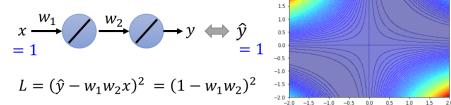


# OPTIMIZATION

- Critical Point → Batch & Momentum
- Learning Rate → Optimizer

## CRITICAL POINT Gradient = 0



$$L(\hat{y}) = (\hat{y} - w_1 w_2 x)^2 = (1 - w_1 w_2)^2$$

$$\frac{\partial L}{\partial w_1} = 2(1 - w_1 w_2)(-w_2) = 0 \quad \text{Critical point: } w_1 = 0, w_2 = 0$$

$$\frac{\partial L}{\partial w_2} = 2(1 - w_1 w_2)(-w_1) = 0$$

**Saddle point**

Local Minima 局部最小值

Saddle Point 轮点, ↗ Hessian 判别

$$H_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} L(\theta')$$

Telling Properties of Critical Points

通过计算 Hessian 矩阵判断是 Saddle Point or Local Minima/Maxima (Eigen Vector & Value)

计算更新方向逃离鞍点  
低维的 Local Minima 高维可能会变成 Saddle Point

H may tell us parameter update direction!

$$u \text{ is an eigen vector of } H \rightarrow u^T H u = u^T (\lambda u) = \lambda \|u\|^2$$

$$\lambda \text{ is the eigen value of } u < 0 < 0$$

$$L(\theta) \approx L(\theta') + \frac{1}{2} (\theta - \theta')^T H (\theta - \theta') \rightarrow L(\theta) < L(\theta')$$

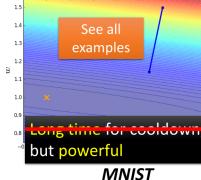
$$\theta - \theta' = u \quad \theta = \theta' + u \quad \text{Decrease } L$$

## Batch Small Batch Size to Escape

Consider 20 examples (N=20)

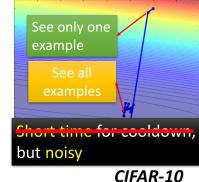
Batch size = N (Full Batch)

Update after seeing all the 20 examples



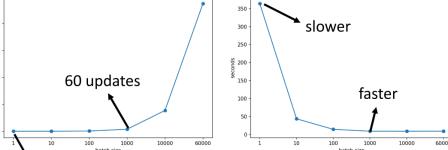
Batch size = 1

Update for each example  
Update 20 times in an epoch



Smaller batch requires longer time for one epoch (longer time for seeing all data once)

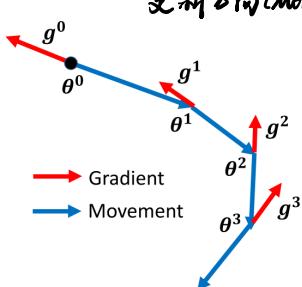
Time for one update



- Smaller batch size has better performance
- What's wrong with large batch size? Optimization Fails

## Momentum 惯性 (Vanilla) Gradient Descent

看直角坐标系



Movement = Negative of  $\partial L / \partial w$

Movement = Negative of  $\partial L / \partial w + \text{Last Movement}$

### Gradient Descent + Momentum

Starting at  $\theta^0$

Compute gradient  $g^0$

Move to  $\theta^1 = \theta^0 - \eta g^0$

Compute gradient  $g^1$

Move to  $\theta^2 = \theta^1 - \eta g^1$

⋮

Movement: movement of last step minus gradient at present

Starting at  $\theta^0$

Compute gradient  $g^0$

Move to  $\theta^1 = \theta^0 - \eta g^0$

Compute gradient  $g^1$

Move to  $\theta^2 = \theta^1 - \eta g^1$

⋮

- Starting at  $\theta^0$
- Movement  $m^0 = 0$
- Compute gradient  $g^0$
- Movement  $m^1 = \lambda m^0 - \eta g^0$
- Move to  $\theta^1 = \theta^0 + m^1$
- Compute gradient  $g^1$
- Movement  $m^2 = \lambda m^1 - \eta g^1$
- Move to  $\theta^2 = \theta^1 + m^2$
- Movement not just based on gradient, but previous movement.

