

Слабая вероятностная аксиоматика и надёжность эмпирических предсказаний

Воронцов Константин Вячеславович
voron@ccas.ru, <http://www.ccas.ru/voron>

Вычислительный Центр РАН,
Москва, Вавилова 40, 119991

Всероссийская конференция
Математические методы распознавания образов (ММРО-13)
30 сентября – 5 октября, 2007
Санкт-Петербург

Содержание

- 1 Теория надёжности эмпирических предсказаний**
 - Слабая вероятностная аксиоматика
 - Задача эмпирического предсказания
 - Примеры
 - Интерпретации, преимущества и недостатки
- 2 Теория обобщающей способности**
 - Классические оценки обобщающей способности
 - Оценки, зависящие от данных (Data-Dependent Bounds)
 - Измерение эффективного локального разнообразия
- 3 Профили обучаемости**
 - Разновидности профилей качества обучения
 - Пример: профиль компактности
 - Система эмпирического измерения качества обучения

Слабая вероятностная аксиоматика

- 1 $X^L = \{x_i\}_{i=1}^L$ — конечная выборка объектов из \mathbb{X} .
- 2 Все разбиения $X^L = X_n^\ell \cup X_n^k$, $n = 1, \dots, N$, где $N = C_L^k$, $L = \ell + k$, имеют равные шансы реализоваться.

Тогда...

Слабая вероятностная аксиоматика

- 1 $X^L = \{x_i\}_{i=1}^L$ — конечная выборка объектов из \mathbb{X} .
- 2 Все разбиения $X^L = X_n^\ell \cup X_n^k$, $n = 1, \dots, N$, где $N = C_L^k$, $L = \ell + k$, имеют равные шансы реализоваться.

Тогда:

- событие A есть функция $A: \{1, \dots, N\} \rightarrow \{0, 1\}$;
- вероятность события есть доля таких разбиений n , для которых $A(n) = 1$:

$$P_n A(n) = \frac{1}{N} \sum_{i=1}^N A(n).$$

Задача эмпирического предсказания

- Реализуется разбиение (X_n^ℓ, X_n^k) , $n \in \{1, \dots, N\}$;
выборка X_n^ℓ — *наблюдаемая*, выборка X_n^k — *скрытая*.
- Задана функция двух выборок $T: \mathbb{X}^k \times \mathbb{X}^\ell \rightarrow R$
- Требуется:
 1. Выбрать функцию $\hat{T}: \mathbb{X}^\ell \rightarrow R$ так, чтобы значение $\hat{T}_n = \hat{T}(X_n^\ell)$ предсказывало бы $T_n = T(X_n^k, X_n^\ell)$.
 2. Оценить точность предсказаний:

$$P_n[d(\hat{T}_n, T_n) > \varepsilon] \leq \eta(\varepsilon),$$

где $d(\hat{r}, r)$ — отклонение предсказания \hat{r} от истины r ,
например $d(\hat{r}, r) = |\hat{r} - r|$.

Слабая вероятностная аксиоматика

- Достаточна для доказательства фундаментальных фактов:
 - закон больших чисел, с точной оценкой скорости сходимости;
 - отклонение эмпирических распределений (критерий Смирнова), с точной оценкой скорости сходимости;
 - непараметрические критерии проверки стат. гипотез;
 - оценки типа Вапника-Червоненкиса.

Слабая вероятностная аксиоматика

- Достаточна для доказательства фундаментальных фактов:
 - закон больших чисел, с точной оценкой скорости сходимости;
 - отклонение эмпирических распределений (критерий Смирнова), с точной оценкой скорости сходимости;
 - непараметрические критерии проверки стат. гипотез;
 - оценки типа Вапника-Червоненкиса.
- Основана на более слабых предположениях:
 - нет определения вероятности как меры на \mathbb{X} ;
 - нет определения вероятности как частоты при $L \rightarrow \infty$;

