

**МАШИННОЕ ОБУЧЕНИЕ
И АНАЛИЗ ДАННЫХ**
(Machine Learning and Data Mining)

Н. Ю. Золотых

<http://www.uic.unn.ru/~zny/ml>

Лекция 17

Теория машинного обучения

Истоки: В. П. Вапник, А. Я. Червоненкис [1971]

Лемма 17.1 Пусть A_1, A_2, \dots, A_q — некоторые события (зависимые или независимые), заданные на одном вероятностном пространстве, тогда

$$\Pr(A_1 \cup A_2 \cup \dots \cup A_q) \leq \Pr A_1 + \Pr A_2 + \dots + \Pr A_q.$$

Лемма 17.2 (С.Н. Бернштейн, Н. Chernoff, W. Hoeffding) Пусть $Z^{(1)}, Z^{(2)}, \dots, Z^{(N)}$ — независимые одинаково распределенные случайные величины.

$$\Pr \left\{ Z^{(i)} = 1 \right\} = \theta, \quad \Pr \left\{ Z^{(i)} = 0 \right\} = 1 - \theta$$

(схема Бернулли). Тогда

$$\Pr(|\hat{\theta} - \theta| > \gamma) \leq 2e^{-2\gamma^2 N},$$

где

$$\hat{\theta} = \frac{1}{N} \sum_{i=1}^N Z^{(i)}.$$

Рассмотрим задачу бинарной классификации.

\mathcal{F} — некоторое множество функций $f : \mathbf{R}^d \rightarrow \{0, 1\}$.

(X, Y) — $(d + 1)$ -мерная с.в. с неизвестной функцией распределения $P(x, y)$

Ожидаемый риск

$$R(f) = \mathbf{E} I(f(X) \neq Y) = \Pr I(f(X) \neq Y).$$

$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})$ — обучающая выборка.

Эмпирический риск (частота ошибки) на данной выборке есть

$$\widehat{R}(f) = \frac{1}{N} \sum_{i=1}^N I(f(x^{(i)}) \neq y^{(i)}).$$

Задача: оценить $R(f)$ в терминах $\widehat{R}(f)$.

Принцип минимизации эмпирического риска:

$$\widehat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \widehat{R}(f)$$

17.1. Случай конечного \mathcal{F}

Пусть $\mathcal{F} = \{f_1, f_2, \dots, f_q\}$.

Рассмотрим $f \in \mathcal{F}$.

Введем семейство случайных величин

$$Z^{(i)} = I(f(X^{(i)}) \neq Y^{(i)}) \quad (i = 1, 2, \dots, N).$$

$Z^{(i)}$ распределены одинаково и независимо — схема испытаний Бернулли.

По лемме 2

$$\Pr(|\widehat{R}(f) - R(f)| > \gamma) \leq 2e^{-2\gamma^2 N}$$

Откуда

$$\Pr(|\widehat{R}(f) - R(f)| \leq \gamma) \geq 1 - 2e^{-2\gamma^2 N}$$

Итак, для *конкретной* f эмпирический риск близок к ожидаемому с *большой* вероятностью.

Нас же интересует более сильное условие:

$$\Pr \left\{ \forall f \in \mathcal{F} : |\widehat{R}(f) - R(f)| \leq \gamma \right\}$$

$$\Pr \left\{ \exists f \in \mathcal{F} : |\widehat{R}(f) - R(f)| > \gamma \right\} = \Pr(A_1 \cup \dots \cup A_q) \leq \sum_{i=1}^q \Pr(A_i) \leq \sum_{i=1}^q 2e^{-2\gamma^2 N} = 2qe^{-2\gamma^2 N}$$

Мы доказали утверждение:

Утверждение 17.3 (О равномерной сходимости \widehat{R} к R в случае конечного \mathcal{F})

$$\Pr \left\{ \forall f \in \mathcal{F} : |\widehat{R}(f) - R(f)| \leq \gamma \right\} \geq 1 - 2qe^{-2\gamma^2 N}$$