oss

$m^0 = 0$

$m^1 = -\eta g^0$

$m^2 = -\lambda \eta g^0 - \eta g^1$

⋮

$m^3 = -\lambda \eta g^1 - \eta g^2$

⋮

$m^4 = -\lambda \eta g^2 - \eta g^3$

⋮

$m^5 = -\lambda \eta g^3 - \eta g^4$

⋮

$m^6 = -\lambda \eta g^4 - \eta g^5$

⋮

$m^7 = -\lambda \eta g^5 - \eta g^6$

⋮

$m^8 = -\lambda \eta g^6 - \eta g^7$

⋮

$m^9 = -\lambda \eta g^7 - \eta g^8$

⋮

$m^{10} = -\lambda \eta g^8 - \eta g^9$

⋮

$m^{11} = -\lambda \eta g^9 - \eta g^{10}$

⋮

$m^{12} = -\lambda \eta g^{10} - \eta g^{11}$

⋮

$m^{13} = -\lambda \eta g^{11} - \eta g^{12}$

⋮

$m^{14} = -\lambda \eta g^{12} - \eta g^{13}$

⋮

$m^{15} = -\lambda \eta g^{13} - \eta g^{14}$

⋮

$m^{16} = -\lambda \eta g^{14} - \eta g^{15}$

⋮

$m^{17} = -\lambda \eta g^{15} - \eta g^{16}$

⋮

$m^{18} = -\lambda \eta g^{16} - \eta g^{17}$

⋮

$m^{19} = -\lambda \eta g^{17} - \eta g^{18}$

⋮

$m^{20} = -\lambda \eta g^{18} - \eta g^{19}$

⋮

$m^{21} = -\lambda \eta g^{19} - \eta g^{20}$

⋮

$m^{22} = -\lambda \eta g^{20} - \eta g^{21}$

⋮

$m^{23} = -\lambda \eta g^{21} - \eta g^{22}$

⋮

$m^{24} = -\lambda \eta g^{22} - \eta g^{23}$

⋮

$m^{25} = -\lambda \eta g^{23} - \eta g^{24}$

⋮

$m^{26} = -\lambda \eta g^{24} - \eta g^{25}$

⋮

$m^{27} = -\lambda \eta g^{25} - \eta g^{26}$

⋮

$m^{28} = -\lambda \eta g^{26} - \eta g^{27}$

⋮

$m^{29} = -\lambda \eta g^{27} - \eta g^{28}$

⋮

$m^{30} = -\lambda \eta g^{28} - \eta g^{29}$

⋮

$m^{31} = -\lambda \eta g^{29} - \eta g^{30}$

⋮

$m^{32} = -\lambda \eta g^{30} - \eta g^{31}$

⋮

$m^{33} = -\lambda \eta g^{31} - \eta g^{32}$

⋮

$m^{34} = -\lambda \eta g^{32} - \eta g^{33}$

⋮

$m^{35} = -\lambda \eta g^{33} - \eta g^{34}$

⋮

$m^{36} = -\lambda \eta g^{34} - \eta g^{35}$

⋮

$m^{37} = -\lambda \eta g^{35} - \eta g^{36}$

⋮

$m^{38} = -\lambda \eta g^{36} - \eta g^{37}$

⋮

$m^{39} = -\lambda \eta g^{37} - \eta g^{38}$

⋮

$m^{40} = -\lambda \eta g^{38} - \eta g^{39}$

⋮

$m^{41} = -\lambda \eta g^{39} - \eta g^{40}$

⋮

$m^{42} = -\lambda \eta g^{40} - \eta g^{41}$

⋮

$m^{43} = -\lambda \eta g^{41} - \eta g^{42}$

⋮

$m^{44} = -\lambda \eta g^{42} - \eta g^{43}$

⋮

$m^{45} = -\lambda \eta g^{43} - \eta g^{44}$

⋮

$m^{46} = -\lambda \eta g^{44} - \eta g^{45}$

⋮

$m^{47} = -\lambda \eta g^{45} - \eta g^{46}$

⋮