Слабая вероятностная аксиоматика

- Достаточна для доказательства фундаментальных фактов:
 - закон больших чисел, с точной оценкой скорости сходимости;
 - отклонение эмпирических распределений (критерий Смирнова), с точной оценкой скорости сходимости;
 - непараметрические критерии проверки стат. гипотез;
 - оценки типа Вапника-Червоненкиса.
- Основана на более слабых предположениях:
 - нет определения вероятности как меры на \mathbb{X} ;
 - нет определения вероятности как частоты при $L \rightarrow \infty$;
 - **нет необходимости определять понятие «вероятность»!**

Пример 1: Закон больших чисел

Опр. Частота события $S \subset \mathbb{X}$ на конечной выборке $U \subset \mathbb{X}$:

$$\nu_S(U) = \frac{1}{|U|} \sum_{u \in U} [u \in S].$$

Положим: $R = \mathbb{R}$, $\hat{T}(U) = T(U) = \nu_S(U)$.

Теорема (известный классический факт)

Частота $\nu_S(X_n^\ell)$ предсказывает частоту $\nu_S(X_n^k)$,
причём справедлива **точная** оценка

$$P_n[\nu_S(X_n^k) - \nu_S(X_n^\ell) \geq \varepsilon] = H\left(\begin{smallmatrix} \ell & s(\varepsilon) \\ L & m \end{smallmatrix}\right),$$

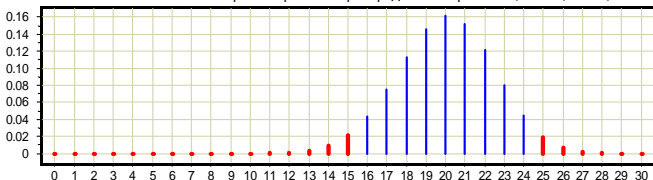
где $H\left(\begin{smallmatrix} \ell & s \\ L & m \end{smallmatrix}\right)$ — левый хвост гипергеометрического
распределения, $s(\varepsilon) = \lfloor \frac{\ell}{L}(m - \varepsilon k) \rfloor$, $m = L\nu_S(X^L)$.

Пример 1: Закон больших чисел

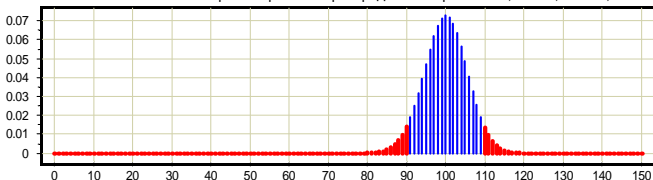
Левый хвост гипергеометрического распределения:

$$H\left(\ell \begin{matrix} s(\varepsilon) \\ L \\ m \end{matrix}\right) = \sum_{t=s_0}^{s(\varepsilon)} \frac{C_m^t C_{L-m}^{\ell-t}}{C_L^\ell}, \quad s_0 = \max\{0, m - k\}$$

H Гипергеометрическое распределение при L=300, k=100, m=30, eta=0.05



H Гипергеометрическое распределение при L=1500, k=500, m=150, eta=0.05



Пример 2: Сходимость эмпирических распределений

Опр. Эмпирическое распределение функции $\xi: \mathbb{X} \rightarrow \mathbb{R}$ на конечной выборке $U \subseteq \mathbb{X}$ есть

$$F_{\xi}(z, U) = \frac{1}{|U|} \sum_{x \in U} [\xi(x) \leq z].$$

Положим: R — пр-во кус.-пост. невозр. функций $F: \mathbb{R} \rightarrow [0, 1]$,
 $T(U) = \hat{T}(U) = F_{\xi}(z, U)$, $d(\hat{r}, r) = \max_{z \in \mathbb{R}} |r(z) - \hat{r}(z)|$.

