

Комбинаторная теория переобучения

К. В. Воронцов
vokov@forecsys.ru

Этот курс доступен на странице вики-ресурса
<http://www.MachineLearning.ru/wiki>
«Теория надёжности обучения по прецедентам
(курс лекций, К. В. Воронцов)»

25 марта 2015 • 1 апреля 2015

1 Основные понятия

- Вероятность переобучения
- Взаимосвязи с SLT и интерпретации
- Теория Вапника–Червоненкиса

2 Эксперименты с переобучением

- Переобучение четырёх модельных семейств
- Переобучение цепей

3 Оценки расслоения–связности

- Граф расслоения–связности
- Порождающие и запрещающие множества
- Основная оценка расслоения–связности

Бинарная функция потерь. Матрица ошибок

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

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

$I(a, x) = [$ алгоритм a ошибается на объекте x];

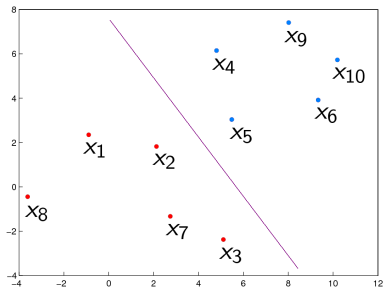
$L \times D$ -матрица ошибок с попарно различными столбцами:

	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^ℓ — наблюдаемая (обучающая) выборка длины ℓ
\dots	0	0	0	0	1	1	\dots	1	
x_ℓ	0	0	1	0	0	0	\dots	0	
$x_{\ell+1}$	0	0	0	1	1	1	\dots	0	X^k — скрытая (контрольная) выборка длины $k = L - \ell$
\dots	0	0	0	1	0	0	\dots	1	
x_L	0	1	1	1	1	1	\dots	0	

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

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

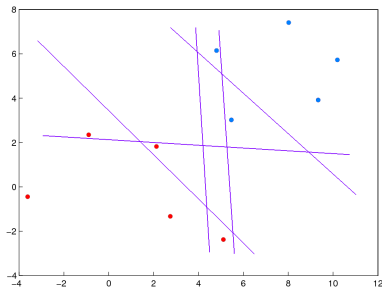
Пример. Матрица ошибок линейных классификаторов



1 вектор с 0 ошибками

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

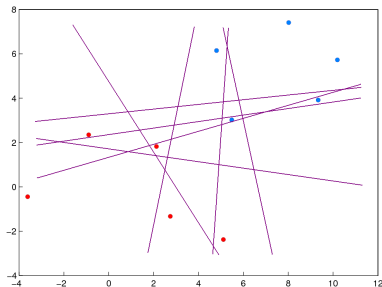
Пример. Матрица ошибок линейных классификаторов



1 вектор с 0 ошибками
5 векторов с 1 ошибкой

x_1	0	1	0	0	0	0
x_2	0	0	1	0	0	0
x_3	0	0	0	1	0	0
x_4	0	0	0	0	1	0
x_5	0	0	0	0	0	1
x_6	0	0	0	0	0	0
x_7	0	0	0	0	0	0
x_8	0	0	0	0	0	0
x_9	0	0	0	0	0	0
x_{10}	0	0	0	0	0	0

Пример. Матрица ошибок линейных классификаторов



1 вектор с 0 ошибками
 5 векторов с 1 ошибкой
 8 векторов с 2 ошибками
 и т. д...

x_1	0	1	0	0	0	0	1	0	0	0	0	1	1	0	...
x_2	0	0	1	0	0	0	1	1	0	0	0	0	0	0	...
x_3	0	0	0	1	0	0	0	1	1	0	0	0	0	1	...
x_4	0	0	0	0	1	0	0	0	1	1	0	0	0	0	...
x_5	0	0	0	0	0	1	0	0	0	1	1	1	0	0	...
x_6	0	0	0	0	0	0	0	0	0	1	0	1	0	0	...
x_7	0	0	0	0	0	0	0	0	0	0	0	0	0	1	...
x_8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
x_9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
x_{10}	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...

Задача обучения по прецедентам

$\mu: X \mapsto a$ — метод обучения, произвольной выборке $X \subset X^L$ ставит в соответствие некоторый алгоритм $a \in A$.

$\mu(X) = \arg \min_{a \in A} \nu(a, X)$ — минимизация эмпирического риска

$\delta(\mu, X^\ell) = \nu(\mu(X^\ell), X^k) - \nu(\mu(X^\ell), X^\ell)$ — переобученность

$\delta(\mu, X^\ell) \geq \varepsilon$ — событие переобучения

Основное вероятностное предположение:

все разбиения $X^\ell \sqcup X^k = X^L$ равновероятны

(слабый вариант **гипотезы независимости** выборки X^L).

Основная задача — оценить **вероятность** переобучения:

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

Обобщающая способность метода обучения

$P \equiv E \equiv \frac{1}{C_L^\ell} \sum_{X^\ell \subset X^L}$ — доля разбиений выборки.

$\delta(\mu, X^\ell) = \nu(\mu(X^\ell), X^k) - \nu(\mu(X^\ell), X^\ell)$ — переобученность.

