

Информатика и её применения

Том 9 Выпуск 1 Год 2015

СОДЕРЖАНИЕ

Моделирование нормальных процессов в стохастических системах со сложными иррациональными нелинейностями И. Н. Сеницын, В. И. Сеницын, Э. Р. Корепанов	2
Метод расчета характеристик интерференции двух взаимодействующих устройств в беспроводной гетерогенной сети Ю. В. Гайдамака, А. К. Самуйлов	9
Heuristic certificates via approximations Sh. Dolev and M. Kogan-Sadetsky	15
Methods and tools for hypothesis-driven research support: A survey L. Kalinichenko, D. Kovalev, D. Kovaleva, and O. Malkov	28
Формальный аксиоматический подход к аспектно-ориентированному расширению технологий программирования С. П. Ковалёв	55
Устойчивые линейные условно оптимальные фильтры и экстраполяторы для стохастических систем с мультипликативными шумами И. Н. Сеницын, Э. Р. Корепанов	70
Выбор оптимальной модели классификации физической активности по измерениям акселерометра М. С. Попова, В. В. Стрижов	76
Оценка погрешности и значимости измерений для линейных моделей С. И. Сливак, О. Г. Кантор, Д. С. Юнусова, С. И. Кузнецов, С. В. Колесов	87
Байесовская рекуррентная модель роста надежности: бета-равномерное распределение параметров Ю. В. Жаворонкова, А. А. Кудрявцев, С. Я. Шоргин	98
К оценке эффективности учебно-познавательной деятельности учащихся с использованием информационных технологий О. М. Корчажкина	106
Об авторах	118
Правила подготовки рукописей	120
Requirements for manuscripts	123

Технический редактор *Л. Кокушкина* Художественный редактор *М. Седакова*
Сдано в набор 12.01.15. Подписано в печать 12.02.15. Формат 60 x 84 / 8
Бумага офсетная. Печать цифровая. Усл.-печ. л. 15,5. Уч.-изд. л. 14,5. Тираж 100 экз.
Заказ № 366к

Издательство «ТОРУС ПРЕСС», Москва 121614, ул. Крылатская, 29-1-43
Отпечатано в Академиздатцентре «Наука» РАН с готовых файлов
Москва 121099, Шубинский пер., д. 6.

МОДЕЛИРОВАНИЕ НОРМАЛЬНЫХ ПРОЦЕССОВ В СТОХАСТИЧЕСКИХ СИСТЕМАХ СО СЛОЖНЫМИ ИРРАЦИОНАЛЬНЫМИ НЕЛИНЕЙНОСТЯМИ*

И. Н. Сеницын¹, В. И. Сеницын², Э. Р. Корепанов³

Аннотация: Рассматриваются дифференциальные стохастические системы (СтС), в том числе и на многообразиях, с винеровскими и пуассоновскими шумами и со сложными иррациональными нелинейностями (СИРН). Такие модели описывают поведение многих современных нано- и квантовооптических технических средств информатики. Приводятся уравнения методов нормальной аппроксимации (МНА) и статистической линеаризации (МСЛ) для аналитического моделирования нестационарных и стационарных нормальных процессов. Рассматриваются методы вычисления типовых интегралов для детерминированных и стохастических одно- и многомерных СИРН скалярного и векторного аргумента. Отмечается возможность использования цилиндрических функций для аналитического расчета интегралов. Обсуждается алгоритмическое обеспечение аналитического и статистического моделирования. Приводится 7 тестовых примеров для типовых СИРН. Рассматривается возможность использования МСЛ для нормализации гиббсовских распределений и распределений с инвариантной мерой для СтС с СИРН.

Ключевые слова: аналитическое и статистическое моделирование; гиббсовское распределение; метод нормальной аппроксимации (МНА); метод статистической линеаризации (МСЛ); распределение с инвариантной мерой; сложные иррациональные нелинейности (СИРН); сложные конечные, дифференциальные и интегральные нелинейности; стохастические системы (СтС); цилиндрические функции

DOI: 10.14357/19922264150101

1 Введение

Вопросам аналитического и статистического моделирования нормальных (гауссовских) стохастических процессов (СтП) в нелинейных СтС (в том числе на многообразиях) на основе МНА и МСЛ посвящена обширная литература (см. обзоры в [1–5]). В [6, 7] дано развитие МНА (МСЛ) для СтС со сложными конечными дифференциальными и интегральными нелинейностями.

Рассмотрим обобщение [6, 7] на случай СтС на многообразиях, содержащих одно- и многомерные детерминированные и стохастические СИРН скалярного и векторного аргумента. Статья состоит из введения, четырех разделов и заключения. В разд. 2 приведены исходные дифференциальные стохастические уравнения (понимаемые в смысле Ито), а также уравнения МНА (МСЛ) для СтС с детерминированными и стохастическими СИРН. Алгоритмическое обеспечение аналитического и статистического моделирования рассмотрено в разд. 3. Типовые интегралы и тестовые примеры для типовых СИРН содержатся в разд. 4. В разд. 5 обсуждаются вопросы моделирования нормальных гиббсо-

вских и распределений с инвариантной мерой в гамильтоновых СтС с линейной диссипацией.

2 Уравнения нормальной аппроксимации (статистической линеаризации) для систем со сложными иррациональными нелинейностями

Как известно [1, 2], уравнения конечномерных непрерывных нелинейных систем со стохастическими возмущениями путем расширения вектора состояния СтС могут быть записаны в виде следующего векторного стохастического дифференциального уравнения Ито:

$$dY_t = a(Y_t, t) dt + b(Y_t, t) dW_0 + \int_{R_0} c(Y_t, t, v) P^0(dt, dv), \quad Y(t_0) = Y_0. \quad (1)$$

* Работа выполнена при поддержке РФФИ (проект 15-07-02244).

¹ Институт проблем информатики Российской академии наук, sinitsin@dol.ru

² Институт проблем информатики Российской академии наук, vsinitsin@ipiran.ru

³ Институт проблем информатики Российской академии наук, ekorepanov@ipiran.ru

Здесь Y_t — $(p \times 1)$ -мерный вектор состояния, $Y_t \in \Delta_y$ (Δ_y — многообразие состояний); $a = a(Y_t, t)$ и $b = b(y_t, t)$ — известные $(p \times 1)$ -мерная и $(p \times m)$ -мерная функции Y_t и t ; $W_0 = W_0(t)$ — $(r \times 1)$ -мерный винеровский СтП интенсивности $\nu_0 = \nu_0(t)$; $c(Y_t, t, v)$ — $(p \times 1)$ -мерная функция Y_t, t и вспомогательного $(q \times 1)$ -мерного параметра v ; $\int_{\Delta} dP^0(t, A)$ — центрированная пуассоновская мера, Δ определяемая следующим образом:

$$\int_{\Delta} dP^0(t, A) = \int_{\Delta} dP(t, A) = \int_{\Delta} \nu_P(t, A) dt.$$

В (1) принято: \int_{Δ} — число скачков пуассоновского СтП в интервале времени $\Delta = (t_1, t_2]$; $\nu_P(t, A)$ — интенсивность пуассоновского СтП $P(t, A)$; A — некоторое борелевское множество пространства R_0^q с выколотым началом. Начальное значение Y_0 представляет собой случайную величину (с.в.), не зависящую от приращений $W_0(t)$ и $P(t, A)$ на интервалах времени, следующих за $t_0, t_0 \leq t_1 \leq t_2$, для любого множества A .

В случае аддитивных нормальных (гауссовских) и обобщенных пуассоновских возмущений уравнение (1) имеет вид:

$$\dot{Y} = a(Y_t, t) + b_0(t)V, \quad V = \dot{W}, \quad Y(t_0) = Y_0. \quad (2)$$

Здесь W — СтП с независимыми приращениями, представляющий собой смесь нормального и обобщенного пуассоновского СтП.

Для компонент $\varphi(Y_t, t) = \{a_h, b_{kj}, c_h\}$ функций a, b и c , являющихся СИРН, примем следующие типовые представления:

$$\varphi(Y_t, t) = \sum_{r=1}^p l_{rt}^{\varphi} |Y_{rt} + d_{rt}|^{\alpha_r^{\varphi}} \rho_r(Y_t); \quad (3)$$

$$\varphi(Y_t, t) = \sum_{r=1}^p l_{r,ht}^{\varphi} \prod_{h=1}^p |Y_{ht} + d_{ht}|^{\alpha_{h,r}^{\varphi}} \rho_{h,r}(Y_t), \quad (4)$$

где $l_{rt}^{\varphi}, l_{r,ht}^{\varphi}, \alpha_r^{\varphi}, \alpha_{h,r}^{\varphi}, d_{rt}, d_{ht}$ и $\rho_r(Y_t), \rho_{h,r}(Y_t)$ — параметры и структурные функции СИРН.

Если предположить существование конечных вероятностных моментов второго порядка для моментов времени t_1 и t_2 , то уравнения МНА примут следующий вид [2, 3]:

— для характеристических функций:

$$g_1^N(\lambda; t) = \exp \left[i\lambda^T m_t - \frac{1}{2} \lambda^T K_t \lambda \right]; \quad (5)$$

$$g_{t_1, t_2}^N(\lambda_1, \lambda_2; t_1, t_2) = \exp \left[i\bar{\lambda}^T \bar{m}_2 - \frac{1}{2} \bar{\lambda}^T \bar{K}_2 \lambda \right], \quad (6)$$

где

$$\bar{\lambda} = [\lambda_1^T \lambda_2^T]^T; \quad \bar{m}_2 = [m_{t_1}^T m_{t_2}^T]^T;$$

$$\bar{K}_2 \begin{bmatrix} K(t_1, t_1) & K(t_1, t_2) \\ K(t_2, t_1) & K(t_2, t_2) \end{bmatrix};$$

— для математических ожиданий m_t , ковариационной матрицы K_t и матрицы ковариационных функций $K(t_1, t_2)$:

$$\dot{m}_t = a_1(m_t, K_t, t), \quad m_0 = m(t_0); \quad (7)$$

$$\dot{K}_t = a_2(m_t, K_t, t), \quad K_0 = K(t_0); \quad (8)$$

$$\left. \begin{aligned} \frac{\partial K(t_1, t_2)}{\partial t_2} &= K(t_1, t_2) a_{21}(m_{t_2}, K_{t_2}, t_2)^T, \\ K(t_1, t_1) &= K_{t_1}. \end{aligned} \right\} \quad (9)$$

Здесь приняты следующие обозначения:

$$a_1 = a_1(m_t, K_t, t) = M_{\Delta_y}^N a(Y_t, t); \quad (10)$$

$$a_2 = a_2(m_t, K_t, t) = a_{21}(m_t, K_t, t) + a_{21}(m_t, K_t, t)^T + a_{22}(m_t, K_t, t); \quad (11)$$

$$a_{21} = a_{21}(m_t, K_t, t) = M_{\Delta_y}^N a(Y_t, t) Y_t^{0T}; \quad (12)$$

$$a_{22} = a_{22}(m_t, K_t, t) = M_{\Delta_y}^N \sigma(Y_t, t); \quad (13)$$

$$\sigma(Y_t, t) = b(Y_t, t) \nu_0(t) b(Y_t, t)^T + \int_{R_0^q} c(Y_t, t, v) c(Y_t, t, v)^T \nu_P(t, dv); \quad (14)$$

$$m_t = M_{\Delta_y}^N Y_t, \quad Y_t^0 = Y_t - m_t;$$

$$K_t = M_{\Delta_y}^N Y_0^0 Y_t^{0T}; \quad K(t_1, t_2) = M_{\Delta_y}^N Y_{t_1}^0 Y_{t_2}^{0T};$$

$M_{\Delta_y}^N$ — символ вычисления математического ожидания для нормальных распределений (5) и (6).

Для стационарных СтС нормальные стационарные СтП — если они существуют, то $m_t = \bar{m}$, $K_t = \bar{K}$, $K(t_1, t_2) = k(\tau)$ ($\tau = t_1 - t_2$), — определяются уравнениями [2, 3]:

$$a_1(\bar{m}, \bar{K}) = 0; \quad a_2(\bar{m}, \bar{K}) = 0; \quad (15)$$

$$\left. \begin{aligned} \dot{k}_{\tau}(\tau) &= a_{21}(\bar{m}, \bar{K}) \bar{K}^{-1} k(\tau), \\ k(0) &= \bar{K} \quad (\forall \tau > 0); \\ k(\tau) &= k(-\tau)^T \quad (\forall \tau < 0). \end{aligned} \right\} \quad (16)$$

При этом необходимо, чтобы матрица $a_{21}(\bar{m}, \bar{K}) = \bar{a}_{21}$ была асимптотически устойчивой.

В случае СтС (2) уравнения МНА переходят в уравнения МСЛ Казакова [2, 3], если принять

$$a(Y_t, t) = a_1(m_t, K_t) + k_1^a(m_t, K_t) Y_t^0; \quad (17)$$

$$\left. \begin{aligned} b(Y_t, t) &= b_0(t); \\ \sigma(Y_t, t) &= b_0(t) \nu(t) b_0(t)^T = \sigma_0(t); \end{aligned} \right\} \quad (18)$$

$$k_1^a(m_t, K_t, t) = \left[\left(\frac{\partial}{\partial m_t} \right) a_0(m_t, K_t, t)^T \right]^T; \quad (19)$$

$$\dot{m}_t = a_1(m_t, K_t, t); \quad m_0 = m(t_0); \quad (20)$$

$$\dot{K}_t = k_1^a(m_t, K_t, t)K_t + K_t k_1^a(m_t, K_t, t)^T + \sigma_0(t), \quad K_0 = K(t_0); \quad (21)$$

$$\frac{\partial K(t_1, t_2)}{\partial t_2} = K(t_1, t_2)K_{t_2} k_1^a(m_{t_2}, K_{t_2}, t_2)^T, \\ K(t_1, t_2) = K_{t_1}. \quad (22)$$

Для стационарных СтС (2) при условии асимптотической устойчивости матрицы $k_1^a(\bar{m}, \bar{K})$ в основе МСЛ лежат уравнения (15) и (16), записанные в виде:

$$a_1(\bar{m}, \bar{K}) = 0; \quad (23)$$

$$k_1^a(\bar{m}, \bar{K})\bar{K} + \bar{K} k_1^a(\bar{m}, \bar{K})^T + \bar{\sigma}_0 = 0; \quad (24)$$

$$\left. \begin{aligned} \dot{k}_\tau(\tau) &= k_1^a(\bar{m}, \bar{K})k(\tau), \quad k(0) = \bar{K} \quad (\forall \tau > 0); \\ k(\tau) &= k(-\tau)^T \quad (\forall \tau < 0). \end{aligned} \right\} \quad (25)$$

Уравнения (5)–(9) лежат в основе МНА для СтС (1), а уравнения (17)–(22) — в основе МСЛ для СтС (2). Для определения стационарных СтП согласно МНА служат соотношения (15) и (16), а МСЛ — (23)–(25).

В задачах практики для СтС со стохастическими СИРН a, b и c являются нормальными случайными функциями состояния Y_t и времени:

$$A = A(Y_t, t); \quad B = B(Y_t, t); \quad C = C(Y_t, t, v).$$

При известных условиях [1–3] они могут быть представлены в виде соответствующих канонических разложений по случайным нормальным скалярным с.в. В этом случае в представлениях (3) и (4) величины l_{rt}^φ и $l_{r,ht}^\varphi$ будут независимыми с.в. L_{rt}^φ и $L_{r,ht}^\varphi$ с известными математическими ожиданиями $m_{rt}^{L_\varphi}$ и $m_{r,ht}^{L_\varphi}$ и дисперсиями $D_{rt}^{L_\varphi}$ и $D_{r,ht}^{L_\varphi}$. Введя расширенный вектор состояния \bar{Y}_t , состоящий из Y_t и с.в. L_{rt}^φ и $L_{r,ht}^\varphi$, и добавив к уравнениям (1) уравнения $dL_{rt}^\varphi = 0$ и $dL_{r,ht}^\varphi = 0$, придем к СтС вида (1).

3 Алгоритмическое обеспечение аналитического и статистического моделирования

Как следует из уравнений (10)–(14), для МНА необходимо уметь вычислять следующие интегралы:

$$I_0^a = I_0^a(m_t, K_t, t) = a_1(m_t, K_t, t) = M_{\Delta_y}^N a(Y_t, t); \quad (26)$$

$$I_1^a = I_1^a(m_t, K_t, t) = a_{21}(m_t, K_t, t) = M_{\Delta_y}^N a(Y_t, t)Y_t^{0T}; \quad (27)$$

$$I_0^{\bar{\sigma}} = I_0^{\bar{\sigma}}(m_t, K_t, t) = a_{22}(m_t, K_t, t) = M_N \bar{\sigma}(Y_t, t). \quad (28)$$

Для МСЛ достаточно вычислить интеграл (26), причем интеграл I_1^a вычисляется по формуле [2–4]:

$$k_1^a = k_1^a(m_t, K_t, t) = \left[\left(\frac{\partial}{\partial m_t} \right) I_0^a(m_t, K_t, t)^T \right]^T.$$

В основе аналитического вычисления интегралов (26)–(28) лежат методы, основанные на свойствах цилиндрических функций [8–10]. Для численных расчетов этих интегралов используются методы [10].

Важно иметь в виду, что уравнения МНА (МСЛ) содержат интегралы I_0^a , I_1^a и $I_0^{\bar{\sigma}}$ в виде соответствующих коэффициентов. Поэтому процедура вычисления интегралов должна быть согласована с методом численного решения обыкновенных дифференциальных уравнений для m_t , K_t и $K(t_1, t_2)$. Эти коэффициенты допускают дифференцирование по m_t и K_t , так как под интегралом стоит сглаживающая нормальная плотность.

В [5] изложены общие алгоритмы аналитического и статистического моделирования распределений, в том числе нормальных в нелинейных СтС на многообразиях. Алгоритмы аналитического, статистического моделирования для СтС с СИРН, а также смешанные алгоритмы различной степени точности относительно шага интегрирования представлены в [5].

4 Тестовые примеры

1. Рассмотрим вычисление интегралов (26) и (27) для одномерных СИРН скалярного аргумента

$$\varphi(Y_t, t) = |Y_t|^{\alpha-1} \text{sign } Y_t \quad (29)$$

(α — нецелый показатель).

Пользуясь (17), представим (29) в виде

$$\varphi(Y_t, t) = \varphi_0(m_t, D_t, t) + k_1^\varphi(m_t, D_t, t)Y_t^0.$$

Здесь введены следующие обозначения:

$$\varphi_0(m_t, D_t, t) = \sqrt{\frac{2}{\pi}} \Gamma(\alpha) D_t^{(\alpha-1)/2} e^{-\xi_t^2/4} D_{-\alpha}(-\xi_t); \quad (30)$$

$$k_1^\varphi(m_t, D_t, t) = \frac{\partial \varphi_0(m_t, D_t, t)}{\partial m_t} = \sqrt{\frac{2}{\pi}} \Gamma(\alpha) D_t^{(\alpha/2)-1} e^{-\xi_t^2/4} D_{1-\alpha}(-\xi_t), \quad (31)$$

где $\Gamma(\alpha)$ — гамма-функция, $\xi_t = m_t/\sqrt{D_t}$ — отношение «сигнал–шум»; $D_{-\alpha}(\xi_t)$ — функция параболического цилиндра [8]. При вычислении были учтены следующие соотношения [8, 10]:

$$\int_0^{\infty} x^{\alpha-1} e^{-\beta x^2 - \gamma x} dx = (2\beta)^{-\alpha/2} \Gamma(\alpha) \exp\left(\frac{\gamma^2}{8\beta}\right) D_{-\alpha}\left(\frac{\gamma}{\sqrt{2\beta}}\right) \quad (\text{Re } \beta > 0, \text{ Re } \alpha > 0); \quad (32)$$

$$\frac{dD_\nu(z)}{dz} = -\frac{z}{2} D_\nu(z) - D_{\nu-1}(z) = \frac{z}{2} D_\nu(z) - D_{\nu+1}(z) \quad (\text{Re } \beta > 0, \text{ Re } \alpha > 0, \nu = -\alpha). \quad (33)$$

Соотношения (32) и (33) могут быть использованы также для вычисления интегралов (28).

Применяя известное асимптотическое разложение для функции параболического цилиндра $D_\nu(z)$ [8, 10]

$$D_\nu(z) \approx e^{-z^2/4} z^\nu \left[1 - \frac{\nu(\nu-1)}{2z^2} + \frac{\nu(\nu-1)(\nu-2)(\nu-3)}{2 \cdot 4z^4} + \dots \right] \quad (|z| \gg 1, |z| \gg \nu),$$

придем к следующим асимптотическим выражениями для $\varphi_0(\xi_t, D_t)$ и $k_1^\varphi(\xi_t, D_t)$:

$$\varphi_0(\xi_t, D_t) = \sqrt{\frac{2}{\pi}} \Gamma(\alpha) D_t^{(\alpha-1)/2} e^{-\xi_t^2/2} |\xi_t|^{-\alpha} \left[1 - \frac{\alpha(\alpha+1)}{2\xi_t^2} + \frac{\alpha(\alpha+1)(\alpha+2)(\alpha+3)}{8\xi_t^4} + \dots \right] > 0;$$

$$k_1^\varphi(\xi_t, D_t) = \frac{\partial \varphi_0(\xi_t, D_t)}{\partial \xi_t} \frac{1}{\sqrt{D_t}}.$$

2. Для одномерной СИРН

$$\varphi(Y_t) = |Y_t|^{\alpha-1} \mathbf{1}(Y_t) \quad (34)$$

имеем

$$\varphi_0(\xi_t, D_t) = \frac{1}{\sqrt{2\pi}} \Gamma(\alpha+1) D_t^{\alpha/2} e^{-\xi_t^2/4} D_{-(\alpha+1)}(-\xi_t);$$

$$k_1^\varphi(\xi_t, D_t) = \frac{\partial \varphi_0(\xi_t, D_t)}{\partial \xi_t} \frac{1}{\sqrt{D_t}}.$$

3. Для нелинейной одномерной дифференциальной СтС с СИРН вида

$$\dot{Y}_t = a_{0t} - a_t |Y_t|^{\alpha-1} \text{sign } Y_t + b_t V \quad (35)$$

(α, a_t, b_t — коэффициенты; V — нормальный белый шум интенсивности ν_t) МСЛ приводит к следующим уравнениям:

$$\dot{m}_t = a_{0t} - a_t \varphi_0(m_t, D_t);$$

$$\dot{D}_t = -2a_t k_1^\varphi(m_t, D_t) + b_t^2 \nu_t.$$

Здесь φ_0 и k_1^φ определены по (30) и (31).

Для статистической системы (35) при $a_t = \bar{a} > 0$ имеет место стационарный режим, для которого $m_t = \bar{m}, D_t = \bar{D}$ определяются из уравнений

$$\bar{a}_0 - \bar{a} \varphi_0(\bar{m}, \bar{D}) = 0; \quad -2\bar{a} k_1^\varphi(\bar{m}, \bar{D}) + \bar{b}^2 \bar{\nu} = 0.$$

4. Для одномерной дифференциальной СтС с СИРН и параметрическим шумом

$$\dot{Y}_t = a_{0t} - a_t Y_t + b_t \left(|Y_t|^{\delta-1} \text{sign } Y_t \right) V$$

с учетом результатов примера 1 согласно МНА имеем:

$$\dot{m}_t = a_{0t} - a_t m_t;$$

$$\dot{D}_t = -2a_t D_t + b_t^2 \nu_t \varphi_0(m_t, D_t).$$

Здесь φ_0 определяется по (30) при $\delta = 2\alpha - 1$.

5. Для некоторых типовых одномерных нелинейных дифференциальных СтС с СИРН вида

$$\dot{Y}_t = \psi(Y_t) + b_t V,$$

допускающих стационарный режим $m_t = \bar{m} = 0, D_t = \bar{D} \neq 0$, выражения для $\bar{\psi}_0 = \bar{\psi}_0(0, \bar{D})$ имеют соответственно вид:

$$\psi(Y_t) = Y_t (Y_t^2 + d^2)^{-1/2};$$

$$\bar{\psi}_0 = \sqrt{\frac{\pi}{\mu}} e^{d^2 \mu} \left[1 - \tilde{\Phi}(d\sqrt{\mu}) \right], \quad (36)$$

где

$$\tilde{\Phi}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt, \quad \mu = \frac{1}{2\bar{D}};$$

$$\psi(Y_t) = Y_t^{\alpha-1} \exp(-\gamma Y_t^{-2});$$

$$\bar{\psi}_0 = 2 \left(\frac{\gamma}{\beta} \right)^{\alpha/4} K_{\alpha/2}(2\sqrt{\beta\gamma}), \quad (37)$$

где $K_z(x)$ — цилиндрическая функция мнимого аргумента, $\beta = 1/(2\bar{D})$;

$$\psi(Y_t) = Y_t^{\mu-1} \sin \gamma Y_t;$$

$$\bar{\psi}_0 = \gamma e^{-\gamma^2/4\beta} \Gamma\left(\frac{1+\mu}{2}\right) {}_1F_1\left(-\frac{\mu}{2}; \frac{3}{2}; \frac{\gamma^2}{4\beta}\right), \quad (38)$$

где ${}_1F_1$ — вырожденная гипергеометрическая функция, $\beta = 1/\bar{D} > 0$, $\text{Re } \mu > -1$;

$$\begin{aligned} \psi(Y_t) &= Y_t^{\mu-1} \cos aY_t; \\ \bar{\psi}_0 &= \frac{\Gamma(\mu/2)}{\beta^{\mu/2}} {}_1F_1\left(\frac{\mu}{2}; \frac{1}{2}; -\frac{a^2}{4\beta}\right) \end{aligned} \quad (39)$$

при $\beta = 1/(2\bar{D})$, $\text{Re } \mu > 0$, $a > 0$;

$$\begin{aligned} \psi(Y_t) &= Y_t^{2\mu-1} \text{sh } \gamma Y_t; \\ \bar{\psi}_0 &= \Gamma(2\mu)(2\beta)^{-\mu} \exp\left(\frac{\gamma^2}{8\beta}\right) \left[D_{-2\mu}\left(-\frac{\gamma}{\sqrt{2\beta}}\right) - \right. \\ &\quad \left. - D_{-2\mu}\left(\frac{\gamma}{\sqrt{2\beta}}\right) \right] \end{aligned} \quad (40)$$

при $\beta = 1/(2\bar{D})$, $\mu > -0,5$;

$$\begin{aligned} \psi(Y_t) &= Y_t^{2\mu-1} \text{ch } \gamma Y_t; \\ \bar{\psi}_0 &= \Gamma(2\mu)(2\beta)^{-\mu} \exp\left(\frac{\gamma^2}{8\beta}\right) \left[D_{-2\mu}\left(-\frac{\gamma}{\sqrt{2\beta}}\right) + \right. \\ &\quad \left. + D_{-2\mu}\left(\frac{\gamma}{\sqrt{2\beta}}\right) \right] \end{aligned} \quad (41)$$

при $\beta = 1/(2\bar{D})$, $\mu > 0$.

6. Статистическая линеаризация одномерной СИРН двумерного аргумента $Y_t = [Y_{1t} Y_{2t}]^T$ для независимых нормальных Y_{it} ($i = 1, 2$) определяется следующими выражениями:

$$\varphi(Y_{1t}, Y_{2t}) = \varphi_0 + k_{11}^\varphi Y_1^0 + k_{12}^\varphi Y_2^0, \quad k_{1i}^\varphi = \frac{\partial \varphi_0}{\partial m_{it}}.$$

Здесь φ_0 , k_{it}^φ зависят от m_{it} и D_{it} . Выражения для φ_0 в случаях (29), (34), (36)–(41) получаются путем перемножения соответствующих выражений для Y_{it} . Аналогично рассматривается случай СИРН n -мерного аргумента ($n > 2$).

7. Пусть одномерная гауссовская стохастическая СИРН аппроксимирована согласно [3] отрезком канонического разложения

$$A(Y_t) = \sum_{j=1}^{n_j} L_j |Y_t|^{\alpha_j-1} \text{sign } Y_t,$$

где L_j — независимые между собой и от Y_t гауссовские с.в. с математическими ожиданиями m^{L_j} и дисперсиями D^{L_j} . Тогда МСЛ приводит к следующему представлению:

$$A(Y_t) = \varphi_0^A + k_1^A Y_t^0 + \sum_{j=1}^{n_j} k_1^{L_j} L_j^0.$$

Здесь коэффициенты φ_0^A , $k_1^{L_j}$ и k_1^A зависят от m_t^y , m^{L_j} , D_t^y и D^{L_j} и вычисляются по формулам примера 1.

5 Применение к моделированию нормальных распределений с инвариантной мерой

В [11] описан ряд многомерных гамильтоновских систем типа Холта, Фокиса–Лагерстрема и др., содержащих СИРН. Для аналитического моделирования нормальных процессов, аппроксимирующих гиббсовские и распределения с инвариантной мерой, если учесть линейные диссипативные силы и нелинейные и параметрические стохастические возмущения, могут быть непосредственно применимы методы [1, 2, 12].

6 Заключение

Рассматриваются дифференциальные СтС, в том числе и на многообразиях, с винеровскими и пуассоновскими шумами и с СИРН. Такие модели описывают поведение многих современных нано- и квантовооптических технических средств информатики.

Приводятся уравнения МНА и МСЛ для аналитического моделирования нестационарных и стационарных нормальных процессов.

Рассматриваются методы вычисления типовых интегралов для детерминированных и стохастических одно- и многомерных СИРН скалярного и векторного аргумента. Отмечается возможность использования цилиндрических функций для аналитического расчета интегралов. Обсуждается алгоритмическое обеспечение аналитического и статистического моделирования.

Приводятся 7 тестовых примеров для типовых СИРН.

Рассматривается возможность использования МСЛ для нормализации гиббсовских распределений и распределений с инвариантной мерой для СтС с СИРН.

Результаты допускают обобщение на случай дискретных, интегродифференциальных и операторных СтС и СИРН, приводимых к дифференциальным, в том числе с автокоррелированными шумами. Как отмечалось в [6, 7], особый интерес представляет развитие МНА (МСЛ) на случай сложных алгебраических и трансцендентных нелинейностей.

Литература

1. Пугачев В. С., Синецын И. Н. Стохастические дифференциальные системы. Анализ и фильтрация. — М.: Наука, 1990. 632 с. (Pugachev V. S., Sinityn I. N.)

- Stochastic differential systems. Analysis and filtering. — Chichester, New York: John Wiley, 1987. 549 p.)
2. Пугачев В. С., Синицын И. Н. Теория стохастических систем. — М.: Логос, 2000; 2004. 1000 с. (Pugachev V. S., Sinitsyn I. N. Stochastic systems. Theory and applications. — Singapore: World Scientific, 2001. 908 p.)
 3. Синицын И. Н. Канонические представления случайных функций и их применение в задачах компьютерной поддержки научных исследований. — М.: ТОРУС ПРЕСС, 2009. 768 с.
 4. Синицын И. Н., Синицын В. И. Лекции по нормальной и эллипсоидальной аппроксимации распределений в стохастических системах. — М.: ТОРУС ПРЕСС, 2013. 488 с.
 5. Синицын И. Н. Параметрическое статистическое и аналитическое моделирование распределений в нелинейных стохастических системах на многообразиях // Информатика и её применения, 2013. Т. 7. Вып. 2. С. 4–16.
 6. Синицын И. Н., Синицын В. И. Аналитическое моделирование нормальных процессов в стохастических системах со сложными нелинейностями // Информатика и её применения, 2014. Т. 8. Вып. 3. С. 2–4.
 7. Синицын И. Н., Синицын В. И., Сергеев И. В., Белюсов В. В., Шоргин В. С. Математическое обеспечение аналитического моделирования стохастических систем со сложными нелинейностями // Системы и средства информатики, 2014. Т. 24. № 3. С. 4–29.
 8. Градштейн И. С., Рыжик И. М. Таблицы интегралов, сумм, рядов и произведений. — М.: ГИФМЛ, 1963. 1100 с.
 9. Справочник по специальным функциям / Под ред. М. Абрамовича и И. Стигана. — М.: Наука, 1979. 832 с.
 10. Попов Б. А., Теслер Г. С. Вычисление функций на ЭВМ: Справочник. — Киев: Наукова Думка, 1984. 599 с.
 11. Переломов А. М. Интегрируемые системы классической механики и алгебры Ли. — М.: Наука, 1990. 240 с.
 12. Синицын И. Н. Аналитическое моделирование распределений с инвариантной мерой в стохастических системах с разрывными характеристиками // Информатика и её применения, 2013. Т. 7. Вып. 1. С. 3–11.

Поступила в редакцию 15.01.15

ANALYTICAL MODELING OF NORMAL PROCESSES IN STOCHASTIC SYSTEMS WITH COMPLEX IRRATIONAL NONLINEARITIES

I. N. Sinitsyn, V. I. Sinitsyn, and E. R. Korepanov

Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation

Abstract: Stochastic systems (including manifolds) with Wiener and Poisson noises and complex irrational nonlinearities (CIRN) are considered. Equations and algorithms of analytical modeling based on the normal approximation method (NAM) and the statistical linearization method (SLM) are given. Typical integrals and software based on cylindrical functions for computing deterministic and stochastic CIRN are presented. Seven test examples for typical CIRN are given. Applications to Gibbs distributions and distributions with invariant measure are discussed.

Keywords: analytical and statistical modeling; complex irrational nonlinearity (CIRN); normal approximation method (NAM); statistical linearization method (SLM); test examples

DOI: 10.14357/19922264150101

Acknowledgments

The research was supported by the Russian Foundation for Basic Research (project 15-07-02244).

References

1. Pugachev, V. S., and I. N. Sinitsyn. 1987. *Stochastic differential systems. Analysis and filtering*. Chichester, New York: John Wiley. 549 p.
2. Pugachev, V. S., and I. N. Sinitsyn. 2001. *Stochastic systems. Theory and applications*. Singapore: World Scientific. 908 p.
3. Sinitsyn, I. N. 2009. Kanonicheskie predstavleniya sluchaynykh funktsiy i ikh primeneniye v zadachakh kom-

- p'yuternoy podderzhki nauchnykh issledovaniy [Canonical expansions of random functions and their application to scientific computer-aided support]. Moscow: TORUS PRESS. 768 p.
4. Sinitsyn, I. N., and V. I. Sinitsyn. 2013. Lektsii po normal'noy i ellipsoidal'noy approksimatsii raspredeleniy v stokhasticheskikh sistemakh [Lectures on normal and ellipsoidal approximation of distributions in stochastic systems]. Moscow: TORUS PRESS. 488 p.
 5. Sinitsyn, I. N. 2013. Parametricheskoe statisticheskoe i analiticheskoe modelirovanie raspredeleniy v nelineynykh stokhasticheskikh sistemakh na mnogoobraznykh [Parametric statistical and analytical modeling of distributions in stochastic systems on manifolds]. *Informatika i ee Primeneniya — Inform. Appl.* 7(2):4–16.
 6. Sinitsyn, I. N., and V. I. Sinitsyn. 2014. Analiticheskoe modelirovanie normal'nykh protsessov v stokhasticheskikh sistemakh so slozhnymi nelineynostyami [Analytical modeling of normal processes in stochastic systems with complex nonlinearities]. *Informatika i ee Primeneniya — Inform. Appl.* 8(3):2–4.
 7. Sinitsyn, I. N., V. I. Sinitsyn, I. V. Sergeev, V. V. Belousov, and V. S. Shorgin. 2014. Matematicheskoe obespechenie analiticheskogo modelirovaniya stokhasticheskikh sistem so slozhnymi nelineynostyami [Mathematical software for analytical modeling of stochastic systems with complex nonlinearities]. *Sistemy i Sredstva Informatiki — Systems and Means of Informatics* 24(3):4–29.
 8. Gradshteyn, I. S., and I. M. Ryzhik. 1963. *Tablitsy integralov, summ, ryadov i proizvedeniy* [Tables of integrals, sums, series, and products]. Moscow: GIFML. 1100 p.
 9. Abramovich, M., and I. Stigan, eds. 1979. *Spravochnik po spetsial'nykh funktsiyam* [Handbook on special functions]. Moscow: Nauka. 832 p.
 10. Popov, B. A., and G. S. Tesler. 1984. *Vychislenie funktsiy na EVM. Spravochnik* [Computing of functions]. Kiev: Naukova Dumka. 599 p.
 11. Perelomov, A. M. 1990. *Integriruemye sistemy klassicheskoy mekhaniki i algebry Li* [Integrable systems of classical mechanics and Li algebra]. Moscow: Nauka. 240 p.
 12. Sinitsyn, I. N. 2013. Analiticheskoe modelirovanie raspredeleniy s invariantnoy meroy v stokhasticheskikh sistemakh s razryvnymi kharakteristikami [Analytical modeling of distributions with invariant measure in stochastic systems with discontinuous characteristics]. *Informatika i ee Primeneniya — Inform. Appl.* 7(1):3–11.

Received January 15, 2015

Contributors

Sinitsyn Igor N. (b. 1940) — Doctor of Science in technology, professor, Honored scientist of Russian Federation, Head of Department, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; sinitsin@dol.ru

Sinitsyn Vladimir I. (b. 1968) — Doctor of Science in physics and mathematics, associate professor, Head of Department, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; vsinitsin@ipiran.ru

Korepanov Eduard R. (b. 1966) — Candidate of Science (PhD) in technology, Head of Laboratory, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; ekorepanov@ipiran.ru

МЕТОД РАСЧЕТА ХАРАКТЕРИСТИК ИНТЕРФЕРЕНЦИИ ДВУХ ВЗАИМОДЕЙСТВУЮЩИХ УСТРОЙСТВ В БЕСПРОВОДНОЙ ГЕТЕРОГЕННОЙ СЕТИ*

Ю. В. Гайдамака¹, А. К. Самуйлов²

Аннотация: Одним из показателей качества функционирования современных беспроводных сетей является отношение сигнала к сумме интерференции и шума (SINR, Signal to Interference plus Noise Ratio) в беспроводных каналах связи. Поскольку значение этой характеристики существенно зависит от расстояния между интерферирующими устройствами, задача оценки значения SINR часто сводится к вычислению длины одной из сторон треугольника, в вершинах которого находятся взаимодействующие устройства. В данной статье решается задача нахождения математического ожидания и среднеквадратического отклонения отношения сигнал/интерференция пары взаимодействующих устройств в достаточно общих предположениях о распределении случайных величин (с.в.) расстояний между интерферирующими устройствами. Предложенный метод может быть использован в качестве основы для анализа интерференции в гетерогенной сети с применением различных беспроводных технологий, включая анализ беспроводных взаимодействий оконечных устройств, на которые интерференция оказывает наиболее сильное воздействие.

Ключевые слова: беспроводная сеть; LTE; интерференция; SINR; взаимодействие устройств; D2D

DOI: 10.14357/19922264150102

1 Постановка задачи

В современных беспроводных сетях, построенных на базе технологии LTE (Long Term Evolution), оценка интерференции между взаимодействующими устройствами является одной из основных задач анализа показателей качества функционирования [1, 2]. Под интерференцией понимается взаимодействие сигналов, передаваемых разными источниками на одном и том же канале. Интерференция вызывает искажение сигнала рассматриваемого источника под воздействием сигнала стороннего источника. В гетерогенных сетях беспроводного взаимодействия оконечных устройств D2D (device-to-device) [3], где плотность интерферирующих объектов высока, интерференция оказывает существенное влияние на принимаемый оконечным устройством сигнал. При анализе беспроводных взаимодействий устройств обычно рассматривается несколько источников сигнала (передатчиков), распределенных на плоскости согласно некоторому закону [4]. Упрощение задачи состоит в том, что, рассмотрев один передатчик и оценив характеристики интерференции на соответствующем ему приемном устройстве (приемнике), можно предположить, что основные показатели

будут идентичны и для остальных пар «передатчик–приемник». В данной статье решается задача нахождения числовых характеристик отношения сигнал/интерференция пары взаимодействующих устройств.

Отношение сигнала к сумме интерференции и шума, SINR, является одной из основных характеристик качества канала в беспроводных сетях связи [5–7]. Отношение сигнала к сумме интерференции и шума на стороне приемника определяется по следующей формуле:

$$\text{SINR} = \frac{S}{\sigma^2 + I}, \quad (1)$$

где S — мощность принимаемого сигнала от соответствующего передатчика; σ^2 — мощность шума; I — мощность принимаемого сигнала от интерферирующих передатчиков. Согласно линейной модели [4]

$$S = gl^{-\alpha}, \quad (2)$$

где g — базовая мощность сигнала передатчика, соответствующего рассматриваемому приемнику; l — расстояние между передатчиком и приемником; α — коэффициент потерь (path loss exponent), принимающий значение от 2 (при условии прямой

* Работа выполнена при финансовой поддержке РФФИ (проекты 14-07-00090 и 15-07-03051).

¹ Российский университет дружбы народов, ygaidamaka@sci.pfu.edu.ru

² Российский университет дружбы народов; Технологический университет г. Тампере, Финляндия, aksamuylov@gmail.com

видимости) до 6 (в худшем случае). Величина I в знаменателе формулы (1) соответствует суммарной мощности сигнала от всех интерферирующих передатчиков, где каждое слагаемое имеет вид (2). Заметим, что принцип повторного использования частот (frequency reuse) в беспроводных сетях связи поколения 4G (4th Generation) позволяет назначать одну и ту же единицу ресурса сети (например, один и тот же ресурсный блок LTE) нескольким парам взаимодействующих устройств, если интерференция не превосходит определенного стандартами уровня.

Рассмотрим случай, когда несколько принимающих устройств (приемников) и одно передающее устройство (передатчик), образующие кластер, расположены на плоскости внутри круга радиуса r_0 , причем передатчик расположен в центре круга. Такой кластер образуется, например, при проведении интерактивного занятия преподавателя с учениками, когда можно предположить, что передатчик располагается в центре круга, а приемники равномерно распределены внутри круга. Для передачи данных на каждую пару взаимодействующих устройств внутри кластера планировщиком распределения радиоресурсов в беспроводной сети 4G назначается по одному ресурсному блоку LTE, и тогда сигналы взаимодействующих пар не интерферируют друг с другом. Но если в соседнем помещении также проходит интерактивное занятие и там использованы те же ресурсные блоки, то пары из соседних кластеров, использующие один и тот же ресурсный блок, будут создавать помехи друг другу. Сведем задачу к анализу взаимодействия двух пар устройств в двух кластерах, как показано на рис. 1.

Пару взаимодействующих устройств, для которой будем рассчитывать показатели качества канала, назовем целевой, а соответствующую ей пару устройств обозначим $TR_0 = \langle Tx_0, Rx_0 \rangle$. Остальные пары, которые создают помехи целевой паре TR_0 ,

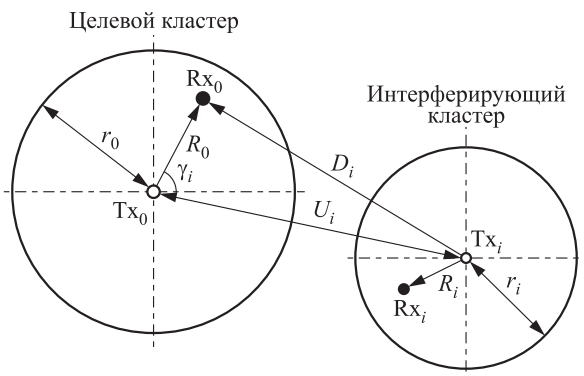


Рис. 1 Схема взаимодействия интерферирующих устройств

обозначим $TR_i = \langle Tx_i, Rx_i \rangle$ и будем называть их интерферирующими. Расстояние между Rx_i и Tx_i обозначим R_i , а расстояние между Tx_0 и Tx_i обозначим U_i . Мощность интерферирующего сигнала от пары TR_i является функцией расстояния между приемником Rx_0 из целевой пары и интерферирующим передатчиком Tx_i , которое обозначим D_i . Угол между прямой, соединяющей целевые передатчик Tx_0 и приемник Rx_0 , и прямой, соединяющей передатчики Tx_0 и Tx_i , обозначим γ_i .

Рассмотрим систему двух кластеров, показанную на рис. 1. В условиях отсутствия шума и одинаковой базовой мощности g сигналов обоих передатчиков искомой характеристикой является отношение сигнал/интерференция SIR для приемника Rx_0 , вычисляемое по формуле:

$$SIR = \left(\frac{D_1}{R_0} \right)^\alpha \tag{3}$$

Будем считать, что R_0 , U_i и γ_i являются с.в. с заданными функциями распределения. Задача состоит в нахождении числовых характеристик с.в. SIR. Для решения задачи в следующем разделе статьи предлагается метод нахождения совместной плотности распределения с.в. R_0 и D_i , что позволяет вычислять начальные моменты $E[SIR^n]$ с.в. SIR.

2 Метод расчета отношения сигнал/интерференция

Как видно из формулы (3), с.в. SIR пропорциональна с.в. D_1 , которая, в свою очередь, зависит от с.в. R_0 . В этом случае для нахождения характеристик с.в. SIR необходимо найти совместное распределение с.в. R_0 и D_1 .

Введем обозначения $\xi_1 := R_0$, $\xi_2 := U_1$, $\xi_3 := \gamma_1$, $\eta_1 := D_1$. Тогда $w_{\xi_1, \xi_2, \xi_3}(x_1, x_2, x_3) := f_{R_0, U_1, \gamma_1}(x_1, x_2, x_3)$ — совместная плотность распределения с.в. R_0 , U_1 и γ_1 , а $W_{\xi_1, \eta_1}(x_1, y_1) := f_{R_0, D_1}(x_1, y_1)$ — искомое совместное распределение с.в. R_0 и D_1 . По теореме косинусов с.в. η_1 является функцией с.в. ξ_1 , ξ_2 и ξ_3 :

$$\eta_1 = \sqrt{\xi_1^2 + \xi_2^2 - 2\xi_1\xi_2 \cos \xi_3} \tag{4}$$

Следуя [8, 9], введя вспомогательную переменную $\eta_2 = \xi_3$, искомое распределение можно найти по следующей формуле:

$$W_{\xi_1, \eta_1}(y_1, y_2) = \sum_{i=1}^2 \int_{Y_{3,j}} w_{\xi_1, \xi_2, \xi_3}(y_1, \varphi_i(y_1, y_2, y_3), y_3) \times \left| \frac{\partial \varphi_j(y_1, y_2, y_3)}{\partial y_2} \right| dy_3, \tag{5}$$

где φ_j — обратное преобразование правой части формулы (4) относительно ξ_2 :

$$\begin{aligned}\varphi_1(y_1, y_2, y_3) &= y_1 \cos y_3 + \sqrt{y_2^2 - y_1^2 + y_1^2 \cos^2 y_3}; \\ \varphi_2(y_1, y_2, y_3) &= y_1 \cos y_3 - \sqrt{y_2^2 - y_1^2 + y_1^2 \cos^2 y_3}.\end{aligned}$$

В формуле (5) области значений $Y_{3,j}$ переменной y_3 для j -й ветви обратного преобразования определяются системой неравенств:

$$\left. \begin{aligned}\varphi_j(y_1, y_2, y_3) &\geq 0; \\ y_1 &\geq 0; \\ y_2 &\geq 0; \\ 0 &\leq y_3 \leq 2\pi.\end{aligned}\right\} \quad (6)$$

Решая систему (6), нетрудно убедиться, что для первой ветви обратного преобразования $Y_{3,1} = Y_{3,1}^1 \cup Y_{3,1}^2 \cup Y_{3,1}^3$, где

$$Y_{3,1}^1 = \left\{ \begin{aligned}0 &\leq y_2 \leq y_1; \\ 0 &\leq y_3 \leq \frac{1}{2} \arccos\left(\frac{y_1^2 - 2y_2^2}{y_1^2}\right); \end{aligned}\right. \quad (7)$$

$$Y_{3,1}^2 = \left\{ \begin{aligned}0 &\leq y_2 \leq y_1; \\ 2\pi - \frac{1}{2} \arccos\left(\frac{y_1^2 - 2y_2^2}{y_1^2}\right) &\leq y_3 \leq 2\pi; \end{aligned}\right. \quad (8)$$

$$Y_{3,1}^3 = \left\{ \begin{aligned}y_2 &\geq y_1; \\ 0 &\leq y_3 \leq 2\pi, \end{aligned}\right. \quad (9)$$

а для второй ветви $Y_{3,2} = Y_{3,2}^1 \cup Y_{3,2}^2$, где

$$Y_{3,2}^1 = \left\{ \begin{aligned}0 &\leq y_2 \leq y_1; \\ 0 &\leq y_3 \leq \frac{1}{2} \arccos\left(\frac{y_1^2 - 2y_2^2}{y_1^2}\right); \end{aligned}\right. \quad (10)$$

$$Y_{3,2}^2 = \left\{ \begin{aligned}0 &\leq y_2 \leq y_1; \\ 2\pi - \frac{1}{2} \arccos\left(\frac{y_1^2 - 2y_2^2}{y_1^2}\right) &\leq y_3 \leq 2\pi. \end{aligned}\right. \quad (11)$$

Таким образом, получена формула для вычисления совместной плотности с.в. R_0 и D_1 :

$$\begin{aligned}W_{\xi_1, \eta_1}(y_1, y_2) &= \\ &= \sum_{i=1}^2 \int_{Y_{3,i}} \frac{w_{\xi_1, \xi_2, \xi_3}(y_1, \varphi_i(y_1, y_2, y_3), y_3) y_2}{\sqrt{y_2^2 - y_1^2 + y_1^2 \cos^2 y_3}} dy_3, \quad (12)\end{aligned}$$

где $Y_{3,j}$ вычисляются по формулам (7)–(11).

В следующем разделе приведен пример численного анализа с использованием формул (7)–(12).

3 Пример численного анализа

В рассматриваемом примере предложенный выше метод использован для расчета начальных моментов $E[SIR^n]$ отношения сигнал/интерференция, которые определяются следующей формулой:

$$\begin{aligned}E[SIR^n] &= \\ &= \int_{0 \leq y_1 \leq r_0} \int_{y_2 \geq 0} \left(\frac{y_2}{y_1}\right)^{n\alpha} W_{\xi_1, \eta_1}(y_1, y_2) dy_2 dy_1. \quad (13)\end{aligned}$$

Рассматривается случай, когда целевой приемник Rx_0 находится внутри круга единичного радиуса ($r_0 = 1$), в центре которого расположен передатчик Tx_0 , а интерферирующий передатчик Tx_1 — в кольце вокруг передатчика Tx_0 с внутренним радиусом r_0 и внешним радиусом h_0 (рис. 2).

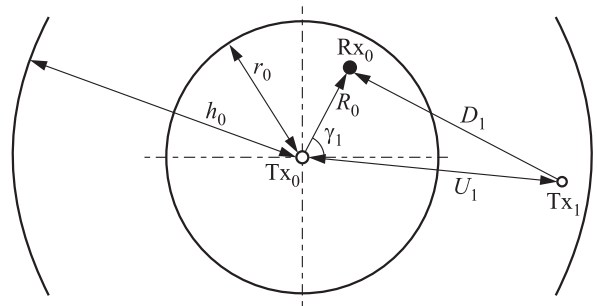


Рис. 2 Пример взаимодействия двух устройств

Тогда с.в. R_0 расстояния от целевого передатчика Tx_0 до соответствующего ему приемника Rx_0 и с.в. U_1 расстояния от целевого передатчика Tx_0 до интерферирующего передатчика Tx_1 имеют распределения

$$\begin{aligned}f_{R_0}(r) &= 2r, \quad 0 \leq r \leq 1; \\ f_{U_1}(u) &= \frac{2u}{h_0^2 - 1}, \quad 1 \leq u \leq h_0.\end{aligned}$$

Будем считать, что с.в. угла γ_1 равномерно распределена на отрезке $[0, 2\pi]$, а коэффициент потерь в формуле (2) принимает значение $\alpha = 2$. Приняты условные единицы измерения: например, расстояние между взаимодействующими устройствами может измеряться в метрах, а величина SIR — в децибелах.

По формулам (7)–(13) рассчитано математическое ожидание отношения сигнал/интерференция $E[SIR]$, представленное в таблице в зависимости от радиуса внешней границы кольца, внутри которого распределены интерферирующие передатчики. В таблице также показаны значения математического ожидания расстояния $E[U_1]$ от целевого передатчика Tx_0 до интерферирующего передатчика Tx_1 .

Математическое ожидание величины SIR

h_0	$E[U_1]$	$E[SIR]$
2	1,56	4,84985
3	2,17	7,41701
4	2,8	9,54562
5	3,44	11,30286

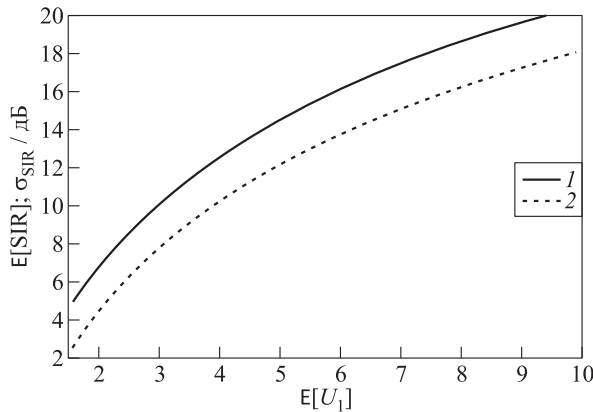


Рис. 3 Числовые характеристики отношения сигнал/интерференция: 1 — $E[SIR]$; 2 — σ_{SIR}

Также были рассчитаны математическое ожидание $E[SIR]$ и среднеквадратическое отклонение $\sigma_{SIR} = \sqrt{E[SIR^2] - E[SIR]^2}$ отношения сигнал/интерференция, показанные на рис. 3 в зависимости от математического ожидания расстояния $E[U_1]$ между целевым передатчиком Tx_0 и интерферирующим передатчиком Tx_1 . Из таблицы и графиков видно, что с ростом расстояния между целевым и интерферирующим передатчиком обе числовые характеристики отношения сигнал/интерференция растут, поскольку мощность интерферирующего сигнала убывает. Вычисления проводились с использованием встроенных средств пакета программ Wolfram Mathematica [10].

4 Заключение

В настоящей статье метод преобразования с.в. применен для анализа основной характеристики качества функционирования беспроводных сетей, а именно: отношения сигнал/интерференция при заданных распределениях расстояний между интерферирующими устройствами. Приведенный пример показывает, что численный анализ является достаточно трудоемким даже в простейших предположениях о распределении исходных с.в., а для оценки характеристик интерференции в условиях

наличия в беспроводной сети нескольких источников интерференции требуется разработка приближенных методов и имитационных моделей, как это сделано, например, в [11]. Задача с несколькими источниками интерференции в беспроводных гетерогенных сетях взаимодействующих устройств представляется особенно актуальной ввиду быстрого развития сетей 4G и принятия в ближайшем будущем стандартов для беспроводных сетей 5G [12].

Авторы выражают благодарность проф. К. Е. Самуйлову за плодотворное обсуждение и ценные советы.

Литература

1. Гайдамака Ю. В., Ефимушкина Т. В., Самуйлов А. К., Самуйлов К. Е. Задачи оптимального планирования межуровневого интерфейса в беспроводных сетях // Информатика и её применения, 2012. Т. 6. Вып. 3. С. 75–81.
2. Basharin G. P., Gaidamaka Yu. V., Samouylov K. E. Mathematical theory of teletraffic and its application to the analysis of multiservice communication of next generation networks // Autom. Control Comp. Sci., 2013. Vol. 47. No. 2. P. 62–69.
3. Andreev S., Pyattaev A., Johnsson K., Galinina O., Koucheryavy Y. Cellular traffic offloading onto network-assisted device-to-device connections // IEEE Commun. Mag., 2014. Vol. 52. No. 4. <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6807935>.
4. Baccelli F., Blaszczyzyn B. Stochastic geometry and wireless networks. Vol. I: Theory. — Boston: NoW Pubs. Inc., 2009. 164 p.
5. Erturk M. C., Mukherjee S., Ishii H., Arslan H. Distributions of transmit power and SINR in device-to-device networks // IEEE Commun. Lett., 2013. Vol. 17. No. 2. <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6472443>.
6. Kim M., Han Y., Yoon Y., Chong Y., Lee H. Modeling of adjacent channel interference in heterogeneous wireless networks // IEEE Commun. Lett., 2013. Vol. 17. No. 9. <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6604524>.
7. Andrews J. G., Singh S., Ye Q., Lin X., Dhillon H. S. An overview of load balancing in hetnets: Old myths and open problems // IEEE Wirel. Commun., 2014. Vol. 21. No. 2. <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6812279>.
8. Левин Б. П. Теоретические основы статистической радиотехники. — 3-е изд. — М.: Радио и связь, 1989. 656 с.
9. Mardia K., Jupp P. Directional statistics. — Wiley Press, 1999. 441 p.

10. Wolfram Mathematica: Программное обеспечение для технических вычислений. <http://www.wolfram.com/mathematica>.
11. Гайдамака Ю. В., Печинкин А. В., Разумчик Р. В., Самуйлов А. К., Самуйлов К. Е., Соколов И. А., Сопин Э. С., Шоргин С. Я. Распределение времени выхода из множества состояний перегрузки в системе $M|M|1\langle L, H \rangle\langle H, R \rangle$ с гистерезисным управлением нагрузкой // Информатика и её применения, 2013. Т. 7. Вып. 4. С. 20–33.
12. Tehrani M., Uysal M., Yanikomeroglu H. Device-to-device communication in 5G cellular networks: Challenges, solutions, and future directions // IEEE Commun. Mag., 2014. Vol. 52. No. 5. <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6815882>.

Поступила в редакцию 20.01.15

METHOD FOR CALCULATING NUMERICAL CHARACTERISTICS OF TWO DEVICES INTERFERENCE FOR DEVICE-TO-DEVICE COMMUNICATIONS IN A WIRELESS HETEROGENEOUS NETWORK

Yu. Gaidamaka¹ and A. Samuylov^{1,2}

¹Peoples' Friendship University of Russia, Applied Probability and Informatics Department, 6 Miklukho-Maklaya Str., Moscow 117198, Russian Federation

²Tampere University of Technology, Department of Electronics and Communications Engineering, 10 Korkeakoulunkatu, Tampere 33720, Finland

Abstract: In wireless networks, one of the key performance metrics is the signal to noise ratio, SINR. As this metric highly depends on the distance between the interfering devices, the problem of SINR estimation is often reduced to the calculation of a triangle's side length, where the vertices represent the interacting devices. This paper addresses the problem of calculating the numerical characteristics of the signal to interference ratio for a pair of interfering devices determined by the known distributions of distances between the entities in question. The proposed method can be used as a basis for analyzing heterogeneous networks, including the analysis of device-to-device (D2D) communications as one of the interference-limited cases.

Keywords: wireless network; LTE; interference; SINR; D2D

DOI: 10.14357/19922264150102

Acknowledgments

The reported study was partially supported by the Russian Foundation for Basic Research, research projects Nos. 14-07-00090 and 15-07-03051.

References

1. Gaidamaka, Yu. V., T. V. Efimushkina, A. K. Samuylov, and K. E. Samouylov. 2012. Zadachi optimal'nogo planirovaniya mezhurovnevnogo interfeysa v besprovodnykh setyakh [Cross-layer optimization planning problems in wireless networks]. *Informatika i ee Primeneniya — Inform. Appl.* 6(3):75–81.
2. Basharin, G. P., Yu. V. Gaidamaka, and K. E. Samouylov. 2013. Mathematical theory of teletraffic and its application to the analysis of multiservice communication of next generation networks. *Autom. Control Comp. Sci.* 47(2):62–69.
3. Andreev, S., A. Pyattaev, K. Johnsson, O. Galinina, and Y. Koucheryavy. 2014. Cellular traffic offloading on-to network-assisted device-to-device connections. *IEEE Commun. Mag.* 52(4). Available at: <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6807935> (accessed January 10, 2015).
4. Baccelli, F., and B. Blaszczyszyn. 2009. *Stochastic geometry and wireless networks*. Vol. 1: Theory. Boston: NoW Publ. Inc. 164 p. January 10, 2015).
5. Erturk, M. C., S. Mukherjee, H. Ishii, and H. Arslan. 2013. Distributions of transmit power and SINR in device-to-device networks. *IEEE Commun. Lett.* 17(2). Available at: <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6472443> (accessed January 10, 2015).
6. Kim, M., Y. Han, Y. Yoon, Y. Chong, and H. Lee. 2013. Modeling of adjacent channel interference in heterogeneous wireless networks. *IEEE Commun. Lett.* 17(9). Available at: <http://ieeexplore.ieee.org/>

- xpl/tocresult.jsp?isnumber=6604524 (accessed January 10, 2015).
7. Andrews, J. G., S. Singh, Q. Ye, X. Lin, and H. S. Dhillon. 2014. An overview of load balancing in hetnets: Old myths and open problems. *IEEE Wirel. Commun.* 21(2). Available at: <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6812279> (accessed January 10, 2015).
 8. Levin, B. R. 1989. *Teoreticheskie osnovy statisticheskoy radiotekhniki* [Theoretical basis of statistical radiotechnics]. 3rd ed. Moscow: Radio and Communications. 656 p.
 9. Mardia, K., and P. Jupp. 1999. *Directional statistics*. 1st ed. Wiley Press. 441 p.
 10. Wolfram mathematica: Software for technical computing. [Free access] Available at: <http://www.wolfram.com/mathematica> (accessed December 1, 2014).
 11. Gaidamaka, Yu. V., A. V. Pechinkin, R. V. Razumchik, A. K. Samuylov, K. E. Samouylov, I. A. Sokolov, E. S. Sopin, and S. Ya. Shorgin. 2013. Raspredelenie vremeni vykhoda iz mnozhestva sostoyaniy peregruzki v sisteme $M|M|1|\langle L, H \rangle|\langle H, R \rangle$ s gisterezisnym upravleniem nagruzkoy [The distribution of the return time from the set of overload states to the set of normal load states in a system $M|M|1|\langle L, H \rangle|\langle H, R \rangle$ with hysteretic load control]. *Informatika i ee Primeneniya — Inform. Appl.* 7(4):20–33.
 12. Tehrani, M., M. Uysal, and H. Yanikomeroglu. 2014. Device-to-device communication in 5G cellular networks: Challenges, solutions, and future directions. *IEEE Commun. Mag.* 52(5). Available at: <http://ieeexplore.ieee.org/xpl/tocresult.jsp?isnumber=6815882> (accessed January 10, 2015).

Received January 20, 2015

Contributors

Gaidamaka Yuliya V. (b. 1971) — Candidate of Science (PhD) in physics and mathematics, associate professor, Applied Probability and Informatics Department, Peoples' Friendship University of Russia, 6 Miklukho-Maklaya Str., Moscow 117198, Russian Federation; ygaidamaka@sci.pfu.edu.ru

Samuylov Andrey K. (b. 1988) — PhD student, Peoples' Friendship University of Russia, Moscow 117198, Russian Federation; researcher, Department of Electronics and Communications Engineering, Tampere University of Technology, 10 Korkeakoulunkatu, Tampere 33720, Finland; aksamuylov@gmail.com

HEURISTIC CERTIFICATES VIA APPROXIMATIONS

Sh. Dolev¹ and M. Kogan-Sadetsky¹

Abstract: This paper suggests a new framework in which the quality of a (not necessarily optimal) heuristic solution is certified by an approximation algorithm. Namely, a result of a heuristic solution is accompanied by a scale obtained from an approximation algorithm. The creation of a scale is efficient while getting a solution from an approximation algorithm is usually concerned with long calculation relatively to heuristics approach. On the other hand, a result obtained by heuristics without scale might be useless. The criteria for choosing an approximation scheme for producing a scale have been investigated. To obtain a scale in practice, not only approximations have been examined by their asymptotic behavior but also relations as a function of an input size of a given problem. For study case only, heuristic and approximation algorithms for the SINGLE KNAPSACK, MAX 3-SAT, and MAXIMUM BOUNDED THREE-DIMENSIONAL MATCHING (MB3DM) NP-hard problems have been examined. The certificates for the heuristic runs have been obtained by using fitting approximations.

Keywords: heuristics; approximation algorithm; optimal solution; approximation preserving reducibility

DOI: 10.14357/19922264150103

1 Introduction

Many important optimization problems are NP-hard. This rules out the possibility of finding in polynomial time an optimal solution (for those that are in NP, unless $P = NP$). Instead, commonly, artificial intelligence (AI) methods are used to try to cope with instances of these problems. Namely, heuristic solutions are used to get a near to optimal solution for a given instance of such problems. However, heuristic may obtain a suboptimal solution such as a local maxima rather than the absolute maximum. Evaluation of the result obtained from the heuristic is the complicated task that is addressed in this work. The present authors suggest using approximation algorithms to obtain a scale that may be used to certify heuristic results.

Heuristic solutions. Heuristic solutions explore various states and solutions for a given problem in a short time². There are well-known heuristics such as: Hill Climbing search that uses local estimates of the distance to the goal and expands the node with the smallest estimate; Branch and Bound search that uses the costs of the already formed partial paths and expands the partial path with the minimum cost, and Genetic Algorithms (GA). In particular, GA will be explored as an example for a heuristic. Genetic algorithms were invented by John Holland [1] and others, finding their inspiration in the evolutionary process occurring in nature. A population of individuals should collectively adapt to conditions of some environment. In order to face this, the reproduction and survival of individuals are promoted by

the elimination of useless and harmful features and by gratifying useful behavior.

The dynamics of the natural evolutionary process are complex and mostly unknown. Nevertheless, the artificial optimization approach is feasible if in comparison to another approaches, a good enough solution can be obtained. One of such different approaches is the solution that uses approximation algorithms rather than heuristics.

Approximation algorithms. A way to cope with optimization tasks that do not have a (known) polynomial solution is by devising an algorithm that finds a solution to the problem in a reasonable time, such that the solution is suboptimal with a known ratio or bound from the (unknown) optimal solution. On one hand, heuristics may perform better than approximations in practice, namely, obtain a better solution with the same processing effort, but on the other hand, heuristics are associated with uncertainty regarding the quality of the obtained result.

The authors' contribution. The authors suggest a new framework in which the quality of a solution that is based on heuristics is certified by an approximate algorithm. In this framework, a heuristic solution is accompanied by a scale obtained from the approximation algorithm. The usefulness of the framework is demonstrated by presenting criteria for choosing a fitting approximation scheme for producing a scale. The SINGLE KNAPSACK, MAX 3-SAT, and MB3DM NP-hard problems are examined.

¹Department of Computer Science, Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel

²The authors do not refer to the A^* heuristic search, in which some optimizations for the exponential time exhaustive search are introduced for obtaining the optimal solution at the price of worst case exponential time.

Using polynomial time approximation schemes. In some cases, it is possible to use polynomial time approximation schemes (PTAS). Let start by demonstrating the framework for such a case, the case of the SINGLE KNAPSACK problem. To solve instances of SINGLE KNAPSACK, galib246 [2] library for genetic algorithms has been used. In order to certificate the solution obtained from the genetic algorithm, the approximation algorithm CKPP (Cardinality Constrained Knapsack Problem) [3], has been used. In the present authors' experiment [4], each item has weight and profit; the restricted total weight was 300 weight units and the optimization task was to maximize the profit while respecting the weight bound. The profit obtained by the genetic algorithm was 577 profit units (with weight of 293 weight units). The total profit of all the given items was much larger – 1,345 profit units. Thus, the authors had no clue whether the obtained profit is a local maxima that is way beyond the possible maximum profit. CKPP approximation certified the heuristic solution as 0.88 of the possible maximum profit. Since CKPP is PTAS, it can be tuned to obtain a refined scale using more and more computation time. Using the performance ratio of 0.75, the profit of 493 profit units (with weight of 293 weight units) running in a few seconds was obtained.

Using approximation preserving reducibility. Obtaining certificates as a part of the practice in using heuristics may require the design of many approximations, which is not a simple task. Luckily, in some cases, approximation preserving reductions may be used to obtain an approximation algorithm from variant of a problem in hand to a known approximation algorithm. Next, the way such a reduction is chosen and used will be demonstrated. The experimental demonstration for this example obtained a 0.64 certificate for the GA result [4]. **MAX 3-SAT.** An instance of the MAX 3-SAT problem is a predicate for which each clause is of not more than 3 literals. A solution is a truth assignment for the predicate variables, such that the number of satisfied clauses is the maximum possible. One may evaluate a given assignment for the predicate variables by the number of clauses it satisfies. Once again, galib246 library [2] was used to solve instances of the MAX 3-SAT problem. In their experimental demonstration, the present authors used a predicate that consisted of 100 clauses. The heuristic algorithm found (in a few seconds) an assignment for the variables that satisfies 58 clauses. Again, the authors had no clue whether the solution was only 0.58 of the optima.

