

# Комбинаторная теория переобучения и поиск логических закономерностей

Воронцов Константин Вячеславович

Вычислительный центр им. А. А. Дородницына РАН

15-я всероссийская научная конференция  
«Математические методы распознавания образов»  
11–17 сентября 2011 • г. Петрозаводск

## Содержание

- 1 **Комбинаторная теория переобучения**
  - Задача оценивания вероятности переобучения
  - Граф расслоения–связности
  - Оценки расслоения–связности
- 2 **Логические алгоритмы классификации**
  - Композиции интерпретируемых информативных правил
  - Критерий предсказанной информативности
  - Эксперименты и выводы

## Задача обучения по прецедентам. Матрица ошибок.

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

$\mathbb{A} = \{a_1, \dots, a_D\}$  — конечное семейство алгоритмов;

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

$L \times D$ -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	
$x_l$	0	0	1	0	0	0	$\dots$	0	
$x_{l+1}$	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	

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

$\nu(a, X) = m(a, X)/|X|$  — частота ошибок  $a$  на выборке  $X$ ;

## Задача оценивания вероятности переобучения

**Опр.** Метод обучения  $\mu: X \mapsto a$  произвольной выборке  $X \subset \mathbb{X}$  ставит в соответствие некоторый алгоритм  $a \in \mathbb{A}$ .

**Пример.** Метод минимизации эмпирического риска:

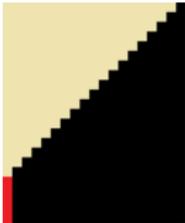
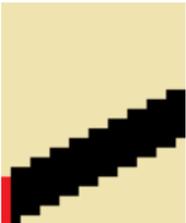
$$\mu X = \arg \min_{a \in \mathbb{A}} \nu(a, X).$$

**Опр.** Вероятность переобучения — доля разбиений  $X \sqcup \bar{X} = \mathbb{X}$ , при которых частота ошибок на контроле  $\bar{X}$  существенно (на  $\varepsilon$ ) превышает частоту ошибок на обучении  $X$ :

$$Q_\varepsilon(\mu, \mathbb{X}) = \mathbf{P}[\nu(\mu X, \bar{X}) - \nu(\mu X, X) \geq \varepsilon].$$

Здесь и далее  $\mathbf{P} \equiv \frac{1}{C_L^\ell} \sum_{X \subset \mathbb{X}: |X|=\ell}$

## Эксперимент: четыре семейства алгоритмов, заданных матрицами ошибок; лучший алгоритм у всех одинаков

	с расслоением по числу ошибок	без расслоения по числу ошибок
каждая пара соседних алгоритмов отличается только на одном объекте (образуется <i>цепь</i> )		
соседние алгоритмы существенно различны, ( <i>цепь</i> не образуется)		

Постепенно добавляя алгоритмы в  $\{a_1, \dots, a_D\}$ , построим зависимости вероятности переобучения  $Q_\epsilon$  от числа  $D$ .

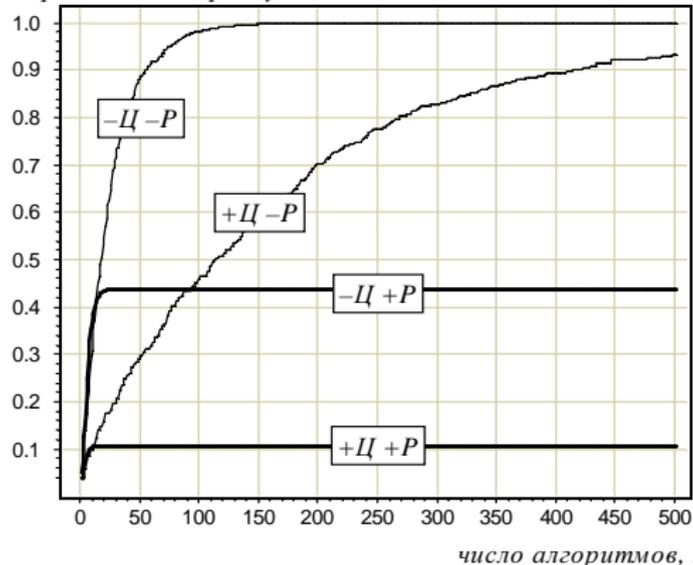
## Вывод: расслоение и связность снижают переобучение

### Эксперимент

с четырьмя  
семействами  
алгоритмов,  
 $\ell = k = 100$ ,  $\varepsilon = 0.05$ :

- +Ц — цепь;
- Ц — не цепь;
- +P — с расслоением;
- P — без расслоения;

Вероятность переобучения



**Вывод:** получение точных оценок вероятности переобучения невозможно без учёта эффектов расслоения и связности.

## Граф расслоения–связности множества алгоритмов

Определим бинарные отношения на множестве алгоритмов  $\mathbb{A}$ :  
частичный порядок  $a \leq b$ :  $I(a, x) \leq I(b, x)$  для всех  $x \in \mathbb{X}$ ;  
предшествование  $a \prec b$ :  $a \leq b$  и  $\|b - a\| = 1$ .

Опр. Граф расслоения–связности  $\langle \mathbb{A}, E \rangle$ :

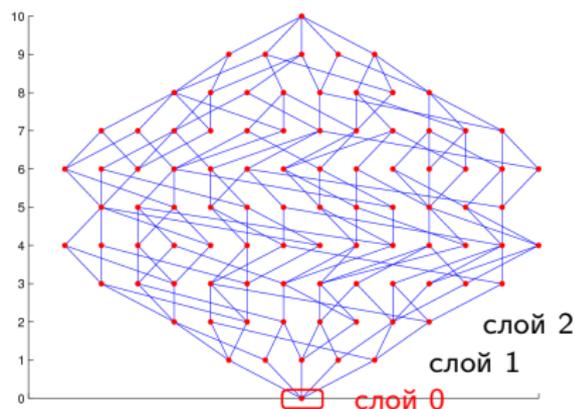
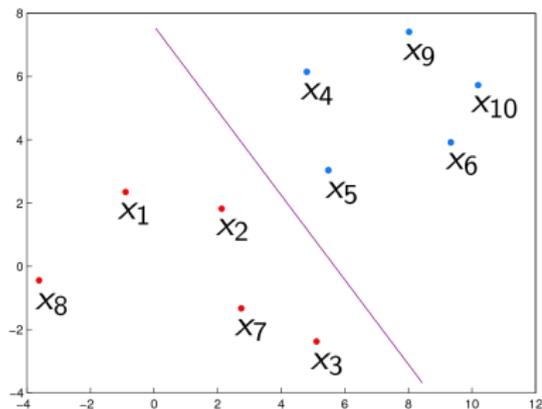
$\mathbb{A}$  — множество попарно различных векторов ошибок;

$E = \{(a, b) : a \prec b\}$ .

