# Error surface is rugged …

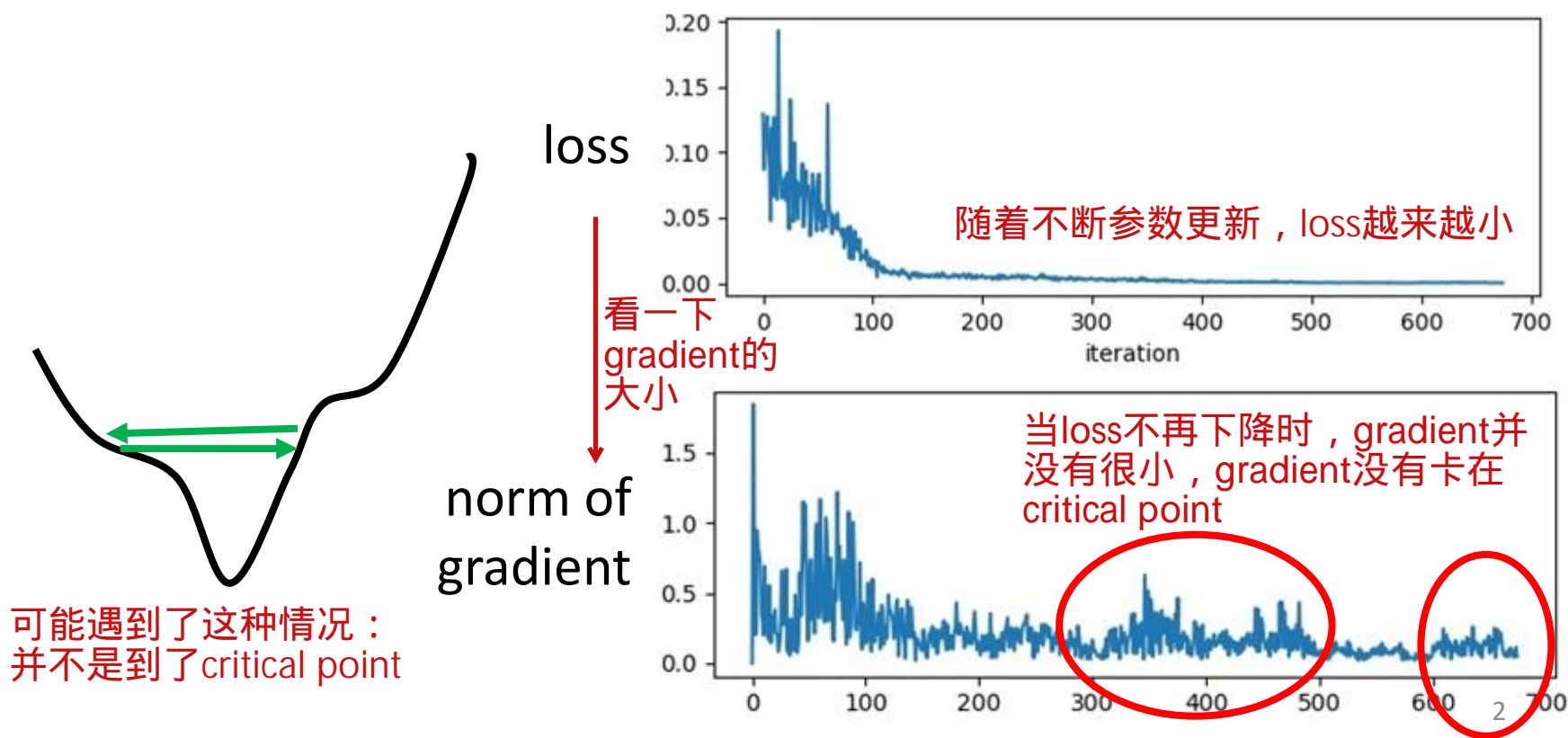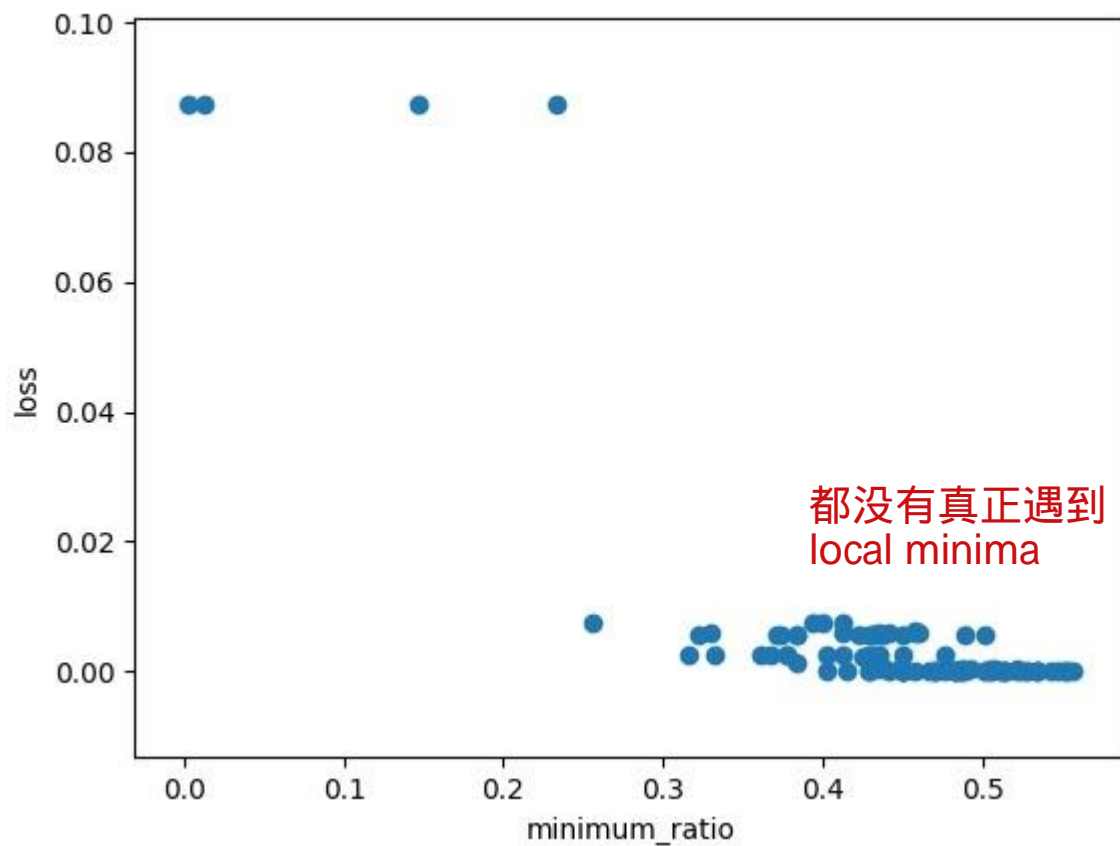Tips for training: **Adaptive** **Learning Rate**

