# Gradient descent.

1) $\Theta_0 = \Theta_{init} \sim N(\cdot, \cdot)$. $\quad (x, y) \in D \quad |D| = N$

2) $\Theta_n = \Theta_{n-1} - \alpha \frac{1}{N} \sum_{i=0}^{N} \nabla_\theta L \left( f(x_i, \Theta_{n-1}), y_i \right)$ ⟶ D

---

1) $\Theta_0 = \Theta_{init} \sim N(\cdot, \cdot)$. SGD
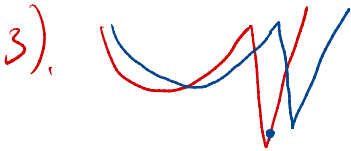
2) $d$ ↙ mini batch. is random subset of $D$ of fixed size k.

3) $\Theta_n = \Theta_{n-1} - \alpha \frac{1}{k} \sum_{|x, y| \in d} \nabla_\theta L \left( f(x_i, \Theta_{n-1}), y_i \right)$

1) От $k$ зависит "уровень шума".

2) $k = N \quad \alpha = \alpha_0$. $\quad k = N/5 \quad \alpha = \alpha_0/5$.

3).

$$m_t = \nabla_{\theta_{t-1}} L (1-\beta) + m_{t-1} \cdot \beta$$

momentum

$$\theta_t = \theta_{t-1} - \alpha m_t$$

$$\beta = 0.9$$

SGD momentum.

---

$$m_t = \nabla_\theta L(\theta_{t-1}) \cdot (1-\beta) + m_{t-1} \cdot \beta$$

$$m_t = \nabla_\theta L(\theta_{t-1} + \beta m_{t-1}) \cdot (1-\beta) + m_{t-1} \cdot \beta$$

lookahead gradient.

Nesterov momentum

---

$$V_t = \beta V_{t-1} + (1-\beta)\left(\nabla_\theta L(\theta_{t-1})\right)^2$$

$$\theta_t = \theta_{t-1} - \alpha \frac{\nabla_\theta L(\theta_{t-1})}{\sqrt{V_t} + \epsilon} \begin{bmatrix} AdaGrad \\ RmsProp \end{bmatrix}$$

AdaGrad + Momentum
    AdaM.

AMSgrad.

$$m_t = \nabla_\theta L(\theta_{t-1}) \cdot (1-\beta) + m_{t-1} \cdot \beta_1$$

$$V_t = \beta_2 V_{t-1} + (1-\beta_2)(\nabla_\theta L(\theta_{t-1}))^2$$

$$\theta_t = \theta_{t-1} - \alpha \frac{m_t}{\sqrt{V_t} + \epsilon}$$

$$\hat{V}_t = \beta_2 V_{t-1} + (1-\beta_2)(\nabla_\theta L(\theta_{t-1}))^2$$
$$m_t$$

$$V_t = \max(\hat{V}_t, V_{t-1})$$

Nadam
AdaBelief.

$$V_t = V_{t-1} \cdot \beta_2 + (1-\beta_2)(m_t - \nabla_\theta L(\theta_{t-1}))^2$$

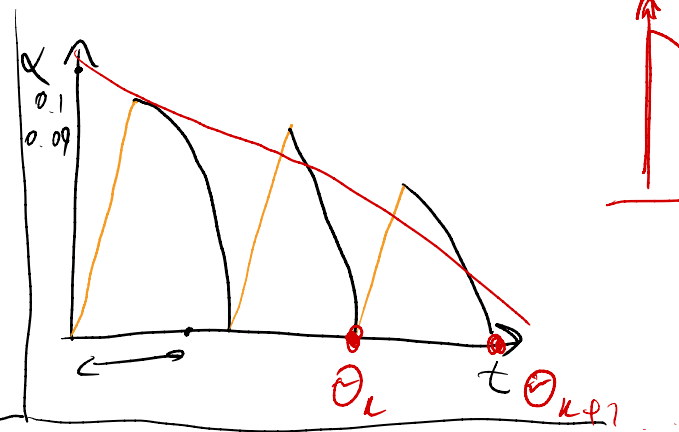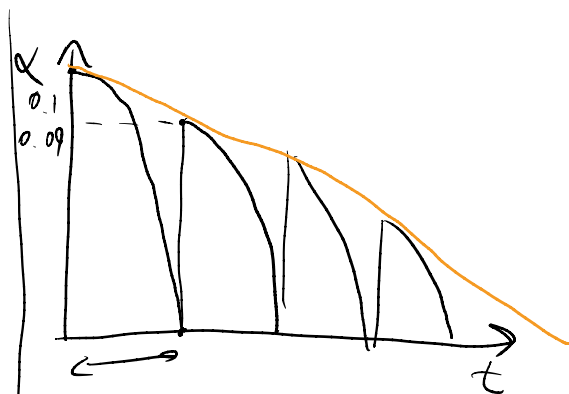n - number of ops.    ⎤  backprop.
~2n      $O(n)$.         ⎦

P - num of params  $|\theta|$

$\underline{2n} + 2p$ — momentum

$\underline{2n} + 3p$ — adam

$2n + 2p$ — Nesterov.



$\alpha_1 = \alpha_0 / 10$





$\theta_i$   $\theta_{k+1}$

$\Theta = \frac{1}{n} \Sigma \theta_i$