Свойства графа расслоения–связности:

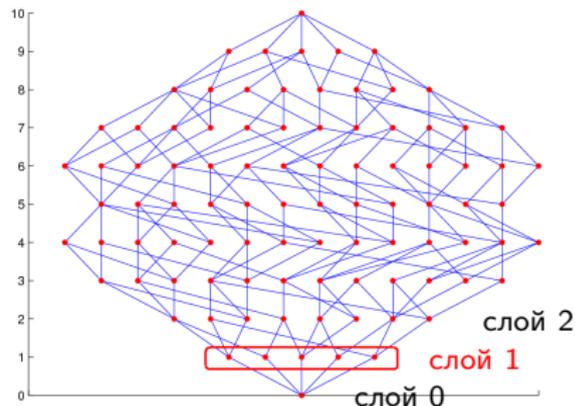
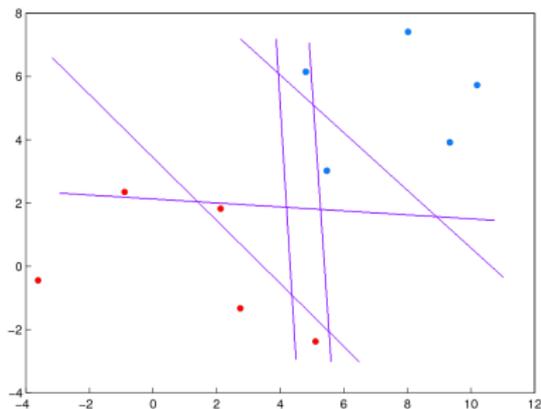
- это подграф графа Хассе отношения порядка  $\leq$  на  $\mathbb{A}$ ;
- каждому ребру  $(a, b)$  соответствует объект  $x_{ab} \in \mathbb{X}$ , такой, что  $I(a, x_{ab}) = 0$ ,  $I(b, x_{ab}) = 1$ ;
- граф является многодольным со слоями  
 $A_m = \{a \in \mathbb{A} : m(a, \mathbb{X}) = m\}$ ,  $m = 0, \dots, L + 1$ ;

## Пример. Семейство линейных алгоритмов классификации



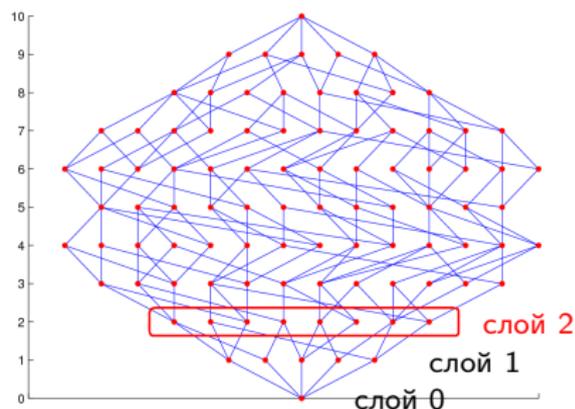
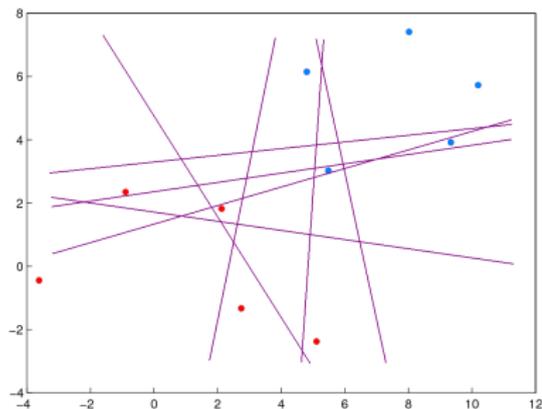
	слой 0
$x_1$	0
$x_2$	0
$x_3$	0
$x_4$	0
$x_5$	0
$x_6$	0
$x_7$	0
$x_8$	0
$x_9$	0
$x_{10}$	0

## Пример. Семейство линейных алгоритмов классификации



	слой 0	слой 1				
X <sub>1</sub>	0	1	0	0	0	0
X <sub>2</sub>	0	0	1	0	0	0
X <sub>3</sub>	0	0	0	1	0	0
X <sub>4</sub>	0	0	0	0	1	0
X <sub>5</sub>	0	0	0	0	0	1
X <sub>6</sub>	0	0	0	0	0	0
X <sub>7</sub>	0	0	0	0	0	0
X <sub>8</sub>	0	0	0	0	0	0
X <sub>9</sub>	0	0	0	0	0	0
X <sub>10</sub>	0	0	0	0	0	0

## Пример. Семейство линейных алгоритмов классификации



	слой 0	слой 1						слой 2								
X <sub>1</sub>	0	1	0	0	0	0	0	1	0	0	0	0	1	1	0	...
X <sub>2</sub>	0	0	1	0	0	0	0	1	1	0	0	0	0	0	0	...
X <sub>3</sub>	0	0	0	1	0	0	0	0	1	1	0	0	0	0	1	...
X <sub>4</sub>	0	0	0	0	1	0	0	0	0	1	1	0	0	0	0	...
X <sub>5</sub>	0	0	0	0	0	0	1	0	0	0	1	1	1	0	0	...
X <sub>6</sub>	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	...
X <sub>7</sub>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	...
X <sub>8</sub>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
X <sub>9</sub>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
X <sub>10</sub>	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...