Функционалы обобщающей способности:

- Полный скользящий контроль (Complete Cross-Validation):

$$CCV(\mu, X^L) = E \nu(\mu(X^\ell), X^k).$$

- Ожидаемая переобученность (Expected OverFitting):

$$EOF(\mu, X^L) = E \delta(\mu, X^\ell).$$

- Вероятность большой частоты ошибок на контроле:

$$R_\varepsilon(\mu, X^L) = P[\nu(\mu(X^\ell), X^k) \geq \varepsilon].$$

- Вероятность переобучения:

$$Q_\varepsilon(\mu, X^L) = P[\delta(\mu, X^\ell) \geq \varepsilon].$$

Отличия от постановки задачи, принятой в SLT

Функционалы обобщающей способности:

- Вероятность равномерного отклонения частоты $\nu(a, X^\ell)$ от вероятности ошибки $P(a)$ [Вапник, Червоненкис, 1971]

$$S_\varepsilon(A, X^L) = P \left[\sup_{a \in A} (P(a) - \nu(a, X^\ell)) \geq \varepsilon \right]$$

- Вероятность переобучения

$$Q_\varepsilon(\mu, X^L) = P[\nu(\mu(X^\ell), X^k) - \nu(\mu(X^\ell), X^\ell) \geq \varepsilon].$$

Основные отличия:

- 1 Отказываемся от завышенной оценки \sup по всему A .
- 2 Стараемся учитывать особенности метода обучения μ .
- 3 Отказываемся оценивать ненаблюдаемую величину $P(a)$.
- 4 Отказываемся от лишних вероятностных допущений.

Связь с кросс-валидацией

Полный скользящий контроль и вероятность переобучения:

$$\text{CCV}(\mu, X^L) = \frac{1}{C_L^\ell} \sum_{X^\ell \subset X^L} \nu(\mu(X^\ell), X^k);$$

$$Q_\varepsilon(\mu, X^L) = \frac{1}{C_L^\ell} \sum_{X^\ell \subset X^L} [\nu(\mu(X^\ell), X^k) - \nu(\mu(X^\ell), X^\ell) \geq \varepsilon].$$

Оценки вероятности методом Монте–Карло — как доли разбиений выборки из случайного подмножества N разбиений, $|N|$ порядка 10^3 – 10^4 :

$$\widehat{\text{CCV}}(\mu, X^L) = \frac{1}{|N|} \sum_{(X^k, X^\ell) \in N} \nu(\mu(X^\ell), X^k);$$

$$\widehat{Q}_\varepsilon(\mu, X^L) = \frac{1}{|N|} \sum_{(X^k, X^\ell) \in N} [\nu(\mu(X^\ell), X^k) - \nu(\mu(X^\ell), X^\ell) \geq \varepsilon].$$

Связь с радемахеровской сложностью

Радемахеровская сложность множества A на X^L

$$\text{RC}(A, X^L) = \mathbf{E}_\sigma \sup_{a \in A} \frac{2}{L} \sum_{i=1}^L \sigma_i l(a, x_i), \quad \sigma_i = \begin{cases} +1, & \text{вер. } \frac{1}{2} \\ -1, & \text{вер. } \frac{1}{2} \end{cases}$$

$\sigma_1, \dots, \sigma_L$ — независимые радемахеровские случ. величины

Ожидаемая переобученность в случае $\ell = k$:

$$\text{EOF}(\mu, X^L) = \mathbf{E} \sup_{a \in A} \frac{2}{L} \sum_{i=1}^L \sigma_i l(a, x_i), \quad \sigma_i = \begin{cases} +1, & x_i \in X^k \\ -1, & x_i \in X^\ell \end{cases}$$

если μ — метод максимизации переобученности:

$$\mu(X) = \arg \max_{a \in A} \left(\nu(a, X^k) - \nu(a, X^\ell) \right)$$

Переход от комбинаторных оценок к SLT

Обычные вероятностные предположения:

X^L — i.i.d. выборка

из вероятностного пространства $\langle \mathcal{X}, \sigma, P \rangle$ над бесконечным \mathcal{X}

Переход от комбинаторной оценки к вероятностной:

- 1 Пусть имеется комбинаторная оценка:

$$P_{X \sim X^L} [\delta(\mu, X) \geq \varepsilon] = Q_\varepsilon(\mu, X^L) \leq \eta(\varepsilon, X^L)$$

- 2 Возьмём матожидание от обеих частей неравенства по X^L :

$$P_{\substack{X \sim \mathcal{X}^l \\ X^k \sim \mathcal{X}^k}} [\delta(\mu, X) \geq \varepsilon] = \mathbf{E}_{X^L} Q_\varepsilon(\mu, X^L) \leq \mathbf{E}_{X^L} \eta(\varepsilon, X^L).$$

Если оценка $\eta(\varepsilon, X^L)$ не зависит от выборки X^L ,
 то она переносится из КТП в SLT непосредственно.

Несколько цитат

А. Н. Колмогоров. Теория информации и теория алгоритмов:

«представляется важной задача освобождения всюду, где это возможно, от *излишних вероятностных допущений*.

На независимой ценности *чисто комбинаторного* подхода к теории информации я неоднократно настаивал в своих лекциях.»

В. Д. Гоппа. Введение в алгебраическую теорию информации:

«Надобность в вероятностной модели отпадает, поскольку теория информации оказывается достаточно интересной и богатой приложениями в алгебраической постановке. Одним из таких приложений является распознавание образов.»