# Training stuck ≠ Small Gradient

critical points

- People believe training stuck because the parameters are around a critical point …



loss

gradient

norm of gradient

critical point

loss

loss

gradient

gradient

critical point

2

# Wait a minute …



local minima

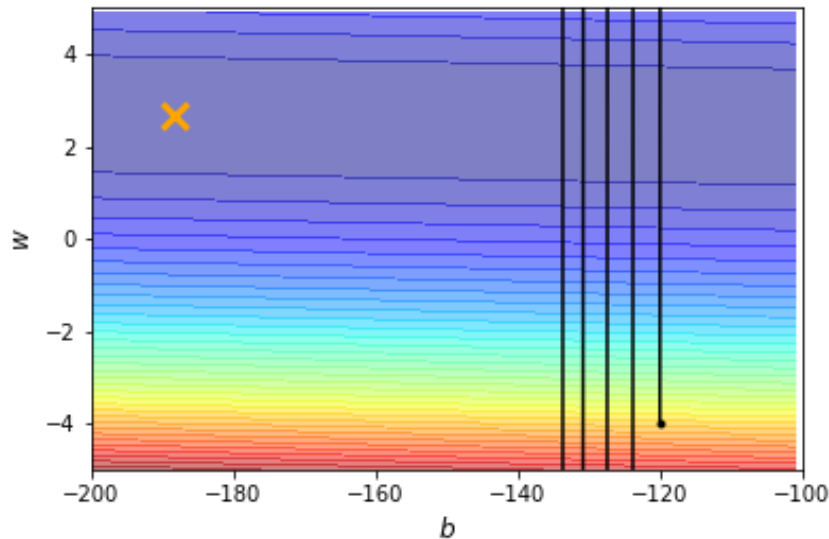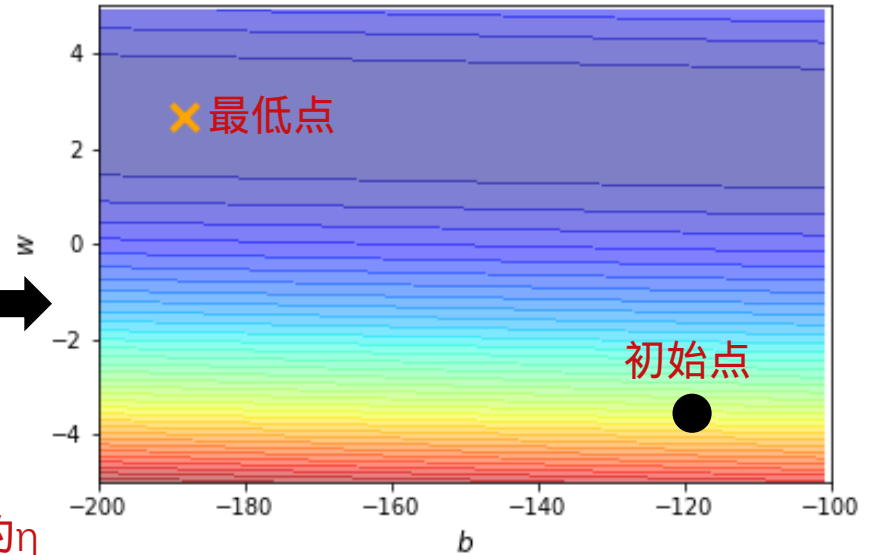# *Training can be difficult even without critical points.*

