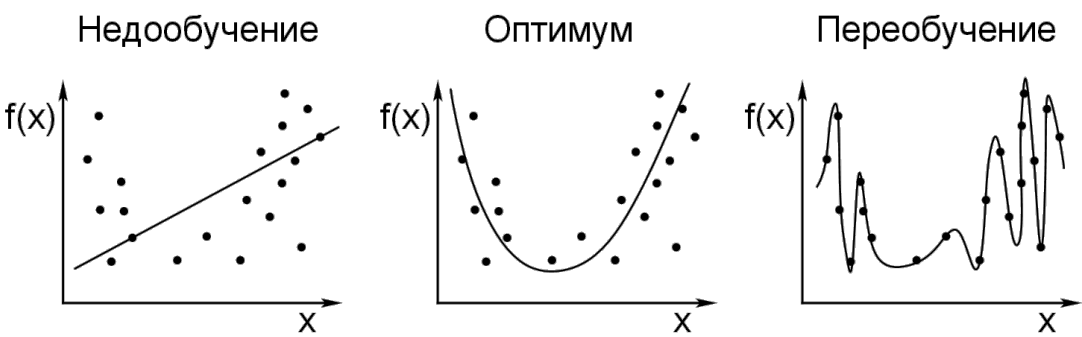


Лекция 3. Оценка обобщающей способности, градиентные методы обучения.

Переобучение.

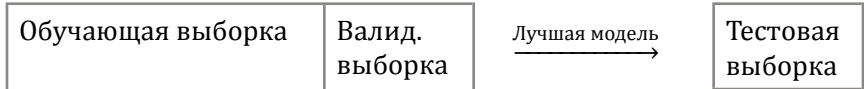


Переобучение (overfitting). Ошибка на новых данных > ошибка на обучающей выборке.

Обобщающая способность (generalization). Качество модели на новых данных не сильно хуже, чем на обучении.

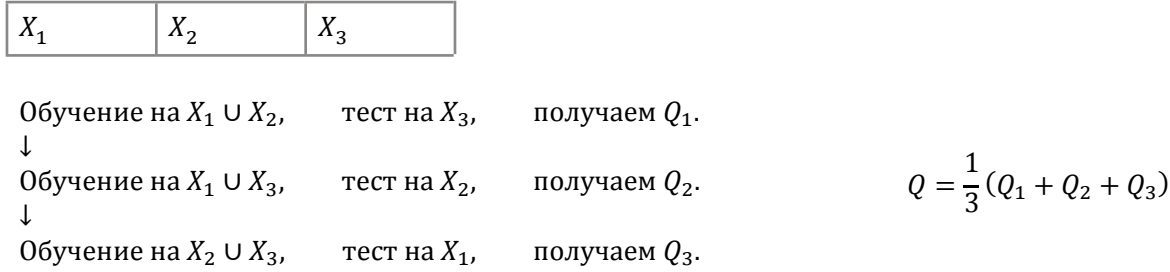
Оценка обобщающей способности.

1. Отложенная выборка (hold-out set).



2. Кросс-валидация (CV).

K-число блоков (folds).



При  $k = \ell \Rightarrow LOO$  (leave-one-out)

Что делать дальше?

★ Обучить модель на всей выборке.

★  $x \rightarrow \frac{1}{3}(a_1(x) + a_2(x) + a_3(x))$

Замечания:

★ Если  $\ell \gg 0$ , то CV вряд ли оправдано.

★ CV требует обучения k моделей.

Обучение линейных моделей.

$$\frac{1}{\ell} \sum_{i=1}^{\ell} (< w, x_i > - y_i)^2 \rightarrow \min_{w \in \mathbb{R}^d} \square, \quad X = \begin{pmatrix} x_{11} & \dots & x_{1d} \\ x_{21} & \dots & x_{2d} \\ \vdots & \ddots & \vdots \\ x_{\ell 1} & \dots & x_{\ell d} \end{pmatrix}, \quad w = \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix}, \quad y = \begin{pmatrix} y_1 \\ \vdots \\ y_{\ell} \end{pmatrix}, \quad X \cdot w = \begin{pmatrix} < w, x_1 > \\ \vdots \\ < w, x_{\ell} > \end{pmatrix}$$

Тогда векторный вид:  $Q(w) = \frac{1}{\ell} ||Xw - y||_2^2 \rightarrow \min_w \square$

$\nabla Q(w) = 0 \Rightarrow \boxed{w_* = (X^T X)^{-1} X^T y}$

Проблемы:

★ Матрица может быть вырожденной.

★ Обращение матрицы за  $\mathcal{O}(d^3)$

★ Если функция потерь  $L(y, z)$  более хитрая, то решить  $\nabla Q(w) = 0$  не выйдет.

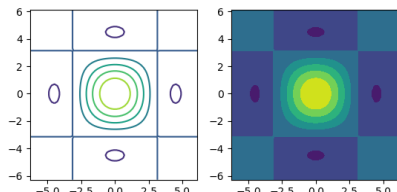
Градиентные методы обучения.

