

# О ПРИМЕНЕНИИ ТОПОЛОГИЧЕСКОГО ПОДХОДА К АНАЛИЗУ ПЛОХО ФОРМАЛИЗУЕМЫХ ЗАДАЧ ДЛЯ ПОСТРОЕНИЯ АЛГОРИТМОВ ВИРТУАЛЬНОГО СКРИНИНГА КВАНТОВО-МЕХАНИЧЕСКИХ СВОЙСТВ ОРГАНИЧЕСКИХ МОЛЕКУЛ. ЧАСТЬ 1: ОСНОВЫ ПРОБЛЕМНО ОРИЕНТИРОВАННОЙ ТЕОРИИ\*

И. Ю. Торшин<sup>1</sup>

**Аннотация:** Топологический подход к анализу плохо формализованных задач и теория хемографов являются расширениями алгебраического подхода к распознаванию, развиваемого в научной школе академика РАН Ю. И. Журавлёва. В первой части статьи предложен проблемно ориентированный формализм для разработки алгоритмов скрининговых оценок квантово-механических (КМ) свойств молекул по их химической структуре. Предложены способы введения метрик на множествах молекул и процедуры порождения «синтетических» признаков описаний, основанные на согласовании значений «экспертной» метрики на множестве значений свойств молекул со значениями настраиваемой метрики на множестве структур молекул.

**Ключевые слова:** алгебраический подход; хемоинформатика; размеченные графы; комбинаторный анализ разрешимости

**DOI:** 10.14357/19922264220106

## 1 Введение. Проблемная область

Алгебраический подход к анализу проблем распознавания, классификации включает формулировку и исследование критериев разрешимости и регулярности задач [1–3], корректности и полноты моделей алгоритмов [4–6]. Исследуются алгоритмы в виде композиций  $A = B \circ C \circ D$ , где  $B$  — распознающий оператор;  $C$  — корректирующая операция (корректор);  $D$  — решающее правило. Исходно в рамках алгебраического подхода наибольшее внимание уделяется именно формулировке и исследованию критериев разрешимости и регулярности задач распознавания/классификации, корректности и полноты моделей алгоритмов [1].

Исследование выполнимости этих критериев и «обучение» алгоритмов подразумевает, что задано формальное описание: *матрица информации* (набор признаков описаний объектов) и *информационная матрица* (отнесение объектов к определенным классам). В то же время во многих прикладных областях (биология, химия, медицина) встречаются

задачи, формальные описания которых могут быть получены очень многими способами.

Прикладные задачи, для которых не существует однозначного метода выделения объектов, определения признаков описаний и классов объектов на основе имеющегося «исходного описания» считаются *плохо формализуемыми* [4]. Адекватная формализация задач позволяет улучшить аккуратность и обобщающую способность соответствующих алгоритмов распознавания. С этой целью в русле алгебраического подхода разрабатываются топологическая и метрическая теории анализа данных [7–9].

К плохо формализуемым относятся и рассматриваемые в настоящей работе задачи типа «структура–свойство», где существенную неоднозначность представляют собой способы порождения признаков описаний молекул.

В статье представлены результаты совместного применения топологического подхода к плохо формализованным задачам [4] и методов теории анализа хемографов [8, 9] для разработки алгоритмов виртуального скрининга КМ свойств молекул.

\* Работа выполнена при поддержке РФФИ (проекты 19-07-00356, 18-07-00944, 20-07-00537).

<sup>1</sup> Федеральный исследовательский центр «Информатика и управление» Российской академии наук, tiy135@yahoo.com

## 2 Введение. О виртуальных скринингах молекул

Квантовая механика — один из ярких примеров успешного применения теории вероятностей, теории операторов, теории групп и функционального анализа в теоретической физике. Тем не менее точное аналитическое решение уравнения Шредингера имеется только для атома водорода, а высокоточные полумпирические схемы КМ-расчетов вычислительно затратны. В то же время задачи поиска перспективных материалов и лекарств подразумевают проведение скринингов пространства всех возможных «малых» органических молекул (более  $10^{11}$  структур [10]).

Для осуществления первоначальных стадий таких виртуальных скринингов необходимы алгоритмы для оценочных расчетов КМ-свойств молекул, позволяющие выделять подмножества молекул с экстремальными значениями КМ-свойств. Практически важны (а) точность, (б) обобщающая способность, (в) производительность вычислений и (г) интерпретируемость таких алгоритмов. В случае оценочных расчетов КМ-свойств молекул необходимо особо акцентировать интерпретируемость таких алгоритмов в терминах теории химической связи на уровне, доступном для химиков-практиков. Одно из возможных решений — алгоритмы, построенные на основании теории хемографов в рамках топологического подхода к анализу данных.

## 3 Основы топологической теории анализа данных

Алгебраический подход требует определения множества *начальных информации* ( $I_i$ ) и множества *конечных информации* ( $I_f$ ) [1]. Исследуемые алгоритмы  $A(\theta) : I_i \rightarrow I_f$  ( $\theta$  — вектор внутренних параметров) «обучаются» на основе конечных множеств прецедентов  $\text{Pr} \subset I_i \times I_f$ . Топологическая теория анализа данных позволяет проводить систематические исследования различных способов определения множеств  $I_i$  и  $I_f$  и так называемого множества оценок (область значений распознающего оператора  $B$ ).

Пусть  $\mathbf{X} = \{x_1, x_2, \dots, x_\alpha, \dots, x_{N_0}\}$  — множество исходных описаний объектов,  $\mathbf{X} \subseteq S$ , где  $S$  — пространство таких описаний. Пусть  $J_{\text{об}}$  — пространство допустимых признаков описаний объектов и определена функция  $D : S \rightarrow J_{\text{об}}$ . Определим  $J_{\text{об}} \subseteq I_1 \times I_2 \times \dots \times I_k \times \dots \times I_{n+1}$ ,  $I_i \subseteq I_1 \times I_2 \times \dots \times I_k \times \dots \times I_n$  и  $I_f \subseteq I_{n+1} \times I_{n+2} \times \dots \times I_{n+l}$  на основе  $I_k = \{\lambda_{k_1}, \lambda_{k_2}, \dots, \lambda_{k_b}, \dots, \lambda_{k_{|I_k|-1}}, \Delta\}$  —

множеств значений  $k$ -х компонент признакового описания,  $k = 1, \dots, n + l$ , где  $\Delta$  — неопределенность;  $n$  — число признаков;  $l$  — число целевых (прогнозируемых) переменных.

