# Линейная регрессия



①

Оценивание обобщ. способности:

1) отлож. выборка

2) кросс-валидация

$K$ - число блоков (fold)

$X^k$ - к-й блок

$X^{\backslash k}$ - все блоки, кроме к-го

$$CV = \frac{1}{K} \sum_{k=1}^{K} Q(a(x; X^{\backslash k}), X^k)$$

обучена на $X^{\backslash k}$

Leave-one-out
$K = l$

Итоговая модель: 1) обучить на всех данных

2) усреднить $K$ моделей

$$c(x) = \frac{1}{K} \sum_{k=1}^{K} a(x, X^{\backslash k})$$

## Борьба с переобучением

В ① $w_j$ очень большие
$10^6, 10^7, ...$

Регуляризация - запрет на большие веса

$(Q(w, X))^{\to min}$ - наш функционал

$$\boxed{Q(w, X) + a \|w\|_2^2 \to \min_w}$$

$$\begin{cases} w = (X'X + aI) X'y \\ \text{(для MSE)} \end{cases}$$

коэф. регуляризатор
регуляриз.

$$\boxed{a \geq 0}$$

**Важно**: в $\|w\|_2^2$ не входит $w_0$ !!!

(иначе $a(x)$ не сможет быть
одного порядка с $y$)

Почему большие веса — плохо?

$$a(x) = 10^6 + \underbrace{3 \cdot 10^7 (\text{площадь})} - 5 \cdot 10^6 \cdot (\text{рассш. до})_{\text{метро.}}$$

площадь + 0.001
$$a(x) + \underbrace{3 \cdot 10^4}_{30.000}$$

Как выбирать $a$?

— по обуч. выборке: выбираем $a$,
  при которой $Q(a, X)$ минимальна
  плохо — оптимально $a = 0$

Гиперпараметр — нельзя подбирать по обуч. выб.
  вводятся для улучш. качества
  на новых данных

  нужно подбирать по CV или отлож. выб.


обуч.   тест

Регуляриз. не обязательно через $L_2$-норму !

$$Q(w, X) + a \|w\|_1 \to \min \qquad a \, \pi$$
$$\vcentcolon= \sum_{j=1}^{d} |w_j| \qquad \boxed{\|w\|_1}$$

Допустим, в $X$ есть лин. зав. призиаки

$$\exists \vartheta \; \forall x \in X \; \langle \vartheta, x \rangle = 0$$

$w_*$ — лучшие по MSE веса

$$\langle \underbrace{w_* + \alpha \vartheta}_{}, x \rangle = \langle w_*, x \rangle + \alpha \underbrace{\langle \vartheta, x \rangle}_{=0} = \langle \underbrace{w_*}_{}, x \rangle$$
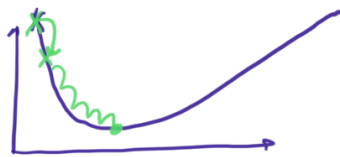
т.е. решений много

$$w_* + \alpha \vartheta \xrightarrow[\alpha \to \infty]{}$$

---

Обучение лин. регр.

для MSE: $\quad w = \underbrace{(X^T X)^{-1} X^T y}_{}$

$\underline{Q(d^3)}$

Градиентное обучение моделей



1) $w^{(0)}$ — начальное приближение

$\nabla_w Q(w)$ — градиент $Q$ по $w$

$-\nabla Q(w)$ — в сторону наискорейш. убывания

2) $w^{(k)} = w^{(k-1)} - \eta_k \nabla Q(w^{(k-1)})$ — шаг град. спуска

↳ длина шага
(learning rate)

3) когда останавливаться?

— когда ошибка на тесте перестает уменьшаться

— $\| w^{(k)} - w^{(k-1)} \| < \varepsilon$

— $| Q(w^{(k)}, X) - Q(w^{(k-1)}, X) | < \varepsilon$

— ...

Сходимость: 1) к $\nabla Q(w) \approx 0$

для лин. моделей с матрицей полного ранга такая точка одна

2) если решений несколько, то можно делать мульти-старт
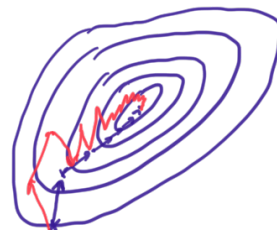


## Улучшения GD

оценки градиента

градиентный шаг

## Оценивание градиента

$$Q(w, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} (\langle w, x_i \rangle - y_i)^2$$

$$Q(w, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} q_i(w)$$

$$\nabla_w Q(w, X) = \frac{1}{\ell} \sum_{i=1}^{\ell} \nabla_w q_i(w)$$

полный градиент



Стохастический GD (SGD): $\nabla Q(w) \approx \nabla q_i(w)$

$$w^{(k)} = w^{(k-1)} - \eta_k \boxed{\nabla q_{i_k}(w^{(k-1)})} \qquad q_{i_k} = (\langle w, x_i \rangle - y_i$$

$i_k$ - случ. номер объекта

$$\rightarrow \sum_{k=1}^{\infty} \eta_k = \infty$$

$$\rightarrow \sum_{k=1}^{\infty} \eta_k^2 = const$$

$$\eta_k = \frac{1}{k}$$

$$\eta_k = \lambda \left( \frac{s_0}{s_0 + k} \right)^p$$

$\lambda, s_0, p$ – параметры

---

**Mini-batch GD:**

$$\nabla Q(w) \approx \frac{1}{n} \sum_{j=1}^{n} \nabla q_{i_{k_j}}(w)$$

---

**Full GD:**

$$Q(w^{(k)}) - Q(w_*) =$$

$$= \underline{O}\left( \frac{1}{k} \right)$$

**SGD:** $\underline{O}\left( \frac{1}{\sqrt{k}} \right)$

**SAG:** $\quad z_i^{(0)} = \nabla q_i(w^{(0)})$

$k$-я итерация:

$$z_i^{(k)} = \begin{cases} \nabla q_i(w^{(k-1)}), & i = i_k \\ z_i^{(k-1)}, & \text{иначе} \end{cases}$$

$$\nabla Q(w^{(k-1)}) \approx \frac{1}{l} \sum_{i=1}^{l} z_i^{(k)}$$

$$\underline{O}\left( \frac{1}{k} \right)$$