This error surface is <mark>convex</mark>. ➡️



η

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



$\eta = 10^{-2}$

learning rate
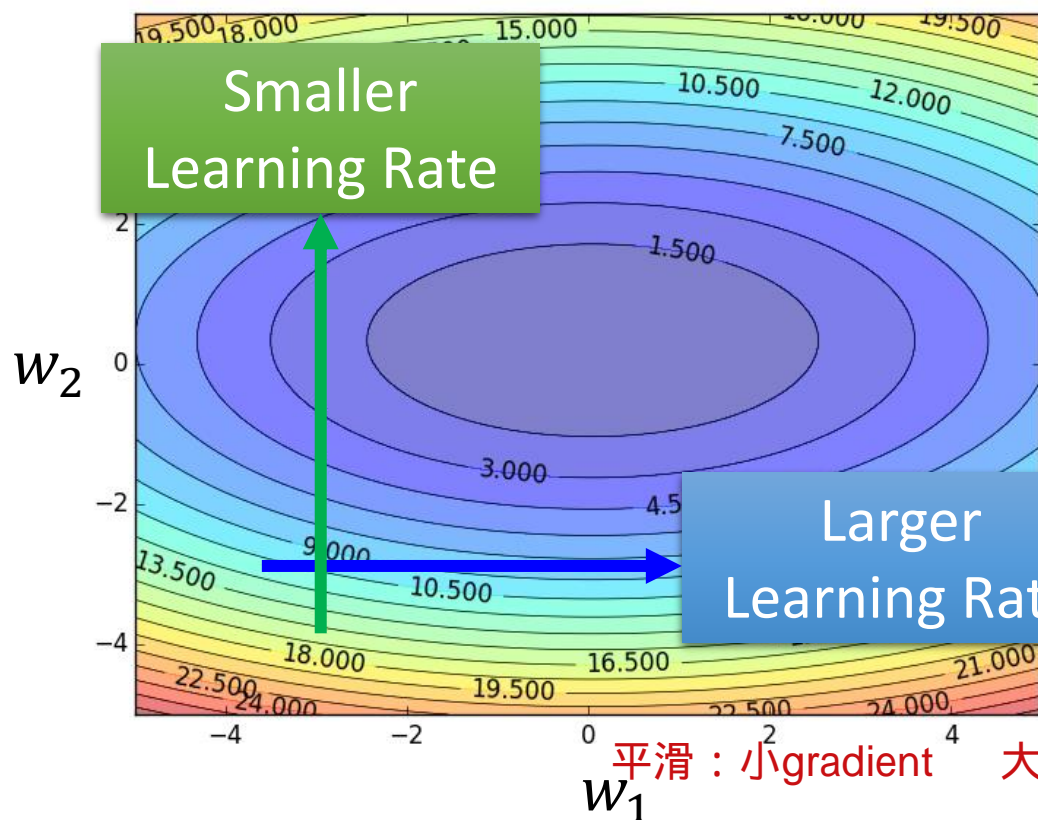


100,000 updates

$\eta = 10^{-7}$

learning rate

4

# Different parameters needs different learning rate

gradient      learning rate

Smaller Learning Rate

Larger Learning Rate

$w_2$

$w_1$

gradient     learning rate

Formulation for **one** parameter:

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\eta}\, \boldsymbol{g}_i^t$$

learning rate

$$\boldsymbol{g}_i^t = \frac{\partial L}{\partial \boldsymbol{\theta}_i}\big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^t}$$

Adaptive learning rate

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\frac{\eta}{\sigma_i^t}}\, \boldsymbol{g}_i^t$$