Теорема (малоизвестный факт)

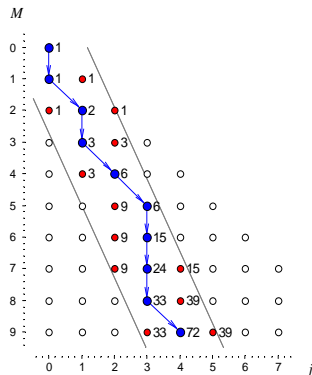
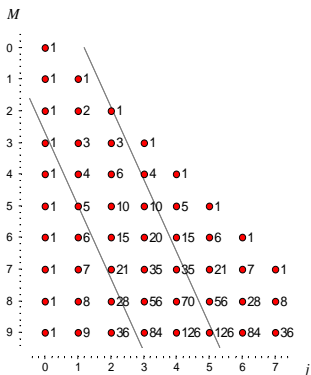
Если значения $\xi(x_i)$ попарно различны на X^L , то

$$P_n \left\{ \max_{z \in \mathbb{R}} |F_{\xi}(z, X_n^k) - F_{\xi}(z, X_n^{\ell})| > \varepsilon \right\} = \frac{G_L^k(\varepsilon)}{C_L^k},$$

где $G_L^k(\varepsilon)$ — значение из усечённого треугольника Паскаля.

Пример 2: Сходимость эмпирических распределений

Усечённый треугольник Паскаля $G_m^j(\varepsilon)$ — между двумя прямыми $j^+(m) = \frac{k}{L}(m + \varepsilon l)$, $j^-(m) = \frac{k}{L}(m - \varepsilon l)$.



Здесь $L = 16$, $k = 7$, $\varepsilon = 0.3$.

Связь с классической аксиоматикой Колмогорова

Теорема (Принцип соответствия)

Если в слабой аксиоматике для некоторой функции $\phi(X^\ell, X^k)$ получена оценка

$$Q_\varepsilon(X^L) = P_n[\phi(X_n^\ell, X_n^k) > \varepsilon] \leq \eta(\varepsilon, X^L)$$

то аналогичная оценка верна и в аксиоматике Колмогорова

$$E_{X^L} Q_\varepsilon(X^L) = P_{X^L}[\phi(X^\ell, X^k) > \varepsilon] \leq E_{X^L} \eta(\varepsilon, X^L)$$

Связь с классической аксиоматикой Колмогорова

Теорема (Принцип соответствия)

Если в слабой аксиоматике для некоторой функции $\phi(X^\ell, X^k)$ получена оценка

$$Q_\varepsilon(X^L) = P_n[\phi(X_n^\ell, X_n^k) > \varepsilon] \leq \eta(\varepsilon, X^L)$$

то аналогичная оценка верна и в аксиоматике Колмогорова

$$E_{X^L} Q_\varepsilon(X^L) = P_{X^L}[\phi(X^\ell, X^k) > \varepsilon] \leq E_{X^L} \eta(\varepsilon, X^L)$$

Трансдукция или индукция?

Если $\eta(\varepsilon, X^L) \equiv \eta(\varepsilon)$,

то оценка справедлива для любой выборки X^L .

Связь с эмпирическим оцениванием

Скольльзящий контроль (Cross-Validation)

Если оценку не удаётся получить теоретически или если теоретическая оценка слишком завышена:

$$Q_\varepsilon(X^L) = P_n[d(\hat{T}_n, T_n) > \varepsilon] \leq \boxed{???},$$

то можно измерить Q_ε эмпирически:

$$Q_\varepsilon(X^L) \leq \frac{1}{|N'|} \sum_{n \in N'} [d(\hat{T}_n, T_n) > \varepsilon] + \underbrace{\delta(N, |N'|)}_{\rightarrow 0 \text{ при } |N'| \rightarrow N},$$

где множество разбиений $N' \subset \{1, \dots, N\}$ выбирается

- либо случайно (метод Монте-Карло),
- либо по блокам (k -fold Cross-Validation)

Слабая вероятностная аксиоматика: за и против

- + Лучше подходит для задач анализа данных и обучения по прецедентам
- ... но хуже — для континуальных стохастических явлений

Слабая вероятностная аксиоматика: за и против

- + Лучше подходит для задач анализа данных и обучения по прецедентам
- ... но хуже — для непрерывных стохастических явлений
- + Даёт не асимптотические, не завышенные оценки
- ... требующие сложных комбинаторных вычислений

