

Разбор летучки

Лекция 7

Выбор моделей

Екатерина Тузова

Задача выбора метода обучения

X - множество объектов

Y - множество классов

Обучающая выборка: $X^l = (x_i, y_i)_{i=1}^l$

Целевая функция: $f : X \rightarrow Y$

Набор моделей алгоритмов $A_t : X \rightarrow Y, t \in T$

Методы обучения $\mu_t : X \times Y \rightarrow A_t, t \in T$

Задача: Найти алгоритм $a \in A_t$ с наилучшей обобщающей способностью

Модель

Неизвестная целевая функция

$$f : X \rightarrow Y$$



Обучающая выборка

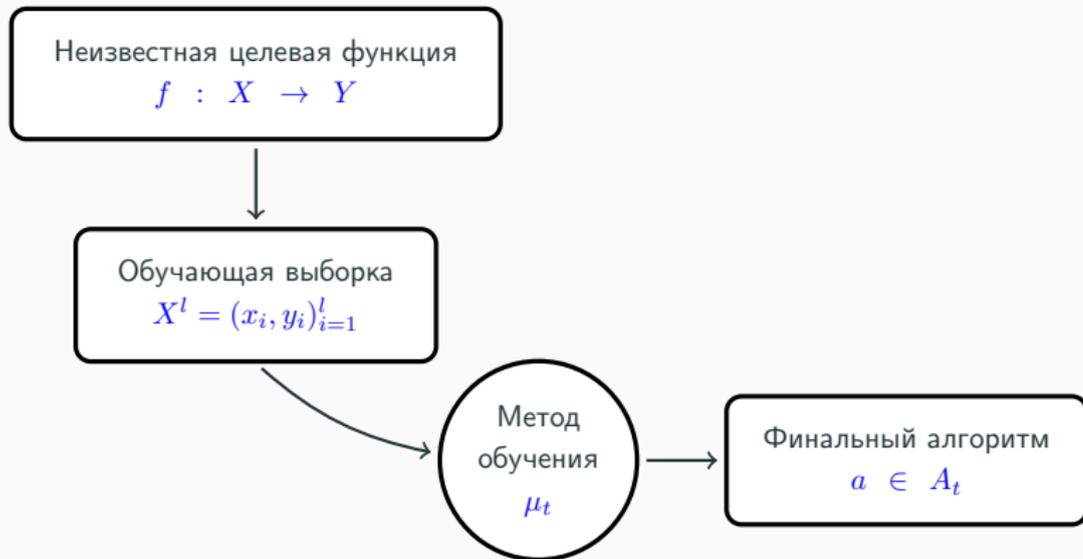
$$X^l = (x_i, y_i)_{i=1}^l$$



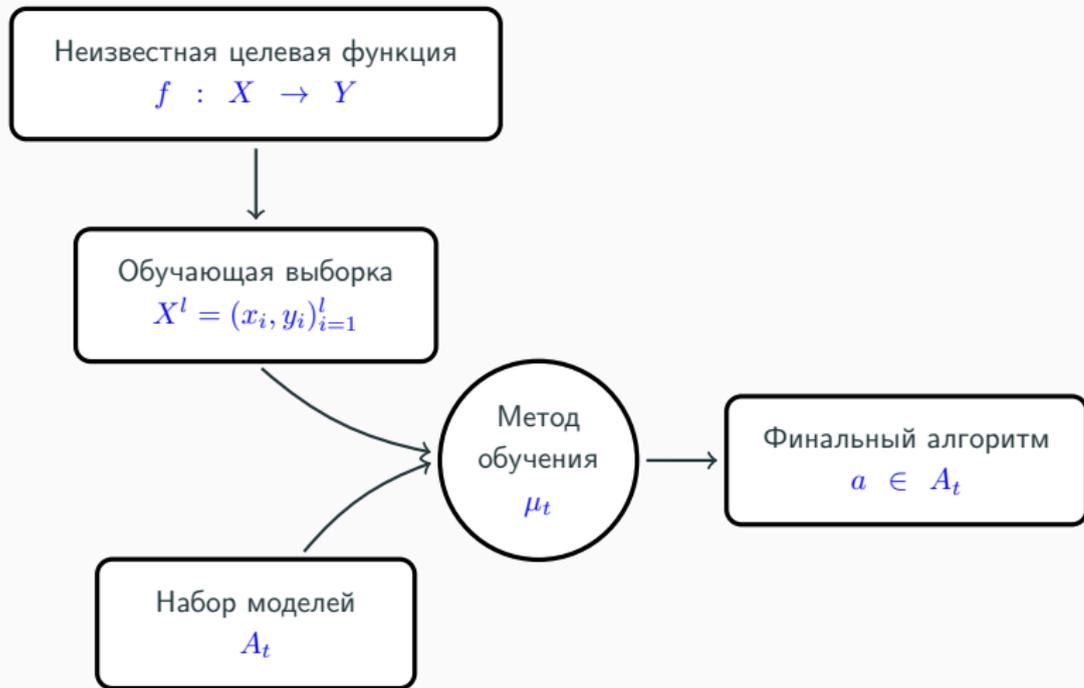
Финальный алгоритм

$$a \in A_t$$

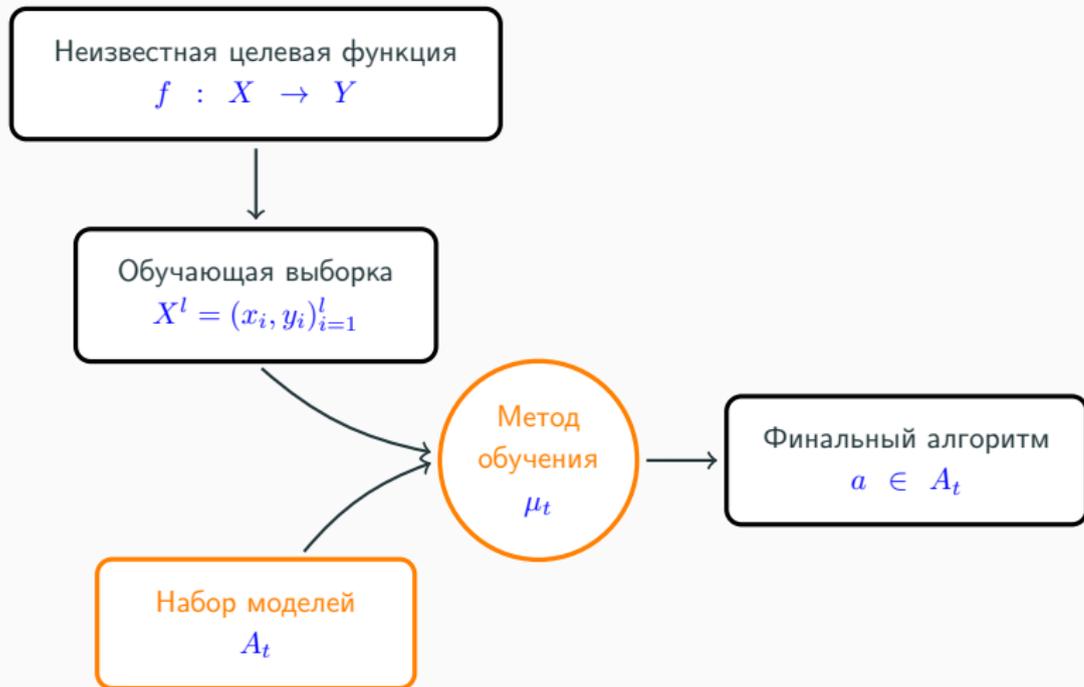
Модель



Модель



Модель



Пример

Набор моделей:

$$a(\mathbf{x}) = \text{sign}\left(\sum_{j=1}^n w_j x^j - w_0\right)$$

Набор моделей:

$$a(\mathbf{x}) = \text{sign}\left(\sum_{j=1}^n w_j x^j - w_0\right)$$

Набор моделей:

$$a(\mathbf{x}) = \text{sign}\left(\sum_{j=1}^n w_j x^j - w_0\right)$$

Метод обучения:

```
1 function PERCEPTRON( $X^l$ )
2   Инициализировать  $w_0, \dots, w_n$ 
3   repeat[пока  $\mathbf{w}$  изменяются]
4     for  $i = 1, \dots, l$  do
5       if  $a(x_i) \neq y_i$  then
6          $\mathbf{w} = \mathbf{w} + y_i \mathbf{x}_i$ 
```

Научиться оценивать метод обучения и обобщающую способность алгоритма

Функция потерь $\mathcal{L}(a, x_i)$ – характеризует величину ошибки алгоритма a на объекте x_i .