## Характеристики **расслоения** и **связности** алгоритма

Множество объектов, соответствующих рёбрам, исходящим из  $a$ :

$$X_a = \{x_{ab} \in \mathbb{X} \mid a \prec b\},$$

Множество объектов, соответствующих всем рёбрам на путях, ведущих в  $a$ :

$$X'_a = \{x \in \mathbb{X} \mid \exists b \in \mathbb{A}: b \prec a, l(b, x) < l(a, x)\}.$$

**Опр.** Характеристики **расслоения** и **связности** алгоритма  $a$ :

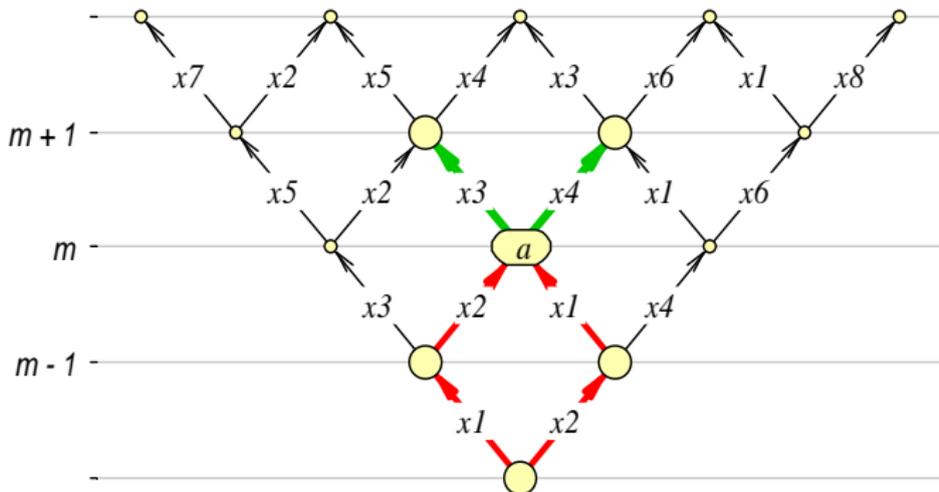
$u(a) = |X_a|$  — **верхняя связность** алгоритма  $a$ ;

$q(a) = |X'_a|$  — **неполноценность** алгоритма  $a$ ;

## Пример: двумерная сеть алгоритмов

Верхняя связность:  $u(a) = \#\{x3, x4\} = 2$ ;

Неполноценность:  $q(a) = \#\{x1, x2\} = 2$ ;



## Монотонные методы обучения

**Опр.** Метод обучения  $\mu$  называется *монотонным*, если

$$\mu X = \arg \min_{a \in \mathbb{A}} K(a, X),$$

где  $K(a, X)$  — строго монотонная функция вектора ошибок  $a$ :

$$\forall X \subset \mathbb{X}, \forall a, b \in \mathbb{A} \quad a < b \rightarrow K(a, X) < K(b, X).$$

**Пример 1.** Взвешенный эмпирический риск с весами  $w_i$ :

$$K(a, X) = \sum_{x_i \in X} w_i l(a, x_i), \quad w_i > 0.$$

**Пример 2.** Стандартные критерии информативности логических закономерностей (рассматриваются далее).

## Верхняя оценка вероятности переобучения

### Теорема (оценка расслоения–связности)

Пусть  $\mu$  — монотонный метод обучения.

Тогда для любых  $\mathbb{X}$ ,  $\mathbb{A}$  и  $\varepsilon \in (0, 1)$

$$Q_\varepsilon \leq \sum_{a \in \mathbb{A}} \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell} \mathcal{H}_{L-u-q}^{\ell-u, m-q} \left( \frac{\ell}{L} (m - \varepsilon k) \right)$$

где  $u = |X_a|$  — верхняя связность алгоритма  $a$ ,

$q = |X'_a|$  — неполноценность алгоритма  $a$ ,

$m = m(a, X)$  — число ошибок алгоритма  $a$ ,

$$\mathcal{H}_L^{\ell, m}(z) = \sum_{s=0}^{\lfloor z \rfloor} \frac{C_m^s C_{L-m}^{\ell-s}}{C_L^\ell}, \quad z = 0, \dots, \ell$$

— функция гипергеометрического распределения:

## Свойства оценки

$$Q_\varepsilon \leq \sum_{a \in \mathbb{A}} \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell} \mathcal{H}_{L-u-q}^{\ell-u, m-q} \left( \frac{\ell}{L} (m - \varepsilon k) \right)$$

- 1 Вклад алгоритма  $a \in \mathbb{A}$  убывает экспоненциально по  $u(a) \Rightarrow$  **связные семейства меньше переобучаются**; по  $q(a) \Rightarrow$  **только нижние слои вносят вклад в  $Q_\varepsilon$** .
- 2 Оценка обращается в равенство в случае многомерных монотонных сетей алгоритмов.
- 3 Вероятность получить алгоритм в результате обучения

$$P[\mu X = a] \leq P_a = \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell}.$$

- 4 Если  $q(a) > k$ , то  $P_a = 0$  и вклад алгоритма  $a$  равен 0  $\Rightarrow$  при малых  $k$  оценка вырождается.
- 5  $\sum_{a \in \mathbb{A}} P_a$  — оценка степени завышенности.