Parameter dependent

i → parameter dependent
t → iteration dependent

5

# Root Mean Square

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

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

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

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

$$\vdots$$

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

AdaGrad

gradient
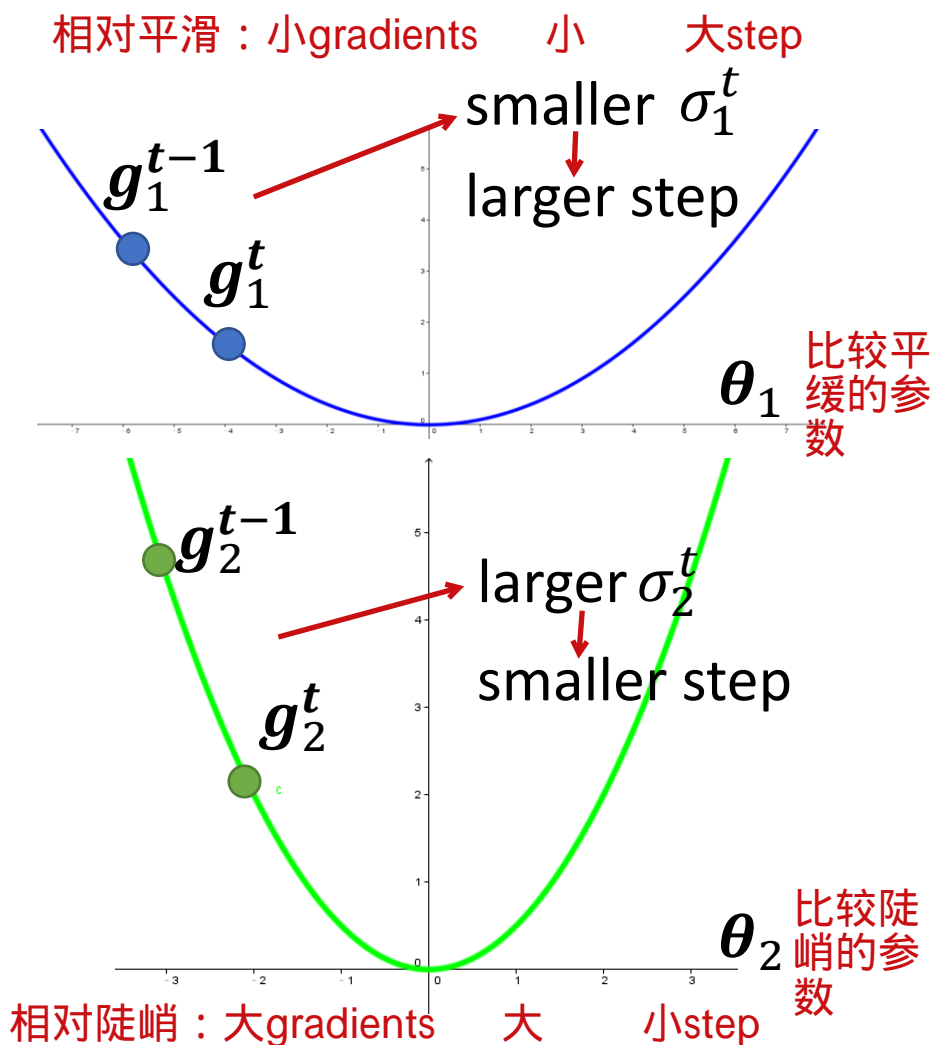
6

# Root Mean Square

smaller $\sigma_1^t$

larger step

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\frac{\eta}{\sigma_i^t}} \boldsymbol{g}_i^t$$

$\boldsymbol{g}_1^{t-1}$

$\boldsymbol{g}_1^t$

$\boldsymbol{\theta}_1$

$$\sigma_i^t = \sqrt{\frac{1}{t+1}\sum_{i=0}^{t}\left(\boldsymbol{g}_i^t\right)^2}$$

$\boldsymbol{g}_2^{t-1}$

larger $\sigma_2^t$

smaller step

$\boldsymbol{g}_2^t$

## Used in **Adagrad**

$\boldsymbol{\theta}_2$

gradients      step

7

# Learning rate adapts dynamically

Error Surface can be very complex.



Smaller Learning Rate

Larger Learning Rate

$w_2$

$w_1$

# RMSProp

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\frac{\eta}{\sigma_i^t}} g_i^t$$

hyperparameter
gradient
gradient
gradient

$$0 < \alpha < 1$$

$$\boldsymbol{\theta}_i^1 \leftarrow \boldsymbol{\theta}_i^0 - \frac{\eta}{\sigma_i^0} \boldsymbol{g}_i^0 \qquad \sigma_i^0 = \sqrt{(\boldsymbol{g}_i^0)^2}$$

$$\boldsymbol{\theta}_i^2 \leftarrow \boldsymbol{\theta}_i^1 - \frac{\eta}{\sigma_i^1} \boldsymbol{g}_i^1 \qquad \sigma_i^1 = \sqrt{\alpha(\sigma_i^0)^2 + (1-\alpha)(\boldsymbol{g}_i^1)^2}$$

$$\boldsymbol{\theta}_i^3 \leftarrow \boldsymbol{\theta}_i^2 - \frac{\eta}{\sigma_i^2} \boldsymbol{g}_i^2 \qquad \sigma_i^2 = \sqrt{\alpha(\sigma_i^1)^2 + (1-\alpha)(\boldsymbol{g}_i^2)^2}$$

$$\vdots$$

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\eta}{\sigma_i^t} \boldsymbol{g}_i^t \qquad \sigma_i^t = \sqrt{\alpha(\sigma_i^{t-1})^2 + (1-\alpha)(\boldsymbol{g}_i^t)^2}$$

Adagrad                gradient

9

# RMSProp

$$g_i^1 \ g_i^2 \ \cdots\cdots \ g_i^{t-1}$$

$$0 < \alpha < 1$$

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

The recent gradient has larger influence, and the past gradients have less influence.

small $\sigma_i^t$

larger step

gradient

step

increase $\sigma_i^t$

smaller step

decrease $\sigma_i^t$

larger step

10

# 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 $\beta_1^t$ and $\beta_2^t$ we denote $\beta_1$ and $\beta_2$ to the power $t$.

