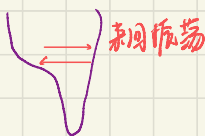




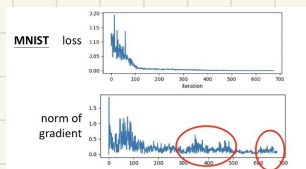
Error surface is fugged ... 误差表面凹凸不平

Adaptive Learning Rate

Training stuck \neq Small Gradient

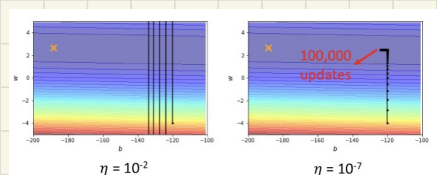
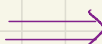
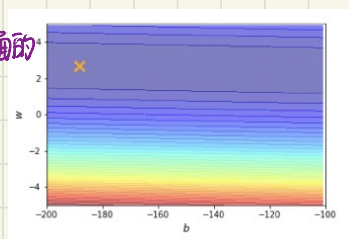


来回振荡



Training stuck without critical points

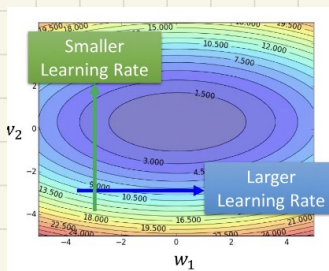
Convex 凸面的



无法靠近 critical points



Different parameters needs different learning rate



update one parameter

$$Q_i^{t+1} \leftarrow Q_i^t - \eta g_i^t$$

$$g_i^t = \left. \frac{dL}{dQ_i} \right|_{Q=Q^t}$$

$$Q_i^{t+1} \leftarrow Q_i^t - \frac{\eta}{\sigma_i^t} g_i^t$$

A. Root Mean Square (均方根)

Adagrad

$$\theta_i^1 \leftarrow \theta_i^0 - \frac{\eta}{\sigma_i^0} g_i^0 \quad \sigma_i^0 = \sqrt{(g_i^0)^2} = |g_i^0|$$

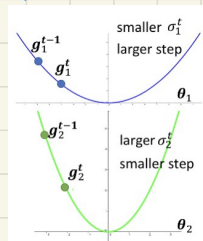
$$\theta_i^2 \leftarrow \theta_i^1 - \frac{\eta}{\sigma_i^1} g_i^1 \quad \sigma_i^1 = \sqrt{\frac{1}{2} [(g_i^0)^2 + (g_i^1)^2]}$$

$$\theta_i^3 \leftarrow \theta_i^2 - \frac{\eta}{\sigma_i^2} g_i^2 \quad \sigma_i^2 = \sqrt{\frac{1}{3} [(g_i^0)^2 + (g_i^1)^2 + (g_i^2)^2]}$$

...

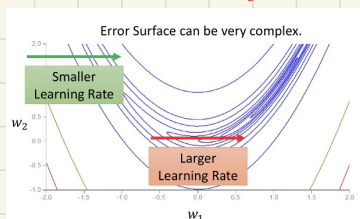
$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t \quad \sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g_i)^2}$$

Used in Adagrad



b. RMSProp

Learning rate adapt dynamically

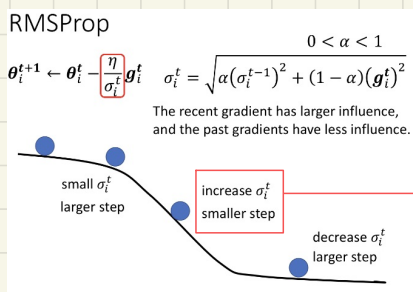


在绿色箭头上 需要 small LR
在红色箭头上 需要 larger LR
需要在同一参数同一方向, LR也能动态调整 \Rightarrow RMSProp

RMSProp

$$\begin{aligned}\theta_i^1 &\leftarrow \theta_i^0 - \frac{\eta}{\sigma_i^0} g_i^0 & \sigma_i^0 &= \sqrt{(g_i^0)^2} & 0 < \alpha < 1 \\ \theta_i^2 &\leftarrow \theta_i^1 - \frac{\eta}{\sigma_i^1} g_i^1 & \sigma_i^1 &= \sqrt{\alpha(\sigma_i^0)^2 + (1-\alpha)(g_i^1)^2} \\ \theta_i^3 &\leftarrow \theta_i^2 - \frac{\eta}{\sigma_i^2} g_i^2 & \sigma_i^2 &= \sqrt{\alpha(\sigma_i^1)^2 + (1-\alpha)(g_i^2)^2} \\ &\vdots & & \\ \theta_i^{t+1} &\leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t & \sigma_i^t &= \sqrt{\alpha(\sigma_i^{t-1})^2 + (1-\alpha)(g_i^t)^2}\end{aligned}$$

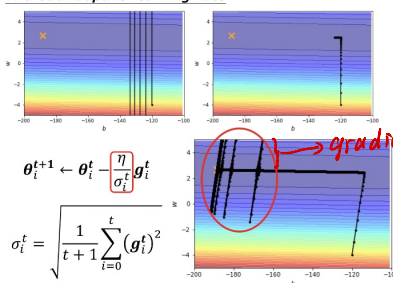
$\alpha \rightarrow 0$ g_i^t 相较于之前计算的 gradient 更重要
 $\alpha \rightarrow 1$ 现在算出的 g_i^t 不重要, 之前的重要



\rightarrow 友列第三个球, 按照 Ada grad 的话 LR 会很大, 导致第三个球飞出; 所以得控制, α 可以用来控制 α

Adam: RMSProp + Momentum

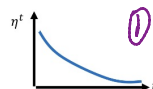
Without Adaptive Learning Rate



$$\begin{aligned}\theta_i^{t+1} &\leftarrow \theta_i^t - \frac{\eta}{\sigma_i^t} g_i^t \\ \sigma_i^t &= \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g_i^t)^2}\end{aligned}$$

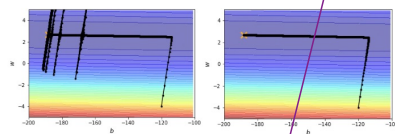
Learning Rate Scheduling

Learning Rate Scheduling



Learning Rate Decay

After the training goes, we are close to the destination, so we reduce the learning rate.



Summary of optimization

(Vanilla) Gradient Descent

$$\theta_i^{t+1} \leftarrow \theta_i^t - \eta g_i^t$$

Various Improvements

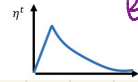
$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta^t}{\sigma_i^t} m_i^t$$

Learning rate scheduling

Momentum: weighted sum of the previous gradients

root mean square of the gradients

为按大小和方向
只考虑 gradient 大小



2

Warm Up

Increase and then decrease?

At the beginning, the estimate of σ_i^t has large variance.