## Свойства оценки

$$Q_\varepsilon \leq \sum_{a \in \mathbb{A}} \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell} \mathcal{H}_{L-u-q}^{\ell-u, m-q} \left( \frac{\ell}{L} (m - \varepsilon k) \right)$$

- 6 При  $|\mathbb{A}| = 1$  это оценка скорости сходимости частот в двух выборках (вариант закона больших чисел):

$$Q_\varepsilon = \mathcal{H}_L^{\ell, m} \left( \frac{\ell}{L} (m - \varepsilon k) \right) \rightarrow 0 \text{ при } \ell, k \rightarrow \infty.$$

- 7 При  $q = u = 0$  и  $\ell = k$  это оценка Вапника-Червоненкиса:

$$Q_\varepsilon \leq \sum_{a \in \mathbb{A}} \mathcal{H}_L^{\ell, m} \left( \frac{\ell}{L} (m - \varepsilon k) \right) \leq |\mathbb{A}| \cdot \exp(-\varepsilon \ell^2).$$

- 8 При замене неполноценности  $q$  на нижнюю связность  $d$  это верхняя оценка функционала равномерной сходимости

$$P \left[ \sup_{a \in \mathbb{A}} (\nu(a, \bar{X}) - \nu(a, X)) \geq \varepsilon \right],$$

которая учитывает связность, но не учитывает расслоение.

## Взвешенное голосование пороговых конъюнкций

Взвешенное голосование правил:

$$a(x) = \arg \max_{y \in Y} \sum_{r \in R_y} w_r r(x),$$

где  $r: X \rightarrow \{0, 1\}$  — правило,  $w_r$  — вес правила  $r$ ,  
 $R_y$  — множество правил класса  $y$ ,  $Y$  — множество классов,

Будем выбирать  $R_y$  из семейства пороговых конъюнкций  $R$ :

$$r(x) = \bigwedge_{j \in J} [f_j(x) \leq \theta_j],$$

где  $f_j(x)$  — числовые признаки,  $\theta_j$  — пороги,  $j = 1, \dots, n$ ;  
 $J \subseteq \{1, \dots, n\}$  — подмножество признаков, обычно  $|J| \lesssim 7$ .

## Критерии информативности правил

Задача поиска закономерностей класса  $y$  в семействе правил  $R$ :

$$\begin{cases} n(r, X) = \#\{x_i \in X \mid r(x_i) = 1, y_i \neq y\} \rightarrow \min_{r \in R}; \\ p(r, X) = \#\{x_i \in X \mid r(x_i) = 1, y_i = y\} \rightarrow \max_{r \in R}. \end{cases}$$

Задача максимизации информативности:  $\mathcal{H}(p, n) \rightarrow \max_{r \in R}$ .

Введём индикатор ошибки для правил класса  $y$ :

$$I(r, x_i) = [r(x_i) \neq [y_i=y]], \quad x_i \in \mathbb{X}.$$

Он индуцирует векторы ошибок и отношения  $\leq$ ,  $\prec$  на  $R$ .

### Теорема

Если  $\mathcal{H}(p, n)$  строго возрастает по  $p$  и строго убывает по  $n$ , то максимизация информативности  $\mathcal{H}(p(r, X), n(r, X))$  является монотонным методом обучения правил.

## Примеры критериев информативности

$P = |\{x_i: y_i = y\}|$  — число «своих» в обучающей выборке;

$N = |\{x_i: y_i \neq y\}|$  — число «чужих» в обучающей выборке.

- точность:  $\mathcal{H}(p, n) = \frac{p + N - n}{P + N}$ ;

- взвешенная точность:  $\mathcal{H}(p, n) = p - \lambda n$ ;

- энтропийный критерий информационного выигрыша:

$$\mathcal{H}(p, n) = h\left(\frac{P}{\ell}\right) - \frac{p+n}{\ell} h\left(\frac{p}{p+n}\right) - \frac{\ell-p-n}{\ell} h\left(\frac{P-p}{\ell-p-n}\right),$$

где  $h(q) = -q \log_2 q - (1 - q) \log_2(1 - q)$ ;

- индекс Джини (Gini impurity): то же, при  $h(q) = 2q(1 - q)$ ;

- точный тест Фишера:  $\mathcal{H}(p, n) = -\log C_P^p C_N^n / C_{P+N}^{p+n}$ ;

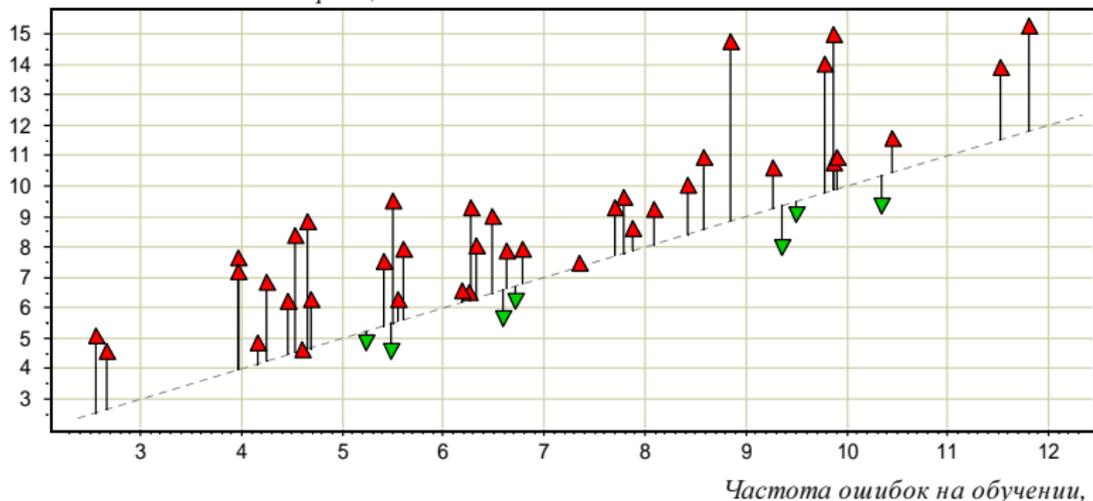
- критерий бустинга:  $\mathcal{H}(p, n) = \sqrt{p} - \sqrt{n}$ ;

## Проблема переобучения закономерностей

Общий недостаток всех критериев информативности — они не учитывают переобучение порогов  $\theta_j$  и наборов  $J$ .

