

# Вопросы эффективности логических корректоров

Дюкова Елена Всеволодовна  
Журавлёв Юрий Иванович  
Прокофьев Пётр Александрович\*

ВЦ РАН, Москва

19 сентября 2015 г.

# Задача корректного распознавания по прецедентам

- $M$  — множество объектов,  $M = K_1 \sqcup \dots \sqcup K_l$
- $\{x_1, \dots, x_n\}$  — система целочисленных признаков
- $(x_1(S), \dots, x_n(S))$  — признаковое описание объекта  $S \in M$
- $T = \{S_1, \dots, S_m\}$  — обучающая выборка (прецеденты)
- $y_i$  — номер класса, которому принадлежит прецедент  $S_i$
- $A_T : M \rightarrow \{0, 1, \dots, l\}$  — алгоритм распознавания
- если  $A_T(S_i) = y_i, \forall S_i \in T$ , то  $A_T$  — *корректный алгоритм*

- 1 Голосование по тестам (*Дмитриев, Журавлёв, Кренделев, 1968*)
- 2 Голосование по представительным наборам (*Вайнцвайг, 1973*)
- 3 Голосование по корректным элементарным классификаторам (*Дюкова, Песков, 2002*).

## Определение

*Элементарным классификатором* (эл.кл.) ранга  $r$ ,  $1 \leq r \leq n$ , называется пара  $(H, \sigma)$ , где  $H = (x_{j_1}, \dots, x_{j_r})$  — набор признаков и  $\sigma = (\sigma^1, \dots, \sigma^r)$  — набор целых чисел.

- $H(S) = (x_{j_1}(S), \dots, x_{j_r}(S))$  — признаковое подписание  $S$
- если  $H(S) = \sigma$ , то эл.кл.  $(H, \sigma)$  *выделяет* объект  $S$

## Определение

Эл.кл.  $(H, \sigma)$  называется *корректным для класса  $K$* , если не существует двух прецедентов  $S_i \in K$  и  $S_t \notin K$ , выделяемых  $(H, \sigma)$ .

- Корректное распознавание на базе произвольных эл.кл. (Дюкова, Журавлёв, Рудаков, 1996).
- $U = ((H_1, \sigma_1), \dots, (H_d, \sigma_d))$  — набор эл.кл.  
 $U(S) = ([H_1(S) = \sigma_1], \dots, [H_d(S) = \sigma_d])$  — отклик

## Определение

- Набор эл.кл.  $U$  называется *корректным для класса  $K$* , если для любой пары прецедентов  $S_i \in K$  и  $S_t \notin K$  выполняется  $U(S_i) \neq U(S_t)$ .
- Булева функция  $F(t_1, \dots, t_d)$ , для которой  $F(U(S_i)) \neq F(U(S_t))$ ,  $S_i \in T \cap K$ ,  $S_t \in T \setminus K$ , называется *корректирующей*.
- *Корректный набор эл.кл.  $U$ , имеющий монотонную корректирующую функцию, называется **МОНОТОННЫМ**.*

# Алгоритмы голосования по корректным наборам эл.кл.

Первое экспериментальное исследование логических корректоров проведено в (Дюкова, Журавлёв, Сотнезов, 2010).

## Простейший логический корректор

- На этапе обучения для каждого класса  $K \in \{K_1, \dots, K_l\}$  строится семейство  $W_K$  корректных для  $K$  наборов эл.кл.
- При распознавании объекта  $S$  для каждого класса  $K \in \{K_1, \dots, K_l\}$  вычисляется оценка принадлежности  $S$  к  $K$

$$\Gamma(S, K) = \sum_{U \in W_K} \frac{1}{|W_K|} \sum_{S_i \in T \cap K} \frac{1}{|T \cap K|} [U(S_i) = U(S)] \quad (1)$$

В случае голосования по монотонным корректным наборам эл.кл. в формуле (1) вместо  $[U(S_i) = U(S)]$  берётся  $[U(S_i) \preceq U(S)]$ , где  $(a_1, \dots, a_n) \preceq (b_1, \dots, b_n) \Leftrightarrow a_j \leq b_j, \forall j \in \{1, \dots, n\}$ .

# Зачем усовершенствовать логические корректоры?

- Эксперименты показывают, что на «простых» задачах монотонные логические корректоры ошибаются чаще, чем алгоритмы голосования по представительным наборам, так как монотонный логический корректор эквивалентен «плохому» голосованию по представительным наборам.
- Этап обучения логических корректоров требует существенных временных затрат.
- Для расширения области применимости логических корректоров предлагается решить следующие задачи:
  - ❶ построить схему логического корректора общего вида;
  - ❷ построить логический корректор, отличный от ранее разработанных логических алгоритмов распознавания;
  - ❸ разработать методики повышения скорости и качества распознавания логических корректоров;

## Определение

- Предикат  $V : M \rightarrow \{0, 1\}$  называется **корректным для класса  $K$** , если множество прецедентов, на которых значение  $V$  равно 1, является подмножеством либо  $T \cap K$ , либо  $T \setminus K$ .
- **Корректный для класса  $K$  предикат** называется **представительным**, если существует прецедент из  $K$ , на котором значение  $V$  равно 1.