Формализация задачи, т. е. переход от множества  $\mathbf{X} \subseteq S$  к множеству прецедентов  $\mathbf{Q} \subseteq J_{\text{об}}$ , заключается в определении таких функций  $\Gamma_k : S \rightarrow I_k$ ,  $k = 1, \dots, n + l$ , что  $D(x_\alpha) = (\Gamma_1(x_\alpha) \times \dots \times \Gamma_k(x_\alpha) \times \dots \times \Gamma_{n+l}(x_\alpha))_\Delta$  [4, 6],  $m_\alpha = (\Gamma_1(x_\alpha) \times \dots \times \Gamma_n(x_\alpha))_\Delta$  — *вектор значений признаков*, а  $t_\alpha = (\Gamma_{n+1}(x_\alpha) \times \dots \times \Gamma_{n+l}(x_\alpha))_\Delta$  — *вектор значений  $l$  целевых переменных*, образующих *матрицу информации*  $\mathbf{M}_I = (m_\alpha)$ ,  $m_\alpha \in I_i$ , и *информационную матрицу*  $\mathbf{M}_F = (t_\alpha)$ ,  $t_\alpha \in I_f$ , соответственно. Функции  $D$  соответствует функция  $\varphi : 2^S \rightarrow 2^{J_{\text{об}}}$ ,  $\varphi(\mathbf{X}) = \{D(x_\alpha) | x_\alpha \in \mathbf{X}\}$ , формирующая множество прецедентов  $\mathbf{Q} = \{q_i | q_i = (m_i, t_i)\} \subseteq I_i \times I_f$ ,  $q_i[k] \in I_k$ ,  $i = 1, \dots, |\mathbf{X}|$ ,  $\mathbf{Q} = \varphi(\mathbf{X})$  [4]. Для  $\Gamma_k$  определена функция полного прообраза значения  $\lambda_{k_b} \in I_k$  в множестве  $\mathbf{X}$ ,  $\Gamma_k^{-1}(\lambda_{k_b}) \subseteq \mathbf{X}$ , а для функции  $D$  — функция полного прообраза объекта,  $D^{-1}(q) = \bigcap_{k=1, n} \Gamma_k^{-1}(q[k])$ .

**Определение 1.** *Регулярным будем называть множество прецедентов  $\mathbf{Q}$ , для которого выполнено  $\forall (q_1, q_2) : D^{-1}(q_1) \neq D^{-1}(q_2)$ . Если множества  $\mathbf{X}$  и  $\mathbf{Q}$  изоморфны, так что  $\forall x \in \mathbf{X} : x = D^{-1}(D(x))$ , то  $\mathbf{X}$  также регулярно [4]. В дальнейшем рассматриваются только регулярные множества  $\mathbf{X}/\mathbf{Q}$ .*

В подавляющем большинстве реальных задач множества  $I_k$  выбираются произвольно. Для получения более информативных «синтетических признаков описаний» [4] множество  $\mathbf{U}(\mathbf{X}) = \{\Gamma_k^{-1}(\lambda_{k_b})\}$  рассматривается как предбаза топологии  $\mathbf{T}(\mathbf{X}) = \{\emptyset, I, a \cup b, a \cap b : a, b \in \mathbf{U}(\mathbf{X})\}$ , где  $I = \{\mathbf{X}\}$  — единичный элемент. Введение отношения порядка на элементах  $\mathbf{T}(\mathbf{X})$  позволяет частично упорядочить элементы топологии  $\mathbf{T}(\mathbf{X})$  в *решетку*  $\mathbf{L}(\mathbf{T}(\mathbf{X})) = \{a \vee b, a \wedge b : a, b \in \mathbf{T}(\mathbf{X}), (a \geq b) \text{ или } (a \leq b)\}$  так, что  $a \leq b \equiv a \subseteq b$  и  $a \vee b = \sup(a, b)$  — объединение, а  $a \wedge b = \inf(a, b)$  — пересечение множеств  $a$  и  $b$ . Показано, что при регулярности  $\mathbf{X}/\mathbf{Q}$  (определение 1) решетка  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$  — булева и элементы  $\mathbf{U}(\mathbf{X})$  являются вершинами  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$  [4]. При этом возникают три фундаментальные разновидности признаков: *булевы* (вершина решетки  $\Gamma_k^{-1}(1)$ ), *категорные* (антицепи решетки) и *числовые* (цепи решетки) [4].

Таким образом,  $k$ -я числовая величина, заданная на  $\mathbf{Q}$  посредством  $I_k$ , соответствует некоторой цепи  $A_k(\mathbf{X})$  в  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$ , образованной множествами  $u(\lambda_{k_b}) = \bigcup_{\beta=1}^b \Gamma_k^{-1}(\lambda_{k_\beta})$ ,  $\lambda_{k_{b-1}} \leq \lambda_{k_b} \leq \lambda_{k_{b+1}}$ . Каждому значению  $\lambda_{k_b}$  соответствует множество объектов  $u(\lambda_{k_b})$  и единственное дополнение  $\neg u(\lambda_{k_b})$

(так как решетка булева). Совокупность точек  $\text{cdf}(A_k(\mathbf{X})) = \{(\lambda_{k_b}, |u(\lambda_{k_b})|/N)\}$  представляет эмпирическую функцию распределения (э.ф.р.)  $k$ -й переменной.