**Пример.** Задача предсказания результата хирургического лечения атеросклероза,  $L = 98$ , точки — закономерности.

*Частота ошибок на контроле, %*



## Идея уменьшения переобучения закономерностей

### План:

- 1 Вывести оценки расслоения–связности для метода обучения порогов  $\theta$  при фиксированном наборе признаков  $J$ , отдельно для:
  - числа ошибок I рода  $m' = P - p$  (на объектах класса  $y$ )
  - числа ошибок II рода  $m'' = n$  (на остальных классах).
- 2 С их помощью получить оценки  $\hat{p}$ ,  $\hat{n}$  на неизвестной контрольной выборке по известным  $p$ ,  $n$  на обучающей выборке.
- 3 Подставить эти оценки в критерий  $\mathcal{H}$  и получить критерий предсказанной информативности  $\mathcal{H}(\hat{p}, \hat{n})$ .
- 4 По этому критерию отбирать подмножества признаков  $J$ .

## Ошибки I и II рода относительно класса $y$

Два определения индикатора ошибки, числа и частоты ошибок:

$$\text{I рода: } l'(r, x_i) = [r(x_i) = 0] [y_i = y], \quad m' = P - p, \quad \nu';$$

$$\text{II рода: } l''(r, x_i) = [r(x_i) = 1] [y_i \neq y], \quad m'' = n, \quad \nu''.$$

Оценки вероятности переобучения для  $r = \mu X$ :

$$Q'_\varepsilon = P[\nu'(r, \bar{X}) - \nu'(r, X) \geq \varepsilon] \leq \eta'_{\ell, k}(\varepsilon);$$

$$Q''_\varepsilon = P[\nu''(r, \bar{X}) - \nu''(r, X) \geq \varepsilon] \leq \eta''_{\ell, k}(\varepsilon).$$

Обращение оценок: с вероятностью не менее  $(1 - \eta)$

$$\nu'(r, \bar{X}) \leq \nu'(r, X) + \varepsilon'_{\ell, k}(\eta),$$

$$\nu''(r, \bar{X}) \leq \nu''(r, X) + \varepsilon''_{\ell, k}(\eta).$$

Критерий предсказанной информативности:

$$\widehat{\mathcal{H}}(p, n) = \mathcal{H}\left(\frac{k}{\ell}p - k\varepsilon'_{\ell, k}\left(\frac{1}{2}\right), \frac{k}{\ell}n + k\varepsilon''_{\ell, k}\left(\frac{1}{2}\right)\right).$$

## Оценка расслоения–связности для ошибок I и II рода

### Теорема

Для монотонного метода обучения  $\mu$  и любых  $\mathbb{X}, R, \varepsilon \in (0, 1)$

$$Q'_\varepsilon \leq \sum_{r \in R} \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell} \mathcal{H}_{L-u-q}^{\ell-u, m'-q'} \left( \frac{\ell}{L} (m' - \varepsilon k) \right);$$

$$Q''_\varepsilon \leq \sum_{r \in R} \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell} \mathcal{H}_{L-u-q}^{\ell-u, m''-q''} \left( \frac{\ell}{L} (m'' - \varepsilon k) \right);$$

где  $u = |X_r|$  — верхняя связность правила  $r$ ,

$q = |X'_r|$ ,  $q' = |X'_r \cap \mathbb{X}_y|$ ,  $q'' = |X'_r \cap \mathbb{X}_{\bar{y}}|$  — неполноценность правила  $r$  относительно индикаторов ошибки  $I, I', I''$ ,

$m' = m'(r, \mathbb{X})$ ,  $m'' = m''(r, \mathbb{X})$  — число ошибок правила  $r$  относительно индикаторов ошибки  $I', I''$ .

## Наблюдаемые и ненаблюдаемые оценки

Мы построили критерий предсказанной информативности:

$$\widehat{\mathcal{H}} = \mathcal{H} \left( \frac{k}{\ell} p(r, X) - k\varepsilon'_{\ell, k} \left( \frac{1}{2} \right), \frac{k}{\ell} n(r, X) + k\varepsilon''_{\ell, k} \left( \frac{1}{2} \right) \right).$$

Но нам нужен критерий, предсказывающий информативность на скрытой выборке  $\bar{X}$  длины  $K$

$$\widehat{\mathcal{H}}_{\text{ненабл}} = \mathcal{H} \left( \frac{K}{L} p(r, \bar{X}) - K\varepsilon'_{L, K} \left( \frac{1}{2} \right), \frac{K}{L} n(r, \bar{X}) + K\varepsilon''_{L, K} \left( \frac{1}{2} \right) \right).$$

Однако мы не можем вычислить  $\varepsilon'_{L, K}$  и  $\varepsilon''_{L, K}$ , т.к.  $\bar{X}$  скрыта.

**Эвристика:** заменив  $\varepsilon'_{L, K}$  на  $\varepsilon'_{\ell, k}$  и  $\varepsilon''_{L, K}$  на  $\varepsilon''_{\ell, k}$ , получим наблюдаемый критерий предсказанной информативности:

$$\widehat{\mathcal{H}}_{\text{набл}} = \mathcal{H} \left( \frac{K}{L} p(r, \bar{X}) - K\varepsilon'_{\ell, k} \left( \frac{1}{2} \right), \frac{K}{L} n(r, \bar{X}) + K\varepsilon''_{\ell, k} \left( \frac{1}{2} \right) \right).$$

Почему мы имеем право на такую замену?