Some maximization problems have trivial bound on solution size which is related to the size of the input. This bound gives one a preliminary Certificate of the heuristic solution. This certificate maybe good enough if the heuristic solution is at least, say, 0.8 of the input size. This implies that an optimal solution for the given

problem is at least 0.8 of the input size, which is not a common case. For instance, in the case of the MAX 3-SAT problem, the trivial upper bound is obvious since all predicate's clauses may be satisfiable by some truth assignment. But it is possible that any truth assignment can satisfy near to half of predicate's clauses, and in this case, the trivial upper bound gives a poor certificate to a heuristic solution. For some problems, it is not easy to determine a trivial upper bound. For example, for the SINGLE KNAPSACK problem, the common case is that only a small part of the items can be put into the knapsack.

To demonstrate the way reductions are chosen, two examples are presented: MAX 3-SAT and MB3DM.

For the first example, two reductions from MAX 3-SAT were used: one to MAX2-SAT and the second to MAX INDEPENDENT SET-B. In this case, both approximations are not PTAS and, therefore, only a single scale can be obtained for a given instance. It turned out that the scale obtained for the particular instance, by using the reduction to MAX 2-SAT, added no information beyond the trivial bound of 100 clauses, while the reduction to MAX INDEPENDENT SET-B resulted in a refined scale of 92 clauses.

Also, two reductions were used from MB3DM: one to minimization of sum of squared machine loads and the second to STAR-GCA-SIMPLE. In this case, the reductions once again, both gave a single scale. It will be shown that according to the all comparable criteria for choosing reduction – quality of approximation preserve, quality of approximation algorithm of a destination problem, inflation of input, and the complexity of reduction (that would be defined later in the paper), the reduction to MINIMIZATION OF SUM OF SQUARED MACHINE LOADS is a better choice for producing a scale and a certificate.

Choosing approximation preserving reductions. There are four main criteria that are suggested using for choosing a specific reduction algorithm from a given problem. The first criterion is the quality of approximation preserve. This is a reduction preservation ratio, which is implied by the parameters of reduction constrains. The second criterion is the quality of approximation algorithm of a destination problem. This means the performance ratio of (best) known approximate algorithm for a destination problem. The third criterion is the inflation of input of a source problem when it is translated to input of a destination problem. The last criterion is the complexity of reduction functions. This parameter represents degree of ease of a reduction usage and also contributes to the general performance of a reduction.

Paper organization. In section 2, the framework will be presented and the way approximation for the SINGLE KNAPSACK problem is used will be demonstrated to certify a solution obtained by a genetic algorithm. An

overview on Genetic Algorithms (the authors' choice for examining heuristics) appears in section 3. Then, in section 4, a short overview on approximation algorithms and approximation preserving reductions will be given including: PTAS, L-reduction, and AP-reduction. Section 5 discusses criteria for choosing an approximation preserving reduction to obtain the best scale. Section 6 concludes the paper.

2 The Framework and the SINGLE KNAPSACK Problem Example

The framework. First, let define some concepts that will be used in the paper:

- *Input*: an instance of hard optimization problem for which there is no (efficient) exact polynomial algorithm or good polynomial approximation.
- *Heuristic certificate*. Given a solution of some heuristic algorithm, the quality of this solution will be determined, i. e., a (closest) relation between the heuristic and the optimal solutions of a given problem instance will be found.
- *Solution scale*. To define a solution scale, a problem solution upper bound obtained from approximation algorithm is used. This scale is used to measure the quality of the heuristic solution.

– Certificate

$$= \frac{\text{Heuristic solution}}{\text{upper bound (obtained from approximations)}}$$

The steps defined by the framework are as follows:

- Choose heuristics, for example, GA to solve the problem instance.
- Find (at least) one approximation for the problem. Sometimes, approximation preserving reductions should be used.
- Choose, using the criteria that will be defined later in the paper, the approximation that leads to the best solution scale.
- Execute the heuristic and the approximation algorithms (possibly, in parallel) and compute the certificate according to the results (possibly, repeatedly, in particular, in the case of PTAS and FULLY PTAS, until the required certificated result is obtained).

SINGLE KNAPSACK problem example. Let consider the SINGLE KNAPSACK problem. Given items of different profits and weights, find the most valuable set of items that fit in a knapsack of fixed weight bound. This problem is chosen since it has PTAS, i. e., the algorithm that for each $\epsilon > 0$ produces a solution that is within ϵ factor of being optimal. The approximation

Table 1 SINGLE KNAPSACK approximation solution

ϵ	Approximation ratio	Solution weight	Solution profit
0.4	2/3	293	493

Table 2 Genetic algorithm solution

Generation id	Solution weight	Solution profit	Optimal solution upper bound	Heuristic certificate
12	293	577	740	0.78

algorithm used is the algorithm CKPP [3]. Its running time is $O(n^{\lceil 1/\epsilon \rceil - 2})$ for $\epsilon < 1/2$, and its guaranteed approximation ratio is $(\lceil 1/\epsilon \rceil - 1) / \lceil 1/\epsilon \rceil$. The test problem instance consists of 30 items, each with integer profit and weight, and of a knapsack with capacity 300. The total profit of all items is 1,345 profit units and the total weight of all items is 1,712 weight units. The run was performed with $\epsilon = 0.4$ and returns in a few seconds with the solution presented in Table 1.

The experimental runs of the GA heuristics were performed by using the galib246 library [2], with a mutation rate $p_m = 0.001$, crossover rate $p_c = 0.6$, and a one-point crossover operator [4]. The initial population consists of a set of maximal (in respect to inclusion) feasible solutions that were obtained by the approximation algorithm during its execution. Therefore, the initial population set contains a subpart of an optimal solution assisting the GA to easily find a near-to-optimal solution. Thus, the use of approximation may result in an additional advantage for the GA. Several runs of the GA were performed, each of 50 generations and takes a few seconds; the best result of these runs is presented in Table 2.

The quality of the solution that is based on heuristics by creating a scale for the problem using the approximate algorithm solution has been determined. The upper bound of the solution can be easily calculated using the approximation algorithm ratio:

$$\text{upper bound} = \frac{\lceil 1/\epsilon \rceil}{\lceil 1/\epsilon \rceil - 1} \times \text{approximation algorithm solution profit.}$$

The performance ratio of the heuristic run is $\text{heuristic solution} / \text{upper bound}$, which equals 0.78 in the considered case. Thus, 0.78 is the heuristic certificate obtained for the heuristic results.

3 Heuristic Choice, Genetic Algorithms (based on [5])

A GA is an iterative search technique which simulates a population of individuals. The search space consists

of candidate solutions to the given problem, each one is encoded by a finite string of symbols called *gnome*. The GA is especially useful for the problems with search space that is too large to be exhaustively explored. Mostly, solutions are represented as binary strings, but other encodings are also possible, e. g., character-based encoding, real-value encoding, tree representation, etc.

The evolution usually starts from a population of individuals which is randomly or heuristically generated. Then, it proceeds by iterative generations – the fitness of every individual in the current population is evaluated according to some predefined objective function, referred to as the *fitness function*. Multiple individuals are selected from the current population based on their fitness. One of the simplest selection procedures is *fitness-proportionate selection*, where individuals are selected with a probability proportional to their relative fitness. This procedure ensures that the expected number of selections of some individual is proportional to its expediency for the population. This policy ensures that high-fitness individuals get higher chance to be promoted and low-fitness individuals are commonly eliminated.

In order to produce new search points, one cannot use a selection only. For this, genetically-inspired operators like *crossover* and *mutation* were used. Selected individuals are modified, recombined, and possibly mutated to form a new population. Crossover is performed as follows: two individuals called *parents* are selected and parts of their gnomes are exchanged between them with probability p_{cross} . This forms two new individuals, called *offsprings*. The simple example of crossover is when substrings are exchanged after a randomly selected crossover point. The mutation is used to maintain genetic diversity from one generation to the next. It prevents impulsive junction to local optima rather than to optimal solution. It is performed by flipping random bits with some usually small probability p_{mut} . The new population is then used in the next iteration of the algorithm. Genetic algorithms are stochastic iterative processes that are not guaranteed to converge. There are a variety of termination conditions, e. g., a solution is found that satisfies minimum criteria, fixed number of generations reached, computational time limitations, etc. In other words, the process stops when some acceptable fitness level is reached.

4 Approximation Algorithms and Approximation Preserving Reductions

An approximation algorithm for an optimization problem generates feasible, but not necessarily optimal, so-

lutions. Since there are no polynomial-time algorithms to get optimal solutions for NP-hard problems (unless $P = NP$), usually, nonoptimal solutions are accepted which can be obtained in polynomial time. Unlike heuristic, the term *approximation algorithm* implies some proven worst case bound on the solution quality in measurable time.

The approximation properties of different problems vary a great deal. Some problems cannot be approximated even with a factor of n , for instance, a problem of finding a maximum clique in a graph. Some other problems can be approximated only with the factor of $O(\log n)$ like problems related to graph separators, or with a predefined constant factor like MAX 3-SAT problem. Only a small quantity of problems can be approximated with any arbitrary constant factor, i. e., have PTAS. One would prefer an approximation algorithm to run in polynomial time and have close to 1 *approximation ratio*, i. e., the worst-case ratio between the solution obtained by the approximation algorithm and the optimal solution. But as the solution given by approximation gets closer to the optimal solution, the time cost gets closer to exponential. Therefore, heuristic algorithm may be preferred to get near-optimal solution with potentially better approximation ratio.

Polynomial time approximation schemes. PTAS for an optimization problem A is a polynomial-time algorithm which input is the problem instance and $\epsilon > 0$ and output is the solution that approximates a given problem within the factor of ϵ . The run time of an algorithm PTAS may depend not just on an input size of a given problem, but also on ϵ .

Approximation preserving reductions (based on [6]).

Definition 1 [6]. Given a class of functions F , an NP-optimization (NPO) problem A that belongs to the class F-APX (an abbreviation of *approximable*) is defined if an $r(n)$ -approximate algorithm T for A exists, for some function $r \in F$. In particular, APX, \log APX, polyAPX, and expAPX will denote the class F-APX with F equal to the set $O(1)$, to the set $O(\log n)$, to the set $O(n^{\log n})$, and to the set $O(2^{n^{O(\log n)}})$, respectively.

Given an instance I of an NPO problem, $|I|$ is used to denote the length of I and $\text{OPT}(I)$ to denote the optimum value for this instance. For any solution S to I , the *objective value* of the solution is denoted by $c(I, S)$.

Definition 2 [6]. Given a solution S to an instance I of an NPO problem, the relative error of S with respect to I is defined as

$$\varepsilon(I, S) := \max \left\{ \frac{c(I, S)}{\text{OPT}(I)}, \frac{\text{OPT}(I)}{c(I, S)} \right\}$$

The above definition is applied to maximization and minimization problems as well.

Definition 3 [6]. An approximation algorithm A for an optimization problem Π has performance ratio $R(n)$ if, given an instance I of Π with $|I| = n$, the solution $A(I)$ satisfies

$$\max \left\{ \frac{c(I, S)}{\text{OPT}(I)}, \frac{\text{OPT}(I)}{c(I, S)} \right\} \leq R(n).$$

A solution of value within a multiplicative factor r of the optimal value is referred to as an r -approximation.

A reduction from a problem A to a problem B specifies some procedure to solve A by means of an algorithm solving B . In content of approximation, the reduction should guarantee that an approximate solution for B can be used to obtain an approximate solution for A . The reduction functions f and g are used where f maps an instance of a problem A to an instance of a problem B and g maps an approximate solution of a problem B into a feasible solution of a problem A .

Definition 4 [6]. Let A and B be two NPO problems, a reduction template (f, g) between them is a tuple of polynomial computable functions such that the following properties hold:

- for any $x \in I_A$, $f(x) \in I_B$;
- for any $x \in I_A$, if $\text{sol}(x) \neq \emptyset$ then also $\text{sol}(f(x)) \neq \emptyset$; and
- for any $y' \in \text{sol}(f(x))$, $g(x, y') \in \text{sol}(x)$.

Typically, polynomial time reductions do not preserve the near-optimality of the solutions. Indeed, all NP-complete problems are equally hard from the viewpoint of obtaining exact solutions. However, from the viewpoint of obtaining near-optimal solutions, they exhibit a rich set of possibilities. An approximation preserving reduction not only has to map instances of a problem A to instances of a problem B , but it also has to be able to obtain good solutions for A from good solutions for B . Approximation preserving reducibilities are defined by imposing some relations between the performance ratios of y' and $g(x, y')$. Several kinds of reducibilities may be found in literature: *Strict reducibility*, *L-reducibility*, *E-reducibility*, *PTAS-reducibility*, *AP-reducibility*. These reducibilities are identical with respect to the overall scheme but differ essentially in the way they preserve approximability: they range from the *Strict reducibility* in which the error cannot increase to the *PTAS-reducibility* where there are basically no restrictions. Below, only two reductions will be described in detail: L-reduction and AP-reduction which are used later in this paper.

L-reducibility [6]. L-reduction stands for *linear reduction* because there is a linear blow-up in the relative approximation error. The L-reduction enforces both optimal and approximation solutions of an instance I of A to be linearly related, respectively, to optimal and

approximation solutions of the instance I' of B to which it is mapped.

Let f and g be polynomial-time reduction functions from optimization problem A to optimization problem B . We say that (f, g) is an *L-reduction* if there are constants $\alpha, \beta > 0$ such that for each instance x of A :

- the optima of x and $f(x)$, $\text{OPT}(x)$ and $\text{OPT}(f(x))$, respectively, satisfy $\text{OPT}(f(x)) \leq \alpha \text{OPT}(x)$; and
- for any solution y' of $f(x)$ with objective value c' , one can find in a polynomial time a solution $y = g(x, y')$ of x with objective value c , so that $|\text{OPT}(x) - c(y)| \leq \beta |\text{OPT}(f(x)) - c'(y')|$.

From the definition above, one gets that $E_A(x, g(x, y)) \leq \alpha \beta E_B(f(x), y)$. This inequality implies that if A is a minimization problem and an r -approximate algorithm for B exists, then a $(1 + \alpha \beta (r - 1))$ -approximate algorithm for A exists. The difficulty of L-reducibility is, mainly, due to the fact that it does not allow the function g to depend on ϵ : as a consequence, this function is forced to map optimum solution into optimum solution. The L-reducibility preserves membership in PTAS but does not preserve membership in APX unless $P = NP \cap \text{co-NP}$. This means that it cannot be blindly used to obtain the existence of approximation algorithms via reductions. The constant β will be usually 1.

AP-reducibility [6]. The AP-reduction lets the functions f and g depend on the expected performance ratio of the reduction. AP-reducibility preserves approximation but not optimal solution. There are two different constraints that are put on the computational time of g and f : the computation time should be polynomial for fixed values of the performance ratio (to preserve membership in PTAS) and the reduction should be efficient even when poor performance ratios are required (to preserve membership in log APX and polyAPX). Thus, the computation time should not increase when the performance ratio decreases.

Definition 5 [6]. Let A and B be two NPO problems, with instances I_A and I_B , respectively. An extended reduction template (f, g) between them is an ordered 2-tuple of functions such that the following properties hold:

- for any $x \in I_A$ and for any rational $r > 1$, $f(x, r) \in I_B$ and is computable in time $t_f(|x|, r)$.
Moreover, if $\text{sol}(x) \neq \emptyset$ then also $\text{sol}(f(x, r)) \neq \emptyset$;
- for any $y' \in \text{sol}(f(x, r))$, $g(x, y', r) \in \text{sol}(x)$ is computable in time $t_g(|x|, |y|, r)$;
- for any fixed r , both $t_f(\cdot, r)$ and $t_g(\cdot, \cdot, r)$ are bounded by a polynomial; and
- for any fixed n and m , both $t_f(n, \cdot)$ and $t_g(n, m, \cdot)$ are nonincreasing functions.

Let A and B be two NPO problems. A is said to be *AP-reducible* to B , denoted by $A \leq_{AP} B$, if an extended template (f, g) between A and B exists and holds 5 properties. The first four properties are of the extended template above. They ensure that f and g reduction functions are polynomial in time, when f maps instances of A into instances of B and g maps solutions for instances of B into solutions for instances of A . The fifth property implies the existence of a positive constant α such that for any $x \in I_A$, for any rational $r > 1$, and for any $y \in \text{sol}_B(f(x, r))$, $R_B(f(x, r), y) \leq r$ implies $R_A(x, g(x, y, r)) \leq 1 + \alpha(r - 1)$. The triple (f, g, α) is said to be an AP-reduction from A to B .

5 Choosing Approximation Preserving Reductions

There are four main criteria that are proposed to use when choosing a specific reduction algorithm from a given problem. Unlike the traditional criteria, constants and not only asymptotic functions are considered.

1. *Quality of approximation preserve.* Some of the reductions preserve approximation by imposing a linear relation between the performance ratios of y' and $g(x, y')$. Others preserve approximation by putting constrains also on the additive errors of y' and $g(x, y')$. The reduction preservation ratio is implied by the parameters of these constrains.
2. *Quality of approximation algorithm of a destination problem.* A performance ratio of the (best) known approximate algorithm for a destination problem.
3. *Inflation of input and the run time complexity of a destination problem.* This parameter is implicitly included in a quality of approximation preserve parameters. Besides, this parameter impacts on a complexity of the required execution of the approximation algorithm of the destination problem.
4. *Complexity of reduction functions f and g .* This parameter represents a degree of ease of a reduction use and also contributes to the general performance of a reduction.

5.1 Approximation preserve reductions of MAX 3-SAT

Let examine two examples of approximation preserving reductions of MAX 3-SAT NPO problem and compare them.

Example 1. MAX 3-SAT \leq_L MAX 2-SAT

Let consider an instance of MAX 3-SAT $\varphi = \varphi_1 \wedge \varphi_2 \wedge \dots \wedge \varphi_k$ with k clauses, each contains at most

3 literals. The reduction f maps φ to an instance φ' of MAX 2-SAT, clause by clause, based on the following roles:

- if φ_i has at most 2 variables, i. e., of the form (x_i^1) or $(x_i^1 \vee x_i^2)$, then $\varphi'_i = \varphi_i$; and
- if φ_i has 3 variables, i. e., of the form $(x_i^1 \vee x_i^2 \vee x_i^3)$, then $\varphi'_i = (x_i^1) \wedge (x_i^2) \wedge (x_i^3) \wedge (y_i) \wedge (\bar{x}_i^1 \vee \bar{x}_i^2) \wedge (\bar{x}_i^1 \vee \bar{x}_i^3) \wedge (\bar{x}_i^2 \vee \bar{x}_i^3) \wedge (x_i^1 \vee \bar{y}_i) \wedge (x_i^2 \vee \bar{y}_i) \wedge (x_i^3 \vee \bar{y}_i)$, where y_i is the new variable.

Suppose a truth assignment τ for φ . If φ_i is of the first form mentioned above and φ_i is satisfied by τ , then, clearly, φ'_i is also satisfied by τ . If a clause φ_i is of the second form and φ_i is satisfied by τ , it can be shown that τ can be extended to a truth assignment τ' for φ'_i which satisfies exactly (and not more) seven clauses of φ'_i . Let consider several cases. In the first case, exactly one literal among x_i^1 , x_i^2 , and x_i^3 is set to *true*. Then, by setting y_i to *false*, let get an assignment that satisfies exactly seven of ten clauses in φ'_i . In the second case, exactly two literals among x_i^1 , x_i^2 , and x_i^3 are set to *true*. Then, by setting y_i indifferently to *true* or *false*, one gets an assignment that satisfies exactly seven of ten clauses in φ'_i . In the third case, all literals x_i^1 , x_i^2 , and x_i^3 are set to *true*. Then, by setting y_i to *false*, one gets an assignment that satisfies exactly seven of ten clauses in φ'_i . Finally, if φ_i is not satisfied by τ , no truth assignment for y_i can satisfy more than six clauses of φ'_i while six are guaranteed by setting y_i to *false*.

Lemma 1 [7]. *Given a propositional formula in conjunctive normal form, at least one half of its clauses can always be satisfied.* (Proof: try some random assignment. If this does not satisfy half the clauses, then its bitwise complement will.)

Denote the number of 3-literal clauses of φ by m . This implies that $\text{opt}(\varphi') = 6m + \text{opt}(\varphi)$. By using the lemma above, one gets that $m \leq 2\text{opt}(\varphi)$ and $\text{opt}(\varphi') = 6m + \text{opt}(\varphi) \leq 13\text{opt}(\varphi)$. This means that $\text{opt}(\varphi') \leq 13\text{opt}(\varphi)$ and $\alpha = 13$. Given a truth assignment τ' for φ' , let consider its restriction $\tau = g(\varphi, \tau')$ on the variables of φ ; for such assignment τ , one has: $m(\varphi, \tau) \geq m(\varphi', \tau') - 6m$. Then $\text{opt}(\varphi) - m(\varphi, \tau) = \text{opt}(\varphi') - 6m - m(\varphi, \tau) \leq \text{opt}(\varphi') - m(\varphi', \tau')$. This means that $\text{opt}(\varphi) - m(\varphi, \tau) \leq \text{opt}(\varphi') - m(\varphi', \tau')$ and $\beta = 1$.

Conclusion 1. According to [7], there is an L-reduction with $\alpha = 13$ and $\beta = 1$ which implies approximation algorithm with ratio of $1/(1 + \alpha\beta(1/0.955 - 1)) = 0.62$ for MAX 3-SAT.

Therefore, the values of the four main criteria mentioned before to analyze the quality of the L-reduction above are:

- (1) *quality of approximation preserve according to* [7]: in L-reduction, this criterion is measured by mul-

tiplication of values of the reduction parameters α and β and equals 13 in this example;

- (2) *quality of approximation algorithm of a destination problem according to [8]*: MAX 2-SAT can be approximated with the ratio 0.955;
- (3) *inflation of input and the run time complexity of a destination problem according to [7]*: recall that a size of φ is k clauses, where m of them have 3 variables. A size of φ' is, thus, $k - m + 10m = k + 9m$ clauses, since every clause with 3 variables is expanded to 10 clauses with 2 and 1 variables. The run time complexity of the used approximation scheme of MAX 2-SAT is much more than linear; and
- (4) *complexity of reduction functions f and g according to [7]*: both functions f and g are linear in a size of input.

Example 2. MAX 3SAT-B \leq_L MAX INDEPENDENT SET-B

Let consider an instance of MAX 3-SAT $\varphi = \varphi_1 \wedge \varphi_2 \wedge \dots \wedge \varphi_k$ with k clauses, where there are at most B occurrences of each variable, for some constant B . The reduction f maps φ to a graph with vertex degree bounded by B , for some constant B (need not be the same as for MAX 3SAT-B instance) in the following way. Construct a graph G with one node for every occurrence of every literal. There is an edge connecting literal occurrences from the same clause — a triangle for every 3-literal clause and a single edge for every 2-literal clause. In addition, there is an edge connecting any two occurrences of complementary literals.

Claim 1 [9, 10]. If every variable of φ occurs $\leq k$ times in the clauses, then the degree of the graph G is $\leq k + 1$.

Every literal has at most 2 edges to the literals that appear in its clause, and at most $k - 1$ edges to its complementary instances in other clauses of φ . Thus, one gets that a degree of the nodes in the graph G is upper bounded by $k + 1$.

Claim 2 [9, 10]. An independent set of size c corresponds to a truth assignment that satisfies at least c clauses of φ .

An independent set cannot select both a literal and its complement and can select at most one literal from each clause; so, a truth assignment can be obtained by setting the variables according to which nodes were in the independent set. For each node in the independent set, there is a satisfied clause. Note that there may be other satisfied clauses of φ .

Conclusion 2 [9, 10]. Let c be an independent set of G , and τ be $g(c)$. Thus,

$$\begin{aligned} \text{opt}(\varphi) - m(\varphi, \tau) \\ = \text{opt}(G) - m(\varphi, \tau) \leq \text{opt}(G) - m(G, c). \end{aligned}$$

This means that $\beta = 1$.

Claim 3 [9, 10]. The size of the maximum independent set in the graph is equal to the maximum number of clauses that can be satisfied.

For each satisfied clause, there is a node in the graph G that can be added to the independent set. This means that $\text{opt}(\varphi) = \text{opt}(G)$ and, thus, $\alpha = 1$.

Conclusion 3. According to [9, 10], there is an L-reduction with $\alpha = \beta = 1$, which implies approximation algorithm with ratio of $\alpha\beta B/6 + o(1) = B/6 + o(1)$ or $\alpha\beta O(B/\log \log B) = O(B/\log \log B)$.

Therefore, the values of the four main criteria mentioned before to analyze the quality of the L-reduction above are:

- (1) *quality of approximation preserve according to [9, 10]*: in L-reduction, this criterion is measured by multiplication of values of the reduction parameters α and β and equals 1 in this example;
- (2) *quality of approximation algorithm of a destination problem according to [11]*: MAX INDEPENDENT SET-B can be approximated with the ratio $B/6 + o(1)$ or $O(B/\log \log B)$ (depends on a value of B);
- (3) *inflation of input and the run time complexity of a destination problem according to [9, 10]*: recall that a size of φ is k clauses where m of them have 3 literals. Denote the number of distinct literals of φ by n . A size of the graph G is, thus, $k - m + 3m + n\lceil k/2 \rceil^2 = k + 2m + n\lceil k/2 \rceil^2$ edges of a graph G , where $3m$ stands for number of edges of a clause with 3 literals, and $n\lceil k/2 \rceil^2$ stands for an upper bound on the number of edges between literals and their complementary literals. Each literal appears at most k times (i. e., in each clause) and $\lceil k/2 \rceil^2$ is the maximal size of full dual graph of k nodes. A number of nodes of graph G is the total number of literals, with duplicates, in all clauses of φ . The run time complexity of the used approximation scheme of MAX INDEPENDENT SET-B is near to linear: $Nk/(j - 1) + \min\{k^2N, N \log N\} + |E|$ where N is number of nodes in the graph G ; E is the number of nodes; and j is the maximal degree of cliques that should be removed from graph G according to the approximation scheme; and
- (4) *complexity of reduction functions f and g according to [9, 10]*: both functions f and g are linear in a size of input.

Comparison of Examples 1 and 2. The approximation ratio of a source problem depends on the parameter k . Note that $6/(k + 1) > 0.509$ leads to $k \leq 10$, and this is the case when one should prefer MAX INDEPENDENT SET-B reduction according to this criterion. The inflation of input is better for MAX 2-SAT, but its run time complexity of destination problem is worse

Table 3 Comparison of Examples 1 and 2

Example	Quality of approximation preserve	Approximation ratio of a destination problem	Inflation of input	Run time of a destination problem	Complexity of f and g	Approximation ratio of a source problem
1	$\alpha\beta = 13$	0.955	$k + 9m$	\gg linear	Linear	0.509
2	$\alpha\beta = 1$	$B/6 + o(1), O(B/\log \log B)$	$k + 2m + n \lceil k/2 \rceil^2$	$N(k/(j-1)) + \min\{k^2N, N \log N\} + E $, (near to linear)	Linear	$B/6 + o(1), O(B/\log \log B)$

than of MAX INDEPENDENT SET-B. Besides, MAX 2-SAT approximation scheme is cumbersome. So, there are trade-offs in choosing of reduction in this case. The present authors suggest to make a choice according to a given instance of MAX 3-SAT problem.

In the following, let examine the reductions above with specific (simplest) approximation algorithms for MAX 2-SAT and MAX INDEPENDENT SET-B. It will be shown how to choose a reduction for getting a scale given an instance of MAX 3-SAT. The comparison parameters are summarized in Table 3.

It seems natural to consider first the reduction to MAX 2-SAT, since MAX 2-SAT and MAX 3-SAT are common problems, and the MAX 2-SAT approximation algorithm is easier to implement [4]. Unfortunately, the reduction to MAX 2-SAT is found useless, and the present authors turned to examine the reduction to MAX INDEPENDENT SET-B which resulted with a nontrivial scale.

Let use the following input $\varphi = (x_1) \wedge (\bar{x}_1) \wedge \dots \wedge (x_{40}) \wedge (x_{40}) \wedge (x_{41} \vee x_{42}) \wedge (x_{43}) \wedge (x_{41} \vee x_{43}) \wedge (x_{41} \vee x_{42} \vee x_{43}) \wedge \dots \wedge (x_{53} \vee x_{54}) \wedge (x_{55}) \wedge (x_{53} \vee x_{55}) \wedge (x_{53} \vee x_{54} \vee x_{55})$ of MAX 3-SAT. There are 100 clauses in φ , 85 of them are one-literal, 10 of them are 2-literal, and 5 of them are 3-literal. The approximation algorithm for MAX 2-SAT is described in [12] and is called *the probabilistic method*. The approximation algorithm for MAX INDEPENDENT SET-B is a Greedy algorithm which is described in [11]. First, it will be shown that for the input φ , ahead can be determine (i. e., without running the reduction code) that the MAX 2-SAT reduction above does not produce a good enough scale (in fact, it produces a trivial scale that is related to the length of the input).

According to the used approximation algorithm for MAX 2-SAT, the translated input φ' is approximated with the ratio $r = 0.37$ of an input size. According to L-reduction properties, this means that the reduction is $1/(1+\alpha\beta(1/r-1)) = 1/(1+13 \cdot 1 \cdot (1/0.37-1)) = 0.04$ -approximative. Assume that the reduction obtains a truth assignment τ that satisfies all 5 3-literal clauses of φ . This means that exactly 7 of 10 clauses of φ' that represent every 3-literal clause of φ , are satisfied. Since there are 5

3-literal clauses in φ , there are $7 \cdot 5 = 35$ satisfied clauses of total 145 clauses of φ' . Since at least 0.37 of φ' clauses should be satisfied, at least 19 2-literal or unit clauses of φ' are also satisfied. These clauses are not transformed when φ is translated to φ' ; so, they are satisfied also in the reduction solution for φ . Thus, one gets the upper bound of $\min\{|\varphi|, 19/0.04\} = \min\{100, 475\} = 100$, which adds no information beyond the trivial upper bound of φ (i. e., equals to the input size) and implies the trivial scale. Otherwise, the reduction obtains a truth assignment τ that does not satisfy all 5 3-literal clauses of φ . In this case, even more 2-literal or unit clauses of φ are satisfied by τ and the reduction once again returns a trivial upper bound for φ .

The MAX INDEPENDENT SET-B reduction [4] results in a few seconds in 0.55 of satisfied clauses of φ and its performance ratio is 0.6. This implies a non-trivial upper bound of $55/0.6 = 92$ clauses.

The run of GA for MAX 3-SAT is performed by using, once again, galib246 library [2] with a mutation rate $p_m = 0.001$, crossover rate $p_c = 0.6$, and a one-point crossover operator [4]. The initial population is random. In a few seconds, a truth assignment was got that satisfies 58 of 100 clauses of φ . The solution scale from MAX INDEPENDENT SET-B reduction certifies the GA solution as $58/92 = 0.64$ of the best possible result.

5.2 Approximation preserve reductions of MB3DM

The MB3DM NPO problem is defined as follows. Given are disjointed sets A , B , and C and a subset of triples $T \subseteq A \times B \times C$. Let consider the restricted version of the MB3DM problem where $|A| = |B| = |C| = q$ and each element of these sets appears in one, two, or three triples of T . The goal is to find a subset $T' \subseteq T$ of maximum cardinality such that no two triples of T' agree in any coordinate. This version of the MB3DM problem has been shown to be APX-hard in [10]. Note that each triple can intersect at most six other triples, which implies that the maximum matching consists of at least $|T|/7$ triples. Let examine two examples of approximation preserving reductions of MB3DM and compare them.

Example 3. MB3DM \leq_L MINIMIZATION OF SUM OF SQUARED MACHINE LOADS

The problem of scheduling unrelated parallel machines is defined as follows. Given are a set of n independent jobs, J_1, \dots, J_n and a set of m parallel machines M_1, \dots, M_m . Each job j can be allocated to one of the machines in a subset $M(j) \subseteq 1, \dots, m$. Each machine can process one job at a time and all machines are available at start time. Denote the time a job j takes to be proceeded on a machine M_i by p_{ij} . In addition, each job j has nonnegative weight w_j . The goal is to arrange the jobs to machines so that the minimal sum of the weighted completion times is achieved.

Denote the sum of weights of jobs assigned to the machine i by the load l_i . Let define Sum of Squared Machine Loads as the l_2 norm of the machines load vector $\vec{l} = (l_1, \dots, l_m)$, where l_2 is defined by $(\sum_{i=1}^m l_i^2)^{1/2}$. Indeed, l_2 is a measure of the quality of a given assignment.

Let consider an instance I of the restricted version of MB3DM. The reduction f maps I into an instance of the Sum of Squared Machine Loads problem in the following way. Let define $3q$ machines, a machine $M(T_i)$ for each triple T_i in T , and $3q - |T|$ dummy machines. Also, let define $5q$ jobs, a job per each element of A , B , and C and $2q$ dummy jobs. On all machines, each element job has a processing time 1 and each dummy job has a processing time 3. The element job can be assigned only to some triple machine $M(T_i)$ and only if the triple assigned with the machine contains the appropriate element.

Note that for the restricted version of MB3DM, the optimal solution $\text{Opt}(I)$ consists of $q = |A| = |B| = |C|$ triples, since each element of A , B , and C appear in at least one triple of T and a feasible solution consists of disjoint triples, the size of optimal solution is q . Let estimate a size of the optimal solution of $f(I)$. The best case is when all three jobs of each triple T_i in $\text{Opt}(I)$ are scheduled to the appropriate machine $M(T_i)$. The $2q$ dummy jobs are scheduled to the $2q$ dummy machines, one job per machine. One gets that every machine has load 3 (three element jobs with load 1, or one dummy job with load 3). Then the objective value of this schedule is the sum of squares of machine load and equals to $3^2 \cdot 3q = 27q$. Thus, $\text{Opt}(f(I)) \leq 27q = 27\text{Opt}(I)$ which means that the parameter α of the L-reduction equals to 27.

Denote by m_k , $k = 0, \dots, 3$, a number of machines in some feasible solution of $f(I)$ that process exactly k element jobs. Then the total number of machines equals $m_0 + m_1 + m_2 + m_3 = 3q$, and $m_1 + 2m_2 + 3m_3 = 3q$ is the total number of element jobs. Note that according to the reduction definition, the objective value $c(g(f(I))) = m_3$.

Lemma 2 [13]. *The objective value $c(s)$ of the feasible solution s of the scheduling instance $f(I)$ satisfies $c(s) \geq 29q - 2m_3$.*

This lemma yields that $|c(g(s)) - \text{Opt}(I)| = q - m_3 = 1/2(29q - 2m_3 - 27q) \leq 1/2|c(s) - \text{Opt}(f(I))|$ which means that the β parameter of the L-reduction equals to $1/2$.

Conclusion 4. According to [13], there is an L-reduction with $\alpha = 27$ and $\beta = 1/2$, which implies approximation algorithm with ratio of $1/(1 + \alpha\beta(1 + \sqrt{2} - 1)) = 1/(1 + 27 \cdot 1/2 \cdot \sqrt{2}) = 0.05$ for MB3DM.

Therefore, the values of the four main criteria mentioned before to analyze the quality of the L-reduction above are:

- (1) *quality of approximation preserve according to [13]:* in L-reduction this criterion is measured by multiplication of values of the reduction parameters α and β and equals 13.5 in this example;
- (2) *quality of approximation algorithm of a destination problem according to [14]:* the SUM OF SQUARED MACHINE LOADS problem can be approximated with the ratio $1 + \sqrt{2}$;
- (3) *inflation of input and the run time complexity of a destination problem according to [13]:* $3q$ machines and $5q$ jobs are defined which implies the inflation of input to be $8q$. The run time complexity of the used approximation scheme of SUM OF SQUARED MACHINE LOADS is quadratic in input size; and
- (4) *complexity of reduction functions f and g according to [13]:* both functions f and g are linear in a size of input.

Example 4. MB3DM \leq_{AP} STAR-GCA-SIMPLE

Let A and B be two NPO problems. A is said to be *AP-reducible* to B , denoted by $A \leq_{AP} B$, if an extended template (f, g) between A and B and a positive constant α exist such that for any $x \in I_A$, for any rational $r > 1$, and for any $y \in \text{sol}_B(f(x, r))$, $R_B(f(x, r), y) \leq r$ implies $R_A(x, g(x, y, r)) \leq 1 + \alpha(r - 1)$. The triple (f, g, α) is said to be an AP-reduction from A to B [6].

The problem of STAR-GCA (General Call Admission control problem in STAR networks) is defined as follows. A star network is the undirected graph $G = (V, E)$, when the node 0 represents a unique central node and rest of nodes are connected to it and called outer nodes. The edge set E consists of the edges $e_i = (i, 0)$ for $i = 1, \dots, n$ according to the network links structure. Each link has a positive capacity $c(e)$.

A request for a connection (call) is defined by a tuple $(u_i, v_i, t_i, d_i, b_i, p_i)$ consisting of a source node $u_i \in V$, a destination node $v_i \in V$, a starting time t_i , a duration d_i , a positive bandwidth requirement b_i , and

a profit p_i . A solution is a set of accepted calls of those arrived. A feasibility of a solution is determined by the sum of bandwidth requirements of simultaneously active calls using the same edge does not exceed the capacity of that edge. The goal is to maximize the sum of profits of the accepted calls.

STAR-GCA-SIMPLE is the restriction of the STAR-GCA problem which is defined as follows. Each call i has one of $t_i \in 0, 1, 2$ starting times, duration $d_i = 2$ and unit profit and the needed bandwidth is also one unit on each call edge. The capacity of each edge is a unit which implies that only a single path per edge can be active in every given moment. Thus, the objective function is a cardinality of accepted calls set.

Let I be an instance of MB3DM problem. The reduction f maps I into instance of STAR-GCA-SIMPLE in the following way. For every element $a_i \in A, b_i \in B$, and $c_i \in C$, a link of this element to the central node 0 is defined. For each triple $t_j = (a_j, b_j, c_j) \in T$, additional three nodes $d_{j,1}, d_{j,2}$, and $d_{j,3}$ are defined and linked to the central node. Also, 5 additional requests are added: $r_1 = (d_{j,1}, d_{j,2})$ with time interval $[0, 2]$; $r_2 = (d_{j,2}, d_{j,3})$ with time interval $[2, 4]$; $r_3 = (d_{j,1}, a_j)$ with time interval $[1, 3]$; $r_4 = (d_{j,2}, b_j)$ with time interval $[1, 3]$; and $r_5 = (d_{j,3}, c_j)$ with time interval $[1, 3]$.

One has to ensure that no more than three requests for one triple are accepted. Note that according to the reduction definition above, either the requests r_1 and r_2 are accepted since they have disjoint time intervals or the requests r_3, r_4 , and r_5 are accepted since they do not share any edge. In any case, the accepted requests of one triple block the rest of the requests of the same triple. Let obtain an approximate solution M_1 for a given instance of MB3DM problem using the following

Lemma 3 [15]. *There is a greedy procedure that computes a 1/3-approximation for the bounded maximum 3-dimensional matching problem.*

Then the reduction function g composes a solution M_2 using the solution obtained for the instance $f(I)$ of the STAR-GCA-SIMPLE problem in the following way. If three calls r_3, r_4 , and r_5 of some triple $t \in T$ are accepted, this triple is added to the solution of the MB3DM problem instance I . Then, let choose the maximal solution obtained from the reduction and from the approximation algorithm for MB3DM, i.e., $|g(f(I))| = \max\{|M_1|, |M_2|\}$. From the lemma above, one gets that $|M_1| \geq |M^*|/3$ where M^* is the maximum matching for I . By the following

Lemma 4 [15]. *Let $T \subseteq A \times B \times C$ be an instance of the maximum three-dimensional matching problem, and let (G, R) be the corresponding instance of STAR-GCA-SIMPLE defined above. There is a feasible solution for (G, R) that accepts $2|T| + k$ requests iff T has a matching of size k .*

one obtains that $|M_2| \geq |g(f(I))| - 2|T|$. Till now, the first four properties of AP-reduction have been satisfied. Let show that the fifth property is also holds with $\alpha = 43$. According to the lemma above, an optimal solution for $f(I)$ consists of $|M^*| + 2|T|$ requests where M^* is the optimal solution of an instance I . Assume that one has an r -approximation algorithm for $f(I)$ that implies a solution Q consists of at least $(|M^*| + 2|T|)/r$ requests. If $r \geq 45/43$, the inequality $|g(f(I))| \geq |M_1| \geq |M^*|/3$ shows that g computes a 1/3-approximation. Since $3 = 1 + 43(45/43 - 1) \leq 1 + 43(r - 1)$, $|g(f(I))| \geq M^*/(1 + 43(r - 1))$ and $1 + 43(r - 1) \geq M^*/|g(f(I))|$. According to the fifth property of the AP-reduction, this means that the value of the reduction parameter α equals 43. Otherwise, $r < 45/43$. From $|Q| \geq (|M^*| + 2|T|)/r$, one gets

$$\begin{aligned} |Q| &\geq \frac{2r|T| + |M^*| - 2(r-1)|T|}{r} \\ &= 2|T| + \frac{|M^*| - 2(r-1)|T|}{r} \\ &\geq 2|T| + \frac{|M^*|(1 - 14(r-1))}{r} \end{aligned}$$

where the property $|T| \leq 7|M^*|$ mentioned before in the last inequality was used. As $|g(f(I))| \geq |M_2| \geq |Q| - 2|T| \geq (1 - 14(r - 1))|M^*|/r$, one gets that

$$\begin{aligned} \frac{|M^*|}{|g(f(I))|} &\leq \frac{|M^*|}{(1 - 14(r - 1))|M^*|/r} \\ &= 1 + \frac{15}{15 - 14r}(r - 1) \leq 1 + 43(r - 1) \end{aligned}$$

where the last inequality holds for $1 < r < 45/43$. Again, the fifth property is fulfilled with $\alpha = 43$.

Conclusion 5. According to [15], there is an AP-reduction with $\alpha = 43$, which implies approximation algorithm with ratio of $1/(1 + \alpha(1/18 - 1)) = 1/(1 + 43(18 - 1)) = 1/732 = 0.001$ for MB3DM.

Therefore, the values of the four main criteria mentioned before to analyze the quality of the AP-reduction above are:

- (1) *quality of approximation preserve according to [15]:* in AP-reduction this criterion is measured by the value of the reduction parameter α and equals 43 in this example;
- (2) *quality of approximation algorithm of a destination problem according to [15]:* the STAR-GCA-SIMPLE problem can be approximated with the ratio 1/18;
- (3) *inflation of input and the run time complexity of a destination problem according to [15]:* suppose the restricted version of MB3DM problem as of the Example 3 reduction, where $|A| = |B| = |C| = q$

Table 4 Comparison of Examples 3 and 4

Example	Quality of approximation preserve	Approximation ratio of a destination problem	Inflation of input	Run time of a destination problem	Complexity of f and g	Approximation ratio of a source problem
3	$\alpha\beta = 13.5$	$1 + \sqrt{2}$	$8q$	Quadratic in input size	Linear	0.05
4	$\alpha = 43$	$1/18$	$11q \leq \dots \leq 27q$	Quadratic in input size	Linear	0.001

and $q \leq |T| \leq 3q$. Let add a vertex for each element of A , B , and C , and for each triple of T , let add 3 vertices and 5 requests to R . Then the inflation of input in this case is $3q + 8|T|$. Since $q \leq |T| \leq 3q$, one gets that $11q \leq 3q + 8|T| \leq 27q$. The run time complexity of the used approximation scheme of STAR-GCA-SIMPLE is quadratic in input size; and

- (4) *complexity of reduction functions f and g according to [15]*: both functions f and g are linear in a size of input.

Comparison of Examples 3 and 4. The approximation ratio of the destination problem in Example 3 is tighter than the approximation ratio of the destination problem in Example 4, namely $1 + \sqrt{2}$ vs. $1/18$. Since these examples use different kinds of reductions, the quality of approximation preserve cannot be compared directly but only the approximation ratios obtained for MB3DM problem. According to this criterion, the reduction of Example 3 is better: 0.05 vs. 0.001. The input inflation for Example 3 is only slightly smaller, and the execution time of both destination problems approximation algorithms is quadratic in input size. Therefore, this criterion does not dominate the choice of a reduction. The complexity of the reduction functions f and g of both reductions is linear in a size of input.

According to all comparable criteria, summarized in Table 4, the reduction of Example 3 is better than the reduction of Example 4.

6 Concluding Remarks

The main goal of this paper is to define a new framework in which the quality of a solution that is based on heuristics is certified by an approximation algorithm. The framework was defined and it was verified that it is feasible, useful, and essential in order to get a certificate for a heuristic.

It was note that it is possible that the approximation ratio stated for a particular approximation algorithm is the worst case ratio over all inputs, one may investigate tight approximation ratio for a given input. This may further motivate the designer of approximation schemes to provide a function from inputs to approximation ratio rather than only the worst case single bound. Also, it

was found that the approximation results are useful to define the initial generation of the GA as an important technique.

The authors hope that the computation of a scale and a certificate will become a standard practice accompanying any heuristics which, in turn, will assist in core AI tasks such as symbolic planning, scheduling, and theorem proving.

Acknowledgments

The authors thank Moshe Sipper and Eitan Bachmat for helpful discussions.

References

- Holland, J. 1971. Genetic algorithms and the optimal allocation of trials. *SIAM J. Comput.* 2:88–105.
- Wall, M., and MIT. 1994–2005. GALib: A C++ library of genetic algorithm components. Available at: <http://lancet.mit.edu/ga/> (accessed February 10, 2015).
- Kellerer, H., U. Pferschy, and D. Pisinger. 2004. *Knapsack problems*. Berlin: Springer. 161–166.
- <http://www.cs.bgu.ac.il/~sadetsky/Thesis/> (accessed February 10, 2015).
- Sipper, M. 1996. A brief introduction to genetic algorithms. Available at: <http://www.cs.bgu.ac.il/~sipper/ga.html> (accessed February 10, 2015).
- Trevisan, L. 1997. Reductions and (non-)approximability. *Universita Degli Studi di Roma 'La Sapienza,' dottorato di Ricerca in Informatica*. IX-97-7:17–35.
- Ausiello, G., and V. Th. Paschos. 2005. Approximability preserving reductions. *Cahier Du Lamsade* 227:12.
- Hastad, J. 1997. Some optimal inapproximability results. *29th ACM Symposium on Theory of Computing Proceedings*. 1–10.
- <http://www.cs.cmu.edu/afs/cs.cmu.edu/academic/class/15451-s00/www/lectures/lect0406post.txt> (accessed February 10, 2015).
- Papadimitriou, C., and M. Yannakakis. 1988. Optimization, approximation, and complexity classes. *20th Annual ACM Symposium on the Theory of Computing Proceedings*. 229–234.
- Halldorsson, M., and J. Radhakrishnan. 1994. Greed is good: Approximating independent sets in sparse and

- bounded-degree graphs. *30th ACM Symposium on Theory of Computing Proceedings*. 439–448.
12. Yannakakis, M. 1994. *On the approximation of maximum satisfiability*. *3rd Annual ACM-SIAM Symposium on Discrete Algorithms Proceedings*. Orlando, FL, USA. 475–502.
 13. Azar, Y., L. Epstein, Y. Richter, and G. Woeginger. 2004. All-norm approximation algorithms. *J. Algorithm*. 52(2):120–133.
 14. Awerbuch, B., Y. Azar, E. Grove, M. Kao, P. Krishman, and J. Vitter. 1995. Load balancing in the L_p norm. *IEEE Symposium on Foundations of Computer Science (FOCS) Proceedings*.
 15. Adamy, U., T. Erlebach, D. Mitsche, I. Schurr, B. Speckmann, and E. Welzl. 2005. Off-line admission control for advance reservations in star networks. *Approximation Online Algorithms* 3351:211–224.

Received January 12, 2015

Contributors

Dolev Shlomi (b. 1958) — Doctor of Science in computer science, professor, Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel; dolev@cs.bgu.ac.il

Kogan-Sadetsky Marina (b. 1977) — PhD student, Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel; sadetsky@cs.bgu.ac.il

ЭВРИСТИЧЕСКИЕ СЕРТИФИКАТЫ ПОСРЕДСТВОМ ПРИБЛИЖЕНИЙ

Ш. Долев¹, М. Коган-Садецкая²

¹Факультет компьютерных наук, Университет Бен-Гурион в Негеве, Израиль, dolev@cs.bgu.ac.il

²Факультет компьютерных наук, Университет Бен-Гурион в Негеве, Израиль, sadetsky@cs.bgu.ac.il

Аннотация: Предложен новый метод, в котором качество (необязательно оптимального) эвристического решения сертифицировано приближенным алгоритмом, а именно: результат эвристического решения сопровождается шкалой, получаемой из приближенного алгоритма. Создание шкалы эффективно, в то время как получение решения от алгоритма аппроксимации обычно требует длительных расчетов относительно эвристического подхода. С другой стороны, результаты, полученные с помощью эвристики без шкалы, могут быть бесполезными. Исследованы критерии для выбора схемы аппроксимации для получения шкалы. Чтобы получить шкалу на практике, приближения рассмотрены не только по их асимптотическому поведению, но также изучены соотношения между ними относительно размера ввода для данной проблемы. Для практического примера рассмотрены эвристические и приближенные алгоритмы для задач SINGLE KNAPSACK, MAX 3-SAT и MAXIMUM BOUNDED THREE-DIMENSIONAL MATCHING, которые являются известными NP-сложными задачами. Получены сертификаты для эвристических запусков с использованием подходящих приближений.

Ключевые слова: эвристика; приближенный алгоритм; оптимальное решение; сводимости сохраняющие приближения

DOI: 10.14357/19922264150103

Литература

1. Holland J. Genetic algorithms and the optimal allocation of trials // *SIAM J. Comput.*, 1971. Vol. 2. P. 88–105.
2. Wall, M., and MIT. 1994–2005. GALib: A C++ library of genetic algorithm components. <http://lancet.mit.edu/ga/>.
3. Kellerer H., Pferschy U., Pisinger D. Knapsack problems. — Berlin: Springer, 2004. P. 161–166.
4. <http://www.cs.bgu.ac.il/~sadetsky/Thesis/>.
5. Sipper M. 1996. A brief introduction to genetic algorithms. <http://www.cs.bgu.ac.il/~sipper/ga.html>.
6. Trevisan L. 1997. Reductions and (non-)approximability. Universita Degli Studi di Roma ‘La Sapienza,’ dottorato di Ricerca in Informatica. Vol. IX-97-7. P. 17–35.
7. Ausiello G., Paschos V. Th. Approximability preserving reductions // *Cahier Du Lamsade*, 2005. Vol. 227. P. 12.
8. Hastad J. Some optimal inapproximability results // *29th ACM Symposium on Theory of Computing Proceedings*, 1997. P. 1–10.
9. <http://www.cs.cmu.edu/afs/cs.cmu.edu/academic/class/15451-s00/www/lectures/lect0406post.txt>.
10. Papadimitriou C., Yannakakis M. Optimization, approximation, and complexity classes // *20th Annual ACM*

- Symposium on the Theory of Computing Proceedings, 1988. P. 229–234.
11. *Halldorsson M., Radhakrishnan J.* Greed is good: Approximating independent sets in sparse and bounded-degree graphs // 30th ACM Symposium on Theory of Computing Proceedings, 1994. P. 439–448.
 12. *Yannakakis M.* On the approximation of maximum satisfiability // 3rd Annual ACM-SIAM Symposium on Discrete Algorithms Proceedings. — Orlando, FL, USA, 1994. P. 475–502.
 13. *Azar, Y., L. Epstein, Y. Richter, and G. Woeginger.* All-norm approximation algorithms // J. Algorithm., 2004. Vol. 52. No. 2. P. 120–133.
 14. *Awerbuch B., Azar Y., Grove E., Kao M., Krishman P., Vitter J.* 1995. Load balancing in the L_p norm. *IEEE Symposium on Foundations of Computer Science (FOCS) Proceedings.*
 15. *Adamy U., Erlebach T., Mitsche D., Schurr I., Speckmann B., Welzl E.* Off-line admission control for advance reservations in star networks // Approximation Online Algorithms, 2005. Vol. 3351. P. 211–224.

Поступила в редакцию 12.01.2015

METHODS AND TOOLS FOR HYPOTHESIS-DRIVEN RESEARCH SUPPORT: A SURVEY*

L. Kalinichenko¹, D. Kovalev¹, D. Kovaleva², and O. Malkov²

Abstract: Data intensive research (DIR) is being developed in frame of the new paradigm of research study known as the Fourth paradigm, emphasizing an increasing role of observational, experimental, and computer simulated data practically in all research domains. The principal goal of DIR is an extraction (inference) of knowledge from data. The intention of this work is to make an overview of the existing approaches, methods, and infrastructures of the data analysis in DIR accentuating the role of hypotheses in such process and efficient support of hypothesis formation, evaluation, and selection in course of the natural phenomena modeling and experiments carrying out. An introduction into various concepts, methods, and tools intended for effective organization of hypothesis-driven experiments in DIR is presented.

Keywords: data intensive research; Fourth paradigm; hypotheses; models; theories; hypothetico-deductive method; hypothesis testing; hypothesis lattice; Galaxy model; connectome analysis; automated hypothesis generation

DOI: 10.14357/19922264150104

1 Hypotheses, Theories, Models and Laws in Data Intensive Science

Data intensive research is being developed in accordance with the Fourth Paradigm [1] of research study (following three previous historical paradigms of the science development (empirical science, theoretical science, and computational science)) emphasizing that science as a whole is becoming increasingly dependent on data as the core source for discovery. Emerging of the Fourth Paradigm is motivated by the huge amount of data coming from scientific instruments, sensors, simulations, as well as from people accumulating data in Web or social nets. The basic objective of DIR is to infer knowledge from the integrated data organized in networked infrastructures (such as warehouses, grids, clouds). At the same time, “Big Data” movement has emerged as a recognition of the increased significance of massive data in various domains. Open access to large volumes of data, therefore, becomes a key prerequisite for discoveries in the XXI century. Data intensive research denotes a crosscut of DIR/IT areas aimed at the creation of effective data analysis technologies for DIR covering scientific and other data intensive domains (including finance, economy, social environment, business, etc.).

Science endeavors to give a meaningful description of the world of natural phenomena using that are known as laws, hypotheses, and theories. Hypotheses, theories,

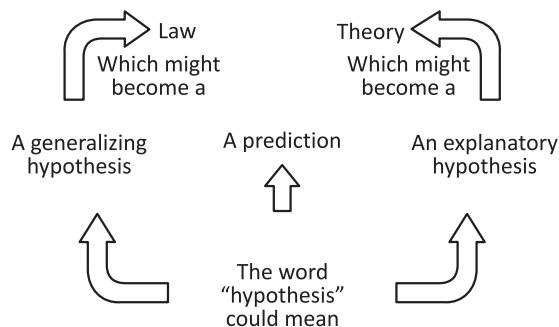


Figure 1 Multiple incarnations of hypotheses

and laws in their essence have the same fundamental character (Fig. 1) [2].

A *scientific hypothesis* is a proposed explanation of a phenomenon which still has to be rigorously tested. In contrast, a *scientific theory* has undergone extensive testing and is generally accepted to be the accurate explanation behind an observation. A *scientific law* is a proposition, which points out any such orderliness or regularity in nature, *the prevalence of an invariable association between a particular set of conditions and particular phenomena*. In the exact sciences, laws can often be expressed in the form of mathematical relationships. Hypotheses explain laws, and well-tested, corroborated hypotheses become theories (see Fig. 1). At the same time, the laws do not cease to be laws, just because they did not appear first as hypotheses and pass through the stage of theories.

*This work has been partially supported by the RFBR grants 13-07-00579 and 14-07-00548.

¹Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation

²Institute of Astronomy, Russian Academy of Sciences, 48 Pyatnitskaya Str., Moscow 119017, Russian Federation

Though theories and laws are different kinds of knowledge, actually, they represent different forms of the same knowledge construct. Laws are generalizations, principles, or patterns in nature, and theories are the explanations of those generalizations. However, classification expressed in Fig. 1 is subjective. Article [3] provides examples showing that the differences between laws, hypotheses, and theories consist only in that they stand at different levels in their claim for acceptance depending on how much empirical evidence is amassed. Therefore, there is no essential difference between constructs used for expressing hypotheses, theories, and laws. Important role of hypotheses in scientific research can scarcely be overestimated. In the edition of M. Poincaré's book [4], it is stressed that *without hypotheses, there is no science*. Thus, it is not surprising that so much attention in the scientific research and the respective publications is devoted to the methods for hypothesis manipulation in experimenting and modeling of various phenomena applying the means of informatics. The idea that the new approaches are needed that can address both *data-* and *hypothesis-driven sciences* runs all through this paper. Such symbiosis alongside with the hypothesis-driven tradition of science (“first hypothesize-then-experiment”) might cause wide application of another one that is typified by “first experiment-then-hypothesize” mode of research. Often, the “first experiment” ordering in DIR is motivated by the necessity of analysis of the existing massive data to generate a hypothesis.

In the course of the present study, paying attention to the issue of inductive and deductive reasoning in hypothesis-driven sciences will be emphasized. In Fig. 2, such ways of knowledge production are shown [2]. Here, “generalization” means any subset of hypotheses, theories, and laws and “Evidence” is any subset of all facts accumulated in a specific DIR.

All researchers collect and interpret empirical evidence through the process called *induction*. This is a technique by which individual pieces of evidence are collected and examined until a law is discovered or a theory is invented. Frances Bacon first formalized induction [5]. The method of (naive) induction (see Fig. 2), he suggested, is, in part, the principal way by which humans traditionally have produced generalizations that permit predictions. The problem with induction is that

it is impossible to collect all observations pertaining to a given situation in all time — past, present, and future.

The formulation of a new law begins through induction as facts are heaped upon other relevant facts. Deduction is useful in checking the validity of a law. Figure 2 shows that a valid law would permit the accurate prediction of facts not yet known. Also an *abduction* [6] is the process of validating a given hypothesis through reasoning by successive approximation. Under this principle, an explanation is valid if it is the best possible explanation of a set of known data. Abductive validation is common practice in hypothesis formation in science. Hypothesis related logic reasoning issues are considered in more details in section 3.

In [4], the useful hypotheses of science are considered to be of two kinds:

- (1) the hypotheses which are valuable *precisely* because they are either verifiable or, else, refutable through a definite appeal to the tests furnished by experience; and
- (2) the hypotheses which, despite the fact that experience suggests them, are valuable *despite*, or even *because*, of the fact that experience can neither confirm nor refute them.

Aspects of science which are determined by the use of the hypotheses of the second kind are considered in [4] as “constituting an essential human way of viewing nature, an interpretation rather than a portrayal or a prediction of the objective facts of nature, an adjustment of our conceptions of things to the internal needs of our intelligence.” According to Poincaré's discussion, the central problem of the logic of science becomes the problem of the relation between the two fundamentally distinct kinds of hypotheses, i. e., between those which cannot be verified or refuted through experience and those which can be empirically tested.

The analysis in this paper will be focused mostly on the modeling of hypotheses of the first kind, leaving issues of analysis of the relations between such two kinds of hypotheses to further study.

The rest of the paper is organized as follows. Section 2 discusses the basic concepts defining the role of hypotheses in the formation of scientific knowledge and the respective organization of the scientific experiments. Approaches for hypothesis formulation, logical reasoning, hypothesis modeling, and testing are briefly introduced in section 3. In section 4, a general overview of the basic facilities provided by informatics for the hypothesis-driven experimentation scenarios, including conceptual modeling, simulations, statistics and machine learning methods is given. In section 5, several examples of organization of hypothesis-driven scientific experiments are included. Concluding remarks summarize the discussion.

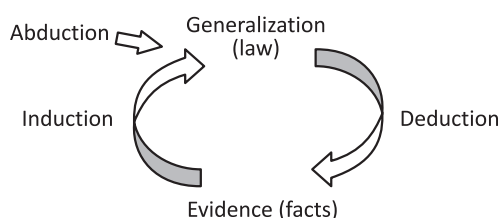


Figure 2 Enhanced knowledge production diagram

2 Role of Hypotheses in Scientific Experiments: Basic Principles

Normally, scientific hypotheses have the form of a mathematical model. Sometimes, one can also formulate them as existential statements, stating that some particular instance of the phenomenon under examination has some characteristic and causal explanations, which have the general form of universal statements, stating that every instance of the phenomenon has a particular characteristic (e. g., *for all x , if x is a swan, then x is white*). Scientific hypothesis considered as a declarative statement identifies the predicted relationship (associative or causal) between two or more variables (independent and dependent). In causal relationship, a change caused by the independent variable is predicted in the dependent variable. Variables are more commonly related in noncausal (associative) way [7].

In experimental studies, the researcher manipulates the independent variable. The dependent variable is often referred to as consequence or the presumed effect that varies with a change of the independent variable. The dependent variable is not manipulated. It is observed and assumed to vary with changes in the independent variable. Predictions are made from the independent variable to the dependent variable. It is the dependent variable that the researcher is interested in understanding, explaining, or predicting [7].

In case when a possible correlation or similar relation between variables is investigated (such as, for example, whether a proposed medication is effective in treating a disease, that is, at least to some extent and for some patients), a few cases in which the tested remedy shows no effect do not falsify the hypothesis. Instead, statistical tests are used to determine how likely it is that the overall effect would be observed if no real relation as hypothesized exists. If that likelihood is sufficiently small, the existence of a relation may be assumed. In statistical hypothesis testing, two hypotheses are compared, which are called the *null hypothesis* and the *alternative hypothesis*. The null hypothesis states that there is no relationship between the phenomena (variables) whose relation is under investigation or, at least, not of the form given by the alternative hypothesis. The alternative hypothesis, as the name suggests, is the alternative to the null hypothesis: it states that there *is* some kind of relation.

Alternative hypotheses are generally used more often than null hypotheses because they are more desirable to state the researcher's expectations. But in any study that involves statistical analysis, the underlying null hypothesis is usually assumed [7]. It is important that the conclusion "do not reject the null hypothesis" does not necessarily mean that the null hypothesis is true. It sug-

gests that there is not sufficient evidence against the null hypothesis in favor of the alternative hypothesis. Rejecting the null hypothesis suggests that the alternative hypothesis may be true.

Any useful hypothesis will enable *predictions by reasoning* (including *deductive reasoning*). It might predict the outcome of an experiment in a laboratory setting or the observation of a phenomenon in nature. The prediction may also invoke statistics assuming that a hypothesis must be *falsifiable* [8] and that one cannot regard a proposition or theory as scientific if it does not admit the possibility of being shown false. The way to demarcate between hypotheses is to call *scientific* those for which we can specify (beforehand) one or more potential falsifiers as the respective experiments. Falsification was supposed to proceed deductively instead of inductively.

Other philosophers of science have rejected the criterion of falsifiability or supplemented it with other criteria, such as verifiability (only statements about the world that are empirically confirmable or logically necessary are cognitively meaningful). They claim that science proceeds by "induction" — that is, by finding confirming instances of a conjecture. Popper treated confirmation as never certain [8]. However, a falsification can be sudden and definitive. Einstein said: "No amount of experimentation can ever prove me right; a single experiment can prove me wrong." To scientists and philosophers outside the Popperian belief [8], science operates mainly by induction (confirmation), and also and less often by disconfirmation (falsification). Its language is almost always one of induction. For this survey both philosophical treatment of hypotheses are acceptable. Sometimes such way of reasoning is called the *hypothetico-deductive method*. According to it, scientific inquiry proceeds by formulating a hypothesis in a form that could conceivably be falsified by a test on observable data. A test that could and does run contrary to predictions of the hypothesis is taken as a falsification of the hypothesis. A test that could but does not run contrary to the hypothesis corroborates the theory.

A scientific method involves experiment to test the ability of some hypothesis to adequately answer the question under investigation. A prediction enabled by hypothesis suggests a test (observation or experiment) for the hypothesis thus becoming testable. If a hypothesis does not generate any observational tests, there is nothing that a scientist can do with it.

For example, not testable hypothesis: "Our universe is surrounded by another, larger universe, with which we can have absolutely no contact;" not verifiable (though testable) hypothesis: "There are other inhabited planets in the universe;" scientific hypothesis (both testable and verifiable): "Any two objects dropped from the same

height above the surface of the earth will hit the ground at the same time as long as air resistance is not a factor” (<http://www.batesville.k12.in.us/physics/phynet/aboutscience/hypotheses.html>).

A *problem (research question)* should be formulated as an issue of what relation exists between two or more variables. The problem statement should be such as to imply possibilities of empirical testing; otherwise, this will not be a scientific problem. Problems and hypotheses being generalized relational statements enable to deduce specific empirical manifestations implied by the problem and hypotheses. In this process, hypotheses can be deduced from theory and from other hypotheses. A problem cannot be scientifically solved unless it is reduced to hypothesis form, because a problem is not directly testable [9].

Most formal hypotheses connect concepts by specifying the expected relationships between *propositions*. When a set of hypotheses are grouped together, they become a type of *conceptual framework*. When a conceptual framework is complex and incorporates causality or explanation, it is generally referred to as a *theory* [10]. In general, hypotheses have to reflect the multivariate complexity of the reality. A scientific theory summarizes a hypothesis or a group of hypotheses that have been supported with repeated testing. A theory is valid as long as there is no evidence to dispute it. *Scientific paradigm* explains the working set of theories under which science operates.

Elements of hypothesis-driven research and their relationships are shown in Fig. 3 [11, 12]. The hypothesis triangle relations, *explains*, *formulates*, and *represents*, are functional in the scientist’s final decision in adopting a particular model m_1 to formulate a hypothesis h_1 , which is meant to explain phenomenon p_1 .

In [12], the lattice structure for hypothesis interconnection is proposed as shown in Fig. 4. A hypothesis lattice is formed by considering a set of hypotheses equipped with *wasDerivedFrom* as a strict order (from the bottom to the top). Hypotheses directly derived from exactly one hypothesis are *atomic*, while those directly derived from at least two hypotheses are *complex*.

The hypothesis lattice is unfolded into model and phenomena isomorphic lattices according to the hypothesis triangle (see Fig. 3) [12]. The lattices are isomorphic if one takes subsets of M (Model), H (Hypotheses), and P (Phenomenon) such that *formulates*, *explains*, and *represents* are both one-to-one and onto mappings (i. e., bijections), seen as structure-preserving mappings (morphisms). Example of the isomorphic lattice is shown in Fig. 5 [12]. This particular lattice corresponds to the case in Computational Hemodynamics considered in [12]. Here, model m_1 formulates hypothesis h_1 , which explains phenomenon p_1 . Similarly, m_2 formulates h_2 , which explains p_2 , and so on. Prop-

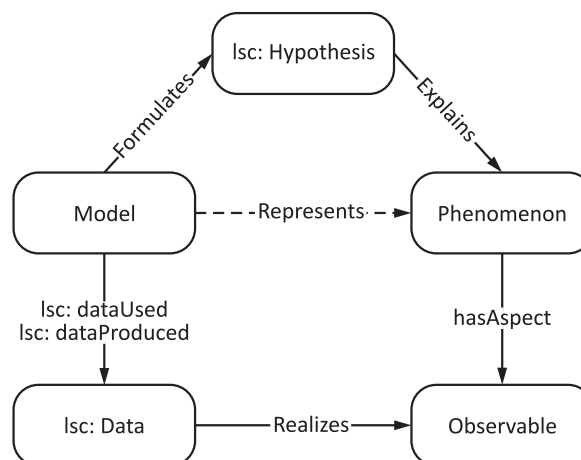


Figure 3 Elements of hypothesis-driven research

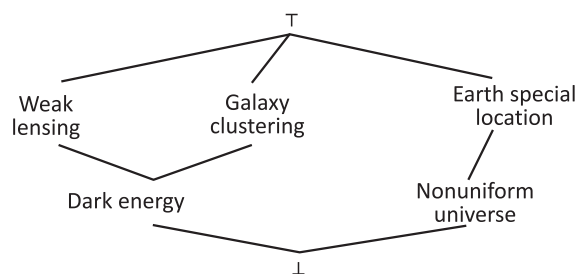


Figure 4 A lattice theoretic representation for hypothesis relationship

erties of the hypothesis lattices and operations over them are considered in [13].

Models are one of the principal instruments of modern science. Models can perform two fundamentally different representational functions: a model can be a representation of a selected part of the world, or a model can represent a theory in the sense that it interprets the laws and hypotheses of that theory.

Here, let consider scientific models to be representations in both senses at the same time. One of the most perplexing questions in connection with models is how they relate to theories. In this respect, models can be considered as a complement to theories, as preliminary theories, can be used as substitutions of theories when the latter are too complicated to handle. Learning about the model is done through experiments, thought experiments, and simulation. Given a set of parameters, a model can generate expectations about how the system will behave in a particular situation. A model and the hypotheses it is based upon are supported when the model generates expectations that match the behavior of its real-world counterpart.

