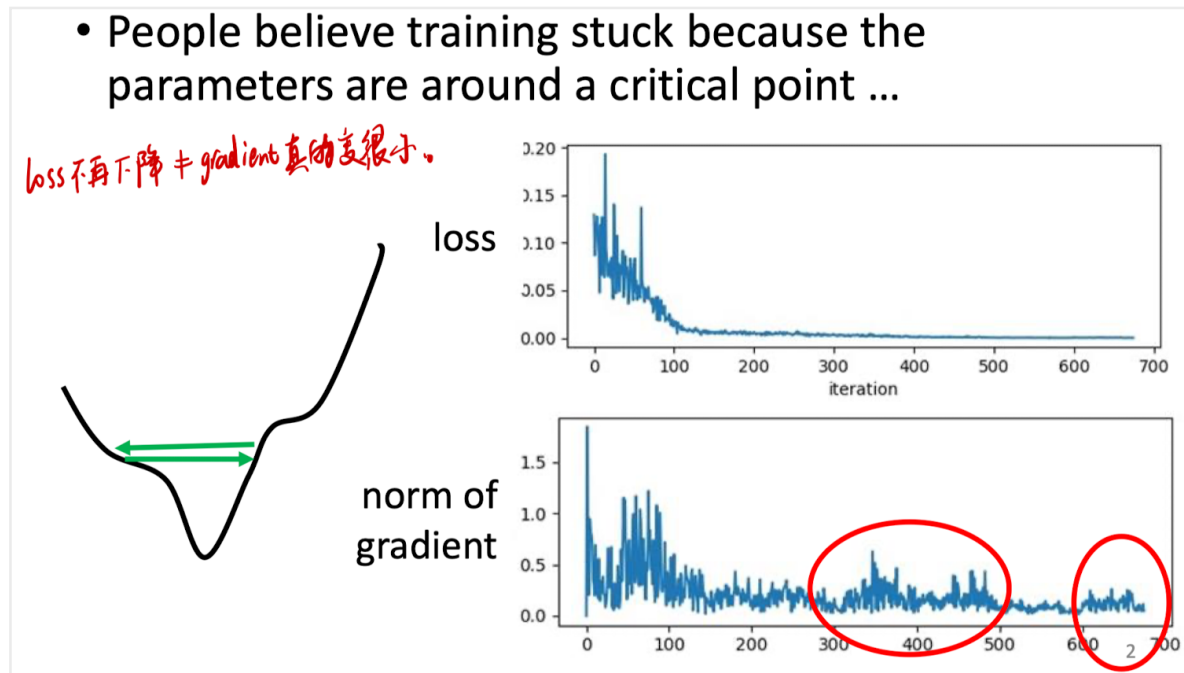


2.3 Adaptive Learning Rate

- Root Mean Square
- RMSProp
- Adam: RMSProp + Momentum

loss 不再下降并不意味着 gradient 真的变得很小，他可能在震荡，卡住的时候可以算 gradient 的 norm。

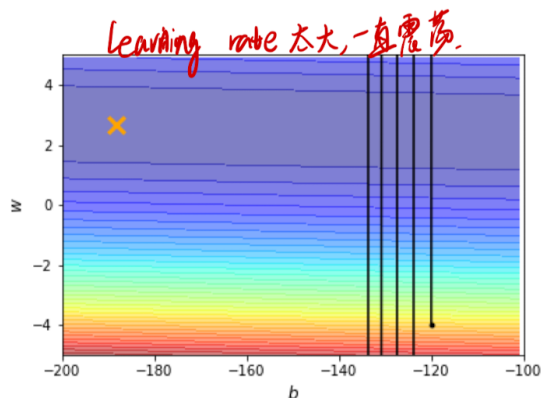


没有 critical point 时候训练也可能是困难。下图左上角 X 是 error surface 的最低点，这个 error surface 是 convex 的。

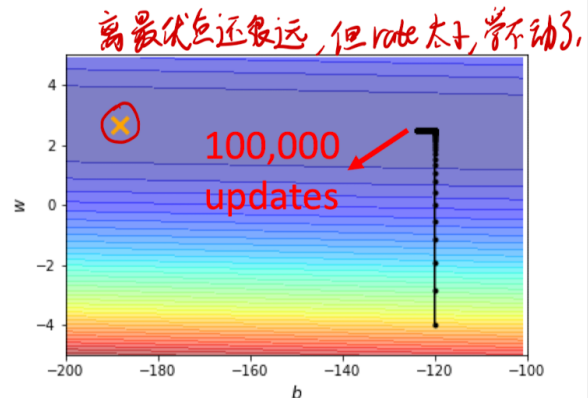
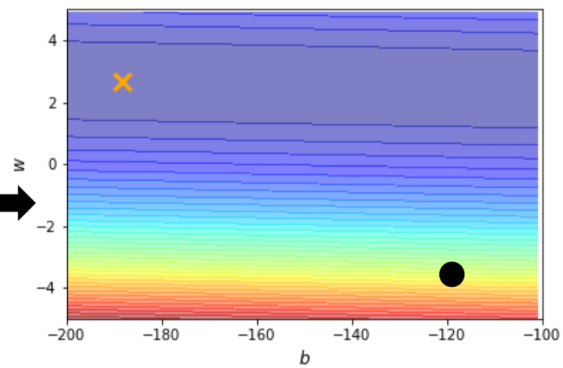
Training can be difficult even without critical points.