Следствие 17.4 Если $N \geq \frac{1}{2\gamma^2} \ln \frac{2q}{\delta}$, тогда

$$\Pr \left\{ \forall f \in \mathcal{F} : |\widehat{R}(f) - R(f)| \leq \gamma \right\} \geq 1 - \delta.$$

Следствие 17.5

$$\Pr \left\{ \forall f \in \mathcal{F} : |\widehat{R}(f) - R(f)| \leq \sqrt{\frac{1}{N} \ln \frac{2q}{\delta}} \right\} \geq 1 - \delta$$

Следствие 17.6 (Принцип минимизации эмпирического риска для конечного \mathcal{F}) Пусть

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} R(f), \quad \hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}(f),$$

тогда

$$\Pr \left\{ R(\hat{f}) \leq R(f^*) + 2\sqrt{\frac{1}{N} \ln \frac{2q}{\delta}} \right\} \geq 1 - \delta.$$

ДОКАЗАТЕЛЬСТВО. Если $\forall f \in \mathcal{F}$

$$|\hat{R}(f) - R(f)| \leq \gamma = \sqrt{\frac{1}{N} \ln \frac{2q}{\delta}},$$

то

$$\begin{aligned} 0 \leq R(\hat{f}) - R(f^*) &= R(\hat{f}) - \hat{R}(\hat{f}) + \hat{R}(\hat{f}) - \hat{R}(f^*) + \hat{R}(f^*) - R(f^*) = \\ &= \underbrace{R(\hat{f}) - \hat{R}(\hat{f})}_{\leq \gamma} + \underbrace{\hat{R}(\hat{f}) - \hat{R}(f^*)}_{\leq 0} + \underbrace{\hat{R}(f^*) - R(f^*)}_{\leq \gamma} \leq 2\gamma. \end{aligned}$$

Поэтому

$$\begin{aligned} \Pr \left\{ 0 \leq R(\hat{f}) - R(f^*) \leq 2\gamma \right\} &\geq \Pr \left\{ R(\hat{f}) - \hat{R}(\hat{f}) \leq \gamma \text{ и } \hat{R}(f^*) - R(f^*) \leq \gamma \right\} \geq \\ &\geq \Pr \left\{ \forall f \in \mathcal{F} : |\hat{R}(f) - R(f)| \leq \gamma \right\} \geq 1 - \delta. \end{aligned}$$

■

17.2. Случай бесконечного \mathcal{F}

Коэффициентом разнообразия (*shatter coefficient*) $\Delta_0(\mathcal{F}, \mathbf{x})$ множества \mathcal{F} на выборке $\mathbf{x} = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$ – количество всевозможных способов разбить \mathbf{x} на 2 подмножества, т. е. количество всевозможных бинарных векторов вида

$$(f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(N)})),$$

порождаемых всеми функциями f из \mathcal{F} .

Функция роста:

$$\Delta(\mathcal{F}, N) = \max_{|\mathbf{x}|=N} \Delta_0(\mathcal{F}, \mathbf{x}).$$

Пример.