И ещё пара цитат

Ю. К. Беляев. Вероятностные методы выборочного контроля:
«возникло глубокое убеждение, что в теории выборочных методов можно получить содержательные аналоги большинства основных утверждений теории вероятностей и математической статистики, которые к настоящему времени найдены в предположении *взаимной независимости* результатов измерений.»

А. Н. Колмогоров. Теория информации и теория алгоритмов:
«чистая математика благополучно развивается как по преимуществу наука о бесконечном. . . Весьма вероятно, что с развитием современной вычислительной техники будет понято, что в очень многих случаях разумно изучение реальных явлений вести, избегая промежуточный этап их стилизации в духе представлений математики бесконечного и непрерывного, *переходя прямо к дискретным моделям.*»

Две альтернативные интерпретации

Задача комбинаторной теории переобучения:

- 1 Оценивание вероятности переобучения, математического ожидания переобученности или частоты ошибок на контроле в *слабой вероятностной аксиоматике*, предполагающей, что все разбиения $X^\ell \sqcup X^k = X^L$ равновероятны.
- 2 Получение комбинаторных формул для эффективного вычисления функционалов полного скользящего контроля Q_ε , CCV и др., основанных на усреднении *по всем* разбиениям заданной конечной выборки X^L на обучающую X^ℓ и контрольную X^k подвыборки, *без явного применения метода обучения μ* .

Простейший, но важный частный случай

Пусть $A = \{a\}$ — одноэлементное множество, $m = n(a, X^L)$.

Тогда вероятность переобучения есть вероятность большого отклонения частот ошибок в двух подвыборках:

$$Q_\varepsilon(a, X^L) = P[\nu(a, X^k) - \nu(a, X^\ell) \geq \varepsilon].$$

Теорема

Для любого X^L , любого $\varepsilon \in [0, 1]$

$$Q_\varepsilon(a, X^L) = \mathcal{H}_L^{\ell, m} \left(\frac{\ell}{L}(m - \varepsilon k) \right),$$

где $\mathcal{H}_L^{\ell, m}(z) = \sum_{s=0}^{\lfloor z \rfloor} \frac{C_m^s C_{L-m}^{\ell-s}}{C_L^\ell}$ — функция гипергеометрического распределения.

Доказательство

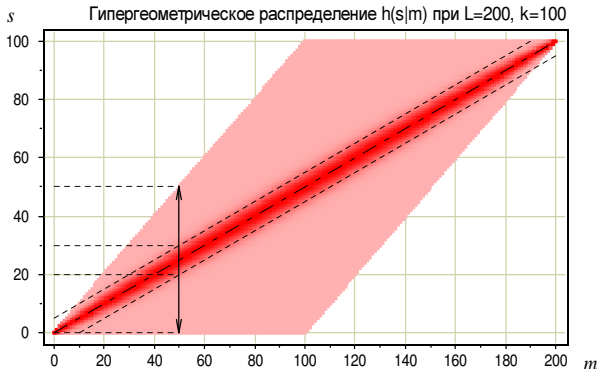
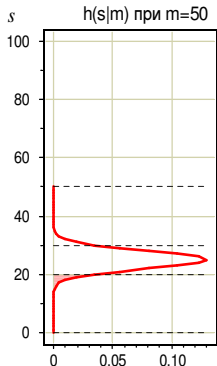
1. Обозначим $s = n(a, X^\ell)$.
2. «Школьная» задача по теории вероятностей:
 в урне L шаров, m из них чёрные; извлекаем ℓ шаров наугад.
 Какова вероятность того, что s из них чёрные?

$$P[n(a, X^\ell) = s] = C_m^s C_{L-m}^{\ell-s} / C_L^\ell.$$

3. Распишем Q_ε , подставив $\nu(a, X^k) = \frac{m-s}{k}$, $\nu(a, X^\ell) = \frac{s}{\ell}$:

$$\begin{aligned} Q_\varepsilon(a, X^\ell) &= P[\nu(a, X^k) - \nu(a, X^\ell) \geq \varepsilon] = \\ &= \sum_{s=0}^{\ell} \underbrace{\left[\frac{m-s}{k} - \frac{s}{\ell} \geq \varepsilon \right]}_{s \leq \frac{\ell}{L}(m-\varepsilon k)} \underbrace{P[n(a, X^\ell) = s]}_{C_m^s C_{L-m}^{\ell-s} / C_L^\ell} = \\ &= \mathcal{H}_L^{\ell, m} \left(\frac{\ell}{L}(m - \varepsilon k) \right). \quad \blacksquare \end{aligned}$$

Гипергеометрическое распределение $h(s|m) = C_m^s C_{L-m}^{\ell-s} / C_L^\ell$



Предсказание числа $m = n(a, X^L)$ по числу $s = n(a, X^\ell)$ возможно благодаря узости гипергеометрического пика, причём при $\ell, k \rightarrow \infty$ он сужается, и $\nu(a, X^\ell) \rightarrow \nu(a, X^k)$ (явление концентрации вероятности, закон больших чисел).

Теория Вапника–Червоненкиса

Рассмотрим общий случай — A произвольное, конечное.