Конструирование корректных (представительных) предикатов:

- 1  $U = ((H_1, \sigma_1), \dots, (H_d, \sigma_d))$  — набор эл.кл.;
- 2  $P = (\rho_1, \dots, \rho_d)$  — набор бинарных отношений на  $\{0, 1\}$ , причём  $\rho_j \in \{[x \leq y], [x \geq y], [x \vee y], [\neg x \vee \neg y]\}$ ;
- 3  $G$  — непустой набор прецедентов класса  $K$ ;
- 4  $B_{(U,P,G)}(S) = \bigvee_{S_i \in G} \bigwedge_{j=1}^d \rho_j([H_j(S_i) = \sigma_j], [H_j(S) = \sigma_j])$ .

Алгоритм голосования по корректным предикатам вида  $B_{(U,P,G)}$

- На этапе обучения для каждого класса  $K$  строятся два семейства  $Z_K$  и  $Z_{\bar{K}}$  корректных для  $K$  предикатов вида  $B_{(U,P,G)}$ .
- При распознавании объекта  $S$  для каждого класса  $K$  вычисляется оценка принадлежности  $S$  к  $K$

$$\Gamma(S, K) = \sum_{B \in Z_K} \alpha_B B(S) - \sum_{B \in Z_{\bar{K}}} \alpha_B B(S).$$

Утверждение

Пусть предикат  $B_{(U,P,G)}$  корректен для  $K$ , набор  $U$  состоит из *различных* эл.кл. и  $G \subseteq T \cap K$ . Тогда  $U$  является представительным для  $K$  и имеет поляризуемую<sup>a</sup> корректирующую функцию.

<sup>a</sup>булева функция  $F(t_1, \dots, t_d)$  называется *поляризуемой*, если для некоторого  $(\alpha_1, \dots, \alpha_d)$  функция  $F(t_1 \oplus \alpha_1, \dots, t_d \oplus \alpha_d)$  монотонна.



# Задача этапа обучения логического корректора POLAR

Поиск корректных для  $K$  предикатов  $B_{(U,P,G)}$  сводится к поиску покрытий булевой матрицы

$$L_K = \left( \begin{array}{ccc} (H, \sigma, \rho) & S_j \in T \cap K & \\ \vdots & \vdots & \\ \neg \rho ([H(S_i) = \sigma], [H(S_t) = \sigma]) & [i = j] & \dots (S_i, S_t), \\ & & S_i \in T \cap K, \\ & & S_t \in T \setminus K \\ \vdots & \vdots & \end{array} \right)$$

Проблема вычислительного характера

Матрица  $L_K$ , как правило, имеет большой размер и много покрытий.

## Ранее используемая на этапе обучения методика

- Поиск покрытий осуществлялся генетическим алгоритмом. Для каждого класса  $K$  отбиралось не более  $N_{max}$  различных неприводимых покрытий матрицы  $L_K$ , соответствующих корректным наборам эл.кл. с распознающей способностью близкой к максимальной.
- Рассматривались только эл.кл. ранга 1 (уменьшение числа столбцов).
- Выборка случайно делилась на базовую и настроечную подвыборки. Матрица  $L_K$  строилась по базовой подвыборке (уменьшение числа строк). По настроечной подвыборке оценивалась распознающая способность корректных наборов эл.кл.

# Формирование семейств предикатов бустингом

- Семейства голосующих предикатов формируются итеративно по принципу *бустинга* (*Freund, Schapire, 1997*).
- Каждый добавляемый предикат соответствует покрытию подматрицы, составленной из сравнительно *небольшой части строк* матрицы  $L_K$  и не обязательно является корректным (подматрица меняется от итерации к итерации).
- Получена оценка числа итераций.

## Теорема

Пусть решается задача распознавания с  $l$  классами и обучающей выборкой размера  $m$ . Если бустинг-алгоритм формирования семейств предикатов выполняет более  $\frac{lm \ln m}{(\sqrt{m}-1)^2}$  итераций, то алгоритм голосования по построенным предикатам корректен.

- При поиске предиката рассматривается подматрица, состоящая из части столбцов матрицы  $L_K$ . Если эта подматрица не содержит нулевой строки, то соответствующие её столбцам тройки  $(H, \sigma, \rho)$  образуют *локальный базис*<sup>1</sup> для  $K$ .
- Использование локального базиса позволяет с одной стороны работать с эл.кл. произвольного ранга, а с другой стороны — отбросить малоинформативные эл.кл.
- Построены и экспериментально исследованы различные алгоритмы и стратегии формирования локальных базисов.