This error surface is convex. →

Learning rate **cannot** be one-size-fits-all



$$\eta = 10^{-2}$$



$$\eta = 10^{-7}$$

4

learning rate 的特制化: learning rate 需要自动根据 gradient 的大小做调整
如何做: 加上一个 dependant on 参数 i 并且 dependant on 于 iteration t .
parameter dependant & iteration dependant.

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

$$g_i^t = \frac{\partial L}{\partial \theta_i} \big|_{\theta = \theta^t}$$

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

Parameter
dependent

parameter dependant 的常见计算方式

1.Root Mean Square

①

Root Mean Square

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

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

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

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

⋮

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

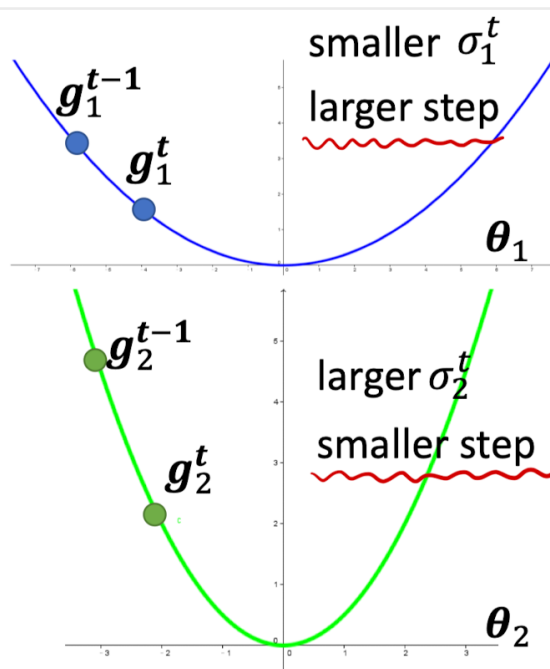
6

why?

参数1的坡度平坦，平均 gradient 小，step 更大；参数2坡度陡峭，平均 gradient 大，step 小

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

$$\sigma_i^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g_i^t)^2}$$

Used in **Adagrad**

2.RMSProp

第一个方法的问题 /；就算是一个参数，它需要的 learning rate 也会随着时间而改变

②

RMSProp

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

$$\theta_i^1 \leftarrow \theta_i^0 - \frac{\eta}{\sigma_i^0} g_i^0 \quad \sigma_i^0 = \sqrt{(g_i^0)^2} \quad \alpha \text{ 要自己调} \quad 0 < \alpha < 1$$

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

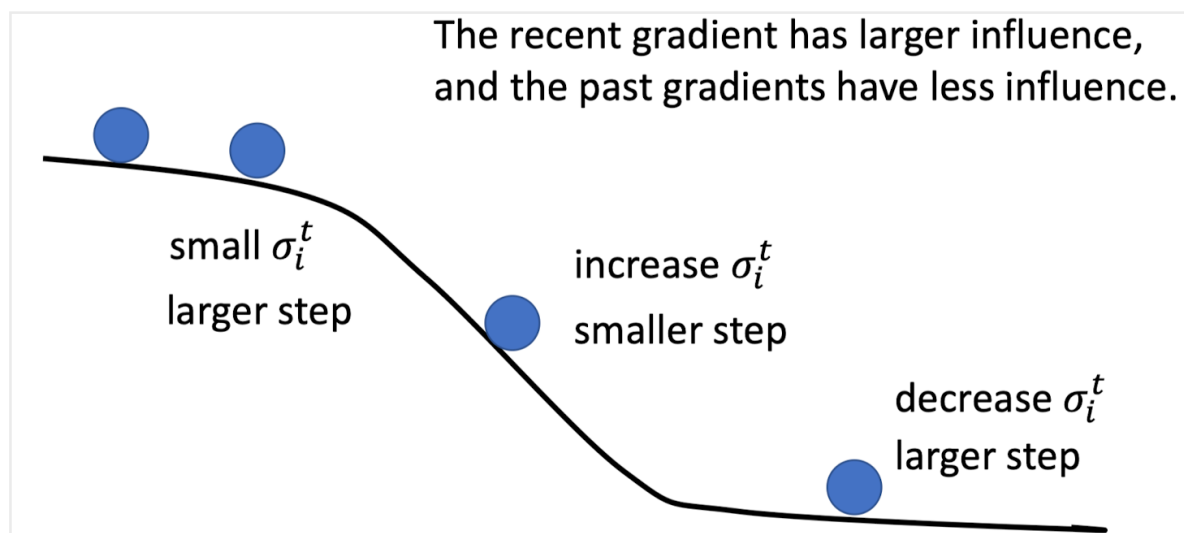
$$\theta_i^3 \leftarrow \theta_i^2 - \frac{\eta}{\sigma_i^2} g_i^2 \quad \sigma_i^2 = \sqrt{\alpha(\sigma_i^1)^2 + (1 - \alpha)(g_i^2)^2}$$

