

Optimization

GD - градиентный спуск

SGD - стохастический градиентный спуск

mini batch SGD - мини батч SGD

batch размер стечки обновлений

Основные обозначения

x, y - \mathcal{L} шумы

l число батчей 1

B - batch size

L - loss function

w/θ - параметры

$a_w(x)$ - модель

\hat{y} - предсказания модели

$$L = \frac{1}{B} \sum_{i=1}^B (y_i - \hat{y}_i)^2 \rightarrow \min_w - MSE$$

$$\hat{y}_i = a_w(x_i)$$

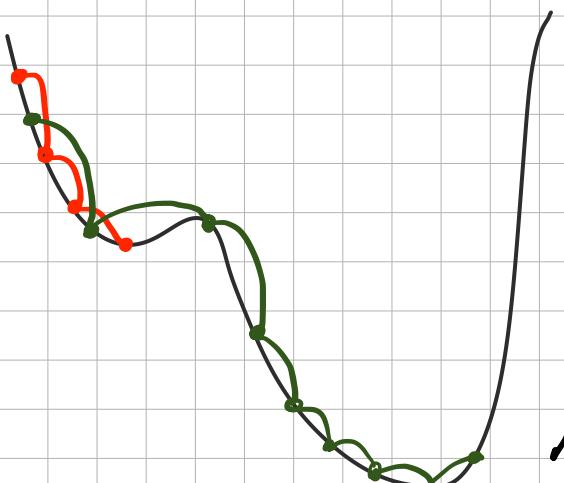
t - целевая оптимизируемая

SGD

$$g_t = \nabla_w L^t(w^{t-1})$$

$$w_t = w^{t-1} - \eta_t g_t$$

SGD + momentum



$$m_t = \mu \cdot m_{t-1} + g_t$$

$$w_t = w_{t-1} - \eta_t m_t$$

μ - значение момента
(инерциальный)

Значение инерциальных

$$\mu \in (1e^{-2}; 1e^{-7}) - \text{з}$$

$$m_0 = 0$$

$$\mu \in [0; 1] - \text{momentum}$$

обычно 0,8 - 0,9

+ weight decay

$$L \rightarrow L^* = L + \lambda \cdot \|w\|_2^2$$

$$\nabla_w \|w\|_2^2 = \frac{\partial}{\partial w} \|w\|_2^2$$

$$g_t \rightarrow g_t^* = \nabla_w L(w_{t-1}) + \underbrace{\lambda \cdot 2}_{\text{объко просимо}} w_{t-1}$$

λ - weight_decay

Adam

Универсальный алгоритм:

- 1) Momentum
- 2) Adaptive Learning Rate

$$g_t = \nabla_w L(w_{t-1}) + \underbrace{\lambda w_{t-1}}_{\text{weight decay}}$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \cdot g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) \cdot g_t^2 \quad \leftarrow \text{нормализация}$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad \leftarrow \text{сдвиги}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$w_t = w_{t-1} - \eta_t \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad \leftarrow \text{запись от градиента к весам}$$

$$m_0 = v_0 = 0$$

$$\beta_1 = 0,9$$

$$\beta_2 = 0,999$$

$$\epsilon = 10^{-9}$$

Разложение Adam

$$m_t = (1 - \beta_1) g_t + \beta_1 m_{t-1} =$$

раскладываем по формуле

$$= (1 - \beta_1) g_t + \beta_1 (1 - \beta_1) \cdot g_{t-1} + \beta_1^2 \cdot m_{t-2} =$$

раскладываем по формуле \times

$$= (1 - \beta_1) g_t + \beta_1 (1 - \beta_1) \cdot g_{t-1} + \beta_1^2 \cdot (1 - \beta_1) \cdot g_{t-2} + \beta_1^3 \cdot m_{t-3} \quad (\equiv)$$

