

# Об основах метрического анализа плохо формализованных задач распознавания и классификации

И. Ю. Торшин, К.В.Рудаков

Московский Физико-Технический Институт

Вычислительный центр им. А.А. Дородницына РАН

*алгебраический подход, теория классификации **значений признаков**, комбинаторная теория разрешимости, метрический анализ данных, метрические сгущения*

# Рамки алгебраического подхода к проблемам распознавания и классификации

- **Заданы**

- матрица информации (набор признакововых описаний объектов)
- информационная матрица (отнесение объектов к определенным классам)

- **Исследовать:**

- Разрешимость/регулярность задач
- Корректность/полноту моделей алгоритмов

# Формализованные задачи – задано формальное описание

- ***Задача формализована*** - заданы матрица информации и информационная матрица
- ***Множество прецедентов*** - совокупность пар соответствующих строк матрицы информации и информационной матрицы



Построение алгоритмов  
распознавания

# Плохо формализованные задачи

- Формальные описания существенно различны, даже могут быть получены существенно различными способами на основании некоторого «исходного описания» задачи в терминах проблемной области.
- На основе имеющегося «исходного описания» не существует однозначного метода определения
  - 1) **Объектов**
  - 2) **Признаковых описаний объектов**
  - 3) **Классов объектов**

*Процедуре формализации той или иной прикладной задачи редко уделяется должное внимание...*

- Адекватная формализация - существенное улучшение аккуратности и обобщающей способности алгоритмов распознавания
- Примеры из различных областей: анализ биомедицинских данных, биоинформатика, хемоинформатика, анализ текстов

# Примеры плохо формализованных задач

Хемоинформатика

Биоинформатика

Задача	Однозначность формализации		
	Объекты	Классы	Признаки
Распознавание вторичной структуры белка	-	-	-
Аннотация генома	+	-	-
Поиск структурно схожих молекул	+	+	-
Поиск релевантных научных публикаций	+	+	-
Анализ биомедицинских данных	+	-	+-

Биомедицина

*При адекватной формализации этих задач в рамках соответствующих проблемно-ориентированных теорий было достигнуто существенное улучшение качества распознавания.*



# Задача распознавания вторичной структуры белка

- распознавание = перевод последовательности из 20-буквенного алфавита в 3-буквенный

VHLPREEKSA...

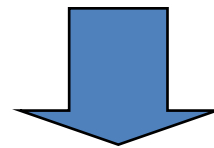


LLHHHHHHHH...

Алфавит «А»

1D

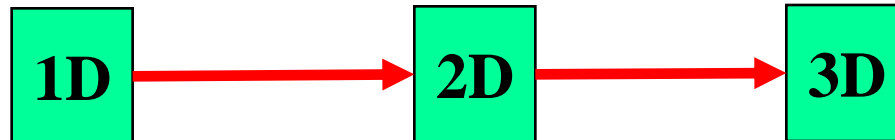
$A = \{A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y\}$



Алфавит «В»