**Require:** $\alpha$: Stepsize
**Require:** $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates
**Require:** $f(\theta)$: Stochastic objective function with parameters $\theta$
**Require:** $\theta_0$: Initial parameter vector
   $m_0 \leftarrow 0$ (Initialize $1^{st}$ moment vector)  → for momentum
   $v_0 \leftarrow 0$ (Initialize $2^{nd}$ moment vector)  → for RMSprop
   $t \leftarrow 0$ (Initialize timestep)
  **while** $\theta_t$ 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)
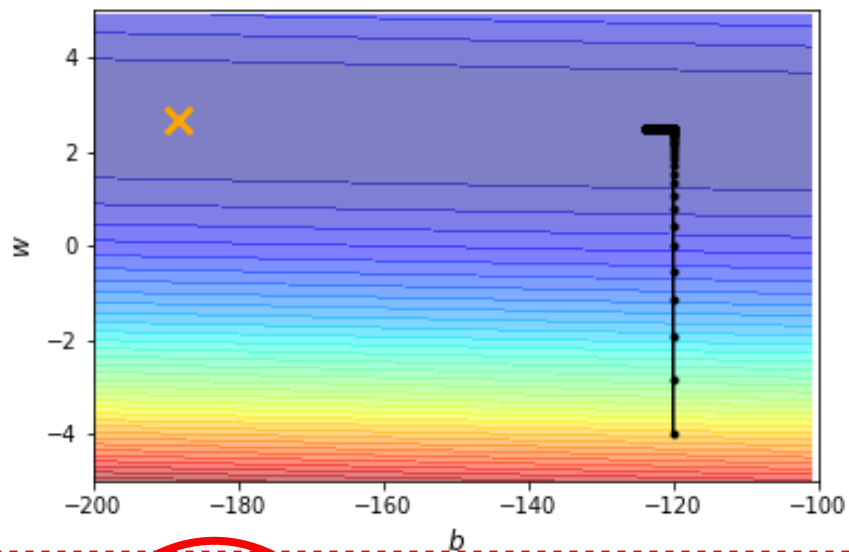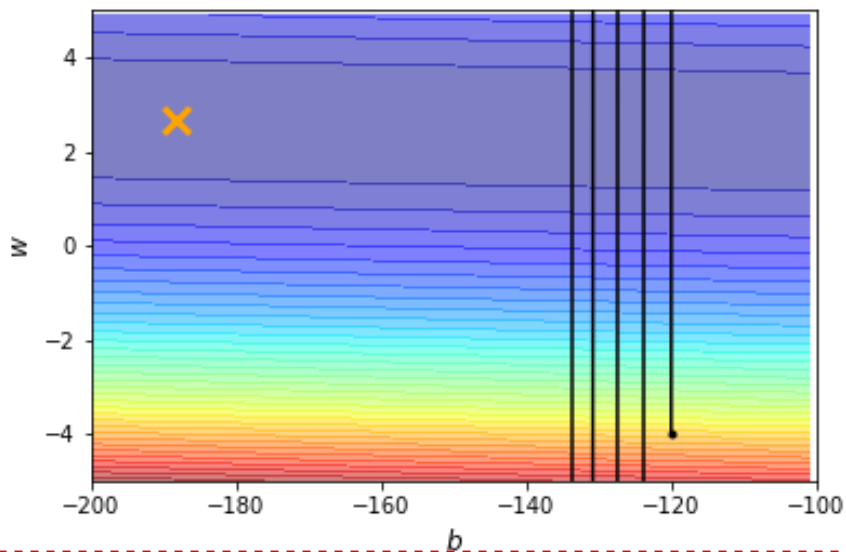    $\widehat{m}_t \leftarrow m_t/(1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
    $\widehat{v}_t \leftarrow v_t/(1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
    $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \widehat{m}_t/(\sqrt{\widehat{v}_t} + \epsilon)$ (Update parameters)
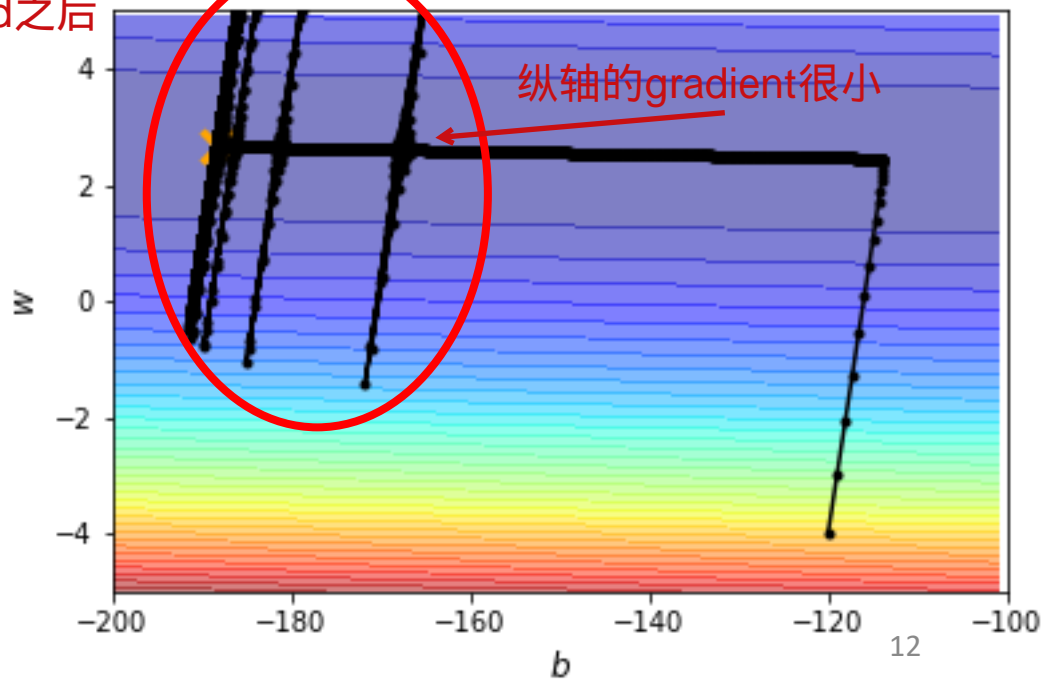  **end while**
  **return** $\theta_t$ (Resulting parameters)