#### 4 Порождение «синтетических» числовых признаков

Булевой решетке  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$  сопоставляется метрическое пространство значений признаков  $\mathbf{M}_{\mathbf{L}}(\mathbf{L}(\mathbf{T}(\mathbf{X})), \rho_{\mathbf{L}})$  с метрикой  $\rho_{\mathbf{L}} : \mathbf{L}^2 \rightarrow \mathbf{R}^+$ . Над  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$  также определяется метрическое пространство объектов  $\mathbf{M}_q[\mathbf{L}(\mathbf{T}(\mathbf{X}))](\mathbf{Q}, \rho_q)$  с метрикой  $\rho_q : \mathbf{Q}^2 \rightarrow \mathbf{R}^+$ . Для определения метрик  $\rho_{\mathbf{L}}$  и  $\rho_q$  вводятся понятия *изотонной оценки* на решетке и *окрестности* элемента в топологии [4].

**Определение 2.** Оценка на решетке  $\mathbf{L}$  есть такая функция  $v : \mathbf{L} \rightarrow \mathbf{R}^+$ , что  $\forall a, b : v[a] + v[b] = v[a \wedge b] + v[a \vee b]$ . Оценка изотонна, если  $\forall a, b : a \subseteq \subseteq b \Rightarrow v[a] \geq v[b]$  (например, высота элемента  $h[x]$ ). Функция  $\rho(x, y) = v[x \vee y] - v[x \wedge y]$ , где  $v$  изотонна, является метрикой типа  $\rho_{\mathbf{L}}$  [4].

**Определение 3.** Окрестность  $u(x)$  точки  $x \in \mathbf{X}$  в топологии  $\mathbf{T}(\mathbf{X})$  есть произвольное  $u \in \mathbf{T}(\mathbf{X})$ ,  $x \in u$ ;  $u(x)$  разделяет  $x, y \in \mathbf{X}$  при  $(x \in u) \neq (y \in u)$ .

Любой объект  $q \in \mathbf{Q}$  представлен в  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$  как набор окрестностей  $\Gamma_k^{-1}(q[k])$ ,  $k = 1, \dots, n$ , что позволяет определить метрику  $\rho_q(q_1, q_2)$  как функцию  $f_q : \mathbf{R}^n \rightarrow \mathbf{R}^+$ ,  $\rho_q(q_1, q_2) = f_q((\rho_{\mathbf{L}}(\Gamma_k^{-1}(q_1[k]), \Gamma_k^{-1}(q_2[k])))$ ,  $k = 1, \dots, n$ .

**Определение 4.** Для аддитивных  $f_q$  расстояние  $\rho_q(q_1, q_2)$  между объектами  $q_1$  и  $q_2$  определяется посредством линейной комбинации  $S \left( \sum_{k=1, \dots, n} \omega_k s(\rho_{\mathbf{L}}(\Gamma_k^{-1}(q_1[k]), \Gamma_k^{-1}(q_2[k]))) \right)$ , где  $S, s : \mathbf{R} \rightarrow \mathbf{R}^+$  — произвольные «сглаживающие» функции, а  $\omega_k$  — вес  $k$ -го признакового описания.

Если  $S = 1$ ,  $s = 1/|\mathbf{Q}|$ ,  $\rho_{\mathbf{L}}(x, y) = (x \neq y)$ , ( $\omega_k = 1$ ), а  $\Gamma_k^{-1}$  и  $\Gamma_k$  определены для  $I_k = [0, 1]$ , то  $\rho_q$  — метрика Хэмминга. При  $S = \sqrt{\cdot}$ ,  $s(x) = x^p$ ,  $\rho_{\mathbf{L}}(x, y) = |x \Delta y|$ ,  $I_k \subset \mathbf{R}$ ,  $\rho_q$  — взвешенная метрика Минковского и т. д. В рамках топологического подхода порождение синтетических признаков осуществляется посредством метрики  $\rho_{\mathbf{L}}$  либо метрики  $\rho_q$  (последнее требует определения  $\rho_{\mathbf{L}}$ ).

#### 5 Синтетические признаки на основании метрики $\rho_{\mathbf{L}}$

**Определение 5.** Расстояние  $\rho_A : \mathbf{A}(\mathbf{X})^2 \rightarrow \mathbf{R}^+$  между цепями  $a = \langle a_1, \dots, a_i, \dots, \mathbf{I} \rangle$  и  $b = \langle b_1, \dots, b_j, \dots, \mathbf{I} \rangle$

есть сумма расстояний между соответствующими элементами:

$$\rho_A(a, b) = \min \left( \sum_{i=1, |a|} \rho_{\mathbf{L}} \left( a_i, \arg \min_{b_j \in b} \rho_{\mathbf{L}}(a_i, b_j) \right), \sum_{i=1, |b|} \rho_{\mathbf{L}} \left( b_j, \arg \min_{a_i \in a} \rho_{\mathbf{L}}(b_j, a_i) \right) \right).$$

Пусть  $\mathbf{A}(\mathbf{X})$  — множество всех цепей  $\mathbf{L}(\mathbf{T}(\mathbf{X}))$ . Алгоритмы прогнозирования  $k$ -й переменной соответствуют цепям в  $\mathbf{A}(\mathbf{X})$ , наиболее близким к цепи  $\mathbf{A}_k(\mathbf{X})$ . Пусть  $\mathbf{A}(\mathbf{X})_{1, n} \subset \mathbf{A}(\mathbf{X})$  соответствует цепям над элементами  $\mathbf{U}(\mathbf{X})$ ,  $k = 1, \dots, n$ . Тогда искомые алгоритмы соответствуют решению задачи

$$aa = \arg \min_{a \in \mathbf{A}(\mathbf{X})_{i, n}} \rho_A(\mathbf{A}_k(\mathbf{X}), a). \quad (1)$$

#### 6 Синтетические признаки на основании метрики $\rho_q$

В рамках данного подхода согласовываются значения настраиваемой «признаковой» метрики  $\rho_q$  с вектором весов ( $\omega_k$ )

$$\arg \min_{\{\omega_k\}} \sum_{m=1}^{|\mathbf{X}|} \sum_{j \neq m}^{|\mathbf{X}|} \mathbf{L}_F(\rho_q((\omega_k), X_m, X_j) - \rho_e(X_m, X_j)), \quad (2)$$

где  $\mathbf{L}_F$  — та или иная функция потерь. Например, в настоящей работе в качестве  $\mathbf{L}_F$  используется функция модуля  $\mathbf{L}_F(x) = |x|$ . Практически важным случаем экспертной метрики  $\rho_e$  является «скалярная» метрика (например, модуль разности значений).