Если $\mathcal{L}(a, x_i) = 0$, то ответ $a(x_i)$ называется корректным.

Задача классификации:

$$\mathcal{L}(a, x_i) = [a(x_i) \neq y_i] - \text{индикатор ошибки}$$

Задача регрессии:

$$\mathcal{L}(a, x_i) = |a(x_i) - y_i| - \text{абсолютное значение ошибки}$$

$$\mathcal{L}(a, x_i) = (a(x_i) - y_i)^2 - \text{квадратичная ошибка}$$

Функционал качества алгоритма a на выборке X^l :

$$Q(a, X^l) = \frac{1}{l} \sum_{i=1}^l \mathcal{L}(a, x_i)$$

Минимизация эмпирического риска:

$$\arg \min_{A_t} Q(a, X^l)$$

$$Q_{\mu}(X^l) = Q(\mu(X^l), X^l)$$

Этот функционал оценивает качество обучения на выборке X^l

$$Q_{\mu}(X^l) = Q(\mu(X^l), X^l)$$

Этот функционал оценивает качество обучения на выборке X^l

Какая с этим может быть проблема?

$$E_{in} = Q_{\mu}(X^l) = Q(\mu(X^l), X^l)$$

Внешний функционал по отложенной выборке:

$$E_{out} = Q_{\mu}(X^t, X^k) = Q(\mu(X^t), X^k)$$

$$E_{in} = Q_{\mu}(X^l) = Q(\mu(X^l), X^l)$$

Внешний функционал по отложенной выборке:

$$E_{out} = Q_{\mu}(X^t, X^k) = Q(\mu(X^t), X^k)$$

Какой здесь недостаток?

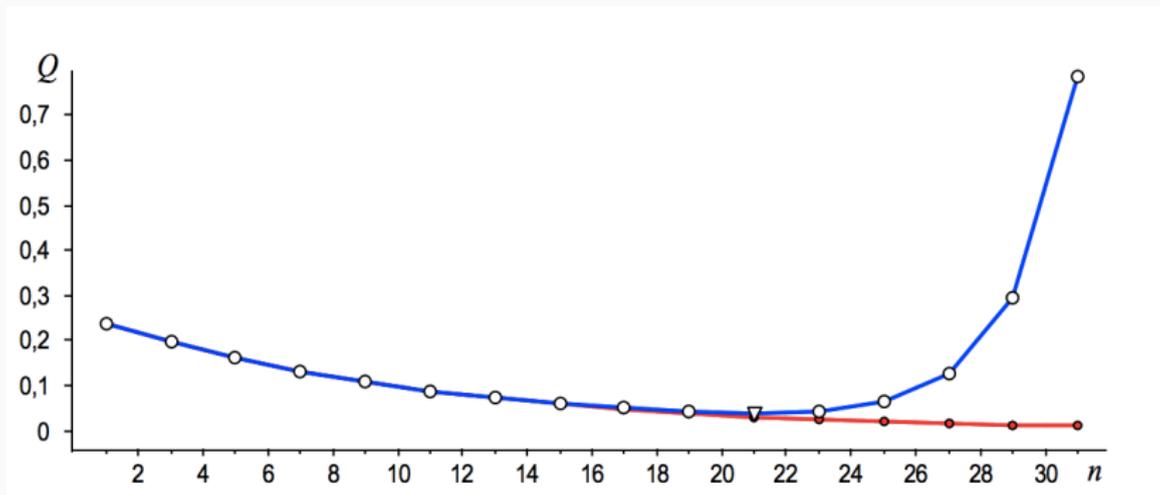
$$E_{in} = Q_{\mu}(X^l) = Q(\mu(X^l), X^l)$$