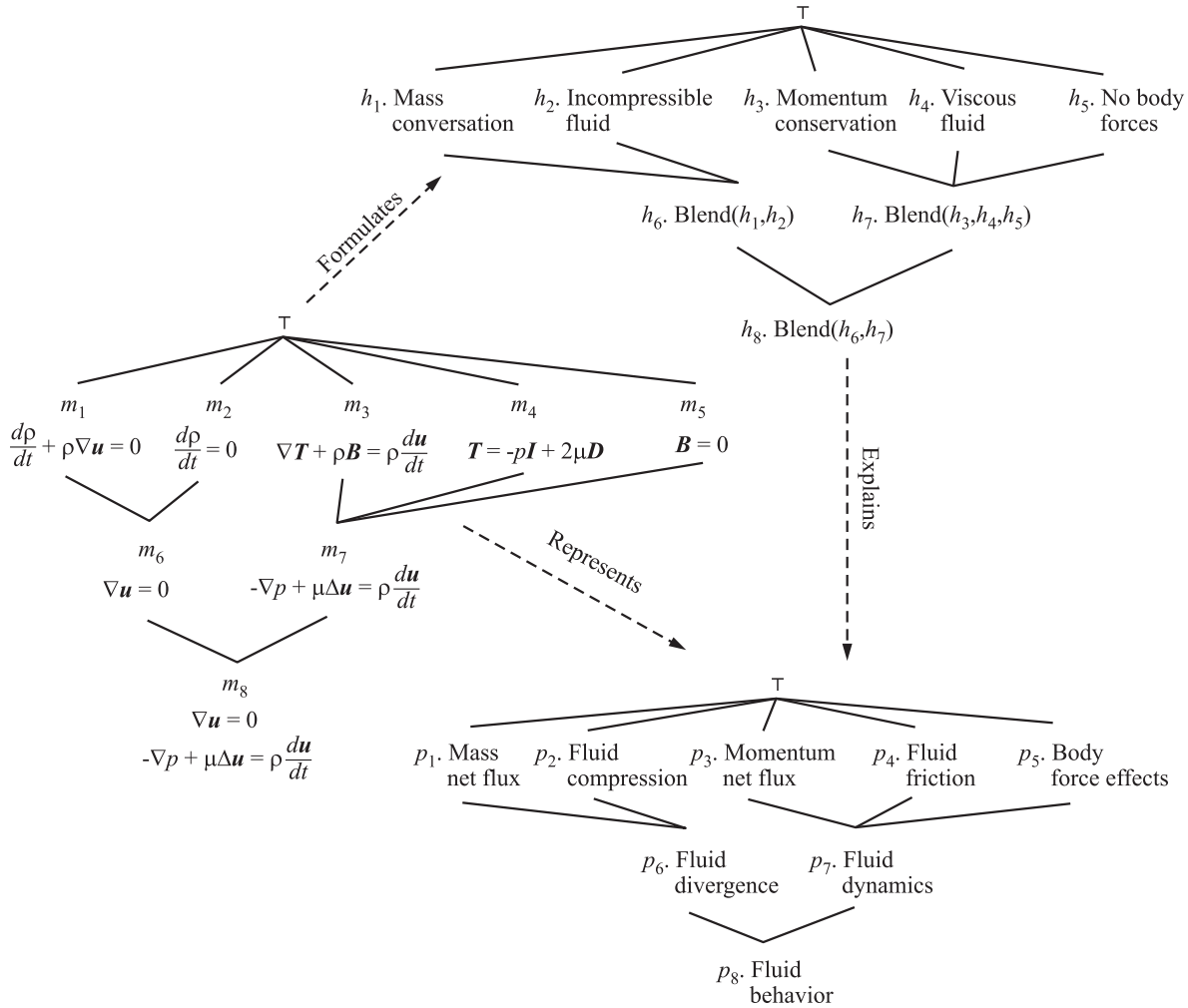


Figure 5 Hypothesis lattice unfolded into model and phenomenon isomorphic lattice

A law generalizes a body of observations. Generally, a law represents a group of related undisputable hypotheses using a handful of fundamental concepts and equations to define the rules governing a set of phenomena. A law does not attempt to explain why something happens — it simply states that it does.

Facilities for support of the hypothesis-driven experimentation will be discussed in the remaining sections.

3 Hypothesis Manipulation in Scientific Experiments

3.1 Hypothesis generation

Researchers that support rationality of scientific discovery presented several methods for hypothesis generation, including discovery as abduction, induction, anomaly detection, heuristics programming, and use of analogies [14].

Discovery as abduction characterizes reasoning processes that take place before a new hypothesis is justified. The abductive model of reasoning that leads to plausible hypotheses formulation is conceptualized as an inference beginning with data. According to [15], an abduction happens as follows:

- (1) some phenomena p_1, p_2, p_3, \dots are encountered for which there is no or little explanation;
- (2) however, p_1, p_2, p_3, \dots would not be surprising if a hypothesis H were added. They would certainly follow from something like H and would be explained by it; and
- (3) therefore, there is a good reason for elaborating a hypothesis H — for proposing it as a possible hypothesis from which the assumption p_1, p_2, p_3, \dots might follow.

The abductive model of reasoning is primarily a process of explaining anomalies or surprising phenomena [16]. The scientists' reasoning proceeds abductively from an

anomaly to an explanatory hypothesis in light of which the phenomena would no longer be surprising. There can be several different hypotheses that can serve as the explanations for phenomena; so, additionally some criteria for choosing among different hypotheses are required.

One way to implement abductive model of reasoning is the abductive logic programming [17]. Hypothesis generation in abduction logical framework is organized as follows. During the experiment, some new observations are encountered. Let B represents the background knowledge and O is the set of facts that represents observations. Both B and O are the logic programs (set of rules in some rule language). In addition, Γ stands for a set of literals representing the set of abducibles, which are candidate assumptions to be added to B for explaining O . Given B , O , and Γ , the hypothesis-generation problem is to find a set H of literals (called a hypothesis) such that:

- (1) B and H entail O ;
- (2) B and H are consistent; and
- (3) H is some subset of Γ .

If all conditions are met, then H is an explanation of O (with respect to B and Γ). Examples of abductive logic programming systems include ACLP [18], A-system [19], ABDUAL [20], and ProLogICA [21]. Abductive logic programming can also be implemented by means of Answer Set Programming systems, e. g., by the DLV system [22].

The example abductive logic program in ProLogICA describes a simple model of the lactose metabolism of the bacterium E.Coli [21]. The background knowledge B describes that E.coli can feed on the sugar lactose if it makes two enzymes permease and galactosidase. Like all enzymes (E), these are made if they are coded by a gene (G) that is expressed. These enzymes are coded by two genes (lac(y) and lac(z)) in cluster of genes (lac(X)) called an operon that is expressed when the amounts (amt) of glucose are low and lactose are high or when they are both at medium level. The abducibles, Γ , declare all ground instances of the predicates “amount” as assumable. This reflects the fact that in the model, it is not known what are the amounts at any time of the various substances. This is incomplete information that should be found out in each problem case that is examined. The integrity constraints state that the amount of a substance (S) can only take one value.

```
## Background Knowledge (B)
feed(lactose):- make(permease),
               make(galactosidase).
make(Enzyme):- code(Gene,Enzyme), express(Gene).
express(lac(X)):- amount(glucose,low),
                 amount(lactose,hi).
express(lac(X)):- amount(glucose,medium),
                 amount(lactose,medium).
```

```
code(lac(y),permease).
code(lac(z),galactosidase).
temperature(low):-amount(glucose,low).
false :- amount(S,V1), amount(S,V2), V1 != V2.
```

```
## Abducibles ( $\Gamma$ )
abducible_predicate(amount).
```

```
## Observation (O)
feed(lactose).
```

```
This goal generates two possible hypotheses:
{amount(lactose,hi), amount(glucose,low)}
{amount(lactose,medium), amount(glucose,medium)}
```

Below, just a couple of another examples of real rule-based systems, where abductive logic programming is used, are presented. Robot Scientist (see subsection 4.4) abductively hypothesizes new facts about the yeast functional biology by inferring what is missing from a model [23]. In [24], both abduction and induction are used to formulate hypotheses about inhibition in metabolic pathways. Augmenting background knowledge is done with abduction; after that, induction is used for learning general rules. Authors of [25] use SOLAR reasoning system to abductively generate hypotheses about the inhibitory effects of toxins on the rat metabolisms.

The process of discovery is deeply connected also with the search of *anomalies*. There are a lot of methods and algorithms to discover anomalies. Anomaly detection is an important research problem in data mining aimed at search of the objects that are considerably dissimilar, exceptional, and inconsistent with respect to the majority data in an input database [26].

Analogies play several roles in science. Not only do they contribute to discovery but they also play a role in the development and evaluation of scientific theories (new hypotheses) by analogical reasoning.

3.2 Hypothesis evaluation

Being testable and falsifiable, a scientific hypothesis provides a solid basis to its further modeling and testing. There are several ways to do it, including the use of statistics, machine learning, and logic reasoning techniques.

3.2.1 Statistical testing of hypotheses

The classical (frequentist) and Bayesian statistic approaches are applicable for hypothesis testing and selection. Brief summary of the basic differences between these approaches are as follows [27].

Classical (frequentist) statistics is based on the following beliefs:

- probabilities refer to relative frequencies of events. They are objective properties of the real world;
- parameters of hypotheses (models) are fixed, unknown constants. Because they are not fluctuating, probability statements about parameters are meaningless; and
- statistical procedures should have well-defined long-run frequency properties.

In contrast, Bayesian approach takes the following assumptions:

- probability describes the degree of subjective belief, not the limiting frequency. Probability statements can be made about things other than data, including hypotheses (models) themselves as well as their parameters; and
- inferences about a parameter are made by producing its probability distribution — this distribution quantifies the uncertainty of our knowledge about that parameter. Various point estimates, such as expectation value, may then be readily extracted from this distribution.

The Bayesian interpretation of probability can be seen as an extension of propositional logic that enables reasoning with hypotheses, i.e., the propositions whose truth or falsity is uncertain.

Bayesian probability belongs to the category of evidential probabilities; to evaluate the probability of a hypothesis, the Bayesian probabilist specifies some prior probability, which is then updated in the light of new, relevant data (evidence) [28]. The Bayesian interpretation provides a standard set of procedures and formulae to perform this calculation.

Hypothesis testing in classical statistic style. After null and alternative hypotheses are stated, some statistical assumptions about data samples should be done, e.g., assumptions about statistical independence or distributions of observations. Failure in providing correct assumptions leads to the invalid test results.

A common problem in classical statistics is to ask whether a given sample is consistent with some hypothesis. For example, one might be interested in whether a measured value x_i , or the whole set $\{x_i\}$, is consistent with being drawn from a Gaussian distribution $N(\mu, \sigma)$. Here, $N(\mu, \sigma)$ is the *null hypothesis*.

It is always assumed that we know how to compute the probability of a given outcome from the null hypothesis: for example, given the cumulative distribution function, $0 \leq H_0(x) \leq 1$, the probability that we would get a value at least as large as x_i is $p(x > x_i) = 1 - H_0(x_i)$ and is called the *p-value*. Typically, a threshold p value is adopted, called *the significance level* α , and the null hypothesis is rejected when $p \leq \alpha$ (e.g., if $\alpha = 0.05$ and

$p < 0.05$, the null hypothesis is rejected at a 0.05 significance level). If one fails to reject a hypothesis, it does not mean that its correctness is proved because it may be that the sample is simply not large enough to detect an effect.

When performing these tests, one can meet with two types of errors, which statisticians call *Type I and Type II errors*. Type I errors are the cases when the null hypothesis is true but incorrectly rejected. In the context of source detection, these errors represent spurious sources or, more generally, false positives (with respect to the alternative hypothesis). The false-positive probability when testing a single datum is limited by the adopted significance level α . Cases when the null hypothesis is false but it is not rejected are called Type II errors (missed sources, or false negatives (again, with respect to the alternative hypothesis)). The false-negative probability when testing a single datum is usually called β and is related to *the power of α test as $(1 - \beta)$* . Hypothesis testing is intimately related to comparisons of distributions.

As the significance level α is decreased (the criterion for rejecting the null hypothesis becomes more conservative), the number of false positives decreases and the number of false negatives increases. Therefore, there is a trade-off to be made to find an optimal value of α , which depends on the relative importance of false negatives and positives in a particular problem. Both the acceptance of false hypotheses and the rejection of true ones are errors that scientists should try to avoid. There is discussion as to what states of affairs is less desirable; many people think that the acceptance of a false hypothesis is always worse than failure to accept a true one and that science should in the first place try to avoid the former kind of error.

When many instances of hypothesis testing are performed, a process called *multiple hypothesis testing*, the fraction of false positives can significantly exceed the value of α . The fraction of false positives depends not only on α and the number of data points, but also on the number of true positives (the latter is proportional to the number of instances when an alternative hypothesis is true).

Depending on data type (discrete vs. continuous random variables) and what one can assume (or not) about the underlying distributions, and the specific question one asks, different statistical tests can be used. The underlying idea of statistical tests is to use data to compute an appropriate statistic and then compare the resulting data-based value to its expected distribution. The expected distribution is evaluated by *assuming that the null hypothesis is true*. When this expected distribution implies that the data-based value is unlikely to have arisen from it by chance (i.e., the corresponding p value is small), the null hypothesis is rejected with some thresh-

old probability α , typically 0.05 or 0.01 ($p < \alpha$). Note again that $p > \alpha$ does *not* mean that the hypothesis is *proven* to be correct.

The number of various statistical tests in the literature is overwhelming and their applicability is often hard to decide (see [29, 30] for variety of statistical methods in SPSS (Statistical Package for the Social Sciences)). When the distributions are not known, tests are called non-parametric, or distribution-free tests. The most popular nonparametric test is the Kolmogorov–Smirnov (K-S) test, which compares the cumulative distribution function, $F(x)$, for two samples, $\{x_{1i}\}$, $i = 1, \dots, N_1$, and $\{x_{2i}\}$, $i = 1, \dots, N_2$. The K-S test is not the only option for nonparametric comparison of distributions. The Cramér – von Mises criterion, the Watson test, and the Anderson–Darling test are similar in spirit to the K-S test, but consider somewhat different statistics. The Mann–Whitney–Wilcoxon test (or the Wilcoxon rank-sum test) is a nonparametric test for testing whether two data sets are drawn from distributions with different location parameters (if these distributions are known to be Gaussian, the standard classical test is called the t test). A few standard statistical tests can be used when it is known, or can be assumed, that both $h(x)$ and $f(x)$ are the Gaussian distributions (e. g., the Anderson–Darling test, the Shapiro–Wilk test) [27]. More on statistical tests can be found in [27, 29, 30, 31].

Hypothesis (model) selection and testing in Bayesian style. The Bayesian approach can be thought of as formalizing the process of continually refining our state of knowledge about the world, beginning with no data (as encoded by the *prior*), then updating that by multiplying in the likelihood once the data are observed to obtain the *posterior*. When more data are taken, then the posterior based on the first data set can be used as the prior for the second analysis. Indeed, the data sets can be different.

The question often arises as to which is the ‘best’ model (hypothesis) to use; ‘*model selection*’ is a technique that can be used when we wish to discriminate between competing models (hypotheses) and identify the best model (hypothesis) in a set, $\{M_1, \dots, M_n\}$, given the data.

Let remind the basic notation. The Bayes theorem can be applied to calculate the posterior probability $p(M_j|d)$ for each model (or hypothesis) M_j representing our state of knowledge about the truth of the model (hypothesis) in the light of the data d as follows:

$$p(M_j|d) = p(d|M_j) \frac{p(M_j)}{p(d)}$$

where $p(M_j)$ is the prior belief in the model (hypothesis) that represents our state of knowledge (or ignorance) about the truth of the model (hypothesis) before the current data have been analyzed; $p(d|M_j)$ is the model

(hypothesis) *likelihood* (represents the probability that some data are produced under the assumption of this model); and $p(d)$ is the normalization constant given by

$$p(d) = \sum_i p(d|M_i)p(M_i).$$

The relative ‘goodness’ of models is given by a comparison of their posterior probabilities; so, to compare two models M_a and M_b , let look at the ratio of the model posterior probabilities:

$$\frac{p(M_a|d)}{p(M_b|d)} = \frac{p(d|M_a)p(M_a)}{p(d|M_b)p(M_b)}.$$

The Bayes factor, B_{ab} , can be computed as the ratio of the model likelihoods:

$$B_{ab} = \frac{p(d|M_a)}{p(d|M_b)}.$$

Empirical scale for evaluating the strength of evidence from the Bayes factor B_{ij} between two models is shown in the table [32].

Strength of evidence for Bayes factor B_{ij} for two models

$ \ln B_{ij} $	Odds	Strength of evidence
< 1.0	$< 3 : 1$	Inconclusive
1.0	$\sim 3 : 1$	Weak evidence
2.5	$\sim 12 : 1$	Moderate evidence
5.0	$\sim 150 : 1$	Strong evidence

The Bayes factor gives a measure of the ‘goodness’ of a model regardless of the prior belief about the model; the higher the Bayes factor, the better the model is. In many cases, the prior belief in each model in the set of proposed models will be equal; so, the Bayes factor will be equivalent to the ratio of the posterior probabilities of the models. The ‘best’ model in the Bayesian sense is the one which gives the best fit to the data with the smallest parameter space.

A special case of model (hypothesis) selection is *Bayesian hypothesis testing* [27, 33]. Taking M_1 to be the “null” hypothesis, one can ask whether the data supports the alternative hypothesis M_2 , i. e., whether one can reject the null hypothesis. Taking equal priors $p(M_1) = p(M_2)$, the odds ratio is

$$B_{21} = \frac{p(d|M_1)}{p(d|M_2)}.$$

The inability to reject M_1 in the absence of an alternative hypothesis is very different from the hypothesis testing procedure in classical statistics. The latter procedure rejects the null hypothesis if it does not provide a good description of the data, that is, when it is very unlikely that the given data could have been generated as prescribed by the null hypothesis. In contrast, the Bayesian approach is based on the posterior rather than

on the data likelihood and cannot reject a hypothesis if there are no alternative explanations for observed data [27].

Comparing classical and Bayesian approaches [27], it is rare for a mission-critical analysis be done in the “fully Bayesian” manner, i. e., without the use of the frequentist tools at the various stages. Philosophy and beauty aside, the reliability and efficiency of the underlying computations required by the Bayesian framework are the main practical issues. A central technical issue at the heart of this is that it is much easier to do optimization (reliably and efficiently) in high dimensions than it is to do integration in high dimensions. Thus, the usable machine learning methods, while there are ongoing efforts to adapt them to Bayesian framework, are almost all rooted in frequentist methods.

Most users of Bayesian estimation methods, in practice, are likely to use a mix of Bayesian and frequentist tools. The reverse is also true — frequentist data analysts, even if they stay formally within the frequentist framework, are often influenced by “Bayesian thinking,” referring to “priors” and “posteriors.” The most advisable position is probably to know both paradigms well, in order to make informed judgments about which tools to apply in which situations [27]. More details on Bayesian style of hypothesis testing can be found in [27, 28, 33].

3.2.2 Logic-based hypothesis testing

According to the hypothetico-deductive approach, the hypotheses are tested by deducing predictions or other empirical consequences from general theories. If these predictions are verified by experiments, this supports the hypothesis. It should be noted that not everything that is logically entailed by a hypothesis can be confirmed by a proper test for it. The relation between hypothesis and evidence is often empirical rather than logical. A clean deduction of empirical consequences from a hypothesis, as it may sometimes exist in physics, is practically inapplicable in biology. Thus, entailment of the evidence by hypotheses under test is neither sufficient nor necessary for a good test. Inference to the best explanation is usually construed as a form of inductive inference (see abduction in subsection 3.1) where hypothesis’ explanatory credentials are taken to indicate its truth [34].

An inductive logic is a system of evidential support that extends deductive logic to less-than-certain inferences. For valid deductive arguments, the premises logically entail the conclusion where the entailment means that the truth of the premises provides a guarantee of the truth of the conclusion. Similarly, in a good inductive argument, the premises should provide some degree of support for the conclusion, where such support means that the truth of the premises indicates with some degree of strength that the conclusion is true. If the logic of good inductive arguments is to be of any real

value, the measure of support it articulates should meet the Criterion of Adequacy (CoA): as evidence accumulates, the degree to which the collection of true evidence statements comes to support a hypothesis, as measured by the logic, should tend to indicate that the hypotheses are probably false or probably true. In [35], the extent to which a kind of logic based on the Bayes theorem can estimate how the implications of hypotheses about evidence claims influences the degree to which hypotheses are supported is discussed in detail. In particular, it is shown how such a logic may be applied to satisfy the CoA: as evidence accumulates, false hypotheses will very probably come to have evidential support values (as measured by their posterior probabilities) that approach 0; and as this happens, a true hypothesis will very probably acquire evidential support values (measured by their posterior probabilities) that approach 1.

3.2.3 Parameter estimation

Models (hypotheses) are typically described by parameters θ whose values are to be estimated from data. The authors describe this process according to [27]. For a particular model M and prior information I , one gets:

$$p(M, \theta|d, I) = \frac{p(d|M, \theta, I)p(M, \theta|I)}{p(d|I)}.$$

The result $p(M, \theta|d, I)$ is called the *posterior* probability density function (pdf) for model M and parameters θ , given data d and other prior information I . This term is a $(k+1)$ -dimensional pdf in the space spanned by k model parameters and the model M . The term $p(d|M, \theta, I)$ is the *likelihood* of data *given* some model M and some fixed values of parameters θ describing it and all other prior information I . The term $p(M, \theta|I)$ is the *a priori* joint probability for model M and its parameters θ in the absence of any of the data used to compute likelihood and is often simply called the *prior*.

In the Bayesian formalism, $p(M, \theta|d, I)$ corresponds to the state of our *knowledge* (i. e., belief) about a model and its parameters, given data d . To simplify the notation, $M(\theta)$ will be substituted by M whenever the absence of explicit dependence on θ is not confusing. A completely Bayesian data analysis has the following conceptual steps.

1. Formulation of the data likelihood $p(d|M, I)$.
2. Choice of the prior $p(\theta|M, I)$, which incorporates all other knowledge that might exist, but is *not* used when computing the likelihood (e. g., prior measurements of the same type, different measurements, or simply an uninformative prior). Several methods for constructing “objective” priors have been proposed. One of them is the *principle of maximum entropy* for assigning uninformative priors by maximizing the entropy over a suitable set of pdfs,

finding the distribution that is least informative (given the constraints). Entropy maximization with no testable information takes place under a single constraint: the sum of the probabilities must be one. Under this constraint, the maximum entropy for a discrete probability distribution is given by the uniform distribution.

3. Determination of the posterior $p(M|d, I)$, using Bayes theorem above. In practice, this step can be computationally intensive for complex multidimensional problems.
4. The search for the best model M parameters, which maximizes $p(M|d, I)$, yielding the *maximum a posteriori* (MAP) estimate. This point estimate is the natural analog to the *maximum likelihood estimate* (MLE) from classical statistics.
5. Quantification of uncertainty in parameter estimates, via *credible regions*. As in MLE, such an estimate can be obtained analytically by doing mathematical derivations specific to the chosen model. The same as in MLE, various numerical techniques can be used to simulate samples from the posterior. This can be viewed as an analogy to the frequentist approach, which can simulate draws of samples from the true underlying distribution of the data. In both cases, various descriptive statistics can then be computed on such samples to examine the uncertainties surrounding the data and estimators of model parameters based on that data.
6. Hypothesis testing as needed to make other conclusions about the model (hypothesis) or parameter estimates.

3.3 Algorithmic generation and evaluation of hypotheses

Two cultures of data analysis (*formulaic modeling*¹ and *algorithmic modeling*) distinguished here in accordance with [36] can be applied to the hypothesis extraction and generation based on data.

Formulaic modeling is a process for estimating the relationships among variables. It includes many techniques for modeling and analyzing several variables, when the focus is on the formulae $y = f(x)$ that give a relation specifying a vector of dependent variables y in terms of a vector of independent variables x . In a statistics experiment (based on various regression techniques), the dependent variable defines the event studied and is expected to change whenever the independent variable (*predictor* variables, extraneous variables) is altered. Such methods as linear regression, logistic regression, and multiple regression are the well-known examples of the representatives of this modeling approach.

¹In [36], instead of “formulaic modeling,” the term “data modeling” is used that looks misleading in the computer science context.

In the *algorithmic modeling* culture, the approach is to find an algorithm that operates on x to predict the responses y . What is observed is a set of x 's that go in and a subsequent set of y 's that come out. Predictive accuracy and properties of the algorithms (such as, for example, their convergence if they are iterative) are the issues to be investigated. *Machine learning algorithms* focus on prediction, based on known properties learned from the training data. Such machine learning algorithms as decision tree, association rule, neural networks, support vector machines as well as other techniques of learning in Bayesian and probabilistic models [37, 38] are examples of the methods that belong to this second culture.

The models that best emulate the nature in terms of predictive accuracy are also the most complex and inscrutable. Nature forms the outputs y from the inputs x by means of a black box with complex and unknown interior. Current accurate prediction methods are also *complex black boxes* (such as neural nets, forests, support vectors). So, we are facing two black boxes, where ours seem only slightly less inscrutable than nature's [36]. In a choice between *accuracy* and *interpretability*, in applications, people sometimes prefer interpretability.

However, the goal of a model is not interpretability (a way of getting information), but getting useful, accurate information about the relation between the response and predictor variables. It is stated in [36] that algorithmic models can give better predictive accuracy than formulaic models, providing also better information about the underlying mechanism. And actually, this is what the goal of statistical analysis is. The researchers should be focused on solving the problems instead of asking what regression model they can create.

An objection to this idea (expressed by Cox) is that prediction without some understanding of underlying process and linking with other sources of information becomes more and more tentative. Due to that, it is suggested to construct the stochastic calculation models that summarize the understanding of the phenomena under study. One of the objectives of such approach might be an understanding and test of hypotheses about underlying process. Given the relatively small sample size, following such direction could be productive. But data characteristics are rapidly changing. In many of the most interesting current problems, the idea of starting with a formal model is not tenable. The methods used in statistics for small sample sizes and a small number of variables are not applicable. Data analytics need to be more pragmatic. Given a statistical problem, find a good solution, whether it is a formulaic model, an algorithmic model, or a Bayesian model or a completely different approach.

In the context of the hypothesis-driven analysis, one should pay attention to the question how far can we

go applying the algorithmic modeling for hypothesis generation and testing. Various approaches to machine learning use related to hypothesis formation and selection can be found in [27, 36, 38].

Besides machine learning, an interesting example of algorithmic generation of hypotheses can be found in the IBM Watson project [39] where the symbiosis of the general-purpose reusable natural language processing (NLP) and knowledge representation and reasoning (KRR) technologies (under the name DeepQA) is exploited for answering arbitrary questions over the existing natural language documents as well as structured data resources. Hypothesis generation takes the results of question analysis and produces candidate answers by searching the available data sources and extracting answer-sized snippets from the search results. Each candidate answer plugged back into the question is considered a hypothesis, which the system has to prove correct with some degree of confidence. After merging, the system must rank the hypotheses and estimate confidence based on their merged scores. A machine-learning approach adopted is based on running the system over a set of training questions with known answers and training a model based on the scores. An important consideration in dealing with NLP-based scorers is that the features they produce may be quite sparse, and so, accurate confidence estimation requires the application of confidence-weighted learning techniques [39] — a new class of online learning methods that maintain a probabilistic measure of confidence in each parameter. It is important to note that instead of statistics based hypothesis testing, contextual evaluation of a wide range of loosely coupled probabilistic question and semantic based content analytics is applied for scoring different questions (hypotheses) and content interpretations. Training different models on different portions of the data in parallel and combining the learned classifiers into a single classifier allow to make the process applicable to the large collections of data. More details on that can be found in [39, 40] as well as in other Watson project related publications.

3.4 Bayesian motivation for discovery

One way for discriminating between competing models of some phenomenon is to use Bayesian model selection approach (see paragraph 3.2.1), the Bayesian evidences for each of the proposed models (hypotheses) can be computed and the models can then be ranked by their Bayesian evidence. This is a good method for identifying which is the best model in a given set of models, but it gives no indication of the *absolute goodness* of the model. Bayesian model selection says nothing about the *overall quality* of the set of models (hypotheses) as a whole —

the best model in the set may merely be the best of in a set of poor models. Knowing that the best model in the current set of models is not particularly good model would provide *motivation to search for a better model* and, hence, may lead to model discovery.

One way of assigning some measure of the absolute goodness of a model is to use the concept of Bayesian doubt first introduced in [41]. Bayesian doubt works by comparing all the known models in a set with an idealized model, which acts as a benchmark model.

An application of the Bayesian doubt method for the cosmological model building is given in [32, 42]. One of the most important questions in cosmology is to identify the fundamental model underpinning the vast amount of observations nowadays available. The so-called ‘cosmological concordance model’ is based on the cosmological principle (i. e., the Universe is isotropic and homogeneous, at least on large enough scales) and on the hot big bang scenario, complemented by an inflationary epoch. This remarkably simple model is able to explain with only half a dozen free parameter observations spanning a huge range of time and length-scales. Since both a cold dark matter (CDM) and a cosmological constant (Λ) component are required to fit the data, the concordance model is often referred to as ‘the Λ CDM model.’

Several different types of explanation are possible for the apparent late time acceleration of the Universe, including different classes of dark energy model such as Λ CDM, w CDM; theories of modified gravity; void models or the back reaction [32]. The methodology of Bayesian doubt which gives an absolute measure of the degree of goodness of a model has been applied to the issue of whether the Λ CDM model should be doubted.

The methodology of Bayesian doubt dictates that an unknown idealized model X should be introduced against which the other models may be compared. Following [41], ‘doubt’ may be defined as the posterior probability of the unknown model:

$$D \equiv p(X|d) = \frac{p(d|X)p(X)}{p(d)}.$$

Here, $p(X)$ is the prior doubt, i. e., the prior on the unknown model, which represents the degree of belief that the list of known models does not contain the true model. The sum of all the model priors must be unity.

The methodology of Bayesian doubt requires a baseline model (the best model in the set of known models), for which, in this application, the Λ CDM has been chosen. The average Bayes factor between Λ CDM and each of the known models is given by:

$$\langle B_{i\Lambda} \rangle \equiv \frac{1}{N} \sum_{i=1}^N B_{i\Lambda}.$$

The ratio R between the posterior doubt and prior doubt, which is called the relative change in doubt, is:

$$R \equiv \frac{D}{p(X)}.$$

For doubt to grow, i. e., the posterior doubt to be greater than the prior doubt ($R \ll 1$), the Bayes factor between the unknown model X and the baseline model must be much greater than the average Bayes factor:

$$\frac{\langle B_{i\Lambda} \rangle}{B_{X\Lambda}} \ll 1.$$

To genuinely doubt the baseline model, Λ CDM, it is not sufficient that $R > 1$, but additionally, the probability of Λ CDM must also decrease such that its posterior probability is greater than its prior probability, i. e., $p(\Lambda|d) < p(\Lambda)$. One can define:

$$R_\Lambda \equiv \frac{p(\Lambda|d)}{p(\Lambda)}.$$

For Λ CDM to be doubted, the following two conditions must be fulfilled:

$$R > 1; \quad R_\Lambda < 1.$$

If these two conditions are fulfilled, then it suggests that the set of known models is incomplete, and gives motivation to search for a better model not yet included, which may lead to model discovery.

In [41], a way of computing an absolute upper bound for $p(d|X)$ achievable among the class of known models has been proposed. Finally, it was found that current cosmic microwave background (CMB), matter power spectrum (mpk), and Type Ia supernovae (SNIa) observations do not require the introduction of an alternative model to the baseline Λ CDM model. The upper bound of the Bayesian evidence for a presently unknown dark energy model against Λ CDM gives only weak evidence in favor of the unknown model. Since this is an absolute upper bound, it was concluded that Λ CDM remains a sufficient phenomenological description of currently available observations.

4 Facilities for the Scientific Hypothesis-Driven Experiment Support

4.1 Conceptualization of scientific experiments

Data intensive research increasingly becomes dependent on computational resources to aid complex researches.

It becomes paramount to offer scientists mechanisms to manage the variety of knowledge produced during such investigations. Specific conceptual modeling facilities [43] are investigated to allow scientists to represent scientific hypotheses, models, and associated computational or simulation interpretations which can be compared against phenomena observations (see Fig. 3). The model allows scientists to record the existing knowledge about an observable investigated phenomenon, including a formal mathematical interpretation of it, if any. Model evolution and model sharing need also to be supported taking either a mathematical or computational view (e. g., expressed by scientific workflows). Declarative representation of scientific model allows scientists to concentrate on the scientific issues to be investigated. Hypotheses can be used also to bridge the gap between an ontological description of studied phenomena and the simulations. Conceptual views on scientific domain entities allow for searching for definitions supporting scientific models sharing among different scientific groups.

In [12], the engineering of hypothesis as linked data is addressed. A semantic view on scientific hypotheses shows their existence apart from a particular statement formulation in some mathematical framework. The mathematical equation is considered as not enough to identify the hypothesis: first, because it must be physically interpreted, and second, because there can be many ways to formulate the same hypothesis. The link to a mathematical expression, however, brings to the hypothesis concept higher semantic precision. Another link, in addition, to an explicit description of the explained phenomenon (emphasizing its “physical interpretation”) can bring forth the intended meaning. By dealing with that hypothesis as a conceptual entity, the scientists make it possible to change its statement formulation or even to assert a semantic mapping to another incarnation of the hypothesis in case someone else reformulates it.

In [43], the following elements related to hypothesis-driven science are conceptualized: a phenomenon observed, a model interpreting this phenomenon, the metadata defining the related computation together with the simulation definition (for simulation, a declarative logic-based language is proposed). In this work, specific attention is devoted to hypothesis definition. The explanation, a scientific hypothesis conveys, is a relationship between the causal phenomena and the simulated one, namely, that the simulated phenomenon is caused by or produced under the conditions set by the causal phenomena. By running the simulations defined by the antecedents in the causal relationship, the scientist aims at providing hypothetical analysis of the studied phenomenon.

Thus, the scientific hypothesis becomes an element of the scientific model that may replace a phenomenon.

When computing a simulation based on a scientific hypothesis, i. e., according to the causal relationship it establishes, the output results may be compared against phenomenon observations to assess the quality of the hypothesis. Such interpretation provides for bridging the gap between qualitative description of the phenomenon domain (scientific hypotheses may be used in qualitative (i. e., ontological) assertions) and the corresponding quantitative valuation obtained through simulations. According to the approach [43], complex scientific models can be expressed as the composition of computation models similarly to database views.

4.2 Hypothesis space browsers

In the HyBrow (Hypothesis Space Browser) project [44], the hypotheses for the biology domain are represented as a set of first-order predicate calculus sentences. In conjunction with an axiom set specified as rules that model known biological facts over the same universe and experimental data, the knowledge base may contradict or validate some of the sentences in hypotheses, leaving the remaining ones as candidates for new discovery. As more experimental data are obtained and rules are identified, discoveries become positive facts or are contradicted. In the case of contradictions, the rules that caused the problems must be identified and eliminated from the theory formed by the hypotheses. In such model-theoretical approach, the validation of hypotheses considers the satisfiability of the logical implications defined in the model with respect to an interpretation. This might be applicable also for simulation-based research, in which validation is solved based on the quantitative analysis between the simulation results and the observations [43]. HyBrow is based on an OWL ontology and application-level rules to contradict or validate hypothetical statements. HyBrow provides for designing hypotheses and evaluating them for consistency with existing knowledge and uses an ontology of hypotheses to represent hypotheses in machine understandable form as relations between objects (agents) and processes [45].

As an upgrade of HyBrow, the HyQue [46] framework adopts linked data technologies and employs Bio2RDF linked data to add to HyBrow semantic interoperability capabilities. HyBrow/HyQue's hypotheses are domain-specific statements that correlate biological processes (seen as events) in the First-Order Logic (FOL). Hypotheses are formulated as instances of the HyQue Hypothesis Ontology and are evaluated through a set of SPARQL queries against biologically-typed OWL and HyBrow data. The query results are scored in terms of how the set of events correspond to background expectations. A score indicates the level of support the data lend the hypothesis. Each event is evaluated independently

in order to quantify the degree of support it provides for the hypothesis posed. Hypothesis scores are linked as properties to the respective hypothesis.

OBI (the Ontology for Biomedical Investigations) project (<http://obi-ontology.org>) aims to model the design of an investigation: the protocols, the instrumentation, and the materials used in experiments and the data generated [47]. Ontologies such as EXPO and OBI enable the recording of the whole structure of scientific investigations: how and why an investigation was executed, what conclusions were made, the basis for these conclusions, etc. As a result of these generic ontology development efforts, the Minimum Information about a Genotyping Experiment (MIGen) recommends the use of terms defined in OBI. The use of a generic or a compliant ontology to supply terms will stimulate cross-disciplinary data-sharing and reuse. As much detail about an investigation as possible in order to make the investigation more reproducible and reusable can be collected [48].

Hypothesis modeling is embedded into the knowledge infrastructures being developed in various branches of science. One example of such infrastructure is considered under the name SWAN — a Semantic Web Application in Neuromedicine [47]. SWAN is a project for developing an integrated knowledge infrastructure for the Alzheimer disease (AD) research community. SWAN incorporates the full biomedical research knowledge lifecycle in its ontological model, including support for personal data organization, hypothesis generation, experimentation, laboratory data organization, and digital prepublication collaboration. The common ontology is specified in an RDF Schema. SWAN's content is intended to cover all stages of the "truth discovery" process in biomedical research, from formulation of questions and hypotheses to capture of experimental data, sharing data with colleagues, and ultimately, the full discovery and publication process.

Several information categories created and managed in SWAN are defined as subclasses of Assertion. They include Publication, Hypothesis, Claim, Concept, Manuscript, DataSet, and Annotation. An Assertion may be made upon any other Assertion, or upon any object specifiable by URL. For example, a scientist can make a Comment upon, or classify, the Hypothesis of another scientist. Linking to objects "outside" SWAN by URL allows one to use SWAN as metadata to organize, for example, all one's PDFs of publications, or the Excel files in which one's laboratory data are stored, or all the websites of tools relevant to Neuroscience. Annotation may be structured or unstructured. Structured annotation means attaching a Concept (tag or term) to an Assertion. Unstructured annotation means attaching free text. Concepts are nodes in controlled vocabularies, which may also be hierarchical (taxonomies).

- (2) use of various resources of biological data as well as human expertise to intelligently generate hypotheses; and
- (3) support for ranking hypotheses and for designing experiments to verify hypotheses.

The extended system is positioned as a prototype of an intelligent research assistant of molecular biologists.

4.4 Hypothesis-driven robots

The Robot Scientist [53] oriented on genomic applications is a physically implemented system which is capable of running cycles of scientific experimentation and discovery in a fully automatic manner: hypothesis formation, experiment selection to test these hypotheses, experiment execution using robotic system, results analysis and interpretation, repeating the cycle (closed-loop in which the results obtained are used for learning from them and feeding the resulting knowledge back into the experimental models). Deduction, induction, and abduction are the types of logical reasoning used in scientific discovery (see section 3). The full automation of science requires ‘closed-loop learning,’ where the computer not only analyses the results, but learns from them and feeds the resulting knowledge back into the next cycle of the process (Fig. 7).

In the Robot Scientist, the automated formation of hypotheses is based on the following key components:

- (1) machine-computable representation of the domain knowledge;
- (2) abductive or inductive inference of novel hypotheses;
- (3) an algorithm for the selection of hypotheses; and

- (4) deduction of the experimental consequences of hypotheses.

Adam, the first Robot Scientist prototype, was designed to carry out microbial growth experiments to study functional genomics in the yeast *Saccharomyces cerevisiae*, specifically to identify the genes encoding ‘locally orphan enzymes.’ Adam uses a comprehensive logical model of yeast metabolism, coupled with a bioinformatic database (Kyoto Encyclopaedia of Genes and Genomes — KEGG) and standard bioinformatics homology search techniques (PSI-BLAST and FASTA) to hypothesize likely candidate genes that may encode the locally orphan enzymes. This hypothesis generation process is abductive.

To formalize Adam’s functional genomics experiments, the LABORS ontology (LABORatory Ontology for Robot Scientists) has been developed. LABORS is a version of the ontology EXPO (as an upper layer ontology) customized for Robot scientists to describe biological knowledge. LABORS is expressed in OWL-DL. LABORS defines various structural research units, e.g., trial, study, cycle of study and replicate as well as design strategy, plate layout, expected actual results. The respective concepts and relations in the functional genomics data and metadata are also defined. Both LABORS and the corresponding database (used for storing the instances of the classes) are translated into Datalog in order to use the SWI-Prolog reasoner for required applications [48].

There were two types of hypotheses generated. The first level links an orphan enzyme, represented by its enzyme class (E.C.) number, to a gene (ORF) that potentially encodes it. This relation is expressed as a two-place predicate where the first argument is the ORF and the second is the E.C. number. An example of

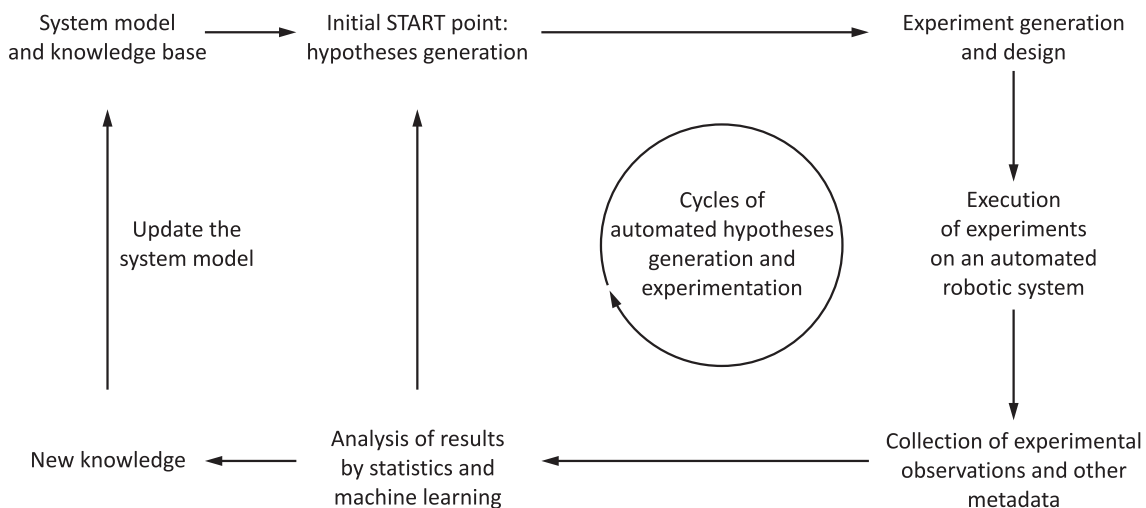


Figure 7 Hypothesis-driven closed-loop learning

hypothesis at this level is: *encodesORFtoEC('YBR166C', '1.1.1.25')*.

The second level of hypothesis involves the association between a specific strain, referenced via the name of its missing ORF, and a chemical compound which should affect the growth of the strain, if added as a nutrient to its environment. This level of hypothesis is derived from the first by logical inference using a specific model of yeast metabolism. An example of such a hypothesis is: *affects growth('C00108', 'YBR166C')*, where the first argument is the compound (names according to KEGG) and the second argument is the strain considered.

Adam then designs the experimental assays required to test these hypotheses for execution on the laboratory robotic system. These experiments are based on a two-factor design that compares multiple replicates of the strains with and without metabolites compared against wild type strain controls with and without metabolites.

Adam follows a hypothetico-deductive methodology (see section 2). Adam abductively hypothesizes new facts about yeast functional biology, then it deduces the experimental consequences of these facts using its model of metabolism, which it then experimentally tests. To select experiments, Adam takes into account the variable cost of experiments, and the different probabilities of hypotheses. Adam chooses its experiments to minimize the expected cost of eliminating all but one hypothesis. This is, in general, an NP complete problem and Adam uses heuristics to find a solution [45].

It is now likely that the majority of hypotheses in biology are computer-generated. Computers are increasingly automating the process of hypothesis formation, for example: machine learning programs (based on induction) are used in chemistry to help design drugs; and in biology, genome annotation is essentially a vast process of (abductive) hypothesis formation. Such computer-generated hypotheses have been necessarily expressed in a computationally amenable way, but it is still not common practice to deposit them into a public database and make them available for processing by other applications [45].

The details describing the software and informatics decisions in the Robot Scientist project can be found in [45, 53] and online at the website <http://www.aber.ac.uk/compsci/Research/bio/robotsci/data/informatics/>. The details for developing the formalization used for Adam's functional genomics investigations can be found in [48, 54]. An ontology-based formalization based on graph theory and logical modeling makes it possible to keep an accurate track of all the result units used for different goals, while preserving the semantics of all the experimental entities involved in all the investigations. It is shown how experimentation and machine learning are used to identify additional knowledge to improve the metabolic model [54].

4.5 Hypotheses as data in probabilistic databases

Another view of hypotheses encoding and management is presented in [55]. Authors use probabilistic database techniques for hypotheses systematic construction and management. MayBMS [56], a probabilistic database management system, is used as a core for hypothesis management. This methodology (called γ -DB) enables researchers to maintain several hypotheses explaining some phenomena and provides evaluation mechanism based on Bayesian approach to rank them.

The construction of γ -DB database comprises several steps. In the first step, phenomenon and hypothesis entities are provided as input to the system. Hypothesis is a set of mathematical equations expressed as functions in W3C MathML-based format and is associated with one or more simulation trial dataset, consisting of tuples with input variables of equation and its corresponding output as functionally dependent variables (the predictions). Phenomenon is represented by at least one empirical dataset similar to simulation trials. In the next step, the system deals with hypotheses and phenomena in the following way:

- (1) researcher has to provide some metadata about hypotheses and phenomena; e. g., hypotheses need to be associated with the respective phenomena and assigned a prior confidence distribution (uniform by default according to the principle of maximum entropy (see paragraph 3.2.3));
- (2) functional dependencies (FD) are extracted from equations in order to obtain database schema to store simulations and experimental data; it should be mentioned that to precisely identify hypothesis formulation, the special attributes for phenomena and hypothesis references are introduced into FD;
- (3) tuples are synthesized from simulation trials and observational data by uncertain pseudotransitive closure and reasoning; and finally,
- (4) the probabilistic γ -DB database is formed.

Once phenomenon and hypothesis (with empirical datasets and simulation trials) are produced, it becomes possible to manipulate them with database tools.

MayBMS provides tools to evaluate competing hypotheses for the explanation of a single phenomenon. With prior probabilities already provided, the system allows to make one or more (if new observational data appears) Bayesian inference steps. In each step, the prior probability is updated to posterior according to Bayes' theorem. As a result, hypotheses which better explain phenomenon get higher probabilities enabling researchers to make more confident decisions (see also paragraph 3.2.1). The γ -DB approach provides a promising way to analyze hypotheses in large-scale

DIR as uncertain predictive database in face of empirical data.

5 Examples of Hypothesis-Driven Scientific Research

5.1 Hypotheses in Besançon Galaxy model

Various models in astronomy heavily rely on hypotheses. One of the most impressive is the Besançon galaxy model (BGM) [57–59] evolving for many years and representing the population and structure synthesis model for the Milky Way. It allows astronomers to test hypotheses on the star formation history, star evolution, and chemical and dynamical evolution of the Galaxy. As the result of simulation process, one can get the following: multidimensional histograms of intrinsic star properties or observable properties, a catalog of pseudoobservations, or the integrated luminosity in a specified photometric band [60]. From the beginning, the aim of the BGM was not only to be able to simulate reasonable star counts but further to test scenarios of Galactic evolution from assumptions on the rate of star formation (SFR), initial mass function (IMF), and stellar evolution.

The model has explicit and implicit hypotheses associated with it. Explicit hypotheses are usually some sets of equations, taken from the literature studies and put as the ingredient of the model. Some of explicit hypotheses are passed as the input of the model, e.g., star formation rate, initial mass function, evolutionary tracks, chemical evolution, atmosphere models, density laws, interstellar extinction model.

The model has some implicit hypotheses as well. For example, it is assumed that no star population comes from the outside of the Galaxy. There are several more implicit hypotheses about disk formation and dark matter assumptions encoded inside the model. It is usually much harder to get all the implicit hypotheses, since many of them are not described in the articles and are difficult to pin from the code.

BGM has not only the large number of explicit and implicit hypotheses, but also a complex interrelations between them. So, some of the hypotheses are being independent, e.g., IMF and SFR; so, it is possible to change them independently. On the other side, some of the hypotheses are connected, e.g., the age distribution, the density laws, and the potential are linked with the age–velocity dispersion via the Boltzmann equation and need to be consistent. Such kind of dependencies make the model hard to be tested and to keep it consistent while varying different parameters during model fitting. Another example of interrelations of hypotheses is competing hypotheses.

BGM has changed drastically over the last 30 years. This has happened because of the appearance of new

data surveys, technologies, and methods of observation development. As an example of such evolution, the model developed in 2014 compared to previous versions handles variations of the SFR, IMF, evolutionary tracks, and atmosphere models. These hypotheses are passed as input parameters to the model; so, the user can vary them.

The second improvement of the model is the implementation of the stellar binarity, being an important change since binaries can account for about 50% of the total stellar content of the Milky Way. The authors of the new version underline the importance of understanding interrelationships between different hypotheses and need for model evolution tools [60]: “In practice, to build a Galaxy from the fundamental building-blocks, we had to reconstruct the previous model and apply important changes in the code arrangement. That required to understand well the underlying relations between all mentioned components.”

It is planned further to focus on the renewed BGM [59], in which authors draw their attention to the Galaxy thin disk treatment and use of Tycho-2 as a testing dataset. The parameters of BGM (such as IMF, SFR and evolutionary track sets) explicitly and model ingredients implicitly can be treated as hypotheses. Model ingredients include the treatment of binarity, the local stellar mass densities of thin disk, extinction model, age-metallicity and age-velocity relations, radial scale length, the age of the Galaxy thin disc, different sets of the star atmosphere models, etc.

Tycho-2 dataset and χ^2 -type statistics test is used to test various versions of these hypotheses in order to choose the most appropriate ones and update model to better fit the provided data. The tests were made by comparing star counts and $(B - V)_T$ color distributions between data and simulations. Two different tests were used to evaluate the adequacy of the stellar densities globally and to test the shape of the color distribution. Other parameters to be tested are: star counts, radio velocity, magnitudes, colors, proper motions, parallax, effective temperatures, gravity, and metallicity. Authors use histograms, 2 goodness of fit (maximum likelihood and χ^2 -test) and for velocity parameter, Kolmogorov–Smirnov and Henderson–Darling tests.

Due to the fact that some ingredients of the model are highly correlated (such as the IMF, SFR, and the local mass density), the authors defined default models as a combination of a new set of ingredients that significantly improve the fit to Tycho data. So, 11 IMF functions, 2 SFR functions, 2 evolutionary track sets, 3 sets of atmosphere models, 3 values for the age of the formation of the thin disk, and 3 sets of values of the thin disk local stellar volume mass density were tested. As a result of testing, the two most appropriate IMF and SFR hypotheses were chosen.

BGM authors have plans to incorporate other star surveys and test the model against them. To do simulations directly comparable with data, the selected magnitudes from the surveys need to be complete in terms of magnitude. Among these surveys, there are the Geneva-Copenhagen survey, SDSS-II/III, SEGUE/SEGUE2, APOGEE, RAVE, LAMOST, Gaia, Gaia-ESO, GALAH LSST, WEAVE, 4MOST, and MOONS surveys [61].

5.2 Hypothesis testing applying connectome data

In the neuroscience community, the development of common paradigms for interrogating the myriad functional systems in the brain remains to be the core challenge. Building on the term “*connectome*,” coined to describe the comprehensive map of neural connections in the human brain, the “functional connectome” denotes the collective set of functional connections in the human brain (its “wiring diagram”) [62]. More broadly, a connectome would include the mapping of all neural connections within an organism’s nervous system. The production and study of connectomes, known as *connectomics*, may range in scale from a detailed map of the full set of neurons and synapses within part or all of the nervous system of an organism to a macroscale description [63] of the functional and structural connectivity between all cortical areas and subcortical structures. The ultimate goal of connectomics is to map the human brain. In functional magnetic resonance imaging (fMRI), associations are thought to represent functional connectivity in the sense that the two regions of the brain participate together in the achievement of some higher-order function, often in the context of performing some task. fMRI has emerged as a powerful tool used to interrogate a multitude of functional circuits simultaneously. This has elicited the interest of statisticians working in that area. At the level of basic measurements, neuroimaging data can be considered to consist typically of a set of signals (usually, time series) at each of a collection of pixels (in two dimensions) or voxels (in three dimensions). Building from such data, various forms of higher-level data representations are employed in neuroimaging. In recent years, a substantial interest in network-based representations has emerged in neuroimaging to use *networks* to summarize relational information in a set of measurements, typically assumed to be reflective of either functional or structural relationships between regions of interest in the brain. With neuroimaging, now, a standard tool in clinical neuroscience, quickly moving towards a time in which we will have available databases composed of large collections of secondary data in the form of *network-based data objects*, is predictable.

One of the most basic tasks of interest in the analysis of such data is the testing of hypotheses in answer to questions such as “Is there a difference between the networks of these two groups of subjects?” Networks are not Euclidean objects and, hence, classical methods of statistics do not directly apply. Network-based analogues of classical tools for statistical estimation and hypothesis testing are investigated in [64, 65]. Such research is motivated by the 1000 Functional Connectomes Project (FCP) launched in 2010 [62]. The 1000 FCP [66] constitutes the largest data set of its kind similarly to large data sets in genetics. Other projects (such as the Human Connectome Project (HCP)) are aimed to build a network map of the human brain in healthy, living adults. The total volume of data produced by the HCP will likely be multiple petabytes [67]. HCP informatics platform includes data management system ConnectomeDB that is based on the XNAT (eXtensive Neuroimaging Archive Toolkit) imaging informatics platform [68], a widely used open source system for managing and sharing imaging and related data.

Now, HCP has information about more than 500 subjects including structural scans (T1w and T2w), resting-state fMRI (rfMRI), task fMRI (tfMRI), and high angular resolution diffusion imaging (dMRI). In addition, some resting-state MEG (rMEG) and/or task MEG (tMEG) data are available.

Data come in several formats: “unprocessed” raw data, “minimally preprocessed,” and “analysis” datasets. Preprocessed datasets have spatial distortions minimized and data have been aligned across modalities and across subjects using appropriate volume-based and surface-based registration methods. HCP consortium recommends to use the preprocessing dataset.

Visualization, processing, and analysis of high-dimensional data such as images often require some kind of preprocessing to reduce the dimensionality of the data and find a mapping from the original representation to a low-dimensional vector space. The assumption is that the original data resides in a low-dimensional subspace or manifold [69], embedded in the original space. This topic of research is called dimensionality reduction, nonlinear dimensionality reduction, including methods for parameterization of data using low-dimensional manifolds as models. Within the neural information processing community, this has become known as manifold learning. Methods for manifold learning are able to find nonlinear manifold parameterizations of datapoints residing in high-dimensional spaces, very much like Principal Component Analysis (PCA) is able to learn or identify the most important linear subspace of a set of data points (projecting data on a n -dimensional linear subspace which maximizes the variance of the data in the new space).

In [64], necessary mathematical properties associated with a certain notion of a ‘space’ of networks used to interpret functional neuroimaging connectome-oriented data are established. Extension of the classical statistics tools to network-based datasets, however, appeared to be highly nontrivial. The main challenge in such an extension is due to the fact that networks are not Euclidean objects (for which classical methods were developed) — rather, they are combinatorial objects, defined through their sets of vertices and edges. In [64], it was shown that networks can be associated with certain natural subsets of Euclidean space and demonstrated that through a combination of tools from geometry, probability on manifolds, and high-dimensional statistical analysis, it is possible to develop a principled and practical framework in analogy to classical tools. In particular, an asymptotic framework for one- and two-sample hypothesis testing has been developed. Key to this approach is the correspondence between an undirected graph and its Laplacian, where the latter is defined as a matrix (associating with a network). Graph Laplacian appeared to be particularly appropriate to be used for such matrices. The space of graph Laplacians is used working in certain subsets of Euclidean space which are some submanifolds of the standard Euclidean space.

The 1000 FCP describes functional neuroimaging data from 1093 subjects, located in 24 community-based centers. The mean age of the participants was 29 years, and all subjects were 18 years old or older. It is of interest to compare the subject-specific networks of males and females in the 1000 FCP data set. In [64], for the 1000 FCP, database comparing networks with respect to the sex of the subjects, over different age group, and over various collection sites is considered. It is shown that it is necessary to compute the means in each subgroup of networks. This was done by constructing the Euclidean mean of the Laplacians for each group of subjects in different age groups. Such group-specific mean Laplacians can then be interpreted as the mean functional connectivity in each group. Such approach provides for building the hypothesis tests about the average of networks or groups of networks to investigate the effect of sex differences on entire networks.

For the 1000 FCP data set, it was tested using the two-sample test for Laplacians whether sex differences were significant to influence patterns of brain connectivity. The null hypothesis of no group differences was rejected with high probability. Similarly for the three different age cohorts, the null hypothesis of no cohort differences also was rejected with high probability.

On such examples, it was shown [64] that the proposed global test has sufficient power to reject the null hypothesis in cases when mass-univariate approach (considered to be the gold standard in fMRI research [70]) fails to detect the differences at the local level. Accord-

ing to the mass-univariate approach, statistical analysis is performed iteratively on all voxels to identify brain regions whose fMRI detected responses display significant statistical effects. Thus, it was shown that a framework for network-based statistical testing is more statistically powerful than a mass-univariate approach.

It is expected that in the near future, there will be a plethora of databases of network-based objects in neuroscience motivating the development and extension of various tools from classical statistics to global network data.

In paper [71] discussing the relationship between neuroimaging and Big Data areas, it is analyzed how modern neuroimaging research represents a multifactorial and broad ranging data challenge, involving the growing size of the data being acquired; sociological and logistical sharing issues; infrastructural challenges for multisite, multitype archiving; and the means by which to explore and mine these data. As neuroimaging advances further, e. g., aging, genetics, and age-related disease, new vision is needed to manage and process this information while marshalling of these resources into novel results. It is predicted that on this way, “big data” can become “big” brain science.

In [72], authors formulate a hypothesis about the brain connectivity and evaluate it against HCP data. They use the task fMRI data, there is specific data about a well-validated task used to probe animate motion detection. The audience was shown short video clips (20 s) of objects (squares, circles, and triangles) either interacting in some way (animate motion) or moving mechanically (inanimate motion). Participants rated the video by selecting if there was any social interaction, no interaction, or not sure for interaction. There were 2 sessions comprised of 5 videoblocks.

Hypothesis states that some regions of the brain (V5 and pSTS) are effectively connected and impacted by animate motion.

To test it, general linear models were used. The time series were modeled with regressors All motion—No motion, Animate—Inanimate motion. Together with regressors about head, tongue, and finger movement, these regressors were used to build general linear model. A group level ANOVA was performed to identify significant regional effects for the All Motion contrast and a contrast for Animate—Inanimate motion. For effective connectivity discovery, Dynamic Causal Modeling (DCM) technique was used. DCM tells about self-, forward, and backward connections between active brain regions during an experiment, enabling to infer the way of brain regions impact each other mostly. As the result of DCM modeling, 16 models were created and passed as the input to Bayesian Model Selection procedure, which chose the winning model among them.

The results show that there is a connectivity between V5 and the pSTS brain regions in both hemispheres, which was independent of the type of motion. Animate motion stimulates the forward and backward connection between V5 and the pSTS in both hemispheres.

5.3 Climate in Australia

Another view on hypothesis representation and evaluation is presented in [73]. Authors argue that as long as in DIR data relevant to some hypotheses get continuously aggregated as time passes, hypotheses should be represented as programs that are executed repeatedly, as new relevant amounts of data get aggregated. Their method and techniques are illustrated by examining hypotheses about temperature trends in Australia during the 20th century. The hypothesis being tested comes from [74], stated that the temperature series is not stationary and is integrated of order 1 ($I(1)$). Nonstationarity means that the level of the time series is not stable in time and can show increasing and decreasing trends; $I(1)$ means that by differentiating the stochastic process, a stationary process (main statistical properties of the series remain unchanged) is obtained. Phillips–Perron test and the Kwiatkowski–Phillips–Schmidt–Shin (KPSS) test are used and both of them are executed in R. Several data sources are crawled: (i) The National Oceanographic and Atmospheric Administration marine and weather information; and (ii) Australian Bureau of Meteorology dataset. The framework consists of R interpreter and R *SPARQL*, *tseries* packages. Authors also used agINFRA for computation and rich semantics to support traditional scientific workflows for natural sciences. Authors received further evidence on different independent dataset that time series is integrated of order 1.

5.4 Financial market

Efficient-market hypothesis (EMH) is one of the most prominent in finance and “*asserts that financial markets are “informationally efficient.”*” In [75], authors test the weak form of EMH, stating that prices on traded assets (e. g., stocks, bonds, or property) already reflect all past publicly available information. The null hypothesis states that successive prices changes are independent (random walk). The alternative hypothesis states that they are dependent. To check if the successive closing prices are dependent of each other, the following statistical tests were used: a serial correlation test, a runs test, an augmented Dickey–Fuller test, and the multiple variance ratio test. Tests were performed on daily closing prices from the six European stock markets (France, Germany, U.K., Greece, Portugal, and Spain) during the period between 1993 and 2007. The result of each test

states whether successive closing prices are dependent of each other.

Test provides evidence that for monthly prices and returns, the null hypothesis should not be rejected for all six markets. If daily prices are concerned, the null hypothesis is not rejected for France, Germany, U.K., and Spain, but this hypothesis is rejected for Greece and Portugal. However, on the 2003–2007 dataset, the null hypothesis for these two countries is not rejected as well.

In [76], Bollen *et al.* use different approach to test EMH. Authors investigate whether public sentiment, as expressed in large-scale collections of daily Twitter posts, can be used to predict the stock market. They build public mood time series by sentiment analysis of tweets from February 28 to December 19, 2008 and try to show that it can predict Dow Jones Index corresponding values. The null hypothesis states that the mood time series do not predict DJIA (Dow Jones Industrial Average) values. Granger causality analysis in which Dow Jones values and mood time series are correlated is used to test the null hypothesis. Granger causality analysis is used to determine if one time series can predict another time series. Its results reject the null hypothesis and claim that public opinion is predictive of changes in DJIA closing values.

5.5 Publication-based automated hypothesis generation in life sciences

Researchers and scientists from leading academic, pharmaceutical, and other research centers have begun deploying IBM’s Watson Discovery Advisor to rapidly analyze and test hypotheses using data in millions of scientific papers available in public databases. A new scientific research paper is published nearly every 30 s, which equals more than a million annually. According to the National Institutes of Health, a typical researcher reads about 23 scientific papers per month, which translates to nearly 300 per year, making it humanly impossible to keep up with the evergrowing body of scientific material available. Building on Watson’s ability to understand nuances in natural language, Watson Discovery Advisor can understand the language of science, such as how chemical compounds interact, making it a uniquely powerful tool for researchers in life sciences and other research and industrial domains. Specifically, the Watson Discovery Advisor for Life Sciences is armed with expertise and understands field-specific lexicon in areas such as clinical trial data, genomics, drugs, and human anatomy.

Recently, scientists of Baylor College of Medicine and IBM using the Baylor Knowledge Integration Toolkit (KnIT), based on Watson technology, identified new enzymes (called kinases) that can modify p53, an important protein related to many cancers [77]. There are

over 240,000 papers that mention one or more of 500+ known human kinases in their Medline abstract. There are over 70,000 papers published on p53 to make their analysis manually is completely unrealistic task. Watson analyzed the scientific articles related to p53 to predict proteins that turn on or off p53's activity. This automated analysis led the Baylor cancer researchers to identify six potential p53 kinases to target for new research. These results are notable, considering that over the last 30 years, scientists averaged one p53 kinase discovery per year. Knowing which proteins are modified by each kinase and, therefore, which kinases would make good drug targets is a difficult and unsolved problem. There are over 500 known human kinases and tens of thousands of possible proteins they can target.

KnIT collects the abstracts to be mined applying queries. A specific kinase name and its synonyms are used in this process. Entity resolution process looks as follows. The words and phrases that make up the document feature space are determined by counting the number of documents in which each word appears and identifying the words with the highest counts. A phrase is considered to be a sequence of two words. Only the N most frequent words and phrases are selected. This becomes the feature space.

Once a feature space is received, a representation of each kinase by averaging the feature vectors of all documents that contain the kinase is created. This is the kinase centroid. Next, a distance matrix is calculated that measures the distance between each kinase and every other kinase in the space.

Finally, a meaningful picture of kinase–kinase relationships is obtained. Thus, it is possible to identify a set of kinases that may modify p53. However, some sort of principled ranking scheme is needed in order to prioritize the kinases for further experimentation. To provide such a scheme, the graph diffusion method [78] was used. Graph diffusion is a semisupervised learning approach for classification based on labeled and unlabeled data. It takes known information (initial labels) and then constrains the new labels to be smooth in respect to a defined structure (e. g., a network). In the case considered, it is known which kinases can modify p53 (initial labels); one would like to know which other proteins can modify p53 (final labels). The distance matrix based on the literature gives the structure of the kinase network. The initial labels are extracted from current knowledge found in review articles.

To test the algorithm, it was first applied in a retrospective analysis to show whether recent annotations of new p53 kinases occurring after a certain date (2003) could be predicted from a model that only took into account papers written before that date, at a time when these discoveries of p53 kinases were still unknown. Next, it was asked whether some variations in

the algorithm could improve p53 kinase prediction as its performance was compared to the common approach used most typically to identify functionally similar proteins in biology. Finally, the analysis was expanded to a larger set of proteins to test scalability.

This research represents the first stage in the IBM–Baylor collaborative effort and as such, it proves the principle that mining past literature is a viable strategy for predicting previously unknown biological events. It was shown that p53 kinases predicted with the text mining methods are supported by laboratory findings. In the future, it should be possible to make many other kinds of predictions on a much larger scale as the infrastructure and capabilities will be increased. In the future, it is planned to focus on a wider area of proteins and functions, building up comprehensive networks of interactions and predicting where new connections ought to exist based on everything else that is known. It is expected that this will ultimately accelerate the pace of cancer discoveries by an order of magnitude and allow scientists to come to a much more complete understanding of the mechanisms behind this disease.

Expanding KnIT to other areas of biology or the physical sciences is not straightforward. For example, to generalize to more proteins and genes is a big problem. In subjects like physics, results tend to be presented using equations and graphs rather than words. However, data-mining groups are working to retrieve information from these, too.

6 Concluding Remarks

The objective of this study is to analyze, collect, and systematize information on the role of hypotheses in the DIR process as well as on support of hypothesis formation, evaluation, selection, and refinement in course of the natural phenomena modeling and scientific experiments. The discussion is started with the basic concepts defining the role of hypotheses in the formation of scientific knowledge and organization of the scientific experiments. Based on such concepts, the basic approaches for hypothesis formulation applying logical reasoning, various methods for hypothesis modeling and testing (including classical statistics, Bayesian hypothesis, and parameter estimation methods, hypothetico-deductive approaches) are briefly introduced. Special attention is given to discussion of the data mining and machine learning methods role in process of generation, selection, and evaluation of hypotheses as well as the methods for motivation of new hypothesis formulation. Facilities of informatics for support of hypothesis-driven experiments, considered in the paper, are aimed at the conceptualization of scientific experiments, hypothesis formulation, and browsing in various domains (includ-

ing biology, biomedical investigations, neuromedicine, and astronomy), automatic organization of hypothesis-driven experiments. Examples of scientific researches applying hypotheses considered in the paper include modeling of population and structure synthesis of the Galaxy, connectome-related hypothesis testing, studying of temperature trends in Australia, analysis of stock markets applying the EMH, as well as algorithmic generation of hypotheses in the collaborative project based on IBM Watson—Baylor Knowledge Integration Toolkit applying the NLP and knowledge representation and reasoning technologies. An introduction into the state of the art of the hypothesis-driven research presented in the paper opens a way for investigation of the generalized approaches for efficient organization of hypothesis-driven experiments applicable for various branches of DIR.

