

| Название                    | Формула  | Недостатки  |
|-----------------------------|--|---|
| Gradient Descent            | $\theta = \theta - \alpha \nabla_{\theta} J(\theta)$   | Градиент вычисл. проходим по всему наб. данным<br>$\Rightarrow$ Треб. много памяти.<br>Также из-за слишком большого градиента алгоритм может не сойтись     |
| Stochastic Gradient Descent | $\theta = \theta - \alpha \nabla_{\theta} J(\theta, \mathcal{S})$ ,<br>где $\mathcal{S}$ - параметры для обучения            | На каждом парам. обуч. меняются веса $\Rightarrow$<br><ul style="list-style-type: none"> <li>• аномалии сильно влияют</li> <li>• долгое обучение</li> </ul> |
| Mini-Batch Gradient Descent | $\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta, N_{\mathcal{S}})$ ,<br>$N_{\mathcal{S}}$ - группы обуч. параметров | Параметр $N$ оказывает основное влияние на обучение $\Rightarrow$ необходимо потратить много времени на его подбор  |
| $\nabla GD$ + Momentum      | $\vartheta = \gamma \cdot \vartheta + \eta \nabla_{\theta} J(\theta)$<br>$\theta = \theta - \alpha \vartheta$                | Также необходимо потратить много времени на ручной подбор параметра.<br>К тому же импульс   |

|                                      |  |  |
|--------------------------------------|--|--|
|                                      |  | <p>может стать слишком большим <math>\Rightarrow</math> начнёт пропускать локал. мин.</p>                      |
| <p>SGD + Momentum + Acceleration</p> | $v = \gamma \cdot v + \eta \nabla_{\theta} J(\theta - \gamma \cdot v)$ $\theta = \theta - \alpha v$  | <p>Аналогично предыдущ. сложно подобрать параметр.</p>   |
| <p>Adagrad</p>                       | $g_t = \nabla_{\theta} J(\theta_t),$ $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot g_t$  | <p>Знаменатель всё время увелич.<br/> <math>\Rightarrow</math> скорость обуч. уменьшается и может останов.</p> |
| <p>Adadelta</p>                      | $\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{E[G_{t,ii}] + \epsilon}} \cdot \nabla_{\theta_{t,i}} J(\theta_{t,i})$   | <p>Вычислительно дорого</p>  |
| <p>Adam</p>                          | $\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{E[G_{t,ii}] + \epsilon}} \cdot E[g_{t,ii}]$   | <p>Вычислительно дорого<br/> Хуже обобщ. данные</p>  |
| <p>Nadam</p>                         | $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[G_t] + \epsilon}} \cdot \left( \beta_1 E g_t + \frac{(1 - \beta_1) \nabla_{\theta} J(\theta_t)}{1 - \beta_1^t} \right)$ | <p>Низкая скорость обучения</p>  |

Accelerate - создан для ускорения импульса, даёт возможность вычислять его относительно предположит. будущего.

**Adagrad:**

меняет скорость обуч. на  $\eta$  для каждого параметра на каждом шаге. Работает на  $\sqrt{\text{ф-ции ошибок}}$ .

$$g_t = \nabla_{\theta} J(\theta_t)$$

$$G_t = G_t + g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot g_t$$

$\eta$  - скорость обуч., которая изменяется для заданного  $\theta_i$ .

**Adadelta:**

Расширение Adagrad. Огранич окно накопленных градиентов до фиксир. размера.

Используется экспоненцирующая средняя

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1-\gamma)g_t^2$$

$$RMSE[g]_t = \sqrt{E[g^2]_t + \epsilon}$$

$$\theta_{t+1} = \theta_t - \frac{RMSE[\theta]_{t-1}}{RMSE[g]_t} \cdot g_t$$

## Adam

Работает с шп. 1-го и 2-го порядка.

Уменьшается скорость во избежание проскаков. min

Сохраняются скользящие средние прошлых  $\nabla$ .

$$\hat{m}_t = \frac{m_t}{1-\beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1-\beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \cdot \hat{m}_t$$