**Теорема 1.** При использовании скалярной метрики  $\rho_e$  условие (1) для метрики Хэмминга  $\rho_q$  соответствует аддитивной схеме учета признаков.

**Доказательство.** Для  $\rho_e$  в виде модуля разности выполнены все три аксиомы метрики (так как они выполнены для любых трех коллинеарных точек). Пусть нулевой элемент входит во все множества  $I_k$ , так что можно определить расстояние от нулевого элемента до любого другого элемента множества  $I_k$  посредством экспертной метрики  $\rho_e$ . В  $\rho$ -конфигурации объектов, образованной метрикой  $\rho_q$ , за нулевую точку с номером  $m_0$  можно, вообще говоря, принять любой объект (например, центральный объект, соответствующий минимальной сумме расстояний до всех остальных). Пусть в  $m$ -сумме в выражении (2) одна из  $j$ -сумм соответствует нулевому элементу  $\mathbf{X}_{m_0}$ .

Использование скалярной метрики  $\rho_e$  в условии (2) соответствует проекции  $|\mathbf{X}|$ -мерной  $\rho_q$ -конфигурации на одномерную числовую ось значений  $I_k$ . Любая из  $j$ -сумм в (1) включает все объекты из  $\mathbf{X}$  и покрывает все множество значений  $I_k$ . Разница между любыми двумя значениями из  $I_k$  (т.е.  $\rho_e(\mathbf{X}_m, \mathbf{X}_j)$ ) также является отрезком на числовой оси  $k$ -й переменной, соответствующим сдвигу значений в множестве  $I_k$  на константу  $\lambda_{k\beta}$ .

Так как каждая из  $j$ -сумм в (2) соответствует одной и той же (с точностью до константы) задаче согласования метрик, то задача (2) с  $\mathbf{L}_F(x) = |x|$  может быть переформулирована через расстояния от нулевого элемента, т.е. произведен переход от оценки попарных расстояний к суммированию по всем объектам:

$$\arg \min_{\{\omega_k\}} \sum_{m \neq m_0}^{|\mathbf{X}|} |\rho_q(\{\omega_k\}, \mathbf{X}_{m_0}, \mathbf{X}_m) - \mathbf{T}_m|, \quad (3)$$

где  $\mathbf{T}_m = \rho_e(\mathbf{X}_{m_0}, \mathbf{X}_m)$  — значение прогнозируемой  $k$ -й числовой переменной для объекта  $\mathbf{X}_m$  (так как, по построению, точка  $\mathbf{X}_{m_0}$  соответствует нулевому значению). Если  $\rho_q$  определена как метрика Хэмминга (см. комментарий к определению 4), то задача в постановке (3) соответствует представлению  $\rho_q$  в виде линейной формы, т.е. аддитивной схеме учета признаков. Теорема доказана.

Таким образом, в случае скалярных  $\rho_e$  расстояния  $\rho_q$  служат своего рода «синтетическими» числовыми признаками объектов в множествах  $\mathbf{X}/\mathbf{Q}$ . Для практического применения (2) необходимо определить функции  $\Gamma_k$ , а для определения метрических пространств  $\mathbf{M}_L$  и  $\mathbf{M}_q$  — метрики  $\rho_L$  и  $\rho_q$ . В случае задач типа «структура—свойство» молекул для этого применяется теория хемографов.

## 7 О теории анализа размеченных графов

Для порождения признаков описаний молекулярных структур в рамках теории анализа размеченных графов вводится понятие хемографа [8]. Признаковые описания порождаются над разметками хемографов [9].

**Определение 6.** Граф  $G = (\mathbf{V}, \mathbf{E})$  — совокупность множества вершин  $\mathbf{V} = \mathbf{V}(G)$  и множества ребер  $\mathbf{E} = \mathbf{E}(G)$ ,  $\mathbf{E} \subset \mathbf{V}^2$ . Хемограф ( $\chi$ -граф) — конечный, связный, неориентированный, размеченный граф без петель, с кликовым числом не более 3. Множество  $\Gamma = \{(\mathbf{V}, \mathbf{E}) | \mathbf{V} \subset \mathbf{N}, \mathbf{E} \subset \mathbf{N}^2\}$ , где  $\mathbf{N}$  — натуральный ряд, есть множество всех графов.

Множество вершин  $\chi$ -графа  $\mathbf{X}$  соответствует множеству атомов молекулы, а множество ребер —

множеству химических связей молекулы. Хемографы строятся на основе «внешних» (декартовых) координат ядер атомов  $\mathbf{R}(X) = \{\vec{R}_i(X) | \vec{R}_i \in \mathbf{R}^3\}$ ,  $i = 1, \dots, |\mathbf{V}(X)|$ , или «внутренних» координат (межатомные расстояния), так что  $\chi$ -графу  $X$  сопоставлен ряд матриц  $\mathbf{M}(X) = (m_{ij}(X))$ ,  $m_{ij} \in \mathbf{R}$ ,  $i, j = 1, \dots, |\mathbf{V}(X)|$  (матрица межатомных расстояний, матрица смежности и др.). Химические формулы в различных «машинных» форматах (SMILES, XYZ и др.) суть упрощенные представления матриц  $\mathbf{M}(X)$  или координат  $\mathbf{R}(X)$ .

