

# Комбинаторные оценки переобучения пороговых классификаторов

Ишкина Шаура Хабировна

Вычислительный центр им. А. А. Дородницына РАН

ММО-17

19 – 25 сентября 2015, г. Светлогорск, Калининградская область

- 1 **Комбинаторная теория переобучения**
  - Задача оценивания вероятности переобучения
- 2 **Переобучение пороговых классификаторов**
  - Случай непрерывного и дискретного признаков
  - Алгоритм вычисления вероятности переобучения
- 3 **Эксперименты**
  - Проверка алгоритма
  - Задача медицинской диагностики
  - Выводы и открытые проблемы

## Основные понятия

$\mathbb{X} = \{x_1, \dots, x_L\}$  — конечное *генеральное множество* объектов;

$A = \{a_1, \dots, a_D\}$  — конечное множество *классификаторов*;

$I(a, x) = [\text{классификатор } a \text{ ошибается на объекте } x];$

$L \times D$ -матрица ошибок с попарно различными столбцами:

	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$\dots$	$a_D$	
$x_1$	1	1	0	0	0	1	$\dots$	1	$X$ — наблюдаемая (обучающая) выборка длины $l$
$\dots$	0	0	0	0	1	1	$\dots$	1	
$x_\ell$	0	0	1	0	0	0	$\dots$	0	
$x_{\ell+1}$	0	0	0	1	1	1	$\dots$	0	$\bar{X}$ — скрытая (контрольная) выборка длины $k = L - l$
$\dots$	0	0	0	1	0	0	$\dots$	1	
$x_L$	0	1	1	1	1	1	$\dots$	0	

### Вероятностная аксиома:

Пусть все разбиения  $X \sqcup \bar{X} = \mathbb{X}$  равновероятны.

## Задача оценивания вероятности переобучения

$n(a, X) = \sum_{x \in X} I(a, x)$  — число ошибок  $a \in A$  на выборке  $X \subset \mathbb{X}$ ;

$\nu(a, X) = \frac{1}{|X|} n(a, X)$  — частота ошибок  $a$  на выборке  $X$ ;

**Опр.** Метод обучения  $\mu: 2^{\mathbb{X}} \rightarrow A$  произвольной выборке  $X \subset \mathbb{X}$  ставит в соответствие некоторый классификатор  $a \in A$ .

**Опр.** Метод минимизации эмпирического риска:

$$\mu X \in A(X) = \underset{a \in A}{\text{Arg min}} n(a, X).$$

$\delta(\mu, X) = \nu(\mu X, \bar{X}) - \nu(\mu X, X)$  — переобученность  $\mu$  на  $X$ .

**Вероятность переобучения:**

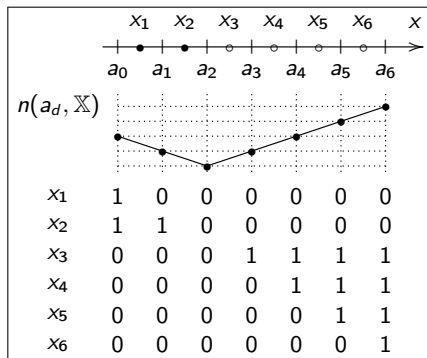
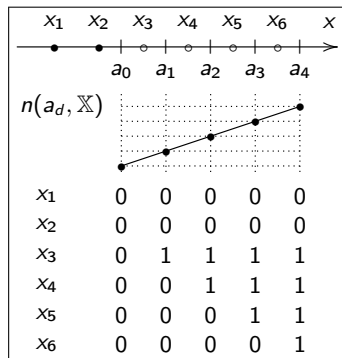
$$Q_\varepsilon(\mu, \mathbb{X}) = P[\delta(\mu, X) \geq \varepsilon] = \frac{1}{C_L^{\mathbb{X}}} \sum_{X \subset \mathbb{X}} [\delta(\mu, X) \geq \varepsilon].$$

## Пороговые классификаторы. Пример

Пусть  $x = (x^1, \dots, x^n) \in \mathbb{X}$ ,  $\mathbb{Y} = \{0, 1\}$ . Пусть  $x^1 \in \mathbb{R}$ .