Внешний функционал по отложенной выборке:

$$E_{out} = Q_{\mu}(X^t, X^k) = Q(\mu(X^t), X^k)$$

Сильная зависимость от разбиения $X^l = X^t \sqcup X^k$

Кривая обучения



Идея: Усреднить по всем C_l^t выборкам $X^l = X^t \sqcup X^k$

$$CCV(\mu, X^l) = \frac{1}{C_l^t} \sum_{X^t} Q_\mu(X^t, X^k)$$

Идея: Усреднить по всем C_l^t выборкам $X^l = X^t \sqcup X^k$

$$CCV(\mu, X^l) = \frac{1}{C_l^t} \sum_{X^t} Q_\mu(X^t, X^k)$$

Во что превратится оценка при $k = 1$?

- Оценка вычислительно слишком сложна
- Не учитывает дисперсию X^k

Идея: Возьмём случайное разбиение $X^l = X_1 \sqcup \dots \sqcup X_k$ на k блоков равной длины.

$$CV_k(\mu, X^l) = \frac{1}{k} \sum_{i=1}^k Q_\mu(X^l \setminus X_i, X_i)$$

Идея: Возьмём случайное разбиение $X^l = X_1 \sqcup \dots \sqcup X_k$ на k блоков равной длины.

$$CV_k(\mu, X^l) = \frac{1}{k} \sum_{i=1}^k Q_\mu(X^l \setminus X_i, X_i)$$

Недостатки:

- Оценка зависит от разбиения на блоки
- Каждый объект только один раз участвует в контроле

Идея: Выборка разбивается t раз случайным образом на k блоков

$$CV_{tk}(\mu, X^l) = \frac{1}{t} \sum_{j=1}^t \frac{1}{k} \sum_{i=1}^k Q_{\mu}(X^l \setminus X_{ji}, X_{ji})$$

Идея: Выборка разбивается t раз случайным образом на k блоков

$$CV_{tk}(\mu, X^l) = \frac{1}{t} \sum_{j=1}^t \frac{1}{k} \sum_{i=1}^k Q_{\mu}(X^l \setminus X_{ji}, X_{ji})$$

- + Выбором t можно улучшать точность оценки
- + Каждый объект участвует в контроле t раз

Критерий непротиворечивости моделей

Идея: Если модель верна, то алгоритмы, настроенные по разным частям данных, не должны противоречить друг другу.

Аналитический подход

1. Получить верхнюю оценку вероятности переобучения R_ε

$$R_\varepsilon(\mu, X^l) = P[E_{out} - E_{in} \geq \varepsilon] \leq \eta(\varepsilon, A)$$

2. Тогда для любых X^l , A , μ и η с вероятностью не менее $1 - \eta$ справедливо

$$E_{out} \leq E_{in} + \varepsilon(\eta, A)$$

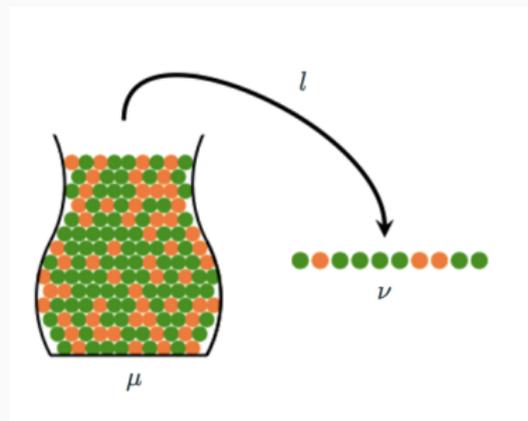
3. Будем оптимизировать

$$E_{in} + \varepsilon(\eta, A) \rightarrow \min_{\mu}$$

Регуляризатор ε — аддитивная добавка к внутреннему критерию, обычно штраф за сложность модели A .

Неравенство Бернштейна-Хёфдинга

$$P[|\nu - \mu| < \varepsilon] \leq 2e^{-2\varepsilon^2 l}$$



ν – доля оранжевых шаров в выборке размера l

μ – истинная доля оранжевых шаров в данных

Какое отношение это имеет к нашим гипотезам?