Признаковые описания  $\chi$ -графов вводятся на основании анализа специальной разновидности подграфов  $\chi$ -графов — цепей и узлов, для описания которых в работах [8, 9] вводится комплекс понятий и соответствующих обозначений: множество всех замкнутых подграфов  $\chi$ -графа  $X$ ,  $\Pi(X)$  (включает множество связных подграфов  $\mathbf{S}(X)$ , множество цепей  $\mathbf{C}(X)$ , множество  $\chi$ -узлов  $\mathbf{K}(X)$ ), операция объединения/пересечения множества подграфов ( $\Pi = \bigcup_{i=1}^{|\Pi|} \pi_i = (\bigcup_{i=1}^{|\Pi|} \mathbf{v}_i, \bigcup_{i=1}^{|\Pi|} \mathbf{e}_i)$ ;  $\hat{\Pi} = \bigcap_{i=1}^{|\Pi|} \pi_i$ ), условие образования графа  $O = X$  множеством  $\mathbf{O} \subseteq \Pi(X)$ , оператор  $\hat{c}^n$  вычисления цепей длины  $n$ , множество смежности (окружение) вершины  $v$ :  $\Gamma(v) = \hat{v} \hat{e} v$ , узел вершины  $v$ ,  $(\Gamma(v), \hat{e} v)$  и др. [8]. Такие подграфы необходимы для исследования изоморфизма  $\chi$ -графов.

**Определение 7.** Графы  $G_1$  и  $G_2$  изоморфны ( $G_1 \simeq G_2$ ) если существует взаимно однозначное соответствие между их вершинами и ребрами. Изоморфизм  $G_1 \simeq G_2$  соответствует существованию функции перенумерации  $\mu_I : \mathbf{N} \rightarrow \mathbf{N}$ , так что из  $G_2 = \mu_I(G_1)$  следует  $G_1 \simeq G_2$ . Определяется класс  $\mathbf{I}(G) = \{g \in \Gamma | \exists \mu_I : G = \mu_I(g)\}$ .

**Определение 8.** Пусть задан алфавит меток  $\mathbf{Y} = \{v_1, v_2, \dots, v_{n(Y)}\}$ . Функция разметки  $\mu_Y : \mathbf{V} \rightarrow \mathbf{Y}$  сопоставляет метку каждой вершине  $\chi$ -графа.

Для комбинаторного анализа свойств изоморфизма графов и для порождения признаков описаний  $\chi$ -графов в теории анализа размеченных графов вводятся множество  $\chi$ -цепей  $\hat{\mathbf{Y}}$  над алфавитом  $\mathbf{Y}$  (в том числе подмножества всех  $\chi$ -цепей длины  $n$ ,  $\hat{\mathbf{Y}}^n \subset \hat{\mathbf{Y}}$ ), множество  $\chi$ -узлов  $\hat{\mathbf{Y}} = \bigcup_{k=2} \hat{\mathbf{Y}}(k)$ ,  $\hat{\mathbf{Y}}(k) = \{\mathbf{Y} \times 2^{\mathbf{Y}^k}\}$ , понятия инварианта графа ( $\iota : \Gamma \rightarrow \mathbf{R}^n, n \in \mathbf{N} : \forall a \in \Gamma : b \in \mathbf{I}(a) \Rightarrow \iota(b) = \iota(a)$ ), кортеж-инварианта  $\iota : \Gamma \rightarrow \mathbf{R}^n, n \geq 2$ , полноты инварианта  $\forall a \in \Gamma : b \in \mathbf{I}(G) \Leftrightarrow \iota(a) = \iota(b)$ , изомерных по инварианту  $\iota$  графов, функции разметки  $\chi$ -цепей  $\mu_c : \mathbf{C} \rightarrow \hat{\mathbf{Y}}$  и  $\chi$ -узлов  $\mu_\kappa : \mathbf{K} \rightarrow \hat{\mathbf{Y}}$ ,  $\chi$ -фрагментов как элементов множеств  $\hat{\mu}_c^{-1}\alpha$  и  $\hat{\mu}_\kappa^{-1}\kappa$ , оператора вхождения множества подграфов  $\pi \subset \Gamma$  в  $\chi$ -граф  $\mathbf{X}$   $\hat{\beta}[\mathbf{X}] : 2^\Gamma \rightarrow \{0, 1\}$ , опе-

ратора числа вхождений множества подграфов в  $\mathbf{X}$   $\hat{\eta}[\mathbf{X}] : 2^\Gamma \rightarrow \mathbf{N}$  и др. [9].

Эти понятия позволяют проводить комбинаторный анализ изоформизма  $\chi$ -графов и порождать признаковые описания  $\chi$ -графов как  $\hat{\eta}[\mathbf{X}]\hat{\mu}_c^{-1}\alpha$ ,  $\hat{\eta}[\mathbf{X}]\hat{\mu}_\kappa^{-1}\kappa$ ,  $\hat{\beta}[\mathbf{X}]\hat{\mu}_c^{-1}\alpha$ ,  $\hat{\beta}[\mathbf{X}]\hat{\mu}_\kappa^{-1}\kappa$  (лемма 1 в работе [7]). Результат последовательного применения  $\hat{\mu}_c^{-1}$  к  $\alpha = \{\alpha \in \tilde{\mathbf{Y}}\}$  обозначим  $\hat{\mu}_c^{-1}\alpha = \{\hat{\mu}_c^{-1}\alpha, \alpha \in \alpha\}$ , оператора  $\hat{\mu}_\kappa^{-1}$  ко множеству  $\chi$ -узлов  $\kappa = \{\kappa \in \tilde{\mathbf{Y}}\} - \hat{\mu}_\kappa^{-1}\kappa = \{\hat{\mu}_\kappa^{-1}\kappa, \kappa \in \kappa\}$ , оператора  $\hat{\beta}$  ко множеству  $\tilde{\pi} = \{\pi_1, \pi_2, \dots, \pi_n\}$ ,  $\tilde{\pi} \subset \Gamma$  обозначим  $\hat{\beta}\tilde{\pi} = \{\hat{\beta}\pi_1, \hat{\beta}\pi_2, \dots, \hat{\beta}\pi_n\}$ .