Рассмотрим семейство

$$A = \{a_\theta(x) = [x^1 \geq \theta], \theta \in \mathbb{R}\}.$$

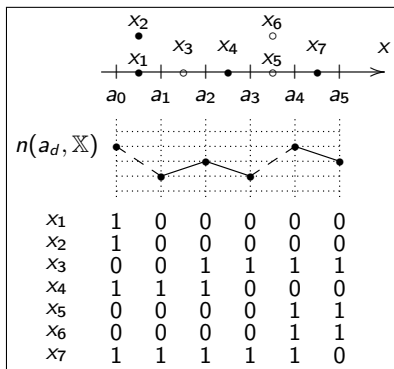


## Прямая цепь

Прямая цепь  $A = \{a_0, \dots, a_D\}$  – семейство классификаторов:

$$\mathbb{G}_d = \{x \in \mathbb{X} \mid I(a_{d-1}, x) \neq I(a_d, x)\}, \quad d = 1, \dots, D.$$

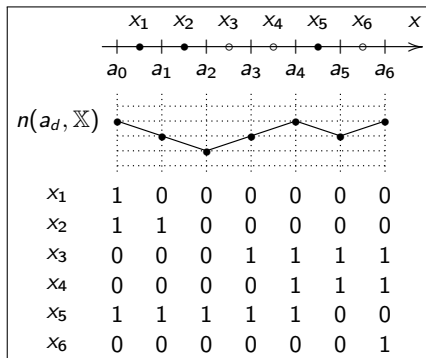
Обозначим  $\mathbb{G} = \mathbb{G}_1 \sqcup \dots \sqcup \mathbb{G}_D = \{x \in \mathbb{X} \mid I(a_0, x) \neq I(a_D, x)\}$



## Непрерывная прямая цепь

Непрерывная прямая цепь  $A = \{a_0, \dots, a_D\}$  – прямая цепь,

$$\forall d \in \{1, \dots, D\} \quad |\mathbb{G}_d| = 1.$$



## Алгоритм вычисления вероятности переобучения

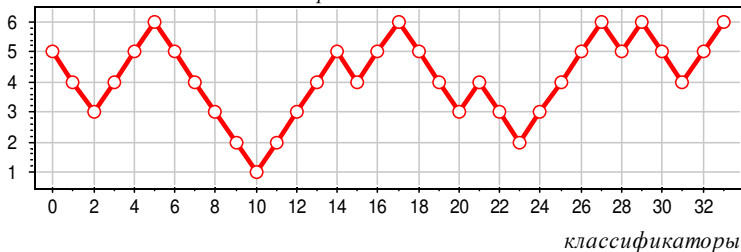
Результат: алгоритм вычисления вероятности переобучения со сложностью

$$O(G^6 + Gl^3 + L) = O(L^6),$$

где  $G = |\mathbb{G}| = \rho(a_0, a_D)$ .

Для непрерывной цепи параметр  $G$  равен длине  $D$ .

число ошибок на полной выборке





## Схема работы алгоритма

$$Q_\varepsilon = P[\delta(\mu X, X) \geq \varepsilon] = \sum_{d=0}^D P[\mu X = a_d][\delta(a_d, X) \geq \varepsilon]$$

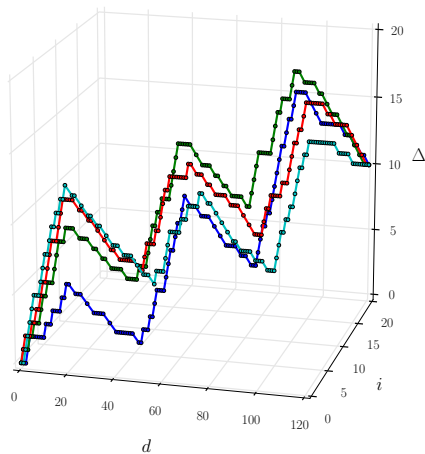
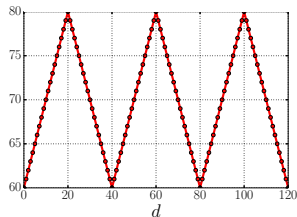
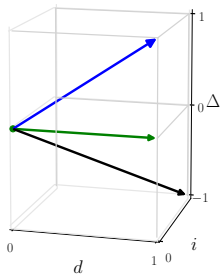
$[\mu X = a_d]$  зависит только от выбора разбиения  $\{X \cap \mathbb{G}, \bar{X} \cap \mathbb{G}\}$ .  
 Зададим параметры  $t = |X \cap \mathbb{G}|$ ,  $e = n(a_d, X \cap \mathbb{G})$ ,

$$Q_\varepsilon = \sum_d \sum_t \sum_e G_d(t, e) N_d(t, e, \varepsilon),$$

где  $G_d(t, e)$  – количество разбиений множества  $\mathbb{G}$ , таких, что  $[\mu X = a_d]$ , множители  $N_d(t, e, \varepsilon)$  вычисляются аналитически.