Каждый шар это объект x из пространства X .
Неизвестная целевая функция f .

Оранжевый шар – модель h верна ($h(x) = f(x)$)

Зелёный шар – модель h не верна ($h(x) \neq f(x)$)

По выборке X^l можем оценить ν – долю объектов, на которых верна модель.

$E_{in}(h) = \nu$ - доля объектов в выборке X^l , на которых h верна
 $E_{out}(h) = \mu$ - доля объектов во всём множестве X , на которых h верна

$$P[|E_{in}(h) - E_{out}(h)| < \varepsilon] \leq 2e^{-2\varepsilon^2 l}$$

Неравенство выполняется для каждой модели.

Какова вероятность, что модель $a \in \mathcal{H}$, наилучшим образом приближающая f по выборке, наилучшим образом приближает f на всём множестве?

С какой вероятностью монета, подброшенная 10 раз, выпадет одной и той же стороной все 10 раз?

С какой вероятностью монета, подброшенная 10 раз, выпадет одной и той же стороной все 10 раз?

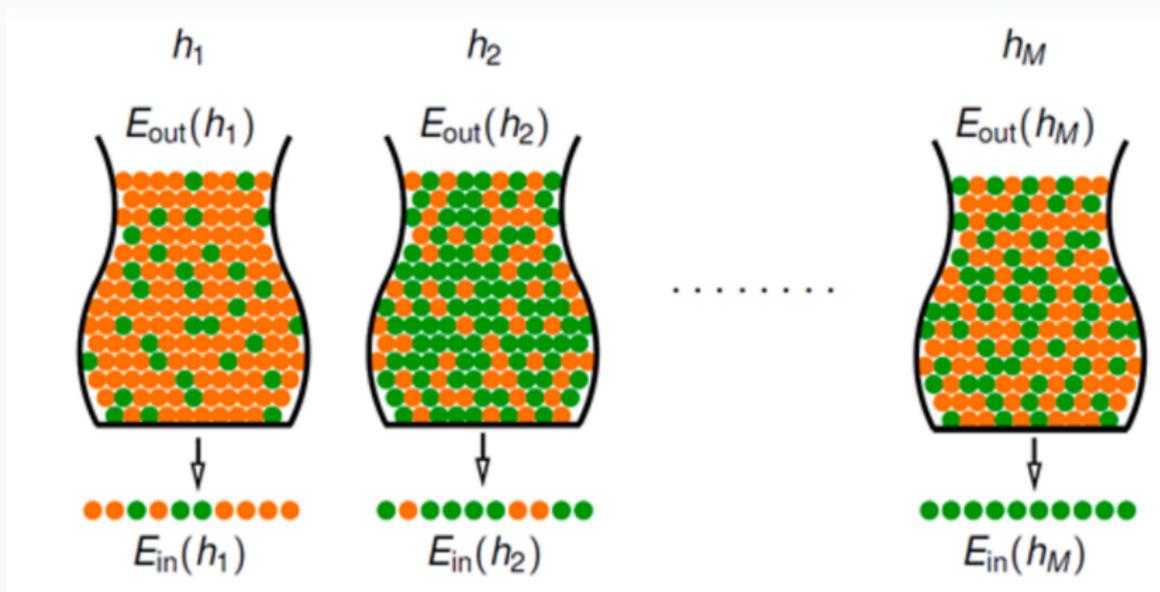
0.001

С какой вероятностью одна из 1000 монет, каждая из которых подброшена 10 раз, выпадет одной и той же стороной все 10 раз?

С какой вероятностью одна из 1000 монет, каждая из которых подброшена 10 раз, выпадет одной и той же стороной все 10 раз?