2D

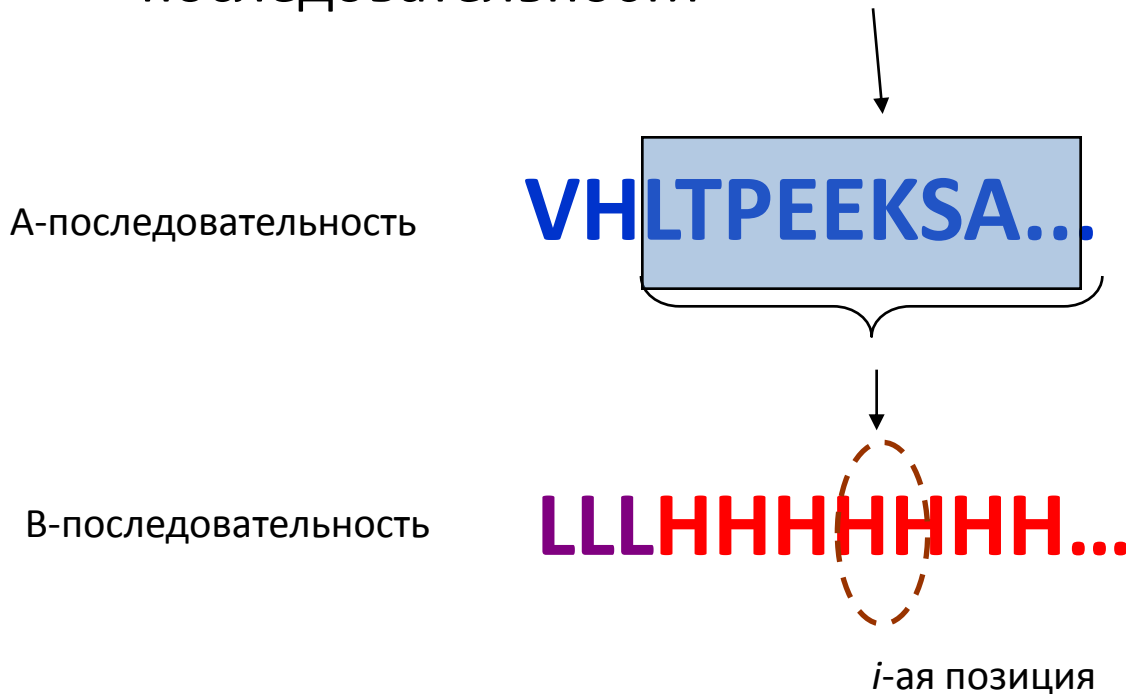
$B = \{H, S, L\}$





# Гипотеза локальной разрешимости задачи позволяет ввести «объекты» задачи распознавания

- В  $i$ -ой позиции, состояние В-последовательности определяется локальным контекстом (**окрестностью**) в А-последовательности



# Формальное описание локальности

слово  $U = \{u_1, u_2, \dots, u_n\}$

ведущая позиция  $i, 1 \leq i \leq n$

«маска»  $\hat{m} = \{\mu_1, \mu_2, \dots, \mu_m\}$

позиции маски  $\mu_i \in \mathbb{Z}$

$\mu_1 < \mu_2 < \dots < \mu_m$

размерность маски

$|\hat{m}| = m$

протяженность маски

$[\hat{m}] = \mu_m - \mu_1 + 1$

оператор выбора под слова  $\eta(i, \hat{m}, U)$

$\eta(i, \hat{m}, U) = \begin{cases} u_{i+\mu_1} u_{i+\mu_2} \dots u_{i+\mu_m}, & \text{если } i + \mu_1 \geq 1, i + \mu_m \leq n, \\ \emptyset & \text{в противном случае.} \end{cases}$