References

- Hey, T., S. Tansley, and K. Tolle, eds. 2009. *The Fourth paradigm: Data-intensive scientific discovery*. Redmond, Microsoft Research. 252 p.
- McComas, W. F. 1998. The principal elements of the nature of science: Dispelling the myths of science. *Nature of science in science education: Rationales and strategies*. Ed. W. F. McComas. Kluwer Academic Pubs. 53–70.
- Lakshmana Rao, J. R. 1998. Scientific ‘Laws,’ ‘Hypotheses’ and ‘Theories’. *Meanings Distinctions Reson.* 3:69–74.
- Poincaré, H. 2012. The foundations of science: Science and hypothesis, the value of science, science and method. *The Project Gutenberg EBook*. No. 39713. 554 p. Available at: <http://www.gutenberg.org/files/39713/39713-8.txt> (accessed February 10, 2015).
- Bacon, F. 1952. The new organon. *Great books of the Western World. Vol. 30. The works of Francis Bacon*. Ed. R. M. Hutchins. Chicago: Encyclopedia Britannica, Inc. 107–195.
- Menzies, T. 1996. Applications of abduction: Knowledge-level modeling. *Int. J. Hum.-Comput. St.* 45(3):305–335.
- Haber, J. 2010. Research questions, hypotheses, and clinical questions. *Evolve resources for nursing research*. 7th ed. Elsevier. 27–55.
- Popper, K. 2005. *The logic of scientific discovery*. London – New York: Routledge, Taylor & Francis. 545 p. Available at: <http://strangebeautiful.com/other-texts/popper-logic-scientific-discovery.pdf> (accessed February 10, 2015).
- Kerlinger, F. N., and H. B. Lee. 1964. *Foundations of behavioral research: Educational and psychological inquiry*. New York: Holt, Rinehart and Winston. 739 p.
- Hempel, C. G. 1952. Fundamentals of concept formation in empirical science. *Int. Encyclopedia Unified Sci.* 2(7). Available at: <http://www.iep.utm.edu/hempel/> (accessed February 10, 2015).
- Porto, F., and S. Spaccapietra. 2011. Data model for scientific models and hypotheses. *Evolution Conceptual Modeling* 6520:285–305.
- Gonçalves, B., and F. Porto. 2013. A lattice-theoretic approach for representing and managing hypothesis-driven research. *25th Conference (International) on Scientific and Statistical Database Management (ACM) Proceedings*. Baltimore. 41.
- Gonçalves, B., F. Porto, and A. M. C. Moura. 2012. On the semantic engineering of scientific hypotheses as linked data. *2nd Workshop (International) on Linked Science Proceedings*. Boston.
- Woodward, J. 2011. Scientific explanation. *The Stanford Encyclopedia of Philosophy*. Available at: <http://plato.stanford.edu/archives/win2011/entries/scientific-explanation/> (accessed February 10, 2015).
- Nickles, T., ed. 1980. *Scientific discovery: Case studies*. Taylor & Francis. 501 p.
- Schickore, J. 2014. Scientific discovery. *The Stanford Encyclopedia of Philosophy*. Available at: <http://plato.stanford.edu/archives/spr2014/entries/scientific-discovery/> (accessed February 10, 2015).
- Kakas, A. C., R. A. Kowalski, and F. Toni. 1993. Abductive logic programming. *J. Logic Comput.* 2(6):719–770.
- Kakas, A. C., A. Michael, and C. Mourlas. 2000. ACLP: Abductive constraint logic programming. *J. Logic Program.* 44(1):129–177.
- Van Nuffelen, B., and A. Kakas. 2001. A-system: Declarative programming with abduction. *Logic programming and nonmonotonic reasoning*. Eds. T. Eiter, W. Faber, and M. Truszczyński. Lecture notes in computer science ser. Berlin–Heidelberg: Springer. 2173:393–397.
- Alferes, J. J., L. M. Pereira, and T. Swift. 2004. Abduction in well-founded semantics and generalized stable models via tabled dual programs. *Theor. Pract. Log. Progr.* 4(4):383–428.
- Ray, O., and A. Kakas. 2006. ProLogICA: A practical system for Abductive Logic Programming. *11th Workshop (International) on Non-Monotonic Reasoning Proceedings*. 304–312.
- Citrigno, S., T. Eiter, W. Faber, G. Gottlob, C. Koch, N. Leone, and F. Scarcello. 1997. The div system: Model generator and application frontends. *12th Workshop on Logic Programming Proceedings*. 128–137.
- King, R. D., M. Liakata, C. Lu, S. G. Oliver, and L. N. Soldatova. 2011. On the formalization and reuse of scientific research. *J. Roy. Soc. Interface* 8(63):1440–1448.
- Tamaddoni-Nezhad, A., R. Chaleil, A. Kakas, and S. H. Muggleton. 2006. Application of abductive ILP to learning metabolic network inhibition from temporal data. *Mach. Learn.* 64:209–230.
- Inoue K., T. Sato, M. Ishihata, Y. Kameya, and H. Nabeshima. 2009. Evaluating abductive hypotheses using and EM algorithm on BDDs. *21st Joint Conference (International) on Artificial Intelligence (IJCAI09) Proceedings*. Pasadena. 810–815.
- Bartha, P. 2013. Analogy and analogical reasoning. *The Stanford Encyclopedia of Philosophy*. Available at: <http://plato.stanford.edu/archives/fall2013/entries/reasoning-analogy/> (accessed February 10, 2015).

27. Ivezić, Ž., A. J. Connolly, J. T. VanderPlas, and A. Gray. 2014. *Statistics, data mining, and machine learning in astronomy: A practical Python guide for the analysis of survey data*. Princeton University Press. 552 p.
28. Sivia, D. S., and J. Skilling. 2006. *Data analysis. A Bayesian tutorial*. New York: Oxford University Press Inc. 264 p.
29. Field, A. 2013. *Discovering statistics using IBM SPSS statistics*. 4th ed. Sage. 915 p.
30. IBM SPSS Statistics for Windows, Version 22.0. 2013. Armonk, N.Y.: IBM Corp. IBM SPSS Statistics base. Available at: https://www.uio.no/tjenester/it/forskning/statistikk/hjelp/programveilednigner/ibm_spss_statistics_brief_guide-2.pdf (accessed February 10, 2015).
31. Ihaka, R., and R. Gentleman. 1996. R: A language for data analysis and graphics. *J. Comput. Graph. Stat.* 5(3):299–314.
32. March, M. C., G. D. Starkman, R. Trotta, and P. M. Vaudrevange. 2011. Should we doubt the cosmological constant? *Mon. Not. Roy. Astron. Soc.* 410(4):2488–2496.
33. Rouder, J. N., P. L. Speckman, D. Sun, R. D. Morey, and G. Iverson. 2009. Bayesian t tests for accepting and rejecting the null hypothesis. *Psychon. Bull. Rev.* 16(2):225–237.
34. Weber, M. 2014. Experiment in biology. *The Stanford Encyclopedia of Philosophy*. Available at: <http://plato.stanford.edu/archives/fall2014/entries/biology-experiment/> (accessed February 10, 2015).
35. Hawthorne, J. 2014. Inductive logic. *The Stanford Encyclopedia of Philosophy*. Available at: <http://plato.stanford.edu/archives/sum2014/entries/logic-inductive/> (accessed February 10, 2015).
36. Breiman, L. 2001. Statistical modeling: The two cultures. *Stat. Sci.* 16(3):199–231.
37. Hastie, T., R. Tibshirani, J. Friedman, and J. Franklin. 2005. The elements of statistical learning: Data mining, inference and prediction. *Math. Intell.* 27(2):83–85.
38. Barber, D. 2010. *Bayesian reasoning and machine learning*. Cambridge University Press. 720 p.
39. Ferrucci, D., E. Brown, J. Chu-Carroll, J. Fan, D. Gondek, A. A. Kalyanpur, and C. Welty. 2010. Building Watson: An overview of the DeepQA project. *AI Mag.* 31(3):59–79.
40. Dredze, M., K. Crammer, and F. Pereira. 2008. Confidence-weighted linear classification. *25th Conference (International) on Machine Learning Proceedings*. Helsinki. 264–271.
41. Starkman, G. D., R. Trotta, and P. M. Vaudrevange. 2008. Introducing doubt in Bayesian model comparison. arXiv preprint arXiv:0811.2415.
42. March, M. C. 2013. Advanced statistical methods for astrophysical probes of cosmology. Springer Theses. Vol. 20. 177 p.
43. Porto, F. 2013. Big data in astronomy. The LINEAR-DEXL case. *EMC Summer School on BIG DATA — NCE/UFRJ*. Available at: <http://www.slideshare.net/fabiomporto/emc-2013-big-data-in-astronomy> (accessed February 10, 2015).
44. Racunas, S. A., N. H. Shah, I. Albert, and N. V. Fedoroff. 2004. Hybrow: A prototype system for computer-aided hypothesis evaluation. *Bioinformatics* 20(1):257–264.
45. Soldatova, L. N., A. Rzhetsky, and R. D. King. 2011. Representation of research hypotheses. *J. Biomed. Semantics* 2(S-2):S9.
46. Callahan, A., M. Duumontier, and N. Shah. 2011. HyQue: Evaluating hypotheses using Semantic Web technologies. *J. Biomed. Semantics* 2(S-2):S3.
47. Gao, Y., J. Kinoshita, E. Wu, E. Miller, R. Lee, A. Seaborne, and T. Clark. 2006. SWAN: A distributed knowledge infrastructure for Alzheimer disease research. *J. Web Semant.* 4(3):222–228.
48. King, R. D., K. E. Whelan, F. M. Jones, P. G. Reiser, C. H. Bryant, S. H. Muggleton, and S. G. Oliver. 2004. Functional genomic hypothesis generation and experimentation by a robot scientist. *Nature* 427(6971):247–252.
49. Porto, F., A. M. C. Moura, B. Gonçalves, R. Costa, and S. A. Spaccapietra. 2012. A scientific hypothesis conceptual model. *Advances in conceptual modeling*. Eds. S. Castano, P. Vassiliadis, L. V. Lakshmanan, and M. Li Lee. Lecture notes in computer science ser. Berlin—Heidelberg: Springer. 7518:101–110.
50. Porto, F., and A. M. C. Moura. 2011. Scientific hypothesis database. Report. Available at: <http://livroaberto.ibict.br/bitstream/1/869/1/Scientific%20Hypothesis%20Database.pdf> (accessed February 10, 2015).
51. Asgharbeygi, N., P. Langley, S. Bay, and K. Arrigo. 2006. Inductive revision of quantitative process models. *Ecol. Model.* 194(1):70–79.
52. Tran, N., C. Baral, V. J. Nagaraj, and L. Joshi. 2005. Knowledge-based integrative framework for hypothesis formation in biochemical networks. *Data integration in the life sciences*. Eds. B. Ludäscher and L. Raschid. Lecture notes in computer science ser. Berlin—Heidelberg: Springer. 3615:121–136.
53. Sparkes, A., W. Aubrey, E. Byrne, A. Clare, M. N. Khan, M. Liakata, and R. D. King. 2010. Towards Robot Scientists for autonomous scientific discovery. *Autom. Exp.* 2(1). Available at: <http://www.aejournal.net/content/2/1/1> (accessed February 10, 2015).
54. Castrillo, J. I., and S. G. Oliver, eds. 2011. *Yeast systems biology: Methods and protocols*. Methods in molecular biology ser. Berlin—Heidelberg: Springer. Vol. 759. 549 p.
55. Plotkin, G. D. 1970. A note on inductive generalization. *Mach. Intell.* 5:153–163.
56. Huang, J., L. Antova, C. Koch, and D. Olteanu. 2009. MayBMS: A probabilistic database management system. *2009 ACM SIGMOD Conference (International) on Management of Data Proceedings*. Rhode Island. 1071–1074.
57. Robin, A., and M. Crézé. 1986. Stellar populations in the Milky Way — a synthetic model. *Astron. Astrophys.* 157:71–90.
58. Robin, A. C., C. Reylé C., S. Derrière, and S. Picaud. 2006. A synthetic view on structure and evolution of the Milky Way. arXiv preprint astro-ph/0401052.

59. Czekaj, M. A., A. C. Robin, F. Figueras, X. Luri, and M. Haywood. 2014. The Besançon Galaxy model renewed-I. Constraints on the local star formation history from Tycho data. *Astron. Astrophys.* 564:A102.
60. Czekaj, M. A. 2012. Galaxy evolution: A new version of the Besançon Galaxy Model constrained with Tycho data. PhD Thesis. Barcelona: Universitet de Barcelona. 167 p.
61. Martins, A. M. M. 2014. Statistical analysis of large scale surveys for constraining the Galaxy evolution. PhD Thesis. Barcelona: Universitet de Barcelona. 221 p.
62. Biswal, B. B., M. Mennes, X. N. Zuo, S. Gohel, C. Kelly, S. M. Smith, and C. Windischberger. 2010. Toward discovery science of human brain function. *Proc. Nat. Acad. Sci. USA* 107(10):4734–4739.
63. Craddock, R. C., S. Jbabdi, C. G. Yan, J. T. Vogelstein, F. X. Castellanos, A. Di Martino, and M. P. Milham. 2013. Imaging human connectomes at the macroscale. *Nat. Methods* 10(6):524–539.
64. Ginestet, C. E., P. Balachandran, S. Rosenberg, and E. D. Kolaczyk. 2014. Hypothesis testing for network data in functional neuroimaging. arXiv preprint arXiv:1407.5525.
65. Ginestet, C. E., A. P. Fournel, and A. Simmons. 2014. Statistical network analysis for functional MRI: Summary networks and group comparisons. *Front. Comput. Neurosci.* 8:51. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4018548/> (accessed February 10, 2015).
66. Yan, C. G., R. C. Craddock, X. N. Zuo, Y. F. Zang, and M. P. Milham. 2013. Standardizing the intrinsic bra towards robust measurement of inter-individual variation in 1000 functional connectomes. *Neuroimage* 80:246–262.
67. Marcus, D. S., J. Harwel, T. Olsen, M. Hodge, M. F. Glasser, F. Prior, and D. C. Van Essen. 2011. Informatics and data mining tools and strategies for the human connectome project. *Front. Neuroinform.* 5. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3127103/> (accessed February 10, 2015).
68. Marcus, D. S., T. R. Olsen, M. Ramaratnam, and R. L. Buckner. 2007. The extensible neuroimaging archive toolkit. *Neuroinformatics* 5(1):11–33.
69. Brun, A. 2006. Manifold learning and representations for image analysis and visualization. Department of Biomedical Engineering, Linköpings Universitet. 104 p.
70. Mahmoudi, A., S. Takerkart, F. Rezagui, D. Bous-saoud, and A. Brovelli. 2012. Multivoxel pattern analysis for fMRI data: A review. *Comput. Math. Methods Med.* Available at: <http://www.hindawi.com/journals/cmmm/2012/961257/> (accessed February 10, 2015).
71. Van Horn, J. D., and A. W. Toga. 2014. Human neuroimaging as a “Big Data” science. *Brain Imaging Behavior* 8(2):323–331.
72. Hillebrandt, H., K. J. Friston, and S. J. Blakemore. 2014. Effective connectivity during animacy perception-dynamic causal modelling of Human Connectome Project data. *Sci. Rep.* 4. Available at: <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4150124/> (accessed February 10, 2016).
73. Lappalainen, J., M. A. Sicilia, and B. Hernández. 2013. Automatic hypothesis checking using eScience Research Infrastructures, ontologies, and linked data: A case study in climate change research. *Procedia Comput. Sci.* 18:1172–1178.
74. Lenten, L. J., and I. A. Moosa. 2003. An empirical investigation into long-term climate change in Australia. *Environ. Modell. Softw.* 18(1):59–70.
75. Borges, M. R. 2010. Efficient market hypothesis in European stock markets. *Eur. J. Financ.* 16(7):711–726.
76. Bollen, J., H. Mao, and X. Zeng. 2011. Twitter mood predicts the stock market. *J. Comput. Sci.* 2(1):1–8.
77. Spangler, S., A. D. Wilkins, B. J. Bachman, et al. 2014. Automated hypothesis generation based on mining scientific literature. *KDD'14 Proceedings*. New York. 1877–1886.
78. Zhou, D., O. Bousquet, T. N. Lal, J. Weston, and B. Schölkopf. 2004. Learning with local and global consistency. *Adv. Neur. Inform. Proc. Syst.* 16(16):321–328.

Contributors

Received February 10, 2015

Kalinichenko Leonid A. (b. 1937) — Doctor of Science in physics and mathematics, professor; Head of Laboratory, Institute of Informatics Problems, Russian Academy of Sciences; 44-2 Vavilov Str., Moscow 119333, Russian Federation; professor, Faculty of Computational Mathematics and Cybernetics, M. V. Lomonosov Moscow State University, 1-52 Leninskiye Gory, GSP-1, Moscow 119991, Russian Federation; leonidandk@gmail.com

Kovalev Dmitry Yu. (b. 1988) — junior scientist, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; dkovalev@ipiran.ru

Kovaleva Dana A. (b. 1973) — Candidate of Science (PhD) in physics and mathematics, scientist, Institute of Astronomy, Russian Academy of Sciences, 48 Pyatnitskaya Str., Moscow 119017, Russian Federation; dana@inasan.ru

Malkov Oleg Yu. (b. 1961) — Doctor of Science in physics and mathematics, associate professor; Head of Department, Institute of Astronomy, Russian Academy of Sciences; 48 Pyatnitskaya Str., Moscow 119017, Russian Federation; professor, Faculty of Physics, M. V. Lomonosov Moscow State University, 1-52 Leninskiye Gory, GSP-1, Moscow 119991, Russian Federation; malkov@inasan.ru

МЕТОДЫ И СРЕДСТВА ПОДДЕРЖКИ ИССЛЕДОВАНИЙ, ДВИЖИМЫХ ГИПОТЕЗАМИ: ОБЗОР

Л. А. Калиниченко¹, Д. Ю. Ковалев², Д. А. Ковалева³, О. Ю. Малков⁴

¹Институт проблем информатики Российской академии наук; leonidandk@gmail.com

²Институт проблем информатики Российской академии наук; dkovalev@ipiran.ru

³Институт астрономии Российской академии наук; dana@inasan.ru

⁴Институт астрономии Российской академии наук; malkov@inasan.ru

Аннотация: Исследования с интенсивным использованием данных (ИИИД), развиваемые в рамках новой парадигмы изучения естественных явлений, именуемой Четвертой парадигмой, придают особое значение все возрастающей роли, которую играют данные, полученные в результате наблюдений, экспериментов или компьютерного моделирования, практически во всех областях анализа и накопления информации. Главной целью ИИИД является извлечение (вывод) знаний из данных. Целью настоящей работы является обзор существующих подходов, методов и инфраструктур анализа данных в ИИИД с акцентом на роли гипотез в процессе анализа информации и эффективной поддержки формирования, оценки и выбора гипотез при моделировании естественных явлений и проведении экспериментов. Статья включает введение в разнообразные понятия, методы и средства эффективной организации движимых гипотезами экспериментов в ИИИД.

Ключевые слова: исследования с интенсивным использованием данных; Четвертая парадигма; гипотезы; модели; теории; гипотетико-дедуктивный метод; проверка гипотез; решетка гипотез; модель Галактики, анализ коннектома; автоматизированная генерация гипотез

DOI: 10.14357/19922264150104

Литература

1. The Fourth paradigm: Data-intensive scientific discovery / Eds. T. Hey, S. Tansley, K. Tolle. — Redmond, Microsoft Research, 2009. 252 p.
2. *McComas W. F.* The principal elements of the nature of science: Dispelling the myths of science // Nature of science in science education: Rationales and strategies / Ed. W. F. McComas. — Kluwer Academic Publs., 1998. P. 53–70.
3. *Lakshmana Rao J. R.* Scientific ‘Laws’, ‘Hypotheses’ and ‘Theories’ // Meanings Distinctions Reson., 1998. Vol. 3. P. 69–74.
4. *Poincaré H.* The foundations of science: Science and hypothesis, the value of science, science and method. The Project Gutenberg EBook, 2012. No. 39713. P. 554. <http://www.gutenberg.org/files/39713/39713-8.txt>.
5. *Bacon F.* The new organon // Great books of the Western World. Vol. 30. The works of Francis Bacon / Ed. R. M. Hutchins. — Chicago, Encyclopedia Britannica, Inc., 1952. P. 107–195.
6. *Menzies T.* Applications of abduction: Knowledge-level modeling // Int. J. Hum.-Comput. St., 1996. Vol. 45. No. 3. P. 305–335.
7. *Haber J.* Research questions, hypotheses, and clinical questions // Evolve resources for nursing research. — 7th ed. — Elsevier, 2010. P. 27–55.
8. *Popper K.* The logic of scientific discovery. — London—New York: Routledge, Taylor & Francis, 2005. 545 p. <http://strangebeautiful.com/other-texts/popper-logic-scientific-discovery.pdf>.
9. *Kerlinger F. N., Lee H. B.* Foundations of behavioral research: Educational and psychological inquiry. — New York: Holt, Rinehart and Winston, 1964. 739 p.
10. *Hempel C. G.* Fundamentals of concept formation in empirical science // Int. Encyclopedia Unified Sci., 1952. Vol. 2. No. 7. <http://www.iep.utm.edu/hempel/>.
11. *Porto F., Spaccapietra S.* Data model for scientific models and hypotheses // Evolution Conceptual Modeling, 2011, Vol. 6520. P. 285–305.
12. *Gonçalves B., Porto F.* A lattice-theoretic approach for representing and managing hypothesis-driven research // AMW, 2013.
13. *Gonçalves B., Porto F., Moura A. M. C.* On the semantic engineering of scientific hypotheses as linked data // 2nd Workshop (International) on Linked Science Proceedings, 2012.
14. *Woodward J.* Scientific explanation // The Stanford Encyclopedia of Philosophy, 2011. <http://plato.stanford.edu/archives/win2011/entries/scientific-explanation/>.
15. Scientific discovery: Case studies / Ed. T. Nickles. — Taylor & Francis, 1980. Vol. 2. 501 p.
16. *Schickore J.* Scientific discovery // The Stanford Encyclopedia of Philosophy, 2014. <http://plato.stanford.edu/archives/spr2014/entries/scientific-discovery/>.
17. *Kakas A. C., Kowalski R. A., Toni F.* Abductive logic programming // J. Logic Comput., 1993. Vol. 2. No. 6. P. 719–770.

18. *Kakas A. C., Michael A., Mourlas C.* ACLP: Abductive constraint logic programming // *J. Logic Program.*, 2000. Vol. 44. No. 1. P. 129–177.
19. *Van Nuffelen B., Kakas A.* A-system: Declarative programming with abduction // *Logic programming and nonmonotonic reasoning* / Eds. T. Eiter, W. Faber, M. Truszczynski. — Lecture notes in computer science ser. — Berlin–Heidelberg: Springer, 2001. Vol. 2173. P. 393–397.
20. *Alferes J. J., Pereira L. M., Swift T.* Abduction in well-founded semantics and generalized stable models via tabled dual programs // *Theor. Pract. Log. Prog.*, 2004. Vol. 4. No. 4. P. 383–428.
21. *Ray O., Kakas A.* ProLogICA: A practical system for Abductive Logic Programming // 11th Workshop (International) on Non-Monotonic Reasoning Proceedings, 2006. P. 304–312.
22. *Citrigno S., Eiter T., Faber W., Gottlob G., Koch C., Leone N., Scarcello F.* The dlv system: Model generator and application frontends // 12th Workshop on Logic Programming Proceedings, 1997. P. 128–137.
23. *King R. D., Liakata M., Lu C., Oliver S. G., Soldatova L. N.* On the formalization and reuse of scientific research // *J. Roy. Soc. Interface*, 2011. Vol. 8. No. 63. P. 1440–1448.
24. *Tamaddoni-Nezhad A., Chaleil R., Kakas A., Muggleton S. H.* Application of abductive ILP to learning metabolic network inhibition from temporal data // *Mach. Learn.*, 2006. Vol. 64. P. 209–230.
25. *Inoue K., Sato T., Ishihata M., Kameya Y., Nabeshima H.* Evaluating abductive hypotheses using an EM algorithm on BDDs // *IJCAI-09 Proceedings*, 2009. P. 810–815.
26. *Bartha P.* Analogy and analogical reasoning // *The Stanford Encyclopedia of Philosophy*, 2013. <http://plato.stanford.edu/archives/fall2013/entries/reasoning-analogy/>.
27. *Ivezic Ž., Connolly A. J., VanderPlas J. T., Gray A.* Statistics, data mining, and machine learning in astronomy: A practical Python guide for the analysis of survey data. — Princeton University Press, 2014. 552 p.
28. *Sivia D. S., Skilling J.* Data analysis. A Bayesian tutorial. — New York: Oxford University Press Inc., 2006. 264 p.
29. *Field A.* Discovering statistics using IBM SPSS statistics. — 4th ed. — Sage, 2013. 915 p.
30. *IBM SPSS Statistics for Windows, Version 22.0.* Armonk, N.Y.: IBM Corp. IBM SPSS Statistics base, 2013. https://www.uio.no/tjenester/it/forskning/statistikk/hjelp/programveilednigner/ibm_spss_statistics_brief_guide-2.pdf.
31. *Ihaka R., Gentleman R. R.* R: A language for data analysis and graphics // *J. Comput. Graph. Stat.*, 1996. Vol. 5. No. 3. P. 299–314.
32. *March M. C., Starkman G. D., Trotta R., Vaudrevange P. M.* Should we doubt the cosmological constant? // *Mon. Not. Roy. Astron. Soc.*, 2011. Vol. 410. No. 4. P. 2488–2496.
33. *Rouder J. N., Speckman P. L., Sun D., Morey R. D., Iverson G.* Bayesian t tests for accepting and rejecting the null hypothesis // *Psychon. Bull. Rev.*, 2009. Vol. 16. No. 2. P. 225–237.
34. *Weber M.* Experiment in biology // *The Stanford Encyclopedia of Philosophy*, 2014. <http://plato.stanford.edu/archives/fall2014/entries/biology-experiment/>.
35. *Hawthorne J.* Inductive logic // *The Stanford Encyclopedia of Philosophy*, 2014. <http://plato.stanford.edu/archives/sum2014/entries/logic-inductive/>.
36. *Breiman L.* Statistical modeling: The two cultures // *Stat. Sci.*, 2001. Vol. 16. No. 3. P. 199–231.
37. *Hastie T., Tibshirani R., Friedman J., Franklin J.* The elements of statistical learning: Data mining, inference and prediction // *Math. Intell.*, 2005. Vol. 27. No. 2. P. 83–85.
38. *Barber D.* Bayesian reasoning and machine learning. — Cambridge University Press, 2010. 720 p.
39. *Ferrucci D., Brown E., Chu-Carroll J., Fan J., Gondek D., Kalyanpur A. A., Welty C.* Building Watson: An overview of the DeepQA project // *AI Mag.*, 2010. Vol. 31. No. 3. P. 59–79.
40. *Dredze M., Crammer K., Pereira F.* Confidence-weighted linear classification // 25th Conference (International) on Machine Learning Proceedings. Helsinki, Finland, 2008. P. 264–271.
41. *Starkman G. D., Trotta R., Vaudrevange P. M.* Introducing doubt in Bayesian model comparison. arXiv preprint arXiv:0811.2415, 2008.
42. *March M. C.* Advanced statistical methods for astrophysical probes of cosmology. Springer Theses, 2013. Vol. 20. 177 p.
43. *Porto F.* Big data in astronomy. The LINEA-DEXL case // Presentation at the EMC Summer School on BIG DATA — NCE/UFRJ, 2013. <http://www.slideshare.net/fabiomporto/emc-2013-big-data-in-astronomy>.
44. *Racunas S. A., Shah N. H., Albert I., Fedoroff N. V.* Hybrow: A prototype system for computer-aided hypothesis evaluation // *Bioinformatics*, 2004. Vol. 20. No. 1. P. 257–264.
45. *Soldatova L. N., Rzhetsky A., King R. D.* Representation of research hypotheses // *J. Biomed. Semantics*, 2011. Vol. 2. No. S-2. P. S9.
46. *Callahan A., Duumontier M., Shah N.* HyQue: Evaluating hypotheses using Semantic Web technologies // *J. Biomed. Semantics*, 2011. Vol. 2. No. S-2. P. S3.
47. *Gao Y., Kinoshita J., Wu E., Miller E., Lee R., Seaborne A., Clark T.* SWAN: A distributed knowledge infrastructure for Alzheimer disease research // *J. Web Semant.*, 2006. Vol. 4. No. 3. P. 222–228.
48. *King R. D., Whelan K. E., Jones F. M., Reiser P. G., Bryant C. H., Muggleton S. H., Oliver S. G.* Functional genomic hypothesis generation and experimentation by a robot scientist // *Nature*, 2004. Vol. 427. No. 6971. P. 247–252.
49. *Porto F., Moura A. M. C., Gonçalves B., Costa R., Spaccapetra S. A.* A scientific hypothesis conceptual model // *Advances in conceptual modeling* / Eds. S. Castano, P. Vassiliadis, L. V. Lakshmanan, M. Li Lee. — Lecture notes in computer science ser. — Berlin–Heidelberg: Springer, 2012. Vol. 7518. P. 101–110.

50. Porto F., Moura A. M. C. Scientific hypothesis database. Report, 2011. <http://livroaberto.ibict.br/bitstream/1/869/1/Scientific%20Hypothesis%20Database.pdf>.
51. Asgharbeygi N., Langley P., Bay S., Arrigo K. Inductive revision of quantitative process models // *Ecol. Model.*, 2006. Vol. 194. No. 1. P. 70–79.
52. Tran N., Baral C., Nagaraj V. J., Joshi L. Knowledge-based integrative framework for hypothesis formation in biochemical networks // *Data integration in the life sciences* / Eds. B. Ludäscher, L. Raschid. — Lecture notes in computer science ser. Berlin–Heidelberg: Springer, 2005. Vol. 3615. P. 121–136.
53. Sparkes A., Aubrey W., Byrne E., Clare A., Khan M. N., Liakata M., King R. D. Towards Robot Scientists for autonomous scientific discovery // *Autom. Exp.*, 2010. Vol. 2. No. 1. <http://www.aejournal.net/content/2/1/1>.
54. Yeast systems biology: Methods and protocols / Eds. J. I. Castrillo, S. G. Oliver. — *Methods in molecular biology* ser. — Berlin–Heidelberg: Springer, 2011. Vol. 759. 549 p.
55. Plotkin G. D. A note on inductive generalization // *Mach. Intell.*, 1970. Vol. 5. P. 153–163.
56. Huang J., Antova L., Koch C., Olteanu D. MayBMS: A probabilistic database management system // 2009 ACM SIGMOD Conference (International) on Management of Data Proceedings, 2009. P. 1071–1074.
57. Robin A., Crézé M. Stellar populations in the Milky Way — a synthetic model // *Astron. Astrophys.*, 1986. Vol. 157. P. 71–90.
58. Robin A. C., Reylé C., Derrière S., Picaud S. A synthetic view on structure and evolution of the Milky Way. arXiv preprint astro-ph/0401052, 2004.
59. Czekaj M. A., Robin A. C., Figueras F., Luri X., Haywood M. The Besançon Galaxy model renewed-I. Constraints on the local star formation history from Tycho data // *Astron. Astrophys.*, 1986. Vol. 564. P. A102.
60. Czekaj M. A. Galaxy evolution: A new version of the Besançon Galaxy Model constrained with Tycho data. PhD Thesis, 2012. Universitet de Barcelona, Spain. 167 p.
61. Martins A. M. M. Statistical analysis of large scale surveys for constraining the Galaxy evolution. PhD Thesis, 2014. Universitet de Barcelona, Spain. 221 p.
62. Biswal B. B., Mennes M., Zuo X. N., Gohel S., Kelly C., Smith S. M., Windischberger C. Toward discovery science of human brain function // *Proc. Nat. Acad. Sci. USA*, 2010. Vol. 107. No. 10. P. 4734–4739.
63. Craddock R. C., Jbabdi S., Yan C. G., Vogelstein J. T., Castellanos F. X., Di Martino A., Milham M. P. Imaging human connectomes at the macroscale // *Nat. Methods*, 2013. Vol. 10. No. 6. P. 524–539.
64. Ginestet C. E., Balachandran P., Rosenberg S., Kolarczyk E. D. Hypothesis testing for network data in functional neuroimaging. arXiv preprint arXiv:1407.5525, 2014.
65. Ginestet C. E., Fournel A. P., Simmons A. Statistical network analysis for functional MRI: Summary networks and group comparisons // *Front. Comput. Neurosci.*, 2014. Vol. 8. P. 51. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4018548/>.
66. Yan C. G., Craddock R. C., Zuo X. N., Zang Y. F., Milham M. P. Standardizing the intrinsic bra towards robust measurement of inter-individual variation in 1000 functional connectomes // *Neuroimage*, 2013. Vol. 80. P. 246–262.
67. Marcus D. S., Harwell J., Olsen T., Hodge M., Glasser M. F., Prior F., Van Essen D. C. Informatics and data mining tools and strategies for the human connectome project // *Front. Neuroinform.*, 2011. Vol. 5. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3127103/>.
68. Marcus D. S., Olsen T. R., Ramaratnam M., Buckner R. L. The extensible neuroimaging archive toolkit // *Neuroinformatics*, 2007. Vol. 5. No. 1. P. 11–33.
69. Brun A. Manifold learning and representations for image analysis and visualization. Department of Biomedical Engineering, Linköpings Universitet, 2006. 104 p.
70. Mahmoudi A., Takerkart S., Regragui F., Boussaoud D., Brovelli A. Multivoxel pattern analysis for fMRI data: A review // *Comput. Math. Methods Med.*, 2012. <http://www.hindawi.com/journals/cmdd/2012/961257/>.
71. Van Horn J. D., Toga A. W. Human neuroimaging as a “Big Data” science // *Brain Imaging Behavior*, 2014. Vol. 8. No. 2. P. 323–331.
72. Hillebrandt H., Friston K. J., Blakemore S. J. Effective connectivity during animacy perception-dynamic causal modelling of Human Connectome Project data // *Sci. Rep.*, 2014. Vol. 4. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4150124/>.
73. Lappalainen J., Sicilia M. Á., Hernández B. Automatic hypothesis checking using eScience Research Infrastructures, ontologies, and linked data: A case study in climate change research // *Procedia Comput. Sci.*, 2013. Vol. 18. P. 1172–1178.
74. Lenten L. J., Moosa I. A. An empirical investigation into long-term climate change in Australia // *Environ. Modell. Softw.*, 2003. Vol. 18. No. 1. P. 59–70.
75. Borges M. R. Efficient market hypothesis in European stock markets // *Eur. J. Financ.*, 2010. Vol. 16. No. 7. P. 711–726.
76. Bollen J., Mao H., Zeng X. Twitter mood predicts the stock market // *J. Comput. Sci.*, 2011. Vol. 2. No. 1. P. 1–8.
77. Spangler S., Wilkins A. D., Bachman B. J., et al. Automated hypothesis generation based on mining scientific literature // *KDD’14 Proceedings*, 2014. P. 1877–1886.
78. Zhou D., Bousquet O., Lal T. N., Weston J., Schölkopf B. Learning with local and global consistency // *Adv. Neur. Inform. Proc. Syst.*, 2004. Vol. 16. No. 16. P. 321–328.

Поступила в редакцию 10.02.2015

ФОРМАЛЬНЫЙ АКСИОМАТИЧЕСКИЙ ПОДХОД К АСПЕКТНО-ОРИЕНТИРОВАННОМУ РАСШИРЕНИЮ ТЕХНОЛОГИЙ ПРОГРАММИРОВАНИЯ*

С. П. Ковалёв¹

Аннотация: Исследуется процедура расширения технологий модульной разработки программных систем приемами аспектно-ориентированного подхода. Расширение описано как обогащение формальных моделей программных модулей разметкой их интерфейсов классами задач, образующими аспектную структуру. Предложен новый подход к разделению ответственности (separation of concerns) путем естественной модуляризации аспектной структуры. В качестве обобщения этого подхода предложена процедура частичной модуляризации аспектной структуры. Для формализации образующихся конструкций на общесистемном уровне, не зависящем от частных парадигм программирования, привлекается теория категорий. Технологиям разработки программ отвечают категории, объектами которых служат формальные модели программ, а морфизмами — технологические операции. Аспектно-ориентированное расширение (АО-расширение) технологии описано аксиоматически как преобразование таких категорий — функтор, обладающий сопряженными подходящего вида как справа, так и слева. В качестве иллюстративного примера АО-расширения приводится событийный подход к моделированию систем.

Ключевые слова: аспектно-ориентированное программирование; трассируемость; теория категорий; формальная технология проектирования; разделение ответственности

DOI: 10.14357/19922264150105

1 Введение

Современный комплексный подход к созданию программных систем требует автоматизировать не только основные процессы предметной деятельности, но и управляющие и обеспечивающие (а также обеспечивающие для обеспечивающих и т.д.) (см., например, [1]). В сложных предметных областях они обладают значительным разнообразием, глубоко погружаются в контекст основной деятельности почти на каждом шаге и в то же время с трудом совмещаются с нею на понятийном уровне.

Многие задачи обеспечивающих процессов рассеиваются по системе, не поддаваясь локализации в рамках программных модулей, автоматизирующих шаги основных. Использование традиционных «модульных» технологий программирования для таких задач приводит к значительным затратам труда, поскольку приходится дублировать одни и те же алгоритмы в контексте различных шагов. Примеры можно найти как среди функциональных задач (ведение информационной модели объекта управления, верификация данных), так и среди программно-технических (защита ин-

формации, ведение журналов функционирования системы и др.).

Для повышения эффективности программной реализации таких рассеянных задач в конце 1990-х гг. была предложена новая парадигма — аспектно-ориентированное программирование (АОП) [2]. Рассеянные задачи оформляются в виде аспектов — особых программных единиц, код которых составляется однократно и затем автоматически вставляется в код основных единиц в точках, явно задаваемых внешним образом. В результате вставки аспект получает полный доступ к контексту. Однако на практике АОП применяется значительно реже, чем модульные подходы, поскольку на концептуальном уровне неясно, как оптимально выделять и соединять аспекты и какие аспекты целесообразно реализовывать модулями [3]. Существующие технологии АОП, такие как AspectJ (АО-расширение языка Java) [4], предлагают лишь частные решения, специфичные для частных парадигм программирования.

В связи с этим целью настоящей работы ставится построение универсальной концептуальной модели расширения «модульных» технологий средствами разработки аспектов — основного метода внедре-

*Статья подготовлена при поддержке Российского гуманитарного научного фонда (грант 13-03-00384).

¹Институт проблем управления им. В. А. Трапезникова Российской академии наук, kovalyov@nm.ru

ния АОП в практику программирования. Расширение описано как обогащение модулей разметкой классами задач, описывающей «историю» их разработки, на уровне интеграционных интерфейсов. Такой подход обусловлен тем, что разметка позволяет непосредственно проводить трассирование задач — операцию, которая фактически является обратной по отношению к рассеиванию и вследствие этого лежит в основе АОП [5]. Классы задач образуют аспектную структуру модулей, так что аспектно-ориентированная (де)композиция проводится согласованно на двух уровнях: модульных основ и аспектных структур.

Чтобы строго сформулировать и верифицировать этот подход, в работе привлекается теория категорий, поскольку ее средствами можно дать универсальное аксиоматическое описание процессов создания программных систем с позиций общей теории систем [6]. Результативность аксиоматического подхода здесь связана с тем, что можно потребовать соблюдения аксиом при организации (и тем более при автоматизации) труда программистов, что обеспечивает применимость и эффективность результатов, выведенных из аксиом. Отправной точкой служит теоретико-категорная конструкция формальной технологии проектирования (architecture school) [7]. Добавление поддержки аспектов в формальную технологию описано в работе как ее преобразование, порождающее категории помеченных системных единиц. Эти категории оснащены функторами выделения модульной основы и аспектной структуры. Приемы АОП формализуются универсальными конструкциями в таких категориях, и их свойства строго доказываются путем вывода из аксиом.

В качестве источника примеров для иллюстрации предлагаемого подхода выбрано событийное моделирование программных систем, поскольку в его понятиях формулируются классические частные семантические модели АОП (см., например, [8]). Вообще в литературе предлагается много подходов к формализации АОП (см., например, [9]), однако все они представлены в контексте тех или иных частных формализмов теоретического программирования (лямбда-исчисление, проверка на моделях и др.) и поэтому могут применяться только в рамках частных парадигм.

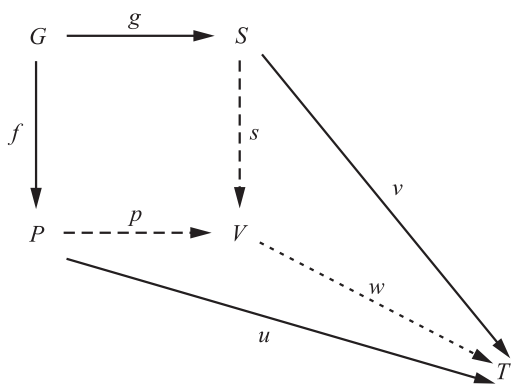
Работа построена следующим образом. В разд. 2 вводится понятие формальной технологии. Раздел 3 посвящен преобразованию модульных технологий в аспектно-ориентированные. В разд. 4 и 5 описаны конструкции полной и частичной модуляризации аспектной структуры соответственно. В заключении подводятся итоги исследования.

2 Теоретико-категорное описание разработки программ

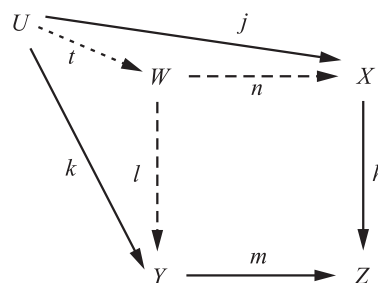
Категория C — это класс абстрактных объектов $\text{Ob } C$, попарно связанных морфизмами (абстрактными аналогами отображений) [10, разд. 1.1]: каждый морфизм f имеет область $\text{dom } f \in \text{Ob } C$ и кообласть $\text{codom } f \in \text{Ob } C$. Соотношения вида $\text{dom } f = A$ и $\text{codom } f = B$ наглядно записываются в форме стрелки $f : A \rightarrow B$, а множество всех морфизмов, удовлетворяющих этим соотношениям, обозначается через $\text{Mor}(A, B)$. Для любой пары морфизмов f, g такой, что $\text{codom } f = \text{dom } g$, определена композиция-морфизм $g \circ f : \text{dom } f \rightarrow \text{codom } g$. Композиция ассоциативна: для любой тройки морфизмов f, g, h если $\text{codom } f = \text{dom } g$ и $\text{codom } g = \text{dom } h$, то $h \circ (g \circ f) = (h \circ g) \circ f$. Наконец, любой объект A обладает тождественным морфизмом $1_A : A \rightarrow A$ таким, что для любого морфизма $f : A \rightarrow B$ выполняется соотношение $f \circ 1_A = 1_B \circ f = f$.

Для формального аксиоматического описания разработки программных систем категории применяются начиная с 1970-х гг. (см., например, [6]). Здесь объекты отвечают компонентам и системам — обычно это формальные модели программ (алгебраические спецификации, графы, термы лямбда-исчисления и т. п.). Морфизмы часто обозначают действия по интеграции компонентов в системы. Композиция морфизмов отвечает конструированию многошаговых действий (процессов), а тождественные морфизмы — «ничегонеделанию». Будем обозначать категорию такого рода через $c\text{-DESC}$. Конфигурации взаимосвязанных компонентов, из которых собираются системы, задаются $c\text{-DESC}$ -диаграммами — ориентированными графами, вершины которых помечены объектами, а ребра — морфизмами категории $c\text{-DESC}$. Актам сборки систем отвечают копределы диаграмм [10, разд. 3.3]. Поясним конструкцию копредела на нескольких примерах. В качестве первого рассмотрим *соединение* компонента P с системой S — прием сборки, состоящий в добавлении промежуточного компонента G , называемого «клеем» (glue), или связкой (connector) [7], который способен интегрироваться как с компонентом, так и с системой. Путем соединения часто строятся системы на базе промежуточного программного обеспечения (middleware). Конфигурация соединения имеет вид пары $c\text{-DESC}$ -морфизмов $f : P \leftarrow G \rightarrow S : g$. Ее копредел, называемый *кодекартовым квадратом*, задается объектом-вершиной V и парой морфизмов-ребер $p : P \rightarrow V \leftarrow S : s$ таких, что $p \circ f = s \circ g$ и выполняется следующее условие универ-

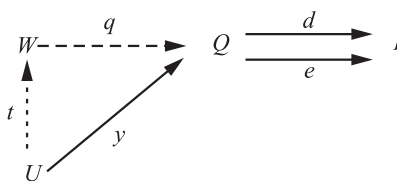
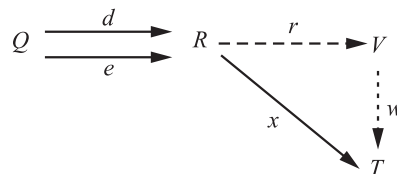
сальности: для любых объекта T и пары морфизмов $u : P \rightarrow T \leftarrow S : v$ если $u \circ f = v \circ g$, то существует единственный морфизм $w : V \rightarrow T$, удовлетворяющий соотношениям $w \circ p = u$ и $w \circ s = v$. Тогда объект V действительно отвечает системе, которая собрана из S и P путем соединения посредством G (и не содержит ничего «лишнего»). Если такой объект V существует, то он определяется однозначно с точностью до изоморфизма — формального представления несущественного различия между моделями. (Если же копредела не существует, то делается вывод, что клей G не способен соединить компонент P с S посредством действий f и g .)



Конструкция копредела, как видно из ее названия, является двойственной по отношению к конструкции предела, которая была введена первоначально для нужд приложений теории категорий в топологии [10, разд. 3.4]. В приложениях к разработке программных систем пределы привлекаются для формализации процедур разложения систем (декомпозиции), обратных по отношению к сборке. Например, в разд. 4 потребуется выделять из моделей части, соответствующие прообразам относительно действия определенных морфизмов. Формально для объекта Y , вложенного в Z посредством мономорфизма (категорного аналога инъекции) $m : Y \hookrightarrow Z$, полный прообраз относительно некоторого морфизма $h : X \rightarrow Z$ строится при помощи *декартова квадрата*. Это конструкция, двойственная к кодекартову квадрату, т.е. предел диаграммы $m : Y \hookrightarrow Z \leftarrow X : h$. Он задается объектом W и парой морфизмов $n : X \leftarrow W \rightarrow Y : l$ таких, что $h \circ n = m \circ l$ и выполняется следующее условие универсальности: для любых объекта U и пары морфизмов $j : X \leftarrow U \rightarrow Y : k$ если $h \circ j = m \circ k$, то существует единственный морфизм $t : U \rightarrow W$, удовлетворяющий соотношениям $n \circ t = j$ и $l \circ t = k$. Искомый прообраз объекта Y выделяется в X морфизмом n , причем он также является мономорфизмом.



Другим важным частным случаем (ко)пределов являются регулярные морфизмы. Рассмотрим произвольную c -DESC-диаграмму вида $d, e : Q \rightrightarrows R$, состоящую из двух параллельных морфизмов. Ее копредел, если он существует, задается морфизмом $r : R \rightarrow V$ таким, что $r \circ d = r \circ e$ и если соотношение $x \circ d = x \circ e$ выполнено для некоторого морфизма $x : R \rightarrow T$, то существует единственный морфизм $w : V \rightarrow T$, удовлетворяющий соотношению $w \circ r = x$. Морфизм r называется *коуравнителем* пары d, e , он является эпиморфизмом (категорным аналогом сюръекции) и задает некоторую факторизацию объекта R (например, если в качестве c -DESC взять категорию множеств, то морфизм r факторизует R по отношению эквивалентности, порожденному множеством пар $\{(d(z), e(z)) | z \in Q\}$). Двойственно, *уравнителем* пары d, e называется морфизм $q : W \rightarrow Q$ такой, что $d \circ q = e \circ q$, и если соотношение $d \circ y = e \circ y$ выполнено для некоторого морфизма $y : U \rightarrow Q$, то существует единственный морфизм $t : U \rightarrow W$, удовлетворяющий соотношению $q \circ t = y$. Уравнитель является мономорфизмом и содержательно задает вложение W в Q в качестве подобъекта (в категории множеств $W \cong \{z | d(z) = e(z)\} \subseteq Q$). Если некоторый морфизм выступает в качестве коуравнителя некоторой пары, то он называется *регулярным эпиморфизмом* [11, определение 7.71], а если в качестве уравнителя — *регулярным мономорфизмом* [11, определение 7.56].



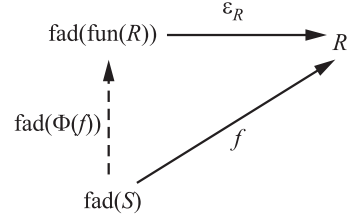
Чтобы описать механизмы формирования конфигураций, формализуется понятие интеграцион-

ного интерфейса — «части» модели, задающей правила интеграции других моделей с нею [7]. Например, у веб-сервиса интерфейсом служит его декларация на языке WSDL (Web Service Description Language). Формальные модели интерфейсов образуют категорию, обозначаемую через SIG, а операция выделения интерфейса у модели программы формализуется как функтор $\text{sig} : c\text{-DESC} \rightarrow \text{SIG}$, называемый сигнатурным. *Функтор* — это отображение категорий, переводящее объекты в объекты, а морфизмы в морфизмы, с сохранением композиции и тождественных морфизмов [10, разд. 1.3]. Поэтому функторами описываются преобразования моделей программ, совместимые с интеграцией систем. Поскольку различные модели могут иметь один и тот же интерфейс, функтор sig не обязан быть инъективным на объектах. Однако sig -образы двух различных действий по интеграции одного и того же компонента в одну и ту же систему должны быть различными: иначе получится, что интерфейсы недостаточно детально описывают интеграционные возможности компонентов. Иными словами, функтор sig должен быть *универсальным* (faithful) [11, определение 3.27(2)], т.е. инъективным на каждом множестве $\text{Mor}(P, S)$, $P, S \in \text{Ob } c\text{-DESC}$.

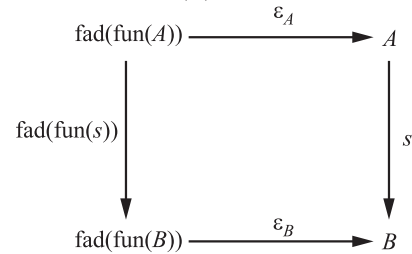
Кроме того, для каждого интерфейса должна существовать хотя бы одна реализация, поддерживающая его интеграционные возможности в полном объеме. Например, WSDL-описание любого веб-сервиса можно реализовать «заглушками» (stubs) — пустыми процедурами: они могут автоматически генерироваться CASE-средствами и позволяют быстро собирать отладочные версии приложений. Формально должен существовать функтор дискретной реализации $\text{sig}^* : \text{SIG} \rightarrow c\text{-DESC}$ такой, что $\text{sig} \circ \text{sig}^* = 1_{\text{SIG}}$ и для любых SIG-объекта I и $c\text{-DESC}$ -объекта S функтор sig сюръективно (следовательно, биективно) отображает множество $\text{Mor}(\text{sig}^*(I), S)$, описывающее все действия по интеграции дискретной реализации интерфейса I в систему S , на множество $\text{Mor}(I, \text{sig}(S))$, определяющее интеграционные возможности интерфейса I . Это означает, что функтор sig^* сопряжен слева к sig , причем единица этого сопряжения тождественна.

Напомним, что для произвольных категорий C , D и функторов $\text{fun} : C \rightarrow D$, $\text{fad} : D \rightarrow C$ сопряжение $\text{fad} \dashv \text{fun}$ — это семейство биекций $\Phi : \text{Mor}(\text{fad}(S), R) \cong \text{Mor}(S, \text{fun}(R))$, $S \in \text{Ob } D$, $R \in \text{Ob } C$, естественное в следующем смысле [10, разд. 4.1]: для любых C -морфизмов $f : \text{fad}(S) \rightarrow R$, $k : R \rightarrow Y$ и D -морфизма $h : X \rightarrow S$ выполняется соотношение $\Phi(k \circ f \circ \text{fad}(h)) = \text{fun}(k) \circ \Phi(f) \circ h$. Семейство C -морфизмов

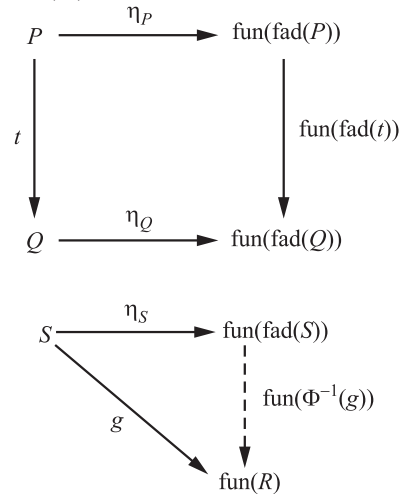
$\varepsilon_A = \Phi^{-1}(1_{\text{fun}(A)}) : \text{fad}(\text{fun}(A)) \rightarrow A$, $A \in \text{Ob } C$, называется *коединицей* сопряжения: имеем $\Phi(1_R \circ \varepsilon_R \circ \text{fad}(\Phi(f))) = \text{fun}(1_R) \circ \Phi(\varepsilon_R) \circ \Phi(f) = \Phi(f)$, откуда получается «треугольное тождество» сопряжения $\varepsilon_R \circ \text{fad}(\Phi(f)) = f$. Это тождество позволяет вычислить действие сопряжения, если известна его коединица.



Естественность сопряжения проявляется в том, что коединица представляет собой естественное преобразование функтора $\text{fad} \circ \text{fun}$ в 1_C , т.е. для любого C -морфизма $s : A \rightarrow B$ выполняется соотношение $\varepsilon_B \circ \text{fad}(\text{fun}(s)) = s \circ \varepsilon_A$. Это соотношение является частным случаем треугольного тождества для морфизма $s \circ \varepsilon_A : \text{fad}(\text{fun}(A)) \rightarrow B$, поскольку $\Phi(s \circ \varepsilon_A) = \Phi(s \circ \varepsilon_A \circ \text{fad}(1_{\text{fun}(A)})) = \text{fun}(s) \circ \Phi(\varepsilon_A) \circ 1_{\text{fun}(A)} = \text{fun}(s)$.



Двойственно, семейство D -морфизмов $\eta_G = \Phi(1_{\text{fad}(G)}) : G \rightarrow \text{fun}(\text{fad}(G))$, $G \in \text{Ob } D$, образует *единицу* сопряжения — естественное преобразование функтора 1_D в $\text{fun} \circ \text{fad}$ ($\eta_Q \circ t = \text{fun}(\text{fad}(t)) \circ \eta_P$ для любого D -морфизма $t : P \rightarrow Q$), порождающее второе треугольное тождество $\text{fun}(\Phi^{-1}(g)) \circ \eta_S = g$ для любого D -морфизма $g : S \rightarrow \text{fun}(R)$.



Если единица состоит из тождественных морфизмов, то $\text{fun} \circ \text{fad} = 1_D$ ввиду естественности, а второе треугольное тождество приобретает вид $\text{fun}(\Phi^{-1}(g)) = g$, откуда $\Phi(f) = \text{fun}(f)$ для любого C -морфизма $f : \text{fad}(S) \rightarrow R$, т.е. биекция сопряжения действует так же, как правый сопряженный функтор (это имеет место для функтора выделения интерфейсов sig). Двойственно, если коединица тождественна, то $\text{fad} \circ \text{fun} = 1_C$ и биекция Φ^{-1} действует как левый сопряженный функтор.

Полноценный процесс создания программных систем в дополнение к интеграции включает трансформации (refinements) — шаги разработки индивидуальных компонентов (уточнение требований, реализация спецификации на языке программирования и др.). Трансформации моделей программ могут быть устроены совершенно иначе, чем действия по интеграции. Поэтому они описываются морфизмами подходящей категории, обозначаемой через r -DESC, которая в общем случае отличается от c -DESC, но обладает таким же классом объектов [7]. Накладываются условия естественности выделения интерфейса и трансформации относительно сборки систем. Получается *формальная технология проектирования* [7] — сложная категорная конструкция, состоящая из категорий c -DESC и r -DESC вместе с классом c -DESC-диаграмм, представляющих конфигурации систем, и функтором выделения интерфейсов $\text{sig} : c\text{-DESC} \rightarrow \text{SIG}$.

Например, в событийном подходе к проектированию [12] в роли основной категории моделей c -DESC выступает категория **Pos** всех частично упорядоченных множеств и всех их монотонных отображений. **Pos**-объекты отвечают сценариям поведения программ — совокупностям событий, частично упорядоченным причинно-следственными связями. Трансформацией сценария X в Y , т.е. морфизмом категории, выступающей в роли r -DESC, служит «раскрытие» событий до подсценариев — произвольное антифункциональное тотальное отношение $R \subseteq X \times Y$, удовлетворяющее условию $\forall x, x' \in X \forall y, y' \in Y (xRy \wedge x'Ry' \wedge x \neq x') \Rightarrow (x \leq x' \Leftrightarrow y \leq y')$. В качестве функтора sig , извлекающего интерфейсы из сценариев, выступает функтор $|-| : \mathbf{Pos} \rightarrow \mathbf{Set} : S \mapsto |S|$, «забывающий» порядок. Левым сопряженным к нему служит функтор дискретного упорядочения $\text{id} : \mathbf{Set} \rightarrow \mathbf{Pos} : I \mapsto \langle I, = \rangle$. Примечательно, что функтор id в свою очередь имеет левый сопряженный — таковым служит функтор распараллеливания $\text{sconn} : \mathbf{Pos} \rightarrow \mathbf{Set}$, сопоставляющий каждому частично упорядоченному множеству множество всех компонент связности его порядка (взаимно независимых подсценариев). Коединица этого сопряжения тождественна, поэтому

для любых сценария S и множества I имеется биекция $\text{sconn} : \text{Mor}(S, \text{id}(I)) \cong \text{Mor}(\text{sconn}(S), I)$, т.е. интерфейсы сценариев в полной мере задают их поведение в роли как компонентов, так и систем. Такое свойство представляет интерес и для произвольных технологий.

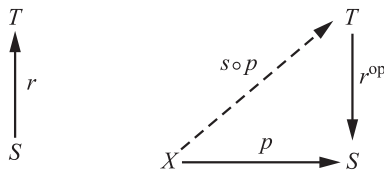
Определение 1. Формальная технология проектирования называется *структурируемой*, если функтор дискретной реализации интерфейсов sig^* имеет левый сопряженный с тождественной коединицей. \square

3 Формальные технологии аспектно-ориентированного проектирования

Среди показателей качества программных систем от рассеивания задач больше всех страдает трассируемость — возможность точно определить, для решения каких задач в систему включен тот или иной фрагмент [13]. В снижении затрат на трассирование фактически состоит назначение АОП: главным мотивом его создателей было отсутствие в традиционных языках программирования конструкций, позволяющих разделять исходный код программ по классам задач [2]. Поэтому, как отмечалось во введении, аспектно-ориентированный подход в целом рассматривается как оснащение модулей разметкой, идентифицирующей классы решаемых ими задач. Универсальная формальная модель АОП основывается на аксиоматическом описании трассирования, которое строится из конструкций в формальных технологиях проектирования следующим образом.

Хорошо известно, что трассируемость легко нарушается при трансформациях (хрестоматийным примером служит реализация алгебраической спецификации программы на алгоритмическом языке программирования). При интеграции, напротив, обычно обеспечивается хотя бы частичное трассирование (здесь примером служит прямая сумма индексированного семейства множеств в категории **Set**, при построении которой элементы каждой компоненты снабжаются ее индексом — «меткой»). Ввиду этого для произвольной трансформации — r -DESC-морфизма $r : S \rightarrow T$ — возникает следующее необходимое условие возможности трассировать вдоль нее результат к источнику [12]: обращение ее направления, т.е. теоретико-категорная дуализация должна превращать эту трансформацию в действие по интеграции — в c -DESC-морфизм $r^{\text{op}} : T \rightarrow S$.

Поскольку в ходе разработки программных систем трансформация перемежается со сборкой систем, необходимо совместно трассировать трансформации и действия по интеграции. Проще всего провести такое совместное трассирование, если трасса $r^{\text{оп}}$ обратима справа [12]. Действительно, существование c -DESC-морфизма $s : S \rightarrow T$ такого, что $r^{\text{оп}} \circ s = 1_S$, эквивалентно тому, что для любого c -DESC-морфизма $p : X \rightarrow S$, задающего интеграцию некоторого компонента X в систему S , существует действие по интеграции X в T , совместимое с трассированием трансформации r в том смысле, что композиция трассы $r^{\text{оп}}$ с этим действием дает p . Таким действием служит $s \circ p$, поскольку $r^{\text{оп}} \circ (s \circ p) = p$.



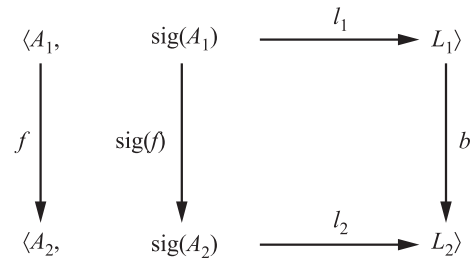
Заметим, что s является регулярным c -DESC-мономорфизмом [11, предложение 7.59(1)], т.е. задает вложение источника трансформации в результат в качестве подобъекта. На практике его построение может быть весьма трудоемким, но он требуется не всегда, поскольку трассированию вдоль процессов сборки систем подлежат в первую очередь интеграционные требования, которые предъявляются к интерфейсам моделей. В таком случае достаточно потребовать обратимости справа не для трассы $r^{\text{оп}}$, а лишь для SIG-морфизма $\text{sig}(r^{\text{оп}})$, представляющего действие трассы на уровне интерфейсов и называемого разметкой [12]. Реализация обращения разметки обычно не требует значительных затрат, поскольку интерфейсы проектируются так, чтобы интегрировать их было «проще», чем сами модели. Получается частный случай известного подхода к снижению затрат на трассирование путем ограничения класса трассируемых требований согласно их значимости (value-based requirements traceability) [14].

Определение 2. r -DESC-морфизм r называется *трассируемой трансформацией*, а двойственный к нему морфизм $r^{\text{оп}}$ — *трассой*, если $r^{\text{оп}}$ принадлежит категории c -DESC и SIG-морфизм $\text{sig}(r^{\text{оп}})$ является ретракцией (т.е. имеет правый обратный). sig -образ трассы называется *разметкой*. \square

Например, в технологии событийного моделирования все трансформации трассируемы, причем класс всех разметок состоит из всех сюръективных отображений множеств (в силу аксиомы выбора любой **Set**-эпиморфизм является ретракцией). Действие трансформации сценариев $t : X \rightarrow Y$,

двойственной к сюръекции $t^{\text{оп}} : |Y| \rightarrow |X|$, интуитивно можно трактовать как реализацию классов задач — точек множества $|X|$ — путем раскрытия (expansion) в их прообразы относительно $t^{\text{оп}}$. По определению трансформация сценариев разбивает свой результат на подмножества, хорошо упорядоченные (well-ordered) в том смысле, что для любых $x, y, z, u \in Y$ таких, что $t^{\text{оп}}(x) = t^{\text{оп}}(y) \neq t^{\text{оп}}(z) = t^{\text{оп}}(u)$, условие $x \leq z$ влечет $y \leq u$.

Наиболее прямым и экономным способом обеспечения полной трассируемости является «запоминание» трасс трансформаций вместе с моделями программ, порожденными ими из классов задач [15]. В контексте АОП интерес представляет в первую очередь влияние трансформаций на интеграционные возможности моделей, поэтому достаточно присоединить к моделям действия трансформаций на уровне интерфейсов, т.е. разметки. Интеграция и трансформация таких обогащенных моделей должна согласованно выполняться на двух уровнях: модульных основ и аспектных структур. Как указано в [16, разд. 7], в теории категорий имеется специальная конструкция, предназначенная для естественного присоединения действий к объектам — категория запятой (comma category) [10, разд. 2.6]. Рассмотрим категорию запятой $\text{sig} \downarrow \text{SIG}$. Ее объектами являются все пары вида $\langle A, l : \text{sig}(A) \rightarrow L \rangle$, где A — c -DESC-объект, l — SIG-морфизм. Морфизмом объекта $\langle A_1, l_1 : \text{sig}(A_1) \rightarrow L_1 \rangle$ в $\langle A_2, l_2 : \text{sig}(A_2) \rightarrow L_2 \rangle$ является любая пара $\langle f : A_1 \rightarrow A_2, b : L_1 \rightarrow L_2 \rangle$ такая, что $b \circ l_1 = l_2 \circ \text{sig}(f)$.



Определение 3. *Аспектно-ориентированной моделью* (АО-моделью) называется любой $(\text{sig} \downarrow \text{SIG})$ -объект $\langle A, l : \text{sig}(A) \rightarrow L \rangle$ такой, что l является разметкой. c -DESC-объект A называется (*модульной*) *основой* АО-модели, SIG-морфизм l — ее (*аспектной*) *разметкой*, SIG-объект L — ее *аспектной структурой*. \square

Будем обозначать через АО полную подкатегорию в $\text{sig} \downarrow \text{SIG}$, класс объектов которой состоит из всех АО-моделей. (Напомним, что полной подкатегорией в произвольной категории C называется категория, состоящая из некоторого класса C -объектов X и объединения всех множеств $\text{Mor}(A, B)$,

$A, B \in X$, оснащенная унаследованными из C операциями). Поясним эту конструкцию на примере событийного моделирования. У АО-модели сценария разметка l — это сюръективное отображение, сопоставляющее каждому элементу множества A (событию) элемент множества L , который можно рассматривать как обозначение класса задач, порождающего это событие [12]. В частности, аспектом естественно называть любую АО-модель, в которой L состоит из одного элемента (такие и только такие АО-модели сценариев удовлетворяют приведенному ниже формальному определению 4) [17]. Поэтому произвольная АО-модель сценария — это частично-упорядоченное мультимножество (multiset), состоящее из аспектов. Подобный подход к моделированию сценариев был предложен еще в 1980-х гг. [18], однако природа меток и способы их синтеза оставались неясными, поскольку они не рассматривались в контексте АОП.

Категория АО снабжена следующими «забывающими» функторами, индуцированными конструкцией категории запятой [12]:

- $\text{mod} : \text{АО} \rightarrow c\text{-DESC} : \langle A, l \rangle \mapsto A, \langle f, b \rangle \mapsto f$ (выделение модульной основы);
- $\text{int} = \text{sig} \circ \text{mod} : \text{АО} \rightarrow \text{SIG} : \langle A, l \rangle \mapsto \text{sig}(A), \langle f, b \rangle \mapsto \text{sig}(f)$ (выделение исходного интерфейса);
- $\text{str} : \text{АО} \rightarrow \text{SIG} : \langle A, l \rangle \mapsto \text{codom } l, \langle f, b \rangle \mapsto b$ (выделение аспектной структуры).

Функтор mod примечателен тем, что он позволяет извлекать из АО-модели модульную основу в форме особого интеграционного интерфейса, т. е. порождает формальную технологию проектирования, поддерживающую аспектно-ориентированный подход. Действительно, функтор mod унивалентен: если $\langle f, b \rangle, \langle f', b' \rangle : \langle A, l \rangle \rightarrow \langle B, k \rangle$ — два произвольных АО-морфизма, то $b' \circ l = k \circ \text{sig}(f) = b \circ l$, откуда $b' = b$, поскольку l обратим справа. Левым сопряженным к функтору mod служит функтор дискретной разметки $\text{mod}^* : c\text{-DESC} \rightarrow \text{АО} : A \mapsto \langle A, 1_{\text{sig}(A)} \rangle, f \mapsto \langle f, \text{sig}(f) \rangle$, единица этого сопряжения тождественна. Естественным образом конструируются трансформации и конфигурации АО-моделей [12]. В результате получается АО-технология — формальная технология проектирования, в которой основной категорией моделей служит АО, а интерфейсы выделяются функтором mod . Важное прикладное значение имеет критерий структурируемости этой технологии, который будет сформулирован и доказан в разд. 5.

Аспектом (aspect) называется элементарный строительный блок аспектно-ориентированной программы (модели), реализующий отдельный класс задач. Аспекты способны сохранять свою

идентичность в составе программы, поэтому их аспектная структура не разрушается при интеграции [17]. Хорошо сохраняют структуру действия по интеграции, обладающие на уровне аспектных структур обратимостью *слева*, т. е. возможностью идентифицировать аспектную структуру компонента в составе системы путем трассирования (см. пояснения перед определением 2). Поэтому АО-морфизм $\langle f, b \rangle$ называется *аспектным*, если SIG-морфизм b обратим слева, и, в частности, *изоаспектным*, если b является изоморфизмом. Например, в АО-технологии событийного моделирования аспектными являются все отображения помеченных сценариев с непустой областью, которые не «склеивают» различные метки.

Определение 4. Аспектно-ориентированная модель A называется *аспектом*, если любой АО-морфизм с областью A является аспектным. \square

Сборка программ из аспектов выходит за рамки традиционной компоновки модулей (linking), поэтому она называется связыванием (weaving). Технологии АОП предлагают разнообразные инструменты связывания: предкомпиляторы исходных текстов программ, преобразователи исполняемого байт-кода, диспетчеры вызова методов. Связывание можно аксиоматически описать универсальной конструкцией в категории АО — копределом диаграммы соединения специального вида [17].

4 Модуляризация аспектов

Наиболее явным способом разметки аспектов, составляющих АО-модель, является их модуляризация, т. е. оформление в виде отдельных единиц модульной архитектуры: объектов, таблиц в базе данных и т. д. Такая модуляризация называется разделением ответственности (separation of concerns) и входит в число основных целей привлечения АОП [19]. Модуляризованные аспекты можно собирать в системы посредством компоновки, не прибегая к связыванию, что позволяет снизить затраты на сборку за счет применения широкодоступных «модульных» технологий программирования. В связи с этим при использовании расширений традиционных технологий средствами АОП, таких как AspectJ, программы, скомпонованные из модулей без применения связывания, тривиальным образом разделяются на модуляризируемые аспекты [4].