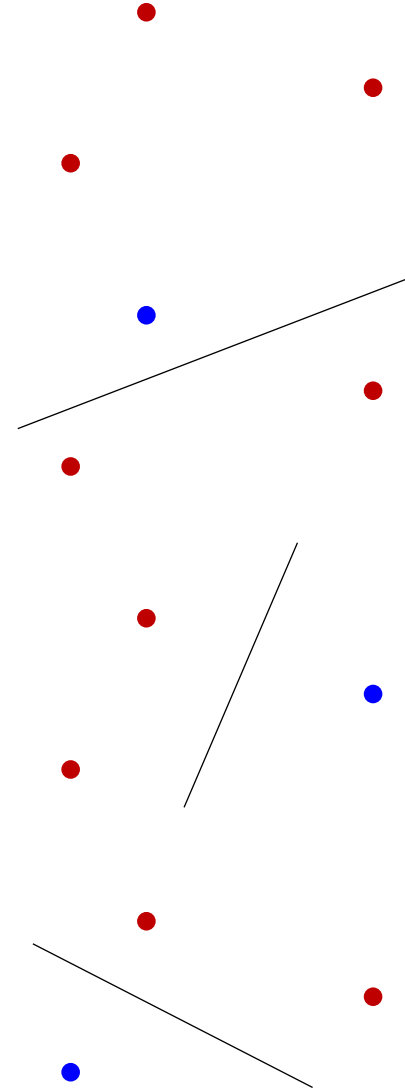
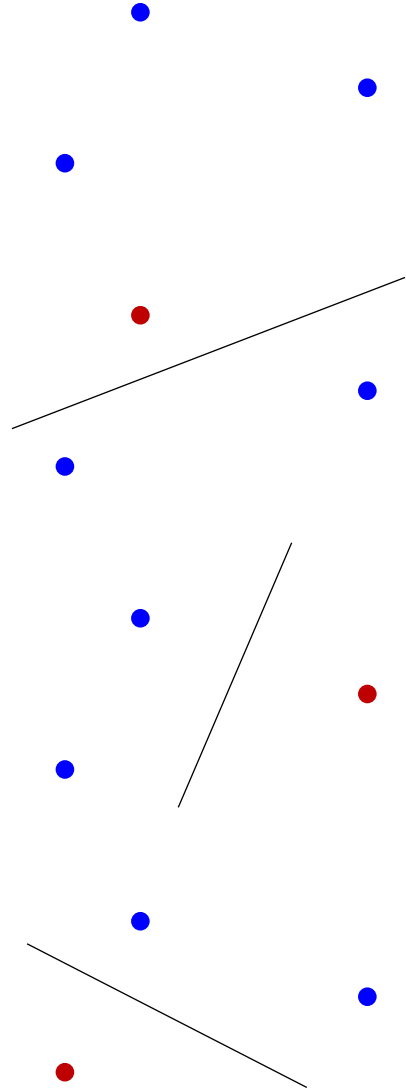
В качестве \mathcal{F} рассмотрим класс пороговых функций (линейных решающих правил)

$$\mathcal{F} = \text{TF}_d = \{f(x) = I(\beta_0 + \langle \beta, x \rangle \leq 0) : \beta_0 \in \mathbf{R}, \beta \in \mathbf{R}^d\}$$



$$\Delta_0(\text{TF}_d, \mathbf{x}) = 8$$

(все возможные разбиения)



Если 3 точки на прямой, то $\Delta_0(\text{TF}_d, \mathbf{x}) = 6$

Если N точек на прямой, то $\Delta_0(\text{TF}_d, \mathbf{x}) = 2N$

Если 4 точки (в общем положении) в \mathbf{R}^2 , то ...

$\Delta_0(\text{TF}_d, \mathbf{x}) = 14$, т. е. все без двух



Утверждение 17.7 *Функция роста класса пороговых функций равна*

$$\Delta(\text{TF}_d, N) = 2 \left(\binom{N-1}{0} + \binom{N-1}{1} + \dots + \binom{N-1}{d} \right).$$

Например,

$$\Delta(\text{TF}_2, 2) = 4, \Delta(\text{TF}_2, 3) = 8, \Delta(\text{TF}_2, 4) = 14, \Delta(\text{TF}_2, 5) = 22$$

$$\Delta(\text{TF}_3, 2) = 4, \Delta(\text{TF}_3, 3) = 8, \Delta(\text{TF}_3, 4) = 16, \Delta(\text{TF}_3, 5) = 30$$

Размерность Ванника–Червоненкиса, или емкость, $VC \mathcal{F}$ – наибольшая мощность множества в \mathcal{X} , разбиваемая (shatter) функциями из \mathcal{F} .