## Эксперимент 1: зависимость $\widehat{\mathcal{H}}_{\text{ненабл}}$ от $\widehat{\mathcal{H}}_{\text{набл}}$

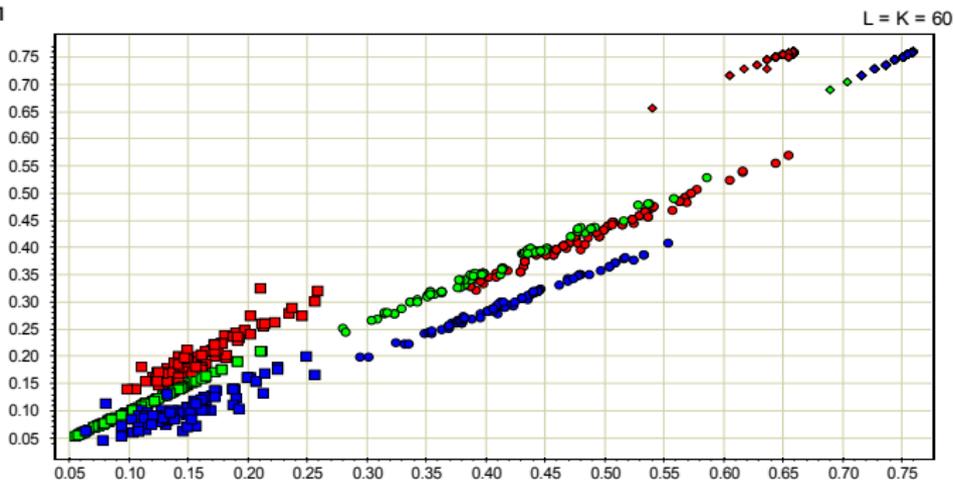
наблюдаемая выборка  $\mathbb{X}$

скрытая выборка  $\bar{\mathbb{X}}$

наблюдаемые  $\varepsilon'_{\ell,k}, \varepsilon''_{\ell,k}$

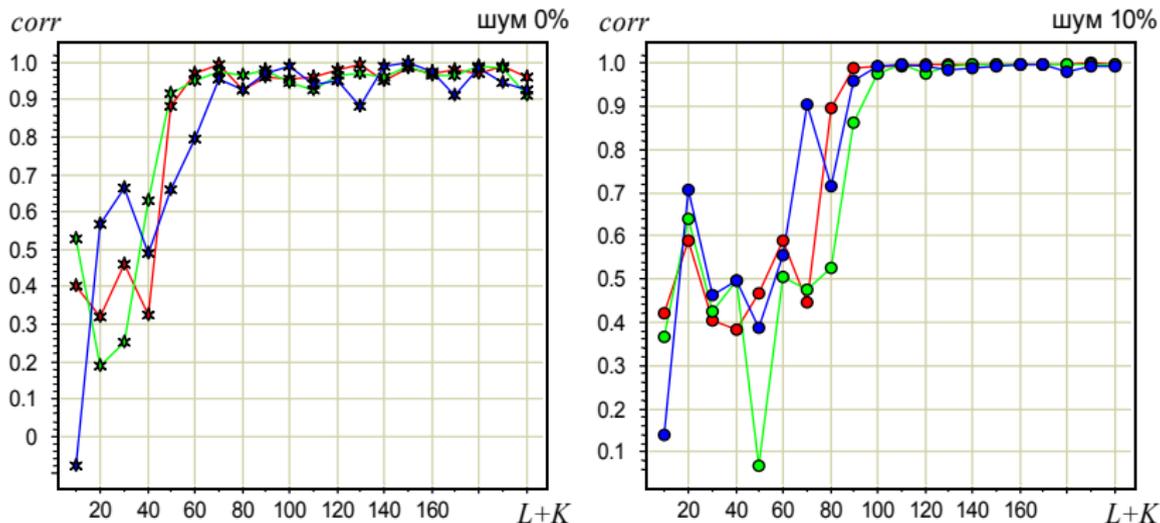
ненаблюдаемые  $\varepsilon'_{L,K}, \varepsilon''_{L,K}$

$\widehat{\mathcal{H}}_{\text{ненабл}}$



$\widehat{\mathcal{H}}_{\text{набл}}$

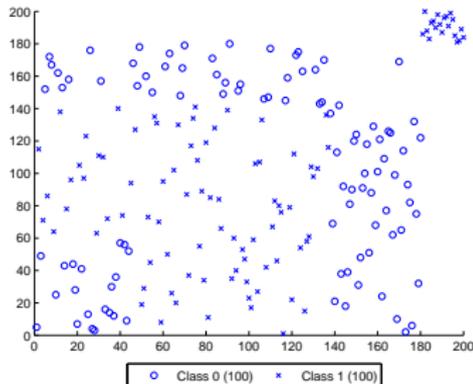
## Эксперимент 1: зависимость корреляции $\hat{\mathcal{H}}_{\text{ненабл}}$ и $\hat{\mathcal{H}}_{\text{набл}}$ от длины выборки и зашумлённости данных



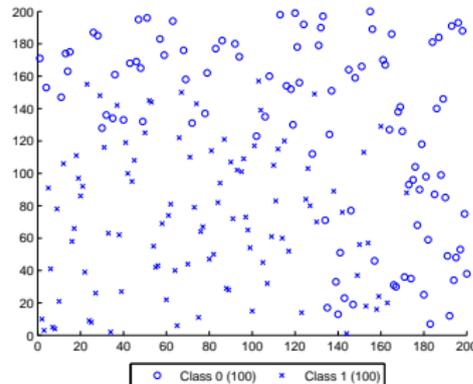
Экспериментальные результаты получены М. Дударенко на модельных данных с энтропийным критерием  $\mathcal{H}$ .

## Эксперимент 2: в каких случаях предсказанная информативность даёт выигрыш?

выборка 1:  $L = 200$ ,  
 шум вдали от границы классов



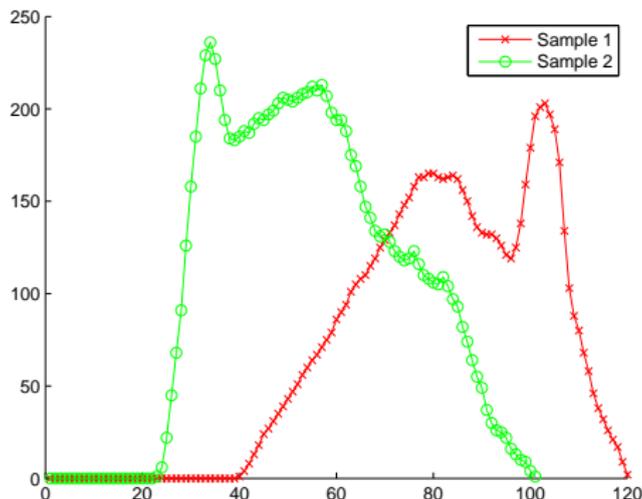
выборка 2:  $L = 200$ ,  
 шум вдоль границы классов



	выборка 1	выборка 2
ошибок лучшего алгоритма, %	20,0	11,5
$\mathcal{H}$ обучение	40,3	40,5
$\mathcal{H}$ контроль	39,1	34,5
наблюдаемый $\widehat{\mathcal{H}}$	34,9	31,2

## Эксперимент 2: Профили расслоения

Профили расслоения  $D(m) = |A_m|$  выборок 1 и 2.