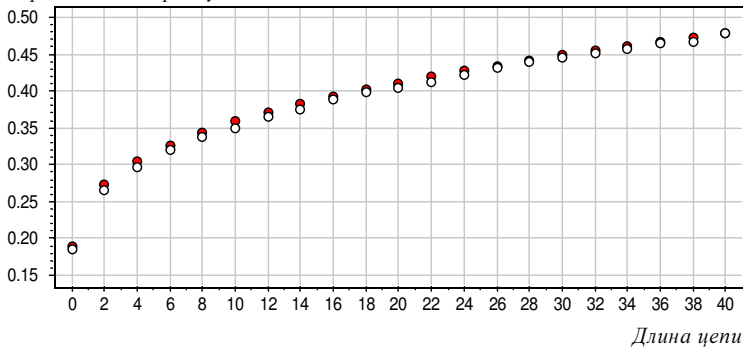
## Разбиение как блуждание по трехмерной сетке

■  $x \in X, I(a_d, x) = 0$    
 ■  $x \in X, I(a_d, x) = 1$    
 ■  $x \in \bar{X}$



## Проверка алгоритма вычисления вероятности переобучения

Вероятность переобучения



● точная оценка    ○ метод Монте-Карло

**Условия эксперимента:** Случайная цепь;

$L = 200$ ,  $\ell = 100$ ,  $m = 40$ ,  $\varepsilon = 0.05$ ; метод Монте-Карло по  $10^5$  случайных разбиений.

## Постановка задачи медицинской диагностики

**Дано:** Выборка  $\mathbb{X}$  пациентов  $S$  – символьных последовательностей в 6-буквенном алфавите. Метки классов: 1 – болен, 0 – здоров.

**Критерий качества:** AUC (Area Under Curve) на контрольной выборке, вычисляемый  $10 \times 10$ -блочной кросс-валидацией.

Признаки – частоты  $p_w$  триграмм  $w$ , т.е.  $6^3 = 216$  признаков.

*Наивный байесовский классификатор (NB)*

$$a(S) = \left[ \ln \frac{\pi_1(S)}{\pi_0(S)} \geq \beta \right] = \left[ \sum_w \gamma_w p_w(S) \geq \beta \right],$$

Для некоторых болезней дает AUC 95 – 99%.

## Отбор признаков

**Старая схема:** top- $K$  по информационному критерию,  $K$  – параметр. Веса  $\gamma_w$  вычисляются по аналитическим формулам.

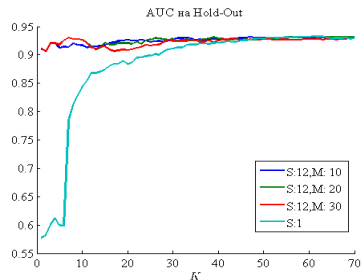
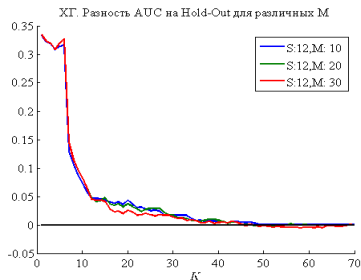
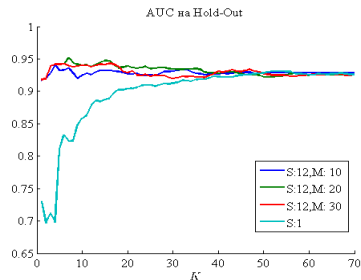
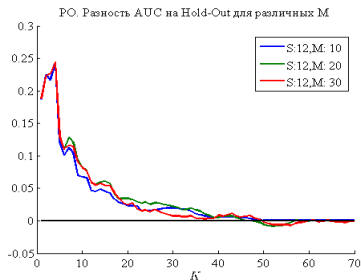
**Новая схема:** Жадный отбор из более широкого множества кандидатов, top- $(K^* + M)$  по информационному критерию.  $K^*$  – оптимальное значение, полученное по старой схеме.

**для всех**  $j = 1, \dots, K^* + M$

Из оставшихся  $K^* + M - j + 1$  кандидатов выбирается тот, который максимально улучшает эмпирический риск на контрольной выборке

Веса  $\gamma_w$  вычисляются по аналитическим формулам.

## Сравнение методов отбора



## Выводы

### Основные результаты

- Разработан алгоритм вычисления вероятности переобучения произвольной прямой цепи, полиномиальный по длине выборки
- Проведены вычислительные эксперименты с применением на практике подхода, основанного на вычислении переобученности

### Открытые проблемы

- Применение в логических алгоритмах классификации
- Сравнение с существующими оценками переобучения в задаче отбора признаков
- Оптимизация весов  $\gamma_w$  признаков с помощью оценок переобучения