Говорят, что \mathcal{F} разбивает множество $\{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$, если функции из \mathcal{F} позволяют разделить это множество на два подмножества всеми 2^N возможными способами, т. е. для любого двоичного набора $y \in \{0, 1\}^N$ найдется $f \in \mathcal{F}$, такое, что $y = (f(x^{(1)}), f(x^{(2)}), \dots, f(x^{(N)}))$.

Эквивалентно: $VC \mathcal{F}$ – это максимальное N , при котором $\Delta(\mathcal{F}, N) = 2^N$

Если такого максимального N не существует, то $VC \mathcal{F} = \infty$

Чтобы доказать, что $VC \mathcal{F} = h$ достаточно проделать две вещи:

- 1) доказать $VC \mathcal{F} \geq h$; для этого достаточно построить $\{x^{(1)}, x^{(2)}, \dots, x^{(h)}\}$, разбиваемое функциями из \mathcal{F} ;
- 2) доказать $VC \mathcal{F} \leq h$; для этого достаточно показать, что никакое множество из $n + 1$ точек нельзя разбить функциями из \mathcal{F} .

Утверждение 17.8 *Размерность Вайника–Червоненкиса класса пороговых функций равна*

$$\text{VC}(\text{TF}_d) = d + 1.$$

Утверждение 17.9 *Если \mathcal{F} конечно, то*

$$\text{VC}(\mathcal{F}) \leq \log_2 |\mathcal{F}|.$$

Теорема 17.10 (Вапник–Червоненкис) (О равномерной сходимости \widehat{R} к R в случае конечной VC \mathcal{F})
Пусть \mathcal{F} задано и VC \mathcal{F} конечно, тогда

$$\Pr \left\{ \forall f \in \mathcal{F} : |\widehat{R}(f) - R(f)| \leq \gamma \right\} \geq 1 - \delta, \quad (1)$$

$$\Pr \left\{ R(\widehat{f}) \leq R(f^*) + 2\gamma \right\} \geq 1 - \delta. \quad (2)$$

где

$$\gamma = O \left(\sqrt{\frac{\text{VC } \mathcal{F}}{N} \ln \frac{N}{\text{VC } \mathcal{F}} + \frac{1}{N} \ln \frac{1}{\delta}} \right). \quad (3)$$

В П.Ч. неравенства $R(\widehat{f}) \leq R(f^*) + 2\gamma$ с ростом VC \mathcal{F} первое слагаемое убывает, а второе возрастает. Сумма достигает своего минимума при некотором оптимальном значении VC \mathcal{F} .

VC \mathcal{F} мало (модель с малой емкостью) – недообучение

VC \mathcal{F} велико (модель с крайне большой емкостью) – переобучение

Также можно доказать, что если VC $\mathcal{F} = \infty$, то равномерной сходимости нет.