1. Вероятность переобучения оценим сверху вероятностью большого *равномерного отклонения* частот: для любых X^L, μ

$$\begin{aligned} Q_\varepsilon(\mu, X^L) &= P[\delta(\mu, X^L) \geq \varepsilon] \leq \\ &\leq P\left[\max_{a \in A} \delta(a, X^L) \geq \varepsilon\right] = \tilde{Q}_\varepsilon(A, X^L). \end{aligned}$$

2. Оценим вероятность объединения событий суммой их вероятностей (неравенство Буля, union bound):

$$\begin{aligned} \tilde{Q}_\varepsilon(A, X^L) &= P \max_{a \in A} [\delta(a, X^L) \geq \varepsilon] \leq \\ &\leq P \sum_{a \in A} [\delta(a, X^L) \geq \varepsilon] = \sum_{a \in A} \underbrace{P[\delta(a, X^L) \geq \varepsilon]}_{Q_\varepsilon(a, X^L)}. \end{aligned}$$

Теория Вапника–Червоненкиса

Таким образом, доказали важную теорему:

Теорема

Для любых X^L , μ , конечного A и $\varepsilon \in [0, 1]$

$$\tilde{Q}_\varepsilon(A, X^L) \leq \sum_{a \in A} \mathcal{H}_L^{\ell, m} \left(\frac{\ell}{L} (m - \varepsilon k) \right),$$

где $m = n(a, X^L)$.

Следствие (Вапник и Червоненкис, 1968)

Для любых X^L , μ , конечного A и $\varepsilon \in [0, 1]$

$$\begin{aligned} \tilde{Q}_\varepsilon(A, X^L) &\leq |A| \cdot \max_m \mathcal{H}_L^{\ell, m} \left(\frac{\ell}{L} (m - \varepsilon k) \right) \leq \\ &\leq |A| \cdot \frac{3}{2} \exp(-\varepsilon^2 \ell), \quad \text{при } \ell = k. \end{aligned}$$

Обобщение на случай бесконечных семейств A

Функция роста $\Delta^A(L)$ семейства A — это максимальное по X^L число различных векторов ошибок $\mathbf{a} = (I(a, x_1), \dots, I(a, x_L))$.
В оценке надо заменить $|A|$ на функцию роста $\Delta^A(L)$.

Ёмкость (размерность Вапника–Червоненкиса) семейства A — это максимальная длина выборки h , для которой $\Delta^A(h) = 2^h$.

Теорема

Если такое h существует, то $\Delta^A(L) \leq C_L^0 + \dots + C_L^h \leq \frac{3}{2} \frac{L^h}{h!}$.

Теорема

Ёмкость семейства линейных классификаторов на два класса

$$a(x) = \text{sign}(w_1 x^1 + \dots + w_n x^n), \quad x = (x^1, \dots, x^n) \in X.$$

равна размерности пространства параметров, $\text{VCdim}(A) = n$.

Обращение оценки Вапника–Червоненкиса (при $\ell = k$)

1. Оценка: $P \left[\max_{a \in A} (\nu(a, X^k) - \nu(a, X^\ell)) \geq \varepsilon \right] \leq \Delta \frac{3}{2} \exp(-\ell \varepsilon^2).$

Тогда для любого $a \in A$ с вероятностью не менее $(1 - \eta)$

$$\nu(a, X^k) \leq \underbrace{\nu(a, X^\ell)}_{\text{эмпирический риск}} + \underbrace{\sqrt{\frac{1}{\ell} \ln \Delta + \frac{1}{\ell} \ln \frac{3}{2\eta}}}_{\text{штраф за сложность}}.$$

2. Оценка: $P \left[\max_{a \in A} (\nu(a, X^k) - \nu(a, X^\ell)) \geq \varepsilon \right] \leq \frac{3}{2} \frac{L^h}{h!} \cdot \frac{3}{2} \exp(-\ell \varepsilon^2).$

Тогда для любого $a \in A$ с вероятностью не менее $(1 - \eta)$

$$\nu(a, X^k) \leq \underbrace{\nu(a, X^\ell)}_{\text{эмпирический риск}} + \underbrace{\sqrt{\frac{h}{\ell} \ln \frac{2e\ell}{h} + \frac{1}{\ell} \ln \frac{9}{4\eta}}}_{\text{штраф за сложность}}.$$

Метод структурной минимизации риска (СМР)

Дано: система вложенных подсемейств возрастающей ёмкости

$$A_0 \subset A_1 \subset \dots \subset A_h \subset \dots$$

Найти: оптимальную ёмкость h^* , такую, что

$$\nu(a, X^k) \leq \underbrace{\min_{a \in A_h} \nu(a, X^\ell)}_{\text{минимизация эмпирического риска}} + \underbrace{\sqrt{\frac{h}{\ell} \ln \frac{2e\ell}{h} + \frac{1}{\ell} \ln \frac{9}{4\eta}}}_{\text{штраф за сложность}} \rightarrow \min_h$$

Недостатки СМР:

- h^* может оказаться заниженной из-за завышенности Q_ε .
- На практике эмпирический CV предпочтительнее этих оценок.

Причины завышенности оценок Вапника-Червоненкиса

- Оценка равномерного отклонения сильно завышена, когда большая часть алгоритмов имеет исчезающе малую вероятность быть результатом обучения.

На практике распределение

$$p(a) = P[\mu(X^\ell) = a], \quad a \in A$$


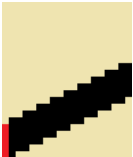


как правило, существенно неравномерно!