Модуляризируемые АО-модели выделяются среди прочих тем, что при интеграции с модулями они ведут себя так же, как модульные единицы. Возможности интеграции модулей в АО-модель определяются ее модульной основой. В свою

очередь, возможности интеграции АО-модели в модули определяются ее аспектной структурой, поскольку при интеграции в модуль аспекты, составляющие АО-модель, выступают в качестве ее элементарных единиц. Формально, модуляризуемые модели образуют полную подкатегорию в АО (которая будет обозначаться через m -АО) такую, что существует АО-расширение модульной технологии — вложение $\text{am} : c\text{-DESC} \mapsto m\text{-АО}$, полностью воспроизводящее интеграционные возможности модулей в следующем смысле. С одной стороны, все способы интеграции модуля $M \in \text{Ob } c\text{-DESC}$ в АО-модель $A \in \text{Ob } m\text{-АО}$ задаются множеством морфизмов $\text{Mor}(\text{am}(M), A)$, поэтому функтор выделения модульного интерфейса mod (точнее, его ограничение на m -АО, обозначаемое далее через am_*) должен устанавливать биекцию между множеством $\text{Mor}(\text{am}(M), A)$ и множеством $\text{Mor}(M, \text{mod}(A))$. С другой стороны, все способы интеграции A в M задаются множеством морфизмов $\text{Mor}(A, \text{am}(M))$, поэтому должен существовать функтор модуляризации аспектной структуры $\text{am}^* : m\text{-АО} \rightarrow c\text{-DESC}$, тривиально действующий на модули ($\text{am}^* \circ \text{am} = 1_{c\text{-DESC}}$) и устанавливающий биекцию между $\text{Mor}(A, \text{am}(M))$ и множеством $\text{Mor}(\text{am}^*(A), M)$. При этом возможность трактовать $c\text{-DESC}$ -объект $\text{am}^*(A)$ как «подъем» аспектной структуры АО-модели A на модульный уровень обеспечивается следующим дополнительным требованием естественности: функтор $\text{sig} \circ \text{am}^*$, выделяющий интеграционный интерфейс из модуляризированной аспектной структуры, должен совпадать с ограничением на m -АО функтора str , определяющего исходную аспектную структуру на уровне интеграционных интерфейсов.

Примером расширения, которым обладает любая формальная технология проектирования, служит изоморфизм между $c\text{-DESC}$ и полной подкатегорией в АО с классом объектов $\{\langle A, 1_{\text{sig}(A)} \rangle \mid A \in \text{Ob } c\text{-DESC}\}$, действующий как функтор mod^* . Например, в технологии событийного моделирования он порождает дискретно размеченные сценарии, которые являются самыми «аспектно-неориентированными» — в них каждый класс задач помечает только одно событие, так что никакого рассеивания не происходит. Назовем такое АО-расширение тривиальным. Можно проверить, что действие любого АО-расширения (с точностью до естественного изоморфизма), по существу, совпадает с его действием: модули (т. е. $c\text{-DESC}$ -объекты) всегда переходят в АО-модели, история получения интерфейсов которых из классов задач путем трансформации утрачена (тривиальна). Таким образом, АО-расширение, по существу, однозначно определяется своей кообластью — классом всех модуля-

ризуемых АО-моделей. Легко видеть, что конструкция АО-расширения устойчива относительно ограничения: любая полная подкатегория в m -АО, содержащая класс $\text{am}(\text{Ob } c\text{-DESC})$, является кообластью АО-расширения, действующего так же, как am . В связи с этим интерес представляют расширения, кообласть которых достаточно «велика».

На языке теории категорий требования к АО-расширению компактно формулируются при помощи конструкции сопряжения функторов.

Определение 5. Функтор $\text{am} : c\text{-DESC} \rightarrow m\text{-АО}$, где $m\text{-АО}$ — некоторая полная подкатегория в АО, называется *АО-расширением* формальной технологии, если он обладает следующими сопряженными функторами:

- правый сопряженный am_* с тождественной единицей, причем $\text{am}_*(f) = \text{mod}(f)$ для любого m -АО-морфизма f ;
- левый сопряженный am^* с тождественной коединицей, причем $\text{sig}(\text{am}^*(f)) = \text{str}(f)$ для любого m -АО-морфизма f .

Аспектно-ориентированное расширение am называется:

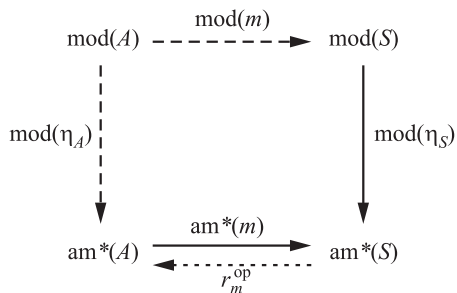
- *тривиальным*, если оно является изоморфизмом;
- *наибольшим*, если любое АО-расширение $\text{am}' : c\text{-DESC} \rightarrow m\text{-АО}'$ удовлетворяет условию $\text{Ob } m\text{-АО}' \subseteq \text{Ob } m\text{-АО}$;
- *полным*, если для любой АО-модели существует m -АО-объект, изоморфный ей. \square

При аксиоматическом описании модуляризации аспектов, составляющих АО-модели, ключевую роль играет единица сопряжения $\text{am}^* \dashv \text{am}$, которая будет обозначаться через η . По определению для любого m -АО-объекта S m -АО-морфизм $\eta_S : S \rightarrow \text{am}(\text{am}^*(S))$ является прообразом $c\text{-DESC}$ -морфизма $1_{\text{am}^*(S)}$ при биекции $\text{am}^* : \text{Mor}(S, \text{am}(\text{am}^*(S))) \cong \text{Mor}(\text{am}^*(S), \text{am}^*(S))$, поэтому $\text{str}(\eta_S) = \text{sig}(\text{am}^*(\eta_S)) = 1_{\text{sig}(\text{am}^*(S))} = 1_{\text{str}(S)}$ ввиду определения 5, т. е. морфизм η_S изоаспектен и его действие нетривиально только на уровне модульной основы. А поскольку $c\text{-DESC}$ -объект $\text{am}^*(S)$ представляет на модульном уровне аспектную структуру АО-модели S , $c\text{-DESC}$ -морфизм $\text{mod}(\eta_S) : \text{mod}(S) \rightarrow \text{am}^*(S)$ можно рассматривать как канонический способ интеграции модульной основы модели в ее модуляризованную аспектную структуру. Этот морфизм является регулярным $c\text{-DESC}$ -эпиморфизмом (см. следствие 1.2 ниже), т. е. задает факторизацию модульной основы на

аспекты. Семейство $\text{mod}(\eta_S)$, $S \in \text{Ob } m\text{-АО}$, образует естественное преобразование функтора выделения модульной основы am_* в функтор выделения аспектной структуры am^* .

Определение 6. (am -)модуляризуемой АО-моделью называется любой m -АО-объект. (am -)модуляризацией (аспектной структуры) модуляризуемой АО-модели S называется c -DESC-морфизм $\text{mod}(\eta_S)$, где η — единица сопряжения $\text{am}^* \dashv \text{am}$. Тождеством (am -)модуляризации произвольного m -АО-морфизма $m : A \rightarrow S$ называется равенство $\text{am}^*(m) \circ \text{mod}(\eta_A) = \text{mod}(\eta_S) \circ \text{mod}(m)$. \square

Естественным способом извлечения модуляризованного аспекта из модульной основы модуляризуемой АО-модели S является трассирование метки аспекта (обозначения классов задач) вдоль модуляризации аспектной структуры модели. Если A — извлекаемый аспект, то должна существовать трассируемая модульная трансформация $r_m : \text{am}^*(A) \rightarrow \text{am}^*(S)$ его (модуляризированной) метки в (модуляризованную) аспектную структуру модели S . Идентификация метки в аспектной структуре производится путем прямого трассирования — морфизма, обратного справа к трассе r_m^{op} (см. пояснения перед определением 2). Таким морфизмом должен служить $\text{am}^*(m)$, если m -АО-морфизм m задает вложение A в S . Поэтому $\text{mod}(m)$ должен выделять в $\text{mod}(S)$ полный прообраз c -DESC-объекта $\text{am}^*(A)$ как подобъекта в $\text{am}^*(S)$ относительно $\text{mod}(\eta_S)$. Следовательно, тождество модуляризации для m -АО-морфизма $m : A \rightarrow S$ должно задавать декартов квадрат. Таким образом, подаспекты — это вложения аспектов, модуляризация которых имеет универсальный характер. Благодаря этому, в частности, подаспекты являются подобъектами в категории АО. Кроме того, аспекты являются атомарными единицами разделения ответственности — они не имеют собственных подаспектов (см. предложение 2 ниже).



Определение 7. m -АО-морфизм $m : A \rightarrow S$ называется (am -)подаспектом АО-модели S , если он удовлетворяет следующим условиям:

- (i) АО-модель A является аспектом;
- (ii) существует трассируемая трансформация $r_m : \text{am}^*(A) \rightarrow \text{am}^*(S)$ такая, что $r_m^{\text{op}} \circ \text{am}^*(m) = 1_{\text{am}^*(A)}$;
- (iii) тождество am -модуляризации для m задает декартов квадрат в категории c -DESC. \square

Предложение 1. Любой подаспект является аспектным регулярным АО-мономорфизмом.

Доказательство. Покажем, что если $m : A \rightarrow S$ — подаспект, то тождество единицы $\text{am}(\text{am}^*(m)) \circ \eta_A = \eta_S \circ m$ задает декартов квадрат в категории АО. Действительно, функтор mod переводит его в тождество модуляризации, которое является декартовым квадратом по условию (iii) определения 7. Функтор int переводит его в тождество $\text{str}(m) \circ \text{int}(\eta_A) = \text{int}(\eta_S) \circ \text{int}(m)$, которое представляет собой sig -образ тождества модуляризации, так что задает декартов квадрат в SIG (ввиду наличия у функтора sig левого сопряженного, он сохраняет все пределы [11, предложение 18.9]). Наконец, функтор str переводит тождество единицы в коммутативный квадрат в SIG, две параллельные стороны которого ($\text{str}(\eta_A)$ и $\text{str}(\eta_S)$) являются тождественными морфизмами, он также декартов по определению. Используя эти факты, можно непосредственно проверить, что тождество единицы задает декартов квадрат. Следовательно, m — уравниватель пары АО-морфизмов η_S , $(\text{am}(\text{am}^*(m)) \circ r_m^{\text{op}}) \circ \eta_S : S \rightrightarrows \text{am}(\text{am}^*(S))$. \square

Предложение 2. Следующие утверждения эквивалентны для любого подаспекта $m : A \rightarrow S$:

- (i) S является модуляризуемым аспектом;
- (ii) m является изоморфизмом;
- (iii) m изоаспектен.

Доказательство.

(i) \Rightarrow (iii). Если S — аспект, то АО-морфизм $t \circ \eta_S : S \rightarrow \text{am}(\text{am}^*(A))$, где $t = \text{am}(r_m^{\text{op}})$, является аспектным. Поскольку η_S изоаспектен, $\text{str}(t)$ обратим слева; следовательно, АО-морфизм t , будучи в свою очередь обратимым справа, изоаспектен. Поэтому SIG-морфизм $w = \text{str}(\text{am}(\text{am}^*(m)))$, правый обратный к $\text{str}(t)$, является изоморфизмом. В свою очередь, применяя функтор str к тождеству единицы $\text{am}(\text{am}^*(m)) \circ \eta_A = \eta_S \circ m$, получаем, что $w = \text{str}(m)$.

(iii) \Rightarrow (ii). Как отмечалось в доказательстве предложения 1, в SIG имеется декартов квадрат, задаваемый тождеством $\text{str}(m) \circ \text{int}(\eta_A) = \text{int}(\eta_S) \circ \text{int}(m)$. По предположению $\text{str}(m)$ — изоморфизм, поэтому $\text{int}(m)$, будучи параллельным ему ребром декартова квадрата, также является

изоморфизмом, в частности эпиморфизмом. А поскольку функтор int , будучи композицией унивалентных функторов mod и sig , сам унивалентен, он отражает эпиморфизмы [11, предложение 7.44]. В свою очередь, если m — эпиморфизм, то в силу предложения 1 он является изоморфизмом.

(ii) \Rightarrow (i). Вытекает из определения 4. \square

Естественным (хотя и не единственным) способом модуляризации аспектной структуры АО-модели является восстановление трассируемой трансформации модулей, индуцирующей ее разметку. Действительно, если для АО-модели $\langle A, l \rangle$ существует подходящая трассируемая трансформация $s : X \rightarrow A$ некоторого c -DESC-объекта X в A , удовлетворяющая условию $\text{sig}(s^{\text{op}}) = l$, то X можно рассматривать как модульную единицу, состоящую из всех классов задач АО-модели, а s^{op} — как «каноническую» модуляризацию ее аспектной структуры. Например, у помеченного сценария трансформация восстанавливается (очевидным образом) из такой и только такой разметки, которая разбивает свою область на хорошо упорядоченную совокупность прообразов точек, причем прообраз любой точки является подаспектом. Очевидными примерами служат любой аспект (единственным подаспектом которого ввиду предложения 2 является он сам) и любой непустой дискретно размеченный сценарий (подаспектом в котором служит любое одноэлементное подмножество). Показательным примером является также сценарий, в котором меткой каждого события выступает содержащая его компонента связности порядка: множество всех подаспектов такого сценария совпадает с множеством взаимно независимых подсценариев, так что с формальной точки зрения распараллеливание (см. конструкцию функтора scopn в конце разд. 2) оказывается частным случаем модуляризации аспектов.

5 Частичная модуляризация

Для формальных технологий, не обладающих полным АО-расширением, можно предложить более слабые подходы к модуляризации аспектной структуры. В частности, аспектную разметку произвольной АО-модели можно представить на модульном уровне частичным морфизмом модульной основы. Напомним, что *частичным морфизмом* объекта X в Y называется пара морфизмов с общим началом, один из которых является мономорфизмом и направлен в X (он выделяет «часть» объекта X , выступающую областью определения частичного морфизма), а другой произволен и направлен в Y (он задает действие частичного морфизма) [11,

определение 28.1(1)]. При некоторых технических ограничениях определена композиция частичных морфизмов. Частичный морфизм можно рассматривать как диаграмму со схемой $* \leftarrow * \rightarrow *$, обозначаемой через V .

Разметка произвольной АО-модели порождает частичный c -DESC-морфизм следующим образом. Воспользуемся тем, что коединица ε сопряжения $\text{sig}^* \dashv \text{sig}$ состоит из мономорфизмов, поскольку $\text{sig}(\varepsilon_A) = 1_{\text{sig}(A)}$, $A \in \text{Ob } c\text{-DESC}$, и любой унивалентный функтор отражает мономорфизмы [11, предложение 7.37(2)]. При помощи коединицы можно выделить в модульной основе АО-модели дискретную «часть» и далее подействовать на нее дискретной реализацией аспектной разметки.

Определение 8. *Частичной модуляризацией* произвольной АО-модели $\langle A, l : \text{sig}(A) \rightarrow L \rangle$ называется следующий частичный c -DESC-морфизм, действующий из ее модульной основы A в дискретную реализацию аспектной структуры L :

$$\varepsilon_A : A \leftarrow \text{sig}^*(\text{sig}(A)) \rightarrow \text{sig}^*(L) : \text{sig}^*(l). \quad \square$$

Предложение 3. Отображение, сопоставляющее каждой АО-модели ее частичную модуляризацию, задает полное вложение категории АО в категорию диаграмм $c\text{-DESC}^V$.

Доказательство. Выберем произвольно АО-модели $S = \langle A, l : \text{sig}(A) \rightarrow L \rangle$ и $S' = \langle A', l' : \text{sig}(A') \rightarrow L' \rangle$, положим $A^* = \text{sig}^*(\text{sig}(A))$, $A'^* = \text{sig}^*(\text{sig}(A'))$. Если $\varepsilon_A = \varepsilon_{A'}$ и $\text{sig}^*(l) = \text{sig}^*(l')$, то $A = \text{codom } \varepsilon_A = \text{codom } \varepsilon_{A'} = A'$ и $l = \text{sig}(\text{sig}^*(l)) = \text{sig}(\text{sig}^*(l')) = l'$. Далее, любой АО-морфизм $\langle f, b \rangle : S \rightarrow S'$ определяет семейство c -DESC-морфизмов $f : A \rightarrow A'$, $f^* : A^* \rightarrow A'^*$, $\text{sig}^*(b) : \text{sig}^*(L) \rightarrow \text{sig}^*(L')$ (здесь $f^* = \text{sig}^*(\text{sig}(f))$), которое является естественным преобразованием диаграмм частичной модуляризации, т. е. $c\text{-DESC}^V$ -морфизмом: по определению коединицы имеем $f \circ \varepsilon_A = \varepsilon_{A'} \circ f^*$, а по определению АО-морфизма — $\text{sig}^*(b) \circ \text{sig}^*(l) = \text{sig}^*(l') \circ f^*$. Ясно, что различные АО-морфизмы определяют различные естественные преобразования такого вида.

$$\begin{array}{ccccc} A & \xleftarrow{\varepsilon_A} & A^* & \xrightarrow{\text{sig}^*(l)} & \text{sig}^*(L) \\ \downarrow f & & \downarrow f^* & & \downarrow \text{sig}^*(b) \\ A' & \xleftarrow{\varepsilon_{A'}} & A'^* & \xrightarrow{\text{sig}^*(l')} & \text{sig}^*(L') \end{array}$$

Таким образом, задано вложение $\text{rmao} : \text{АО} \hookrightarrow c\text{-DESC}^V$. Чтобы убедиться в его полноте,

рассмотрим произвольное естественное преобразование диаграммы $\text{pmao}(S)$ в $\text{pmao}(S')$ — тройку c -DESC-морфизмов $p : A \rightarrow A'$, $q : A^* \rightarrow A'^*$, $s : \text{sig}^*(L) \rightarrow \text{sig}^*(L')$, удовлетворяющих соотношениям $p \circ \varepsilon_A = \varepsilon_{A'} \circ q$ и $s \circ \text{sig}^*(l) = \text{sig}^*(l') \circ q$. Поскольку ε состоит из мономорфизмов, из первого равенства вытекает, что $q = \text{sig}^*(\text{sig}(p))$. С учетом этого, применяя функтор sig ко второму равенству, получаем, что $\text{sig}(s) \circ l = l' \circ \text{sig}(p)$. Следовательно, существует АО-морфизм $\langle p, \text{sig}(s) \rangle : S \rightarrow S'$. \square

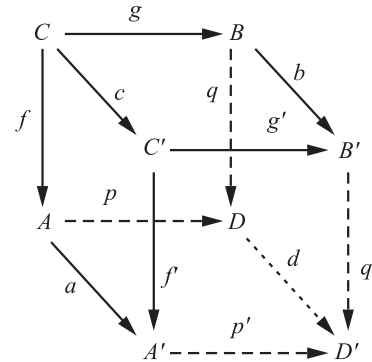
Конечно, частичная модуляризация АО-модели не отражает трансформацию, порождающую модульную основу АО-модели из аспектной структуры, поскольку действие частичной модуляризации не обязано быть трассой трансформации. Тем не менее, можно построить «аппроксимацию» модуляризации аспектной структуры, вычисляя копредел диаграммы частичной модуляризации (если он существует: ребро копредела частичного морфизма, параллельное его действию, можно рассматривать как универсальное расширение частичного морфизма на всю свою область). По существу, речь идет о соединении (в смысле разд. 2) модульной основы с дискретной реализацией аспектной структуры АО-модели. Оказывается, таким способом можно вычислить модуляризацию аспектной структуры любой модуляризуемой АО-модели. При некотором техническом ограничении отсюда выводится существование наибольшего АО-расширения. Оно модуляризирует в точности все АО-модели, аспектную разметку которых можно «поднять» на модульный уровень в виде регулярного эпиморфизма.

Таким образом, переход к обобщенной частичной процедуре модуляризации аспектной структуры позволяет по-новому взглянуть на ее исходную полную форму и установить ряд ее важных свойств. Отправным пунктом служит следующий факт: существование копредела у любой диаграммы частичной модуляризации является критерием структурируемости АО-технологии — необходимого условия существования полного АО-расширения (см. следствие 1.3 ниже). Например, этому критерию удовлетворяет технология событийного моделирования, поскольку любая **Pos**-диаграмма имеет копредел.

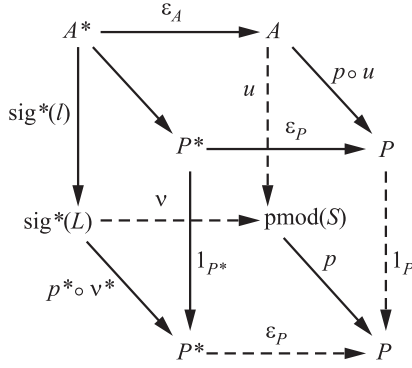
Теорема 1. *Аспектно-ориентированная технология структурируема тогда и только тогда, когда диаграмма частичной модуляризации любой АО-модели имеет копредел.*

Доказательство. Рассмотрим произвольные АО-диаграммы $f : A \leftarrow C \rightarrow B : g$ и $f' : A' \leftarrow C' \rightarrow B' : g'$, обладающие кодекартовыми квадратами $p \circ f = q \circ g$ и $p' \circ f' = q' \circ g'$ с вершинами D

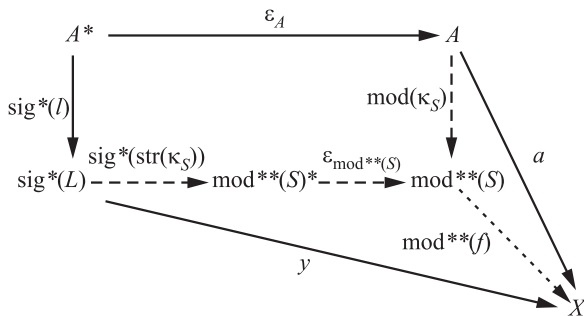
и D' соответственно. Если АО-морфизмы $a : A \rightarrow A'$, $c : C \rightarrow C'$, $b : B \rightarrow B'$ образуют естественное преобразование этих диаграмм, то существует единственный c -DESC-морфизм $d : D \rightarrow D'$, удовлетворяющий условиям $d \circ p = p' \circ a$ и $d \circ q = q' \circ b$ (т.е. делающий куб в c -DESC, составленный из двух параллельных кодекартовых квадратов и морфизмов a, c, b, d , коммутативным). Обозначим морфизм d через $\text{colim}(a, c, b)$. Легко проверить, что для любой подкатегории в c -DESC^V, все объекты которой обладают кодекартовыми квадратами, отображение colim является функцией морфизмов функтора, действующего из этой подкатегории в c -DESC, переводя любой ее объект в вершину кодекартова квадрата над ним. В частности, если все диаграммы из класса $\text{pmao}(\text{Ob AO})$ имеют копредел, то можно определить функтор $\text{pmod} = \text{colim} \circ \text{pmao} : \text{AO} \rightarrow c\text{-DESC}$.



Проверим, что функтор pmod является левым сопряженным к функтору mod^* с тождественной коединицей, так что согласно определению 1 АО-технология структурируема. Поскольку любая АО-модель из класса $\text{mod}^*(\text{Ob } c\text{-DESC})$ имеет тождественный морфизм в качестве разметки, функтор colim можно выбрать так, чтобы выполнялось соотношение $\text{pmod} \circ \text{mod}^* = 1_{c\text{-DESC}}$. Покажем, что для любых c -DESC-объекта P , АО-модели $S = \langle A, l : \text{sig}(A) \rightarrow L \rangle$ и c -DESC-морфизма $p : \text{pmod}(S) \rightarrow P$ существует единственный АО-морфизм $r : S \rightarrow \text{mod}^*(P)$ такой, что $\text{pmod}(r) = p$. Действительно, пусть $u : A \rightarrow \text{pmod}(S) \leftarrow \text{sig}^*(L) : v$ — ребра копредела диаграммы $\text{pmao}(S)$, в частности $u \circ \varepsilon_A = v \circ \text{sig}^*(l)$. Отсюда $p \circ u \circ \varepsilon_A = p \circ v \circ \text{sig}^*(l)$, следовательно, $\text{sig}(p \circ u) = \text{sig}(p \circ v) \circ l$, так что имеется АО-морфизм $\langle p \circ u, \text{sig}(p \circ v) \rangle : S \rightarrow \text{mod}^*(P)$. Ввиду предложения 3 под действием функтора pmod он переходит в p , поэтому является искомым r (его единственность вытекает из того, что 1_P и ε_P — мономорфизмы, как видно на диаграмме).



Обратно, предположим, что функтор mod^* обладает левым сопряженным $\text{mod}^{**} : \text{AO} \rightarrow c\text{-DESC}$ с тождественной коединицей. Обозначим через κ единицу этого сопряжения. Для произвольной АО-модели $S = \langle A, l : \text{sig}(A) \rightarrow L \rangle$ положим $T = \text{mod}^{**}(S)$, $s = \text{mod}(\kappa_S) : A \rightarrow T$. По определению АО-морфизма имеем $\text{str}(\kappa_S) \circ l = 1_{\text{sig}(T)} \circ \text{sig}(s) = \text{sig}(s)$, поэтому тождество коединицы сопряжения $\text{sig}^* \dashv \text{sig}$ для s можно записать в виде $s \circ \varepsilon_A = \varepsilon_T \circ \text{sig}^*(\text{sig}(s)) = w \circ \text{sig}^*(l)$, где $w = \varepsilon_T \circ \text{sig}^*(\text{str}(\kappa_S))$. Проверим, что оно определяет кокартов квадрат — копредел диаграммы $\text{ptao}(S)$. Выберем произвольно пару $c\text{-DESC}$ -морфизмов $a : A \rightarrow X \leftarrow \text{sig}^*(L) : y$ такую, что $a \circ \varepsilon_A = y \circ \text{sig}^*(l)$. Применяя функтор sig к этому равенству, получаем $\text{sig}(a) = \text{sig}(y) \circ l$, так что имеется АО-морфизм $f = \langle a, \text{sig}(y) \rangle : S \rightarrow \text{mod}^*(X)$. Тождество единицы для него имеет вид $f = \text{mod}^*(q) \circ \kappa_S$, где $q = \text{mod}^{**}(f) : T \rightarrow X$. Применяя к этому тождеству функтор mod , получаем, что $a = q \circ s$. Применяя к нему же функтор str , получаем, что $\text{sig}(y) = \text{str}(\text{mod}^*(q)) \circ \text{str}(\kappa_S) = \text{sig}(q \circ w)$, откуда $y = q \circ w$ ввиду унивалентности функтора sig . Поэтому q — стрелка копредела (ее единственность вытекает из определения сопряжения: для любого $c\text{-DESC}$ -морфизма q' соотношения $a = q' \circ s$ и $y = q' \circ w$ влекут $f = \text{mod}^*(q') \circ \kappa_S$, откуда $q' = q$). \square



Следствие 1.1. Диаграмма частичной модуляризации любой модуляризуемой АО-модели имеет копредел, ребро которого, параллельное действию

частичной модуляризации, задает модуляризацию аспектной структуры модели. \square

Следствие 1.2. Модуляризация аспектной структуры любой модуляризуемой АО-модели является регулярным эпиморфизмом. \square

Следствие 1.3. Если некоторая формальная технология проектирования обладает полным АО-расширением, то АО-технология над ней структурируема. \square

Следствие 1.4. Если категория $c\text{-DESC}$ подвижна, то формальная технология проектирования обладает наибольшим АО-расширением, причем АО-модель $\langle A, l \rangle$ содержится в его кообласти тогда и только тогда, когда существует регулярный $c\text{-DESC}$ -эпиморфизм с областью A , переходящий в l под действием функтора sig . \square

Напомним, что категория $c\text{-DESC}$, рассматриваемая вместе с функтором sig , называется подвижной (transportable) [11, определение 5.1], если для любых $c\text{-DESC}$ -объекта A и SIG -изоморфизма i с областью $\text{sig}(A)$ существует $c\text{-DESC}$ -изоморфизм j с областью A , удовлетворяющий условию $\text{sig}(j) = i$. Например, категория **Pos** с «забывающим» функтором $|-| : \mathbf{Pos} \rightarrow \mathbf{Set}$ очевидным образом является подвижной, поэтому ввиду следствия 1.4 технология событийного моделирования обладает наибольшим АО-расширением, которое строится следующим образом.

Обозначим через **HAOSM** полную подкатеорию в категории всех помеченных сценариев, состоящую из всех объектов, которые допускают частичное упорядочение своих множеств меток, превращающее разметку в монотонное отображение. Вложение категории **Pos** в **HAOSM**, задающее дискретную разметку сценария, является АО-расширением: правым сопряженным к нему служит функтор, «забывающий» разметку, а левым сопряженным — функтор, переводящий произвольный **HAOSM**-объект $\langle X, \leq, l : X \rightarrow L \rangle$ в частично упорядоченное множество $\langle L, \preceq \rangle$, где \preceq — пересечение всех частичных порядков на L , превращающих l в монотонное отображение. Ясно, что категория **HAOSM** содержит все помеченные сценарии, разметка которых может быть восстановлена до трансформации модулей (т. е. непомеченных сценариев), как описано в конце разд. 4. Есть в **HAOSM** и помеченные сценарии, не обладающие этим свойством, например следующий трехэлементный сценарий с двумя несравнимыми элементами, размеченный двумя метками: $\bullet \leftarrow \bullet \rightarrow \circ$. В свою очередь, линейно упорядоченный сценарий $\bullet \rightarrow \circ \rightarrow \bullet$ («интерливинг») не содержится в **HAOSM**, так что технология событийного моделирования не обладает полным АО-расширением.

6 Заключение

Эффект от применения теории категорий к формализации процессов создания программных систем обусловлен возможностью явно выразить основные положения общей теории систем [6]. Поэтому теоретико-категорное аксиоматическое описание концепции АОП как расширения модульного подхода к интеграции систем формально выражает его суть с общесистемной точки зрения, которая согласно [20] представляет собой совокупность принципов обеспечения абстрагируемости, модульности, удобства сборки систем в условиях рассеивания задач. В работе удалось строго сформулировать и доказать путем вывода из аксиом ряд утверждений, характеризующих ключевые возможности и ограничения АОП как такового, поскольку его описание не зависит от частных парадигм программирования.

В результате появилась возможность рационально применять аспектно-ориентированный подход при проектировании сложных гетерогенных программных систем, эффективно комбинируя его с традиционным объектно-ориентированным. Например, для систем технологического управления был формально смоделирован и затем реализован программный модуль ведения единого журнала событий, с которым связаны (в смысле АОП) многочисленные аспекты обработки событий: модули изменения паспорта объекта автоматизации, расчета разнообразных показателей, проверки соответствия фактического хода автоматизируемых процессов плановому, оповещения участников процессов и т. д. [21]. Такое динамическое связывание аспектов позволяет тонко настраивать порядок выполнения рассеянных задач непосредственно в функционирующей системе, без дублирования их программного кода и вообще без участия программистов.

Литература

1. Репин В. В., Елиферов В. Г. Процессный подход к управлению. Моделирование бизнес-процессов. — М.: Манн, Иванов и Фербер, 2013. 544 с.
2. Kiczales G., Lamping J., Mendhekar A., Maeda C., Lopes C. V., Loingtier J.-M., Irwin J. Aspect-oriented programming // 11th Conference (European) on Object-Oriented Programming Proceedings / Eds. M. Aksit, S. Matsuoka. — Lecture notes in computer science ser. — Springer, 1997. Vol. 1241. P. 220–242.
3. Steimann F. The paradoxical success of aspect-oriented programming // Conference (International) OOPSLA'06 Proceedings. — Portland, 2006. P. 481–497.
4. Colyer A., Clement A., Harley G., Webster M. Eclipse AspectJ: Aspect-oriented programming with AspectJ and the Eclipse AspectJ development tools. — Addison-Wesley, 2004. 504 p.
5. Rashid A., Chitchyan R. Aspect-oriented requirements engineering: A roadmap // 13th Workshop (International) on Early Aspects Proceedings. — Leipzig, 2008. P. 35–41.
6. Goguen J. Categorical foundations for General Systems Theory // Advances in cybernetics and systems research. — London: Transcripta Books, 1973. P. 121–130.
7. Fiadeiro J. L., Lopes A., Wermelinger M. A mathematical semantics for architectural connectors // Generic programming — advanced lectures / Eds. R. C. Backhouse, J. Gibbons. — Lecture notes in computer science ser. — Springer, 2003. Vol. 2793. P. 190–234.
8. Douence R., Fradet P., Südholt M. Trace-based aspects // Aspect-oriented software development. — Reading: Addison Wesley, 2004. P. 201–218.
9. Jagadeesan R., Pitcher C., Riely J. Open bisimulation for aspects // Conference (International) AOSD'07 Proceedings. — Vancouver, Canada, 2007. P. 107–120.
10. Маклейн С. Категории для работающего математика / Пер. с англ. — М.: Физматлит, 2004. 352 с. (*Mac Lane S. Categories for the working mathematician. — Berlin — Heidelberg — New York: Springer, 1978. 317 p.*)
11. Adámek J., Herrlich H., Strecker G. Abstract and concrete categories. — New York: Wiley and Sons, 1990. 482 p.
12. Ковалёв С. П. Формальный подход к аспектно-ориентированному моделированию сценариев // Сиб. журн. индустр. математики, 2010. Т. 13. № 3. С. 30–42.
13. Gotel O., Finkelstein A. An analysis of the requirements traceability problem // 1st Conference (International) on Requirements Engineering Proceedings. — Colorado Springs, 1994. P. 94–101.
14. Egyed A., Grünbacher P., Heindl M., Biffl S. Value-based requirements traceability: Lessons learned // Design requirements engineering: A ten-year perspective / Eds. K. Lyytinen, P. Loucopoulos, J. Mylopoulos, B. Robinson. — Lecture notes in business information processing ser. — Springer, 2009. Vol. 14. P. 240–257.
15. Aizenbud-Reshef N., Nolan B., Rubin J., Shaham-Gafni Y. Model traceability // IBM Syst. J., 2006. Vol. 45. No. 3. P. 515–526.
16. Goguen J. A categorical manifesto // Math. Struct. Comp. Sci., 1991. Vol. 1. No. 1. P. 49–67.
17. Ковалёв С. П. Семантика аспектно-ориентированного моделирования данных и процессов // Информатика и её применения, 2013. Т. 7. Вып. 3. С. 70–80.
18. Pratt V. R. Modeling concurrency with partial orders // Int. J. Parallel Prog., 1986. Vol. 15. No. 1. P. 33–71.
19. Sutton S. M., Rouvellou I. Concern modeling for aspect-oriented software development // Aspect-oriented software development. — Reading: Addison Wesley, 2004. P. 479–505.
20. Rashid A., Moreira A. Domain models are not aspect free // 9th Conference (International) on Model Driven

- Engineering Languages and Systems Proceedings / Eds. O. Nierstrasz, J. Whittle, D. Harel, G. Reggio. — Lecture notes in computer science ser. — Springer, 2006. Vol. 4199. P. 155–169.
21. Андриюшкевич С. К., Ковалёв С. П. Динамическое связывание аспектов в крупномасштабных системах технологического управления // Вычисл. технологии, 2011. Т. 16. № 6. С. 3–12.

Поступила в редакцию 25.08.14

FORMAL AXIOMATIC APPROACH TO ASPECT-ORIENTED EXTENSION OF PROGRAMMING TECHNOLOGIES

S. P. Kovalyov

Institute of Control Problem, Russian Academy of Sciences, 65 Profsoyuznaya Str., Moscow 117997, Russian Federation

Abstract: The procedure of extending modular software systems design technologies by aspect-oriented techniques is considered. The extension is described as enrichment of formal module models by labeling their interfaces by concerns they handle which comprise aspect structure. A novel approach to separation of concerns based on the natural modularizing aspect structure is proposed. Partial modularization of the aspect structure is proposed to generalize this approach. In order to formalize these constructs at the general systems level independently of particular programming paradigms, the category theory is employed. Software engineering technologies are represented as categories with formal models of programs as objects and technological operations as morphisms. The aspect-oriented extension of the technology is axiomatically described as a functor between such categories that has appropriate right and left adjoints. The event-based approach to system modeling is employed as an illustrative case of the aspect-oriented extension.

Keywords: aspect-oriented programming; traceability; category theory; architecture school; separation of concerns

DOI: 10.14357/19922264150105

Acknowledgments

The research was financially supported by the Russian Foundation for Humanities (grant 13-03-00384).

References

1. Repin, V.V., and V.G. Eliferov. 2013. *Protsessnyy podkhod k upravleniyu. Modelirovanie biznes-protsessov* [Process approach to control. Business process modeling]. Moscow: Mann, Ivanov and Ferber. 544 p.
2. Kiczales, G., J. Lamping, A. Mendhekar, C. Maeda, C. V. Lopes, J.-M. Loingtier, and J. Irwin. 1997. Aspect-oriented programming. *11th Conference (European) on Object-Oriented Programming Proceedings*. Eds. M. Aksit and S. Matsuoka. Lecture notes in computer science ser. Springer. 1241:220–242.
3. Steimann, F. 2006. The paradoxical success of aspect-oriented programming. *Conference (International) OOPSLA'06 Proceedings*. Portland. 481–497.
4. Colyer, A., A. Clement, G. Harley, and M. Webster. 2004. *Eclipse AspectJ: Aspect-oriented programming with AspectJ and the Eclipse AspectJ development tools*. Addison-Wesley. 504 p.
5. Rashid, A., and R. Chitchyan. 2008. Aspect-oriented requirements engineering: A roadmap. *13th Workshop (International) on Early Aspects Proceedings*. Leipzig. 35–41.
6. Goguen, J. 1973. Categorical foundations for General Systems Theory. *Advances in cybernetics and systems research*. London: Transcripta Books. 121–130.
7. Fiadeiro, J.L., A. Lopes, and M. Wermelinger. 2003. A mathematical semantics for architectural connectors. *Generic programming — advanced lectures*. Eds. R. C. Backhouse and J. Gibbons. Lecture notes in computer science ser. Springer. 2793:190–234.
8. Douence, R., P. Fradet, and M. Südholt. 2004. Trace-based aspects. *Aspect-oriented software development*. Reading: Addison Wesley. 201–218.
9. Jagadeesan, R., C. Pitcher, and J. Riely. 2007. Open bisimulation for aspects. *Conference (International) AOSD'07 Proceedings*. Vancouver, Canada. 107–120.
10. Mac Lane, S. 1978. *Categories for the working mathematician*. Berlin – Heidelberg – New York: Springer. 317 p.
11. Adámek, J., H. Herrlich, and G. Strecker. 1990. *Abstract and concrete categories*. New York: Wiley and Sons. 482 p.
12. Kovalyov, S.P. 2010. Formal'nyy podkhod k aspektno-orientirovannomu modelirovaniyu stsensariiev [Formal approach to aspect-oriented scenario modeling]. *Sibirskiy*

- Zhurnal Industrial'noy Matematiki* [J. Appl. Industrial Math.] 13(3):30–42.
13. Gotel, O., and A. Finkelstein. 1994. An analysis of the requirements traceability problem. *1st Conference (International) on Requirements Engineering Proceedings*. Colorado Springs. 94–101.
 14. Egyed, A., P. Grünbacher, M. Heindl, and S. Biffl. 2009. Value-based requirements traceability: Lessons learned. *Design requirements engineering: A ten-year perspective*. Eds. Lyytinen, K., P. Loucopoulos, J. Mylopoulos, and B. Robinson. Lecture notes in business information processing ser. Springer 14:240–257.
 15. Aizenbud-Reshef, N., B. Nolan, J. Rubin, and Y. Shaham-Gafni. 2006. Model traceability. *IBM Syst. J.* 45(3):515–526.
 16. Goguen, J. 1991. A categorical manifesto. *Math. Struct. Comp. Sci.* 1(1):49–67.
 17. Kovalyov, S. P. 2013. Semantika aspektno-orientirovanogo modelirovaniya dannyykh i protsessov [Semantics of aspect-oriented modeling of data and processes]. *Informatika i ee Primeneniya — Inform. Appl.* 7(3):70–80.
 18. Pratt, V. R. 1986. Modeling concurrency with partial orders. *Int. J. Parallel Prog.* 15(1):33–71.
 19. Sutton, S. M., and I. Rouvellou. 2004. Concern modeling for aspect-oriented software development. *Aspect-oriented software development*. Reading: Addison Wesley. 479–505.
 20. Rashid, A., and A. Moreira. 2006. Domain models are not aspect free. *9th Conference (International) on Model Driven Engineering Languages and Systems Proceedings*. Eds. Nierstrasz, O., J. Whittle, D. Harel, and G. Reggio. Lecture notes in computer science ser. Springer. 4199:155–169.
 21. Andryushkevich, S. K., and S. P. Kovalyov. 2011. Dinamicheskoe svyazyvanie aspektov v krupnomasshtabnykh sistemakh tekhnologicheskogo upravleniya [Dynamic aspect weaving in large-scale manufacturing control systems]. *Vychislitel'nye Tekhnologii* [J. Comput. Technol.]. 16(6):3–12.

Received August 25, 2014

Contributor

Kovalyov Sergey P. (b. 1972) — Doctor of Science in physics and mathematics, senior scientist, Institute of Control Problem, Russian Academy of Sciences, 65 Profsoyuznaya Str., Moscow 117997, Russian Federation; kovalyov@nm.ru

УСТОЙЧИВЫЕ ЛИНЕЙНЫЕ УСЛОВНО ОПТИМАЛЬНЫЕ ФИЛЬТРЫ И ЭКСТРАПОЛЯТОРЫ ДЛЯ СТОХАСТИЧЕСКИХ СИСТЕМ С МУЛЬТИПЛИКАТИВНЫМИ ШУМАМИ

И. Н. Сеницын¹, Э. Р. Корепанов²

Аннотация: Статья посвящена теории аналитического синтеза непрерывных равномерно асимптотически устойчивых условно оптимальных (по среднеквадратическому критерию) линейных фильтров и экстраполяторов (ЛУОФ и ЛУОЭ) для линейных дифференциальных стохастических систем (СтС) с линейными мультипликативными шумами. Предполагается, что наблюдение входит как в уравнение состояния, так и в уравнение наблюдения. Белые шумы в уравнениях наблюдения и состояния предполагаются заданными априори в виде производных по времени от произвольных процессов с независимыми приращениями. Доказаны теоремы, лежащие в основе теории непрерывных устойчивых ЛУОФ и ЛУОЭ. Достаточные условия равномерной асимптотической устойчивости сформулированы в виде требований положительной определенности и равномерной стохастической ограниченности некоторых матриц, отражающих свойства наблюдаемости и управляемости. Приведен иллюстративный пример. Сформулированы некоторые обобщения.

Ключевые слова: мультипликативный белый шум; равномерная асимптотическая устойчивость; стохастическая система (СтС); точность; уравнение Риккати; линейный условно оптимальный фильтр и экстраполятор (ЛУОФ и ЛУОЭ)

DOI: 10.14357/19922264150106

1 Введение

В настоящее время теория ЛУОФ и ЛУОЭ Пугачёва для непрерывных линейных СтС с аддитивными гауссовскими и негауссовскими шумами подробно разработана и реализована в универсальных и специальных программных средствах (см., например, [1–5]).

Для гауссовских СтС с мультипликативными шумами в уравнениях состояния и наблюдения в [1–3] получены соответствующие матричные уравнения Риккати для оценки точности ЛУОФ и ЛУОЭ. При этом их устойчивость для стохастических процессов (СтП) в таких СтС не рассматривалась.

В настоящей статье на основе известных свойств матричных уравнений Риккати получены достаточные условия равномерной асимптотической устойчивости непрерывных ЛУОФ и ЛУОЭ и их теория. Сформулированы дальнейшие обобщения.

2 Равномерно асимптотически устойчивые линейные условно оптимальные фильтры

Рассмотрим наблюдаемую непрерывную (дифференциальную) СтС с мультипликативными и ад-

дитивными (в общем случае негауссовскими) белыми шумами, описываемую следующими стохастическими дифференциальными уравнениями Ито [1–3]:

$$\left. \begin{aligned} \dot{X}_t &= aY_t + a_1X_t + a_0 + \left(c_{10} + \sum_{r=1}^{n_y} c_{1r}Y_r + \right. \\ &\quad \left. + \sum_{r=1}^{n_x} c_{1,n_y+r}X_r \right) V, \quad X_{t_0} = X_0; \\ \dot{Y}_t &= bY_t + b_1X_t + b_0 + \left(c_{20} + \sum_{r=1}^{n_y} c_{2r}Y_r + \right. \\ &\quad \left. + \sum_{r=1}^{n_x} c_{2,n_y+r}X_r \right) V, \quad Y_{t_0} = Y_0, \end{aligned} \right\} (1)$$

где X_t и Y_t — векторы состояния и наблюдения ($\dim X_t = n_x$, $\dim Y_t = n_y$); $V = \dot{W}$ — белый шум ($\dim V = n_v$); W — векторный СтП с независимыми приращениями вида

$$W(t) = W_0(t) + \int_{R_0^a} c(\rho) P^0(t, d\rho). \quad (2)$$

Здесь $W_0(t)$ — винеровский СтП интенсивности $\nu_0(t)$; $c(\rho)$ — некоторая векторная функция той

¹Институт проблем информатики Российской академии наук, sinitsin@dol.ru

²Институт проблем информатики Российской академии наук, ekorepanov@ipiran.ru

же размерности, что и $W(t)$, q -мерного аргумента, а интеграл при любом $t \geq t_0$ представляет собой стохастический интеграл по центрированной пуассоновской мере $P^0(t, B)$ интенсивности $\nu_P(t, \rho)$, независимой от винеровского СтП $W_0(t)$ и имеющей независимые значения на непересекающихся множествах; B — борелевское множество пространства R_0^q с выколотым началом. При этом интенсивность СтП $W(t)$ определяется по формуле [1, 3, 5]

$$\nu = \nu(t) = \nu_0(t) + \int_{R_0^q} c(\rho)c(\rho)^T \nu_P(\tau, \rho) d\rho.$$

В (1) a_0, b_0, a, a_1, b, b_1 и c_{ij} ($i = 1, 2, j = 1, \dots, n_x$) — векторно-матричные функции t , не зависящие от $X_t = [X_1 \dots X_{n_x}]^T$ и $Y_t = [Y_1 \dots Y_{n_y}]^T$. Следуя [1, 3, 5], класс допустимых ЛУОФ зададим линейным уравнением

$$\begin{aligned} \dot{\hat{X}}_t &= aY_t + a_1\hat{X}_t + a_0 + \\ &+ \beta_t \left[\dot{Y}_t - (bY_t + b_1\hat{X}_t + b_0) \right]. \end{aligned} \quad (3)$$

Здесь β_t и σ_{ij} ($i, j = 1, 2$) определяются уравнениями:

$$\begin{aligned} \beta_t &= (R_t b_1^T + \sigma_{12}) \sigma_{22}^{-1}; \\ \sigma_{12} &= \left(c_{10} + \sum_{r=1}^{n_y+n_x} c_{1r} m_r \right) \nu \left(c_{20}^T + \sum_{r=1}^{n_y+n_x} c_{2r}^T m_r \right) + \\ &+ \sum_{r,s=1}^{n_y+n_x} c_{1r} \nu c_{2s}^T k_{rs}; \\ \sigma_{22} &= \left(c_{20} + \sum_{r=1}^{n_y+n_x} c_{2r} m_r \right) \nu \left(c_{20}^T + \sum_{r=1}^{n_y+n_x} c_{2r}^T m_r \right) + \\ &+ \sum_{r,s=1}^{n_y+n_x} c_{2r} \nu c_{2s}^T k_{rs}. \end{aligned} \quad (4)$$

Ковариационная матрица R_t ошибки фильтрации $\tilde{X}_t = \hat{X}_t - X_t$ удовлетворяет матричному уравнению Риккати:

$$\begin{aligned} \dot{R}_t &= a_1 R_t + R_t a_1^T - \\ &- (R_t b_1^T + \sigma_{12}) \sigma_{22}^{-1} (b_1 R_t + \sigma_{21}) + \sigma_{11}, \end{aligned} \quad (5)$$

где

$$\begin{aligned} \sigma_{21} &= \left(c_{20} + \sum_{r=1}^{n_y+n_x} c_{2r} m_r \right) \nu \left(c_{10}^T + \sum_{r=1}^{n_y+n_x} c_{1r}^T m_r \right) + \\ &+ \sum_{r,s=1}^{n_y+n_x} c_{2r} \nu c_{1s}^T k_{rs}; \end{aligned}$$

$$\begin{aligned} \sigma_{11} &= \left(c_{10} + \sum_{r=1}^{n_x+n_y} c_{1r} m_r \right) \nu \left(c_{10}^T + \sum_{r=1}^{n_x+n_y} c_{1r}^T m_r \right) + \\ &+ \sum_{r,s=1}^{n_y+n_x} c_{1r} \nu c_{1s}^T k_{rs}. \end{aligned}$$

При этом вероятностные моменты первого и второго порядка $m_t = [m_r]$, $K_t = [k_{rs}]$ вектора $Q_t = [X_t Y_t]^T$ ($r, s = 1, \dots, n_y + n_x$) определяются следующими уравнениями:

$$\dot{m}_t = \bar{a} m_t + \bar{a}_0, \quad m_{t_0} = m_0; \quad (6)$$

$$\begin{aligned} \dot{K}_t &= \bar{a} K_t + K_t \bar{a}^T + c_0 \nu c_0^T + \sum_{r=1}^{n_y+n_x} (c_0 \nu c_r^T + c_r \nu c_0^T) m_r + \\ &+ \sum_{r,s=1}^{n_y+n_x} c_r \nu c_s^T (m_r m_s + k_{rs}), \quad K_{t_0} = K_0 \end{aligned} \quad (7)$$

$$\begin{aligned} \left(\bar{a} = \begin{bmatrix} b & b_1 \\ a & a_1 \end{bmatrix}, \quad \bar{a}_0 = \begin{bmatrix} b_0 \\ a_0 \end{bmatrix}, \right. \\ \left. c_r = \begin{bmatrix} c_{2r} \\ a_{1r} \end{bmatrix}, \quad r = 0, 1, \dots, n_y + n_x \right). \end{aligned}$$

Таким образом, в основе теории синтеза ЛУОФ лежит следующее утверждение.

Теорема 2.1. Пусть процессы X_t и Y_t в дифференциальной СтС (1) обладают конечными вероятностными моментами второго порядка, а матрица (4) не вырождена ($\det \sigma_{22} \neq 0$). Тогда ЛУОФ определяется уравнением (3), а его точность — матричным уравнением Риккати (5).

Замечание 2.1. Фильтр, определяемый уравнением (3), оптимален в классе всех линейных фильтров, причем ЛУОФ является оптимальным и в более широком классе всех линейных фильтров. В частном случае линейной СтС (1) без мультипликативных шумов уравнения (3) и (5) совпадают с уравнениями теории оптимальной линейной фильтрации [1, 3]. Для винеровского процесса $W_0(t)$ при отсутствии мультипликативных шумов ЛУОФ оказывается оптимальным и в классе всех возможных фильтров.

Замечание 2.2. Особенностью ЛУОФ является то обстоятельство, что $m_r = m_{rt}$, $K_r = K_{rt}$ и $R = R_t$ могут быть вычислены заранее во время синтеза фильтра, так как не требуют знания результатов текущих наблюдений.

Далее, применяя известное матричное неравенство [6, 7] к уравнению Риккати (5), получим $\forall t \geq t_0$

$$0 \leq R_t(R_0, t_0) \leq u_A(t, t_0)R_0 u^T(t, t_0) + \int_{t_0}^t u_A(t, \tau) \bar{\sigma}_{11} u_A(t, \tau)^T d\tau = u_A(t, t_0)R_0 u(t, t_0)^T + \mathcal{W}(t_0, t).$$

Здесь введены следующие обозначения: $R_0 = R_{t_0}$ — неотрицательно определенная матрица; $u_A(t, t_0)$ — фундаментальная матрица однородного уравнения, полученного из уравнения Риккати (5):

$$\dot{Z} = AZ, \quad (8)$$

где

$$A = a_1 - \beta_t b_1 = a_1 - R_t b_1^T \sigma_{22}^{-1} b_1 - \sigma_{12} \sigma_{22}^{-1} b_1;$$

$$\mathcal{W}(t_0, t) = \int_{t_0}^t u_A(t, \tau) \bar{\sigma}_{11}(\tau) u_A(t, \tau)^T d\tau; \quad (9)$$

$$\bar{\sigma}_{11} = \sigma_{11} - \sigma_{12} \sigma_{22}^{-1} \sigma_{21}. \quad (10)$$

Наконец, обобщая понятия равномерной наблюдаемости и управляемости [6, 7], введем понятие равномерной ограниченности

$$0 < \alpha_1 I_n \leq \mathcal{W}(t - t', t) \leq \alpha_2 I_n \quad \forall t \geq t_0 + t', \quad (11)$$

где α_1, α_2, t' — постоянные; $n = n_y + n_x$; I_n — единичная ($n \times n$)-матрица.

Правая часть (5) удовлетворяет условию Липшица, поэтому имеет место локальное существование и единственность решения (5). Более того, верхняя граница решения (5) позволяет определить постоянную Липшица на любом конечном интервале времени. А это значит, что уравнение (5) имеет глобальное единственное решение. Методом детерминированных функций Ляпунова вида $\mathcal{L} = \xi_t^T R_t \xi$ для уравнения (8) аналогично [6, 7] устанавливается следующая теорема.

Теорема 2.2. Пусть в условиях теоремы 2.1 система уравнений (1) равномерно вполне наблюдаема и равномерно вполне управляема, т.е. матрица (9) положительно определена и выполнены условия (11). Тогда ЛУОФ (3) равномерно асимптотически устойчив, т.е. тривиальное решение образуемого из (3) при $\dot{Y}_t = 0, Y_t = 0$ уравнения (8) равномерно асимптотически устойчиво.

Случай независимых аддитивных и мультипликативных негауссовских белых шумов V_1 и V_2 соответственно в уравнениях состояния и наблюдения получается на основе уравнений (1), если принять

$$V = [V_1^T V_2^T]^T; \quad \nu = \begin{bmatrix} \nu_1 & 0 \\ 0 & \nu_2 \end{bmatrix}; \quad (12)$$

$$\left(c_{i0} + \sum_{r=1}^{n_y} c_{ir} Y_r + \sum_{s=1}^{n_x} c_{i, n_y+r} X_r \right) V = \left(c'_{i0} + \sum_{r=1}^{n_y} c'_{ir} Y_r + \sum_{s=1}^{n_x} c'_{i, n_y+r} X_r \right) V_i; \quad (13)$$

$$\sigma_{ii} = \left(c'_{i0} + \sum_{r=1}^{n_y+n_x} c'_{ir} m_r \right) \nu_i \left(c_{i0}^T + \sum_{r=1}^{n_y+n_x} c_{ir}^T m_r \right),$$

$$\sigma_{ij} = \sigma_{ji} = 0 \quad (i \neq j);$$

$$\bar{\sigma}_{11} = \sigma_{11}; \quad (14)$$

$$\beta_t = R_t b_1^T \sigma_{22}^{-1}. \quad (15)$$

Теорема 2.3. Пусть процессы X_t и Y_t в системе уравнений (1) при условиях (12) и (13) обладают конечными вероятностными моментами второго порядка, а матрица σ_{22} не вырождена ($\det \sigma_{22} \neq 0$). Тогда ЛУОФ определяется уравнением (3) при условии (15), а его точность — матричным уравнением (5). Для обеспечения равномерной асимптотической устойчивости ЛУОФ (3) достаточно положительной определенности матрицы (9) при условии (14) и равномерной ограниченности (11).

3 Равномерно асимптотически устойчивые линейные условно оптимальные экстраполяторы

Обобщая [1, 3] на случай белых шумов V_1 и V_2 вида (2) и с учетом теоремы 2.1, придем к следующим утверждениям.

Теорема 3.1. Пусть процессы X_t, Y_t в уравнениях

$$\dot{X}_t = a_1 X_t + a_0 + \left(c_{10} + \sum_{r=1}^{n_x} c_{1, n_y+r} X_r \right) V_1,$$

$$V_1 = \dot{W}_1, \quad X_{t_0} = X_0,$$

$$\dot{Y}_t = b Y_t + b_1 X_t + b_0 +$$

$$+ \left(c_{20} + \sum_{r=1}^{n_y} c_{2r} Y_r + \sum_{r=1}^{n_x} c_{2, n_y+r} X_r \right) V_2,$$

$$V_2 = \dot{W}_2, \quad Y_{t_0} = Y_0,$$

где $W_1 = W_1(t), W_2 = W_2(t)$ — независимые процессы с независимыми приращениями, обладают конечными вероятностными моментами второго порядка. Тогда уравнения ЛУОЭ имеют вид:

$$\hat{X}_t = a_1(t + \Delta) \hat{X}_t + a_0(t + \Delta) + \beta_t \left[\dot{Y}_t - (b Y_t + b_1 \varepsilon_t^{-1} \hat{X}_t + b_0 - b_1 \varepsilon_t^{-1} h_t) \right]; \quad (16)$$

$$\varepsilon_t = u(t + \Delta, t), \quad (17)$$

где $u(t, \tau)$ — фундаментальная матрица уравнения:

$$\begin{aligned} \dot{u}_t &= a_1(t)u_t; \\ h_t &= h(t) = \int_t^{t+\Delta} u(t+\Delta, \tau)a_0(\tau) d\tau. \end{aligned} \quad (18)$$

Необходимые для вычисления матриц σ_{22} и β_t находятся согласно теореме 2.3 для составного вектора $[Y_t^T X_t^T \hat{X}_t^T]^T$. Роль матриц c_{1r} , c_{2r} играют матрицы $[0 \ c_{1r}]$ и $[c_{2r} \ 0]$, а матрица ν — диагональна. При этом точность экстраполяции определяется путем интегрирования следующего уравнения:

$$\begin{aligned} \dot{R}_t &= a_1(t+\Delta)R_t + R_t a_1(t+\Delta)^T - \\ &- \beta_t \left[\left(c_{20} + \sum_{r=1}^{n_y+n_x} c_{2r} m_r \right) \nu_1 \left(c_{20}^T + \sum_{r=1}^{n_y+n_x} c_{2r}^T m_r \right) + \right. \\ &\quad \left. + \sum_{r=1}^{n_y+n_x} c_{2r} \nu_1 c_{2s}^T k_{rs} \right] \beta_t^T + \left[c_{10}(t+\Delta) + \right. \\ &\quad \left. + \sum_{r=n_y+1}^{n_y+n_x} c_{1r}(t+\Delta) m_r(t+\Delta) \right] \nu_1(t+\Delta) \times \\ &\times \left[c_{10}(t+\Delta)^T + \sum_{r=n_y+1}^{n_y+n_x} c_{1r}(t+\Delta)^T m_r(t+\Delta) \right] + \\ &\quad + \sum_{s=n_y+1}^{n_y+n_x} c_{1r}(t+\Delta) \nu_1(t+\Delta) c_{1s}(t+\Delta)^T k_{rs} = \\ &= A_\Delta R_t + R_t A_\Delta^T + R_t B_\Delta R_t + C_\Delta, \end{aligned}$$

причем для обеспечения равномерной асимптотической устойчивости достаточно выполнения условий

$$0 < \alpha_1 I_n \leq \mathcal{W}_\Delta(t-t', t) \leq \alpha_2 I_n,$$

где

$$\mathcal{W}_\Delta(t-t', t) = \int_{t-t'}^t u_\Delta(t, \tau) C(\tau) u_\Delta(t, \tau)^T d\tau.$$

Здесь $u_\Delta(t, \tau)$ — фундаментальная матрица для уравнения $\dot{Z} = A_\Delta Z$.

Из (16) видно, что ЛУОЭ можно представить в виде последовательного ЛУОФ, усилителя с коэффициентом усиления ε_t (17) и сумматора, вводящего неслучайное слагаемое h_t (18). Найденный ЛУОЭ оптимален в классе всех линейных экстраполяторов [1–3].

4 Пример

Рассмотрим случай, когда скалярные уравнения (1) и (2) содержат независимые белые шумы V_1 и V_2 :

$$\begin{aligned} \dot{X}_t &= aY_t + a_1X_t + (c_{11}X_t + c_{12}Y_t)V_1, \\ \dot{Y}_t &= bY_t + b_1X_t + (c_{21}X_t + c_{22}Y_t)V_2, \\ \nu &= \begin{bmatrix} \nu_1 & 0 \\ 0 & \nu_2 \end{bmatrix}. \end{aligned}$$

Обратим внимание на то, что c_{1r} и c_{2r} здесь не те, что в уравнениях (1). Они представляют собой соответственно первые и вторые элементы матриц-строк, на которые умножается вектор $[V_1 \ V_2]^T$ в (1). Для простоты оставляем для них обозначения c_{1r} и c_{2r} . Тогда имеем:

$$\begin{aligned} \beta_t &= \sigma_{22}^{-1} b_1 R_t; \\ \sigma_{22} &= \nu_2 (c_{20} + c_{21} m_1 + c_{22} m_2)^2 + \\ &\quad + \nu_2 (c_{21}^2 k_{11} + 2c_{22} c_{21} k_{21} + c_{22}^2 k_{22}). \end{aligned}$$

Уравнения (6), (7) и (5), определяющие m_1 , m_2 , k_{11} , k_{12} , k_{21} , k_{22} и R_t , имеют следующий вид:

$$\dot{m}_1 = a_1 m_1 + a m_2; \quad \dot{m}_2 = b_1 m_1 + b m_2; \quad (19)$$

$$\left. \begin{aligned} \dot{k}_{11} &= 2(a_1 k_{11} + a k_{12}) + \\ &\quad + \nu_1 (c_{11} m_1 + c_{12} m_2 + c_{10})^2 + \\ &\quad + \nu_1 (c_{11}^2 k_{11} + 2c_{12} c_{11} k_{12} + c_{12}^2 k_{22}); \\ \dot{k}_{12} &= (a_1 + b) k_{12} + b_1 k_{11} + a k_{22}; \\ \dot{k}_{22} &= 2(b_1 k_{12} + b k_{22}) + \\ &\quad + \nu_2 (c_{21} m_1 + c_{22} m_2 + c_{20})^2 + \\ &\quad + \nu_2 (c_{21}^2 k_{11} + 2c_{22} c_{21} k_{12} + c_{22}^2 k_{22}), \end{aligned} \right\} \quad (20)$$

$$\begin{aligned} \dot{R}_t &= 2a_1 R_t - \sigma_{22}^{-1} b_1^2 R_t^2 + \nu_1 (c_{11} m_1 + c_{12} m_2 + c_{10})^2 + \\ &\quad + \nu_1 (c_{11}^2 k_{11} + 2c_{11} c_{12} k_{12} + c_{12}^2 k_{22}). \end{aligned} \quad (21)$$

Согласно (10) $\bar{\sigma}_{11}$, входящая в условия устойчивости ЛУОФ (теорема 3.1), определяется формулой

$$\begin{aligned} \bar{\sigma}_{11} &= \sigma_{11} = \nu_1 (c_{11} m_1 + c_{12} m_2 + c_{10}) + \\ &\quad + \nu_1 (c_{11}^2 k_{11} + 2c_{11} c_{12} k_{12} + c_{12}^2 k_{22}). \end{aligned}$$

При $a = a_0 = c_{12} = 0$, $a_1 = \text{const}$ ЛУОЭ представляет собой последовательное соединение ЛУОФ и усилителя с коэффициентом $\varepsilon_t = \exp(a_1 \Delta)$.

Приравнивая правые части уравнений (19)–(21) нулю, найдем условия для синтеза стационарных ЛУОФ и ЛУОЭ с постоянными $\beta_t = \beta^*$, $R_t = R^*$, $\varepsilon_t = \varepsilon^*$:

$$\begin{aligned} m_1^* &= 0; \quad m_2^* = 0; \\ 2(a_1 k_{11}^* + a k_{12}^*) + \\ &\quad + \nu_1^* (c_{11}^2 k_{11}^* + 2c_{12} c_{11} k_{12}^* + c_{12}^2 k_{22}^*) = 0; \\ (a_1 + b) k_{12}^* + b_1 k_{11}^* + a k_{22}^* &= 0; \end{aligned}$$

$$\begin{aligned}
& 2(b_1 k_{12}^* + b k_{22}^*) + \\
& \quad + \nu_2^*(c_{21}^2 k_{11}^* + 2c_{22}c_{21}k_{12}^* + c_{22}^2 k_{22}^*) = 0; \\
& 2a_1 R^* - \sigma_{22}^{*-1} b_1^2 R^{*2} + \\
& \quad + \nu_1^*(c_{11}^2 k_{11}^* + 2c_{12}c_{11}k_{12}^* + c_{12}^2 k_{22}^*) = 0; \\
& \quad \beta^* = \sigma_{22}^{*-1} b_1 R^*; \\
& \bar{\sigma}_{11}^* = \nu_1^*(c_{11}^2 k_{11}^* + 2c_{11}c_{12}k_{12}^* + c_{12}^2 k_{22}^*); \\
& \sigma_{22}^* = \nu_2^*(c_{21}^2 k_{11}^* + 2c_{22}c_{21}k_{12}^* + c_{22}^2 k_{22}^*).
\end{aligned}$$

Для равномерной асимптотической устойчивости ЛУОФ достаточно отрицательности коэффициента

$$A = a_1 - \frac{b_1^2 R^*}{\sigma_{22}^*} < 0.$$

5 Заключение

Получено обобщение известных результатов по теории синтеза непрерывных равномерно асимптотически устойчивых ЛУОФ и ЛУОЭ для случая наблюдаемых дифференциальных СтС с мультипликативными негауссовскими белыми шумами. Приводится иллюстративный пример.

Результаты могут быть использованы и для теории синтеза дискретных ЛУОФ и ЛУОЭ для непрерывных и дискретных СтС, если воспользоваться численными методами приведения нелинейных стохастических дифференциальных уравнений к разностным на основе обобщенной формулы Ито [2, 8].

Практический интерес представляет задача синтеза устойчивых ЛУОФ и ЛУОЭ по критериям

устойчивости, отличным от критерия равномерной асимптотической устойчивости, а также для исследования вопросов эквивалентности различных широкополосных шумов (в том числе автокоррелированных).

Литература

1. Пугачев В. С., Синицын И. Н. Стохастические дифференциальные системы. Анализ и фильтрация. — М.: Наука, 1990. 632 с. (Pugachev V. S., Sinitsyn I. N. Stochastic differential systems. Analysis and filtering. — Chichester, New York: John Wiley, 1987. 549 p.)
2. Пугачев В. С., Синицын И. Н. Теория стохастических систем. — М.: Логос, 2000. 1000 с. (Stochastic systems. Theory and applications. — Singapore: World Scientific, 2001. 908 p.)
3. Синицын И. Н. Фильтры Калмана и Пугачева. — 2-е изд. — М.: Логос, 2007. 776 с.
4. Korepanov E. R. Стохастические информационные технологии на основе фильтров Пугачева // Информатика и её применения, 2011. Т. 5. Вып. 2. С. 36–57.
5. Синицын И. Н., Синицын В. И. Лекции по нормальной и эллипсоидальной аппроксимации в стохастических системах. — М.: ТРУС ПРЕСС, 2013. 476 с.
6. Kalman R. A new approach to linear filtering and prediction problems // J. Basic Eng. (ASME Trans.), 1960. Vol. 82D. P. 35–45.
7. Ройтенберг Я. Н. Автоматическое управление. — 3-е изд., перераб. и доп. — М.: Наука, 1992. 576 с.
8. Синицын И. Н. Параметрическое статистическое и аналитическое моделирование распределений в нелинейных стохастических системах на многообразиях // Информатика и её применения, 2013. Т. 7. Вып. 2. С. 4–16.

Поступила в редакцию 22.09.14

STABLE LINEAR CONDITIONALLY OPTIMAL FILTERS AND EXTRAPOLATORS FOR STOCHASTIC SYSTEMS WITH MULTIPLICATIVE NOISES

I. N. Sinitsyn and E. R. Korepanov

Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation

Abstract: The applied theory of analytical synthesis of linear conditionally optimal filters and extrapolators in linear differential stochastic systems with white multiplicative non-Gaussian noises is presented. Efficient criteria of unique asymptotic stability of conditionally optimal filters and extrapolators are formulated in terms of special positive definite integral forms and unique boundedness of controllability and observability matrices. White noises are assumed to be derivatives of additive and multiplicative non-Gaussian arbitrary stochastic processes with independent increments. An illustrative example is given. Some generalizations are discussed.