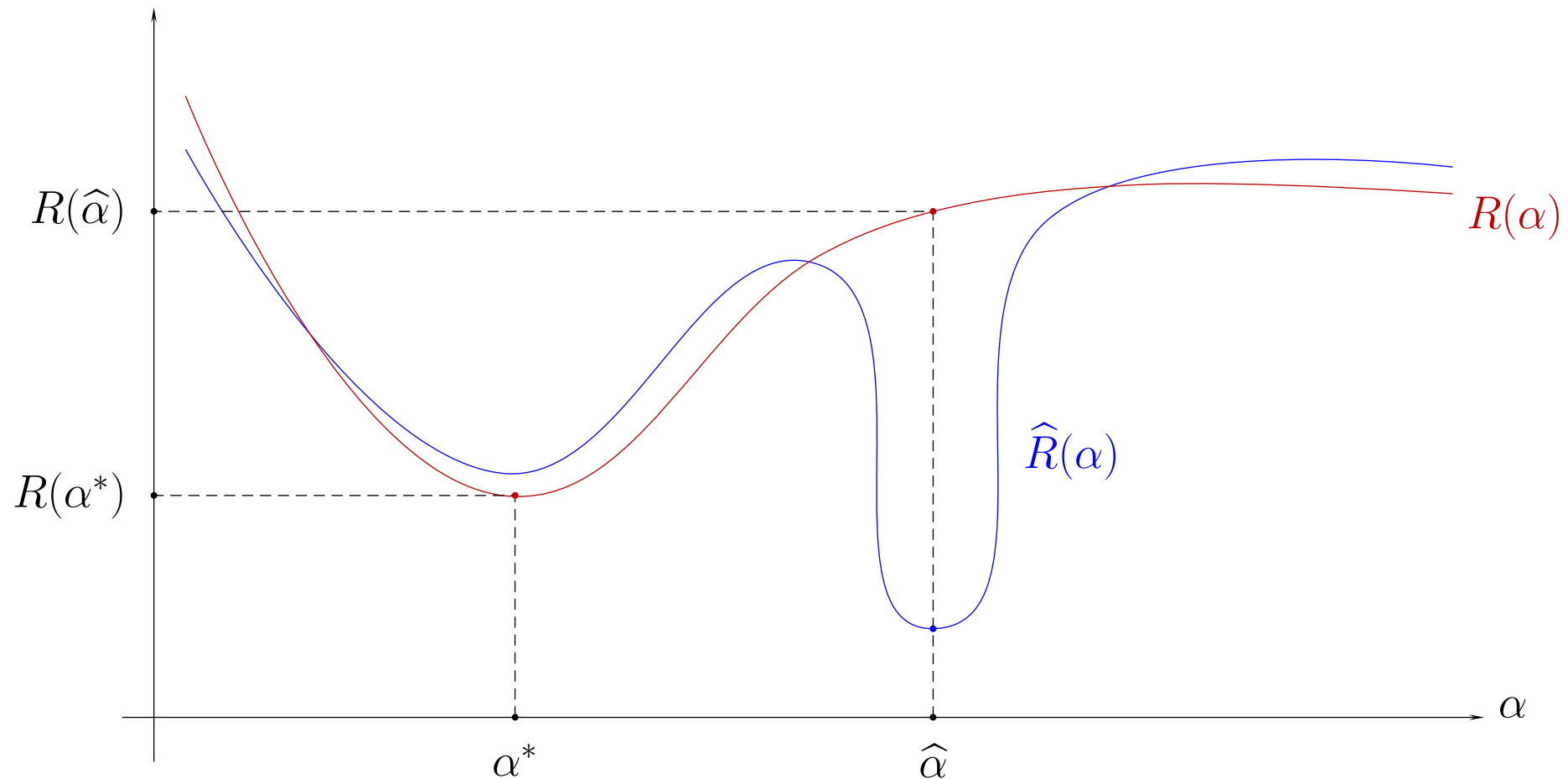
Как бороться с переобучением:

уменьшать VC \mathcal{F}

увеличивать N

Пусть $\mathcal{F} = \{f : f(x, \alpha), \alpha \in [0, 1]\}$ — класс решающих правил