Слабая вероятностная аксиоматика: за и против

- + Лучше подходит для задач анализа данных и обучения по прецедентам
- ... но хуже — для континуальных стохастических явлений
- + Даёт не асимптотические, не завышенные оценки
- ... требующие сложных комбинаторных вычислений
- + Удовлетворяет «принципу соответствия»
- ... однако не все классические теоремы имеют аналоги в слабой аксиоматике

Слабая вероятностная аксиоматика: за и против

- + Лучше подходит для задач анализа данных и обучения по прецедентам
- ... но хуже — для континуальных стохастических явлений
- + Даёт не асимптотические, не завышенные оценки
- ... требующие сложных комбинаторных вычислений
- + Удовлетворяет «принципу соответствия»
- ... однако не все классические теоремы имеют аналоги в слабой аксиоматике

Открытая проблема:

Какая часть математической статистики может быть построена в рамках слабой аксиоматики?

Задача обучения по прецедентам

- \mathbb{X} — множество объектов, \mathbb{Y} — множество ответов.
- Бинарная функция потерь $\mathcal{L}: \mathbb{X} \times \mathbb{Y} \rightarrow \{0, 1\}$.
 - Для задач классификации: $\mathcal{L}(x, y) = [y \neq y^*(x)]$,
 - Для задач регрессии: $\mathcal{L}(x, y) = [|y - y^*(x)| > \delta]$,
где $y^*(x)$ — неизвестная целевая функция.
- Обучающая выборка $X^\ell = \{x_i\}_{i=1}^\ell$ с известными $\mathcal{L}(x_i, y)$.
- Метод обучения $\mu: X^\ell \mapsto f$.
- Частота ошибок алгоритма $f: \mathbb{X} \rightarrow \mathbb{Y}$ на выборке $U \subset \mathbb{X}$:
$$\nu(f, U) = \frac{1}{|U|} \sum_{u \in U} \mathcal{L}(u, f(u)).$$
- Обобщающая способность:
средняя ошибка $\nu(\mu X^\ell, U)$ должна быть **достаточно мала**
для **большинства** выборок $U \in \mathbb{X}^*$.

Оценки Вапника-Червоненкиса [1971]

- Для семейства алгоритмов F в сильной аксиоматике:

$$P_\varepsilon(F) = P_{X^L} \left[\sup_{f \in F} (\nu(f, X^k) - \nu(f, X^\ell)) > \varepsilon \right] \\ \leq \Delta^F(L) \cdot 1.5 e^{-\varepsilon^2 \ell} \quad (\text{при условии } \ell = k);$$

- $\Delta^F(L)$ — функция роста (shatter coefficient) семейства F — макс. число функций $f \in F$, попарно различимых на X^L ;
 $\Delta^F(L) \leq 1.5 \frac{L^h}{h!}$, $h = VCdim(F)$ — ёмкость семейства F .

Оценки Вапника-Червоненкиса [1971]

- Для семейства алгоритмов F в сильной аксиоматике:

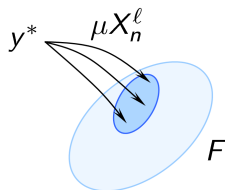
$$\begin{aligned} P_\varepsilon(F) &= P_{X^L} \left[\sup_{f \in F} (\nu(f, X^k) - \nu(f, X^\ell)) > \varepsilon \right] \\ &\leq \Delta^F(L) \cdot 1.5 e^{-\varepsilon^2 \ell} \quad (\text{при условии } \ell = k); \end{aligned}$$

- $\Delta^F(L)$ — функция роста (shatter coefficient) семейства F — макс. число функций $f \in F$, попарно различимых на X^L ;
 $\Delta^F(L) \leq 1.5 \frac{L^h}{h!}$, $h = VCdim(F)$ — ёмкость семейства F .
- Оценка крайне завышена, т. к. она не учитывает:
 - особенности выборки X^ℓ ;
 - особенности целевой зависимости $y^*(x)$;
 - особенности метода обучения μ ;
 - степень различности получаемых алгоритмов;
 - что $1.5 e^{-\varepsilon^2 \ell}$ — лишь асимптотическая аппроксимация.