система масок  $M = \{\hat{m}_1, \hat{m}_2, \dots, \hat{m}_N\}$

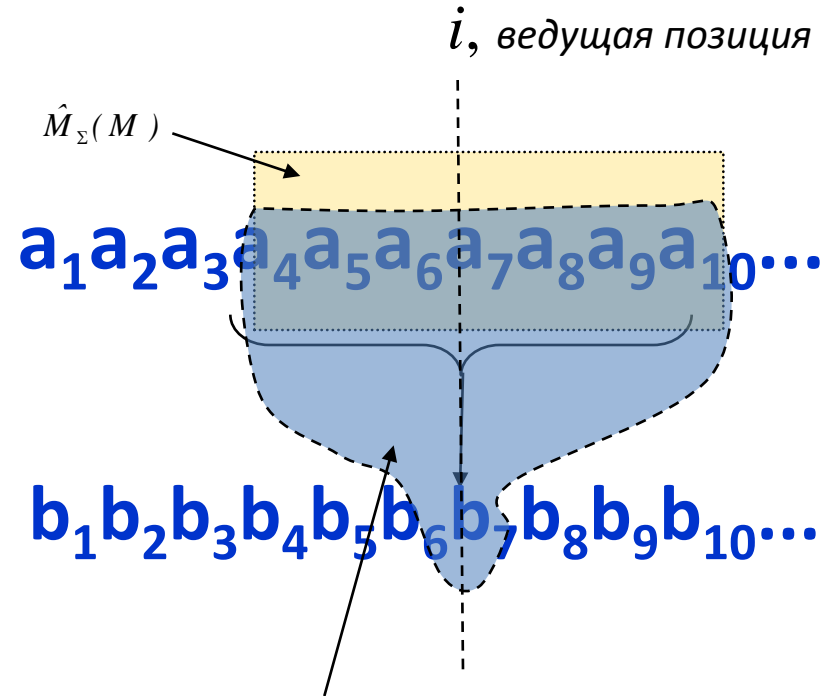
$\mathbf{M}_n^m$

объединенная маска

$M_\Sigma(M) = \bigcup_{k=1}^{|M|} \hat{m}_k$

**Мотив** - пара из под слова и соответствующей маски

$\kappa_\alpha = (\hat{m}_\alpha, \vec{v}_\alpha)$



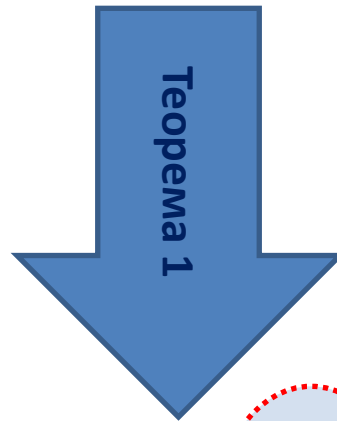
**Элементарный объект  $q$**

$q = (\eta(i, \hat{M}_\Sigma(M), V), W_i) \in Q$

$Q = Q(Pr, M)$  - множество объектов

# Условие локальной разрешимости

$$\forall_{Pr} (V^1, W^1), (V^2, W^2) : (V^1 = V^2) \Rightarrow (W^1 = W^2)$$



$$\hat{m} = \{\mu_1, \mu_2, \dots, \mu_m\}$$

$$\eta(i, \hat{m}, U)$$

$$M = \{\hat{m}_1, \hat{m}_2, \dots, \hat{m}_N\}$$

$$Q = Q(Pr, M)$$

$$\kappa_\alpha = (\hat{m}_\alpha, \vec{v}_\alpha)$$

$$\forall_{Q(Pr, M)} (i, j) : w_i \neq w_j \Rightarrow \exists_{\alpha=1}^{/K_1/} \alpha : (\kappa_\alpha \in \bar{v}_i) \neq (\kappa_\alpha \in \bar{v}_j)$$



# Примеры плохо формализованных задач

Хемоинформатика

Задача	Однозначность формализации		
	Объекты	Классы	Признаки
Распознавание вторичной структуры белка	-	-	-
Аннотация генома	+	-	-
Поиск структурно схожих молекул	+	+	-
Поиск релевантных научных публикаций	+	+	-
Анализ биомедицинских данных	+	-	+ -

Биомедицина

процедуры преобразования  
признаковых описаний

# О множествах начальных и конечных информации в плохо формализованных задачах

Хемоинформатика

Задача	Однозначность формализации		
	Объекты	Классы	Признаки
Распознавание вторичной структуры белка	-	-	-
Аннотация генома	+	-	-
Поиск структурно схожих молекул	-	+	-
Поиск релевантных научных публикаций	+	+	-
Анализ биомедицинских данных	+	-	+/-

Биомедицина

$I_i$

$I_f$

$I_f$ 

# Формализация задачи поиска

## «ОПТИМАЛЬНЫХ» $I_f$

$$Pr \subseteq I_i \times I_f, Pr \neq \emptyset, Pr = \{(\alpha, \beta)\}$$

*разрешимость задачи*

$$\forall_{Pr} (\alpha_1, \beta_1), (\alpha_2, \beta_2) : \alpha_1 = \alpha_2 \Rightarrow \beta_1 = \beta_2$$

*множество признаков описаний  
прецедентов*

$$Q(Pr) = \{\alpha_1, \alpha_2, \dots, \alpha_N\}, Q(Pr) \subseteq I_i$$

*множество описаний классов*

$$B(Pr) = \{\beta \mid \exists \alpha, (\alpha, \beta) \in Pr\} \cup \Delta$$

**Теорема 1.** *Задача разрешима тогда и только тогда, когда совокупность  $\beta$ -подмножеств признаков описаний над  $B(Pr)$  является разбиением  $Q(Pr)$ .*