Keywords: accuracy and unique asymptotic stability of filters; differential stochastic systems; linear conditionally optimal filters and extrapolators; multiplicative white noises; Riccati equation

DOI: 10.14357/19922264150106

References

1. Pugachev, V. S., and I. N. Sinitsyn. 1987. *Stochastic differential systems. Analysis and filtering*. Chichester, New York: John Wiley. 549 p.
2. Pugachev, V. S., and I. N. Sinitsyn. 2001. *Stochastic systems. Theory and applications*. Singapore: World Scientific. 908 p.
3. Sinitsyn, I. N. 2007. *Fil'try Kalmana i Pugacheva* [Kalman and Pugachev filters]. 2-e izd. Moscow: Logos. 776 s.
4. Korepanov, E. R. 2011. Stokhasticheskie informatsionnye tekhnologii na osnove fil'trov Pugacheva [Stochastic informational technologies based on Pugachev filters]. *Informatika i ee Primeneniya — Inform Appl.* 5(2):36–57.
5. Sinitsyn, I. N., and V. I. Sinitsyn. 2013. Lektsii po normal'noy i ellipsoidal'noy approksimatsii v stokhasticheskikh sistemakh [Lectures on normal and ellipsoidal approximation of distributions in stochastic systems]. Moscow: TORUS PRESS. 476 p.
6. Kalman, R. 1960. A new approach to linear filtering and prediction problems. *J. Basic Eng. (ASME Trans.)* 82D:35–45.
7. Roytenberg, Ya. N. 1992. *Avtomaticheskoe upravlenie* [Automatic control]. 3rd ed. Moscow: Nauka. 576 p.
8. Sinitsyn, I. N. 2013. Parametricheskoe statisticheskoe i analiticheskoe modelirovanie raspredeleniy v nelineynykh stokhasticheskikh sistemakh na mnogoobraznykh [Parametric statistical and analytical modeling of distributions in nonlinear stochastic systems on manifolds]. *Informatika i ee Primeneniya — Inform Appl.* 7(2):4–16.

Received September 22, 2014

Contributors

Korepanov Eduard R. (b. 1966) — Candidate of Science (PhD) in technology, Head of Laboratory, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; Ekorepanov@ipiran.ru

Sinitsyn Igor N. (b. 1940) — Doctor of Science in technology, professor, Honored scientist of Russian Federation, Head of Department, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; sinitsin@dol.ru

ВЫБОР ОПТИМАЛЬНОЙ МОДЕЛИ КЛАССИФИКАЦИИ ФИЗИЧЕСКОЙ АКТИВНОСТИ ПО ИЗМЕРЕНИЯМ АКСЕЛЕРОМЕТРА*

М. С. Попова¹, В. В. Стрижов²

Аннотация: Решается проблема построения оптимальных устойчивых моделей в задаче классификации физической активности человека. Каждый тип физической активности конкретного человека описывается набором признаков, сгенерированных по временным рядам с акселерометра. В условиях мультиколлинеарности признаков выбор устойчивых моделей классификации затруднен из-за необходимости оценки большого числа параметров этих моделей. Оценка оптимального значения параметров также затруднена в связи с тем, что функция ошибок имеет большое количество локальных минимумов в пространстве параметров. В работе исследуются модели, принадлежащие классу двуслойных нейронных сетей. Ставится задача нахождения Парето-оптимального фронта на множестве допустимых моделей. Предлагаются критерии оптимального, последовательного и устойчивого прореживания нейронной сети, критерий наращивания сети, а также строится стратегия пошаговой модификации модели с использованием предложенных критериев. В вычислительном эксперименте модели, порождаемые предложенной стратегией, сравниваются по трем критериям качества — сложности, точности и устойчивости.

Ключевые слова: классификация; нейронные сети; сложность; устойчивость; оптимальность по Парето; критерии прореживания и наращивания

DOI: 10.14357/19922264150107

1 Введение

Для получения точного и устойчивого прогноза физической активности человека необходимы методы, позволяющие выбирать адекватные модели из некоторого множества допустимых моделей-претендентов. Проблема выбора моделей обсуждается в работах [1–3]. Настройка параметров универсальной модели является нетривиальной многоэкстремальной оптимизационной задачей. Предлагается упростить эту задачу, рассматривая наборы последовательно порождаемых устойчивых моделей заданной сложности. Модели порождаются путем модификации структуры искусственной нейронной сети. Решается задача последовательной модификации нейронной сети. Требуется получить нейронную сеть с небольшим числом связей между нейронами, которая достаточно точно решала бы задачу классификации физической активности человека по показаниям акселерометра и обладала бы устойчивостью к возмущениям данных. Ввиду этого возникает задача минимизации сложности модели без потери точности классификации [4].

Существует два базовых подхода к решению задачи выбора сетей оптимальной структуры: *наращивание структуры сети* (network growing) [5] и *прореживание структуры сети* (network pruning) [6–8].

Согласно первому подходу в качестве начальной модели выбирается сеть недостаточной сложности, решающая поставленную задачу с большим значением функции ошибки, после чего в сеть добавляются новые нейроны и связи между ними. В [5] описаны некоторые методы наращивания, приведен сравнительный анализ генетических алгоритмов с алгоритмом байесовской оптимизации. В алгоритмах метода прореживания модифицируется многослойная сеть с избыточным числом нейронов и связей между ними. Классическими алгоритмами прореживания нейронных сетей являются «optimal brain damage» [7] и «optimal brain surgery» [8], основанные на вычислении вторых производных функции ошибки. Также получили развитие *гибридные алгоритмы*, в которых объединяются оба упомянутых выше подхода [9–11].

В данной работе предлагается стратегия пошаговой модификации нейронной сети, комбини-

* Работа поддержана Skolkovo Institute of Science and Technology (Skoltech) в рамках SkolTech/MITInitiative.

¹ Московский физико-технический институт, maria_popova@phystech.edu

² Вычислительный центр Российской академии наук им. А. А. Дородницына, strijov@ccas.com

рующая этапы добавления и удаления параметров [12, 13]. Стратегия включает в себя критерии прореживания и наращивания структуры сети, критерии останова этапов добавления и удаления параметров, а также критерий останова процедуры модификации. Согласно предложенной стратегии процедура модификации начинается с нейронной сети избыточной сложности и чередует шаги удаления и добавления параметров до тех пор, пока этот процесс не стабилизируется согласно критерию останова процедуры модификации. Критерии прореживания и наращивания позволяют на каждом шаге процедуры модификации выбирать параметр, добавление или удаление которого улучшит качество нейронной сети. Качество сети оценивается по трем критериям: сложности, точности и устойчивости [2, 14, 15]. Также предлагается рассматривать процедуру пошаговой модификации нейронной сети как путь в многомерном кубе.

В вычислительном эксперименте определяются значения критериев качества для нейронных сетей, порождаемых предложенной стратегией. В качестве тестового примера рассматривается задача классификации физической активности человека по измерениям акселерометра [16].

2 Постановка задачи

Дана выборка $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{t}_i)\}, i \in \mathcal{I} = \{1 \dots m\}$, состоящая из m объектов \mathbf{x} , каждый из которых описывается n признаками, $\mathbf{x}_i \in \mathbb{R}^n$, и принадлежит одному из z классов $\mathbf{t}_i \in \{0, 1\}^z$. Также задано разбиение множества индексов выборки $\mathcal{I} = \mathcal{L} \sqcup \mathcal{T}$ на обучающую $(\mathbf{x}_i, \mathbf{t}_i), i \in \mathcal{L}$, и контрольную $(\mathbf{x}_i, \mathbf{t}_i), i \in \mathcal{T}$. Требуется выбрать устойчивую модель классификации оптимальной сложности.

Определение 1. Моделью назовем отображение

$$\mathbf{f} : \left(\begin{array}{c} \mathbf{w} \\ k \times 1 \end{array}, \begin{array}{c} \mathbf{x} \\ 1 \times n \end{array} \right) \mapsto \begin{array}{c} \mathbf{y} \\ 1 \times z \end{array},$$

где $\mathbf{w} = [w_1, \dots, w_j, \dots, w_k]^T, j \in \mathcal{J} = \{1, \dots, k\}$, — вектор параметров модели; $\mathbf{x} \in \mathbb{R}^{n \times m}$ — матрица плана; $\mathbf{y} \in \{0, 1\}^z$ — зависимая переменная.

Предполагается, что переменная \mathbf{y} — мультиномиально распределенная случайная величина, а переменная \mathbf{w} имеет нормальное распределение с нулевым математическим ожиданием:

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \quad (1)$$

где \mathbf{A}^{-1} — ковариационная матрица параметров общего вида, положительно-определенная: $\mathbf{w}^T \mathbf{A} \mathbf{w} > 0$ для любого $\mathbf{w} \in \mathbb{R}^k$.

В данной работе рассматриваются модели, принадлежащие классу двухслойных нейронных сетей с функциями активации \tanh и softmax :

$$\mathbf{a}(\mathbf{x}) = \mathbf{W}_2^T \tanh(\mathbf{W}_1^T \mathbf{x}); \quad (2)$$

$$\mathbf{f}(\mathbf{x}) = \frac{\exp(\mathbf{a}(\mathbf{x}))}{\sum_{j=1}^n \exp(a_j(\mathbf{x}))}.$$

Вектор \mathbf{f} интерпретируется как вектор вероятностей: f_ξ есть вероятность того, что вектор \mathbf{x} принадлежит классу с номером ξ :

$$\mathbf{f}(\mathbf{x}) = \{f_\xi\}, \quad 0 \leq f_\xi \leq 1, \quad \sum f_\xi = 1, \quad \xi = 1 \dots z.$$

Под вектором параметров двухслойной нейронной сети будем понимать $\mathbf{w} = \text{vec}(\mathbf{W}_1^T | \mathbf{W}_2^T)$, где \mathbf{W}_1 и \mathbf{W}_2 — матрицы весов первого и второго слоя нейронной сети (2). Вектор $\mathbf{y} = [y_1, \dots, y_\xi, \dots, y_z]^T$ определим следующим образом:

$$y_\xi = \begin{cases} 1, & \text{если } \xi = \arg \max_{\xi \in \{1, \dots, z\}} (p_\xi); \\ 0 & \text{иначе.} \end{cases}$$

Вектор \mathbf{y} — это вектор метки класса, полученный для объекта \mathbf{x} с помощью построенной модели, в то время как вектор \mathbf{t} — это вектор метки класса объекта \mathbf{x} из выборки \mathcal{D} .

Под *структурным* параметром двухслойной нейронной сети будем понимать количество нейронов в скрытом слое нейронной сети — N_h . Матрица весов первого слоя имеет размерность $n \times N_h$, матрица весов второго слоя имеет размерность $N_h \times z$. Далее будем считать, что структурный параметр фиксирован и одинаков для всех рассматриваемых моделей.

Определение 2. Параметр w_j модели \mathbf{f} назовем активным, если $w_j \neq 0$.

Определение 3. Структурой \mathcal{A} модели \mathbf{f} назовем множество индексов активных параметров этой модели $\mathcal{A} = \{j : w_j \neq 0\} \subseteq \mathcal{J}$.

Каждая структура $\mathcal{A} \subseteq \mathcal{J}$ задает некоторую модель

$$\mathbf{f}_{\mathcal{A}} : \hat{\mathbf{w}}_{\mathcal{A}} \in \mathbb{R}^k,$$

где $\mathbf{f}_{\mathcal{A}}$ — модель со структурой \mathcal{A} , а $\hat{\mathbf{w}}_{\mathcal{A}} \in \mathbb{R}^k$ — оптимальный вектор параметров модели $\mathbf{f}_{\mathcal{A}}$, определение которому будет дано ниже. Объединение всех $\mathbf{f}_{\mathcal{A}}$ назовем множеством допустимых моделей

$$\mathfrak{F} = \bigcup_{\mathcal{A} \subseteq \mathcal{J}} \{\mathbf{f}_{\mathcal{A}}\}. \quad (3)$$

Оптимальную модель $\hat{\mathbf{f}}_{\mathcal{A}}$ будем выбирать из множества допустимых моделей \mathfrak{F} .

Согласно гипотезе (1) о распределении многомерных случайных величин \mathbf{y} и \mathbf{w} в качестве функции ошибки выберем функцию

$$S(\mathbf{w}|\mathcal{K}) = - \sum_{i \in \mathcal{K}} \sum_{\xi=1}^z t_{i\xi} \ln(f_{\xi}(\mathbf{x}_i, \mathbf{w})),$$

максимизирующую логарифм правдоподобия случайной величины \mathbf{y} и заданную на разбиении выборки \mathfrak{D} , определенном некоторым множеством индексов $\mathcal{K} \subseteq \mathcal{I}$, $\mathbf{t}_i = [t_{i1}, \dots, t_{i\xi}, \dots, t_{iz}]^T$.

Определение 4. Оптимальным вектором параметров модели $\mathbf{f}_{\mathcal{A}}$ назовем такой вектор $\hat{\mathbf{w}}_{\mathcal{A}}$, который является решением следующей задачи оптимизации:

$$\hat{\mathbf{w}}_{\mathcal{A}} = \arg \min_{\mathbf{w}_{\mathcal{A}} \in \mathbb{R}^k} S(\mathbf{w}_{\mathcal{A}}|\mathcal{L}). \quad (4)$$

Для оценки качества моделей и сравнения их друг с другом введем три критерия качества — сложность, устойчивость и точность.

Определение 5. Сложностью $C = C(\hat{\mathbf{w}})$ модели \mathbf{f} с вектором параметров $\hat{\mathbf{w}} = [w_1, \dots, w_k]$ назовем мощность множества активных параметров этой модели

$$C(\mathbf{w}) = \sum_{i=1}^k [w_i \neq 0] = |\mathcal{A}|.$$

Чем больше мощность множества активных параметров, тем сложнее модель. Максимально возможная сложность модели равна размерности пространства параметров k .

Определение 6. Устойчивостью $\eta = \eta(\hat{\mathbf{w}})$ модели \mathbf{f} с вектором параметров \mathbf{w} назовем число η , равное числу обусловленности матрицы \mathbf{A} (1), т. е.

$$\eta(\hat{\mathbf{w}}) = \frac{\lambda_{\max}}{\lambda_{\min}},$$

где λ_{\max} — максимальное, а λ_{\min} — минимальное собственное число матрицы \mathbf{A} .

Чем лучше обусловлена матрица \mathbf{A} , тем более устойчива модель. У абсолютно устойчивой модели $\lambda_{\min} = \lambda_{\max}$, $\eta = 1$.

Определение 7. Под точностью S модели \mathbf{f} с вектором параметров $\hat{\mathbf{w}}$ будем понимать величину функции ошибки (3) на контрольной выборке.

Чем больше значение функции ошибки, тем меньше точность модели.

Введем на множестве допустимых моделей \mathcal{F} отношение доминирования. Будем говорить, что модель \mathbf{f}' доминирует над моделью \mathbf{f} и обозначать $\mathbf{f}' \succ \mathbf{f}$, если

$$C' \leq C; \quad \eta' \leq \eta; \quad S' \leq S,$$

где C , η , S и C' , η' , S' — сложность, устойчивость и точность моделей \mathbf{f} и \mathbf{f}' .

Определение 8. Модель $\mathbf{f} \in \mathcal{F}$ назовем оптимальной по Парето, если не существует $\mathbf{f}' \in \mathcal{F}$ такой, что $\mathbf{f}' \succ \mathbf{f}$.

Определение 9. Множество оптимальных по Парето моделей назовем Парето-оптимальным фронтом $\text{POF}_{\mathfrak{F}}$ множества допустимых моделей \mathfrak{F} .

Задача выбора оптимальной модели состоит в том, чтобы найти Парето-оптимальный фронт $\text{POF}_{\mathfrak{F}}$ множества допустимых моделей \mathfrak{F} .

3 Стратегия пошаговой модификации модели

Определение 10. Стратегией пошаговой модификации модели называется процедура последовательного изменения модели, в которой на каждом шаге решается оптимизационная задача вида

$$\hat{j} = \arg \text{opt}_{j \in \mathcal{A}} Q(\hat{\mathbf{w}}_{\mathcal{A}}),$$

где Q — один из вышеприведенных критериев качества или их Парето-оптимальный набор.

Стратегия задается следующими математическими объектами:

- набором критериев оптимизации — сложностью, точностью, устойчивостью $\{C, S, \eta\}$,
- набором ограничений на структуру и параметры модели $\mathcal{A} \subseteq \mathcal{J}$, $\mathbf{w} = \hat{\mathbf{w}}_{\mathcal{A}}$ из (4),
- критериями останова шагов удаления (см. (9)) и добавления (см. (10)),
- критерием останова процедуры выбора модели (см. (11)).

Действуя согласно стратегии, будем изменять структуру модели, удаляя из нее элементы и добавляя их согласно (11).

Для определения индекса параметра \hat{j} , который должен быть удален из модели или добавлен в нее, ниже предлагается несколько критериев оптимизации модели.

3.1 Критерий оптимального прореживания

Этот критерий позволяет выяснить индекс параметра, удаление которого приведет к минимизации приращения функции ошибки (3). Для функции ошибки используется локальная аппроксимация

вблизи некоторого локального минимума вектора параметров \mathbf{w}_0 :

$$S(\mathbf{w}_0 + \Delta\mathbf{w}) = S(\mathbf{w}_0) + \mathbf{g}^T(\mathbf{w}_0)\Delta\mathbf{w} + \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H} \Delta\mathbf{w} + O(\|\Delta\mathbf{w}\|^3),$$

где $\Delta\mathbf{w}$ — возмущение вектора параметров в данной точке \mathbf{w}_0 ; $\mathbf{g}(\mathbf{w}_0)$ — вектор градиента, вычисленный в точке \mathbf{w}_0 ; $\mathbf{H} = \mathbf{H}(\mathbf{w}_0)$ — матрица вторых производных функции ошибки. Предполагается, что матрица вторых производных $\mathbf{H} = \mathbf{H}(\mathbf{w})$ — диагональная, а функция ошибки в окрестности глобального или локального минимума является квадратичной. На основании этих гипотез аппроксимация функции ошибки записывается в следующем виде:

$$\Delta S = S(\mathbf{w}_0 + \Delta\mathbf{w}) - S(\mathbf{w}_0) = \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H} \Delta\mathbf{w}.$$

Пусть w_j — некоторый параметр. Удаление этого параметра (присваивание ему нулевого значения) эквивалентно выполнению условия

$$\mathbf{e}_j^T \Delta\mathbf{w} + w_j = 0,$$

где \mathbf{e}_j^T — вектор, все элементы которого равны нулю, за исключением j -го, который равен единице. Таким образом, получаем задачу условной минимизации

$$\Delta S = \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H} \Delta\mathbf{w} \rightarrow \min; \quad \mathbf{e}_j^T \Delta\mathbf{w} + w_j = 0.$$

Для решения этой задачи строим лагранжиан

$$L = \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H} \Delta\mathbf{w} - \lambda_j (\mathbf{e}_j^T \Delta\mathbf{w} + w_j).$$

Продифференцировав L по $\Delta\mathbf{w}$, получаем значение лагранжиана L_j для элемента w_j :

$$L_j = \frac{w_j^2}{2[\mathbf{H}^{-1}]_{j,j}},$$

где \mathbf{H}^{-1} — матрица, обратная гессияну \mathbf{H} ; $[\mathbf{H}^{-1}]_{j,j}$ — j -й диагональный элемент этой матрицы. Значение лагранжиана L_j называется выпуклостью w_j . Выпуклость L_j описывает рост среднеквадратичной ошибки, вызываемый удалением параметра w_j .

Критерию оптимального прореживания отвечает параметр w_j , соответствующий минимальному значению выпуклости:

$$\hat{j} = \arg \min_{j \in \mathcal{A}} L_j. \quad (5)$$

3.2 Критерий последовательного прореживания

В качестве второго критерия предлагается простой критерий последовательного удаления параметров w_j — компонент вектора \mathbf{w} . Основной идеей этого критерия является принцип локально-оптимального выбора — критерию отвечает параметр w_j , без которого функция ошибки (3) оказывается минимальной.

Для нахождения параметра, отвечающего этому критерию, решается задача

$$\hat{j} = \arg \min_{j \in \mathcal{A}} S(\mathbf{w}_{\mathcal{A}} \setminus w_j | T). \quad (6)$$

3.3 Критерий устойчивого прореживания

Помимо вышеописанных критериев предлагается критерий устойчивого прореживания, основанный на модификации метода Белсли [17, 18].

Пусть \mathbf{W} — матрица реализаций оптимального вектора параметров $\hat{\mathbf{w}}$, определенного в (4) и рассматриваемого согласно (3) как многомерная случайная величина. Пусть эта матрица имеет размерность $r \times k$. Выполним ее сингулярное разложение:

$$\mathbf{W} = \mathbf{U} \mathbf{S} \mathbf{V}^T, \quad (7)$$

где \mathbf{U} и \mathbf{V} — ортогональные матрицы размера $r \times r$ и $k \times k$, при этом r — количество оценок, а k — размерность вектора параметров \mathbf{w} ; $\mathbf{\Lambda}$ — матрица, на диагонали которой стоят сингулярные числа матрицы \mathbf{W} .

По определению ковариационная матрица вектора параметров \mathbf{w} вычисляется как

$$\mathbf{A}^{-1} = \text{cov}(\mathbf{W}) = \mathbf{E}(\mathbf{W}^T \mathbf{W}) - \mathbf{E}(\mathbf{W})\mathbf{E}(\mathbf{W}^T) = \mathbf{E}(\mathbf{W}^T \mathbf{W}).$$

Последнее равенство выполняется в силу предположения о том, что математическое ожидание вектора параметров равно нулю: $\mathbf{E}(\mathbf{w}) = \mathbf{0}$. По матрице реализаций \mathbf{W} многомерной случайной величины \mathbf{w} ковариационная матрица может быть оценена следующим образом:

$$\mathbf{A}^{-1} = \frac{1}{r} \mathbf{W} \mathbf{W}^T.$$

У ковариационной матрицы есть нулевые строки с индексами из множества $\mathcal{J} \setminus \mathcal{A}$, где \mathcal{J} — множество индексов всех параметров модели, а \mathcal{A} — множество индексов активных параметров. Таким образом, ковариационная матрица является неполноранговой.

Используя сингулярное разложение (7) матрицы \mathbf{W} , получим выражение для матрицы \mathbf{A}^{-1} :

$$\begin{aligned}\mathbf{A}^{-1} &= (\mathbf{W}\mathbf{W}^T) = (\mathbf{U}\mathbf{\Lambda}\mathbf{V}^T\mathbf{V}\mathbf{\Lambda}^T\mathbf{U}^T) = \\ &= (\mathbf{U}\mathbf{\Lambda}\mathbf{\Lambda}^T\mathbf{U}^T) = \mathbf{U}\mathbf{\Lambda}^2\mathbf{U}^T.\end{aligned}$$

Индексом обусловленности η_ζ назовем отношение максимального элемента λ_{\max} матрицы $\mathbf{\Lambda}$ к ζ -му по величине элементу λ_ζ этой матрицы:

$$\eta_\zeta = \frac{\lambda_{\max}}{\lambda_\zeta}.$$

Так как ковариационная матрица \mathbf{A}^{-1} неположительная, то некоторые значения индексов обусловленности не определены. Чтобы избежать этой проблемы, исключим из рассмотрения параметры с дисперсией, меньшей некоторого порога α , и добавим к каждому элементу, стоящему на диагонали ковариационной матрицы, небольшое число τ .

Оценками дисперсии параметров будут диагональные элементы \mathbf{A}^{-1} :

$$\sigma(w_\zeta) = \mathbf{A}_{\zeta\zeta}^{-1}.$$

Долевой коэффициент $q_{\zeta j}$ определим как вклад j -го признака в дисперсию ζ -го элемента вектора параметров \mathbf{w} :

$$q_{\zeta j} = \frac{u_{\zeta j}^2 \lambda_{jj}^2}{\sigma(w_\zeta)}.$$

Находим индексы обусловленности и долевы коэффициенты для набора активных параметров \mathcal{A} . Большие значения индексов обусловленности указывают на зависимость между признаками. Поэтому для нахождения параметра, отвечающего этому критерию прорезивания, находим максимальный индекс обусловленности

$$\hat{\zeta} = \operatorname{argmax}_{\zeta \in \mathcal{A}} \eta_\zeta.$$

Затем находим максимальный долевы коэффициент, соответствующий найденному максимальному индексу обусловленности $\eta_{\hat{\zeta}}$:

$$\hat{j} = \operatorname{argmax}_{j \in \mathcal{A}} q_{\hat{\zeta} j}. \quad (8)$$

Параметр $w_{\hat{j}}$ и есть параметр, отвечающий критерию устойчивого прорезивания.

3.4 Критерий последовательного наращивания

Критерий последовательного добавления параметров, как и критерий (6), основан на принципе

локально-оптимального выбора — критерию отвечает параметр, при добавлении которого в сеть функция ошибки (3) минимальна.

Для нахождения параметра, отвечающего этому критерию, решается задача

$$\hat{j} = \operatorname{argmin}_{j \in \mathcal{J} \setminus \mathcal{A}} S(\mathbf{w}_{\mathcal{A}} \cup w_j | T).$$

3.5 Описание базовой стратегии

Стратегия пошаговой модификации модели состоит из двух этапов — Del и Add. Перед началом процедуры модификации все параметры модели активны.

Этап Del. Ищем параметр с индексом \hat{j} , отвечающий одному из критериев прорезивания (5), (6) или (7), и удаляем его из множества активных параметров:

$$\mathcal{A} = \mathcal{A} \setminus \hat{j}.$$

Этап Del повторяем до тех пор, пока ошибка $S(\mathbf{w}_{\mathcal{A}} | T)$ не превысит свое минимальное значение на данном этапе более чем на некоторое заданное значение δS_1 . Критерием останова шага Del является следующее условие:

$$S(\hat{\mathbf{w}}_{\mathcal{A}} | T) \geq S_{\min} + \delta S_1, \quad (9)$$

где S_{\min} — некоторое заданное значение.

Этап Add. В модели ищем параметр \hat{j} , отвечающий критерию наращивания (8), и добавляем найденный параметр во множество активных параметров:

$$\mathcal{A} = \mathcal{A} \cup \hat{j}.$$

Критерием останова шага Add является выполнение условия

$$S(\hat{\mathbf{w}}_{\mathcal{A}} | T) \geq S_{\min} + \delta S_2, \quad (10)$$

где S_{\min} — некоторое заданное значение. На рис. 1 приведен график, демонстрирующий изменение функции ошибки при удалении параметров из модели. Аналогичным образом ведет себя функция ошибки при добавлении параметров в модель. Из графика видно, что эта зависимость имеет минимум, а значит модели с большим числом параметров не являются наиболее точными. На рис. 2 показано, как согласно критериям останова (9) и (10) сменяются шаги удаления и добавления.

Процедура модификации продолжается до тех пор, пока процесс не стабилизируется. В качестве критерия стабилизации предлагается использовать энтропию изменения структуры модели:

$$H(\mathcal{A}, \mathcal{A}') = - \sum_{j=1}^k \rho(a_j, a'_j) \ln(\rho(a_j, a'_j)) \quad (11)$$

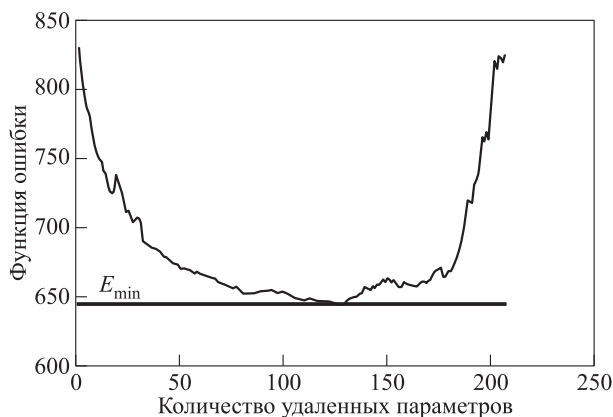


Рис. 1 Изменение функции ошибки при удалении параметров из модели

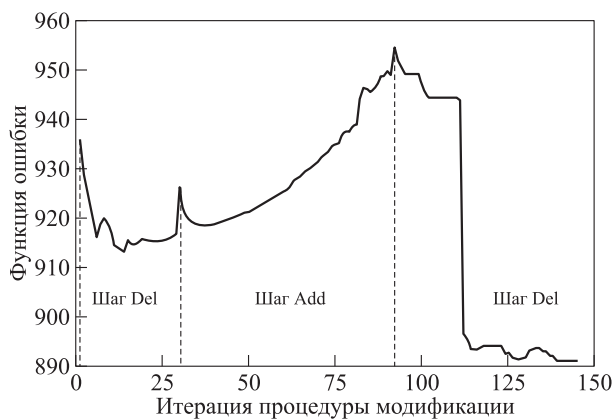


Рис. 2 Смена шагов Del и Add

множества попарных нормированных расстояний Хэмминга между элементами наборов $\mathcal{A} = \{a_1, \dots, a_k\}$ и $\mathcal{A}' = \{a'_1, \dots, a'_k\}$, полученных на двух последовательных итерациях алгоритма следующим образом:

$$a_j = \begin{cases} 1, & \text{если } w_j \neq 0; \\ 0, & \text{если } w_j = 0. \end{cases}$$

Процесс считается стабильным, если энтропия $H(\mathcal{A}, \mathcal{A}')$ не превосходит заданного порога.

4 Путь в k -мерном кубе

В данной задаче будем иметь дело с вектором параметров размерности k . Это означает, что существует 2^k вариантов структуры модели. Из этих 2^k возможных вариантов структуры выбираются оптимальные. Все варианты можно представить в виде вершин k -мерного куба \mathcal{V} . И тогда стратегия задает путь \mathbf{V} по его вершинам. Этот путь

заканчивается в некоторой вершине \hat{v} , к которой сходится процедура модификации. Будем искать оптимальные модели в некоторой окрестности вершины \hat{v} . Так как охватить все возможные варианты слишком трудоемко, то в качестве окрестности \hat{v} будем рассматривать ведущий к \hat{v} путь по вершинам куба, полученный по описанной выше стратегии.

Пример 1. В этом примере использовалась выборка $\{x_i, y_i\}, i \in \{1, \dots, 177\}$. Каждый объект выборки описывался 6 признаками χ_1, \dots, χ_6 и принадлежал одному из трех классов. Схематично взаимное расположение векторов χ_1, \dots, χ_6 изображено на рис. 3.

Для классификации такой выборки модифицировалась двухслойная нейронная сеть с одним нейроном в скрытом слое. Совокупное число параметров такой сети равно девяти. Нейронная сеть модифицировалась за 11 итераций. На рис. 4 изображен путь по вершинам девятимерного куба. По вертикали отложен номер параметра, по горизонтали — номер итерации. Черная клетка означает, что параметр с индексом j — активный, белая клетка — параметр неактивный. Например, на пятой итера-

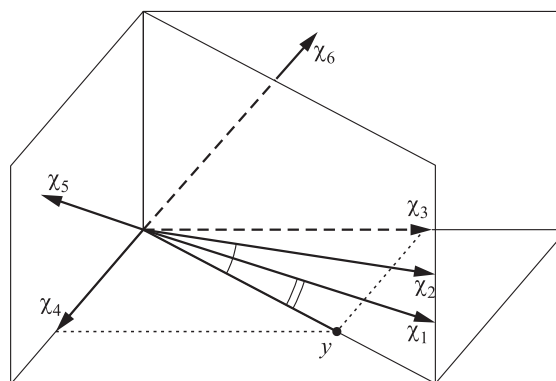


Рис. 3 Данные

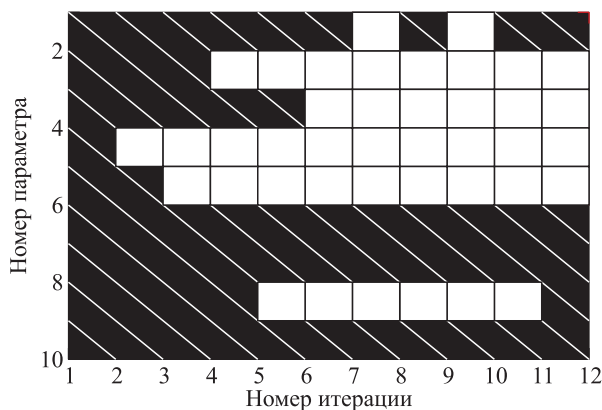


Рис. 4 Путь в кубе

ции из сети был удален параметр 9, а на одиннадцатой итерации этот параметр был снова добавлен в сеть.

5 Вычислительный эксперимент

С целью получить значение критериев качества описанной стратегии был проведен вычислительный эксперимент. Использовались данные акселерометра мобильного телефона. Данные состояли из 5418 векторов признаков, которые были получены в результате обработки соответствующих временных рядов. Было выделено 43 признака и 6 классов физической активности: ходьба, бег, сидение, стояние, подъем и спуск. Временные ряды записывались акселерометром мобильного телефона, который находился в кармане у человека, выполняющего один из типов физической активности. Для выделения признаков временные ряды разделялись на десятисекундные сегменты. Из этих сегментов извлекались признаки, такие как проекции среднего ускорения на координатные оси, среднеквадратические отклонения от проекций среднего ускорения на каждую из трех координатных осей, время между пиками синусоидального сигнала в миллисекундах и др. С более подробным описанием признаков и процессом их генерации можно ознакомиться в [16].

В вычислительном эксперименте оптимизировалась двухслойная нейронная сеть с пятью нейронами в скрытом слое. Размерность вектора параметров такой модели $k = 245$. Нейронная сеть оптимизировалась по стратегии, описанной в разд. 3. Был получен набор из 771 модели. В процедуре модификации использовался каждый из трех критериев прореживания — оптимального, последовательного и устойчивого. Для всех моделей были вычислены значения критериев качества. Был построен Парето-оптимальный фронт трех критериев. На рис. 5 изображены все полученные модели. Пустыми значками обозначены модели, которые были получены по стратегии с применением критерия устойчивого прореживания, серыми значками — критерия последовательного прореживания, черными значками — оптимального прореживания. Парето-оптимальные модели обозначены черными крестиками. Из рис. 5, *a* видно, что самые устойчивые модели получаются при использовании критерия устойчивого прореживания. В таблице приведены значения критериев качества моделей, которые являются точками останова процедуры модификации для каждого из трех критериев прореживания.

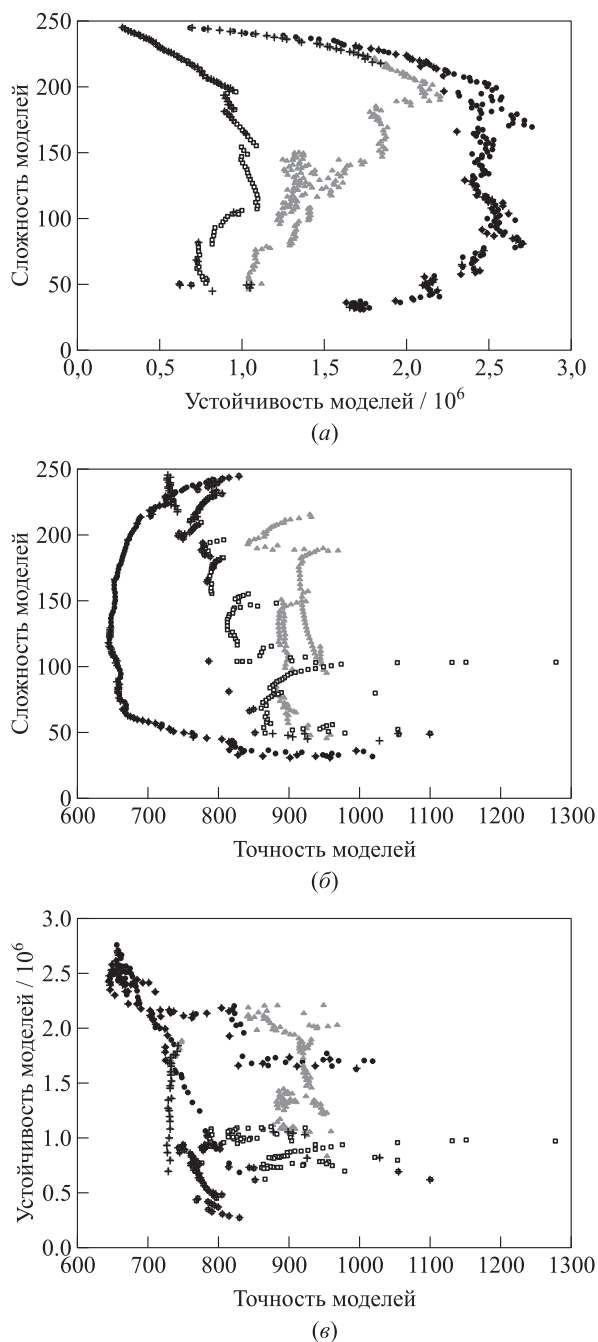


Рис. 5 Множество моделей в координатах «устойчивость—сложность» (*a*), «точность—сложность» (*b*) и «точность—устойчивость» (*v*)

На рис. 6 приведена интерпретация полученных результатов. В верхней области графика Парето-оптимальные модели не интересны для рассмотрения, так как в этой области имеет место недообучение — модели излишне сложны. Парето-оптимальные модели с незначительной сложностью находятся в нижней области графика.

Сложность, точность и устойчивость моделей

Стратегия	Сложность	Точность	Устойчивость
Оптимальное прореживание	50	877	$1,2 \cdot 10^6$
Последовательное прореживание	36	870	$2,0 \cdot 10^6$
Устойчивое прореживание	50	866	$6, \cdot 10^5$

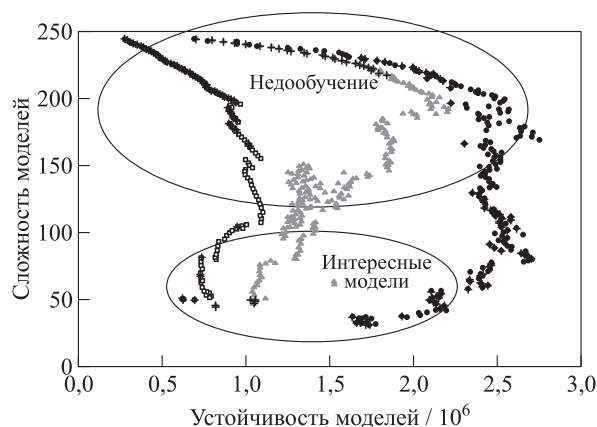
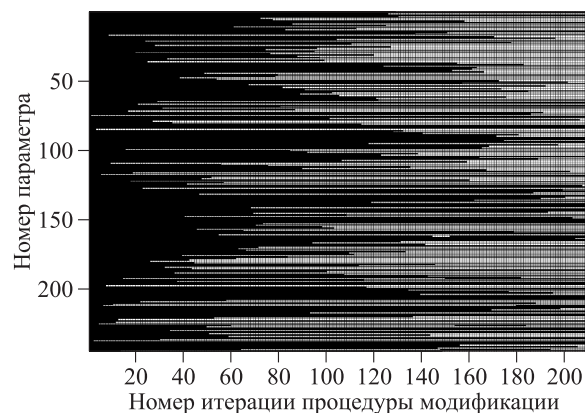


Рис. 6 Интерпретация результатов

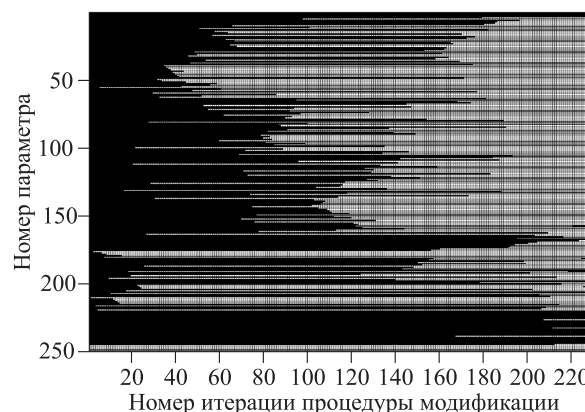
Также была визуализирована процедура пошаговой модификации модели как путь в k -мерном кубе. На рис. 7, так же как и в примере 1, по вертикали отложен номер параметра, по горизонтали — номер итерации. Черная клетка означает, что параметр активный, белая клетка — параметр неактивный. На рис. 6, 7, а и 7, б указана последовательность, в которой параметры удалялись из модели и добавлялись в нее. Из рис. 7, б и 7, в видно, что стратегия с критериями оптимального и последовательного прореживания, которые выбирают для удаления параметр, минимизирующий функцию ошибки, оставляет в моделях параметры с номерами с 216 по 245. Это связано с тем, что параметры с такими номерами относятся ко второму слою нейронной сети, а удаление большого числа параметров второго слоя приводит к росту функции ошибки.

6 Заключение

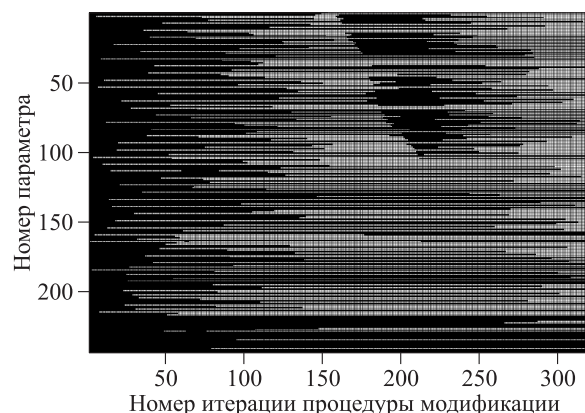
В работе была предложена стратегия пошаговой модификации моделей классификации согласно трем критериям качества — сложности, точности и устойчивости. В рамках стратегии были предложены критерии добавления и удаления параметров в модель, критерии останова шагов добавления и удаления, а также критерий останова процедуры модификации. Процедура пошаговой модификации модели была рассмотрена и визуализи-



(а)



(б)



(в)

Рис. 7 Путь в кубе: устойчивое (а); последовательное (б) и оптимальное (в) прореживания

зирована как путь в многомерном кубе. Был проведен вычислительный эксперимент, в ходе которого был получен набор моделей и найден Парето-оптимальный фронт критериев качества этого набора. Вычислительный эксперимент показал, что наилучшие по рассматриваемым критериям качества модели получаются при использовании критерия устойчивого прореживания. Это связано с тем, что критерий устойчивого прореживания позволяет получать более устойчивые модели, удаляя коррелирующие параметры и тем самым повышая устойчивость и обобщающую способность модели классификации. Программная реализация стратегии пошаговой модификации нейронной сети в среде разработки MatLab находится в свободном доступе [19].

Литература

1. Визильтер Ю. В., Горбацевич В. С., Карамеев С. Л., Костромов Н. А. Обучение алгоритмов выделения кожи на цветных изображениях лиц // Информатика и её применения, 2012. Т. 6. Вып. 1. С. 109–113.
2. Токмакова А. А., Стрижов В. В. Оценка гиперпараметров линейных и регрессионных моделей при отборе шумовых и коррелирующих признаков // Информатика и её применения, 2012. Т. 6. Вып. 4. С. 66–75.
3. Хапланов А. Ю. Асимптотическая нормальность оценки параметров многомерной логистической регрессии // Информатика и её применения, 2013. Т. 7. Вып. 2. С. 69–74.
4. Myung I. J. The importance of complexity in model selection // J. Math. Psychol., 2000. Vol. 44. No. 1. P. 190–204.
5. MacLeod C., Maxwell M. Incremental evolution in ANNs: Neural nets which grow // Artif. Intell. Rev., 2001. Vol. 16. No. 3. P. 201–224.
6. Karnin E. D. A simple procedure for pruning backpropagation trained neural networks // IEEE Trans. Neural Networks, 1990. Vol. 1. No. 2. P. 239–242.
7. LeCun Y., Denker L. S., Solla S. A. Optimal brain damage // Adv. Neur. Inform. Processing Syst., 1990. Vol. 2. No. 2. P. 598–605.
8. Hassibi B., Stork D. G., Woff G. J. Optimal brain surgeon and general network pruning // IEEE Conference (International) on Neural Networks Proceedings, 1993. Vol. 1. P. 293–299.
9. Hong-Gui H., Qi-li C., Jun-Fei Q. An efficient self-organizing RBF neural network for water quality prediction // Neural Networks, 2011. Vol. 24. No. 7. P. 717–725.
10. Yang S., Chen Y. An evolutionary constructive and pruning algorithm for artificial neural networks and its prediction applications // Neurocomputing, 2012. Vol. 86. P. 140–149.
11. Pu X., Pengfei Sun P. A new hybrid pruning neural network algorithm based on sensitivity analysis for stock market forecast // J. Inform. Comput. Sci., 2013. Vol. 3. P. 883–892.
12. Knerr S., Personnaz L., Dreyfus G. Single-layer learning revisited: A stepwise procedure for building and training a neural network // Neurocomputing Algorithms Architectures Applications, 1990. Vol. 68. No. 1. P. 41–50.
13. Strijov V., Krymova E., Weber S. V. Evidence optimization for consequently generated models // Math. Comput. Modell., 2010. Vol. 57. No. 1–2. P. 50–56.
14. Леонтьева Л. Н. Последовательный выбор признаков при восстановлении регрессии // Машинное обучение и анализ данных, 2012. Т. 1. № 3. С. 335–346.
15. Зайцев А. А., Токмакова А. А. Оценка гиперпараметров линейных регрессионных моделей методом максимального правдоподобия при отборе шумовых и коррелирующих признаков // Машинное обучение и анализ данных, 2012. Т. 1. № 3. С. 347–353.
16. Kwapisz J. R., Weiss G. M., Moore S. Activity recognition using cell phone accelerometers // SIGKDD Explorations, 2010. Vol. 12. No 2. P. 74–82.
17. Belsley D. A., Kuh E., Welsch R. E. Regression diagnostics: Identifying influential data and sources of collinearity. — New York: John Wiley and Sons, 2005. 302 p.
18. Сандуляну Л. Н., Стрижов В. В. Выбор признаков в авторегрессионных задачах прогнозирования // Информационные технологии, 2012. Т. 7. С. 11–15.
19. Попова М. С. Реализация стратегии пошаговой модификации нейронной сети // Algorithms Machine Learning, 2014. <http://sourceforge.net/p/mlalgorithms/code/HEAD/tree/Group174/Popova2014OptimalModeISelection/code/main.m>.

Поступила в редакцию 10.08.14

SELECTION OF OPTIMAL PHYSICAL ACTIVITY CLASSIFICATION MODEL USING MEASUREMENTS OF ACCELEROMETER

M. Popova¹ and V. Strijov²

¹Moscow Institute of Physics and Technology, 9 Institutskiy Per., Dolgoprudny, Moscow Region 141700, Russian Federation

²Dorodnicyn Computing Center, Russian Academy of Sciences, 40 Vavilov Str., Moscow 119333, Russian Federation

Abstract: The paper solves the problem of selecting optimal stable models for classification of physical activity. Each type of physical activity of a particular person is described by a set of features generated from an accelerometer time series. In conditions of feature's multicollinearity, selection of stable models is hampered by the need to evaluate a large number of parameters of these models. Evaluation of optimal parameter values is also difficult due to the fact that the error function has a large number of local minima in the parameter space. In the paper, the optimal models from the class of two-layer artificial neural networks are chosen. The problem of finding the Pareto optimal front of the set of models is solved. The paper presents a stepwise strategy of building optimal stable models. The strategy includes steps of deleting and adding parameters, criteria of pruning and growing the model and criteria of breaking the process of building. The computational experiment compares the models generated by the proposed strategy on three quality criteria — complexity, accuracy, and stability.

Keywords: classification; artificial neural networks; complexity; accuracy; stability; Pareto efficiency; growing and pruning criteria

DOI: 10.14357/19922264150107

Acknowledgments

The research was supported by Skolkovo Institute of Science and Technology (Skoltech) in the frame of SkolTech/MITInitiative.

References

1. Vizilter, Y., V. Gorbachevich, S. Karateev, and N. Kostromov. 2012. Obuchenie algoritmov vydeleniya kozhni na tsvetnykh izobrazheniyakh lits [Teaching of skin extraction algorithms for human face color images]. *Informatika i ee Primeneniya — Inform. Appl.* 6(1):109–113.
2. Tokmakova, A. A., and V. V. Strizhov. 2012. Otsenivanie giperparametrov lineynykh i regressionnykh modeley pri otbore shumovykh i korrelirovannykh priznakov [Estimation of linear model hyperparameters for noise or correlated feature selection problem]. *Informatika i ee Primeneniya — Inform. Appl.* 6(4):66–75.
3. Khaplanov, A. Yu. 2013. Asimptoticheskaya normal'nost' otsenki parametrov mnogomernoy logisticheskoy regressii [Asymptotic normality of the estimation of the multivariate logistic regression]. *Informatika i ee Primeneniya — Inform. Appl.* 7(2):69–74.
4. Myung, I. J. 2000. The importance of complexity in model selection. *J. Math. Psychol.* 44(1):190–204.
5. MacLeod, C., and M. Maxwell. 2001. Incremental evolution in ANNs: Neural nets which grow. *Artif. Intell. Rev.* 16(3):201–224.
6. Karnin, E. D. 1990. A simple procedure for pruning back-propagation trained neural networks. *IEEE Trans. Neural Networks* 1(2):239–242.
7. LeCun, Y., L. S. Denker, and S. A. Solla. 1990. Optimal brain damage. *Adv. Neur. Inform. Processing Syst.* 2(2):598–605.
8. Hassibi, B., D. G. Stork, and G. J. Woff. 1993. Optimal brain surgeon and general network pruning. *IEEE Conference (International) on Neural Networks Proceedings.* 293–299.
9. Hong-Gui, H., C. Qi-li, and Q. Jun-Fei. 2011. An efficient self-organizing RBF neural network for water quality prediction. *Neural Networks* 24(7):717–725.
10. Yang, S., and Y. Chen. 2012. An evolutionary constructive and pruning algorithm for artificial neural networks and its prediction applications. *Neurocomputing* 86(1):140–149.
11. Pu, X., and P. Pengfei-Sun. 2013. A new hybrid pruning neural network algorithm based on sensitivity analysis for stock market forecast. *J. Inform. Comput. Sci.* 3(1):883–892.
12. Knerr, S., L. Personnaz, and G. Dreyfus. 1990. Single-layer learning revisited: A stepwise procedure for building

- and training a neural network. *Neurocomputing Algorithms Architectures Applications* 68(1):41–50.
13. Strijov, V., E. Krymova, and S. Weber. 2013. Evidence optimization for consequently generated models. *Math. Comput. Modell.* 57(1-2):50–56.
 14. Leont'eva, L. N. 2012. Posledovatel'nyy vybor priznakov pri vosstanovlenii regressii [Feature selection in autoregression forecasting]. *J. Machine Learning Data Analysis* 1(3):335–346.
 15. Zaytsev, A. A., and A. A. Tokmakova. 2012. Otsenka giperparametrov lineynykh regressionnykh modeley metodom maksimal'nogo pravdopodobiya pri otbore shumovykh i korreliuyushchikh priznakov [Estimation regression model hyperparameters using maximum likelihood]. *J. Machine Learning Data Analysis* 1(3):347–353.
 16. Kwapisz, J. R., G. M. Weiss, and S. Moore. 2010. Activity recognition using cell phone accelerometers. *SIGKDD Explorations* 12(2):74–82.
 17. Belsley, D. A., E. Kuh, R. E. Welsch. 2005. *Regression diagnostics: Identifying influential data and sources of collinearity*. New York: John Wiley and Sons. 302 p.
 18. Sanduljanu, L. N., and V. V. Strizhov. 2012. Vybor priznakov v avtoregressionnykh zadachakh prognozirovaniya [Feature selection in autoregression forecasting]. *Information Technologies* 7:11–15.
 19. Popova, M. S. 2014. Realizatsiya strategii poshagovoy modifikatsii neyronnoy seti [Realization of a stepwise strategy for neural network modification]. Available at: <http://sourceforge.net/p/mlalgorithms/code/HEAD/tree/Group174/Popova2014OptimalModelSelection/code/main.m> (accessed February 10, 2015).

Received August 10, 2014

Contributors

Popova Maria S. (b. 1994) — student, Moscow Institute of Physics and Technology, 9 Institutskiy Per., Dolgoprudny, Moscow Region 141700, Russian Federation; maria_popova@phystech.edu

Strijov Vadim V. (b. 1967) — Candidate of science (PhD) in physics and mathematics; associate professor, Moscow Institute of Physics and Technology, 9 Institutskiy Per., Dolgoprudny, Moscow Region 141700, Russian Federation; leading scientist, Dorodnicyn Computing Center, Russian Academy of Sciences, 40 Vavilov Str., Moscow 119333, Russian Federation; strijov@ccas.com

ОЦЕНКА ПОГРЕШНОСТИ И ЗНАЧИМОСТИ ИЗМЕРЕНИЙ ДЛЯ ЛИНЕЙНЫХ МОДЕЛЕЙ*

С. И. Спивак¹, О. Г. Кантор², Д. С. Юнусова³, С. И. Кузнецов⁴, С. В. Колесов⁵

Аннотация: Решение задач восстановления линейных зависимостей в тех случаях, когда точное решение, полученное стандартными методами, не удовлетворяет объективным требованиям, обуславливает разработку специальных подходов для их численной реализации. В статье приводится описание метода получения приближенных значений параметров линейных зависимостей по экспериментальным данным, в основе которого лежит использование методологии линейного программирования и теории двойственности. Разработанный метод позволяет не только получать приближенные решения, обеспечивающие выполнение всех предъявляемых требований к самой восстанавливаемой зависимости и ее параметрам, но и проводить оценку погрешности измерений и их значимости. А это важно для совершенствования процедуры построения функциональных зависимостей на стадии планирования экспериментов в части уточнения экспериментальных данных или их исключения из рассмотрения как не удовлетворяющих критериям достоверности. Приведены результаты апробации предложенного метода для задач, связанных с исследованиями химических и социально-экономических систем.

Ключевые слова: задачи восстановления линейных зависимостей; погрешность измерений; значимость измерений; двойственные оценки

DOI: 10.14357/19922264150108

1 Введение

Рассматривается задача определения параметров линейных математических моделей по экспериментальным данным, которые не могут быть рассчитаны стандартными методами в силу некоторых объективных причин (например, в виду ограниченного количества имеющихся измерений или отсутствия информации об их статистических характеристиках). Предметом рассмотрения в данной работе являются системы линейных алгебраических уравнений вида

$$AX = B. \quad (1)$$

В системе (1) $A = (a_{ij})$ и $B = (b_i)$ — экспериментальные данные ($i = 1, \dots, m, j = 1, \dots, n$), а $X = (x_1, x_2, \dots, x_n)^T$ — искомые параметры модели.

К решению таких задач сводятся многие задачи восстановления линейных зависимостей по экспериментальным данным, возникающие при исследованиях в различных областях научной и практической деятельности. Примеры задач подобного рода широко представлены в рамках таких хорошо

изученных разделов, как определение регрессионных зависимостей и моделирование временных рядов. Основу подавляющего большинства подходов к решению задачи (1) составляют методы математической статистики, представляющие на сегодняшний день группу подробно изученных и хорошо зарекомендовавших себя на практике инструментов определения параметров линейных зависимостей [1–3].

Существенным препятствием к применению статистических методов является обязательное наличие достаточно большого количества экспериментальных данных, что не всегда достижимо на практике.

При равенстве числа наблюдений и числа оцениваемых параметров могут быть использованы и классические методы решения квадратных систем линейных уравнений (метод Гаусса, метод обратной матрицы и пр.). Однако говорить о статистической значимости искомых параметров в этом случае нельзя, и основное назначение такого подхода заключается в установлении точного вида функциональной связи между исследуемыми величинами по результатам конкретных наблюдений.

* Работа выполнена при поддержке РФФИ (проект 13-01-00749).

¹ Башкирский государственный университет, semen.spivak@mail.ru

² Институт социально-экономических исследований Уфимского научного центра Российской академии наук, o_kantor@mail.ru

³ Башкирский государственный университет, kazakova_d_s@mail.ru

⁴ Институт органической химии Уфимского научного центра Российской академии наук, chemorg@anrb.ru

⁵ Институт органической химии Уфимского научного центра Российской академии наук, kolesovservic@rambler.ru

Применение большинства методов восстановления линейных зависимостей по экспериментальным данным (в том числе и перечисленных выше) сопряжено с еще одной группой проблем: найденные значения параметров могут не удовлетворять некоторым условиям, вытекающим из их физического смысла (например, применение статистических методов или метод решения линейных алгебраических уравнений не обеспечивают неотрицательность искомых величин).

В этой связи актуальным является определение приближенного решения системы (1) и оценка величины погрешности измерений, под которой будем понимать расхождение значений расчетных и экспериментальных величин не в каждом отдельном наблюдении, а в целом по всей совокупности наблюдений. В свою очередь, это обуславливает необходимость изучения способов формализации таких задач.

2 Описание подхода к определению погрешности измерений

Без ограничения общности рассуждений будем предполагать, что на параметры модели X наложены условия неотрицательности:

$$X \geq 0. \quad (2)$$

Достаточно часто на значения параметров накладываются ограничения, выражающие их принадлежность какому-либо множеству значений, поэтому будем считать, что параметры модели также удовлетворяют системе ограничений

$$CX \geq D, \quad (3)$$

где C — это матрица, состоящая из коэффициентов при параметрах модели в системе ограничений; D — концы промежутков значений, которым принадлежат параметры модели.

Тогда задача определения приближенного решения системы (1) с учетом ограничений (2) и (3) может быть сведена к задаче линейного программирования

$$\left. \begin{array}{l} \varepsilon \rightarrow \min; \\ |AX - B| \leq \varepsilon; \\ CX \geq D; X \geq 0. \end{array} \right\} \quad (4)$$

Здесь $A = (a_{ij})$ и $B = (b_i)$ — экспериментальные данные ($i = 1, \dots, m$, $j = 1, \dots, n$); $X = (x_1, x_2, \dots, x_n)^T$ — искомые параметры модели;

$C = (c_{lj})$ — это матрица, состоящая из коэффициентов при параметрах модели в системе ограничений ($l = 1, \dots, k$, $j = 1, \dots, n$), $D = (d_l)$ — вектор-столбец, элементы которого — концы промежутков значений параметров модели ($l = 1, \dots, k$); ε — параметр, характеризующий величину погрешности измерений.

Отметим, что вместо ограничения $|AX - B| \leq \varepsilon$ может использоваться условие вида $|AX - B| = E$, в котором элементы матрицы $E = (\varepsilon_i)$ ($i = 1, \dots, m$) представляют собой параметры, характеризующие величину ошибки в описании i -го эксперимента. Основная сложность такого перехода от неравенств к равенствам сопряжена с ростом числа неизвестных: вместо одной неизвестной величины ε определению подлежат m величин ε_i . При этом, если ввести обозначение $\varepsilon = \max_i \varepsilon_i$, легко убедиться, что задача в постановке (4) более предпочтительна не только по причине меньшего числа неизвестных, но и в силу того, что разность между AX и B , т. е. левые части ограничений $|AX - B| \leq \varepsilon$, позволяют определить величины ε_i ($i = 1, \dots, m$).

В результате решения задачи (4) должны быть определены параметры модели, удовлетворяющие требуемым ограничениям, и величина погрешности измерений.

Следует отметить, что матрицы A и B формируются по результатам наблюдений, а потому не могут рассматриваться как абсолютно точные, так как результаты их измерений неминуемо сопряжены с некоторыми ошибками. Такие ошибки приводят к отклонениям измеряемых значений величин A и B от их истинных значений. Причем очевидно, что исследователю эти отклонения заранее не известны, но, предполагая их наличие, целесообразно ставить задачу определения параметров модели, удовлетворяющих ограничениям (2) и (3), при условии внесения изменений в экспериментальные данные A и B . Будем предполагать незначительные отклонения измеряемых значений величин A и B от их истинных значений, что соответствует ситуации, при которой грубые ошибки при получении экспериментальной информации исключаются. В данных условиях целесообразно рассмотреть способы постановки задач определения параметров модели, обеспечивающих минимально возможные отклонения от экспериментальных данных, для различных случаев вариации матриц A и B .

В случае предполагаемых ошибок в матрице B , различных для каждого отдельного измерения $i = 1, \dots, m$, задача определения параметров модели, обеспечивающих минимальное отклонение от экспериментальных величин (b_i) может быть сведена к задаче линейного программирования:

$$\left. \begin{aligned} \varepsilon &\rightarrow \min; \\ AX &= \Delta B; \\ CX &\geq D; X \geq 0; \\ |\delta_i - 1| &\leq \varepsilon \quad \forall i = 1, \dots, m. \end{aligned} \right\} \quad (5)$$

Здесь ε — это величина погрешности измерений, а Δ — диагональная матрица вида

$$\Delta = \begin{pmatrix} \delta_1 & \dots & 0 \\ \dots & \dots & \dots \\ 0 & \dots & \delta_m \end{pmatrix},$$

где δ_i ($i = 1, \dots, m$) — это неизвестные величины, отождествляемые с параметрами, характеризующими ошибки в измерениях элементов матрицы B . Величины ε и δ_i связаны очевидным соотношением $\varepsilon = \max_i |\delta_i - 1|$.

Решение задачи (5) обеспечивает определение неотрицательных параметров X , при которых наибольшая погрешность измерений элементов матрицы B минимальна, поскольку каждый элемент этой матрицы умножается на число, отличающееся от единицы не больше чем на ε в силу условия $|\delta_i - 1| \leq \varepsilon \quad \forall i = 1, \dots, m$. При такой постановке разницы между левыми и правыми частями системы (1) соответственно равны $(1 - \delta_i)b_i$. Следовательно, справедливы соотношения $|AX - B| \leq \varepsilon$, а это означает, что для системы (1) максимальная разница правых и левых частей по модулю не превысит ε .

В случае предполагаемых ошибок в каждом элементе матрицы A , что отражает случайный характер ошибок в каждом измерении экспериментальных данных $A = (a_{ij})$ и формально соответствует умножению каждого элемента матрицы на некоторое число, задача (1) может быть формализована в следующем виде:

$$\left. \begin{aligned} \varepsilon &\rightarrow \min; \\ A'X &= B; \\ CX &\geq D; X \geq 0; \\ |\gamma_{ij} - 1| &\leq \varepsilon \quad \forall i = 1, \dots, m, \quad \forall j = 1, \dots, n. \end{aligned} \right\} \quad (6)$$

Здесь матрица A' имеет вид:

$$A' = \begin{pmatrix} \gamma_{11}a_{11} & \dots & \gamma_{1n}a_{1n} \\ \dots & \dots & \dots \\ \gamma_{m1}a_{m1} & \dots & \gamma_{mn}a_{mn} \end{pmatrix}.$$

Задача (6) является нелинейной, поскольку в выражении $A'X = B$ элементы вектора X — неизвестные величины, а A' — матрица, которая зависит от $m \times n$ неизвестных величин. Данное обстоятельство создает существенные проблемы для численной реализации модели (6) и не является предметом

рассмотрения данной работы. Вместе с тем относительно элементов матрицы A может быть известна информация, которая позволит упростить процесс получения решения. Примером тому могут служить ситуации, при которых либо во всех наблюдениях применительно к каждой введенной в рассмотрение величине совершаются однотипные ошибки, либо одни и те же погрешности допускаются в рамках каждого наблюдения и по отношению ко всем рассматриваемым величинам. Первая ситуация отражает наличие индивидуальных систематических ошибок при измерении каждой из наблюдаемых величин, а вторая — индивидуальные ошибки каждого отдельного наблюдения. Первая ситуация может возникать, например, ввиду особенностей приборов, используемых для фиксации значений наблюдаемых величин, а вторая — в случае зависимости от условий проведения эксперимента (температурных, временных и пр.).

Первая из перечисленных выше ситуаций отражается в пропорциональных изменениях всех элементов каждого из столбцов, а вторая — строк. Формально первая ситуация соответствует умножению матрицы A на диагональную матрицу Γ справа ($A' = A\Gamma$, $\Gamma = (\gamma_{jj})$, $j = 1, \dots, n$), а вторая — слева ($A' = \Gamma A$, $\Gamma = (\gamma_{ii})$, $i = 1, \dots, m$).

Предполагаемые ошибки в матрице A , выражающиеся в пропорциональных изменениях элементов столбцов, никак не отражаются на погрешности измерений ε , а влияют только на параметры X . Действительно, рассмотрим ограничение задачи (6):

$$A'X = B.$$

Матрица $A' = A\Gamma$; следовательно,

$$A\Gamma X = B.$$

Если обозначить ΓX через X' , получим систему, идентичную исходной:

$$AX' = B.$$

В случае предполагаемых ошибок в матрице A , выражающихся в пропорциональных изменениях строк исходной матрицы A , задача определения неизвестных параметров эквивалентна задаче (5). Покажем это. Соотношения для определения параметров X на основании имеющейся информации имеют вид:

$$A'X = B.$$

Так как матрица $A' = \Gamma A$, то

$$\Gamma AX = B.$$

Умножим обе части уравнения на Γ^{-1} :

$$\Gamma^{-1}\Gamma AX = \Gamma^{-1}B.$$

С учетом $\Gamma^{-1}\Gamma = E$ получим

$$AX = \Gamma^{-1}B.$$

Если Γ^{-1} обозначить через Δ , то последнее соотношение эквивалентно задаче (5). Отметим, что обратная матрица Γ^{-1} всегда существует, поскольку матрица Γ — это диагональная матрица с отличными от нуля элементами на главной диагонали, поэтому и ее определитель также отличен от нуля. А такая матрица, как известно, всегда имеет обратную.

Основным преимуществом предложенного подхода к определению параметров линейных математических моделей (1) является возможность учета всех дополнительных требований, предъявляемых к параметрам X , уже на стадии формализации модели, что позволяет исключить возможность получения заведомо неприемлемых результатов.

3 Методика определения значимости измерений

Важное практическое значение имеет оценка влияния погрешности экспериментальных данных модели на погрешность измерений ε [4, 5], что позволяет осуществлять анализ информационной ценности измерений и, как следствие, выявлять те, которые следует рассматривать как наиболее недоуверенные или значимые и пр. Результатами такого анализа могут быть, например, выводы о необходимости, при наличии соответствующих возможностей, уточнения некоторых экспериментальных данных или рекомендации об их исключении из рассмотрения при непосредственном построении функциональных зависимостей.

Известно, что при решении задач линейного программирования для этих целей используется теория двойственности [6–8]. Согласно третьей теореме двойственности компоненты оптимального решения двойственной задачи равны частным производным целевой функции прямой задачи по соответствующим параметрам, в качестве которых выступают свободные члены системы ограничений исходной задачи. В силу того, что применительно к задаче восстановления зависимостей такие параметры являются величинами, наблюдаемыми в ходе эксперимента, предполагая их малые изменения, с помощью анализа оптимальных значений двойственных переменных можно оценить значимость каждого отдельного наблюдения.

Применительно к исследуемой в настоящей работе задаче решение соответствующих двойственных задач позволяет оценить влияние элементов

матрицы экспериментальных данных B и матрицы ограничений на параметры модели D в прямых задачах линейного программирования на величину минимального значения погрешности измерений ε . Это предоставляет возможность выявлять те элементы матриц B и D , которые вносят наибольший вклад в значение погрешности измерений ε и количественно его оценить.

С этих позиций для рассмотренных выше задач линейного программирования целесообразно рассмотреть двойственные к ним. Двойственная задача для задачи линейного программирования (4) имеет вид:

$$\left. \begin{aligned} (B, y^1) - (B, y^2) + (D, y^3) &\rightarrow \max; \\ A^T y^1 - A^T y^2 + C^T y^3 &\leq 0; \\ \sum_{i=1}^m y_i^1 + \sum_{i=1}^m y_i^2 &\leq 1; \\ y^1 \geq 0; y^2 \geq 0; y^3 &\geq 0. \end{aligned} \right\} \quad (7)$$

Здесь $y^1 = (y_i^1)$, $y^2 = (y_i^2)$ ($i = 1, \dots, m$) и $y^3 = (y_l^3)$ ($l = 1, \dots, k$) — векторы оптимального решения двойственной задачи.

Аналогичным образом может быть выписана двойственная задача для задачи линейного программирования (5):

$$\left. \begin{aligned} -\sum_{i=1}^m y_i^2 + \sum_{i=1}^m y_i^3 + (D, y^4) &\rightarrow \max; \\ A^T y^1 + C^T y^4 &\leq 0; \\ -y_i^1 b_i - y_i^2 + y_i^3 &\leq 0 \quad \forall i = 1, \dots, m; \\ \sum_{i=1}^m y_i^2 + \sum_{i=1}^m y_i^3 &\leq 1; \\ y^2 \geq 0; y^3 \geq 0; y^4 &\geq 0. \end{aligned} \right\} \quad (8)$$

Здесь $y^1 = (y_i^1)$, $y^2 = (y_i^2)$, $y^3 = (y_i^3)$ ($i = 1, \dots, m$) и $y^4 = (y_l^4)$ ($l = 1, \dots, k$) — векторы оптимального решения двойственной задачи.

Заметим, что в задаче (4) для оценки степени влияния i -го соотношения из системы неравенств $|AX - B| \leq \varepsilon$ на значение погрешности измерений ε необходимо рассмотреть соответствующие компоненты векторов $y^1 = (y_i^1)$ и $y^2 = (y_i^2)$, являющихся решением задачи (7), и выбрать из них максимальный. Аналогичную процедуру следует провести с векторами $y^2 = (y_i^2)$ и $y^3 = (y_i^3)$, являющимися решением задачи (8), для оценки степени влияния сводного члена i -го соотношения из системы неравенств $|\delta_i - 1| \leq \varepsilon$ на погрешность измерений ε в задаче (5).

