

Минимизация эмпирического риска:

$$Q(w) = \sum_{i=1}^{\ell} \mathcal{L}(\langle x_i, w \rangle y_i) = \sum_{i=1}^{\ell} \mathcal{L}_i(w) \rightarrow \min_w.$$

Численная минимизация методом *градиентного спуска*:

$w^{(0)}$  := начальное приближение;

$$w^{(t+1)} := w^{(t)} - h \cdot \nabla Q(w^{(t)}), \quad \nabla Q(w) = \left( \frac{\partial Q(w)}{\partial w_j} \right)_{j=0}^n,$$

где  $h$  — *градиентный шаг*, называемый также *темпом обучения*.

$$w^{(t+1)} := w^{(t)} - h \sum_{i=1}^{\ell} \nabla \mathcal{L}_i(w^{(t)}).$$

**Идея ускорения сходимости:**

брать  $(x_i, y_i)$  по одному и сразу обновлять вектор весов.

**Вход:** выборка  $X^\ell$ , темп обучения  $h$ , темп забывания  $\lambda$ ;

**Выход:** вектор весов  $w$ ;

1 инициализировать веса  $w_j$ ,  $j = 1, \dots, n$ ;

2 инициализировать оценку функционала:  $\bar{Q} := \frac{1}{\ell} \sum_{i=1}^{\ell} \mathcal{L}_i(w)$ ;

3 **повторять**

4     выбрать объект  $x_i$  из  $X^\ell$  случайным образом;

5     вычислить потерю:  $\varepsilon_i := \mathcal{L}_i(w)$ ;

6     сделать градиентный шаг:  $w := w - h \nabla \mathcal{L}_i(w)$ ;

7     оценить функционал:  $\bar{Q} := (1 - \lambda) \bar{Q} + \lambda \varepsilon_i$ ;

8 **пока** значение  $\bar{Q}$  и/или веса  $w$  не сойдутся;

---

*Robbins, H., Monro S. A stochastic approximation method // Annals of Mathematical Statistics, 1951, 22 (3), p. 400–407.*

**Проблема:** после каждого шага  $w$  по одному объекту  $x_i$ , не хотелось бы оценивать  $Q$  по всей выборке  $x_1, \dots, x_\ell$ .

**Решение:** использовать рекуррентную формулу.

Среднее арифметическое  $\bar{Q}_m = \frac{1}{m} \sum_{i=1}^m \varepsilon_i$ :

$$\bar{Q}_m = (1 - \frac{1}{m})\bar{Q}_{m-1} + \frac{1}{m}\varepsilon_m.$$

*Экспоненциальное скользящее среднее*

$$\bar{Q}_m := (1 - \lambda)\bar{Q}_{m-1} + \lambda\varepsilon_m;$$

$$\bar{Q}_m = \lambda\varepsilon_m + \lambda(1 - \lambda)\varepsilon_{m-1} + \lambda(1 - \lambda)^2\varepsilon_{m-2} + \lambda(1 - \lambda)^3\varepsilon_{m-3} + \dots$$

Чем больше  $\lambda$ , тем быстрее забывается предыстория ряда.

Параметр  $\lambda \approx \frac{1}{m}$  называется *темпом забывания*.