Будем называть это **эффектом расслоения семейства A** .

- **Неравенство Буля** сильно завышено, когда среди бинарных векторов ошибок есть много похожих.

Будем называть это **эффектом сходства алгоритмов**.

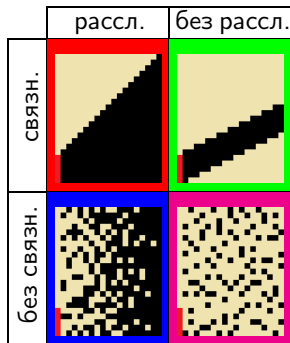
Эксперимент: четыре семейства алгоритмов, заданных матрицами ошибок; лучший алгоритм у всех одинаков

	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)		

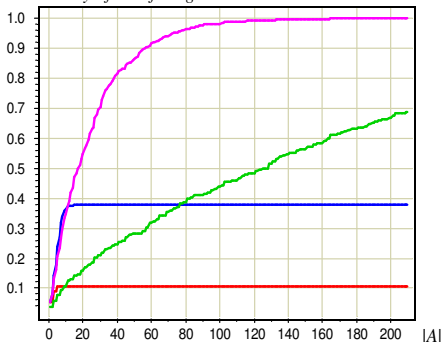
Постепенно добавляя алгоритмы в $\{a_1, \dots, a_D\}$, построим зависимости вероятности переобучения Q_ϵ от числа D .

Эксперимент с «трижды испорченной» монотонной цепью

$\ell = k = 100$, $\varepsilon = 0.05$, $N = 1000$ разбиений Монте-Карло.



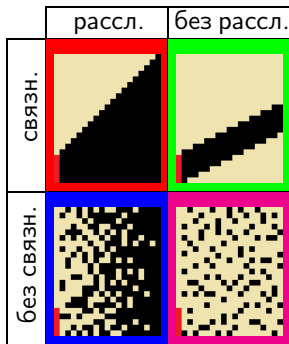
Probability of overfitting



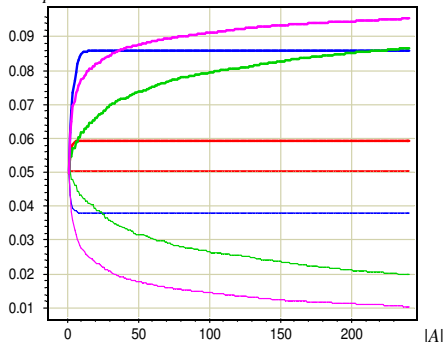
- Огромные семейства с P&C могут почти не переобучаться
- Без P&C даже 30 алгоритмов могут сильно переобучаться

Эксперимент с «трижды испорченной» монотонной цепью

$\ell = k = 100$, $\varepsilon = 0.05$, $N = 1000$ разбиений Монте-Карло.



Complete Cross-Validation



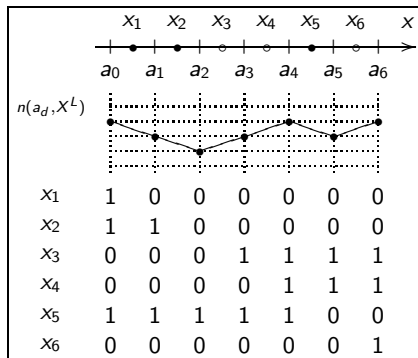
$E \nu(\mu(X^\ell), X^k)$ — CCV, жирные линии
 $E \nu(\mu(X^\ell), X^\ell)$ — тонкие линии

Цель, порождаемая семейством пороговых классификаторов

Пусть $f(x)$ — вещественный признак,

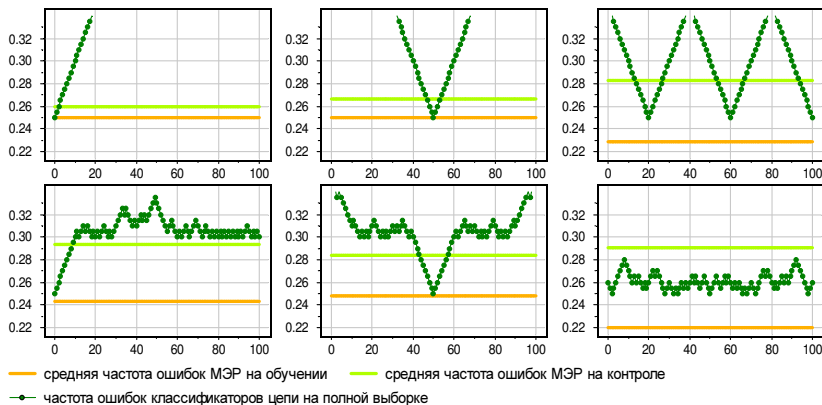
$$A = \{a(x) = [f(x) \geq \theta] : \theta \in \mathbb{R}\}.$$

Пример:



Эксперимент. Переобучение цепей с различным расслоением

Условия эксперимента: $L = 100$, $\ell = 50$, $m = 25$, $\varepsilon = 0.05$,
 метод Монте-Карло по $N = 100000$ случайных разбиений.