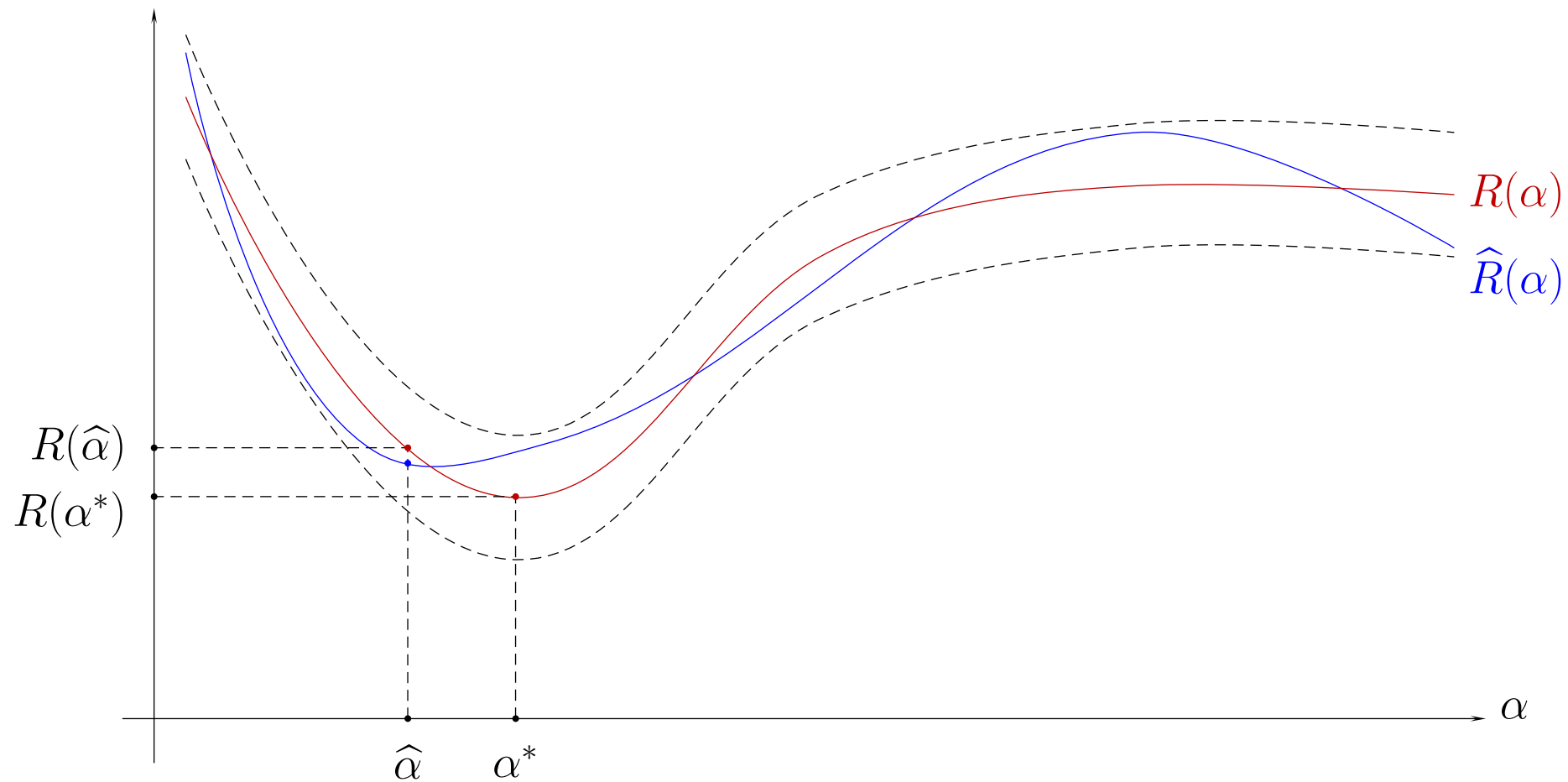
$R(\alpha)$ — средний риск, $\widehat{R}(\alpha)$ — эмпирический риск на функции $f(x, \alpha)$



$R(\widehat{\alpha})$ далеко от минимального значения $R(\alpha^*)$, хотя

$$\lim_{N \rightarrow \infty} \Pr \left\{ |\widehat{R}(f) - R(f)| \leq \gamma \right\} = 1.$$

Равномерная сходимость



$$\lim_{N \rightarrow \infty} \Pr \left\{ \forall f \in \mathcal{F} : |\hat{R}(f) - R(f)| \leq \gamma \right\} = 1.$$

Принцип *структурной минимизации риска* предназначен для отыскания оптимальной емкости.

Пусть в \mathcal{F} выделена некоторая цепочка подсемейств

$$\mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots \subset \mathcal{F}_M = \mathcal{F}.$$

Необходимо выбрать семейство с минимальным значением верхней оценки в (2).