Выборка 2 содержит почти вдвое лучший алгоритм, однако найти его по случайной подвыборке практически невозможно из-за мощных нижних слоёв и переобучения.

## Эксперимент 2: выводы

Выводы:

- Избавиться от переобучения порогов нельзя.
- Но можно предсказать, когда (для каких наборов признаков) оно будет меньше, проанализировав структуру графа расслоения–связности:
- когда профиль расслоение растёт медленнее,
- то есть когда граница классов наименее зашумлена.

## Эксперимент 3: на реальных данных

Задачи из репозитория UCI:

задачи	объектов	признаков
australian	690	14
echo cardiogram	74	10
heart disease	294	13
hepatitis	155	19
labor relations	40	16
liver	345	6

Методы обучения:

- WV — weighted voting (boosting);
- DL — decision list;
- LR — logistic regression.

Методика тестирования: 10-кратный скользящий контроль.

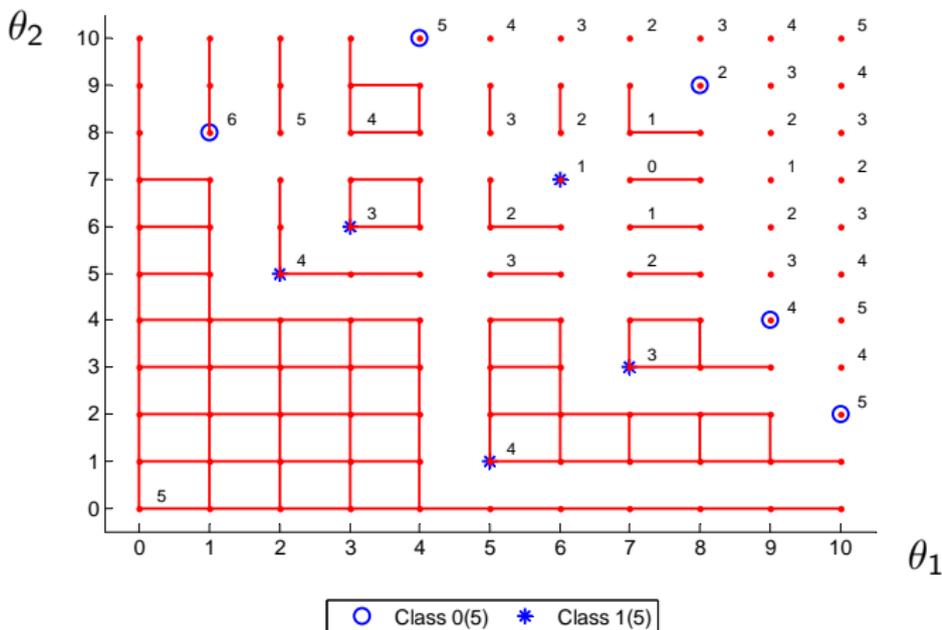
## Результаты эксперимента на реальных данных

методы	задачи					
	austr	echo	heart	hepa	labor	liver
RIPPER-opt	15.5	2.97	19.7	20.7	18.0	32.7
RIPPER+opt	15.2	5.53	20.1	23.2	18.0	31.3
C4.5(Tree)	14.2	5.51	20.8	18.8	14.7	37.7
C4.5(Rules)	15.5	6.87	20.0	18.8	14.7	37.5
C5.0	14.0	4.30	21.8	20.1	18.4	31.9
SLIPPER	15.7	4.34	19.4	17.4	12.3	32.2
LR	14.8	4.30	19.9	18.8	14.2	32.0
WV	14.9	4.37	20.1	19.0	14.0	32.3
DL	15.1	4.51	20.5	19.5	14.7	35.8
WV+CS	14.1	3.2	19.3	18.1	13.4	30.2
DL+CS	14.4	3.6	19.5	18.6	13.6	32.3

По каждой задаче выделено два лучших результата.

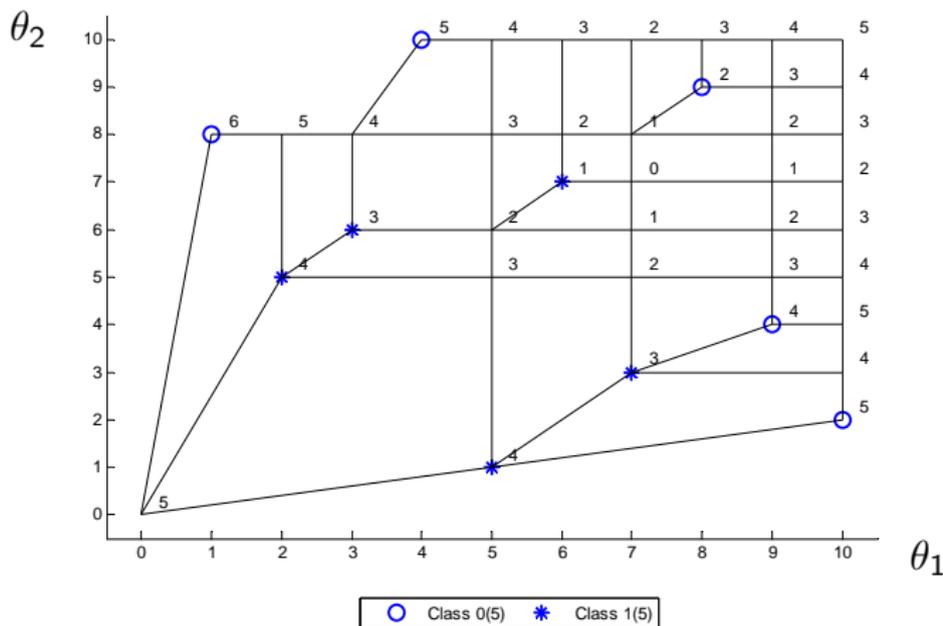
## Классы эквивалентности правил (точки — правила)

Пример: разделимая 2-мерная выборка,  $L = 10$ , два класса.  
 правила:  $r(x) = [f_1(x) \leq \theta_1] \cdot [f_2(x) \leq \theta_2]$ .