Комбинаторный анализ разрешимости задачи анализа изоморфизма  $\chi$ -графов сводится к установлению локальной полноты некоторых кортеж-инвариантов [9]. Пусть  $\iota : \Gamma \rightarrow \mathbf{R}^n$ ,  $n \geq 2$ , — кортеж-инвариант, построенный над некоторым множеством из  $n$  элементарных инвариантов  $\iota_e \subset \mathbf{E}$ . Пусть для графа  $G$  выражение  $\iota_i(G) = \{\iota_i(G), i = 1, \dots, n\}$  будет означать множество значений инвариантов из  $\iota_e$ . Определим функцию нумерации элементарных инвариантов  $\lambda : \iota_e \rightarrow \mathbf{N}$  и оператор формирования кортеж-инварианта. Оператор формирования кортеж-инварианта  $\hat{\iota} : 2^{\mathbf{E}} \rightarrow \mathbf{R}^n$  по заданному  $\iota_e$  определим как  $\hat{\iota}_{\iota_e} = (\iota_j, \iota_k, \dots, \iota_l)$ ,  $\iota_j, \iota_k, \dots, \iota_l \in \iota_e$ ,  $\lambda(\iota_j) = 1 \leq \lambda(\iota_k) < \dots < \lambda(\iota_l) = n$ , а значение  $i$ -го элемента кортежа  $\hat{\iota}_{\iota_e}$  обозначим  $\hat{\iota}[i]\iota_e(G) = \iota(G)|\lambda(i) = i$ . В теореме 2 показана взаимосвязь между изоморфизмом  $\chi$ -графов и полнотой кортеж-инвариантов.

**Теорема 2.**  $\forall a, b \in \Gamma : |\mathbf{I}(a) \cap \mathbf{I}(b)| > 0 \Leftrightarrow \exists_{i=1, n} \hat{\iota}[i]\iota_e(a) \neq \hat{\iota}[i]\iota_e(b)$ , так что  $\hat{\iota}_{\iota_e}$  — полный кортеж-инвариант.

Доказательство приведено в работе [7]. Следствия теоремы важны для анализа множества прецедентов Pr.

**Следствие 1.** Полные инварианты могут быть образованы над подмножествами множеств  $\tilde{\mathbf{Y}}$  и  $\tilde{\mathbf{X}}$ , если для каждого  $\chi$ -графа  $\mathbf{X}$  множества прецедентов Pr существуют  $\kappa'(\mathbf{X}) = \{\kappa \in \kappa | \hat{\eta}[\mathbf{X}]\hat{\mu}_\kappa^{-1}\kappa = 1\}$  и  $\alpha'(\mathbf{X}) = \{\alpha \in \alpha | \hat{\eta}[\mathbf{X}]\hat{\mu}_c^{-1}\alpha = 1\}$  такие, что при  $\pi'(\mathbf{X}) = \hat{\mu}_\kappa^{-1}\kappa'(\mathbf{X}) \cup \hat{\mu}_c^{-1}\alpha'(\mathbf{X})$   $\tilde{\pi}(\mathbf{X}) = \mathbf{X}$ .

**Следствие 2.** Пусть в конечном Pr  $\subset \Gamma$  каждый граф  $G$  помечен меткой изоморфности  $\text{iso}(G) : \mathbf{I}(\Gamma) \rightarrow \mathbf{N}$ . Разрешимость

$$\forall_{a, b \in \text{Pr}} \text{iso}(a) \neq \text{iso}(b) \Rightarrow \iota(a) \neq \iota(b)$$

эквивалентна полноте

$$\forall_{a, b \in \text{Pr}} \text{iso}(a) \neq \text{iso}(b) \Rightarrow \exists_{i=1, |\mathbf{X}|} \hat{\iota}[i]\chi(a) \neq \hat{\iota}[i]\chi(b).$$

**Следствие 3.** Задача распознавания изоморфных графов разрешима тогда и только тогда, когда

$\sum_{a \in \text{Pr}} |\mathbf{i}\mu(a, \iota, \text{Pr}) \setminus \mathbf{i}(a, \text{Pr})| = 0$ , где  $\mathbf{i}\mu(G, \iota, \text{Pr})$  — множество изомерных  $G$  графов;  $\mathbf{i}(G, \text{Pr})$  — множество изоморфных  $G$  графов.

**Следствие 4.** Инвариант  $\iota$  полон при  $r_\iota(\iota, \text{Pr}) = 1 - (1/|\text{Pr}|^2) \sum_{a \in \text{Pr}} |\mathbf{i}\mu(a, \iota, \text{Pr}) \setminus \mathbf{i}(a, \text{Pr})| = 1$ .

**Следствие 5.** Пусть  $\phi(\hat{\iota}\chi, i, \text{Pr}) = \{(a, b) | a, b \in \text{Pr}, \hat{\iota}[i]\chi(a) \neq \hat{\iota}[i]\chi(b)\}$  и  $\phi(\hat{\iota}\chi, i, \text{Pr}) \cap \phi(\hat{\iota}\chi, j, \text{Pr}) = \emptyset$ . Тогда  $r_\iota(\iota, \text{Pr}) = \sum_{i=1}^{|\mathbf{X}|} \varphi_\iota(\hat{\iota}\chi, i, \text{Pr})$ , где  $\varphi_\iota(\hat{\iota}\chi, i, \text{Pr}) = |\phi(\hat{\iota}\chi, i, \text{Pr})|/|\text{Pr}|^2$ .