К сожалению, оценки (1)–(3) являются слишком завышенными, чтобы их можно было использовать для практических вычислений. Однако они представляют большой интерес с теоретической точки зрения и являются стимулом для дальнейших исследований в поиске подобных оценок, в частности в теории вероятностно–приближенно–корректного обучения (probably approximately correct learning, PAC–learning).

Изложенные здесь результаты распространяются на задачу классификации с бóльшим числом классов и на задачу восстановления регрессии.

Пусть \mathcal{F} – некоторый класс функций $f : \mathcal{X} \rightarrow \mathcal{Y}$. Размерностью Вапника–Червоненкиса для класса \mathcal{F} называется $VC \mathcal{F} = VC \mathcal{F}'$, где

$$\mathcal{F}' = \{I(f(x) - y) : f \in \mathcal{F}, y \in \mathcal{Y}\}.$$

17.2.1.* Функция роста для линейных классификаторов

Система из N неравенств с неизвестными $\beta_0, \beta_1, \dots, \beta_d$:

$$\begin{cases} \beta_0 + \langle \beta, x \rangle \leq 0, & \text{если } f(x) = 0 \\ \beta_0 + \langle \beta, x \rangle > 0, & \text{если } f(x) = 1 \end{cases}$$

Не нарушая общности, можем считать, что точка 0 — один из объектов.

Множество всех функций разобьется на два равномоощных класса:

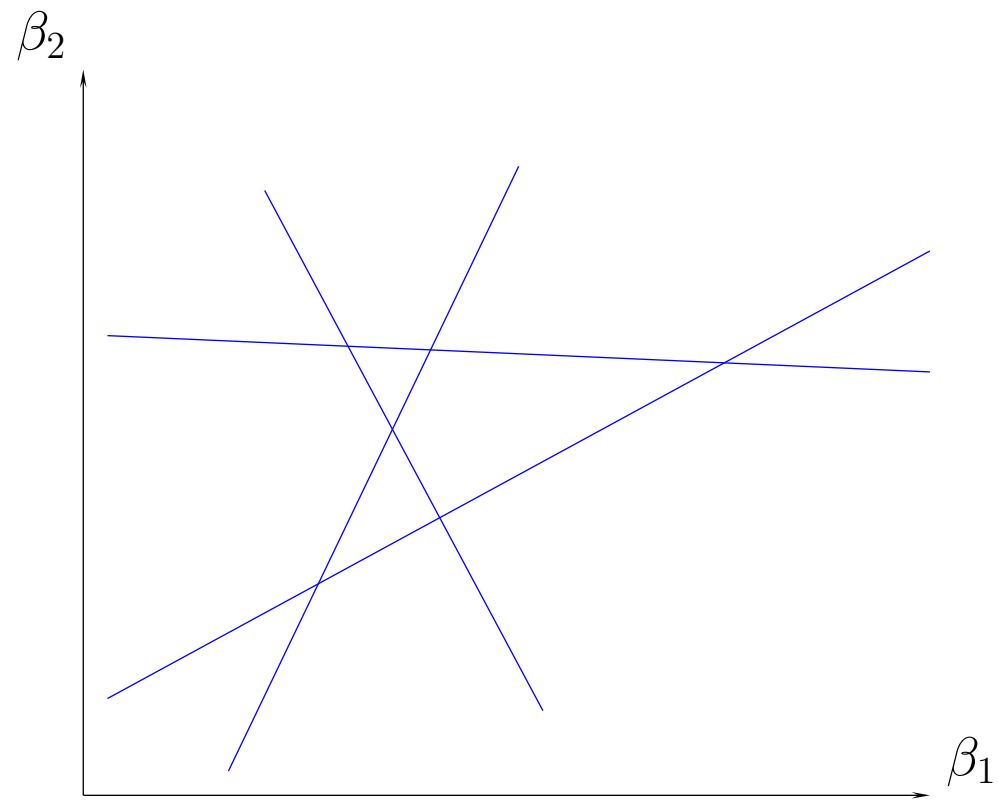
для которых $f(0) = 0$

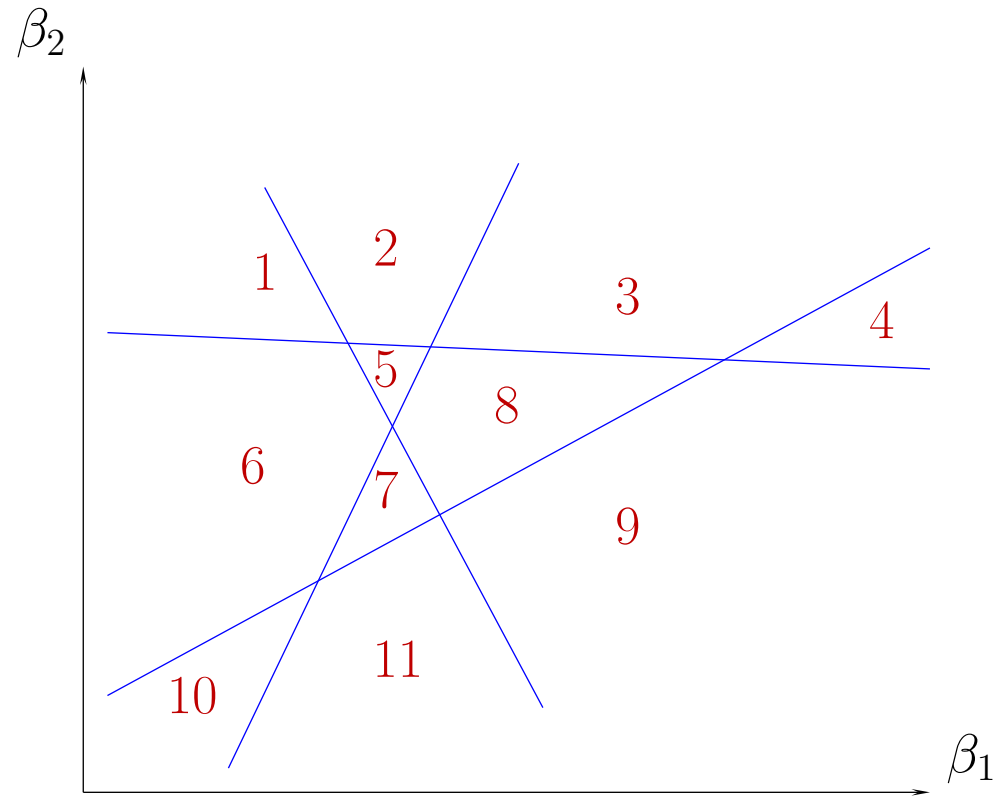
и для которых $f(0) = 1$

Будем считать только первые, а затем результат умножим на 2.

Для них можно считать, что $\beta_0 = -1$:

$$\begin{cases} -1 + \langle \beta, x \rangle \leq 0, & \text{если } f(x) = 0 \\ -1 + \langle \beta, x \rangle > 0, & \text{если } f(x) = 1 \end{cases}$$





Найдем максимальное количество $\Gamma(d, N)$ областей, на которые можно разбить \mathbf{R}^d N гиперплоскостями.
 Справедливо рекуррентное соотношение

$$\Gamma(d, N) = \Gamma(d, N - 1) + \Gamma(d - 1, N - 1),$$

откуда

$$\Gamma(d, N) = \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{d}.$$

Теорема 17.11 (Лемма Зауэра; Sauer, 1972; Shelah, 1972) Пусть класс \mathcal{F} имеет конечную размерность Вайника–Червоненкиса $VC(\mathcal{F}) = h$, тогда

$$\Delta(\mathcal{F}, N) \leq \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{h} < \left(\frac{eN}{h}\right)^h$$