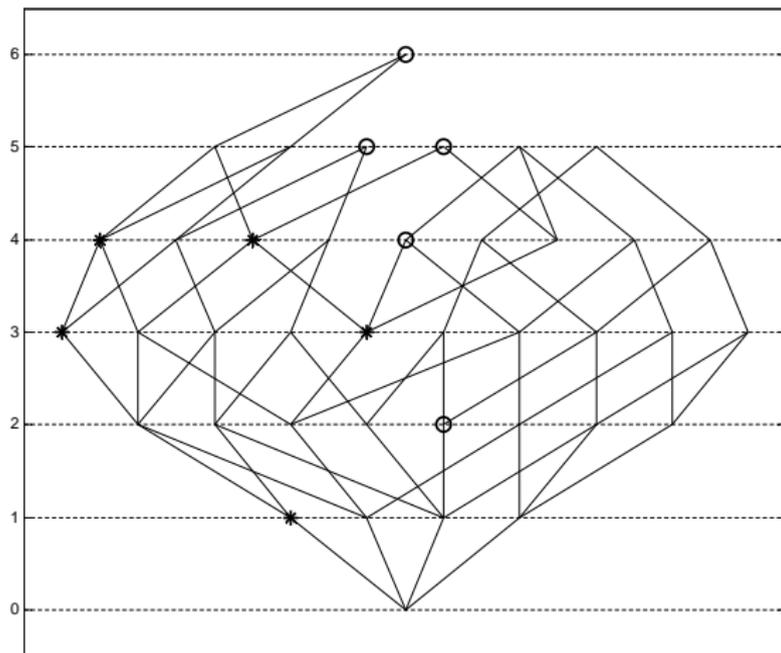
## Классы эквивалентности правил (точки — классы)

Пример: разделимая 2-мерная выборка,  $L = 10$ , два класса.  
 правила:  $r(x) = [f_1(x) \leq \theta_1 \text{ and } f_2(x) \leq \theta_2]$ .



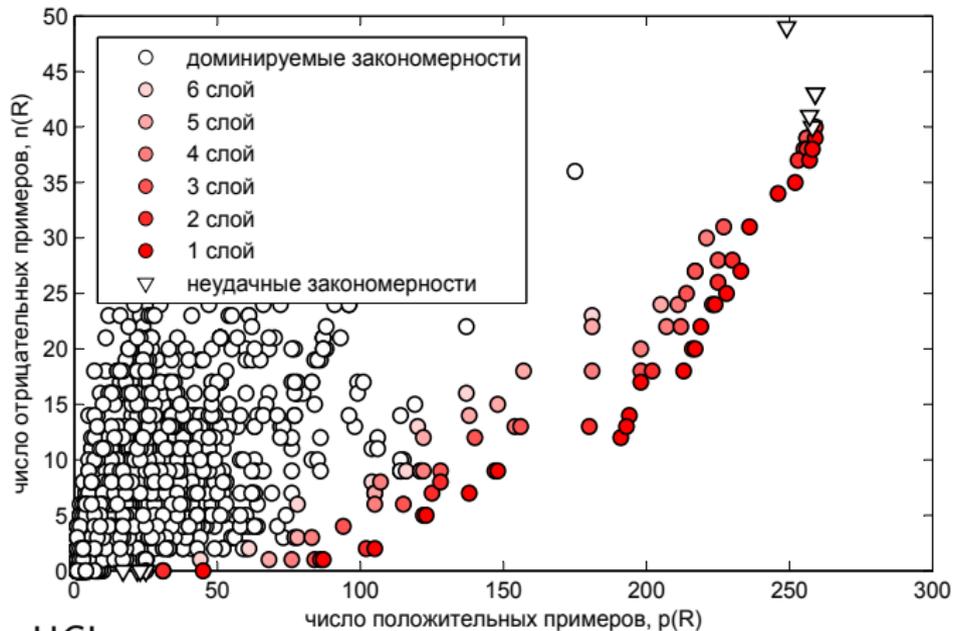
## Классы эквивалентности правил (граф P-C)

**Пример:** Граф расслоения–связности, изоморфный графу классов эквивалентности правил с предыдущего слайда.



## Парето-расслоение множества правил в плоскости $p, n$

**Парето-фронт** — множество недоминируемых закономерностей (точка  $(p, n)$  недоминируема, если правее и ниже точек нет)



задача UCI:german

## Выводы

- 1 Комбинаторные оценки вероятности переобучения учитывают расслоение и связность семейства алгоритмов.
- 2 Для многомерных сетей алгоритмов эти оценки точные.
- 3 Встраивание этих оценок в критерий информативности уменьшает частоту ошибок логических алгоритмов классификации на 1–2%.

Новые результаты:

- 1 Оценка вероятности переобучения обобщена на класс *монотонных* методов обучения.
- 2 Получены отдельные оценки для ошибок I и II рода и предложен критерий предсказанной информативности.
- 3 Обоснована (пока только экспериментально) замена ненаблюдаемой оценки наблюдаемой.

Спасибо за внимание!

Воронцов Константин Вячеславович  
[vokov@forecsys.ru](mailto:vokov@forecsys.ru)

Страницы на [www.MachineLearning.ru](http://www.MachineLearning.ru):

- Участник:Vokov
- Слабая вероятностная аксиоматика
- Расслоение и сходство алгоритмов (виртуальный семинар)