*оператор выделения  $\beta$ -  
подмножества признаков описаний*

$$\hat{\alpha}(\beta, Pr) = \{\alpha_i \mid (\alpha_i, \beta_i) \in Pr, \beta_i = \beta\}$$

*оператор покрытия*

$$\check{\alpha}(B, Pr) = \bigcup_{k=1}^{|B|} \hat{\alpha}(\beta_k, Pr)$$

*$(\xi, \gamma)$ -разбиение – учет  
противоречивости реальных  
множеств прецедентов.*

$$B = \{\beta_1, \beta_2, \dots, \beta_m\}$$

$$\begin{cases} |Q(Pr) \cap \check{\alpha}(B, Pr)| \geq \xi \cdot |Q(Pr)| \\ \forall_{k=1}^{|B|} |\hat{\alpha}(\beta_k, Pr) \cap \check{\alpha}(B \setminus \beta_k, Pr)| \leq \gamma \cdot |Q(Pr)| \end{cases}$$

# Принцип максимального покрытия

Принципом максимального покрытия назовём утверждение о максимальном вкладе  $\alpha$ -го элемента словаря  $B_j \in I_{j,\delta(j)}(\tilde{B})$ ,  $\alpha > 1$ , в покрытие  $Q(Pr')$

$$\bigvee_{i=j}^{j+\delta(j)} \beta_i, \quad \left| \tilde{\alpha}(I(B_j, 1.. \alpha), Pr') \right| \geq \left| \tilde{\alpha}(I(B_j, 1.. \alpha - 1), Pr') \cup \hat{\alpha}(\beta_i, Pr') \right|$$

**Lm 4.**  $\sum_{i=1}^{|\mathcal{B}|} \left| \hat{\alpha}(\beta_i, Pr') \right| \geq |Q(Pr')|$

- необходимое условие покрытия

**Th. 5.** Принцип максимального покрытия ведущих позиций – необходимое условие оптимальности словаря.

Упорядоченный список  $(\tilde{B}(n)) = \{\beta_1, \beta_2, \dots, \beta_j \dots\}$   
 $|\hat{\alpha}(\beta_1, Pr')| \geq |\hat{\alpha}(\beta_2, Pr')| \geq |\hat{\alpha}(\beta_j, Pr')| \geq \dots$

$\delta(\xi)$ -покрытие  $I_{j,\delta}(\tilde{B}(n)) \subseteq I(\tilde{B}(n)) = \{\beta_j, \beta_{j+1}, \dots, \beta_{j+\delta-1}\}$ :

(7)  $\left| \tilde{\alpha}(I_{j,\delta}(\tilde{B}), Pr') \right| \geq \xi \cdot |Q(Pr')|$

$\xi$ -тупиковое покрытие  $B^T$  - (7) нарушено для любого подмножества  $B^T$ .

**Th. 6.**  $\delta(\xi)$ -покрытия содержат  $\xi$ -тупиковые покрытия.

**Th. 7.** Принцип максимального покрытия – необходимое, а при  $\xi'(B_j, Pr') \leq \xi$ ,  $\gamma'(B_j, Pr') \leq \gamma$  достаточное условие  $\xi, \gamma$ -разбиения

↳ **Теорема 1**

**Th. 8.** Множество оптимальных словарей есть подмножество множества  $\Gamma(\xi)$ , каждый элемент которого удовлетворяет принципу максимального покрытия.

Еще одна форма записи критерия разрешимости...

