

# Концентрация меры в комбинаторных оценках обобщающей способности

Н. К. Животовский

Московский физико-технический институт  
Факультет управления и прикладной математики  
Кафедра "Интеллектуальные Системы"

Научный руководитель д.ф.-м.н., с.н.с. ВЦ РАН К. В. Воронцов

20 июня 2013 г.

## Содержание

- 1 Введение**
  - Определения
  - Вероятностные предположения
- 2 Оценки ожидаемой переобученности**
  - Равномерная оценка
  - Энтропийная оценка
  - Учёт метода обучения
  - Критерий точности
- 3 Концентрация меры**
  - Равномерный случай
  - Случай минимизации эмпирического риска
- 4 Выводы**

## Определения

- Задана конечная генеральная выборка объектов:  
 $\mathbb{X} = \{x_1, \dots, x_L\}$ .
- Задано множество  $A(\mathbb{X})$  алгоритмов (классификаторов), сопоставляющих объекту его класс.
- Бинарная функция потерь  $l: A \times \mathbb{X} \rightarrow \{0, 1\}$ , где  $(l(a, x) = 1) \leftrightarrow (a \text{ допускает ошибку на объекте } x)$ .
- Метод обучения  $\mu: \mathbb{X} \rightarrow A$ ,  $X \subseteq \mathbb{X}$ .
- Разбиение генеральной выборки на обучение и контроль  $\mathbb{X} = X \sqcup \bar{X}$ .
- Отношения порядка  
 $\forall a, b \in A, (a \leq b) \leftrightarrow (l(a, x) \leq l(b, x), \forall x \in \mathbb{X})$ , причем  
 $(a < b) \leftrightarrow (a \neq b, a \leq b)$ ,  $(a \prec b) \leftrightarrow (a < b, \rho(a, b) = 1)$
- Пессимистичная минимизация эмпирического риска  
 $\mu X = \arg \max_{a \in A(X)} n(a, \bar{X})$ , где  $A(X) = \arg \min_{a \in A} n(a, X)$ .

## Вероятностное предположение

### Предположение

*Все разбиения  $\mathbb{X} = X \cup \bar{X}$ ,  $|X| = \ell$ ,  $|\bar{X}| = k$  равновероятны.*

### Цель

*Для метода обучения  $\mu$  и обучающей выборки  $X$  оценить частоту ошибок на контрольной выборке  $\nu(\mu X, \bar{X})$ .*

Существующие оценки вероятности переобучения

$$Q_t(\mu, \mathbb{X}) = \mathbb{P}[\nu(\mu X, \bar{X}) - \nu(\mu X, X) \geq t]$$

существенно используют информацию о семействе алгоритмов на всей генеральной выборке.

## Возникающие задачи

### Цель

*Получение верхних оценок обобщающей способности, зависящих только от наблюдаемой обучающей выборки.*

### Определение

Функция

$${}_2F_1(a, b, c, z) = 1 + \sum_{k=1}^{\infty} \left[ \prod_{l=0}^{k-1} \frac{(a+l)(b+l)}{(1+l)(c+l)} \right] z^k$$

*с параметрами  $a, b, c$ , определённая в круге  $|z| < 1$ , называется гипергеометрической.*

## Оценка ожидаемой переобученности

### Теорема

Пусть  $\ell = k = \frac{L}{2}$  и  $\max_{a \in A} m(a) \leq \frac{L}{2}$ , тогда

$$\begin{aligned} \mathcal{E} \mathcal{O} \mathcal{F}_{\max}(\mathbb{X}) &= \mathbb{E} \max_{a \in A} (\nu(a, \bar{X}) - \nu(a, X)) \leq \\ &\min_{\lambda > 0} \frac{1}{\lambda} \left( \ln(\cosh(\lambda)) + \frac{1}{\ell} \ln \left( \sum_{s=0}^{\ell} \Delta_s \varphi(s, \ell, \lambda) \right) \right) \leq \\ &\sqrt{\frac{2 \ln(|A(\mathbb{X})|)}{\ell}}, \end{aligned}$$

где  $\Delta_s$  — число алгоритмов в  $s$ -ом слое семейства алгоритмов,

$$\varphi(s, \ell, \lambda) = {}_2F_1 \left( \frac{1-s}{2}, -\frac{s}{2}, \frac{1}{2} - \ell, (\tanh(\lambda))^2 \right)$$

## Энтропия семейства алгоритмов

Пусть  $\pi$  – некоторая вероятностная мера на  $A(\mathbb{X})$ .

