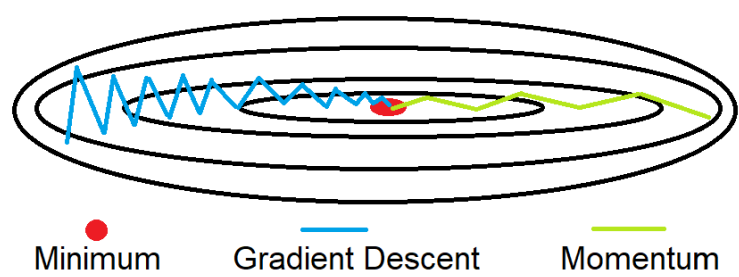


Лекция 4. Модификации градиентного спуска, регуляризация.

Модификации градиентного спуска.

Метод инерции (*momentum*).



$w^{(0)}$ – инициализация весов, h_0 – вектор инерции

Шаг:
$$h_k = \alpha \cdot h_{k-1} + \underbrace{\eta_k \nabla_w Q(w^{(k-1)})}_{\substack{\text{можно} \\ \text{заменять} \\ \text{на оценку град.}}} \Rightarrow w^{(k)} = w^{k-1} - h_k$$

Адаптивный шаг (*AdaGrad*). в основном для разреженных данных.

G_{0j} , j -номер признака.

$$G_{kj} = G_{k-1,j} + \left(\nabla Q(w^{(k-1)})\right)_j^2$$
 – насколько сильно уже обучена w_j

Шаг:
$$w_j^{(k)} = w_j^{(k-1)} - \frac{\eta_k}{\sqrt{G_{kj} + \varepsilon}} \cdot \left(\nabla Q(w^{(k-1)})\right)_j$$

Проблема: G_{kj} только растёт и может быстро остановить оптимизацию.

Модификация: *RMSProp*
$$G_{kj} = \alpha G_{k-1,j} + (1 - \alpha) \cdot \left(\nabla Q(w^{(k-1)})\right)_j^2, \quad \alpha \in (0,1)$$

Adam (*momentum* + *AdaGrad*) одновременно и инерция, и адаптивность.

Регуляризация.

Известный экспериментальный факт:

Линейная модель переобучена \Leftrightarrow большие веса.
по слухам

Почему?

Объяснение №1. Пусть есть линейно зависимые признаки $\exists v \in \mathbb{R}^d: \forall x \in \mathbb{X} \quad \langle v, x \rangle = 0$

w_\star – решение
$$\frac{1}{\ell} \sum_{i=1}^{\ell} (\langle w, x_i \rangle - y_i)^2 \rightarrow \min_w$$

$\alpha > 0: \langle w_\star + \alpha v, x \rangle = \langle w_\star, x \rangle + \alpha \underbrace{\langle v, x \rangle}_0 = \langle w_\star, x \rangle,$ т.е. $w_\star + \alpha v$ – решение (еще одно)

Объяснение №2. $a(x) = 10^8 \cdot \text{площадь} - 10^9 \cdot \text{этаж} + 10^{11} \cdot \text{район}.$ Добавим совсем чуть чуть к признаку

$$10^8(\text{площадь} + 0,001) = 10^8 \cdot \text{площадь} + \underbrace{10^8 \cdot 0,001}_{10^5}$$

Гиперчувствительность к изменениям в признаках – не соответствует тому, как работает мир.

! Запретить большие веса !

$$Q(w) + \underbrace{\alpha}_{\substack{\text{коэф.} \\ \text{регуляр.}}} \cdot \underbrace{R(w)}_{\text{регуляризатор}} \rightarrow \min$$

$$R(w) = \begin{cases} \|w\|_1 = \sum_{j=1}^d |w_j| & \text{— это } L_1 \text{ регуляризатор.} \\ \|w\|_2^2 = \sum_{j=1}^d w_j^2 & \text{— это } L_2 \text{ регуляризатор.} \end{cases}, \quad \begin{cases} \alpha \gg 0 & \Rightarrow w_\star = 0 \\ \alpha = 0 & \Rightarrow \text{нет рег.} \end{cases}$$

α нельзя подбирать по обучающей выборке – гиперпараметр. Подбираем по новым данным (отложенной выборке, ...)

Стратегии подбора α :

- ★ *Grid Search*
- ★ *Random Search*
- ★ *AutoML*

Ridge-регрессия
$$\frac{1}{\ell} \sum_{i=1}^{\ell} (\langle w, x_i \rangle - y_i)^2 + \alpha \|w\|_2^2 \rightarrow \min_w, \quad \text{невыврожденная матрица}$$

LASSO
$$\frac{1}{\ell} \sum (\langle w, x_i \rangle - y_i)^2 + \alpha \|w\|_1 \rightarrow \min, \quad \text{зануляет часть весов}$$