$m^{48} = -\lambda \eta g^{46} - \eta g^{47}$

⋮

$m^{49} = -\lambda \eta g^{47} - \eta g^{48}$

⋮

$m^{50} = -\lambda \eta g^{48} - \eta g^{49}$

⋮

$m^{51} = -\lambda \eta g^{49} - \eta g^{50}$

⋮

$m^{52} = -\lambda \eta g^{50} - \eta g^{51}$

⋮

$m^{53} = -\lambda \eta g^{51} - \eta g^{52}$

⋮

$m^{54} = -\lambda \eta g^{52} - \eta g^{53}$

⋮

$m^{55} = -\lambda \eta g^{53} - \eta g^{54}$

⋮

$m^{56} = -\lambda \eta g^{54} - \eta g^{55}$

⋮

$m^{57} = -\lambda \eta g^{55} - \eta g^{56}$

⋮

$m^{58} = -\lambda \eta g^{56} - \eta g^{57}$

⋮

$m^{59} = -\lambda \eta g^{57} - \eta g^{58}$

⋮

$m^{60} = -\lambda \eta g^{58} - \eta g^{59}$

⋮

$m^{61} = -\lambda \eta g^{59} - \eta g^{60}$

⋮

$m^{62} = -\lambda \eta g^{60} - \eta g^{61}$

⋮

$m^{63} = -\lambda \eta g^{61} - \eta g^{62}$

⋮

$m^{64} = -\lambda \eta g^{62} - \eta g^{63}$

⋮

$m^{65} = -\lambda \eta g^{63} - \eta g^{64}$

⋮

$m^{66} = -\lambda \eta g^{64} - \eta g^{65}$

⋮

$m^{67} = -\lambda \eta g^{65} - \eta g^{66}$

⋮

$m^{68} = -\lambda \eta g^{66} - \eta g^{67}$

⋮

$m^{69} = -\lambda \eta g^{67} - \eta g^{68}$

⋮

$m^{70} = -\lambda \eta g^{68} - \eta g^{69}$

⋮

$m^{71} = -\lambda \eta g^{69} - \eta g^{70}$

⋮

$m^{72} = -\lambda \eta g^{70} - \eta g^{71}$

⋮

$m^{73} = -\lambda \eta g^{71} - \eta g^{72}$

⋮

$m^{74} = -\lambda \eta g^{72} - \eta g^{73}$

⋮

$m^{75} = -\lambda \eta g^{73} - \eta g^{74}$

⋮

$m^{76} = -\lambda \eta g^{74} - \eta g^{75}$

⋮

$m^{77} = -\lambda \eta g^{75} - \eta g^{76}$

⋮

$m^{78} = -\lambda \eta g^{76} - \eta g^{77}$

⋮

$m^{79} = -\lambda \eta g^{77} - \eta g^{78}$

⋮

$m^{80} = -\lambda \eta g^{78} - \eta g^{79}$

⋮

$m^{81} = -\lambda \eta g^{79} - \eta g^{80}$

⋮

$m^{82} = -\lambda \eta g^{80} - \eta g^{81}$

⋮

$m^{83} = -\lambda \eta g^{81} - \eta g^{82}$

⋮

$m^{84} = -\lambda \eta g^{82} - \eta g^{83}$

⋮

$m^{85} = -\lambda \eta g^{83} - \eta g^{84}$

⋮

$m^{86} = -\lambda \eta g^{84} - \eta g^{85}$

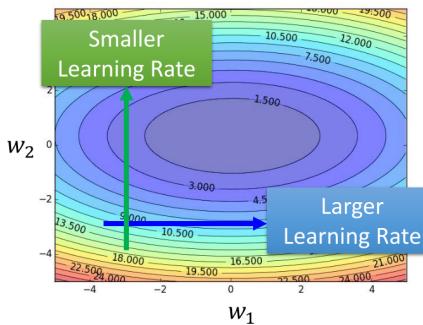
⋮

$m^{87} = -\lambda \eta g^{85} - \eta g^{86}$

⋮

## → Learning Rate 和 Gradient 做调整:

Formulation for one parameter:



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

*t*th Iteration  
 $g_i^t = \frac{\partial L}{\partial \theta_i}|_{\theta=\theta^t}$   
*i*th parameter