$Q(w_1, \dots, w_d) \rightarrow \min_w \square, \quad Q - \text{дифференцируемый.}$

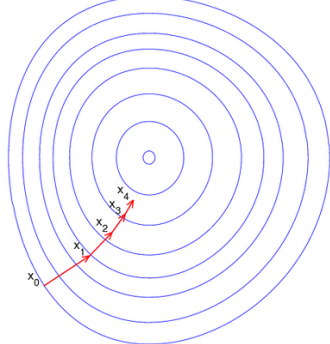
Важные свойства градиента:

★  $\nabla Q(w)$  показывает направление наискорейшего роста в этой точке  $w$ ,  $(-\nabla Q(w))$  наискорейшее убывание)

★  $\nabla Q(w)$  ортогонален линии уровня



Градиентный спуск.



$w^{(0)}$  — инициализация (в линейной модели можно брать случайно)

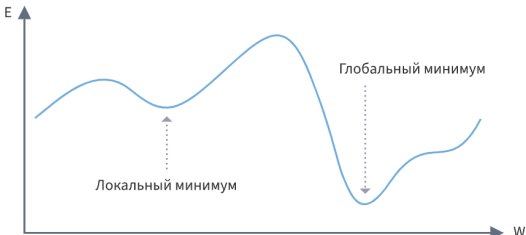
Шаг спуска:  $\boxed{w^{(k)} = w^{(k-1)} - \eta \nabla Q(w^{(k-1)})}$ ,  $\eta$  — длина шага (learning rate)

Останавливаемся, если:

- ★  $||w^{(k)} - w^{(k-1)}|| < \varepsilon$
- ★  $|Q(w^{(k)}) - Q(w^{(k-1)})| < \varepsilon$
- ★  $||\nabla Q(w^{(k)})|| < \varepsilon$
- ★  $k > N$
- ★ ошибка на отложенной выборке перестала убывать.

Некоторые наблюдения:

★ Локальные минимумы:



Решение: мультистарт (сомнительно для большого количества весов)

★ Много условий сходимости:

$$\begin{cases} Q(w) - \text{выпуклая + дифференцируемая} \\ \nabla Q(w) - \text{липшицева} (||\nabla Q(w_1) - \nabla Q(w_2)|| \leq L \cdot ||w_1 - w_2||) \\ \eta \text{ не очень большая } \left( \eta \leq \frac{1}{L} \right) \end{cases} \Rightarrow \text{градиентный спуск сходится к минимуму.}$$

★ С линейными моделями и адекватными функциями потерь  $Q(w)$  всегда (почти всегда) выпуклые.

★  $Q(w^{(k)}) - Q(w_*) = \mathcal{O}\left(\frac{1}{k}\right)$

Оценивание градиента.

$Q(w) = \frac{1}{\ell} \sum_{i=1}^{\ell} L(y_i, a(x_i, w)), \quad Q(w) = \frac{1}{\ell} \sum_{i=1}^{\ell} q_i(w), \quad \nabla Q(w) = \frac{1}{\ell} \sum_{i=0}^{\ell} \nabla q_i(w), \quad \ell = 10^6 \Rightarrow 10^6 \text{ градиентов 🍌🍌🍌}$

Стохастический градиентный спуск (SGD).

$\nabla Q \approx \nabla q_i(w) \Rightarrow$  шаг:  $i_k$  — индекс случайного объекта.

$w^{(k)} = w^{(k-1)} - \eta \nabla q_{i_k}(w^{(k-1)})$

Проблема: если  $||w^{(k)} - w_*|| \gg 0$ , то  $\nabla q_{i_k}(w) \approx \nabla Q(w)$   
если  $||w^{(k)} - w_*|| \approx 0$ , то  $\nabla q_{i_k}(w) \neq \nabla Q(w)$

Идея:  $w^{(k)} = w^{(k-1)} - \eta \nabla q_{i_k}(w^{(k-1)})$ , где  $\sum_{k=1}^{\infty} \eta_k = \infty$  и  $\sum_{k=1}^{\infty} \eta_k^2 < \infty \Rightarrow SGD$  сойдется к минимуму (если угадаем  $w^{(0)}$ )  
условия Роббинса-Монро

$\eta_k = \lambda \left( \frac{s_0}{s_0 + k} \right)^p$ ,  $\lambda, s_0, p$  — надо подбирать.

Скорость сходимости:  $\mathbb{E}[Q(w^{(k)}) - Q(w_*)] = \mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$

Замечания:

★ mini-batch GD:  $\nabla Q(w) \approx \frac{1}{\ell} \sum_{j=1}^{\ell} \nabla q_{i_j}(w)$

★ SGD хорош для онлайн-обучения (можно обучаться на огромных выборках).

Stochastic average gradient (SAG)

$Q(w) = \frac{1}{\ell} \sum_{i=1}^{\ell} q_i(w), \quad z_i^{(0)} = \nabla q_i(w^{(0)})$

Итерация SAG:  $i_k \sim \{1, \dots, \ell\}, \quad z_i^{(k)} = \begin{cases} \nabla q_i(w^{(k-1)}), & i = i_k \\ z_i^{(k-1)}, & i \neq i_k \end{cases}$

$\nabla Q(w^{(k-1)}) \approx \frac{1}{\ell} \sum_{i=1}^{\ell} z_i^{(k)}, \quad w^{(k)} = w^{(k-1)} - \eta_k \cdot \frac{1}{\ell} \sum_{i=1}^{\ell} z_i^{(k)}$

$\mathbb{E}[Q(w^{(k)}) - Q(w_*)] = \mathcal{O}\left(\frac{1}{k}\right)$