Если продолжать раскладывать до m_0 ,

то получим:

$$\equiv (1 - \beta_1) g_t + \beta_1 (1 - \beta_1) \cdot g_{t-1} + \beta_1^2 \cdot (1 - \beta_1) \cdot g_{t-2} + \dots$$

$$\dots + \beta_1^{t-1} (1 - \beta_1) \cdot g_1 + \cancel{\beta_1^t m_0}^0 =$$

$$= \sum_{i=0}^{t-1} \beta_1^i (1 - \beta_1) \cdot g_{t-i}$$

$$\hat{M}_t = \frac{m_t}{1 - \beta_1^t} = \frac{1 - \beta_1}{1 - \beta_1^t} \cdot \sum_{i=0}^{t-1} \beta_1^i g_{t-i} =$$

$$= \frac{1 - \beta_1}{(1 - \beta_1) \left(\sum_{k=0}^{t-1} \beta_1^k \right)} \cdot \sum_{i=0}^{t-1} \beta_1^i g_{t-i} \quad (\equiv)$$

$$= \frac{\sum_{i=0}^{t-1} \beta_i^i g_{t-i}}{\sum_{k=0}^{t-1} \beta_k^k} = \sum_{i=0}^{t-1} d_i g_{t-i} =$$

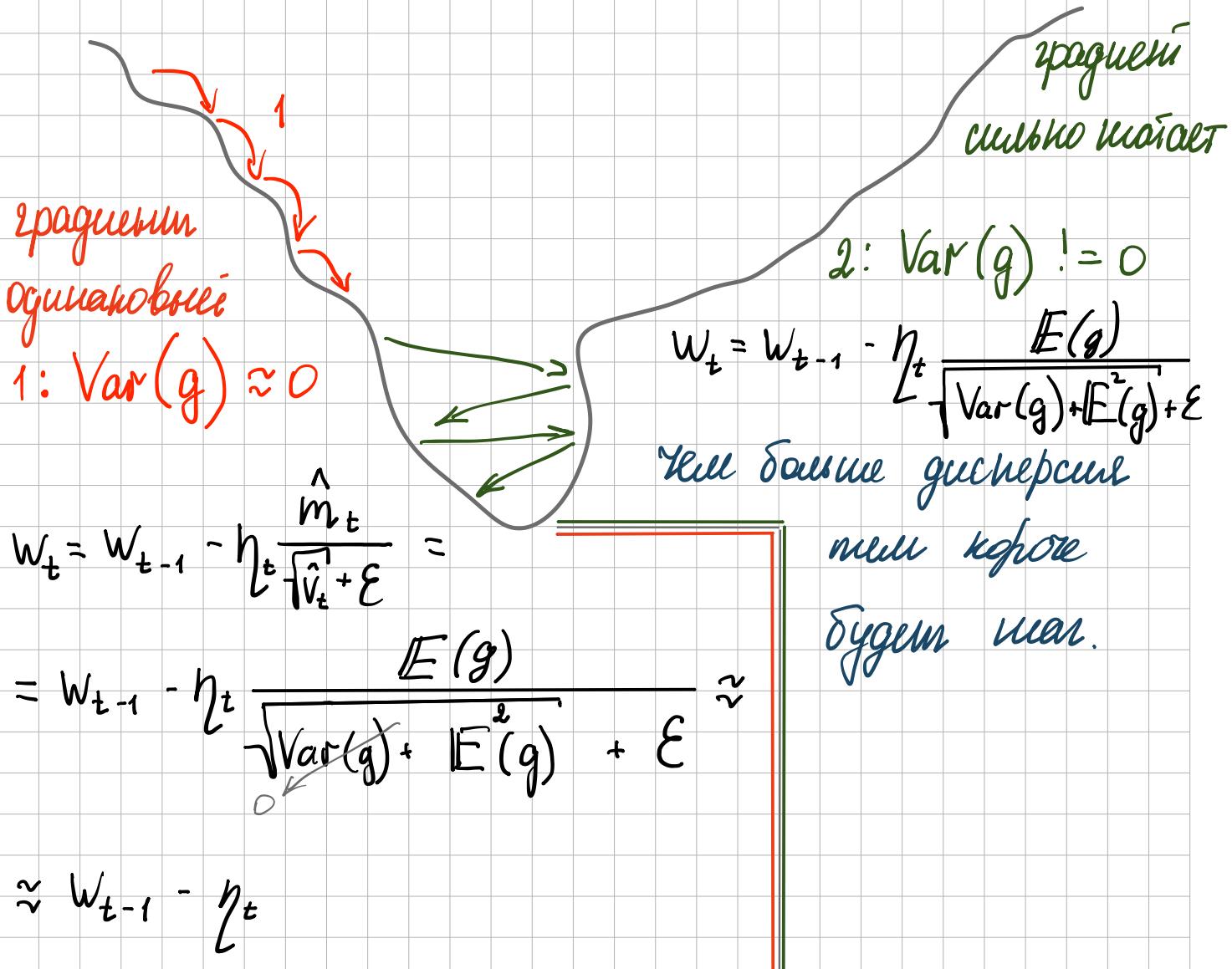
$$d_i = \frac{\beta_i^i}{\sum_{k=0}^{t-1} \beta_k^k}$$