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

与 iteration 有关  
和 parameters 有关  
Parameter dependent

### • AdaGrad

#### Root Mean Square

$g_i^{t-1}$

$g_i^t$

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

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

Used in AdaGrad

RMSProp

$g_i^1, g_i^2, \dots, g_i^{t-1}$

$0 < \alpha < 1$

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

$\sigma_i^t = \sqrt{\alpha(\sigma_i^{t-1})^2 + (1-\alpha)(g_i^t)^2}$

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

small  $\sigma_i^t$   
larger step

increase  $\sigma_i^t$   
smaller step

decrease  $\sigma_i^t$   
larger step

( $t-1$ )

$\alpha \rightarrow 0, g_i^t$  相对于之前梯度更重要

$\alpha \rightarrow 1,$ 之前的梯度相对更重要

→ 通过调整  $\alpha$  改变更新速度

• Adam RMSProp + Momentum

应用最广泛的 Optimizer

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:  $\alpha: \text{Stepsize}$

Require:  $\beta_1, \beta_2 \in [0, 1]: \text{Exponential decay rates for the moment estimates}$

Require:  $f(\theta): \text{Stochastic objective function with parameters } \theta$

Require:  $\theta_0: \text{Initial parameter vector}$

$m_0 \leftarrow 0$  (Initialize 1<sup>st</sup> moment vector)

$v_0 \leftarrow 0$  (Initialize 2<sup>nd</sup> moment vector)

$t \leftarrow 0$  (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)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  (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  $\theta_t$  (Resulting parameters)

11

CONCLUSIONS :

• Adaptive Learning Rate

↳ AdaGrad

↳ RMSProp

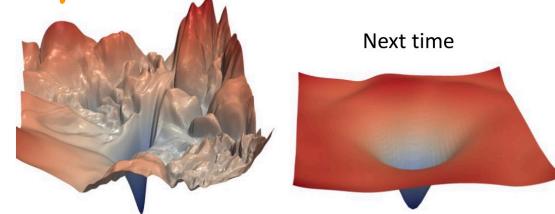
↳ Adam

• 蜜罐

→ Learning Rate Scheduling

↳ Learning Rate Decay

↳ Warm Up

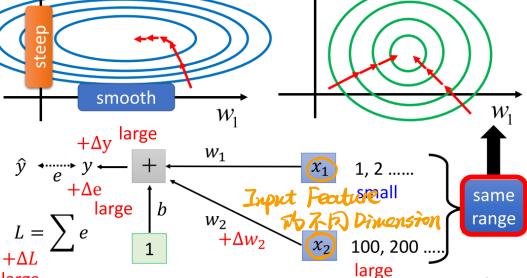


Better optimization strategies:  
If the mountain won't move,  
build a road around it.

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

从崎岖的山坡中找最优路径  
→ 将崎岖的山坡变平滑 →

Changing Landscape 椭圆→圆



## BATCH NORMALIZATION

Input Feature ( $x_1, x_2, \dots$ ) 每个 Dimension 的值差别很大  
→ 产生椭圆形的 Error Surface 不同方向斜率坡度差异较大  
→ Feature Normalization 将不同维度的 Feature 规范到相同范围  
→ 能够加速收敛

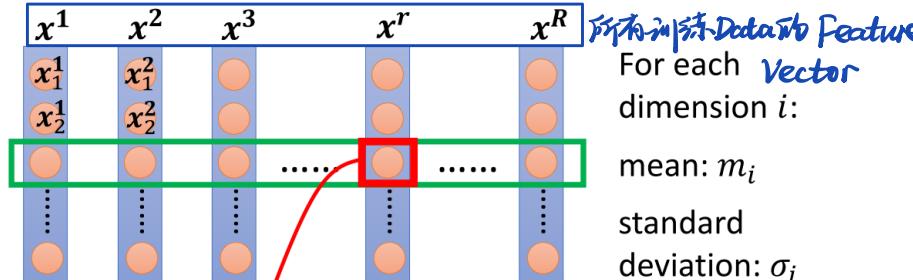
## Feature Normalization (Normalize Feature Vector)