⋮

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

可以自己调整现在的 gradient 的权重，从而动态调整 learning rate。

如下图，当走到下坡中间的时候，如果是 adagrad，它反应会比较慢，step 会很大；但如果用 RMSProp，把 α 的值设小一点，让新的刚看到的 gradient 影响比较大的话，就可以很快让分母变大，就可以很快的让步伐变小。



3.Adam: RMSProp + Momentum

Algorithm 1: *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector) \rightarrow for momentum

$v_0 \leftarrow 0$ (Initialize 2nd moment vector) \rightarrow for RMSprop

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

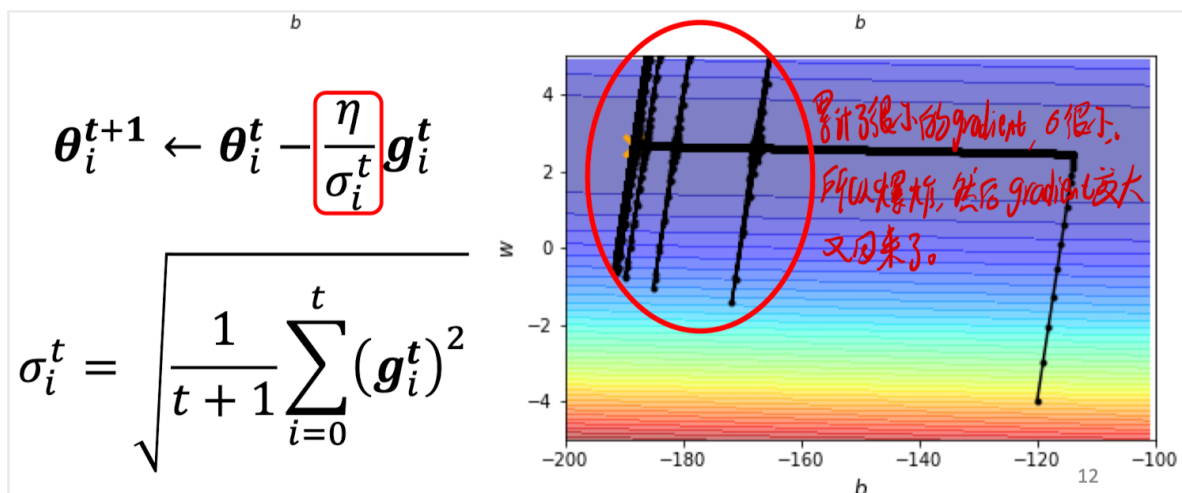
$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

11

adagrad 存在的的问题：会突然爆炸



Learning rate scheduling

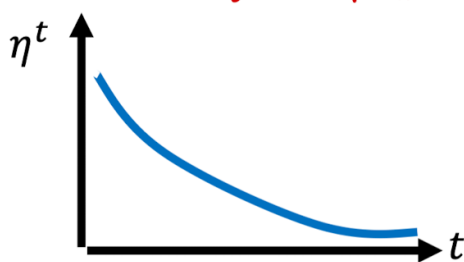
η 不是固定的值，而是和时间有关的，不是常数。

- Learning Rate Decay
- Warm up

Learning Rate Scheduling

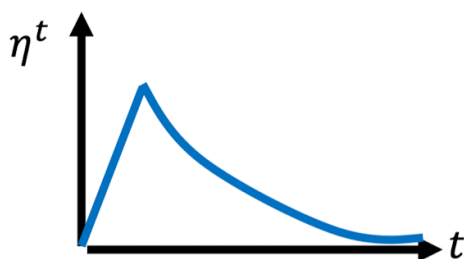
$$\theta_i^{t+1} \leftarrow \theta_i^t - \frac{\eta^t}{\sigma_i^t} g_i^t$$

η与时间有关, 越后面越小。



① Learning Rate Decay

As the training goes, we are closer to the destination, so we reduce the learning rate.



② Warm Up

Increase and then decrease?

learning rate 先变大后变小。

Summary

(Vanilla) Gradient Descent

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

Various Improvements

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

③ Learning rate scheduling

④ Momentum: weighted sum of the previous gradients
考虑大小与方向

② root mean square of the gradients
只考虑大小

Consider direction

only magnitude