0.63

К нашей задаче



На h_M модели наблюдаем переобучение.

$$\begin{aligned} P[|E_{in}(a) - E_{out}(a)| < \varepsilon] &\leq \sum_{m=1}^M P[|E_{in}(h) - E_{out}(h)| < \varepsilon] \\ &\leq 2Me^{-2\varepsilon^2 l} \end{aligned}$$

$$E_{in}(h, X^l) = \frac{1}{l} \sum_{i=1}^l \mathcal{L}(h, x_i)$$

$$E_{out} = \mathbb{E}_x \mathcal{L}(h, x)$$

В идеале хотим достичь $E_{out}(a) = 0$

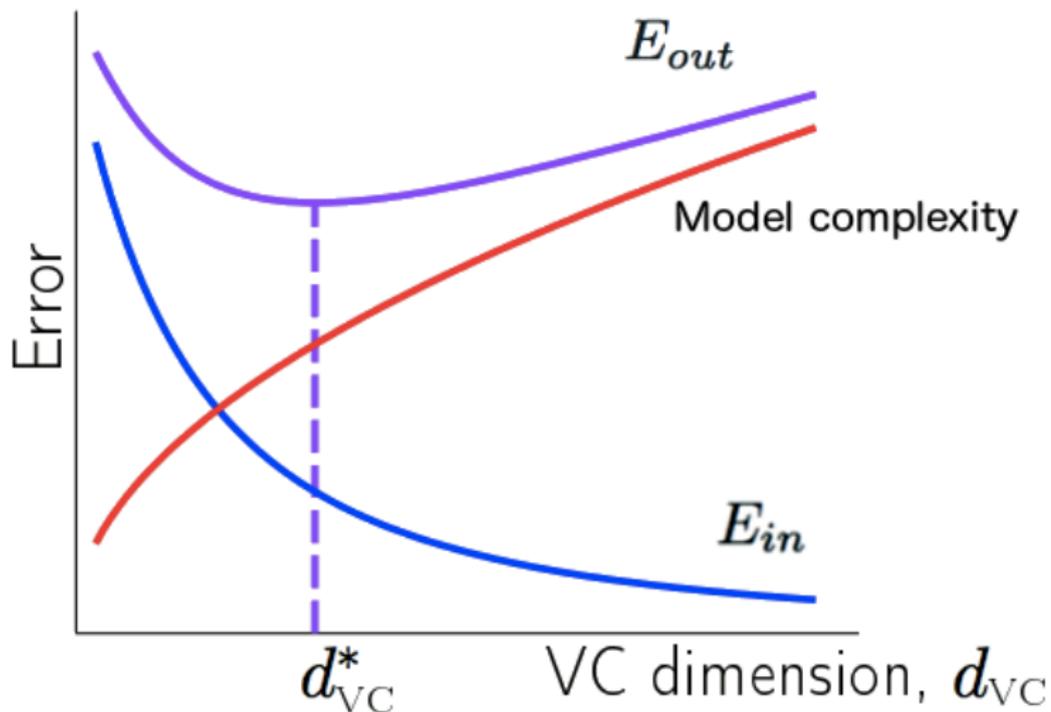
Пытаемся достичь $E_{in}(a) \approx E_{out}(a)$ и $E_{in}(a) = 0$

В идеале хотим достичь $E_{out}(a) = 0$

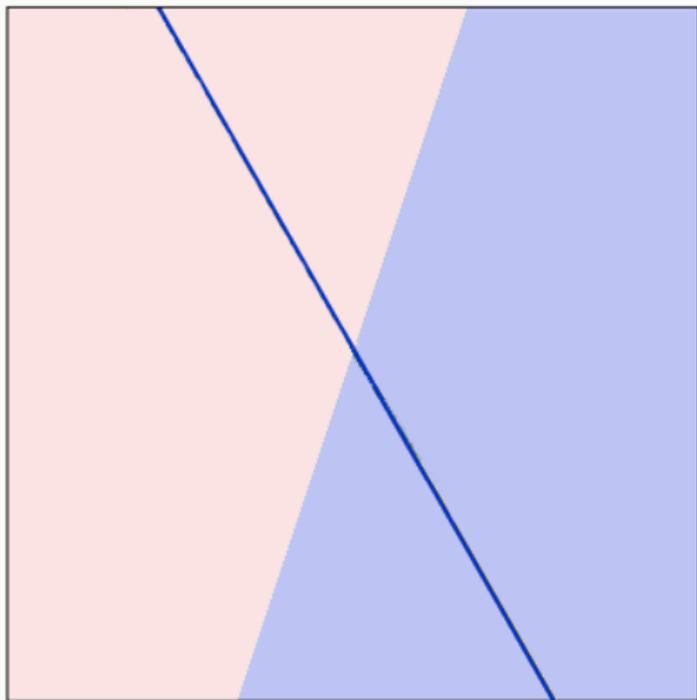
Пытаемся достичь $E_{in}(a) \approx E_{out}(a)$ и $E_{in}(a) = 0$