# *Without Adaptive Learning Rate*



Adagrad

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \boxed{\frac{\eta}{\sigma_i^t}}\boldsymbol{g}_i^t$$

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

gradient

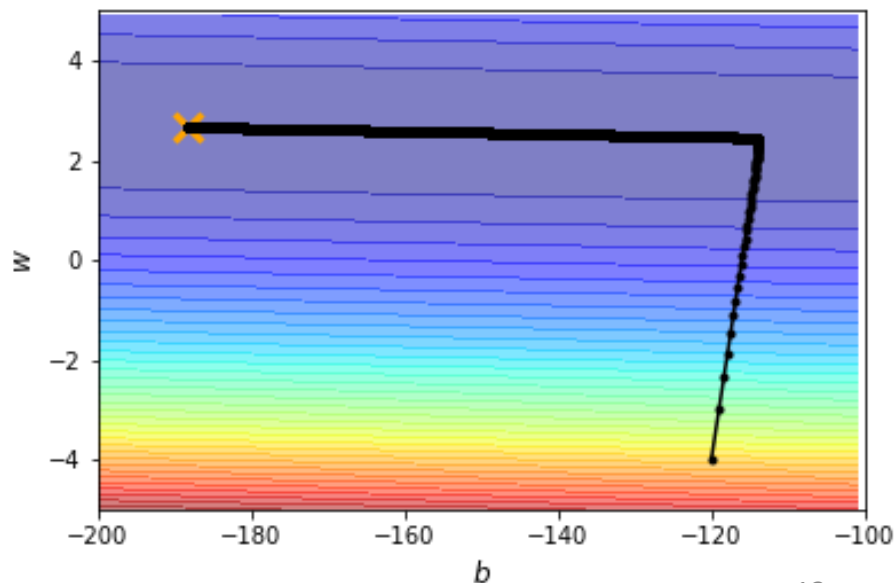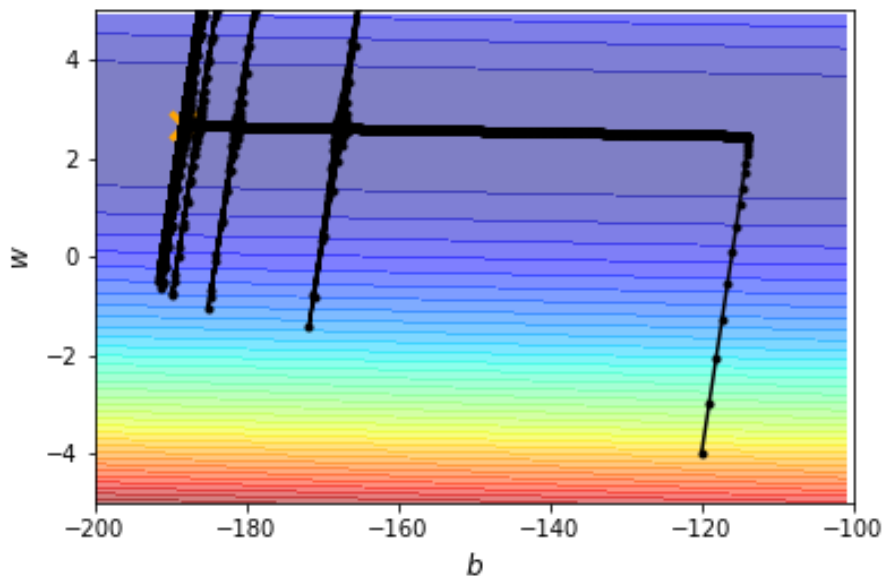# *Learning Rate Scheduling*