## 8 Топологическая теория анализа $\chi$ -графов

На основании алфавита  $\tilde{\mathbf{Y}}$  строятся различные множества  $\chi$ -цепей длины  $m$ ,  $\tilde{\mathbf{Y}}^m$ ,  $m = 1, \dots, m_{\max}$ , множество всех  $\chi$ -узлов  $\tilde{\mathbf{Y}}$  и затем соответствующие элементарные  $\chi$ - и кортеж-инварианты. Для  $\chi$ -графа  $X \in \mathbf{X}$  вычисляются множество всех  $\chi$ -цепей  $X$ ,  $\tilde{\mathbf{Y}}(X)$ , и множество всех  $\chi$ -узлов  $X$ ,  $\tilde{\mathbf{Y}}(X)$ . Для любых  $\alpha \in \tilde{\mathbf{Y}}(X)$ ,  $\kappa \in \tilde{\mathbf{Y}}(X)$  пусть  $\hat{\beta}[X]\hat{\mu}_c^{-1}\alpha$ ,  $\hat{\beta}[X]\hat{\mu}_\kappa^{-1}\kappa$  — булевы инварианты. Тогда для выбранных  $\alpha \subseteq \tilde{\mathbf{Y}}$  и  $\kappa \subseteq \tilde{\mathbf{Y}}$  таких, что  $|\alpha| + |\kappa| = n$ , вычисляются множества  $\hat{\mu}_c^{-1}\alpha$  и  $\hat{\mu}_\kappa^{-1}\kappa$ . Если множества  $\mathbf{I}_k$  определены как  $[0, 1]$ , то множество элементарных инвариантов  $\iota_e = \hat{\beta}\hat{\mu}_c^{-1}\alpha \cup \hat{\beta}\hat{\mu}_\kappa^{-1}\kappa$ , а если  $\mathbf{I}_k \subset \mathbf{N}$ , то  $\iota_e = \hat{\eta}\hat{\mu}_c^{-1}\alpha \cup \hat{\eta}\hat{\mu}_\kappa^{-1}\kappa$ . Определив функцию нумерации инвариантов  $\lambda : \iota_e \rightarrow \mathbf{N}$ , определяем функцию  $\Gamma_k$  как  $\Gamma_k(X) = \hat{\iota}[k]\iota_e(X)$ , а функцию  $D$  — как  $D(X) = \hat{\iota}_{\iota_e}(X)$ .

Таким образом, предбаза  $\mathbf{U}(\mathbf{X})$  топологии  $\mathbf{T}(\mathbf{X})$  состоит из подмножеств  $\chi$ -графов, соответствующих  $\chi$ -цепям или  $\chi$ -узлам. Для  $D(\mathbf{X}) = \hat{\iota}_{\iota_e}(\mathbf{X})$  задача (2) реализуется как хемометрический анализ [7]:

$$\arg \min_{\{\omega_k\}} \sum_{m_0=1}^{|\mathbf{X}|} \sum_{m \neq m_0}^{|\mathbf{X}|} \left| \sum_{k=1}^n \omega_k \hat{\iota}[k] \hat{\beta}[X_{m_0}] \pi \oplus \hat{\iota}[k] \hat{\beta}[X_m] \pi - |T_m - T_{m_0}| \right|,$$

где  $\pi \subseteq \hat{\mu}_c^{-1}\alpha \cup \hat{\mu}_\kappa^{-1}\kappa$  — некоторое «опорное» множество подграфов;  $T_m$  — оцениваемое свойство  $m$ -й молекулы, а выражение под знаком суммы по  $k$  соответствует  $\rho_q$ -метрике Хэмминга. Применяя теорему 1 и «сглаживающую» функцию  $S : \mathbf{R} \rightarrow \mathbf{R}^+$  (определение 4), получаем:

$$\arg \min_{\{\omega_k\}} \sum_{m=1}^{|\mathbf{X}|} \left| S \left( \sum_{k=1}^n \omega_k s \left( \hat{\iota}[k] \hat{\beta}[X_m] \pi \right) \right) - T_m \right|.$$

## 9 Заключение

Для построения проблемно ориентированных метрик и алгоритмов прогнозирования КМ свойств

молекул по их химической структуре разработанные признаковые описания хемографов анализируются посредством теории топологического анализа данных. В результате анализа топологии  $T(X)$  и решетки  $L(T(X))$  с использованием условий (1) и (2) могут быть получены различные алгоритмы прогнозирования произвольной числовой переменной на основе структуры хемографов. Эти алгоритмы одновременно являются алгоритмами порождения «синтетических» числовых признаков, информативных относительно этой переменной. Такие «синтетические» признаки могут использоваться в алгоритмических конструкциях алгебраического подхода и в обычных методах машинного обучения (регрессия, нейронные сети, метрические методы и др.). Сопоставление полученных моделей алгоритмов с формализмом квантовой механики, равно как и экспериментальная апробация соответствующих алгоритмов, будут представлены во второй части статьи.

## Литература

1. Журавлев Ю. И. Об алгебраическом подходе к решению задач распознавания или классификации // Проблемы кибернетики, 1978. № 33. С. 5–68.
2. Рудаков К. В., Торшин И. Ю. Вопросы разрешимости задачи распознавания вторичной структуры белка // Информатика и её применения, 2010. Т. 4. Вып. 2. С. 25–35.
3. Рудаков К. В., Торшин И. Ю. Анализ информативности мотивов на основе критерия разрешимости в задаче распознавания вторичной структуры белка // Информатика и её применения, 2011. Т. 5. Вып. 4. С. 40–50.
4. Torshin I. Yu., Rudakov K. V. On the theoretical basis of metric analysis of poorly formalized problems of recognition and classification // Pattern Recognition Image Analysis, 2015. Vol. 25. No. 4. P. 577–587.
5. Torshin I. Yu., Rudakov K. V. Combinatorial analysis of the solvability properties of the problems of recognition and completeness of algorithmic models. Part 2: Metric approach within the framework of the theory of classification of feature values // Pattern Recognition Image Analysis, 2017. Vol. 27. No. 2. P. 184–199.
6. Torshin I. Yu., Rudakov K. V. On the procedures of generation of numerical features over partitions of sets of objects in the problem of predicting numerical target variables // Pattern Recognition Image Analysis, 2019. Vol. 29. No. 4. P. 654–667. doi: 10.1134/S1054661819040175.
7. Ruddigkeit L., van Deursen R., Blum L., Reymond J. Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17 // J. Chem. Inf. Model., 2012. Vol. 52. No. 11. P. 2864–2875. doi: 10.1021/ci300415d.
8. Torshin I. Yu., Rudakov K. V. On the application of the combinatorial theory of solvability to the analysis of chemographs. Part 1: Fundamentals of modern chemical bonding theory and the concept of the chemograph // Pattern Recognition Image Analysis, 2014. Vol. 24. No. 1. P. 11–23.
9. Torshin I. Yu., Rudakov K. V. On the application of the combinatorial theory of solvability to the analysis of chemographs. Part 2: Local completeness of invariants of chemographs in view of the combinatorial theory of solvability // Pattern Recognition Image Analysis, 2014. Vol. 24. No. 2. P. 196–208.
10. Torshin I. Yu., Rudakov K. V. Topological data analysis in materials science: The case of high-temperature cuprate superconductors // Pattern Recognition Image Analysis, 2020. Vol. 30. No. 2. P. 262–274. doi: 10.1134/S1054661820020157.