• 取不同训练 Data 的同一 Dimension 的 Mean 和 Variance

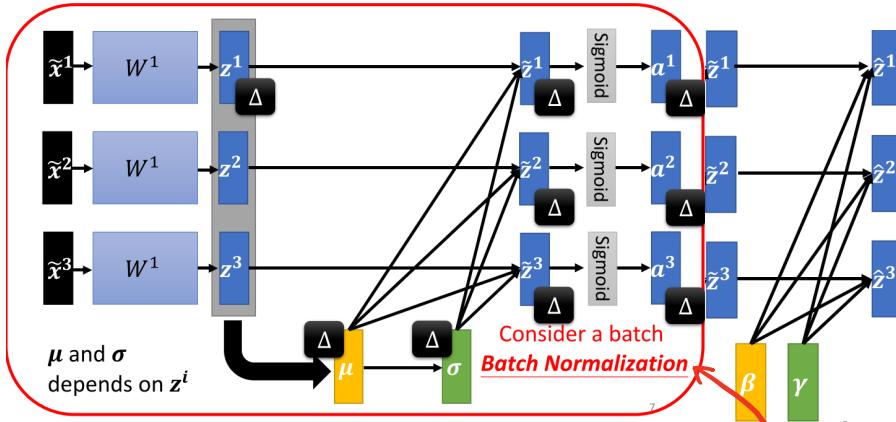
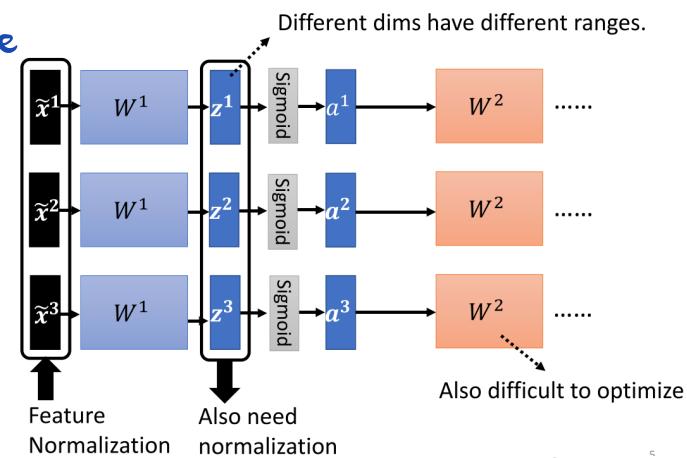
• 上标不同 → 不同训练 Data (不同样本)

下标不同 → Feature Vector 的 Dimension 不同

3



每个 Feature Vector 都不同  
Dimension 不同  
The means of all dims are 0,  
and the variances are all 1

$$\tilde{x}_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$


$$\mu = \frac{1}{3} \sum_{i=1}^3 z^i$$

$$\sigma = \sqrt{\frac{1}{3} \sum_{i=1}^3 (z^i - \mu)^2}$$

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

Element-wise

一个样本的变化会对其他样本的结果产生影响  
→ A Large Network 一个 Network 一次处理一个样本  
Large Network 同时处理所有样本

### Inference 测试阶段

Testing

$$\tilde{z} = \frac{z - \mu}{\sigma}$$

$\mu, \sigma$  are from batch?

Testing 时不依赖到先验的 Batch

We do not always have batch at testing stage.

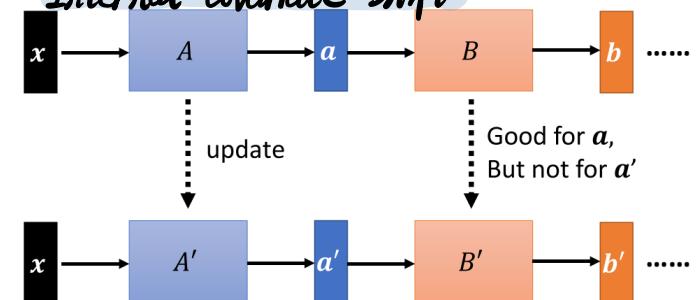
Computing the moving average of  $\mu$  and  $\sigma$  of the batches during training. 训练时计算 Moving Average ( $\bar{\mu}, \bar{\sigma}$ )

$$\mu^1 \quad \mu^2 \quad \mu^3 \quad \dots \quad \mu^t$$

$$\bar{\mu} \leftarrow p\bar{\mu} + (1-p)\mu^t$$

### Why Batch Normalization?

- Internal Covariate Shift



Batch normalization make  $a$  and  $a'$  have similar statistics.

Experimental results do not support the above idea.

How Does Batch Normalization Help Optimization?

<https://arxiv.org/abs/1805.11604>

- Smoother Error Surface