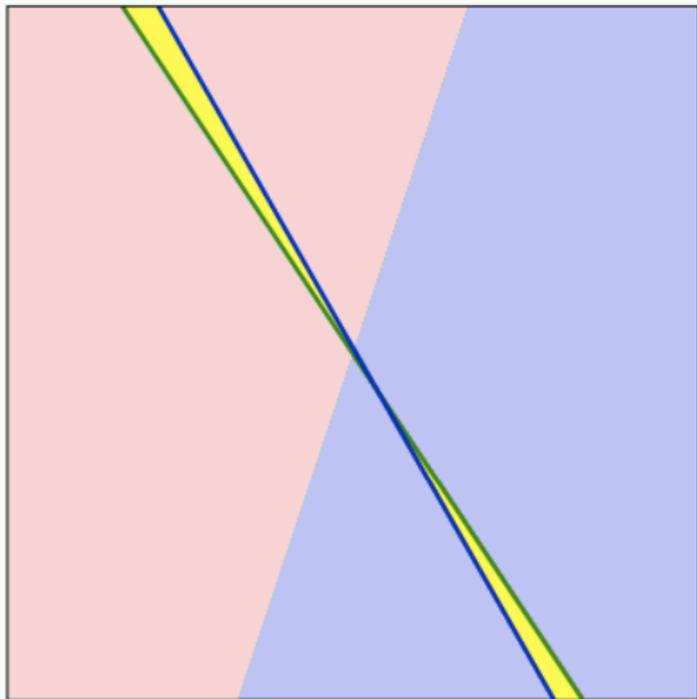
Два вопроса:

1. Можно ли судить о E_{out} по E_{in} ?
2. Как уменьшить E_{in} ?



Какое количество гипотез в
нашем пространстве \mathcal{H} ?



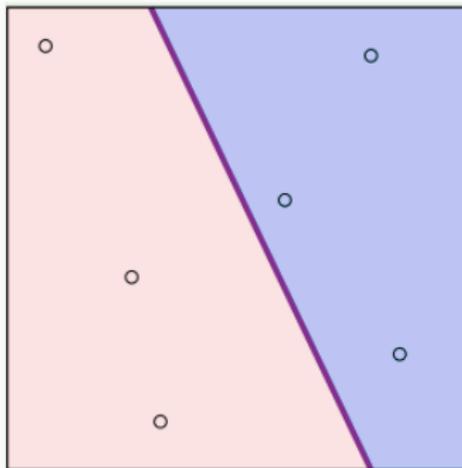


ΔE_{out} = площадь жёлтой области

ΔE_{in} = изменение меток объектов жёлтой области в выборке

$$|E_{in}(h_1) - E_{out}(h_1)| \approx |E_{in}(h_2) - E_{out}(h_2)|$$

Обучающая выборка x_1, \dots, x_l и набор бинарных значений меток y_1, \dots, y_l .



Сколько вариантов y_1, \dots, y_l ?

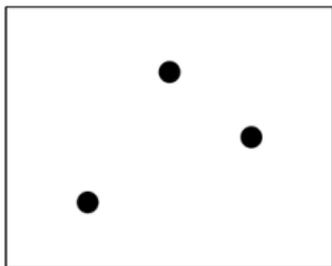
$$P(|E_{in}(a) - E_{out}(a)| > \varepsilon) \leq 2Me^{-2\varepsilon^2 l}$$

Наш набор моделей \mathcal{H} может породить $|\mathcal{H}(x_1, \dots, x_l)|$ дихотомий.