Граф расслоения–связности множества алгоритмов

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

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

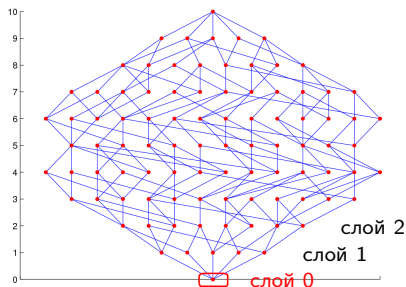
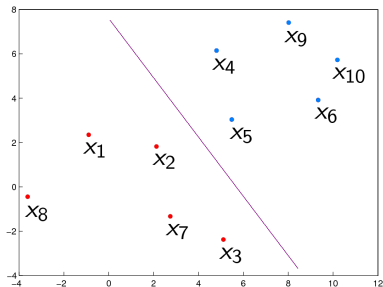
A — множество попарно различных векторов ошибок;

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

Свойства графа расслоения–связности:

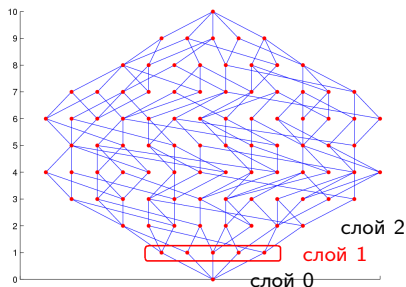
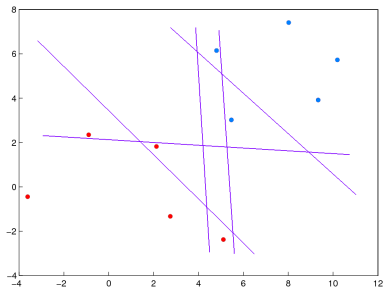
- это подграф графа Хассе отношения порядка \leq на A ;
- каждому ребру (a, b) соответствует объект $x_{ab} \in X^L$, такой, что $I(a, x_{ab}) = 0$, $I(b, x_{ab}) = 1$;
- граф является многодольным со слоями
 $A_m = \{a \in A : n(a, X^L) = m\}$, $m = 0, \dots, L$;

Пример 1. Семейство линейных алгоритмов классификации



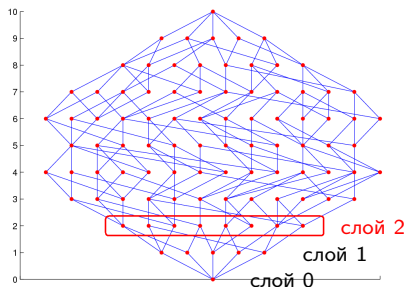
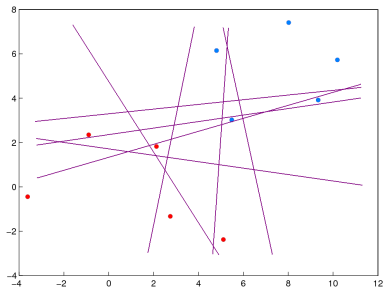
	слой 0
x ₁	0
x ₂	0
x ₃	0
x ₄	0
x ₅	0
x ₆	0
x ₇	0
x ₈	0
x ₉	0
x ₁₀	0

Пример 1. Семейство линейных алгоритмов классификации



	слой 0	слой 1				
x_1	0	1	0	0	0	0
x_2	0	0	1	0	0	0
x_3	0	0	0	1	0	0
x_4	0	0	0	0	1	0
x_5	0	0	0	0	0	1
x_6	0	0	0	0	0	0
x_7	0	0	0	0	0	0
x_8	0	0	0	0	0	0
x_9	0	0	0	0	0	0
x_{10}	0	0	0	0	0	0

Пример 1. Семейство линейных алгоритмов классификации



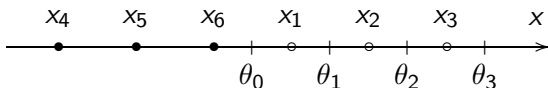
	слой 0	слой 1						слой 2								
X ₁	0	1	0	0	0	0	0	1	0	0	0	0	1	1	0	...
X ₂	0	0	1	0	0	0	0	1	1	0	0	0	0	0	0	...
X ₃	0	0	0	1	0	0	0	0	1	1	0	0	0	0	1	...
X ₄	0	0	0	0	1	0	0	0	0	1	1	0	0	0	0	...
X ₅	0	0	0	0	0	0	1	0	0	0	1	1	0	0	0	...
X ₆	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	...
X ₇	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	...
X ₈	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
X ₉	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...
X ₁₀	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	...

Пример 2. Монотонная цепь

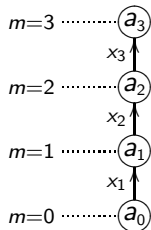
Опр. Монотонная цепь алгоритмов: $a_0 \prec a_1 \prec \dots \prec a_D$.