*I<sub>i</sub>?*

# Формализм для проведения метрического анализа плохо формализованных задач

- Аксиома соответствия
- Топологии и решётки над множеством исходных описаний
- Метрические пространства
- Приложения:
  - *введение метрик на множествах признаков описаний,*
  - *метрики на множествах объектов*
  - *анализ «взаимодействий» разнородных признаков описаний.*



# Об аксиоме (однозначного) соответствия

- $X = \{x_i\}$  – конечное множество исходных описаний объектов в проблемной области
- $\{m_i\}$  матрица информации,  $\{l_i\}$  – информационная матрица
- $Q = \{q_i | q_i = (m_i, l_i)\}$  – множество «объектов ИАД»
- Формализация задачи – нахождение функции  $\varphi: X \rightarrow Q$
- Аксиома соответствия – ограничения на  $\varphi$ :
  - Слабая форма:  $\varphi$  – инъективна
  - Сильная форма:  $\varphi$  – биекция

# О топологиях над множеством исходных описаний

- Процесс нахождения  $\phi$ :
  - От множества  $X$  к некоторой топологии  $\tau(X)$ ,
  - От  $\tau(X)$  к  $\tau$ -пространству  $T(X)=(X,\tau(X))$
  - От  $T(X)$  к решётке  $L(X)=L(T(X))$
  - От  $L(X)$  к метрическому пространству  $M(\tau(X))=(\tau(X),\rho_{L(X)})$ .
  - От  $M(\tau(X))$  к  $Q(X)$
- По сильной форме аксиомы соответствия - существуют биекции между всеми этими представлениями  $X$

# Пример. $X$ - множество описаний пациентов

- Разбиения на неупорядоченные подмножества
  - описаний мужчин и женщин,
  - подмножества, соответствующие тому или иному симптому,
  - числовые и «категориальные» признаки: подмножества, соответствующие *значениям признаков*
  - ...

# Формализуем...

- Над множеством исходных описаний  $X$  задается  $\pi(X)$  - система подмножеств  $X$   $\pi(X) = \{\emptyset, a_1, a_2, \dots, a_n \mid a_i \subseteq X\}$ 
  - объединение элементов  $\pi(X)$   $\check{\pi}(X) = \{\bigcup_{i=1..n} a_i\}$
  - пересечение элементов  $\pi(X)$   $\hat{\pi}(X) = \{\bigcap_{i=1..n} a_i\}$
- Считая все подмножества  $\pi(X)$  «открытыми», формируем множества  $\underline{\tau}(X) = \{\hat{a} \mid a \subseteq \pi(X)\}$   $\tilde{\tau}(X) = \{\check{b} \mid b \subseteq \underline{\tau}(X)\}$
- **Теорема 1.** Система множеств  $\tau(X)$  - топология над  $X$ , а  $\pi(X)$  - предбаза  $\tau(X)$  тогда и только тогда, когда  $\check{\pi}(X) = X$ 
  - Следствие.  $\check{\pi}(X) = X$  - необходимое условие выполнения сильной формы аксиомы соответствия.
  - Следствие.  $\check{\pi}(X) = X$  - необходимое условие разрешимости и регулярности задач распознавания (при выполнении сильной формы аксиомы).

# Выполнимость условия $\check{\pi}(X) = X$ - критерий оценки набора исходных признаков описаний

- При  $\check{\pi}(X) \subset X$  ни одна задача распознавания, которая могла бы быть сформулирована для элементов множества  $X$ , *не является формализуемой.*
- *Условие  $\check{\pi}(X) = X$ , наряду с разрешимостью и регулярностью, является важным конструктивным критерием качества постановки задачи.*

*BigDATA*

# Анализ аксиом отделимости в

## топологическом пространстве $T(X)=(X,\tau(X))$

- $T_0$  (аксиома Колмогорова)

$$\forall_{x,y \in X} \exists_{u \in \tau(X)} : (x \in u) \neq (y \in u)$$

- $T_1$

$$\forall_{x,y \in X} \exists_{u,v \in \tau(X)} : (x \in u) \wedge (y \in v) \wedge (x \notin v) \wedge (y \notin u)$$

- $T_2$  (аксиома Хаусдорфа):

$$\forall_{x,y \in X} \exists_{u,v \in \tau(X)} : (x \in u) \wedge (y \in v) \wedge (u \cap v = \emptyset)$$

- $T_3$  (аксиома Тихонова):

$$\forall_{x \in X} \forall_{v \in \tau(X) \mid x \in v} \exists_{u \in \tau(X)} : (x \in u) \wedge (u \subset v)$$

