

## Gradient descent methods in NN

Cost function:  $J(\theta)$ ,  $\theta \in \mathbb{R}^d$  NN params

Gradient:  $\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}$

Generally speaking,  $\theta^{(t+1)} = \theta^{(t)} - \Delta^{(t)}$

### ① Gradient descent

$$\Delta^{(t)} = \eta \nabla_{\theta} J(\theta^{(t)}) \quad \leftarrow \eta: \text{learning rate}$$

### ② Momentum

$$\Delta^{(t)} = \gamma \Delta^{(t-1)} + \eta \nabla_{\theta} J(\theta^{(t)})$$

$\gamma$ : decay factor; 0.9

### ③ Adagrad

$$\Delta^{(t)} = \frac{\eta}{\sqrt{G^{(t)} + \epsilon I}} \nabla_{\theta} J(\theta^{(t)})$$

$$G^{(t)} = \sum_{i=1}^t \left[ \frac{\partial J}{\partial \theta_i} \right]^2$$

$$i=1 \quad \left[ \frac{\partial J}{\partial \theta_d} \right]_{\theta=\theta^{(i)}}$$

④ RMS prop

$$\Delta^{(t)} = \frac{n}{\sqrt{G^{(t)} + \epsilon I}} \nabla_{\theta} J(\theta^{(t)})$$

$$G^{(t)} = \gamma G^{(t-1)} + (1-\gamma) \begin{bmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix}_{\theta=\theta^{(t)}}^T$$

0.9 ↗

⑤ Adam

$$m^{(t)} = \beta_1 m^{(t-1)} + (1-\beta_1) \nabla_{\theta} J(\theta^{(t)})$$

$$v^{(t)} = \beta_2 v^{(t-1)} + (1-\beta_2) (\nabla_{\theta} J(\theta^{(t)}))^2$$

$$\hat{m}^{(t)} = \frac{m^{(t)}}{1-\beta_1^t}, \quad \hat{v}^{(t)} = \frac{v^{(t)}}{1-\beta_2^t}$$

element wise

$$\Delta^{(t)} = \frac{n}{\sqrt{G^{(t)} + \epsilon}} \hat{m}^{(t)}$$

Typical values:  $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$