Поступила в редакцию 30.03.21

---

# ON THE APPLICATION OF A TOPOLOGICAL APPROACH TO ANALYSIS OF POORLY FORMALIZED PROBLEMS FOR CONSTRUCTING ALGORITHMS FOR VIRTUAL SCREENING OF QUANTUM-MECHANICAL PROPERTIES OF ORGANIC MOLECULES. PART 1: THE BASICS OF THE PROBLEM-ORIENTED THEORY

I. Yu. Torshin

Federal Research Center “Computer Science and Control” of the Russian Academy of Sciences, 44-2 Vavilov Str., Moscow 119333, Russian Federation

**Abstract:** The topological approach to the analysis of poorly formalized problems and the theory of chemographs are extensions of Zhuravlev’s algebraic approach to recognition. In the first part of the article, a problem-oriented

formalism is proposed aimed at development of algorithms for screening assessments of the quantum-mechanical properties of molecules on the basis of their chemical structure. Methods for introducing metrics on sets of molecules and procedures for generating “synthetic” feature descriptions are proposed. The latter are generated by matching the values of some “expert” metric on the set of molecular properties to a tunable metric on the set of molecular structures.

**Keywords:** algebraic approach; chemoinformatics; labeled graphs; combinatorial solvability analysis

**DOI:** 10.14357/19922264220106

## Acknowledgments

This work was supported in part by RFBR grants 19-07-00356, 18-07-00944, 20-07-00537.

## References

- Zhuravlev, Yu. I. 1978. Ob algebraicheskom podkhode k resheniyu zadach raspoznavaniya ili klassifikatsii [On algebraic approach to recognition and classification problems]. *Problemy kibernetiki* [Cybernetic Problems] 33:5–68.
- Rudakov, K. V., and I. Yu. Torshin. 2010. Voprosy razreshimosti zadachi raspoznavaniya vtorichnoy struktury belka [Questions of solvability of the problem of recognition of the secondary structure of a protein]. *Informatika i ee Primeneniya — Inform Appl.* 4(2):25–35.
- Rudakov, K. V., and I. Yu. Torshin. 2011. Analiz informativnosti motivov na osnove kriteriya razreshimosti v zadache raspoznavaniya vtorichnoy struktury belka [Analysis of the informativeness of motives based on the criterion of solvability in the problem of recognizing the secondary structure of a protein]. *Informatika i ee Primeneniya — Inform Appl.* 5(4):40–50.
- Torshin, I. Yu., and K. V. Rudakov. 2015. On the theoretical basis of metric analysis of poorly formalized problems of recognition and classification. *Pattern Recognition Image Analysis* 25(4):577–587.
- Torshin, I. Yu., and K. V. Rudakov. 2017. Combinatorial analysis of the solvability properties of the problems of recognition and completeness of algorithmic models. Part 2: Metric approach within the framework of the theory of classification of feature values. *Pattern Recognition Image Analysis* 27(2):184–199.
- Torshin, I. Yu., and K. V. Rudakov. 2019. On the procedures of generation of numerical features over partitions of sets of objects in the problem of predicting numerical target variables. *Pattern Recognition Image Analysis* 29(4):654–667. doi: 10.1134/S1054661819040175.
- Ruddigkeit, L., R. van Deursen, L. Blum, and J. Raymond. 2012. Enumeration of 166 billion organic small molecules in the chemical universe database GDB-17. *J. Chem. Inf. Model.* 52(11):2864–2875. doi: 10.1021/ci300415d.
- Torshin, I. Yu., and K. V. Rudakov. 2014. On the application of the combinatorial theory of solvability to the analysis of chemographs. Part 1: Fundamentals of modern chemical bonding theory and the concept of the chemograph. *Pattern Recognition Image Analysis* 24(1):11–23.
- Torshin, I. Yu., and K. V. Rudakov. 2014. On the application of the combinatorial theory of solvability to the analysis of chemographs. Part 2: Local completeness of invariants of chemographs in view of the combinatorial theory of solvability. *Pattern Recognition Image Analysis* 24(2):196–208.
- Torshin, I. Yu., and K. V. Rudakov. 2020. Topological data analysis in materials science: The case of high-temperature cuprate superconductors. *Pattern Recognition Image Analysis* 30(2):262–274. doi: 10.1134/S1054661820020157.

Received March 30, 2021

## Contributor

**Torshin Ivan Yu.** (b. 1972) — Candidate of Science (PhD) in physics and mathematics, Candidate of Science (PhD) in chemistry, senior scientist, A. A. Dorodnicyn Computing Center, Federal Research Center “Computer Science and Control” of the Russian Academy of Sciences, 40 Vavilov Str., Moscow 119333, Russian Federation; tiy135@yahoo.com