4 Результаты апробации

Ниже представлены результаты определения параметров и оценки значимости используемых измерений на примере задачи нахождения распределения мольных долей фрагментов фуллерена с различным количеством заместителей в макроцепях полимеров, рассмотренной в работах [9, 10]. Для нахождения этого распределения составляется система уравнений Бугера–Ламберта, представляющая собой систему линейных алгебраических уравнений вида (1), с равным числом уравнений и неизвестных. Элементами квадратной матрицы A являются молярные экстинкции ядер несвязанного фуллерена и ядер, ковалентно связанных одной, двумя и n связями с заместителями (заместителями, фрагментами инициатора, макроцепями) соответственно, а в роли B — значения оптических плотностей, измеряемых спектрофотометрически в ультрафиолетовой/видимой области, для растворов фуллеренсодержащих продуктов (полимеров, смесей специально химически синтезированных индивидуальных замещенных фуллеренов). Определению подлежали параметры X , представляющие собой концентрации содержания фрагментов фуллерена в макроцепях полимеров. На основании параметров X могут быть рассчитаны и мольные доли концентраций, для чего необходимо разделить значение каждой концентрации на сумму всех концентраций.

Матрицы A и B формировались на основании экспериментальных данных:

$$A = \begin{pmatrix} 54\,000 & 30\,800 & 35\,800 & 28\,500 & 30\,900 \\ 30\,900 & 22\,800 & 28\,300 & 27\,900 & 28\,800 \\ 19\,600 & 21\,000 & 21\,800 & 18\,500 & 16\,050 \\ 50\,500 & 24\,370 & 17\,630 & 15\,070 & 11\,800 \\ 60\,780 & 24\,150 & 15\,350 & 13\,000 & 11\,700 \end{pmatrix};$$

$$B = \begin{pmatrix} 2,453 \\ 2,001 \\ 1,475 \\ 1,435 \\ 1,408 \end{pmatrix}.$$

Решение, полученное методом обратной матрицы, следующее:

$$X = (5,067 \cdot 10^{-1} \quad 3,554 \cdot 10^{-5} \quad 1,211 \cdot 10^{-5} \quad -3,096 \cdot 10^{-6} \quad 3,19 \cdot 10^{-5})^T.$$

Такое решение не имеет смысла, так как не все элементы вектора X неотрицательны. Именно поэтому определялось приближенное решение посредством сведения исходной проблемы к задаче линейного программирования вида (4):

$$\left. \begin{aligned} \varepsilon &\rightarrow \min; \\ |AX - B| &\leq \varepsilon; \\ X &\geq 0. \end{aligned} \right\} \quad (9)$$

В данной задаче ограничение $CX \geq D$ отсутствует, поскольку нет дополнительных ограничений на параметры X , кроме ограничения неотрицательности.

Решение задачи (9):

$$X = (9,504 \cdot 10^{-7} \quad 3,333 \cdot 10^{-5} \quad 1,274 \cdot 10^{-5} \quad 0 \quad 2,964 \cdot 10^{-5})^T,$$

параметр ε , характеризующий величину погрешности измерений, в этом случае равен 0,002786.

В предположении существования ошибок в матрице B рассматривалась задача линейного программирования вида (5), которая в обозначениях поставленной задачи имеет вид:

$$\left. \begin{aligned} \varepsilon &\rightarrow \min; \\ AX &= \Delta B; \\ X &\geq 0; \\ |\delta_i - 1| &\leq \varepsilon \quad \forall i = 1, \dots, 5. \end{aligned} \right\} \quad (10)$$

Решение задачи (10):

$$X = (1,031 \cdot 10^{-6} \quad 3,301 \cdot 10^{-5} \quad 1,379 \cdot 10^{-5} \quad 0 \quad 2,859 \cdot 10^{-5})^T.$$

Параметры, характеризующие ошибки в экспериментальных данных (в элементах матрицы B): $\delta_1 = 0,9986$; $\delta_2 = 0,9986$; $\delta_3 = 0,9986$; $\delta_4 = 1,0014$; $\delta_5 = 0,9986$. Погрешность измерений в этом случае составила $\varepsilon = 0,001396$.

Заметим, что в такой постановке, когда каждый элемент матрицы B умножается на параметр, характеризующий ошибку в этом элементе матрицы, погрешность измерений ε допускает удобное представление в процентах (в данном случае $\varepsilon = 0,1396\%$).

Полученные результаты позволили полагать, что отсутствие физического смысла в точном решении рассматриваемой задачи может быть связано с наличием ошибки в экспериментальных данных. Для оценки влияния предполагаемой погрешности экспериментальных данных в матрице B на погрешность измерений ε были выписаны двойственные задачи к задачам (9) и (10). Двойственная задача к задаче (9) имеет вид:

$$\left. \begin{aligned} (B, y^1) - (B, y^2) &\rightarrow \max; \\ A^T y^1 - A^T y^2 &\leq 0; \\ \sum_{i=1}^5 y_i^1 + \sum_{i=1}^5 y_i^2 &\leq 1; \\ y^1 &\geq 0; \quad y^2 \geq 0. \end{aligned} \right\} \quad (11)$$

Здесь $y^1 = (y_i^1)$, $y^2 = (y_i^2)$ ($i = 1, \dots, 5$).
Решение задачи (11):

$$\begin{aligned} y^1 &= (0,044 \ 0 \ 0,177 \ 0 \ 0,284)^T; \\ y^2 &= (0 \ 0,098 \ 0 \ 0,397 \ 0)^T. \end{aligned}$$

Компоненты решения двойственной задачи (11) показывают, что наибольший вклад в значение погрешности измерений ε в задаче (9) вносит четвертый элемент матрицы B , так как выражение $\max_i \{y_i^1; y_i^2\}$ достигает наибольшего значения (0,397) на четвертой компоненте оптимального решения y^2 , а наименьший (0,044) — первый, так как минимальное значение выражение $\min_i \{y_i^1; y_i^2\}$ принимает на первой компоненте оптимального решения y^1 .

Двойственная задача к задаче (10) имеет вид:

$$\left. \begin{aligned} -\sum_{i=1}^5 y_i^2 + \sum_{i=1}^5 y_i^3 &\rightarrow \max; \\ A^T y^1 &\leq 0; \\ -y_i^1 b_i - y_i^2 + y_i^3 &\leq 0 \quad \forall i = 1, \dots, 5; \\ \sum_{i=1}^5 y_i^2 + \sum_{i=1}^5 y_i^3 &\leq 1; \\ y^2 &\geq 0; \quad y^3 \geq 0. \end{aligned} \right\} \quad (12)$$

Здесь $y^1 = (y_i^1)$, $y^2 = (y_i^2)$, $y^3 = (y_i^3)$ ($i = 1, \dots, 5$).
Решение двойственной задачи (12):

$$\begin{aligned} y^1 &= (0,029 \ -0,064 \ 0,115 \ -0,259 \ 0,109)^T; \\ y^2 &= (0 \ 0,128 \ 0,371 \ 0)^T; \\ y^3 &= (0,07 \ 0 \ 0,169 \ 0 \ 0,261)^T. \end{aligned}$$

Решение двойственной задачи (12) также показывает, что наибольший вклад в значение погрешности измерений ε в задаче (10) определяется погрешностью четвертого элемента матрицы B , а наименьший — погрешностью первого. Из этого следует, что наименее достоверным следует полагать измерение величины B в четвертом эксперименте.

Аналогичный подход был применен к определению степени влияния погрешностей наблюдаемых величин на погрешность измерения при решении задачи моделирования численности населения Российской Федерации методом системной динамики [10–17]. Общий вид исследованной модели

системной динамики в терминах разностных уравнений следующий:

$$\left. \begin{aligned} \Delta N &= a_1 N^{\alpha_1} D^{\beta_1} I^{\gamma_1} - a_2 N^{\alpha_2} D^{\beta_2} I^{\gamma_2}; \\ \Delta D &= a_3 N^{\alpha_3} D^{\beta_3} I^{\gamma_3} - a_4 N^{\alpha_4} D^{\beta_4} I^{\gamma_4}; \\ \Delta I &= a_5 N^{\alpha_5} D^{\beta_5} I^{\gamma_5} - a_6 N^{\alpha_6} D^{\beta_6} I^{\gamma_6}, \end{aligned} \right\} \quad (13)$$

где N — численность населения РФ; D — душевые доходы за год; I — индекс потребительских цен. Информационную базу настоящего исследования составили данные официальной статистической отчетности за период с 1998 по 2010 гг.

По результатам специально организованного численного эксперимента, описанного в работах [12, 15–17], на основании данных за 1998–2009 гг. была получена модель, с достаточно высокой точностью описывающая экспериментальные данные:

$$\left. \begin{aligned} \frac{dN}{dt} &= 8,139 \cdot 10^{-22} \frac{N^{2,05} D^2}{I^2} - 64,1 \frac{N^{0,33} D^{0,3}}{I^{0,3}}; \\ \frac{dD}{dt} &= 560 D^{0,35} - 9900 I; \\ \frac{dI}{dt} &= 0,131 I^{-0,4} - 0,0072 \frac{N^{0,092} D^{0,092}}{I^{0,092}}. \end{aligned} \right\} \quad (14)$$

Для непосредственной реализации численного эксперимента в среде программирования Delphi был разработан модуль, позволяющий учитывать ряд вытекающих из смысла решаемой задачи дополнительных условий и давать наглядную интерпретацию проводимых расчетов. Концепция модели базировалась на следующих предположениях:

- численность населения влияет как на прирост, так и на убыль населения (данное предположение очевидно в силу того, что чем больше людей, тем чаще рождаются дети и тем больше смертных случаев от естественных и прочих причин);
- реальный годовой доход жителей страны влияет на изменение численности населения как положительно (чем больше реальный доход семьи, тем проще решиться на рождение ребенка), так и отрицательно (семьи с большими доходами, как правило, имеют по 1–2 ребенка);
- на изменение душевых доходов сами доходы сказываются положительно (как правило, государство, осуществляющее социально-направленную политику, само индексирует доходы граждан и не позволяет работодателям снижать заработную плату);
- индекс потребительских цен снижает доходы;
- высокий уровень индекса потребительских цен заставляет население экономить, что приводит к снижению потребления товаров и услуг и, как

Таблица 1 Исходные данные для модели (13)

Год	Численность населения РФ N , чел.	Душевые доходы D , руб./чел. в год	Индекс потребительских цен I , доля ед.
1998	147 802 133	12 122,4	1,844
1999	147 539 426	19 906,8	1,365
2000	146 890 128	27 373,2	1,202
2001	146 303 611	36 744,0	1,186
2002	145 649 334	47 366,4	1,151
2003	144 963 650	62 044,8	1,120
2004	144 168 205	76 923,6	1,117
2005	143 474 219	97 342,8	1,109
2006	142 753 551	122 352,0	1,090
2007	142 220 968	151 232,4	1,119
2008	142 008 800	179 287,2	1,133
2009	141 904 000	202 282,8	1,088
2010	141 914 509	226 572,0	1,088

следствие, к снижению скорости изменения самого индекса потребительских цен;

- высокий реальный годовой доход населения обуславливает устойчивый спрос на товары и услуги, что обеспечивает благоприятные условия работы для производителей, а следовательно, замедляет скорость изменения индекса потребительских цен.

Вместе с тем, как показывает практика использования статистической информации, исходные данные для модели (13) (табл. 1) нельзя рассматривать как абсолютно точные. Это обусловлено рядом причин объективного характера: запаздывание сбора и обработки данных, неточности (а иногда и умышленные искажения) предоставляемой информации, аккумулируемой органами статистики, и пр. В этой связи, несмотря на хорошую точность построенной модели, целесообразной представляется оценка значимости имеющихся экспериментальных данных. И одним из возможных способов реализации этого является описанный в настоящей работе подход.

Ниже приводится описание предложенного подхода применительно к первому уравнению модели (13). Для этого потребовалось осуществить линеаризацию данного уравнения, что было осуществлено посредством использования разложения в ряд Тейлора с центром в точке $\{a_1^0, a_2^0, \alpha_{1,2} = 0, \beta_{1,2} = 0, \gamma_{1,2} = 0\}$:

$$\Delta N \approx a_1 + a_1^0 \ln N \cdot \alpha_1 + a_1^0 \ln D \cdot \beta_1 + a_1^0 \ln I \cdot \gamma_1 - a_2 - a_2^0 \ln N \cdot \alpha_2 - a_2^0 \ln D \cdot \beta_2 - a_2^0 \ln I \cdot \gamma_2. \quad (15)$$

Были введены следующие обозначения: $x_1 = a_1$; $x_2 = a_2$; $x_3 = \alpha_1$; $x_4 = \alpha_2$; $x_5 = \beta_1$; $x_6 = \beta_2$; $x_7 = \gamma_1$; $x_8 = \gamma_2$.

На основании исходной информации (см. табл. 1) для определения параметров $\{a_i, \alpha_i, \beta_i, \gamma_i\}$, $i = \overline{1,2}$, на основании соотношений (15) были сформированы матрицы A и B :

$$A = \begin{pmatrix} 1 & -1 & 0,019 & 1205,810 & 0,09 & 602,720 \\ 1 & -1 & 0,019 & 1205,696 & 0,010 & 634,514 \\ 1 & -1 & 0,019 & 1205,413 & 0,010 & 654,930 \\ 1 & -1 & 0,019 & 1205,157 & 0,011 & 673,802 \\ 1 & -1 & 0,019 & 1204,869 & 0,011 & 690,079 \\ 1 & -1 & 0,019 & 1204,567 & 0,011 & 707,383 \\ 1 & -1 & 0,019 & 1204,214 & 0,011 & 721,161 \\ 1 & -1 & 0,019 & 1203,905 & 0,011 & 736,252 \\ 1 & -1 & 0,019 & 1203,582 & 0,012 & 750,910 \\ 1 & -1 & 0,019 & 1203,342 & 0,012 & 764,493 \\ 1 & -1 & 0,019 & 1203,247 & 0,012 & 775,433 \\ 1 & -1 & 0,019 & 1203,199 & 0,012 & 783,713 \end{pmatrix}; \quad B = \begin{pmatrix} 0,000 & 0 \\ 0,000 & -19,945 \\ 0,000 & -11,794 \\ 0,000 & -10,935 \\ 0,000 & -9,014 \\ 0,000 & -7,264 \\ 0,000 & -7,092 \\ 0,000 & -6,632 \\ 0,000 & -5,524 \\ 0,000 & -7,207 \\ 0,000 & -8,004 \\ 0,000 & -5,406 \end{pmatrix}; \quad (16)$$

(Значения параметров a_1^0 и a_2^0 полагались равными 0,001 и 64,1 соответственно.)

Все требования к параметрам $\{a_i, \alpha_i, \beta_i, \gamma_i\}$, $i = \overline{1,2}$, были учтены в виде дополнительных ограничений:

$$\begin{aligned} 0 \leq x_1 \leq 3; & \quad 50 \leq x_2 \leq 100; \\ 1,948 \leq x_3 \leq 2,050; & \quad 0 \leq x_4 \leq 1; \\ 0 \leq x_5 \leq 5; & \quad 0 \leq x_6 \leq 5; \\ -5 \leq x_7 \leq 0; & \quad -5 \leq x_8 \leq 0. \end{aligned}$$

Таким образом, задача определения приближенного значения параметров $\{a_i, \alpha_i, \beta_i, \gamma_i\}$, $i = \overline{1, 2}$, была формализована в следующем виде:

$$\left. \begin{aligned} \varepsilon \rightarrow \min; \\ |AX - B| \leq \varepsilon; \\ 0 \leq x_1 \leq 3; \quad 50 \leq x_2 \leq 400\,000; \\ 1,948 \leq x_3 \leq 2,050; \quad 0 \leq x_4 \leq 1; \\ 0 \leq x_5 \leq 5; \quad 0 \leq x_6 \leq 5; \\ -5 \leq x_7 \leq 0; \quad -5 \leq x_8 \leq 0. \end{aligned} \right\} \quad (17)$$

(Матрицы A и B имеют вид (16).)

Решение прямой задачи (17):

$$X = (3 \ 396\,128,4 \ 2,05 \ 0 \ 5 \ 5 \ 0 \ 0)^T;$$

погрешность измерений $\varepsilon = 402\,796,7$; решение двойственной к ней (см. модель (7)):

$$\begin{aligned} y^1 &= (0\ 0\ 0\ 0\ 0,5\ 0\ 0\ 0\ 0\ 0)^T; \\ y^2 &= (0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0,5)^T; \\ y^3 &= (0\ 0\ 0\ 0\ 0\ 0\ 3,82\ 0\ 0)^T. \end{aligned}$$

Аналогичные действия были проведены и с остальными уравнениями системы (13). В качестве центров разложения в рядах Тейлора для параметров a_i , $i = \overline{3, 6}$, были выбраны значения этих же параметров в модели (14). Диапазоны вариации переменных, необходимые для формирования моделей (17) по каждому уравнению системы (13), приведены в табл. 2.

Полученные результаты (см. табл. 2) показывают, что на погрешность измерений ε второго и третьего уравнения системы (13) заданные диапазоны вариаций параметров каждого из этих уравнений не оказывают влияния (об этом свидетельствуют нулевые значения двойственных переменных y^3). Это означает, что величины погрешностей измерений

Таблица 2 Оценки погрешностей измерений для уравнений модели (13) и значимости ограничений на их параметры

Параметр	Диапазон вариации	Решение прямой задачи	Двойственные оценки y^3	
			по левой границе	по правой границе
Уравнение 1 ($\varepsilon = 402\,796,7$)				
a_1	[0, 3]	3,0	—	0
a_2	[50, 400 000]	396 128,4	0	0
α_1	[1,948, 2,050]	2,05	0	0
α_2	[0, 5]	0,0	—	0
β_1	[0, 5]	5,0	—	0
β_2	[0, 5]	5,0	—	3,82
γ_1	[-5, 0]	0,0	0	—
γ_2	[-5, 0]	0,0	0	—
Уравнение 2 ($\varepsilon = 394,56$)				
a_3	[500, 600]	600,0	0	$6,6 \cdot 10^{-5}$
a_4	[9000, 11 000]	9000,0	—	0
α_3	[0, 1]	1,0	—	$1,4 \cdot 10^{-4}$
α_4	[0, 1]	0,019	—	0
β_3	[0, 1]	0,0	—	0
β_4	[0, 1]	0,05	—	0
γ_3	[0, 1]	0,0	—	0
γ_4	[0, 2]	0,187	—	0
Уравнение 3 ($\varepsilon = 0,26$)				
a_5	[0, 1]	0,0	—	0
a_6	[0, 1]	0,023	—	0
α_5	[0, 1]	0,0	—	0
α_6	[0, 1]	0,0	—	0
β_5	[0, 1]	0,0	—	0
β_6	[0, 1]	0,0	—	0
γ_5	[-1, 0]	0,0	0	0
γ_6	[-1, 0]	0,0	0	0

Таблица 3 Оценки значимости измерений для уравнений модели (13)

Номер эксперимента (измерения)	Уравнение 1		Уравнение 2		Уравнение 3	
	y^1	y^2	y^1	y^2	y^1	y^2
1	0	0	0,0345	0	0	0,5
2	0	0	0	0	0,5	0
3	0	0	0	0	0	0
4	0	0	0	0	0	0
5	0	0	0	0	0	0
6	0,5	0	0,1323	0	0	0
7	0	0	0	0	0	0
8	0	0	0	0	0	0
9	0	0	0	0,5	0	0
10	0	0	0	0	0	0
11	0	0	0	0	0	0
12	0	0,5	0,3333	0	0	0

Замечание: y^1 и y^2 — верхняя и нижняя границы, определяемые соответствующим ограничением группы $|AX - B| \leq \varepsilon$.

для данных уравнений в первую очередь определяются измерениями (табл. 3). Так, погрешность измерений второго уравнения определяется в большей степени 9-м и 12-м измерениями (9-я компонента вектора y^2 и 12 компонента вектора y^1 равны соответственно 0,5 и 0,333) и в меньшей — 1-м и 6-м. На погрешность измерений ε третьего уравнения в равной степени влияют 1-е и 2-е измерения. Среди параметров первого уравнения есть один (β_2), верхняя граница значений которого оказывает существенное влияние на погрешность измерений ε (соответствующее значение двойственной оценки y^3 равно 3,82). Также на значение ε в равной степени оказывают влияние 6-е и 12-е измерения, о чем свидетельствуют 6-я и 12-я компоненты векторов y^1 и y^2 соответственно.

Данная информация может быть использована при планировании дальнейших экспериментов с моделями (13) или (14). При этом должен соблюдаться следующий принцип: в каждом последующем эксперименте информативность новых измерений должна быть не менее значима, что должно отражаться при формировании соответствующих условий. Например, сужение или увеличение диапазона значений параметра β_2 целесообразно проводить за счет вариации верхней границы интервала его значений.

Литература

1. *Кендалл М., Стьюарт А.* Многомерный статистический анализ и временные ряды. — М.: Наука, 1976. 736 с.
2. *Айвазян С. А., Мхитарян В. С.* Прикладная статистика и основы эконометрики. — М.: ЮНИТИ, 1998. 1022 с.
3. *Эконометрика / Под ред. И. И. Елисейевой.* — М.: Финансы и статистика, 2006. 576 с.
4. *Марчук Г. И.* Методы вычислительной математики. — М.: Наука, 1977. 457 с.
5. *Сливак С. И., Тимошенко В. И., Слинько М. Г.* Методы построения кинетических моделей стационарных реакций // Химическая промышленность, 1979. № 3. С. 33–36.
6. *Канторович Л. В.* Экономический расчет наилучшего использования ресурсов. — М.: Изд-во Академии наук СССР, 1960. 347 с.
7. *Зуховицкий С. И., Авдеева Л. И.* Линейное и выпуклое программирование. — М.: Наука, 1967. 460 с.
8. *Канторович Л. В., Горстко А. Б.* Оптимальные решения в экономике. — М.: Наука, 1972. 231 с.
9. *Кузнецов С. И., Юмагулова Р. Х., Медведева Н. А., Хамидуллин Ф. Ф., Колесов С. В.* Фуллеренсодержащие полимеры. Уф-спектроскопическое исследование // Высокомолекулярные соединения. Сер. А, 2012. Т. 54. № 6. С. 859–864.
10. *Кузнецов С. И., Хамидуллин Ф. Ф., Юмагулова Р. Х., Медведева Н. А., Лебедев Ю. А., Колесов С. В.* Самоорганизация функционализированных фуллереном C_{60} макромолекул полиметилметакрилата и полистирола // Высокомолекулярные соединения. Сер. А, 2012. Т. 54. № 10. С. 1527–1531.
11. *Сливак С. И., Кантор О. Г.* Оценка параметров моделей системной динамики // Журнал СВМО, 2011. Т. 13. № 3. С. 107–113.
12. *Сливак С. И., Кантор О. Г., Салахов И. Р.* О программе, корректирующей систему уравнений // Журнал СВМО, 2011. Т. 13. № 4. С. 87–93.
13. *Сливак С. И., Кантор О. Г.* Качество моделей математической обработки наблюдений социально-экономических систем // Системы управления и информационные технологии, 2012. № 2(48). С. 44–49.

14. Спивак С. И., Кантор О. Г. Оценка качества спецификации моделей системной динамики // Журнал СВМО, 2012. Т. 14. № 2. С. 34–39.
15. Спивак С. И., Кантор О. Г., Салахов И. Р. Вычислительная реализация оценки управляющих параметров модели системной динамики // Вестник Башкирского университета, 2012. Т. 17. № 4. С. 1658–1660.
16. Спивак С. И., Кантор О. Г., Салахов И. Р. Алгоритм получения прогнозируемых параметров социально-экономических систем // Системы управления и информационные технологии, 2013. № 4(54). С. 43–45.
17. Спивак С. И., Кантор О. Г. Построение моделей системной динамики в условиях ограниченной экспертной информации // Информатика и её применения, 2014. Т. 8. Вып. 2. С. 112–122.

Поступила в редакцию 04.07.14

EVALUATION OF MEASUREMENT ACCURACY AND SIGNIFICANCE FOR LINEAR MODELS

S. I. Spivak¹, O. G. Kantor², D. S. Yunusova¹, S. I. Kuznetsov³, and S. V. Kolesov³

¹Bashkir State University, 32 Valdy Str., Ufa 450076, Russian Federation

²Institute of Social and Economic Research, Ufa Scientific Center, Russian Academy of Sciences; 71 Av. Oktyabrya, Ufa 450054, Russian Federation

³Institute of Organic Chemistry, Ufa Scientific Center, Russian Academy of Sciences, 71 Av. Oktyabrya, Ufa 450054, Russian Federation

Abstract: Identification of a linear dependency, when exact solution obtained by standard methods does not meet the objective requirements, determines development of specific approaches for their numerical realization. A method to obtain approximate values of linear models parameters on experimental data, which is based on the use of the linear programming methodology and the duality theory, is presented. This method makes it possible to obtain approximate solutions that fulfill all requirements to the model and its parameters and to evaluate accuracy and significance of measurements. It is important for improving the procedure of construction of functional dependencies on the stage of planning experiments if they do not satisfy the authenticity criteria. The results of testing the proposed method for problems connected with research of chemical and socioeconomic systems are given.

Keywords: problems of linear dependencies recovering; measurement accuracy; measurement significance; dual estimates

DOI: 10.14357/19922264150108

Acknowledgments

The research was financially supported by the Russian Foundation for Basic Research (project 13-01-00749).

References

1. Kendall, M., and A. St'yuart. 1976. *Mnogomernyy statisticheskiy analiz i vremennyye ryady* [Multivariate statistical analysis and time series]. Moscow: Nauka. 736 p.
2. Ayvazyan, S. A., and V. S. Mkhitarian. 1998. *Prikladnaya statistika i osnovy ekonometriki* [Applied statistics and bases of econometrics]. Moscow: YUNITI. 1022 p.
3. Eliseeva, I. I., ed. 2008. *Ekonometrika* [Econometrics]. Moscow: Finance and Statistics. 576 p.
4. Marchuk, G. I. 1977. *Metody vychislitel'noy matematiki* [Methods of calculus mathematics]. Moscow: Nauka. 457 p.
5. Spivak, S. I., V. I. Timoshenko, and M. G. Slin'ko. 1979. *Metody postroeniya kineticheskikh modeley statsionarnykh reaktsiy* [Methods of creation of kinetic models of stationary reactions]. *Khimicheskaya Promyshlennost'* [Chemical Industry] 3:33–36.
6. Kantorovich, L. V. 1960. *Ekonomicheskiy raschet nailuchshego ispol'zovaniya resursov* [Economic calculation of the best use of resources]. Moscow: Publishing House of Academy of Sciences of the USSR. 347 p.
7. Zukhovitskiy, S. I., and L. I. Avdeeva. 1967. *Lineynoe i vypukloe programmirovaniye* [Linear and convex programming]. Moscow: Nauka. 460 p.
8. Kantorovich, L. V., and A. B. Gorstko. 1972. *Optimal'nyye resheniya v ekonomike* [Optimal solutions in economics]. Moscow: Nauka. 231 p.
9. Kuznetsov, S. I., R. Kh. Yumagulova, N. A. Medvedeva, F. F. Khamidullin, and S. V. Kolesov. 2012. Fulleren-

- soderzhashchie polimery. UF-spektroskopicheskoe issledovanie [Fullerene polymers. UV spectroscopic research]. *Vysokomolekulyarnye Soedineniya. Ser. A* [Macromolecular Compounds A] 6:859–864.
10. Kuznetsov, S. I., F. F. Khamidullin, R. Kh. Yumagulova, N. A. Medvedeva, Yu. A. Lebedev, and S. V. Kolesov. 2012. Samoorganizatsiya funktsionalizirovannykh fullerenom C₆₀ makromolekul polimetilmetakrilata i polistirola [Self-organization functionalized with fullerenes C₆₀ of macromolecules polymethylmethacrylate and polystyrene]. *Vysokomolekulyarnye Soedineniya. Ser. A* [Macromolecular Compounds A] 10:1527–1531.
 11. Spivak, S. I., and O. G. Kantor. 2011. Otsenka parametrov modeley sistemnoy dinamiki [Estimation of parameters of system dynamics models]. *Zh. SVMO* [J. SVMO] 3:107–113.
 12. Spivak, S. I., O. G. Kantor, and I. R. Salakhov. 2011. O programme, korrektiruyushchey sistemu uravneniy [About the program correcting system of the equations]. *Zh. SVMO* [J. SVMO] 4:87–93.
 13. Spivak, S. I., and O. G. Kantor. 2012. Kachestvo modeley matematicheskoy obrabotki sotsial'no-ekonomicheskikh sistem [Quality of mathematical processing observations models of socioeconomic systems]. *Sistemy Upravleniya i Informatsionnye Tekhnologii* [Management and Information Technology] 2(48):44–49.
 14. Spivak, S. I., and O. G. Kantor. 2012. Otsenka kachestva spetsifikatsii modeley sistemnoy dinamiki [Estimation of quality of specification of system dynamics models]. *Zh. SVMO* [J. SVMO] 2:34–39.
 15. Spivak, S. I., O. G. Kantor, and I. R. Salakhov. 2012. Vychislitel'naya realizatsiya otsenki upravlyayushchikh parametrov modeli sistemnoy dinamiki [Computational realization of estimation of the control parameters of system dynamics model]. *Vestnik Bashkirskogo Universiteta* [Bashkir University Bulletin] 4:1658–1660.
 16. Spivak, S. I., O. G. Kantor, and I. R. Salakhov. 2013. Algoritm polucheniya prognoziruemyykh parametrov sotsial'no-ekonomicheskikh sistem [Algorithm of obtaining predicted parameters of social and economic systems]. *Sistemy Upravleniya i Informatsionnye Tekhnologii* [Management and Information Technology] 4(54):43–45.
 17. Spivak, S. I., and O. G. Kantor. 2014. Postroenie modeley sistemnoy dinamiki v usloviyakh ogranichennoy ekspertnoy informatsii [Construction of system dynamics models in the conditions of limited expert information]. *Informatika i ee Primeneniya — Inform. Appl.* 2:112–122.

Received July 4, 2014

Contributors

Spivak Semen I. (b. 1945) — Doctor of Science in physics and mathematics, professor, Bashkir State University, 32 Zaki Validi Str., Ufa 450074, Russian Federation; semen.spivak@mail.ru

Kantor Olga G. (b. 1971) — Candidate of Science (PhD) in physics and mathematics, senior scientist, Institute of Social and Economic Research, Ufa Scientific Centre, Russian Academy of Sciences, 71 October Av., Ufa 450054, Russian Federation; o.kantor@mail.ru

Yunusova Darya S. (b. 1989) — PhD student, Bashkir State University, 32 Zaki Validi Str., Ufa 450074, Russian Federation; kazakova_d_s@mail.ru

Kuznetsov Sergey I. (b. 1955) — scientist, Institute of Organic Chemistry, Ufa Scientific Center, Russian Academy of Sciences, 71 October Av., Ufa 450054, Russian Federation; chemorg@anrb.ru

Kolesov Sergey V. (b. 1951) — Head of Laboratory, Institute of Organic Chemistry, Ufa Scientific Center, Russian Academy of Sciences, 71 October Av., Ufa 450054, Russian Federation; kolesovservic@rambler.ru

БАЙЕСОВСКАЯ РЕКУРРЕНТНАЯ МОДЕЛЬ РОСТА НАДЕЖНОСТИ: БЕТА-РАВНОМЕРНОЕ РАСПРЕДЕЛЕНИЕ ПАРАМЕТРОВ*

Ю. В. Жаворонкова¹, А. А. Кудрявцев², С. Я. Шоргин³

Аннотация: Прогнозирование надежности сложных модифицируемых информационных систем (СМИС) является в настоящее время одной из актуальных задач теории массового обслуживания. Любая впервые созданная сложная система, предназначенная для переработки или передачи информационных потоков, как правило, не обладает требуемой надежностью. Такие системы подвергаются модификациям в ходе разработки, опытной эксплуатации и штатного функционирования. Целью этих модификаций является увеличение надежности информационных систем. В связи с этим возникает необходимость формализации понятия надежности модифицируемых информационных систем и разработки методов и алгоритмов оценивания и прогнозирования различных надежностных характеристик. Одним из подходов к определению надежности системы является вычисление вероятности того, что на сигнал, поданный на вход системы в определенный момент времени, система отреагирует корректно. В статье рассматривается экспоненциальная рекуррентная модель роста надежности, в которой вероятность надежности системы представляется как линейная комбинация параметров «дефективности» и «эффективности» средства, исправляющего недостатки системы. Предполагается, что исследователь не имеет точных сведений об исследуемой системе, а лишь знаком с характеристиками класса, из которого берется данная система. В рамках байесовского подхода предполагается, что один из показателей «дефективности» и «эффективности» имеет бета-распределение, а другой — равномерное распределение. Вычисляется средняя предельная надежность системы. Приводятся численные результаты для модельных примеров.

Ключевые слова: модифицируемые информационные системы; теория надежности; байесовский подход; бета-распределение; равномерное распределение

DOI: 10.14357/19922264150109

1 Постановка задачи

Задача прогнозирования надежности СМИС была сформулирована в [1], а в дальнейшем более подробно рассмотрена в [2]. В статье [3] дано подробное описание класса моделей, в рамках которых возникает необходимость использования байесовского подхода к анализу роста надежности СМИС, таких как новая программная система для компьютера, новая информационно-вычислительная сеть или новая административно-информационная система, которые, как правило, изначально не обладают требуемой надежностью.

Исследуемые СМИС подвергаются периодическим изменениям (модификациям) с целью увеличения надежности информационных систем. В статье рассматривается описанная в книге [2] модель роста надежности, обычно используемая, когда удобно иметь дело непосредственно с параметром, интерпретируемым как надежность системы.

Рассмотрим произвольную систему, на вход которой подаются некоторые сигналы (например, команды оператора или внешние воздействия). Реакция системы на поданные сигналы может быть либо правильной (корректной), либо неправильной (некорректной).

В каждый момент времени t надежность системы можно характеризовать параметром $p(t)$ — вероятностью того, что на сигнал, поданный на вход системы в момент t , система отреагирует правильно. По смыслу такая характеристика надежности ближе всего к традиционно используемому коэффициенту готовности. В случайные моменты времени $0 = Y_0 \leq Y_1 \leq Y_2 \leq \dots$ система подвергается (мгновенной) модификации, в результате чего изменяется параметр $p(t)$.

Следует обратить внимание на то обстоятельство, что ниже рассматривается непрерывное время, без привязки напрямую процесса модифицирования системы к процессу ее тестирования. Пред-

* Исследование выполнено при поддержке Российского научного фонда (проект 14-11-00397).

¹ ООО Спутник, juliana-zh@yandex.ru

² Московский государственный университет им. М. В. Ломоносова, факультет вычислительной математики и кибернетики, nubigena@mail.ru

³ Институт проблем информатики Российской академии наук, sshorgin@ipiran.ru

положим, что траектории процесса $p(t)$ непрерывны справа и кусочно-постоянны, так что $p(t) = p(Y_j)$ при $Y_j \leq t < Y_{j+1}$.

Обозначим $p_j = p(Y_j)$. Рассмотрим поведение p_j в зависимости от изменения j . Другими словами, будем изучать изменение надежности системы в зависимости от номера модификации. В книге [2] рассматривается, в частности, следующая рекуррентная модель роста надежности. Пусть $\{(\theta_j, \eta_j)\}$, $j \geq 1$, — последовательность независимых одинаково распределенных двумерных случайных векторов таких, что $0 < \eta_1 < 1$; $0 < \theta_1 < 1$ почти наверное.

Задав начальную надежность p_0 , рассмотрим модель, определяемую рекуррентным соотношением

$$p_{j+1} = \eta_{j+1}p_j + \theta_{j+1}(1 - p_j).$$

Эта модель названа дискретной экспоненциальной моделью. В такой модели случайные величины η_j (параметры «дефективности») описывают возможное уменьшение надежности из-за некачественных модификаций, в ходе которых вместо исправления существующих дефектов в систему могут быть внесены новые, в то время как величины θ_j (параметры «эффективности») описывают повышение надежности за счет исправления дефектов.

Обозначим $\lambda = 1 - E\theta_1$, $\mu = E\eta_1$. В [2] доказано, что при условии $\lambda + \mu \neq 1$

$$p = \lim_{j \rightarrow \infty} E p_j = \frac{\mu}{\lambda + \mu}.$$

Изучение предельного значения средней величины $E p_j$ представляет значительный интерес, поскольку эта величина характеризует асимптотическое значение надежности системы в рамках некоторой рекуррентной модели, задаваемой набором $\{(\theta_j, \eta_j)\}$. Из результатов [2] следует, что это асимптотическое значение зависит только от средних значений величин $\{(\theta_j, \eta_j)\}$, $j \geq 1$.

В [3, 4] исследовалась ситуация, при которой рассматривается набор однотипных сложных модифицируемых объектов (МО), каждый из которых обслуживается собственной ремонтной бригадой (РБ). Исследователю хотелось бы определить усредненное значение p по всем МО. Для решения этой задачи в указанной работе предложена так называемая байесовская постановка. Предполагается, что рассматривается целая группа однотипных МО и группа им соответствующих однотипных РБ. Пусть $m = 1, 2, \dots$ — номера этих объектов. Для каждого МО (вместе с его РБ) существует собственный набор $\{(\theta_j^m, \eta_j^m)\}$, $j \geq 1$, $m \geq 1$, независимых одинаково распределенных при каждом фиксированном j двумерных случайных векторов таких, что

$0 < \eta_1^m < 1$; $0 < \theta_1^m < 1$ почти наверное. Но средние значения величин θ_j^m , η_j^m , $j \geq 1$, $m \geq 1$, не предполагаются известными; более того, они не предполагаются даже одинаковыми. Вводится предположение, что величины $\lambda = 1 - E\theta_j^m$, $\mu = E\eta_j^m$ сами по себе являются случайными, т. е. на вероятностном пространстве, в которое в качестве элементарных событий входят все рассматриваемые в рамках данной постановки МО вместе с их РБ, заданы случайные величины λ и μ (которые полагаем независимыми), имеющие смысл $\lambda = 1 - E\theta_j^m$, $\mu = E\eta_j^m$, где m — случайный номер МО. Принимаемые исследователем за основу распределения величин λ и μ будем называть априорными.

Подлежащие вычислению характеристики такой «рандомизированной» группы МО, естественно, являются рандомизацией аналогичных характеристик «отдельно взятой» МО с учетом априорного распределения параметров λ и μ , взятого исследователем за основу. Наиболее естественной и удобной для изучения характеристикой является усредненное по всем МО значение предельной вероятности надежности, т. е.

$$p_{\text{сред}} = E p = E \frac{\mu}{\lambda + \mu},$$

где усреднение ведется по совместному распределению случайных величин (λ, μ) .

В рассматриваемой ситуации величины η_j^m и θ_j^m удовлетворяют ограничениям $0 < \eta_j^m < 1$, $0 < \theta_j^m < 1$. Значит, и значения λ и μ величин $1 - E\theta_j^m$ и $E\eta_j^m$ соответственно также находятся на отрезке $[0, 1]$. Поэтому в качестве априорных распределений параметров λ и μ следует выбирать только распределения, сосредоточенные на $[0, 1]$.

В работах [3, 4] были рассмотрены независимые случайные параметры λ и μ , имеющие одновременно равномерное или бета-распределение соответственно. В настоящей статье исследования байесовской рекуррентной модели роста надежности продолжены для ситуации, когда один из параметров имеет бета-распределение, а другой — равномерное распределение.

2 Основные результаты

Введем следующие обозначения. Через $B(m, n)$, $m, n > 0$, будем обозначать бета-функцию. Через $R(a, b)$ и $\beta(m, n)$ обозначим соответственно равномерное распределение и бета-распределение. Пусть

$$(\alpha)_i = \alpha(\alpha + 1) \cdots (\alpha + i - 1), \quad (\alpha)_0 = 1.$$

Несмотря на то что $(\alpha)_i$ имеет смысл неполного факториала, нигде далее не требуется, чтобы α

было положительным. Рассмотрим классическую гипергеометрическую функцию Гаусса

$$G(\alpha, \beta, \gamma; x) = \sum_{i=0}^{\infty} \frac{(\alpha)_i (\beta)_i}{(\gamma)_i i!} x^i.$$

По аналогии с 9.180.1, 9.180.3 и 9.14 п. 1 из [5] введем в рассмотрение две обобщенные гипергеометрические функции двух переменных:

$$G_{s,t}^{p,q}(\alpha, \beta_1, \dots, \beta_p, \beta'_1, \dots, \beta'_q; \gamma, \delta_1, \dots, \dots, \delta_s, \delta'_1, \dots, \delta'_t; x, y) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(\alpha)_{i+j} (\beta_1)_i \dots (\beta_p)_i (\beta'_1)_j \dots (\beta'_q)_j}{(\gamma)_{i+j} (\delta_1)_i \dots (\delta_s)_i (\delta'_1)_j \dots (\delta'_t)_j} \times \frac{x^i y^j}{i! j!}; \quad (1)$$

$$H_{s,t}^{p,q}(\beta_1, \dots, \beta_p, \beta'_1, \dots, \beta'_q; \gamma, \delta_1, \dots, \delta_s, \delta'_1, \dots, \dots, \delta'_t; x, y) = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(\beta_1)_i \dots (\beta_p)_i (\beta'_1)_j \dots (\beta'_q)_j}{(\gamma)_{i+j} (\delta_1)_i \dots (\delta_s)_i (\delta'_1)_j \dots (\delta'_t)_j} \times \frac{x^i y^j}{i! j!}. \quad (2)$$

Теорема 1. Пусть случайные величины λ и μ независимы и имеют соответственно распределения $R(a, b)$ и $\beta(m, n)$, где $0 \leq a < b \leq 1, m, n > 0$. Тогда

$$p_{\text{сред}} = \frac{B(m+1, n)}{(b-a)B(m, n)} \ln \left(\frac{b+1}{a+1} \right) + \frac{b(m+1)^{-2}}{(b-a)B(m, n)(b+1)} \times H_{1,0}^{3,2} \left(1-n, m+1, m+1, 1, 1; m+2, m+2; \frac{1}{b+1}, 1 \right) - \frac{a(m+1)^{-2}}{(b-a)B(m, n)(a+1)} \times H_{1,0}^{3,2} \left(1-n, m+1, m+1, 1, 1; m+2, m+2; \frac{1}{a+1}, 1 \right). \quad (3)$$

Доказательство. Найдем плотность $f_p(x)$ случайной величины p . Имеем

$$f_p(x) = \int_a^b \frac{y}{(1-x)^2} f_{\mu} \left(\frac{x}{1-x} y \right) f_{\lambda}(y) dy.$$

Используя замену переменной $z = xy/(1-x)$, по формуле 8.391 из [5] для $0 < x < 1/(b+1)$ имеем

$$f_p(x) = \int_a^b \frac{(1-x)^{-2} y}{(b-a)B(m, n)} \times \left(\frac{xy}{1-x} \right)^{m-1} \left(1 - \frac{xy}{1-x} \right)^{n-1} dy = \frac{(m+1)^{-1} x^{-2}}{(b-a)B(m, n)} \left(\left(\frac{bx}{1-x} \right)^{m+1} \times \right. \\ \left. \times G \left(1-n, m+1, m+2; \frac{bx}{1-x} \right) - \left(\frac{ax}{1-x} \right)^{m+1} G \left(1-n, m+1, m+2; \frac{ax}{1-x} \right) \right) \equiv S_1(x),$$

а для $1/(b+1) \leq x < 1/(a+1)$

$$f_p(x) = \int_a^{(1-x)/x} \frac{(1-x)^{-2} y}{(b-a)B(m, n)} \left(\frac{xy}{1-x} \right)^{m-1} \times \left(1 - \frac{xy}{1-x} \right)^{n-1} dy = \frac{x^{-2}}{(b-a)B(m, n)} \times \left(B(m+1, n) - \frac{1}{m+1} \left(\frac{ax}{1-x} \right)^{m+1} \times \right. \\ \left. \times G \left(1-n, m+1, m+2; \frac{ax}{1-x} \right) \right) \equiv S_2(x).$$

Таким образом,

$$E p = \int x f_p(x) dx = \int_0^{1/(b+1)} x S_1(x) dx + \int_{1/(b+1)}^{1/(a+1)} x S_2(x) dx. \quad (4)$$

Вычислим отдельно первый интеграл из правой части (4). Имеем

$$\int_0^{1/(b+1)} x S_1(x) dx = \int_0^{1/(b+1)} \frac{(m+1)^{-1} x^{-1}}{(b-a)B(m, n)} \left(\left(\frac{bx}{1-x} \right)^{m+1} \times \right. \\ \left. \times G \left(1-n, m+1, m+2; \frac{bx}{1-x} \right) - \left(\frac{ax}{1-x} \right)^{m+1} \times \right. \\ \left. \times G \left(1-n, m+1, m+2; \frac{ax}{1-x} \right) \right) dx \equiv U_1 - U_2.$$

Для первого слагаемого имеем:

$$U_1 = \frac{b^{m+1}(m+1)^{-1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i(m+1)_i b^i}{(m+2)_i i!} \times$$

$$\times \int_0^{1/(b+1)} \frac{x^{m+i} dx}{(1-x)^{m+i+1}} = -\frac{b^{m+1}(m+1)^{-1}}{(b-a)B(m,n)} \times$$

$$\times \sum_{i=0}^{\infty} \frac{(1-n)_i(m+1)_i b^i}{(m+2)_i i!} \int_0^{1/b} \frac{((1/b)-z)^{m+i}}{(z-(b+1)/b)} dz.$$

Вспользуемся формулой 3.196.1 из [5]:

$$U_1 = \frac{b^{m+1}(m+1)^{-1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i(m+1)_i b^i}{(m+2)_i i!} \times$$

$$\times \frac{b}{b+1} \left(\frac{1}{b}\right)^{m+i+1} \frac{G(1, 1, m+i+2, 1/(b+1))}{m+i+1} =$$

$$= \frac{b}{(b-a)B(m,n)} \times$$

$$\times \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(m+i+1)^{-2} j!}{i!(m+i+2)_j} \left(\frac{1}{b+1}\right)^{j+1}.$$

Аналогично

$$U_2 = \frac{a^{m+1}}{(b-a)B(m,n)} \times$$

$$\times \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(m+i+1)^{-2} a^i j!}{i!(m+i+2)_j} \times$$

$$\times \left(\frac{1}{b}\right)^{m+i} \left(\frac{1}{b+1}\right)^{j+1}.$$

Аналогично получаем для второго слагаемого из (4):

$$\int_{1/(b+1)}^{1/(a+1)} x S_2(x) dx = \frac{B(m+1,n)}{(b-a)B(m,n)} \ln \left(\frac{b+1}{a+1}\right) -$$

$$- \int_{1/(b+1)}^{1/(a+1)} \frac{(m+1)^{-1} x^{-1}}{(b-a)B(m,n)} \left(\frac{ax}{1-x}\right)^{m+1} \times$$

$$\times G\left(1-n, m+1, m+2; \frac{ax}{1-x}\right) dx \equiv U_3 - U_4.$$

Разбив вычитаемое в последнем выражении на два интеграла, получим

$$U_4 = \int_0^{1/(a+1)} \frac{(m+1)^{-1} x^{-1}}{(b-a)B(m,n)} \left(\frac{ax}{1-x}\right)^{m+1} \times$$

$$\times G\left(1-n, m+1, m+2; \frac{ax}{1-x}\right) dx - U_2 =$$

$$= \frac{a}{(b-a)B(m,n)} \times$$

$$\times \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(m+i+1)^{-2} j!}{i!(m+i+2)_j} \left(\frac{1}{a+1}\right)^{j+1} - U_2.$$

Подытоживая все вспомогательные выкладки, получаем для (4)

$$E p = U_1 - U_2 + U_3 - U_4 =$$

$$= \frac{B(m+1,n)}{(b-a)B(m,n)} \ln \left(\frac{b+1}{a+1}\right) + \frac{b}{(b-a)B(m,n)} \times$$

$$\times \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(m+i+1)^{-2} j!}{i!(m+i+2)_j} \left(\frac{1}{b+1}\right)^{j+1} -$$

$$- \frac{a}{(b-a)B(m,n)} \times$$

$$\times \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(m+i+1)^{-2} j!}{i!(m+i+2)_j} \left(\frac{1}{a+1}\right)^{j+1}. \quad (5)$$

Используя элементарные преобразования для (5), по определению (2) получаем (3), что завершает доказательство теоремы.

Теперь рассмотрим симметричный случай для распределений случайных величин λ и μ .

Теорема 2. Пусть случайные величины λ и μ независимы и имеют соответственно распределения $\beta(m, n)$ и $R(a, b)$, где $0 \leq a < b \leq 1$, $m, n > 0$. Тогда

$$p_{\text{сред}} = \frac{B(m+1,n)}{B(m,n)} \left(1 + \frac{1}{b-a} \ln \left(\frac{a+1}{b+1}\right)\right) +$$

$$+ \frac{bm^{-1}(m+1)^{-1}}{(b-a)B(m,n)} \times$$

$$\times G_{1,0}^{2,1} \left(m, 1-n, m+1, 1; m+1, m+2; 1, -\frac{1}{b}\right) -$$

$$- \frac{am^{-1}(m+1)^{-1}}{(b-a)B(m,n)} \times$$

$$\times G_{1,0}^{2,1} \left(m, 1-n, m+1, 1; m+1, m+2; 1, -\frac{1}{a}\right). \quad (6)$$

Доказательство. Найдем плотность $f_p(x)$ случайной величины p . Имеем

$$f_p(x) = \int_0^1 \frac{y}{(1-x)^2} f_\mu\left(\frac{x}{1-x}y\right) f_\lambda(y) dy.$$

По формуле 8.391 из [5] для $a/(a+1) < x < b/(b+1)$ имеем

$$\begin{aligned} f_p(x) &= \int_{a(1-x)/x}^1 \frac{(1-x)^{-2}}{(b-a)B(m,n)} y^m(1-y)^{n-1} dy = \\ &= \frac{B(m+1,n)(1-x)^{-2}}{(b-a)B(m,n)} - \\ &- \frac{a^{m+1}(m+1)^{-1}(1-x)^{m-1}}{(b-a)B(m,n)x^{m+1}} \times \\ &\times G\left(1-n, m+1, m+2, \frac{a(1-x)}{x}\right) \equiv T_1(x), \end{aligned}$$

а для $b/(b+1) \leq x < 1$

$$\begin{aligned} f_p(x) &= \frac{(1-x)^{-2}}{(b-a)B(m,n)} \left(\int_0^{b(1-x)/x} y^m(1-y)^{n-1} dy - \right. \\ &- \left. \int_0^{a(1-x)/x} y^m(1-y)^{n-1} dy \right) = \\ &= \frac{b^{m+1}(m+1)^{-1}(1-x)^{m-1}}{(b-a)B(m,n)x^{m+1}} \times \\ &\times G\left(1-n, m+1, m+2, \frac{b(1-x)}{x}\right) - \\ &- \frac{a^{m+1}(m+1)^{-1}(1-x)^{m-1}}{(b-a)B(m,n)x^{m+1}} \times \\ &\times G\left(1-n, m+1, m+2, \frac{a(1-x)}{x}\right) \equiv T_2(x). \end{aligned}$$

Таким образом,

$$E_p = \int x f_p(x) dx = \int_{a/(a+1)}^{b/(b+1)} x T_1(x) dx + \int_{b/(b+1)}^1 x T_2(x) dx. \quad (7)$$

Вычислим отдельно первый интеграл из правой части (7). Имеем

$$\begin{aligned} \int_{a/(a+1)}^{b/(b+1)} x T_1(x) dx &= \\ &= \frac{B(m+1,n)}{B(m,n)} \left(1 + \frac{1}{b-a} \ln\left(\frac{a+1}{b+1}\right) \right) - \end{aligned}$$

$$- \int_{a/(a+1)}^{b/(b+1)} \frac{a^{m+1}(m+1)^{-1}(1-x)^{m-1}}{(b-a)B(m,n)x^m} \times$$

$$\times G\left(1-n, m+1, m+2, \frac{a(1-x)}{x}\right) dx \equiv V_1 - V_2.$$

Для вычисления V_2 воспользуемся формулой 3.196.1 из [5]. Имеем

$$\begin{aligned} V_2 &= \frac{a^{m+1}(m+1)^{-1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i(m+1)_i a^i}{(m+2)_i i!} \times \\ &\times \int_{a/(a+1)}^{b/(b+1)} \frac{(1-x)^{m+i-1}}{x^{m+i}} dx = \\ &= \frac{a^{m+1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i a^i}{(m+i+1)!} \times \\ &\times \int_{1/(b+1)}^{1/(a+1)} \frac{y^{m+i-1}}{(1-y)^{m+i}} dy = \\ &= \frac{a^{m+1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i a^i}{(m+i+1)!} \times \\ &\times \left(\int_0^{1/(a+1)} \frac{(1/(a+1)-x)^{m+i-1}}{(x+a/(a+1))^{m+i}} dx - \right. \\ &- \left. \int_0^{1/(b+1)} \frac{(1/(b+1)-x)^{m+i-1}}{(x+b/(b+1))^{m+i}} dx \right) = \\ &= \frac{a^{m+1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i a^i}{(m+i+1)!} \times \\ &\times \left(\frac{a^{-m-i}}{m+i} G\left(1, m+i, m+i+1, -\frac{1}{a}\right) - \frac{b^{-m-i}}{m+i} \times \right. \\ &\times \left. G\left(1, m+i, m+i+1, -\frac{1}{b}\right) \right) = \\ &= \frac{a^{m+1}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \frac{(1-n)_i a^i}{(m+i+1)!} \times \\ &\times \left(\sum_{j=0}^{\infty} \frac{a^{-m-i}(-a)^{-j}}{m+i+j} - \sum_{j=0}^{\infty} \frac{b^{-m-i}(-b)^{-j}}{m+i+j} \right). \end{aligned}$$

Для второго слагаемого в (7) аналогично имеем

$$\begin{aligned} \int_{b/(b+1)}^1 x T_2(x) dx &= \int_{b/(b+1)}^1 \frac{b^{m+1}(m+1)^{-1}(1-x)^{m-1}}{(b-a)B(m,n)x^m} \times \\ &\times G\left(1-n, m+1, m+2, \frac{b(1-x)}{x}\right) dx - \end{aligned}$$

$$\begin{aligned}
 & - \int_{b/(b+1)}^1 \frac{a^{m+1}(m+1)^{-1}(1-x)^{m-1}}{(b-a)B(m,n)x^m} \times \\
 & \times G\left(1-n, m+1, m+2, \frac{a(1-x)}{x}\right) dx \equiv V_3 - V_4, \\
 & + \frac{b}{(b-a)B(m,n)} \times \\
 & \times \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(-b)^{-j}}{(m+i+j)(m+i+1)i!}. \quad (8)
 \end{aligned}$$

где

$$\begin{aligned}
 V_3 &= \\
 &= \frac{b}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(-b)^{-j}}{(m+i+j)(m+i+1)i!};
 \end{aligned}$$

$$\begin{aligned}
 V_4 &= \\
 &= \frac{a^{m+1}b^{-m}}{(b-a)B(m,n)} \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i a^i b^{-i} (-b)^{-j}}{(m+i+j)(m+i+1)i!}.
 \end{aligned}$$

Объединяя полученные результаты, получаем из (7)

$$\begin{aligned}
 E_p &= V_1 - V_2 + V_3 - V_4 = \\
 &= \frac{B(m+1, n)}{B(m, n)} \left(1 + \frac{1}{b-a} \ln\left(\frac{a+1}{b+1}\right)\right) - \\
 &- \frac{a}{(b-a)B(m, n)} \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} \frac{(1-n)_i(-a)^{-j}}{(m+i+j)(m+i+1)i!} +
 \end{aligned}$$

Используя элементарные преобразования для (8), по определению (1) получаем (6), что завершает доказательство теоремы.

Замечание 1. При $a = 0$, очевидно, выражение (6) остается справедливым, если последнее слагаемое положить равным нулю.

Замечание 2. Выражения (3) и (6) служат для компактной записи $p_{\text{сред}}$. Для практического использования (непосредственного вычисления) имеет смысл представлять $p_{\text{сред}}$ в виде рядов типа (5) и (8), которые несложно вычисляются с любой наперед заданной точностью.

В качестве иллюстрации приведем несколько таблиц со значениями $p_{\text{сред}}$. Для удобства сравнения с результатами, опубликованными в [3] и [4], значения параметров равномерного и бета-распределения взяты из этих статей. Таблицы 1 и 2 соответствуют распределениям λ и μ из теоремы 1, а табл. 3 и 4 соответствуют теореме 2.

Таблица 1 Частные значения средней надежности ($\lambda \sim R(a, b)$, $\mu \sim \beta(m, n)$)

[a, b]	m; n									
	1; 7	1; 3	3; 5	1; 1	6; 10	2; 2	10; 6	5; 3	3; 1	7; 1
[0, 1/4]	0,47	0,60	0,74	0,75	0,75	0,78	0,84	0,83	0,85	0,88
[0, 1/2]	0,34	0,48	0,61	0,63	0,62	0,66	0,73	0,73	0,76	0,79
[0, 3/4]	0,28	0,40	0,53	0,56	0,54	0,58	0,65	0,65	0,68	0,72
[0, 1]	0,24	0,35	0,47	0,50	0,48	0,52	0,59	0,59	0,62	0,66
[1/4, 1/2]	0,22	0,35	0,48	0,52	0,49	0,54	0,62	0,62	0,66	0,70
[1/4, 3/4]	0,19	0,30	0,42	0,46	0,43	0,48	0,56	0,55	0,60	0,64
[1/4, 1]	0,16	0,27	0,38	0,42	0,38	0,44	0,51	0,51	0,55	0,59
[1/2, 3/4]	0,15	0,26	0,36	0,40	0,37	0,42	0,50	0,49	0,54	0,58
[1/2, 1]	0,13	0,23	0,32	0,37	0,33	0,38	0,46	0,45	0,49	0,54
[3/4, 1]	0,12	0,20	0,29	0,33	0,29	0,35	0,41	0,41	0,45	0,50

Таблица 2 Частные значения средней надежности ($\lambda \sim R(a, b)$, $\mu \sim \beta(m, n)$)

[a, b]	m; n									
	1; 7	1, 1; 3, 3	1, 2; 2, 0	1, 3; 1, 3	1, 4; 2, 33	1, 5; 1, 5	1, 6; 0, 96	1, 7; 1, 02	1, 8; 0, 6	1, 9; 0, 27
[0, 1/4]	0,47	0,61	0,70	0,76	0,71	0,77	0,81	0,81	0,85	0,87
[0, 1/2]	0,34	0,48	0,57	0,64	0,58	0,65	0,70	0,71	0,75	0,78
[0, 3/4]	0,28	0,41	0,50	0,57	0,50	0,57	0,63	0,63	0,67	0,71
[0, 1]	0,24	0,36	0,44	0,51	0,45	0,51	0,57	0,57	0,62	0,66
[1/4, 1/2]	0,22	0,35	0,45	0,53	0,46	0,53	0,60	0,60	0,65	0,69
[1/4, 3/4]	0,19	0,31	0,40	0,47	0,40	0,47	0,54	0,54	0,59	0,63
[1/4, 1]	0,16	0,27	0,36	0,43	0,36	0,43	0,49	0,49	0,54	0,59
[1/2, 3/4]	0,15	0,26	0,34	0,41	0,35	0,42	0,48	0,48	0,53	0,58
[1/2, 1]	0,13	0,23	0,31	0,38	0,31	0,38	0,44	0,44	0,49	0,53
[3/4, 1]	0,12	0,20	0,28	0,34	0,28	0,34	0,40	0,40	0,45	0,49

Таблица 3 Частные значения средней надежности ($\lambda \sim \beta(m, n), \mu \sim R(a, b)$)

$m; n$	$[a, b]$									
	$[0, 1/4]$	$[0, 1/2]$	$[0, 3/4]$	$[0, 1]$	$[1/4, 1/2]$	$[1/4, 3/4]$	$[1/4, 1]$	$[1/2, 3/4]$	$[1/2, 1]$	$[3/4, 1]$
1; 7	0,53	0,79	0,81	0,74	0,72	0,74	0,75	0,71	0,66	0,67
1; 3	0,40	0,52	0,57	0,65	0,73	0,69	0,75	0,77	0,78	0,75
3; 5	0,09	0,39	0,72	0,52	0,69	0,73	0,74	0,83	0,85	0,27
1; 1	0,25	0,35	0,43	0,50	0,51	0,55	0,58	0,61	0,63	0,66
6; 10	0,25	0,38	0,46	0,54	0,70	0,73	0,74	0,83	0,85	0,86
2; 2	0,23	0,32	0,41	0,48	0,39	0,47	0,55	0,59	0,62	0,64
10; 6	0,16	0,27	0,35	0,41	0,55	0,54	0,55	0,65	0,67	0,74
5; 3	0,17	0,26	0,39	0,41	0,55	0,50	0,48	0,87	0,18	0,87
3; 1	0,15	0,25	0,32	0,38	0,33	0,41	0,44	0,49	0,51	0,54
7; 1	0,13	0,22	0,27	0,34	0,30	0,36	0,38	0,42	0,47	0,49

Таблица 4 Частные значения средней надежности ($\lambda \sim \beta(m, n), \mu \sim R(a, b)$)

$m; n$	$[a, b]$									
	$[0, 1/4]$	$[0, 1/2]$	$[0, 3/4]$	$[0, 1]$	$[1/4, 1/2]$	$[1/4, 3/4]$	$[1/4, 1]$	$[1/2, 3/4]$	$[1/2, 1]$	$[3/4, 1]$
1; 7	0,36	0,79	0,81	0,74	0,72	0,74	0,75	0,71	0,69	0,68
1, 1; 3, 3	0,53	0,57	0,61	0,60	0,78	0,67	0,66	0,69	0,71	0,69
1, 2; 2, 0	0,34	0,50	0,54	0,58	0,62	0,59	0,63	0,68	0,69	0,75
1, 3; 1, 3	0,25	0,43	0,48	0,50	0,45	0,53	0,53	0,54	0,60	0,59
1, 4; 2, 33	0,37	0,45	0,56	0,50	0,50	0,56	0,59	0,60	0,64	0,63
1, 5; 1, 5	0,23	0,35	0,42	0,49	0,46	0,53	0,57	0,58	0,60	0,65
1, 6; 0, 96	0,20	0,30	0,36	0,44	0,42	0,46	0,51	0,54	0,59	0,62
1, 7; 1, 02	0,15	0,29	0,35	0,41	0,40	0,45	0,49	0,52	0,58	0,59
1, 8; 0, 6	0,16	0,28	0,32	0,42	0,34	0,40	0,45	0,45	0,49	0,51
1, 9; 0, 27	0,13	0,21	0,28	0,37	0,35	0,39	0,40	0,40	0,46	0,48

3 Заключение

Полученные результаты могут применяться, например, для вычисления других моментов и построения доверительных интервалов для характеристики p . В дальнейшем предполагается расширить класс совместных распределений параметров (λ, μ) , разработать соответствующие расчетные алгоритмы для вычисления величины $p_{\text{сред}}$ и провести тестовые расчеты.

Литература

1. Gnedenko B. V., Korolev V. Yu. Random summation: Limit theorems and applications. — Boca Raton, FL: CRC Press, 1996. 288 p.
2. Королев В. Ю., Соколов И. А. Основы математической теории надежности модифицируемых систем. — М.: ИПИ РАН, 2006. 108 с.
3. Кудрявцев А. А., Соколов И. А., Шоргин С. Я. Байесовская рекуррентная модель роста надежности: равномерное распределение параметров // Информатика и её применения, 2013. Т. 7. Вып. 2. С. 55–59.
4. Жаворонкова Ю. В., Кудрявцев А. А., Шоргин С. Я. Байесовская рекуррентная модель роста надежности: бета-распределение параметров // Информатика и её применения, 2014. Т. 8. Вып. 2. С. 48–54.
5. Градштейн И. С., Рыжик И. М. Таблицы интегралов, сумм, рядов и произведений. — М.: Наука, 1971. 1108 с.

Поступила в редакцию 26.01.15

BAYESIAN RECURRENT MODEL OF RELIABILITY GROWTH: BETA-UNIFORM DISTRIBUTION OF PARAMETERS

Iu. V. Zhavoronkova¹, A. A. Kudryavtsev², and S. Ya. Shorgin³

¹Sputnik Ltd., 8/2 Prishvina Str., Moscow 127549, Russian Federation

²Faculty of Computational Mathematics and Cybernetics, M. V. Lomonosov Moscow State University, 1-52 Leninskiye Gory, GSP-1, Moscow 119991, Russian Federation

³Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation

Abstract: Forecasting reliability of complex modifiable information systems is one of the topical problems of the mass service theory nowadays. Any first established complex system designed for processing or transmission of information flows, as a rule, does not possess the required reliability. Such systems are subject to modifications during development, testing, and regular functioning. The purpose of such modifications is to increase reliability of information systems. In this connection, there is a necessity to formalize the concept of reliability of modifiable information systems and to develop methods and algorithms of estimation and forecasting of various reliability characteristics. One approach to determine system reliability is to compute the probability that the signal fed to the input of the system at a given point of time will be reacted to correctly by the system. The article considers the exponential recurrent growth model of reliability, in which the probability of system reliability is represented as a linear combination of “defectiveness” and “efficiency” parameters of tools correcting the deficiencies in the system. It is assumed that the researcher does not have exact information about the system under study and is only familiar with the characteristics of the class from which this system is taken. In the framework of the Bayesian approach, it is assumed that one of the indicators of “defectiveness” and “efficiency” has the beta-distribution and the other one has the uniform distribution. Average marginal system reliability is calculated. Numerical results for model examples are obtained.

Keywords: modifiable information systems; theory of reliability; Bayesian approach; beta-distribution; uniform distribution

DOI: 10.14357/19922264150109

Acknowledgments

This work was financially supported by the Russian Science Foundation (grant No. 14-11-00397).

References

1. Gnedenko, B. V., and V. Yu. Korolev. 1996. *Random summation: Limit theorems and applications*. Boca Raton, FL: CRC Press, 1996. 288 p.
2. Korolev, V. Yu., and I. A. Sokolov. 2006. *Osnovy matematicheskoy teorii nadezhnosti modifitsiruemykh system* [Fundamentals of mathematical theory of modified systems reliability]. Moscow: IPI RAN, 2006. 108 p.
3. Kudryavtsev, A. A., I. A. Sokolov, and S. Ya. Shorgin. 2013. Bayesovskaya rekurrentnaya model' rosta nadezhnosti: Ravnomernoe raspredelenie parametrov [Bayesian recurrent model of reliability growth: Uniform distribution of parameters]. *Informatika i ee Primeneniya — Inform. Appl.* 7(2):55–59.
4. Zhavoronkova, Iu. V., A. A. Kudryavtsev, and S. Ya. Shorgin. 2013. Bayesovskaya rekurrentnaya model' rosta nadezhnosti: Beta-raspredelenie parametrov [Bayesian recurrent model of reliability growth: Beta-distribution of parameters]. *Informatika i ee Primeneniya — Inform. Appl.* 8(2):48–54.
5. Gradshteyn, I. S., and I. M. Ryzhik. 1971. *Tablitsy integralov, summ, ryadov i proizvedeniy* [Tables of integrals, sums, series, and products]. Moscow: Nauka. 1108 p.

Received January 26, 2015

Contributors

Zhavoronkova Iuliia V. (b. 1990) — software developer, Sputnik Ltd., 8/2 Prishvina Str., Moscow 127549, Russian Federation; juliana-zh@yandex.ru

Kudryavtsev Alexey A. (b. 1978) — Candidate of Science (PhD) in physics and mathematics, associate professor, Faculty of Computational Mathematics and Cybernetics, M. V. Lomonosov Moscow State University, 1-52 Leninskiye Gory, GSP-1, Moscow 119991, Russian Federation; nubigena@mail.ru

Shorgin Sergey Ya. (b. 1952) — Doctor of Science in physics and mathematics, professor, Deputy Director, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; sshorgin@ipiran.ru

К ОЦЕНКЕ ЭФФЕКТИВНОСТИ УЧЕБНО-ПОЗНАВАТЕЛЬНОЙ ДЕЯТЕЛЬНОСТИ УЧАЩИХСЯ С ИСПОЛЬЗОВАНИЕМ ИНФОРМАЦИОННЫХ ТЕХНОЛОГИЙ

О. М. Корчажкина¹

Аннотация: Рассматриваются проблемы измерения эффективности учебно-познавательной деятельности (УПД) учащихся как показателя соответствия планируемых и достигнутых образовательных результатов. Этот показатель выражается в терминах конкретных продуктов УПД, получаемых в ходе выполнения мыслительных операций. Обсуждается вопрос совмещения стиля учения и методов обучения в условиях интеграции педагогических и новых информационных технологий при выполнении заданий различных типов. Приводится пример вербализации достигнутых результатов УПД с использованием мобильных устройств, основанный на таксономии мыслительных операций Бенджамина Блума. Установлено, что уровень эффективности УПД с использованием информационно-коммуникационных технологий (ИКТ) определяется способностью учителя организовать совместную работу с учащимися, ориентированную на развитие форм мыслительной деятельности, приводящих к созданию интегрированного персонального познавательного стиля каждого.