Пример: 1D пороговый классификатор $a_d(x) = [x - \theta_d]$;

2 класса $\{\bullet, \circ\}$
 6 объектов



Граф семейства:



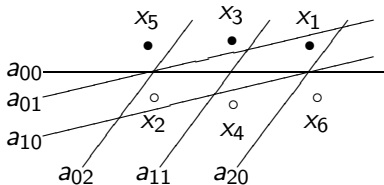
Матрица ошибок:

	a_0	a_1	a_2	a_3
x_1	0	1	1	1
x_2	0	0	1	1
x_3	0	0	0	1
x_4	0	0	0	0
x_5	0	0	0	0
x_6	0	0	0	0

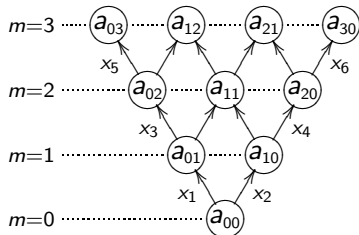
Пример 3. Двумерная сеть классификаторов

Пример:

2D линейный классификатор,
 2 класса {●, ○},
 6 объектов



Граф семейства:



Матрица ошибок:

	a_{00}	a_{01}	a_{10}	a_{02}	a_{11}	a_{20}	a_{03}	a_{12}	a_{21}	a_{30}
x_1	0	1	0	1	1	0	1	1	1	0
x_2	0	0	1	0	1	1	0	1	1	1
x_3	0	0	0	1	0	0	1	1	0	0
x_4	0	0	0	0	0	1	0	0	1	1
x_5	0	0	0	0	0	0	1	0	0	0
x_6	0	0	0	0	0	0	0	0	0	1

Порождающее и запрещающее множества алгоритмов

Определение

Верхняя связность $u(a)$ алгоритма a — это число всех рёбер, исходящих из вершины a :

$$u(a) = |X_a|, \quad X_a = \{x_{ab} \in X^L \mid a < b\};$$

X_a называется *порождающим множеством* алгоритма a .

Определение

Неполноценность $q(a)$ алгоритма a — это число различных объектов, соответствующих всем рёбрам на путях, ведущих в a :

$$q(a) = |X'_a|, \quad X'_a = \{x \in X^L \mid \exists b \in A: b < a, I(b, x) < I(a, x)\};$$

X'_a называется *запрещающим множеством* алгоритма a .

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

Верхняя связность $u(a) = \#\{x_{ab} \in X^L \mid a \prec b\}$

Нижняя связность $d(a) = \#\{x_{ba} \in X^L \mid b \prec a\}$

Неполноценность $q(a) = \#\{x \in X^L \mid \exists b \in A: b \prec a, l(b, x) < l(a, x)\}$

Число ошибок $m(a) = n(a, X^L)$.

Утв.

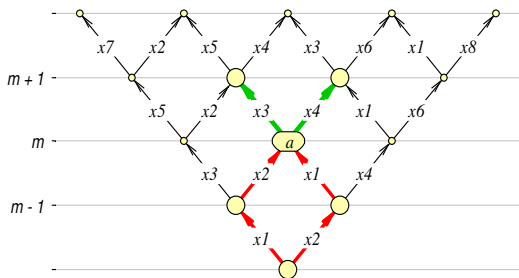
$d(a) \leq q(a) \leq m(a)$

Пример: двумерная
 сеть алгоритмов

$u(a) = \#\{x_3, x_4\} = 2$

$d(a) = \#\{x_1, x_2\} = 2$

$q(a) = \#\{x_1, x_2\} = 2$



Теорема о порождающих и запрещающих множествах

Если $[\mu(X^\ell) = a] = [X_a \subseteq X^\ell][X'_a \subseteq X^k]$, то

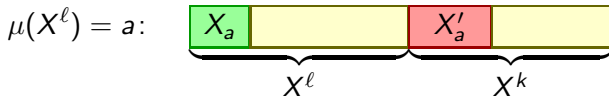
$$Q_\varepsilon(\mu, X^L) = \sum_{a \in A} P_a \mathcal{H}_{L_a}^{\ell_a, m_a}(s_a(\varepsilon)), \quad P_a = \mathbb{P}[\mu(X^\ell) = a] = \frac{C_{L_a}^{\ell_a}}{C_L^\ell},$$

$$L_a = L - |X_a| - |X'_a|,$$

$$\ell_a = \ell - |X_a|,$$

$$m_a = n(a, X^L \setminus X_a \setminus X'_a);$$

$$s_a(\varepsilon) = \frac{\ell}{L} (n(a, X^L) - \varepsilon k) - n(a, X_a).$$



Монотонные методы обучения

Опр. Метод обучения μ называется *монотонным*, если

$$\mu(X^\ell) \in A(X^\ell) = \operatorname{Arg} \min_{a \in A} K(a, X^\ell),$$

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

$$\forall X \subset X^L, \forall a, b \in A \text{ если } a < b, \text{ то } K(a, X) < K(b, X).$$

Опр. Метод μ называется *пессимистичным*, если

$$\mu(X^\ell) = \arg \max_{a \in A(X^\ell)} \delta(a, X^\ell).$$

Лемма

Если метод обучения μ монотонный и пессимистичный, то

$$[\mu(X^\ell) = a] \leq [X_a \subseteq X^\ell] [X'_a \subseteq X^k].$$

Верхняя оценка расслоения–связности

Теорема (Воронцов, Ивахненко, Решетняк, 2010)

Для любого монотонного метода μ , любых X^L , A и $\varepsilon \in (0, 1)$

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

где $u = |X_a|$ — верхняя связность алгоритма a ,
 $q = |X'_a|$ — неполноценность алгоритма a ,
 $m = n(a, X^L)$ — 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:

Идея доказательства

1. Пусть μ — произвольный монотонный метод обучения, $\bar{\mu}$ — монотонный пессимистичный метод обучения. Тогда

$$Q_\varepsilon(\mu, X^L) \leq Q_\varepsilon(\bar{\mu}, X^L).$$

2. Если $\bar{\mu}(X^\ell) = a$, то $\begin{cases} X_a \subseteq X^\ell & \text{в силу пессимистичности } \bar{\mu}, \\ X'_a \subseteq X^k & \text{в силу монотонности } \bar{\mu}. \end{cases}$

$$3. P[\bar{\mu}(X^\ell) = a] \leq P[\underbrace{X_a \subseteq X^\ell \text{ и } X'_a \subseteq X^k}_{S(a, X^\ell)}] = \frac{C_{L-|X_a|-|X'_a|}^{\ell-|X_a|}}{C_L^\ell} = \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell}.$$

4. По формуле полной вероятности:

$$Q_\varepsilon(\bar{\mu}, X^L) = \sum_{a \in A} \underbrace{P[S(a, X^\ell)]}_{C_{L-u-q}^{\ell-u} / C_L^\ell} \cdot \underbrace{P[\delta(a, X^\ell) \geq \varepsilon \mid S(a, X^\ell)]}_{\mathcal{H}_{L-u-q}^{\ell-u, m-q} \left(\frac{\ell}{L} (m - \varepsilon k) \right)}. \quad \blacksquare$$

Свойства оценки

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

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

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

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

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

- 6 При $|A| = 1$ это вероятность большого отклонения частот в двух выборках:

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

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

$$Q_\varepsilon \leq \sum_{a \in A} \mathcal{H}_L^{\ell, m} \left(\frac{\ell}{L} (m - \varepsilon k) \right) \leq |A| \cdot \frac{3}{2} \exp(-\varepsilon \ell^2).$$

- 8 При замене неполноценности q на нижнюю связность d это верхняя оценка функционала равномерного отклонения

$$\tilde{Q}_\varepsilon(A, X^L) = P \left[\sup_{a \in A} (\nu(a, X^k) - \nu(a, X^\ell)) \geq \varepsilon \right],$$

который учитывает связность, но не учитывает расслоение.

Верхние оценки средней частоты ошибок на контроле

Теорема

Для любого монотонного метода μ , любых X^L и A

$$\text{CCV}(\mu, X^L) \leq \sum_{a \in A} \frac{C_{L-u-q}^{\ell-u}}{C_L^\ell} \left(\frac{m}{k} - \frac{(m-q)(\ell-u)}{k(L-u-q)} \right).$$

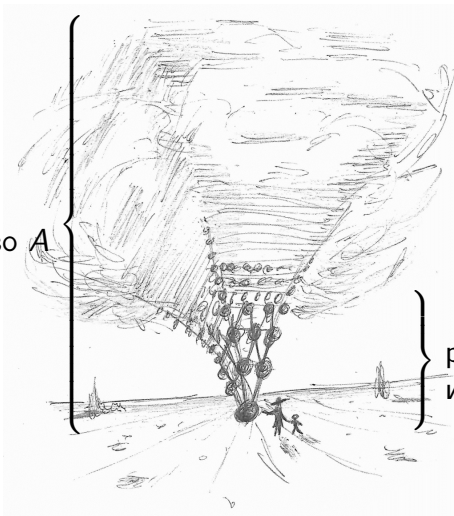
где $u = |X_a|$ — верхняя связность алгоритма a ,
 $q = |X'_a|$ — неполноценность алгоритма a ,
 $m = n(a, X^L)$ — число ошибок алгоритма a .

Преимущество:

оценка CCV вычисляется намного проще, чем оценки Q_ε и R_ε .

Идея использования оценок расслоения–связности

всё семейство A



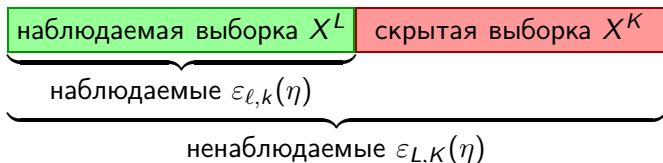
реально
используемая часть A

Резюме

- Свойства расслоения и связности уменьшают переобучение.
- На практике семейства, как правило, ими обладают. Иначе вероятность переобучения была бы близка к 1 уже при $|A|$ порядка нескольких десятков.
- Практическое применение комбинаторных оценок:
 - 1) оценить $\eta = Q_\varepsilon(\mu, X^L)$ по нескольким нижним слоям;
 - 2) применив обращение, оценить ε через η ;
 - 3) использовать оценку $\nu(X^L) + \varepsilon(\eta)$ как внешний критерий для выбора модели, метода или отбора признаков.
- Информацию о нижних слоях можно получать в ходе перебора порогов [Ивахненко 2010, Ботов 2011], случайным блужданием по A [Соколов 2013], в ходе итераций метода обучения,...

Нерешённые проблемы комбинаторной теории переобучения

- Проблема вычислительной неэффективности.
На больших выборках оценки расслоения–связности избыточно детализируют структуру семейства алгоритмов.
- Обоснование переноса оценки на скрытую выборку.
Оценки вычисляются для разбиения $X^L = X^\ell \sqcup X^k$, но применяются к разбиению $X^L \sqcup X^K$ со скрытой X^K .



В экспериментах это не приводит к переобучению, оценки устойчивы к изменению длины выборки при $L > 60$.

Но теоретического обоснования пока нет.