learning rate

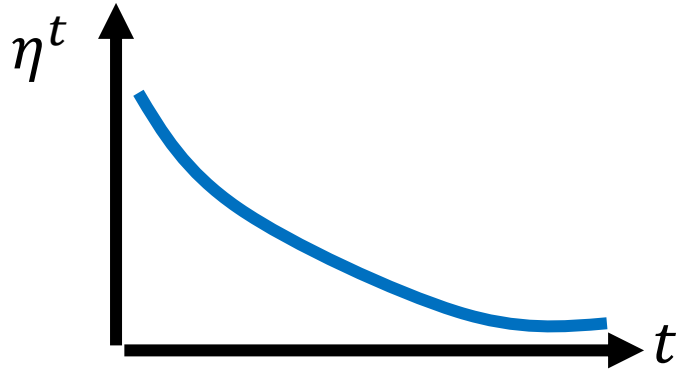$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\eta^t}{\sigma_i^t} g_i^t$$



1. *Learning Rate Decay*

   As the training goes, we are closer to the destination, so we reduce the learning rate.
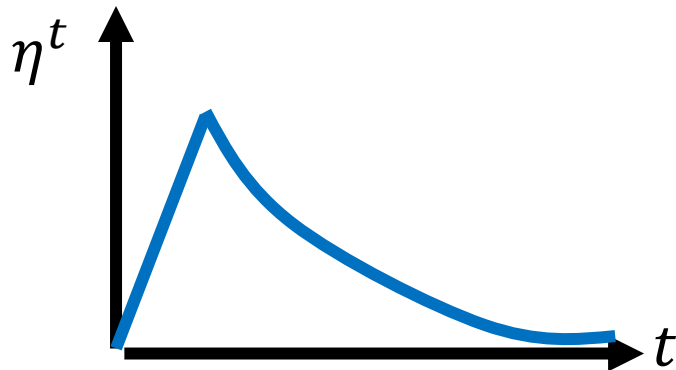
# Learning Rate Scheduling

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\boxed{\eta^t}}{\sigma_i^t} \boldsymbol{g}_i^t$$



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?

warmup

We further explore $n = 18$ that leads to a 110-layer ResNet. In this case, we find that the initial learning rate of 0.1 is slightly too large to start converging[5]. So we use 0.01 to warm up the training until the training error is below 80% (about 400 iterations), and then go back to 0.1 and continue training. The rest of the learning schedule is as done previously. This 110-layer network converges well (Fig. 6, middle). It has *fewer* parameters than other deep and thin

[5]With an initial learning rate of 0.1, it starts converging (<90% error) after several epochs, but still reaches similar accuracy.

**Residual Network**

https://arxiv.org/abs/1512.03385
2015   12

## 5.3   Optimizer

We used the Adam optimizer [17] with $\beta_1 = 0.9$, $\beta_2 = 0.98$ and $\epsilon = 10^{-9}$. We varied the learning rate over the course of training, according to the formula:

$$lrate = d_{\text{model}}^{-0.5} \cdot \min(step\_num^{-0.5}, step\_num \cdot warmup\_steps^{-1.5}) \qquad (3)$$

This corresponds to increasing the learning rate linearly for the first $warmup\_steps$ training steps, and decreasing it thereafter proportionally to the inverse square root of the step number. We used $warmup\_steps = 4000$.

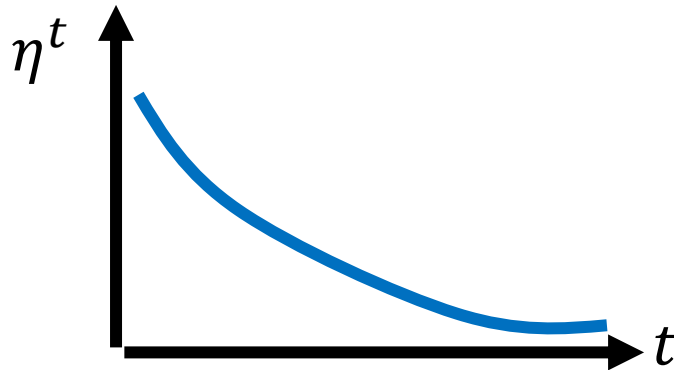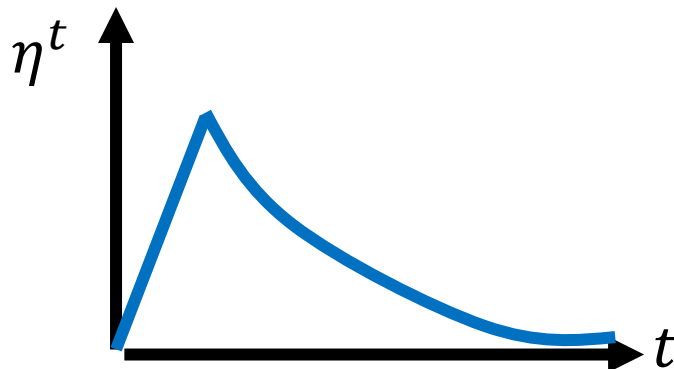**Transformer**   https://arxiv.org/abs/1706.03762

15

# Learning Rate Scheduling

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\boxed{\eta^t}}{\sigma_i^t} \boldsymbol{g}_i^t$$

## Learning Rate Decay

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

## Warm Up  warmup

Increase and then decrease?