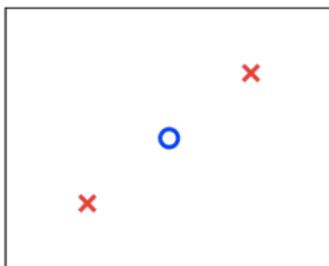
$$|\mathcal{H}(x_1, \dots, x_l)| \leq 2^l$$

При этом сам набор моделей может быть бесконечным.

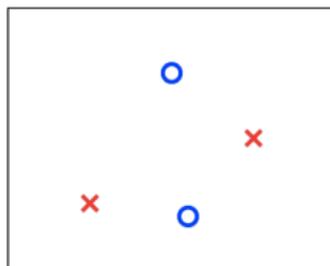
Функция роста $m_{\mathcal{H}(l)}$



$l = 3$



$l = 3$



$l = 4$

$$m_{\mathcal{H}(l)} = \max_{x_1, \dots, x_l} |\mathcal{H}(x_1, \dots, x_l)|$$

$$m_{\mathcal{H}(l)} \leq 2^l$$

Гипотезы:

1. Луч на множестве точек: $m_{\mathcal{H}}(l) = l + 1$
2. Интервал на множестве точек: $m_{\mathcal{H}}(l) = C_{l+1}^2 + 1$
3. Выпуклое множество на множестве точек: $m_{\mathcal{H}}(l) = 2^l$

Если для некоторого k выполняется $m_{\mathcal{H}}(k) < 2^k$, то k называется точкой разрыва.

Наличие точки разрыва означает наличие полиномиального ограничения на функцию роста $m_{\mathcal{H}}(l)$

Пример

	# of rows	x_1	x_2	...	x_{l-1}	x_l
S_1	α	+1	+1	...	+1	+1
		-1	+1	...	+1	-1
		⋮	⋮	⋮	⋮	⋮
		+1	-1	...	-1	-1
		-1	+1	...	-1	+1
S_2	β	+1	-1	...	+1	+1
		-1	-1	...	+1	+1
		⋮	⋮	⋮	⋮	⋮
		+1	-1	...	+1	+1
		-1	-1	...	-1	+1
S_2^-	β	+1	-1	...	+1	-1
		-1	-1	...	+1	-1
		⋮	⋮	⋮	⋮	⋮
		+1	-1	...	+1	-1
		-1	-1	...	-1	-1

Разделим на 2 ситуации относительно x_l :
либо есть только один вариант (+1 или -1), либо оба (+1 и -1)

$B(l, k)$ – максимальное число дихотомий для выборки размера l при наличии точки разрыва k .

- $B(l, k) = \alpha + 2\beta$

- $\alpha + \beta \leq B(l - 1, k)$, т.к. $\alpha + \beta$ – число дихотомий для $l - 1$

- $\beta \leq B(l - 1, k - 1)$

$$\implies B(l, k) \leq B(l - 1, k) + B(l - 1, k - 1)$$

Число дихотомий

$$B(l, k) \leq \sum_{i=0}^{k-1} C_l^i$$

Доказывается по индукции.

Индукционный шаг:

$$\sum_{i=0}^{k-1} C_l^i = \sum_{i=0}^{k-1} C_{l-1}^i + \sum_{i=0}^{k-2} C_l^i$$

$$m_{\mathcal{H}(l)} \leq \sum_{i=0}^{k-1} C_l^i$$

Доказывается по индукции.

Индукционный шаг:

$$\sum_{i=0}^{k-1} C_l^i = \sum_{i=0}^{k-1} C_{l-1}^i + \sum_{i=0}^{k-2} C_l^i$$

$$P(E_{in}(a) - E_{out}(a) > \varepsilon) \leq 4m_{\mathcal{H}}(2l)e^{-\frac{1}{8}\varepsilon^2 l}$$

Размерность d_{VC} = наибольшее l , для которого $m_{\mathcal{H}}(l) = 2^l$.
 $d_{VC} = k - 1$.

Вопросы?

Что почитать по этой лекции

- Professor Yaser Abu-Mostafa MOOC
- Tom Mitchell "Machine Learning" Chapter 4

На следующей лекции

- Задача максимизации зазора - аналог классификатора с регуляризацией
- Двойственная задача
- Что такое опорный вектор
- Регуляризация
- Решение для неразделимых выборок
- Kernel trick