= $E(g)$

$$d_i > 0$$

$$\sum_{i=0}^{t-1} d_i = 1$$

$$\hat{V}_t = E(g^2) = \text{Var}(g) + E^2(g)$$



Adam W

$$g_t = \nabla_w L^t(w_{t-1})$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

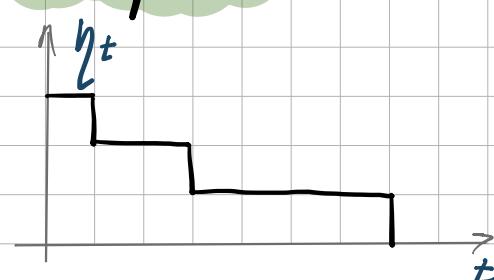
$$w_t = w_{t-1} \left(1 - \underbrace{\eta_t \lambda}_{\text{weight_decay}} \right) - \eta_t \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t + \epsilon}}$$

Изменение шага обновлений

Constant

$$\eta_t = \eta_0$$

Step LR



$$\eta_0 : 1 \leq t < t_0$$

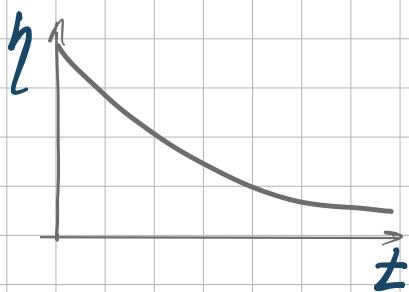
$$\eta_1 : t_0 \leq t < t_1$$

:

Exponential LR

$$\eta_t = \gamma \eta_{t-1}$$

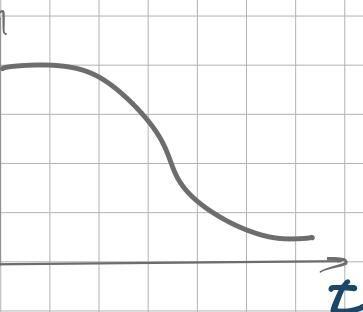
$$0 < \gamma < 1$$



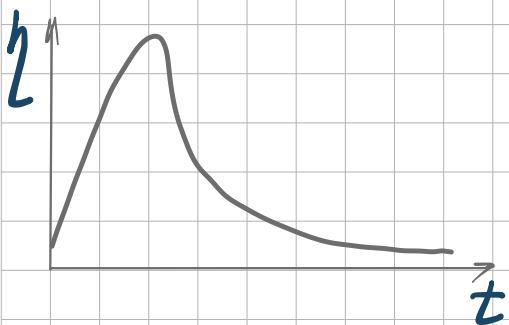
Cosine LR

$$\eta_t = \eta_0 \cdot \frac{1}{2} \left(1 + \cos \frac{\pi t}{T} \right)$$

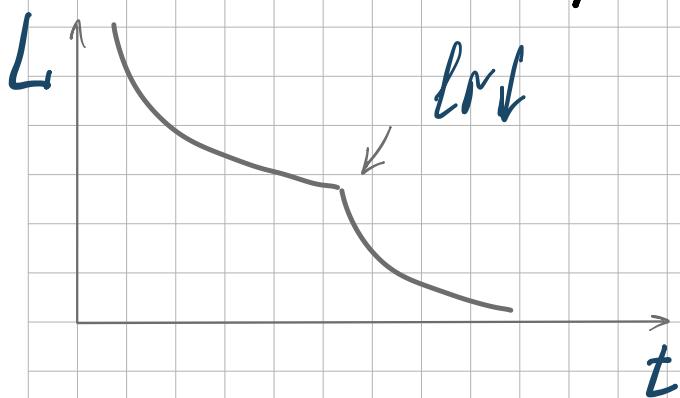
T = epochs



Linear warm up



Reduce LR on plateau



Cosine with Restarts