Алгоритм  $a_0$  – корректный на  $\mathbb{X}$ .

$B(a, \varepsilon)$  – замкнутый шар с центром в алгоритме  $a$  и радиусом  $\varepsilon$  по метрике  $\sqrt{\rho(a, b)}$ .

### Определение

Энтропия семейства  $A(\mathbb{X})$

$$\frac{1}{3} \ln \left( \frac{2}{\min_{a \in A \cup a_0} \pi(B(a, \frac{d}{2}))} \right) + \frac{4}{3} \sum_{k=2}^{\infty} 2^{-k} \ln \left( \frac{2}{\min_{a \in A \cup a_0} \pi(B(a, \frac{d}{2^k}))} \right),$$

где  $d$  – диаметр семейства  $A(\mathbb{X}) \cup a_0$ .

## Энтропийная оценка переобученности

### Теорема

Пусть  $\ell = k = \frac{L}{2}$  и  $\max_{a \in A} m(a) \leq \frac{\ell}{2}$ , тогда

$$\mathcal{E}OF_{\max}(\mathbb{X}) \leq 3\sqrt{\frac{2Q(A(\mathbb{X}))}{\ell}}$$

Показано, что данная оценка не хуже, чем предыдущая (при оптимальном выборе меры  $\pi$ ).

Показано, что в случае, когда ёмкость семейства алгоритмов конечна, энтропия зависит лишь от ёмкости.



## Комбинаторные характеристики семейства алгоритмов

Для алгоритма  $a$  введем

- Порождающее множество

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

- Запрещающее множество

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

- Верхняя связность  $u(a) = |X_a|$ .

- Неполноценность  $q(a) = |X'_a|$ .

- Нижняя связность

$$d(a) = |\{x \in \mathbb{X} \mid \exists b \in A: b \prec a, I(b, x) < I(a, x)\}|.$$

- Графом расслоения–связности множества алгоритмов  $A$  будем называть направленный граф  $\langle A, E \rangle$  с множеством рёбер  $E = \{(a, b): a \prec b\}$ .

## Оценка, учитывающая метод обучения

### Теорема

Пусть метод обучения  $\mu$  – ПМЭР,  $\ell = k = \frac{\ell}{2}$  и  $\max_{a \in A} m(a) \leq \frac{\ell}{2}$

$$\mathcal{E} \mathcal{O} \mathcal{F}_{\mu}(\mathbb{X}) = \mathbb{E}(\nu(\mu X, \bar{X}) - \nu(\mu X, X)) \leq$$

$$\min_{\lambda > 0} \frac{1}{\lambda} \left( \ln(\cosh(\lambda)) + \frac{1}{\ell} \ln \left( \frac{\sum_{a \in A} \varphi(\ell, m, q, u, \lambda)}{C_{2\ell}^{\ell}} \right) \right).$$

$$\varphi(\ell, m, q, u, \lambda) =$$

$$\sum_{j=0}^{m-q} \sum_{i=0}^j (-1)^i C_j^i C_{2\ell-j-u-q}^{\ell-q+i-j} C_{m-q}^j (1 + \tanh(\lambda))^q \tanh^j(\lambda).$$

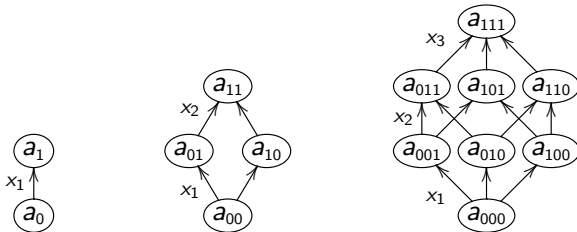
## Факторы завышенности

### Замечание

Предыдущая оценка не превосходит  $\sqrt{\frac{2 \ln(|A(\mathbb{X})|)}{\ell}}$

Завышенность предыдущей оценки появляется из-за использование неравенства

$$[\mu X = a] \leq [X_a \subset X][X'_a \subset \bar{X}].$$



## Критерий в терминах графа расслоения–связности

### Теорема

Пусть  $\ell \geq 2 \max_{a \in A} |X_a|$  и  $k \geq 2 \max_{a \in A} |X'_a|$ ,  $G$  — граф расслоения–связности множества  $A$ . Условие

$$[\mu X = a] = [X_a \subset X][X'_a \subset \bar{X}].$$

выполнено тогда и только тогда, когда

- 1 граф  $G$  имеет только один сток и один исток,
- 2 если  $a$  — произвольная вершина  $G$ ,  $x_1, \dots, x_n$  — объекты  $\mathbb{X}$ , соответствующие выходящим из нее ребрам, то  $G$  содержит направленный  $n$ -мерный куб с нижней вершиной  $a$ , построенный с помощью ребер, соответствующих  $x_1, \dots, x_n$ .