At the beginning, the estimate of $\sigma_i^t$ has large variance.

Please refer to **RAdam**  https://arxiv.org/abs/1908.03265
Rectified Adam

# Summary of Optimization

## (Vanilla) Gradient Descent

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \eta \boldsymbol{g}_i^t$$

## Various Improvements

$$\boldsymbol{\theta}_i^{t+1} \leftarrow \boldsymbol{\theta}_i^t - \frac{\eta^t}{\sigma_i^t} \boldsymbol{m}_i^t$$

Learning rate scheduling

Momentum: weighted sum of the previous gradients

Consider direction

root mean square of the gradients

only magnitude

# To Learn More ……

optimizers for deep learning



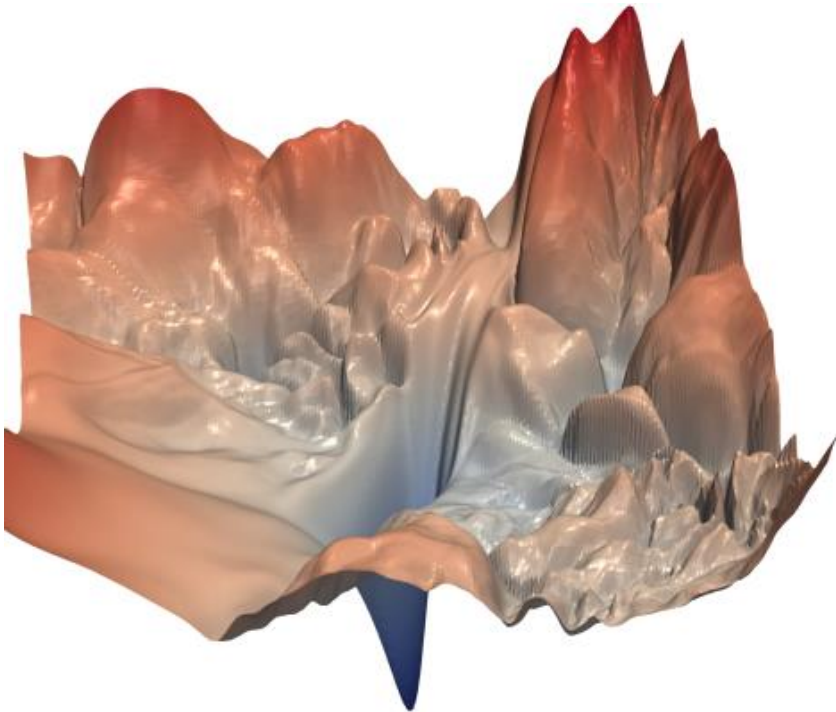https://youtu.be/4pUmZ8hXlHM

(in Mandarin)



https://youtu.be/e03YKGHXnL8

(in Mandarin)

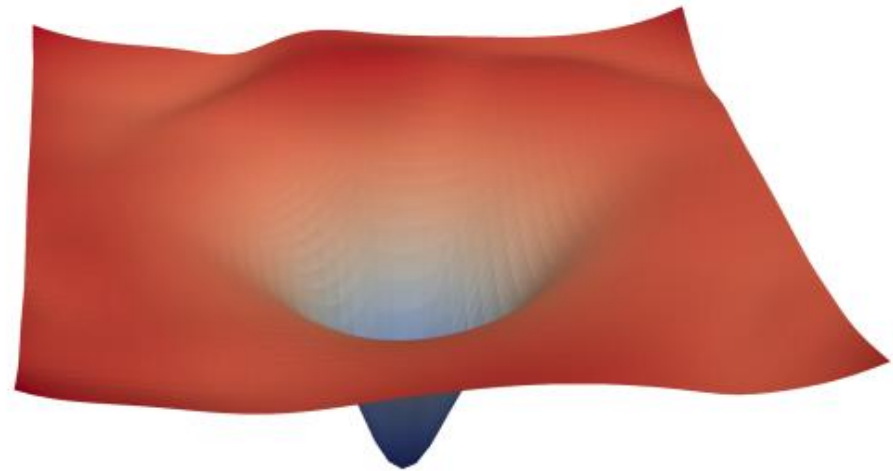# Next Time

Next time



Better optimization strategies:
If the mountain won't move,
build a road around it.

Can we change the error
surface?
Directly move the mountain!