Ключевые слова: эффективность обучения; планируемые образовательные результаты; достигнутые образовательные результаты; мобильные устройства; мыслительные операции; индивидуальный стиль учения; методы обучения; LOA-технология

DOI: 10.14357/19922264150110

1 Введение

Эффективность обучения как соотношение полезного результата и затратных факторов образовательного процесса может выражать различные его стороны, поэтому, когда встает вопрос об измерении эффективности использования ИКТ в образовании, необходимо определиться, какую эффективность следует рассматривать и что конкретно понимать под словом «эффективность»:

- эффективность использования ИКТ (или, более узко, электронных образовательных ресурсов — ЭОР) в учебном процессе;
- эффективность учебного процесса с использованием ИКТ;
- эффективность УПД учащихся с использованием ИКТ.

Очевидно, что следует делать акцент не на экономическую эффективность внедрения ИКТ, т.е. считать затратами при достижении определенного результата не материальные вложения государства в информатизацию образования или затраты конкретного учебного заведения на приобретение компьютеров, поддержку сетей и оплату труда педагогов или технического персонала. За скобками

следует оставить также временные и психологические затраты учителя и учащихся и рассматривать эффективность только как **уровень достижения образовательного результата**.

Измерение эффективности в первой трактовке — как уровня достижения образовательного результата при использовании ИКТ или ЭОР в учебном процессе — должно осуществляться с привлечением очень многих показателей, например, согласно [1, с. 298], *ценности учебного материала, мотивации и компетентности учащихся*.

Среди предложенных авторами [1] показателей эффективности особые сложности для измерения вызывает *ценность учебного материала*, которую, основываясь на положениях теории информации Шеннона и теории статистических решений, где есть показатель «ценность информации», можно вслед за самим Клодом Шенноном определить как «максимальную пользу, которую данное количество информации способно принести в деле уменьшения средних потерь». Однако в рассматриваемом случае под потерями следует, видимо, подразумевать неусвоенное знание, недополученные навыки и умения, которые опять же очень сложно измерить количественно. Кроме того, такой показатель, как «количество учебной информации» в определении Шеннона, нужно корректно оценить при по-

¹Институт проблем информатики Российской академии наук, olgakomax@gmail.com

мощи соответствующих статистических методов — распределения Шеннона, Хартли или Больцмана. Таким образом, первая трактовка «эффективность использования ИКТ в учебном процессе» неоднозначна, а сама эффективность в этом понимании сложна для измерения и, следовательно, не может быть принята как практический вариант.

Вторая трактовка эффективности — эффективности учебного процесса с использованием ИКТ — слишком сложная и многоаспектная категория, и ее оценивание лежит в области так называемых нетривиальных педагогических измерений, что требует привлечения сложного математического аппарата.

Третий вариант трактовки эффективности, напрямую связанный с результатами УПД учащихся с использованием ИКТ, лежит в плоскости практической деятельности учителя и может быть с успехом им реализован, поскольку связан непосредственно с организацией учебного процесса, ориентированного на достижение учащимися планируемых образовательных результатов, декларированных в Федеральных государственных стандартах (далее — ФГОС) второго поколения [2].

2 Эффективность обучения как соотношение планируемых и достигнутых образовательных результатов

В процедуре оценивания результатов УПД учащихся различают три уровня достижения этих результатов:

- (1) планируемый уровень L_P — тот, что декларируется во ФГОС второго поколения и находит воплощение в учебниках и учебно-методических пособиях;
- (2) реализуемый уровень L_R , характеризующий результаты, определяемые учителем в зависимости от своих профессиональных предпочтений и условий обучения;
- (3) достигнутый уровень L_A — уровень объективных, реальных достижений учащихся.

В связи с этим требуется определить, разницу между какими уровнями достижения образовательных результатов: $\delta_{PR} = (L_P - L_R)/L_P$, $\delta_{RA} = (L_R - L_A)/L_R$ или $\delta_{PA} = (L_P - L_A)/L_P$ — необходимо минимизировать, чтобы можно было говорить об объективном измерении эффективности УПД.

Очевидно, что наиболее достоверные результаты измерения эффективности УПД дает показатель δ_{PA} , поскольку он учитывает приближенность

уровня объективных достижений учащихся к уровню объективного абсолюта, т. е. того предела, на который ориентируют всех участников образовательного процесса нормативные документы федерального уровня. Именно поэтому при измерении эффективности УПД учащихся целесообразно говорить об *оценке достижения планируемых образовательных результатов*.

Таким образом, система оценки достижения планируемых образовательных результатов, с одной стороны, направлена на реализацию требований ФГОС, а с другой — способствует объективному измерению эффективности УПД на основе реальных достижений учащихся.

Основной функцией системы оценки, ориентирующей образовательный процесс на достижение планируемых результатов, является обеспечение результативной обратной связи, позволяющей осуществлять управление образовательным процессом, особенно в части принятия педагогических мер, которые способствовали бы повышению эффективности УПД учащихся [3, с. 133].

В систему оценки достижения планируемых образовательных результатов включаются следующие необходимые компоненты [3, с. 133]:

- формулировка основных направлений и целей оценочной деятельности;
- описание объектов и содержания оценки;
- задание критериев, описание процедур и состава инструментария оценивания;
- описание форм представления результатов, условий и границ применения системы оценивания;
- привлечение разнообразных методов и форм оценивания, взаимно дополняющих друг друга.

3 Роль информационно-коммуникационных технологий в повышении эффективности учебно-познавательной деятельности учащихся

С расширением процесса информатизации образования и внедрением в учебный процесс новых ИКТ ожидалось повышение его эффективности. Однако скоро стало очевидным, что использование ИКТ в русле традиционного обучения не столь существенно влияет на уровень обученности, т. е. на эффективность обучения, как предполагалось ранее [4, с. 22, 23], однако значительно

повышает энергоёмкость труда учителя за счет необходимости освоения новых техник и технологий в сжатые сроки. Более того, строгие рамки постоянного мониторинга учебного процесса, в которые поставлены учителя при освоении новых ИКТ, а также повышение интенсивности труда учителя на первых этапах процесса информатизации образования приводили многих учителей к резкому неприятию инновационной педагогической деятельности в области информатизации образования, а часть из них — даже к фрустрации и профессиональному выгоранию.

По прошествии последних десяти лет информатизации российского образования, когда большинство учителей-предметников средней школы вольно или невольно с разной степенью успешности включились в этот процесс, по-прежнему остается без ответа вопрос: почему же ИКТ не повлияли существенным образом на эффективность обучения?

Для того чтобы понять это, нужно проследить, к каким изменениям в учебном процессе привели новые ИКТ с точки зрения технической, социально-педагогической и психолого-педагогической.

Техническая сторона информатизации заключается в обновлении как аппаратного, так и программного обеспечения сферы образования, причем в последние несколько лет — 2010–2014 годы — идет процесс так называемой *электронизации* [5, с. 167], которая характеризуется распространением мобильных электронных устройств различного типа и назначения, использованием мощных персональных компьютеров, быстродействующих накопителей большой емкости, облачных серверов, новых информационных и телекоммуникационных технологий, мультимедиа-технологий и виртуальной реальности, появлением 3D-принтеров, первого поколения электронных учебников.

С *социально-педагогической* точки зрения изменения произошли как в той роли, которую стали играть в учебном процессе учитель и учащийся, так и в изменении ожиданий и приоритетов учащихся по отношению к образовательному процессу.

Если при традиционном обучении учитель является единственным поставщиком готовых знаний и контролером их усвоения, то в ходе внедрения ИКТ в учебный процесс все существеннее ощущается крен в сторону автономии учащихся. При этом преобладающее влияние приобретают следующие формы обучения:

- регулируемое (направляемое) обучение (учитель — консультант, навигатор, содействующий целенаправленной УПД учащихся по освоению планируемых компетенций);

- самообучение (учитель — тьютор, модератор, в обязанности которого входит как содействие раскрытию потенциальных способностей учащихся, так и налаживание контактов между ними для организации совместной работы);
- саморегулируемое обучение (учитель — фасилитатор, работающий в парадигме личностно-ориентированной педагогики и способствующий наиболее эффективному учебному взаимодействию) [6, с. 408–424].

Получая большую автономию, учащиеся самостоятельно начинают выбирать приоритеты, среди которых в плане формы обучения основное место занимают:

- мобильность — желание получать знания не только в учебной аудитории в рамках классно-урочной системы;
- планшетизация — использование планшетных и иных мобильных устройств для обучения;
- коллаборация и краудсорсинг — сетевое сотрудничество;
- гейметизация — использование интерактивных игровых технологий в обучении.

Еще в 1990 г. в проекте Концепции информатизации отечественного образования очередного периода [7, с. 4] указывалось, что с *психолого-педагогической* точки зрения ИКТ в образовании способствуют:

- раскрытию, сохранению и развитию индивидуальных способностей обучаемых;
- формированию у учащихся познавательных способностей, стремления к самосовершенствованию;
- обеспечению комплексности изучения явлений действительности, неразрывности взаимосвязи между естественными и гуманитарными науками;
- постоянному динамическому обновлению содержания, форм и методов процесса обучения и воспитания.

Однако за прошедшую почти четверть века мы так и не смогли ответить на главный вопрос: что сделано для того, чтобы ИКТ способствовали раскрытию, сохранению, формированию, обеспечению, обновлению и т.п.? Или это должно было произойти автоматически, без каких-либо усилий со стороны участников образовательного процесса?

Практика показала, что преимущества, которые предоставляют ИКТ не в технической, а, главным образом, в психолого-педагогической сфере, могут привести к коренным изменениям в учебном

процессе. Остается только построить учебный процесс таким образом, чтобы были созданы условия для реализации этих преимуществ. Именно это обстоятельство позволит впоследствии говорить об эффективности использования ИКТ.

Итак, вопросы, на которые должен ответить учитель, заинтересованный в повышении эффективности своей работы средствами новых ИКТ, можно сформулировать следующим образом:

- С внедрением ИКТ в учебный процесс как изменилось мышление учащихся и в каком направлении следует развивать их мыслительные способности?
- С внедрением ИКТ в учебный процесс как изменилась УПД учащихся на уроке и в информационно-образовательной среде (ИОС) — какие новые виды УПД возникают, какие типы заданий предлагать учащимся и как строить урок в классе и в ИОС школы?

Ответы на эти вопросы лежат в плоскости когнитивной психологии, а непременным условием повышения учебных достижений учащихся, т.е. эффективности обучения, специалисты называют проблему совмещения стилей учения учащихся с методами обучения [8, с. 263]. Эта проблема возникает в связи с таким, казалось бы, непреложным фактом, что обучение происходит тем более эффективно, чем более оно соответствует познавательным стилям учащихся, которые называют еще стилями учения¹.

Тогда встает еще один вопрос: а что означает такое соответствие? Означает ли это, что средства обучения — методы обучения, формы предъявления учебного материала, используемые техники и технологии обучения и пр. — должны подстраиваться под каждого конкретного учащегося, тем самым формируя его образовательную траекторию? Или же, наоборот, задача учителя — создать такие условия обучения, такую образовательную среду, в которой каждый учащийся, носитель своего персонального познавательного стиля, не только сможет выбрать свою линию обучения, но и интеллектуально развиваться, осваивая новые для себя способы познания окружающей действительности [8, с. 266, 267]?

¹Стиль учения — это индивидуальная характеристика личности, психическое образование, которое является многомерным по своим проявлениям в различных видах УПД, иерархическим по устройству, включающему разные уровни стилового поведения, интегральным по своим механизмам, являясь продуктом интеграции разных форм индивидуального ментального опыта, и гибким по своим адаптационным возможностям, что способствует формированию интегрированного персонального познавательного стиля [8, с. 269].

²Психологи-когнитивисты различают следующие уровни базовых механизмов стилового поведения познающей личности: уровень стилей кодирования информации, основанных на разных модальностях опыта (кинестетической, визуальной, словесно-речевой, сенсорно-эмоциональной); уровень стилей переработки информации (интеллектуальные, импульсивные, рефлексивные, неуспешные); уровень стилей постановки и решения проблем (вариации в наборе приемов решения задачи от адаптивного к смыслообразующему); уровень стилей познавательного отношения к миру (чувственное, рациональное, сверхчувственное познание и др.) [8, с. 270].

Очевидно, что первый способ не только трудно осуществим, но и приводит к закреплению у учащегося определенного стиля усвоения учебной информации, ограничивая его интеллектуальное развитие. Второй же способ, напротив, стимулирует формирование интегрированного персонального познавательного стиля каждого учащегося и взаимообогащение стилей учащихся при сотрудничестве, что создает условия для их дальнейшего интеллектуального воспитания и развития.

М.А. Холодная, приводя эти рассуждения, ссылается на целесообразность стилового подхода к обучению, ориентирующегося на внутреннюю дифференциацию как одну из двух форм индивидуализации образовательного процесса. В противовес внешней дифференциации, когда производится отбор детей под определенный тип обучения с целью создания гомогенных классов, имеющих однонаправленную специализацию методов обучения, внутренняя дифференциация предполагает «учет индивидуальных познавательных возможностей каждого ребенка в рамках общего для всех гетерогенного образовательного пространства — вариативного с точки зрения своего содержания и видов учебной деятельности (в том числе с использованием современных педагогических и информационных технологий)» [8, с. 265, 266]. И далее: «Правильнее говорить не об учете индивидуальных познавательных стилей детей, а о формировании у каждого ребенка персонального познавательного стиля на основе актуализации и обогащения всех механизмов стилового поведения» [8, с. 271].

Таким образом, для полноценного интеллектуального развития учащихся, т.е. для развития их мыслительных способностей и повышения эффективности обучения, необходимо организовать УПД в такой образовательной среде, которая была бы вариативной за счет многообразных инструментов, формирующих недостающие механизмы стилового поведения учащихся².

Очевидно, что организовать и обеспечить развитие многообразных познавательных стилей и механизмов стилового поведения учащихся при традиционном обучении весьма трудно, поскольку

автономия учащихся ограничена, существует недостаток в новейших интерактивных учебных материалах и средствах обучения, а учитель и учащиеся поставлены в жесткие рамки классно-урочной системы. Тогда как современная ИОС, предоставляя учащимся разнообразные мобильные технологические инструменты — от приложений для мобильных устройств (см. разд. 4) до электронных учебников нового поколения, создает условия как для реального, так и для продуктивного виртуального учебного взаимодействия. Эти функции ИОС позволяют учителю направить УПД учащихся в русло развития многообразных форм их мыслительной деятельности, что необходимо приводит к формированию интегрированного персонального познавательного стиля каждого учащегося и в перспективе — к повышению эффективности обучения. При этом следует учитывать следующие уровни интеграции информационных и педагогических технологий при осуществлении учащимися УПД с использованием ИКТ.

Уровень 1. Занятия смешанного типа, когда средствами ИКТ вводится новый материал, осуществляется его отработка и контроль усвоения. Учитель использует отдельные элементы готовых или авторских ЭОР и неинтерактивные интернет-ресурсы для визуализации традиционной работы в классе без привлечения учащихся к непосредственной работе с ЭОР.

Уровень 2. Занятия смешанного типа, когда средствами ИКТ вводится новый материал, осуществляется его отработка и контроль усвоения. Электронные образовательные ресурсы и неинтерактивные интернет-ресурсы используются как учителем, так и учащимися для иллюстрации учебного материала и в виде справочных источников (в том числе онлайн-словарей, предметных справочников и энциклопедий, языковых корпусов и онтологий данных); для работы в поисковых системах, проведения онлайн-тестирования и опросов; с привлечением личных электронных учебных блокнотов и заметок при традиционной работе в классе и/или дома.

Уровень 3. Занятия смешанного типа с использованием электронных конструкторов, виртуальных сред и/или интернет-сервисов Web 2.0, позволяющих осуществлять простейшую визуализацию и преобразование учебного материала (определение зависимостей, отношений, построение чертежей, диаграмм, графиков, создание образов, статистическая обработка данных, оформление в виде электронных таблиц, интеллект-карт, облака ключевых слов, интерактивных рабочих листов, электронных

каталогов понятий) с целью его усвоения при традиционной работе в классе и/или дома.

Уровень 4. Занятия смешанного типа, включающие как обязательный компонент совместную работу учащихся в учебных сетевых сообществах с использованием интернет-сервисов Web 2.0 и/или приложений для мобильных устройств (электронная стена; сервис для создания рабочих групп; пространство для создания заметок и совместной работы с ними в группе), позволяющих осуществлять простейшие преобразования учебного материала с целью достижения коллективного учебного результата.

Уровень 5. Занятия смешанного типа в среде программно-методических комплексов в виде виртуальных предметных сред (лабораторий и сред, позволяющих осуществлять алгоритмизацию и моделирование изучаемых явлений и процессов по данному предмету с использованием встроенных функций системы) при индивидуальной или совместной работе учащихся в классе и/или дома, использование мобильных устройств и дистанционных многофункциональных приложений для совместной работы над проектами, а также дистанционных технологий, в том числе для проведения видеоконференций.

Уровень 6. Комплексные задания 3-, 4- и 5-го уровней в ИОС учебного заведения, универсальных рабочих пространствах (в том числе ИОС электронных учебников) и интегральных образовательных платформах Web 2.0.

Приведенная классификация показывает, что чем выше уровень интеграции педагогических и информационных технологий, тем более широкие возможности предоставляются учащимся для формирования и развития интегрированного персонального познавательного стиля, что, как следствие, способствует повышению эффективности УПД с использованием ИКТ.

4 Вербализация целей учебно-познавательной деятельности учащихся

Выше за показатель эффективности обучения был принят уровень соответствия планируемых и достигнутых образовательных результатов (минимум показателя $\delta_{РА}$ — см. разд. 2). Следовательно, эффективность УПД с использованием ИКТ должна выражаться, по меньшей мере, в терминах достигнутых образовательных результатов, т. е. ее конкретных продуктов, получаемых в ходе

выполнения мыслительных операций. Это будет задавать вектор оценочной деятельности учителя, работающего в русле компетентностного и системно-деятельностного подхода. Поэтому учителю при определении эффективности УПД учащихся необходимо ориентироваться на конечные результаты этой деятельности, которые будут понятны ему и учащимся. Кроме того, необходима формулировка целей обучения, а они при компетентностном и системно-деятельностном подходе как раз совпадают с результатами этой деятельности.

Как выразить эти цели-результаты таким образом, чтобы можно было не только однозначно сопоставить то, что было запланировано и что достигнуто, но и выработать адекватные критерии оценки этих достигнутых результатов? Для этого необходима определенного рода вербализация этих целей, т. е. формальное представление их с помощью языковых средств.

В качестве удобной последовательности когнитивных педагогических целей, поддающейся формализации на вербальном уровне, американским психологом методов обучения Бенджамином Блумом в 1956 г. была предложена таксономия (иерархия) этих целей в виде перечня мыслительных (когнитивных) операций, или умственных действий, представленная в виде пирамиды [9] (рис. 1), в основании которой «лежит» *знание* как базовый уровень, а самой высшей по степени сложности и развития мыслительной операцией является *оценка*, находящаяся на шестом уровне.

В приведенной таблице (см. с. 112) для каждой мыслительной операции дается набор смысловых глаголов, соответствующих различным учебным задачам. Пользуясь этой таблицей, учитель может соотносить название мыслительной операции с их содержанием, обеспечить концентрацию усилий на главных аспектах УПД, наметить первоочередные задачи и перспективы дальнейшей работы, создать



Рис. 1 Пирамида Блума

возможности для разъяснения учащимся ориентиров УПД, сформировать эталоны оценки результатов обучения, обеспечивающие надежность и объективность [11].

Таким образом, глаголы мыслительных операций, с одной стороны, помогают учителю в постановке целей и задач УПД, а с другой — позволяют производить описание и оценку результатов этой деятельности в сравнении с поставленными целями. С их помощью учитель выявляет наличие и характер отклонений от запланированных целей УПД на основе образовательного мониторинга, определяет их причины [3, с. 134] и вносит соответствующие коррективы.

В русле развития новых ИКТ в 2000-х гг. появилось так называемое «Педагогическое колесо», которое позволяет установить соответствие между глаголами мыслительных операций Блума, видами УПД учащихся и инструментами мобильных ИКТ [12] (рис. 2), т. е. показывает воплощение процесса УПД в ее продуктивный конкретный результат, достигаемый с помощью современных мобильных инструментов. На рис. 2 во внешней части «Педагогического колеса» приведены приложения для iPad Apple, которые являются инструментами реализации мыслительных операций, собранных в его центральной части.

Например, для раздела «Синтез (Create)» могут быть использованы следующие приложения:

- **Aurasma** — приложение, позволяющее создавать дополненную реальность;
- **Creative Book Builder** — приложение, позволяющее создавать, редактировать и публиковать книги;
- **Easy Release** — приложение для создания и редактирования информационных сообщений;
- **Fotobabble** — приложение для озвучивания изображений (создание «говорящих» фотографий);
- **Garageband** — виртуальный самоучитель игры на музыкальных инструментах;
- **iMovie** — приложение для создания видеороликов;
- **Interview Assistant** — виртуальный микрофон и запись звука;
- **iTimeLapse Pro** — приложение для съемки серии изображений и их компиляции в видео;
- **Nearpod** — интегральная платформа, позволяющая учителю осуществлять совместную работу с учащимися и ее оценку в реальном времени;

Таксономия мыслительных операций по Б. Блуму (1956 г.)*. Строки таблицы соответствуют уровням «Пирамиды Блума» [10]

Название мыслительной операции	Содержание мыслительной операции	Глагольное выражение мыслительной операции
1. ЗНАНИЕ	Учащийся знает употребляемые термины, конкретные факты, методы и процедуры, основные понятия, правила и принципы	Упорядочи, определи, продублируй, составь список, соотнеси, запомни, назови, проранжируй, опознай, отнеси, вспомни, повтори, воспроизведи
2. ПОНИМАНИЕ	Учащийся понимает правила, факты и принципы, интерпретирует словесный материал, схемы, графики, диаграммы, преобразует словесный материал в математические выражения и наоборот, предположительно оценивает будущие события, последствия, вытекающие из имеющихся данных	Классифицируй, опиши, обсуди, объясни, вырази, осознай, укажи, расположи, распознай, сообщи, подтверди, сделай обзор, отбери, отсортируй, расскажи, переведи, проэкстраполируй
3. ПРИМЕНЕНИЕ	Учащийся использует понятия и принципы в новых ситуациях, применяет законы и теории в конкретных практических ситуациях, демонстрирует правильное применение метода или процедуры	Примени, выбери, продемонстрируй, инсценируй, привлекли, проиллюстрируй, проинтерпретируй, произведи операции, приготовь, выполни, осуществи, отработай, составь план/программу, набросай, реши, используй
4. АНАЛИЗ	Учащийся выделяет скрытые (неявные) предположения, видит ошибки и упущения в логике рассуждений, проводит различия между фактами и следствиями, оценивает значимость данных	Проанализируй, оцени, рассчитай, категоризируй, сравни, сопоставь, выскажи критику, составь диаграмму, различи, распознай, найди отличия, исследуй, проэкспериментируй, подведи итог, проясни, опробуй
5. СИНТЕЗ	Учащийся пишет небольшое творческое сочинение, предлагает план проведения эксперимента, использует знания из разных областей, чтобы составить план решения той или иной проблемы	Организуй, собери, скомпонуй, сочини, построй, создай, спроектируй, разработай, овладей, организуй, спланируй, подготовь, предложи, установи, синтезируй, напиши
6. ОЦЕНКА	Учащийся оценивает логику представления материала в виде письменного текста, оценивает соответствие вывода имеющимся данным, оценивает значимость того или иного продукта деятельности исходя из внутренних или внешних критериев	Оцени, поспорь, осуществи экспертизу, выбери, сравни, защити, выскажи суждение, взвесь «за» и «против», сделай вывод, спрогнозируй, проранжируй, выставь оценку, выбери, поддержи, оцени значимость/значение

* В 1990-х гг. группа американских психологов-когнитивистов, возглавляемая бывшим учеником Б. Блума Лорином Андерсоном, предложила обновленную версию таксономии Блума применительно к реалиям XXI в. В пирамиде Блума существительные, называющие мыслительные операции, были заменены на герундий — часть речи, описывающую процесс выполнения мыслительных операций (*knowledge — remembering; comprehension — understanding; application — applying; analysis — analyzing*), а *synthesis — creating* и *evaluation — evaluating* поменялись местами [10].

- **Prezi** — приложение для создания презентаций;
- **ScreenChomp** — интерактивная цифровая доска для создания набросков и заметок;
- **Toontastic** — приложение для создания мультфильмов и анимированных изображений;
- **Voicethread** — интерактивное приложение для одновременного выполнения нескольких операций с документами: правка, обсуждение, создание схем и заметок;
- **Wordpress** — система управления содержимым сайта, блога.

5 Системы оценки эффективности учебно-познавательной деятельности учащихся как уровня достижения планируемых образовательных результатов

В контексте требований новых ФГОС следует уделить внимание двум основным видам оценки

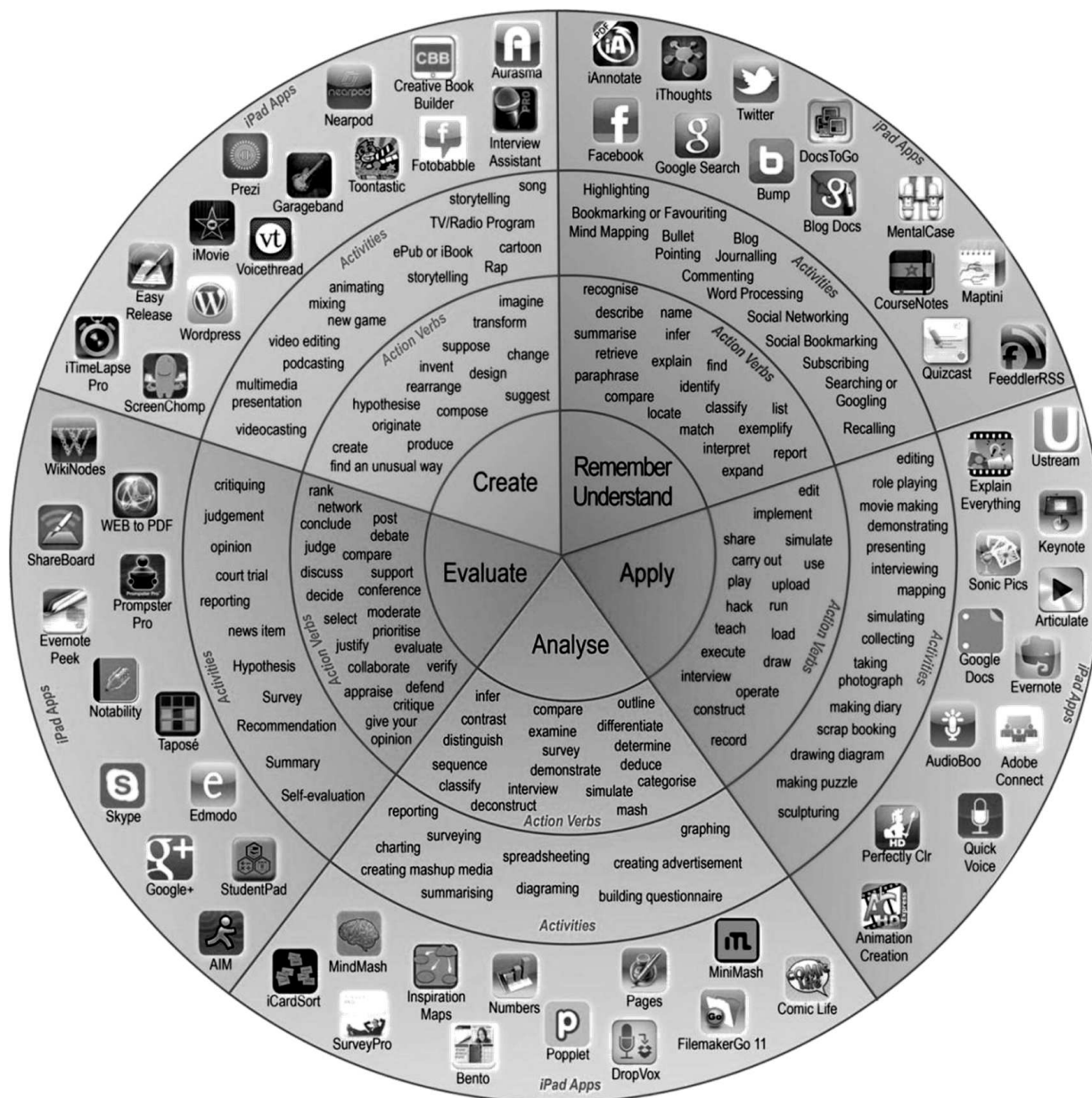


Рис. 2 Оригинальная версия «Педагогического колеса» — современной интерпретации «Пирамиды Блума» с точки зрения интеграции педагогических и мобильных ИКТ

достижения планируемых образовательных результатов: **суммирующему (итоговому) оцениванию и формирующему (процессуальному) оцениванию.**

При *суммирующем оценивании* определяется качество усвоения некоторого объема учебного материала в течение определенного временного (четверти, полугодия, учебного года) или информационного (урока, раздела, модуля, группы модулей, цикла уроков, курса и т.д.) этапа обучения, который рассматривается как некоторый итог.

Суммирующее оценивание имеет целью выставление частично (промежуточной) или полностью итоговой отметки и производится, как правило, в числовом выражении.

Суммирующее оценивание является своего рода объектом презентации передового опыта, отчета учителя и/или учащихся о проделанной работе в течение указанного этапа образовательной деятельности (периода учебного времени или объема учебной информации, подлежащей усвоению). По-

этому суммирующее оценивание можно рассматривать лишь как простую констатацию факта того, были или нет достигнуты запланированные образовательные результаты, и если «да», то какой уровень был достигнут. Это определяет место суммирующего оценивания как бы вовне процесса освоения знаний или приобретения компетенций [11].

Значимость *формирующего оценивания* состоит в том, что оно, будучи комплексной, интегральной процедурой, включенной в сам процесс УПД, позволяет получить «пошаговые» данные об уровне развития мыслительных способностей или компетенций учащихся на сколь угодно мелких промежуточных этапах. Эта процессуальная функция формирующего оценивания позволяет рассматривать его как необходимую составную часть эффективного обучения, требующего, однако, дополнительных затрат учебного времени.

Формирующее оценивание базируется на трех основополагающих принципах педагогической технологии, получившей название LOA (learning-oriented assessment) — «оценивание, направленное на обучение» [13, с. 57, 59–60]:

- (1) постановке задач оценивания как учебных задач;
- (2) привлечению учащихся к оцениванию работы своих товарищей и самооценке;
- (3) осуществлению обратной связи, направленной не на пройденный материал, а на материал, который предстоит освоить.

В названии LOA-технологии ключевым является слово «обучение», а не «оценивание», поскольку вне процесса обучения преимущества формирующего оценивания теряют смысл. Более того, при сбалансированном подходе и обучение, и оценивание не просто движутся в едином русле достижения эффективных результатов обучения, но эти результаты являются *запланированными* результатами, поскольку ориентируются на конкретную учебную цель, поставленную при формулировке учебной задачи, а именно: подобная цель является объектом УПД. В этом смысле формирующее оценивание самым оптимальным образом ориентирует процесс обучения на достижение запланированных результатов обучения как его конечной или промежуточной цели.

Первый принцип LOA-технологии — постановка задач оценивания как учебных задач — предполагает совмещение задач оценивания и задач обучения. Этот принцип состоит в пошаговом оценивании перспективы успешности/неуспешности решения задачи на промежуточных, более мелких этапах, т. е. прогнозирование этой успешности/неуспешности

с целью корректировки алгоритма ее решения. При этом учащимся с целью прогнозирования ситуации поневоле приходится осуществлять действия, связанные с экстраполяцией учебного материала, т. е. с углублением не в пройденный учебный материал, а с обращением к материалу, подлежащему дальнейшему усвоению. Эти шаги способствуют развитию многообразных форм мыслительной деятельности учащихся.

Второй принцип LOA-технологии — привлечение учащихся к оцениванию работы своих товарищей и самооценке — способствует не только развитию навыков рефлексии и самооценки, являющихся важнейшими регулятивными характеристиками как предметных, так и метапредметных компетенций. Он приучает их к осуществлению экспертной оценки в соответствии с критериями, разработанными ими самими, к принятию ответственных решений, влияющих на конечный результат, к учебному сотрудничеству. Кроме того, прозрачность оценки, выносимой в результате совместного обсуждения, служит залогом понимания учащимися конечной цели своей УПД.

Третий принцип LOA-технологии — осуществление обратной связи, направленной не на пройденный материал, а на материал, который предстоит освоить. Обратная связь сама по себе не побуждает учащихся к дальнейшему изучению предмета, однако, выступая как основной способ анализа результатов на отдельных этапах решения учебной задачи с целью корректировки путей ее решения, она тем самым нацеливает учащихся на дальнейшее изучение материала.

Процедура интеграции оценивания и обучения в рамках LOA-технологии сопоставима с одним из принципов формирования операционного стиля мышления, выдвинутых академиком А. П. Ершовым еще в 1980-х гг.: планирование структуры целенаправленных действий в определенных условиях с помощью заданного набора средств. Реализация этого принципа предполагает, что учащийся должен не только представлять себе ситуацию, в которой будет осуществляться решение поставленной задачи, но и уметь анализировать ее, выявляя имеющиеся средства, доступные резервы и предполагаемые трудности. Анализ этой ситуации необходим для выстраивания верной стратегии решения — иными словами, создания адекватной задаче структуры целенаправленных умственных действий (алгоритма), осуществление которых согласно принятому плану поможет привести к успешному результату, что само по себе и предполагает формирующее оценивание.

В процессе анализа исходной ситуации учащиеся подбирают ряд более простых целевых ситуаций,

выстраивают их в определенную иерархию, не противоречащую исходной, хотя и упрощающую ее на некоторых этапах, и тем самым шаг за шагом движутся в направлении нужного решения.

Кроме того, такие пошаговые процедуры сужают поле поиска решения и тем самым упрощают его, делая посильным. Таким образом, деление сложных задач на более простые, элементарные, «пооперациональные» задачи, во-первых, сужает поле поиска вероятного решения, а во-вторых, структурирует траекторию поиска, разделяя ее на шаги или этапы, так что движение осуществляется дозированно, а на каждом этапе решается некоторая элементарная задача, приближающая учащегося к решению исходной более сложной задачи.

Очень важно подчеркнуть, что деление первоначальной задачи на более мелкие и простые дозированные этапы, осуществляемое в результате анализа исходной ситуации, происходит именно путем формирующего оценивания, которое направлено на корректировку стратегической траектории решения задачи (более подробно см. [14, с. 28–30]). При этом результаты формирующего оценивания, выраженные в баллах, демонстрируют не только правильность решения задачи, но и рациональность выбранной стратегической траектории движения к искомому решению.

Необходимо отметить, что формирующее оценивание требует одновременной, «сиюминутной» вовлеченности в процесс обучения и учителя, и учащихся, что имеет место при выполнении заданий, направленных на максимальную кооперацию всех участников процесса УПД. Одним из форматов урока, способствующего реализации концепции формирующего оценивания, направленного на достижение планируемых образовательных результатов при решении конкретных образовательных задач, служит, например, технология «перевернутого» урока. Она предполагает самостоятельную работу учащихся с электронным контентом дома, а основное интерактивное общение ориентирует на выполнение совместных заданий в ИОС учебного заведения и решение задач повышенной трудности в классе при очном общении учителя и учащихся.

Оценка достижения предметных результатов обучения осуществляется, как правило, традиционными балльными методами. Однако в ряде случаев, особенно при текущем или промежуточном оценивании, которое может вестись в формате формирующего оценивания, полученные результаты целесообразно сохранять с помощью так называемой накопительной системы оценивания (например, в форме портфолио) и затем учитывать при определении итоговой оценки вкпе с результатами суммирующего оценивания.

6 Заключение

Уровень эффективности УПД с использованием ИКТ определяется способностью учителя организовать совместную работу с учащимися, ориентированную на развитие форм мыслительной деятельности, приводящих к созданию интегрированного персонального познавательного стиля каждого. Такую возможность предоставляют учителю педагогические и новые информационные технологии, объединенные в целостный дидактический процесс, реализуемый в ИОС учебного заведения.

Литература

1. Капранов В. К., Капранова М. Н. ЭОР от Интернета до учителя // Информационные технологии в образовании XXI века: Сб. науч. тр. II Всеросс. науч.-практич. конф. — М.: НИЯУ МИФИ, 2012. Т. 2. С. 297–300.
2. Федеральный государственный образовательный стандарт основного общего образования / Минобрнауки РФ. — М.: Просвещение, 2011. 48 с.
3. Ривкин Е. Ю. Профессиональная деятельность учителя в период перехода на ФГОС основного общего образования: Теория и технологии. — Волгоград: Учитель, 2014. 183 с.
4. Информационные и коммуникационные технологии в образовании / Под. ред. Б. Дендева. — М.: ИИТО ЮНЕСКО, 2013. 320 с.
5. Вихрев В. В., Христочевская А. С., Христочевский С. А. О новой концепции информатизации образования // Системы и средства информатики, 2014. Т. 24. № 4. С. 162–172.
6. Петти Д. Современное обучение: Практическое руководство / Пер. с англ. П. Кириллова. — М.: Ломоносовъ, 2010. 624 с. (Прикладная психология). (Pet-ti D. Teaching today: A practical guide. — 4th ed. — Cheltenham: Nelson Thornes, 2009. 624 p.)
7. Концепция информатизации образования // Информатика и образование, 1990. № 1. С. 3–9.
8. Холодная М. А. Когнитивные стили. О природе индивидуального ума. — М.: ПЕР СЭ, 2002. 304 с.
9. Bloom B. Developing talent in young people. — New York: Ballantine Books, 1985. 558 p.
10. Bloom's Taxonomy and the Pedagogy Wheel. <https://www.gadsdenstate.edu/academics/elearning/pdf/First%20Friday%20Tech%20Tip%20Aug%202013.pdf>.
11. Чайка В. М. Таксономия целей обучения. http://uchebnikionline.com/pedagogika/osnovi_didaktiki-_chayka_vm/taksonomiya_tsiley_navchannya.htm.
12. The Pedagogy Wheel. <http://www.unity.net.au/padwheel/padwheelposter.pdf>, <http://elearningstuff.net/wp-content/uploads/2013/06/padagogy-wheel.jpg>.

13. Carless D. Learning-oriented assessment: conceptual bases and practical implications // *Innov. Educ. Teach. Int.*, 2007. Vol. 44. No. 1. P. 57–66.
14. Корчажкина О. М. Операционный стиль мышления: взгляд четверть века спустя // *Информатика и образование*, 2010. № 5. С. 28–36.

Поступила в редакцию 04.01.15

ON ACCESS TO THE EFFICIENCY OF STUDENTS' COGNITIVE ACTIVITIES WHILE USING THE NEW INFORMATION TECHNOLOGIES

O. M. Korchazhkina

Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation

Abstract: The paper considers a problem of how to measure the efficiency of students' cognitive activities as the planned outcomes in compliance with the achieved ones, both expressed in terms of specific products of learning and cognitive activities that are obtained while performing mental tasks. Combining the style of teaching and learning methods with the use of pedagogical and new information technologies integrated while performing various types of tasks is discussed. An example of how to verbalize the results achieved during the learning activities with the use of mobile devices is given. The way of verbalizing is based on Bloom's taxonomy action verbs. It is found out that the level of how well students perform cognitive tasks with the use of information and communication technologies depends on their teacher's ability to collaborate with them while developing all forms of their mental activity, which leads to building an integrated personal cognitive style for each student.

Keywords: efficiency of training; planned educational results; achieved educational outcomes; mobile devices; cognitive/mental tasks; individual style of learning; teaching methods; LOA-technology

DOI: 10.14357/19922264150110

References

1. Kapranov, V. K., and M. N. Kapranova. 2012. EOR ot Interneta do uchitelya [Digital educational resources: From the Internet to the teacher]. *Informatsionnye Tekhnologii v Obrazovanii XXI veka: Sb. nauch. tr. II Vseross. nauch.-praktich. konf.* [2nd All-Russian Scientific-Practical Conference "Information Technologies in Education of the XXI Century"]. Moscow. 2:297–300.
2. Minobrnauki RF [Department of Education and Science of the Russian Federation]. 2011. Federal'nyy gosudarstvennyy obrazovatel'nyy standart osnovnogo obshchego obrazovaniya [Federal State Educational Standard of basic general education]. Moscow: Prosveshchenie Publ. House. 48 p.
3. Rivkin, E. Yu. 2014. *Professional'naya deyatel'nost' uchitelya v period perekhoda na FGOS osnovnogo obshchego obrazovaniya: Teoriya i tekhnologii* [The teacher's professional activity in transition to the Federal State Educational Standard of basic general education]. Volgograd: Uchitel' Publ. House. 183 p.
4. Dende, B., ed. 2013. *Informatsionnye i kommunikatsionnye tekhnologii v obrazovanii* [Information and communication technologies in education]. Moscow: IITE UNESCO. 320 p.
5. Vikhrev, V. V., A. S. Christochevskaya, and S. A. Christochevsky. 2014. O novoy kontseptsii informatizatsii obrazovaniya [On a new conception of informatization of education]. *Sistemy i Sredstva Informatiki — Systems and Means of Informatics* 24(4):162–172.
6. Petti, D. 2009. *Teaching today: A practical guide*. 4th ed. Cheltenham: Nelson Thornes. 624 p.
7. Kontseptsiya informatizatsii obrazovaniya [The conception of informatization of education]. 1990. *Informatika i Obrazovanie* [Informatics and Education] 1:3–9.
8. Kholodnaya, M. A. 2002. *Kognitivnye stili. O prirode individual'nogo uma* [Cognitive styles. On the nature of the individual mind]. Moscow: PERSY Publ. House. 304 p.
9. Bloom, B. 1985. *Developing talent in young people*. New York: Ballantine Books. 558 p.
10. Bloom's taxonomy and the Pedagogy Wheel. Available at: <https://www.gadsdenstate.edu/academics/elearning/pdf/First%20Friday%20Tech%20Tip%20Aug%202013.pdf> (accessed December 1, 2014).
11. Chayka, V. M. Taksonomiya tseley obucheniya [Taxonomy of educational objectives]. Available at: http://uchebnikionline.com/pedagogika/osnovi_didaktiki-_chayka_vm/taksonomiya_tsiley_navchannya.htm (accessed November 3, 2014).
12. Pedagogy Wheel, The. Available at: <http://www.unity.net.au/padwheel/padwheelposter.pdf>; <http://elearnings>

- tuff.net/wp-content/uploads/2013/06/padagogy-wheel.jpg (accessed December 11, 2014).
13. Carless, D. 2007. Learning-oriented assessment: conceptual bases and practical implications. *Innov. Educ. Teach. Int.* 44(1):57–66.
 14. Korchazhkina, O. M. 2010. Operatsionnyy stil' myshleniya: Vzglyad chetvert' veka spustya [The operational style of thinking: A sight in a quarter of the century]. *Informatika i Obrazovanie* [Informatics and Education] 5:28–36.

Received January 4, 2015

Contributor

Korchazhkina Olga M. (b. 1953) — Candidate of Science (PhD) in technology, senior scientist, Institute of Informatics Problems, Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation; olgakomax@gmail.com

Гайдамака Юлия Васильевна (р. 1971) — кандидат физико-математических наук, доцент Российского университета дружбы народов

Долев Шломи (р. 1958) — доктор наук по информатике, профессор Университета им. Бен-Гуриона в Негаве, Беэр-Шева, Израиль

Жаворонкова Юлия Вадимовна (р. 1990) — программист-разработчик ООО «Спутник»

Калиниченко Леонид Андреевич (р. 1937) — доктор физико-математических наук, профессор, заведующий лабораторией Института проблем информатики Российской академии наук; профессор факультета вычислительной математики и кибернетики Московского государственного университета им. М. В. Ломоносова

Кантор Ольга Геннадиевна (р. 1971) — кандидат физико-математических наук, старший научный сотрудник Института социально-экономических исследований Уфимского научного центра Российской академии наук

Ковалев Дмитрий Юрьевич (р. 1988) — младший научный сотрудник Института проблем информатики Российской академии наук

Ковалева Дана Александровна (р. 1973) — кандидат физико-математических наук, научный сотрудник Института астрономии Российской академии наук

Ковалёв Сергей Протасович (р. 1972) — доктор физико-математических наук, старший научный сотрудник Института проблем управления им. В. А. Трапезникова

Коган-Садецкая Марина (р. 1977) — аспирант Университета им. Бен-Гуриона в Негаве, Беэр-Шева, Израиль

Колесов Сергей Викторович (р. 1951) — доктор химических наук, заведующий лабораторией Института органической химии Уфимского научного центра Российской академии наук

Корепанов Эдуард Рудольфович (р. 1966) — кандидат технических наук, заведующий сектором Института проблем информатики Российской академии наук

Корчажкина Ольга Максимовна (р. 1953) — кандидат технических наук, старший научный сотрудник Института проблем информатики Российской академии наук

Кудрявцев Алексей Андреевич (р. 1978) — кандидат физико-математических наук, доцент кафедры математической статистики факультета вычислительной математики и кибернетики Московского государственного университета им. М. В. Ломоносова

Кузнецов Сергей Иванович (р. 1955) — научный сотрудник Института органической химии Уфимского научного центра Российской академии наук

Малков Олег Юрьевич (р. 1961) — доктор физико-математических наук, доцент, заведующий отделом Института астрономии Российской академии наук; профессор физического факультета Московского государственного университета им. М. В. Ломоносова

Попова Мария Сергеевна (р. 1994) — студент Московского физико-технического института

Самуйлов Андрей Константинович (р. 1988) — аспирант Российского университета дружбы народов; исследователь Технологического университета г. Тампере, Финляндия

Синицын Владимир Игоревич (р. 1968) — доктор физико-математических наук, доцент, заведующий отделом Института проблем информатики Российской академии наук

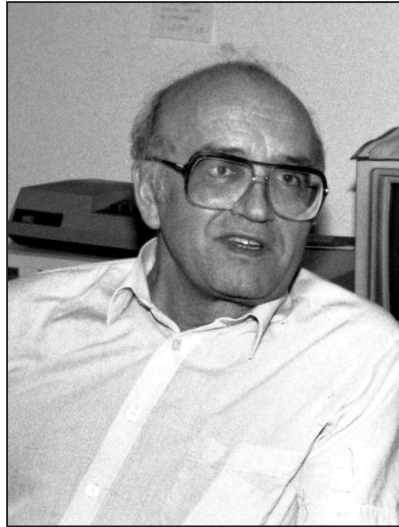
Синицын Игорь Николаевич (р. 1940) — доктор технических наук, профессор, заслуженный деятель науки РФ, заведующий отделом Института проблем информатики Российской академии наук

Спивак Семен Израилевич (р. 1945) — доктор физико-математических наук, профессор Башкирского государственного университета

Стрижов Вадим Викторович (р. 1967) — кандидат физико-математических наук, доцент Московского физико-технического института, ведущий научный сотрудник Вычислительного центра РАН им. Дородницына

Шоргин Сергей Яковлевич (р. 1952) — доктор физико-математических наук, профессор, заместитель директора Института проблем информатики Российской академии наук

Юнусова Дарья Сергеевна (р. 1989) — аспирант Башкирского государственного университета



Профессор Игорь Алексеевич Ушаков

22.01.1935—27.02.2015

Редакционный совет и редакционная коллегия журнала «Информатика и её применения» с глубоким прискорбием извещают, что 27 февраля 2015 г. после тяжелой и продолжительной болезни скончался Игорь Алексеевич Ушаков — доктор технических наук, профессор, член редколлегии журнала «Информатика и её применения».

Игорь Алексеевич Ушаков окончил Московский авиационный институт, в 1963 г. защитил кандидатскую, а в 1968 г. — докторскую диссертацию. С 1958 по 1989 гг. работал в ряде научно-исследовательских организаций СССР, в том числе руководил отделами в НИИ АА и ВЦ АН СССР; с 1969 по 1989 гг. преподавал в МФТИ (был профессором, а затем заведующим кафедрой) и в МЭИ. С 1989 г. — в США: являлся профессором университета Дж. Вашингтона, университета Дж. Мэйсона и Калифорнийского университета, сотрудником компаний MCI, Qualcomm и Hughes.

И. А. Ушаков с момента основания журнала «Надежность и контроль качества» был заместителем ответственного редактора, а затем на протяжении многих лет членом редколлегии. В 2006 г. основал электронный международный журнал “Reliability: Theory & Application”, главным редактором которого оставался до конца жизни.

Учебниками и справочниками по теории надежности, написанными И. А. Ушаковым, пользовались и пользуются несколько поколений ученых и специалистов в разных странах мира.

Игорь Алексеевич всегда уделял огромное внимание работе с молодежью; более 50 его учеников защитили докторские и кандидатские диссертации.

И. А. Ушаков вел активную научно-просветительскую деятельность. В частности, он был одним из организаторов и руководителей Московского кабинета качества и надежности при Политехническом музее (целью этого Кабинета было оказание консультаций работникам промышленных предприятий и чтение курсов лекций для инженеров, занимающихся проблемой надежности). Находясь в США, И. А. Ушаков создал международный интернет-форум им. Б. В. Гнеденко, объединивший около 400 видных специалистов по приложениям теории вероятностей и математической статистики, преимущественно в области теории надежности и анализа риска, из десятков стран мира; коллективным членом этого Форума является и наш журнал. Цели Форума — содействие контактам между специалистами из разных стран, организация обмена профессиональными новостями и информацией (новые публикации, предстоящие события и др.). Также необходимо отметить большое число научно-популярных работ, опубликованных И. А. Ушаковым.

И. А. Ушаков обладал большим личным обаянием, имел широкий круг интересов. Все знавшие И. А. Ушакова всегда будут помнить его как замечательного ученого и прекрасного человека.

Редакционный совет и редакционная коллегия журнала «Информатика и её применения» выражают глубокие соболезнования родным и близким покойного, всем, кто его знал и работал с ним.

Правила подготовки рукописей для публикации в журнале «Информатика и её применения»

Журнал «Информатика и её применения» публикует теоретические, обзорные и дискуссионные статьи, посвященные научным исследованиям и разработкам в области информатики и ее приложений.

Журнал издается на русском языке. По специальному решению редколлегии отдельные статьи могут печататься на английском языке.

Тематика журнала охватывает следующие направления:

- теоретические основы информатики;
- математические методы исследования сложных систем и процессов;
- информационные системы и сети;
- информационные технологии;
- архитектура и программное обеспечение вычислительных комплексов и сетей.

1. В журнале печатаются статьи, содержащие результаты, ранее не опубликованные и не предназначенные к одновременной публикации в других изданиях.

Публикация не должна нарушать закон об авторских правах.

Направляя рукопись в редакцию, авторы сохраняют все права собственников данной рукописи и при этом передают учредителям и редколлегии неисключительные права на издание статьи на русском языке (или на языке статьи, если он отличен от русского) и на ее распространение в России и за рубежом. Авторы должны представить в редакцию письмо в следующей форме:

Соглашение о передаче права на публикацию:

«Мы, нижеподписавшиеся, авторы рукописи «. . .», передаем учредителям и редколлегии журнала «Информатика и её применения» неисключительное право опубликовать данную рукопись статьи на русском языке как в печатной, так и в электронной версиях журнала. Мы подтверждаем, что данная публикация не нарушает авторского права других лиц или организаций, а также не содержит сведений, запрещенных к опубликованию в открытой печати.

Подписи авторов: (ф. и. о., дата, адрес)».

Это соглашение может быть представлено в бумажном виде или в виде отсканированной копии (с подписями авторов).

Редколлегия вправе запросить у авторов экспертное заключение о возможности публикации представленной статьи в открытой печати.

2. К статье прилагаются данные автора (авторов) (см. п. 8). При наличии нескольких авторов указывается фамилия автора, ответственного за переписку с редакцией.
3. Редакция журнала осуществляет экспертизу присланных статей в соответствии с принятой в журнале процедурой рецензирования.

Возвращение рукописи на доработку не означает ее принятия к печати.

Доработанный вариант с ответом на замечания рецензента необходимо прислать в редакцию.

4. Решение редколлегии о публикации статьи или ее отклонении сообщается авторам.
Редколлегия может также направить авторам текст рецензии на их статью. Дискуссия по поводу отклоненных статей не ведется.
5. Редактура статей высылается авторам для просмотра. Замечания к редакции должны быть присланы авторами в кратчайшие сроки.
6. Рукопись предоставляется в электронном виде в форматах MS WORD (.doc или .docx) или ЛАТЭК (.tex), дополнительно — в формате .pdf, на дискете, лазерном диске или электронной почтой. Предоставление бумажной рукописи необязательно.
7. При подготовке рукописи в MS Word рекомендуется использовать следующие настройки.

Параметры страницы: формат — А4; ориентация — книжная; поля (см): внутри — 2,5, снаружи — 1,5, сверху — 2, снизу — 2, от края до нижнего колонтитула — 1,3.

Основной текст: стиль — «Обычный», шрифт — Times New Roman, размер — 14 пунктов, абзацный отступ — 0,5 см, 1,5 интервала, выравнивание — по ширине.

Рекомендуемый объем рукописи — не свыше 20 страниц указанного формата.

Сокращения слов, помимо стандартных, не допускаются. Допускается минимальное количество аббревиатур.

Все страницы рукописи нумеруются.

Шаблоны примеров оформления представлены в Интернете: <http://www.ipiran.ru/journal/template.doc>

8. Статья должна содержать следующую информацию на *русском и английском языках*:

- название статьи;
- Ф.И.О. авторов, на английском можно только имя и фамилию;
- место работы, с указанием почтового адреса организации и электронного адреса каждого автора;
- сведения об авторах, в соответствии с форматом, образцы которого представлены на страницах:
http://www.ipiran.ru/journal/issues/2013_07_01_rus/authors.asp и
http://www.ipiran.ru/journal/issues/2013_07_01_eng/authors.asp;
- аннотация (не менее 100 слов на каждом из языков). Аннотация — это краткое резюме работы, которое может публиковаться отдельно. Она является основным источником информации в информационных системах и базах данных. Английская аннотация должна быть оригинальной, может не быть дословным переводом русского текста и должна быть написана хорошим английским языком. В аннотации не должно быть ссылок на литературу и, по возможности, формул;
- ключевые слова — желательно из принятых в мировой научно-технической литературе тематических тезаурусов. Предложения не могут быть ключевыми словами;
- источники финансирования работы (ссылки на гранты, проекты, поддерживающие организации и т. п.).

9. Требования к спискам литературы.

Ссылки на литературу в тексте статьи нумеруются (в квадратных скобках) и располагаются в каждом из списков литературы в порядке первых упоминаний.

Списки литературы представляются в двух вариантах:

- (1) **Список литературы к русскоязычной части.** Русские и английские работы — на языке и в алфавите оригинала;
- (2) **References.** Русские работы и работы на других языках — в латинской транслитерации с переводом на английский язык; английские работы и работы на других языках — на языке оригинала.

Необходимо для составления списка “References” пользоваться размещенной на сайте <http://translit.ru/> бесплатной программой транслитерации русского текста в латиницу, при этом в закладке «варианты. . . » следует выбрать опцию BGN.

Список литературы “References” приводится полностью отдельным блоком, повторяя все позиции из списка литературы к русскоязычной части, независимо от того, имеются или нет в нем иностранные источники. Если в списке литературы к русскоязычной части есть ссылки на иностранные публикации, набранные латиницей, они полностью повторяются в списке “References”.

Ниже приведены примеры ссылок на различные виды публикаций в списке “References”.

Описание статьи из журнала:

Zagurenko, A. G., V. A. Korotovskikh, A. A. Kolesnikov, A. V. Timonov, and D. V. Kardymon. 2008. Tekhniko-ekonomicheskaya optimizatsiya dizayna gidrorazryva plasta [Technical and economic optimization of the design of hydraulic fracturing]. *Neftyanoe hozyaystvo [Oil Industry]* 11:54–57.

Zhang, Z., and D. Zhu. 2008. Experimental research on the localized electrochemical micromachining. *Rus. J. Electrochem.* 44(8):926–930. doi:10.1134/S1023193508080077.

Описание статьи из электронного журнала:

Swaminathan, V., E. Lepkoswka-White, and B. P. Rao. 1999. Browsers or buyers in cyberspace? An investigation of electronic factors influencing electronic exchange. *JCMC* 5(2). Available at: <http://www.ascusc.org/jcmc/vol5/issue2/> (accessed April 28, 2011).

Описание статьи из продолжающегося издания (сборника трудов):

Astakhov, M. V., and T. V. Tagantsev. 2006. Eksperimental'noe issledovanie prochnosti soedineniy “stal”–kompozit [Experimental study of the strength of joints “steel–composite”]. *Trudy MGTU “Matematicheskoe modelirovanie slozhnykh tekhnicheskikh sistem” [Bauman MSTU “Mathematical Modeling of Complex Technical Systems” Proceedings]*. 593:125–130.

Описание материалов конференций:

Usmanov, T. S., A. A. Gusmanov, I. Z. Mullagalin, R. Ju. Muhametshina, A. N. Chervyakova, and A. V. Sveshnikov. 2007. Osobennosti proektirovaniya razrabotki mestorozhdeniy s primeneniem gidrorazryva plasta [Features of the design of field development with the use of hydraulic fracturing]. *Trudy 6-go Mezhdunarodnogo Simpoziuma "Novye resursoberegayushchie tekhnologii nedropol'zovaniya i povysheniya neftegazootdachi"* [6th Symposium (International) "New Energy Saving Subsoil Technologies and the Increasing of the Oil and Gas Impact" Proceedings]. Moscow. 267–272.

Описание книги (монографии, сборники):

Lindorf, L. S., and L. G. Mamikonians, eds. 1972. *Ekspluatatsiya turbogeneratorov s neposredstvennym okhlazhdeniem* [Operation of turbine generators with direct cooling]. Moscow: Energy Publs. 352 p.

Latyshev, V. N. 2009. *Tribologiya rezaniya. Kn. 1: Friksionnye protsessy pri rezanii metallov* [Tribology of cutting. Vol. 1: Frictional processes in metal cutting]. Ivanovo: Ivanovskii State Univ. 108 p.

Описание переводной книги (в списке литературы к русскоязычной части необходимо указать: / Пер. с англ. — после названия книги, а в конце ссылки указать оригинал книги в круглых скобках):

1. В русскоязычной части:

Тимошенко С. П., Янг Д. Х., Уивер У. Колебания в инженерном деле / Пер. с англ. — М.: Машиностроение, 1985. 472 с. (*Timoshenko S. P., Young D. H., Weaver W. Vibration problems in engineering. — 4th ed. — N.Y.: Wiley, 1974. 521 p.*)

2. В англоязычной части:

Timoshenko, S. P., D. H. Young, and W. Weaver. 1974. *Vibration problems in engineering*. 4th ed. N.Y.: Wiley. 521 p.

Описание неопубликованного документа:

Латыпов, А. Р., М. М. Хасанов, и В. А. Байков. 2004. Geology and production (NGT GiD). Certificate on official registration of the computer program No. 2004611198. (In Russian, unpubl.)

Описание интернет-ресурса:

Pravila tsitirovaniya istochnikov [Rules for the citing of sources]. Available at: <http://www.scribd.com/doc/1034528/> (accessed February 7, 2011).

Описание диссертации или автореферата диссертации:

Semenov, V. I. 2003. *Matematicheskoe modelirovanie plazmy v sisteme kompaktnyy tor* [Mathematical modeling of the plasma in the compact torus]. D.Sc. Diss. Moscow. 272 p.

Kozhunova, O. S. 2009. *Tekhnologiya razrabotki semanticheskogo slovarya informatsionnogo monitoringa* [Technology of development of semantic dictionary of information monitoring system]. PhD Thesis. Moscow: IPI RAN. 23 p.

Описание ГОСТа:

GOST 8.586.5-2005. 2007. *Metodika vypolneniya izmereniy. Izmerenie raskhoda i kolichestva zhidkostey i gazov s pomoshch'yu standartnykh suzhayushchikh ustroystv* [Method of measurement. Measurement of flow rate and volume of liquids and gases by means of orifice devices]. Moscow: Standardinform Publs. 10 p.

Описание патента:

Bolshakov, M. V., A. V. Kulakov, A. N. Lavrenov, and M. V. Palkin. 2006. *Sposob orientirovaniya po krenu letatel'nogo apparata s opticheskoy golovkoy samonavedeniya* [The way to orient on the roll of aircraft with optical homing head]. Patent RF No. 2280590.

10. Присланные в редакцию материалы авторам не возвращаются.
11. При отправке файлов по электронной почте просим придерживаться следующих правил:
 - указывать в поле subject (тема) название журнала и фамилию автора;
 - использовать attach (присоединение);
 - в состав электронной версии статьи должны входить: файл, содержащий текст статьи, и файл(ы), содержащий(е) иллюстрации.
12. Журнал «Информатика и её применения» является некоммерческим изданием. Плата за публикацию не взимается, гонорар авторам не выплачивается.

Адрес редакции журнала «Информатика и её применения»:

Москва 119333, ул. Вавилова, д. 44, корп. 2, ИПИ РАН

Тел.: +7 (499) 135-86-92 Факс: +7 (495) 930-45-05

e-mail: rust@ipiran.ru (Сейфуль-Мулюков Рустем Бадриевич)

<http://www.ipiran.ru/journal/issues/>

Requirements for manuscripts submitted to Journal “Informatics and Applications”

Journal “Informatics and Applications” (Inform. Appl.) publishes theoretical, review, and discussion articles on the research and development in the field of informatics and its applications.

The journal is published in Russian. By a special decision of the editorial board, some articles can be published in English.

The topics covered include the following areas:

- theoretical fundamentals of informatics;
 - mathematical methods for studying complex systems and processes;
 - information systems and networks;
 - information technologies; and
 - architecture and software of computational complexes and networks.
1. The Journal publishes original articles which have not been published before and are not intended for publication in other editions. An article submitted to the Journal must not violate the Copyright law. Sending the manuscript to the Editorial Board, the authors retain all rights of the owners of the manuscript and transfer the nonexclusive rights to publish the article in Russian (or the language of the article, if not Russian) and its distribution in Russia and abroad to the Founders and the Editorial Board. Authors should submit a letter to the Editorial Board in the following form:

Agreement on the transfer of rights to publish:

“We, the undersigned authors of the manuscript “. . .”, pass to the Founder and the Editorial Board of the Journal “Informatics and Applications” the nonexclusive right to publish the manuscript of the article in Russian (or in English) in both print and electronic versions of the Journal. We affirm that this publication does not violate the Copyright of other persons or organizations.

Author(s) signature(s): (name(s), address(es), date).

This agreement should be submitted in paper form or in the form of a scanned copy (signed by the authors).

2. A submitted article should be attached with **the data on the author(s)** (see item 8). If there are several authors, the contact person should be indicated who is responsible for correspondence with the Editorial Board and other authors about revisions and final approval of the proofs.
3. The Editorial Board of the Journal examines the article according to the established reviewing procedure. If the authors receive their article for correction after reviewing, it does not mean that the article is approved for publication. The corrected article should be sent to the Editorial Board for the subsequent review and approval.
4. The decision on the article publication or its rejection is communicated to the authors. The Editorial Board may also send the reviews on the submitted articles to the authors. Any discussion upon the rejected articles is not possible.
5. The edited articles will be sent to the authors for proofread. The comments of the authors to the edited text of the article should be sent to the Editorial Board as soon as possible.
6. The manuscript of the article should be presented electronically in the MS WORD (.doc or .docx) or \LaTeX (.tex) formats, and additionally in the .pdf format. All documents may be sent by e-mail or provided on a CD or diskette. A hard copy submission is not necessary.

7. The recommended typesetting instructions for manuscript.

Pages parameters: format A4, portrait orientation, document margins (cm): left — 2.5, right — 1.5, above — 2.0, below — 2.0, footer 1.3.

Text: font — Times New Roman, font size — 14, paragraph indent — 0.5, line spacing — 1.5, justified alignment.

The recommended manuscript size: not more than 20 pages of the specified format.

Use only standard abbreviations. Avoid abbreviations in the title and abstract. The full term for which an abbreviation stands should precede its first use in the text unless it is a standard unit of measurement.

All pages of the manuscript should be numbered.

The templates for the manuscript typesetting are presented on site: <http://www.ipiran.ru/journal/template.doc>.

8. The articles should enclose data both in **Russian and English:**

- title;
- author’s name and surname;
- affiliation — organization, its address with ZIP code, city, country, and official e-mail address;
- data on authors according to the format: (see site)

http://www.ipiran.ru/journal/issues/2013_07_01/authors.asp and

http://www.ipiran.ru/journal/issues/2013_07_01_eng/authors.asp;

- abstract (not less than 100 words) both in Russian and in English. Abstract is a short summary of the article that can be published separately. The abstract is the main source of information on the article and it could be included in leading information systems and data bases. The abstract in English has to be an original text and should not be an exact translation of the Russian one. Good English is required. In abstracts, avoid references and formulae;
 - indexing is performed on the basis of keywords. The use of keywords from the internationally accepted thematic Thesauri is recommended.
Important! Keywords must not be sentences;
 - Acknowledgments.
9. References. Russian references have to be presented both in English translation and Latin transliteration (refer <http://www.translit.ru>, option BGN).
Please take into account the following examples of Russian references appearance:
- Article in journal:**
Zhang, Z., and D. Zhu. 2008. Experimental research on the localized electrochemical micromachining. *Rus. J. Electrochem.* 44(8):926–930. doi:10.1134/S1023193508080077.
- Journal article in electronic format:**
Swaminathan, V., E. Lepkoswka-White, and B. P. Rao. 1999. Browsers or buyers in cyberspace? An investigation of electronic factors influencing electronic exchange. *JCMC* 5(2). Available at: <http://www.ascusc.org/jcmc/vol5/issue2/> (accessed April 28, 2011).
- Article from the continuing publication (collection of works, proceedings):**
Astakhov, M. V., and T. V. Tagantsev. 2006. Eksperimental’noe issledovanie prochnosti soedineniy “stal’–kompozit” [Experimental study of the strength of joints “steel–composite”]. *Trudy MGTU “Matematicheskoe modelirovanie slozhnykh tekhnicheskikh sistem” [Bauman MSTU “Mathematical Modeling of Complex Technical Systems” Proceedings]*. 593:125–130.
- Conference proceedings:**
Usmanov, T. S., A. A. Gusmanov, I. Z. Mullagalin, R. Ju. Muhametshina, A. N. Chervyakova, and A. V. Sveshnikov. 2007. Osobennosti proektirovaniya razrabotki mestorozhdeniy s primeneniem gidrorazryva plasta [Features of the design of field development with the use of hydraulic fracturing]. *Trudy 6-go Mezhdunarodnogo Simpoziuma “Novye resursoberegayushchie tekhnologii nedropol’zovaniya i povysheniya neftegazootdachi” [6th Symposium (International) “New Energy Saving Subsoil Technologies and the Increasing of the Oil and Gas Impact” Proceedings]*. Moscow. 267–272.
- Books and other monographs:**
Lindorf, L. S., and L. G. Mamikonians, eds. 1972. *Ekspluatatsiya turbogeneratorov s neposredstvennym okhlazhdeniem [Operation of turbine generators with direct cooling]*. Moscow: Energy Publs. 352 p.
- Dissertation and Thesis:**
Kozhunova, O. S. 2009. Tekhnologiya razrabotki semanticheskogo slovarya informatsionnogo monitoringa [Technology of development of semantic dictionary of information monitoring system]. PhD Thesis. Moscow: IPI RAN. 23 p.
- State standards and patents:**
GOST 8.586.5-2005. 2007. Metodika vypolneniya izmereniy. Izmerenie raskhoda i kolichestva zhidkostey i gazov s pomoshch’yu standartnykh suzhayushchikh ustroystv [Method of measurement. Measurement of flow rate and volume of liquids and gases by means of orifice devices]. M.: Standardinform Publs. 10 p.
Bolshakov, M. V., A. V. Kulakov, A. N. Lavrenov, and M. V. Palkin. 2006. Sposob orientirovaniya po krenu letatel’nogo apparata s opticheskoy golovkoy samonavedeniya [The way to orient on the roll of aircraft with optical homing head]. Patent RF No. 2280590.
- References in Latin transcription are presented in the original language.
References in the text are numbered according to the order of their first appearance; the number is placed in square brackets. All items from the reference list should be cited.
10. Manuscripts and additional materials are not returned to Authors by the Editorial Board.
11. Submissions of files by e-mail must include:
- the journal title and author’s name in the “Subject” field;
 - an article and additional materials have to be attached using the “attach” function;
 - an electronic version of the article should contain the file with the text and a separate file with figures.
12. “Informatics and Applications” journal is not a profit publication. There are no charges for the authors as well as there are no royalties.

Editorial Board address:

IPI RAN, Vavilova Str., 44, block 2, Moscow 119333, Russia
Ph.: +7 (499) 135 86 92, Fax: +7 (495) 930 45 05
e-mail: rust@ipiran.ru (to Prof. Rustem Seyful-Mulyukov)
<http://www.ipiran.ru/english/journal.asp>