## Неравенство концентрации в равномерном случае.

### Теорема

При  $\ell = k = \frac{l}{2}$  для  $t > 0$  имеет место неравенство

$$\begin{aligned} P \left( \max_{a \in A(\mathbb{X})} (\nu(a, \bar{X}) - \nu(a, X)) \geq \mathcal{E} \mathcal{O} \mathcal{F}_{\max}(\mathbb{X}) + \frac{8\sqrt{\pi}\sigma}{\ell} + t \right) \\ \leq 2 \exp \left( -\frac{t^2 \ell^2}{16\sigma^2} \right), \text{ где } \sigma = \sqrt{\max_{a \in A(\mathbb{X})} m(a)} \leq \sqrt{2\ell} \end{aligned}$$

### Замечание

- $\mathcal{E} \mathcal{O} \mathcal{F}_{\max}(\mathbb{X})$  оценивается сверху величинами, не зависящими от ненаблюдаемой выборки.
- Локализация (уменьшение  $\max_{a \in A(\mathbb{X})} m(a)$ ) улучшает оценку.

## Неравенство концентрации с методом обучения.

- Для разбиения  $\mathbb{X} = X \sqcup \bar{X}$  и метода обучения  $\mu$  функционал переобученности  $\delta(\mu, X \sqcup \bar{X}) = \nu(\mu X, \bar{X}) - \nu(\mu X, X)$ .
- Для разбиения  $X \sqcup \bar{X}$  разбиение  $X' \sqcup \bar{X}'$  получается перестановкой некоторого элемента из  $X$  с элементом  $\bar{X}$ .
- Для чётных  $\ell$  естественным образом вводится  $\mathcal{E}\mathcal{O}\mathcal{F}_\mu(X)$

### Определение

ПМЭР  $\mu$  вместе с семейством  $A(\mathbb{X})$  называется  $(\alpha, \beta)$ -устойчивым, если для всех пар  $\{X \sqcup \bar{X}, X' \sqcup \bar{X}'\}$

$$|\delta(\mu, X \sqcup \bar{X}) - \delta(\mu, X' \sqcup \bar{X}')| \leq \frac{2\beta}{\ell},$$

$$\mathcal{E}\mathcal{O}\mathcal{F}_\mu(\mathbb{X}) - \mathcal{E}\mathcal{O}\mathcal{F}_\mu(X) \leq \alpha.$$

## Оценка ненаблюдаемой частоты ошибок.

### Теорема

Пусть  $\ell = k = \frac{L}{2}$ ,  $\mu$  —  $(\alpha, \beta)$ -устойчивый ПМЭР,  $t > 0$

$$\mathbb{P}\left(\nu(\mu X, \bar{X}) - \nu(\mu X, X) \geq \varepsilon \mathcal{O}\mathcal{F}_\mu(X) + \alpha + t\right) \leq \exp\left(\frac{-(\ell+1)t^2}{4\beta^2}\right).$$

### Следствие

С вероятностью не меньшей  $1 - \delta$

$$\nu(\mu X, \bar{X}) < \nu(\mu X, X) + \varepsilon \mathcal{O}\mathcal{F}_\mu(X) + \alpha + 2\beta \sqrt{\frac{\ln\left(\frac{1}{\delta}\right)}{\ell+1}}.$$

### Замечание

Оценка вычислима только по наблюдаемой выборке.

## Выводы

- Получены оценки ожидаемой переобученности с нелинейным вкладом алгоритмов.
- Исследованы степени вкладов алгоритмов в эти оценки.
- Получен критерий точности комбинаторных оценок обобщающей способности.
- Получены необходимые для комбинаторного подхода неравенства концентрации меры.
- Получена оценка ненаблюдаемой частоты ошибок в равномерном по семейству алгоритмов случае.
- Получена оценка ненаблюдаемой частоты ошибок для устойчивой минимизации эмпирического риска.