- $T_4$  (аксиома нормальности):

$$\forall_{u,v \in \tau(X) : u \cap v = \emptyset} \Rightarrow \exists_{u_1, v_1 \in \tau(X)} : (u \subset u_1, v \subset v_1) \wedge (u_1 \cap v_1 = \emptyset)$$

# Регулярность (по Журавлёву) – достаточное условие разрешимости задачи

- *множество прецедентов*  $Q \subseteq I_i \times I_f$ ,  $Q \neq \emptyset$ ,  $Q = \{(\alpha_j, \beta_j)\}$
- *разрешимость задачи – существует*  $f : I_i \rightarrow I_f$   
 $\bigvee_Q (\alpha_1, \beta_1), (\alpha_2, \beta_2) : \alpha_1 = \alpha_2 \Rightarrow \beta_1 = \beta_2$
- *регулярность задачи*  
 $\bigvee_Q (\alpha_1, \beta_1), (\alpha_2, \beta_2) : \alpha_1 \neq \alpha_2$
- **Теорема 2.** Пусть выполнена аксиома соответствия, так что задано  $Q = \phi(X)$ . Аксиома отделимости  $T_1$  выполнена для  $\tau(X)$  тогда и только тогда, когда выполнено условие регулярности множества прецедентов  $Q$ .
  - Следствие. Регулярность следует из аксиомы  $T_2$ .
  - Следствие. При выполнении условия регулярности  $T(X) = (X, \tau(X))$  - дискретное топологическое пространство.

# Решётки над множеством ИСХОДНЫХ ОПИСАНИЙ

- Формально, топология  $\tau(X)$  - неупорядоченное множество. В то же время, элементы  $\tau(X)$  могут являться подмножествами друг друга, так что между ними возникает естественное упорядочение по включению
- **Теорема 3.** *Решётка  $L(X)$  над  $L(X)$  является булевой при выполнении аксиомы соответствия и условия регулярности для  $X$ .*
  - *Следствие.  $L(X)$  – дистрибутивна*



# Метрическое пространство $M(\tau(X))$

- *Оценка на решётке  $L(X)$  – функция  $v : L(X) \rightarrow \mathbb{R}^+$ , для которой выполнено условие оценки  $\forall a, b \in L(X) : v[a] + v[b] = v[a \wedge b] + v[a \vee b]$*
- *Изотонная оценка:  $v$ - монотонна*
- **Лемма 5 (Биркгофа).** *Пусть  $L$  - дистрибутивная решётка,  $v$  - оценка на  $L$ ,  $d(x, y) = v[x \vee y] - v[x \wedge y]$ . Тогда, для произвольных выполнено  $d(a \vee x, a \vee y) + d(a \wedge x, a \wedge y) = d(x, y)$*
- **Теорема 6.** *В решётке  $L(X)$  с изотонной оценкой  $v[\ ]$  функция  $\rho(x, y) = v[x \vee y] - v[x \wedge y]$  - метрика, образующая метрическое пространство  $M(\tau(X)) = (\tau(X), \rho(x, y))$* 
  - *Следствие. При  $v[x] = h[x]$   $\rho(x, y) = |x \Delta y|$*
  - *Следствие. В конечной  $L(X)$  с заданной изотонной оценкой  $v[\ ]$  определено  $\varepsilon_{min} > 0$  - минимальное, не равное нулю значение расстояния  $\rho(x, y)$  между элементами  $(x, y)$ .*
  - *Следствие. Произвольное подпространство  $\tilde{m} \subseteq M(\tau(X))$  также является метрическим пространством.*

# Выводы

- Разработан концептуальный аппарат для анализа математической структуры исследуемых наборов *разнородных признаков описаний* в произвольной задаче распознавания или классификации.
- Регулярность (по Журавлеву) множества объектов позволяет переносить результаты анализа метрического пространства на элементы множества  $X$ , т.е. в изучаемую проблемную область.

# Выводы

- Условие  $\check{\pi}(X) = X$  - конструктивный критерием качества постановки задачи, допускающий комбинаторное тестирование
- Разрабатываются приложения аппарата к
  - введению метрик на множествах признаков
  - анализу «значимости взаимодействий» между разнородными признаковыми описаниями.
  - введению метрик на множествах объектов