Оценки, зависящие от данных

- *Эффект локализации:*

Если зависимость y^* , метод μ и выборка X^L фиксированы, то не все функции из F могут быть результатом обучения.



- Равномерная сходимость [Вапник, Червоненкис, 1969]
- Theory of learnable (PAC-learning) [Valiant, 1982]
- Concentration inequalities [Talagrand, 1995]
- Data-dependent bounds [Haussler, 1992 – Bartlett, 1998 – ...]
- Self-bounding learning algorithms [Freund, 1998]
- PAC-Bayesian model averaging [McAllester, 1999]
- Microchoice bounds [Langford, Blum, 2001]
- Algorithmic luckiness [Herbrich, Williamson, 2002]
- Tight sample complexity bounds [Langford, PhD, 2002]

Оценки, зависящие от данных

- В слабой вероятностной аксиоматике:

$$\begin{aligned} Q_\varepsilon(\mu, X^L) &= P_n \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \varepsilon \right] \\ &\leq \Delta_L^\ell(\mu, X^L) \cdot \max_m H_{Lm}^{\ell s(\varepsilon)} \\ &(\leq \Delta^F(L) \cdot 1.5 e^{-\varepsilon^2 \ell}); \end{aligned}$$

где алгоритм $f_n = \mu X_n^\ell$ — результат обучения на X_n^ℓ ;

- $\Delta_L^\ell(\mu, X^L)$ — локальный коэффициент разнообразия (*local shatter coefficient*) множества алгоритмов $\{f_n = \mu X_n^\ell \mid n = 1, \dots, N\}$, которые могут быть результатом обучения для данной задачи $\langle \mu, X^L \rangle$:

Оценки, зависящие от данных: дальнейшее уточнение

- **Идея:** скалярная характеристика разнообразия содержит слишком мало информации о процессе обучения.
- Вводится *профиль разнообразия (shatter profile)* $\{D_m\}_{m=0}^L$:

$$\Delta_L^\ell(\mu, X^L) = \sum_{m=1}^L D_m(\mu, X^L),$$

где $D_m(\mu, X^L)$ — локальный коэффициент разнообразия множества алгоритмов $\{f_n \mid \nu(f_n, X^L) = \frac{m}{L}, n = 1, \dots, N\}$.

Оценки, зависящие от данных: дальнейшее уточнение

- **Идея:** скалярная характеристика разнообразия содержит слишком мало информации о процессе обучения.
- Вводится профиль разнообразия (shatter profile) $\{D_m\}_{m=0}^L$:

$$\Delta_L^\ell(\mu, X^L) = \sum_{m=1}^L D_m(\mu, X^L),$$

где $D_m(\mu, X^L)$ — локальный коэффициент разнообразия множества алгоритмов $\{f_n \mid \nu(f_n, X^L) = \frac{m}{L}, n = 1, \dots, N\}$.

Теорема

В слабой аксиоматике справедлива оценка:

$$Q_\varepsilon(\mu, X^L) \leq \sum_{m=1}^L D_m(\mu, X^L) \cdot H\binom{\ell}{L m}^{s(\varepsilon)};$$

Эффективный локальный профиль разнообразия

- Обратная задача:
при каком профиле D_m оценка была бы точной?
- Теорема

$$Q_{\varepsilon, m}(\mu, X^L) = P_n \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \varepsilon \right] \left[\nu(f_n, X^L) = \frac{m}{L} \right] \\ \leq D_m(\mu, X^L) \cdot H(L_m^{\ell s(\varepsilon)});$$

Заменяем здесь « \leq » на « $=$ » и выразим D_m :

Эффективный локальный профиль разнообразия

- Обратная задача:
при каком профиле D_m оценка была бы точной?
- Теорема

$$Q_{\varepsilon, m}(\mu, X^L) = P_n \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \varepsilon \right] \left[\nu(f_n, X^L) = \frac{m}{L} \right] \\ \leq D_m(\mu, X^L) \cdot H_{Lm}^{\ell s(\varepsilon)};$$