---

<sup>1</sup>термин введён в (*Воронцов, 1998*)

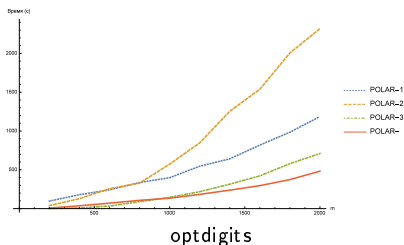
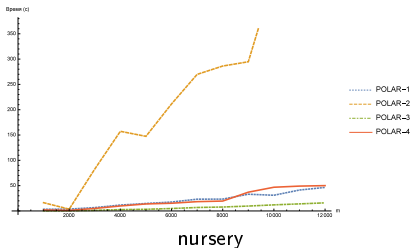
№	Название	$l$	$m$	$n$	$z$
1.	audiology	24	226	69	161
2.	balance scale	3	625	4	20
3.	breast cancer	2	699	9	90
4.	car	4	1728	6	21
5.	dermatology	4	366	34	192
6.	house votes	2	435	16	48
7.	kr vs kp	2	3196	36	73
8.	monks-2	2	601	6	17
9.	nursery	5	12960	8	27
10.	soybean large	19	307	35	132
11.	tic tac toe	2	958	9	27
12.	optdigits	10	5620	64	914
13.	letter recognition	26	20000	16	256
14.	lenses	3	24	4	9
15.	soybean small	4	47	35	72

# Результаты экспериментов

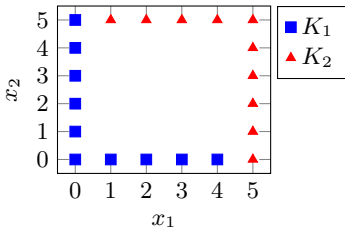
№	Задача	T*	ПН*	МОН*	P.-1	P.-2	P.-3	P.-4
1.	audiology	0.03	0.03	0.03	0.03	0.03	0.02	0.03
2.	b. scale	0.25	0.2	0.19	0.18	0.23	0.23	0.25
3.	b. cancer	0.046	0.04	0.057	0.061	0.059	0.065	0.059
4.	car	0.061	0.032	0.033	0.013	0.027	0.022	0.01
5.	dermat.	0.41	0.4	0.4	0.39	0.42	0.44	0.43
6.	h. votes	0.07	0.05	0.07	0.05	0.06	0.07	0.08
7.	kr-vs-kp	0.008	0.004	0.003	0.008	0.007	0.004	0.003
8.	monks-2	0.37	0.55	0.42	0.04	0.44	0.56	0.36
9.	nursery	0.027	0.003	0.005	0.002	—	0.002	0.004
10.	soybean l.	0.075	0.06	0.072	0.078	0.106	0.083	0.075
11.	tic-tac-toe	0.011	0.002	0.005	0.028	0.002	0.001	0.007
12.	letter r.	0.21	0.16	0.25	—	—	0.23	0.25
13.	optdigits	0.25	0.23	0.17	0.15	—	0.27	0.14
14.	lenses	0.42	0.25	0.29	0.33	0.29	0.38	0.25
15.	soybean s.	0	0	0	0	0.02	0.04	0

# Приложение А. Время обучения при различных стратегиях построения локальных базисов

- POLAR-1 — один локальный базис из одноранговых эл. кл. с отношениями  $[x \leq y]$  и  $[x \geq y]$ .
- POLAR-2 — один локальный базис строится бустингом над представительными наборами.
- POLAR-3 — локальный базис модифицируется на каждой итерации классическим голосованием по представительным наборам.
- POLAR-4 — локальный базис модифицируется на каждой итерации бустингом над представительными наборами.



# Приложение Б. Преимущество поляризуемой корректирующей функции



**Модельная задача  
распознавания с двумя  
классами и двумя  
признаками**

- Наименьшая мощность монотонного корректного для  $K_1$  набора равна 6, например,  $U_1 = ([x_2=0], [x_1=0], [x_1=1], [x_1=2], [x_1=3], [x_1=4])$ .  
Предикаты  $[U_1(S_i) \preceq U_1(S)], S_i \in T \cap K_1$ :  
 $[x_1=0], [x_1=1 \wedge x_2=0], \dots, [x_1=4 \wedge x_2=0]$ .
- Корректный набор, имеющий поляризуемую корректирующую функцию, можно составить из 3 эл.кл., например,  $U_2 = ([x_1 = 5], [x_2 = 0], [x_1 = 0])$ .  
При  $P = ([x \geq y], [x \leq y], [x \leq y])$  и  $S_i \in T \cap K_1$  предикат  $B_{(U_2, P, \{S_i\})}$  корректен для  $K_1$ . Каждый такой предикат имеет вид  $[x_1 \neq 5 \wedge x_2 = 0]$  или  $[x_1 = 0]$ .