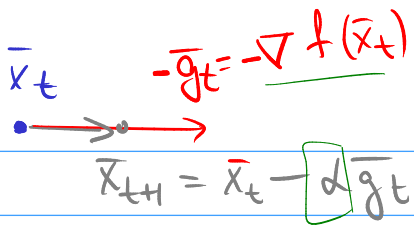
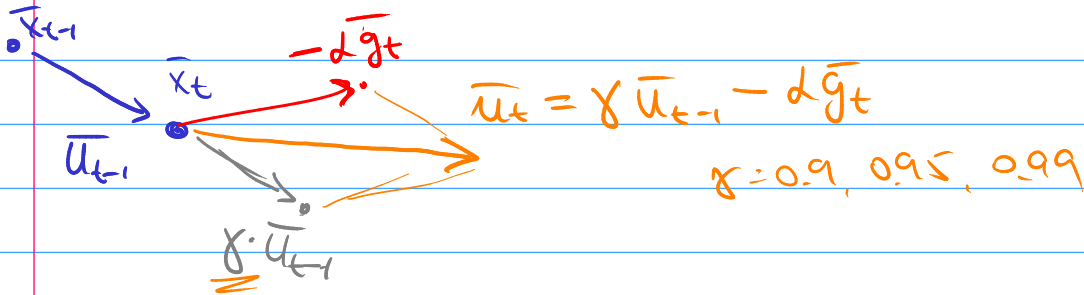


SGD

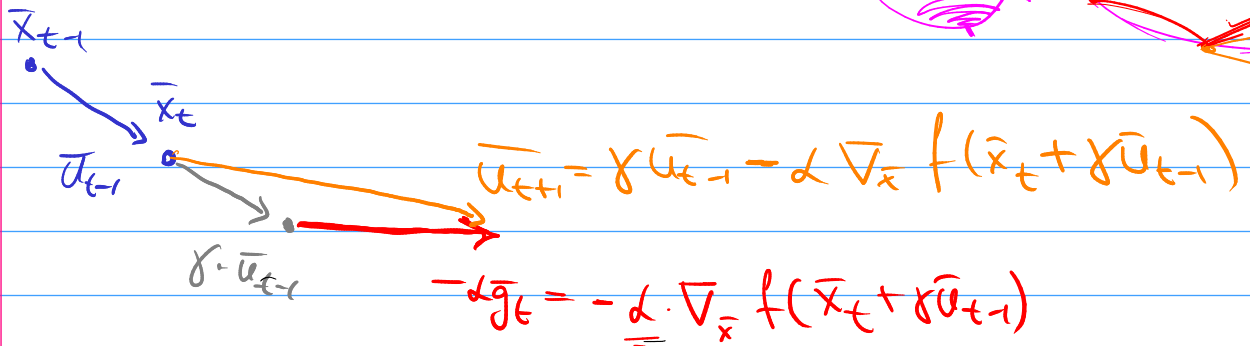


$$\hat{g} = \frac{1}{m} \sum \dots$$

SGD with Momentum



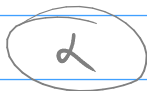
Nesterov accelerated gradients (NAG)



Adaptive SGD

$$\bar{x}_{t+1} = \bar{x}_t - \frac{\alpha}{(m/c)^T} \bar{g}_t$$

Adagrad



$$\bar{G}_0 = 0$$

$$G_{t,i} = G_{t-1,i} + g_{t,i}^2, \quad \bar{g}_t = \nabla_x f(x_t)$$

$$\bar{x}_{t+1} = \bar{x}_t - \frac{\alpha}{\sqrt{G_t + \epsilon}} \bar{g}_t$$

$$c = c - \frac{1}{\sqrt{m/d}} \cdot m/c$$

$$c = c - 1 \cdot m/c$$

RMSprop

$$G_{t,i} = \gamma G_{t-1,i} + (1-\gamma) g_{t,i}^2$$

$$\bar{x}_{t+1} = \bar{x}_t - \frac{\alpha}{\sqrt{G_t + \epsilon}} \bar{g}_t$$

Adadelta

$$\bar{x}_{t+1} = \bar{x}_t - \alpha \cdot \frac{\sqrt{R_{t+1} \epsilon}}{\sqrt{G_{t+1} \epsilon}} \cdot \bar{g}_t$$

$$R_{t,i} = \rho R_{t-1,i} + (1-\rho) \cdot \bar{u}_{t,i}^2$$

Adam

$$\bar{x}_{t+1} = \bar{x}_t - \frac{\alpha}{\sqrt{G_{t+1} \epsilon}} \bar{m}_t, \text{ use}$$

$$G_t = \beta_2 G_{t-1} + (1-\beta_2) \bar{g}_t^2$$

$$\bar{m}_t = \beta_1 \bar{m}_{t-1} + (1-\beta_1) \bar{g}_t$$

$\beta_1 = 0.9$
 $\beta_2 = 0.999$
 $\epsilon = 10^{-8}$

AMSGrad

$$\sqrt{G_{t+1} \epsilon}, \text{ use } G'_t = \max(G_t, G'_{t-1})$$

Nadam

AdamW

1980-e: weight decay \bar{w}_k

$$\bar{w}_{k,t+1} = (1-\beta) \bar{w}_k + \alpha \nabla_w F(\bar{w}_k) =$$

$$= \bar{w}_k - \alpha \left(\nabla_w F + \frac{\beta}{\alpha} \bar{w}_k \right)$$

$$= \bar{w}_k - \alpha \cdot \nabla_w \left(F + \frac{\beta}{2\alpha} \bar{w}^T \bar{w} \right)$$

L_2 -penalty.