Заменим здесь « \leq » на « $=$ » и выразим D_m :

- Эффективный локальный профиль разнообразия. $\{\hat{D}_m\}_{m=0}^L$:

$$\hat{D}_m = \frac{\frac{1}{|N'|} \sum_{n \in N'} \left[\nu(f_n, X_n^k) - \nu(f_n, X_n^\ell) > \varepsilon \right] \left[\nu(f_n, X^L) = \frac{m}{L} \right]}{H_{Lm}^{\ell s(\varepsilon)}}.$$

- Эффективный локальный коэффициент разнообразия:
 $\hat{\Delta}_L^\ell = \hat{D}_0 + \dots + \hat{D}_L.$

Эмпирическое исследование

- Сравнить численно причины завышенности:
 - пренебрежение эффектом локализации
 - свёртка профиля в коэффициент разнообразия
 - пренебрежение степенью различности алгоритмов
- На примере логических алгоритмов классификации [см. доклад А. А. Ивахненко]
 - известна функция роста $\Delta^F(L)$
 - легко оценить снизу локальный профиль $\Delta_L^\ell(\mu, X^L)$
 - легко оценить эффективный локальный профиль

Разновидности профилей качества обучения

- Профиль разнообразия (или сложности) [см. выше]
- Профиль делимости (margin distribution)
 - [Mason, Bartlett, Baxter, 1998]
 - [Shawe-Taylor, Cristianini, 1999]
 - [Garg 2002, 2003]
- Профиль монотонности [Воронцов, 2002]
- Профиль компактности [Воронцов, 2004]
- Профиль устойчивости [???
- Профиль чего ещё ???

Определение профиля компактности

- Профиль компактности выборки X^L — доля объектов, у которых m -й сосед лежит в другом классе:

$$R(m) = \frac{1}{L} \sum_{i=1}^L [y_i \neq y_i^{(m)}], \quad m = 1, \dots, L-1.$$

- Функционал полного скользящего контроля:

$$Q_c(\mu, X^L) = \frac{1}{N} \sum_{n=1}^N \nu(\mu X_n^\ell, X_n^k).$$

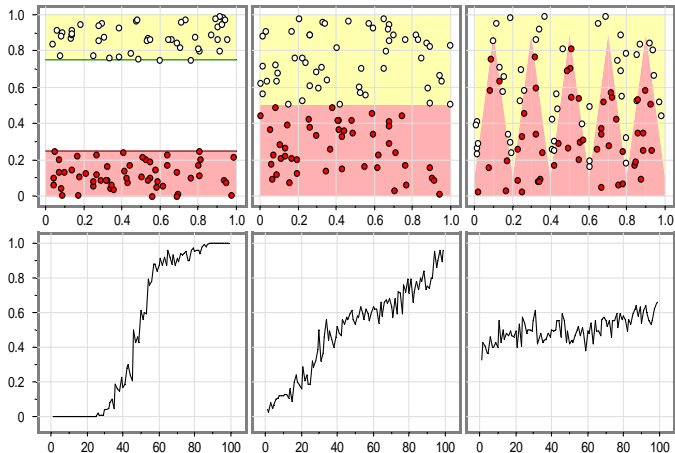
Теорема

Справедлива **точная** оценка:

$$Q_c(\mu, X^L) = \sum_{m=1}^k R(m) \frac{C_{L-1-m}^{\ell-1}}{C_{L-1}^\ell}.$$

Пример: профиль компактности

Серия модельных выборок, $L = 100$:



Методика эмпирического измерения качества обучения

- Оценки скользящего контроля
- Профиль разнообразия
- Распределения вариации и смещения по объектам
- Профиль представительности объектов
- Профиль устойчивости
- Профиль делимости
- Распределение признаков по информативности
- Временные показатели эффективности обучения

[см. доклад А. В. Лисицы]