



When gradient is small ... (optimization)

### ① Training Fails because

- local minima No way to get out
- saddle point escape

Which one?

Taylor Series Approximation

$L(\theta)$  around  $\theta = \theta'$  Can be approximated below

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

Gradient  $g$  is a vector:  $g = \nabla L(\theta')$ ,  $g_i = \frac{\partial L(\theta')}{\partial \theta_i}$

Hessian  $H$  is a matrix:  $H_{ij} = \frac{\partial^2 L(\theta')}{\partial \theta_i \partial \theta_j}$

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

a.  $V^T H V > 0 \Rightarrow$  Around  $\theta'$ :  $L(\theta) > L(\theta') \Rightarrow$  local minima

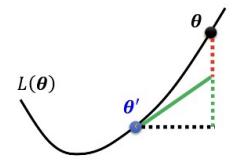
$\Rightarrow H$  is positive definite (凸的)  $\Rightarrow$  All eigen values are positive

b.  $V^T H V < 0 \Rightarrow$  Around  $\theta'$ :  $L(\theta) < L(\theta') \Rightarrow$  local maxima

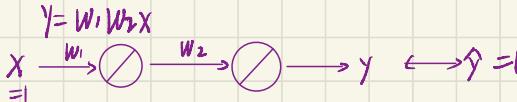
$\Rightarrow H$  is negative definite  $\Rightarrow$  All eigen values are negative

c. Sometimes  $V^T H V > 0$ , sometimes  $V^T H V < 0 \Rightarrow$  saddle point

$\Rightarrow$  some eigen values are positive. Some are negative



Example:



$$L = (y - w_1 w_2 x)^2 = (1 - w_1 w_2)^2$$

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= 2(1 - w_1 w_2)(-w_2) && \text{Critical point: } w_1 = 0, w_2 = 0 \\ &= 0 \\ \frac{\partial L}{\partial w_2} &= 2(1 - w_1 w_2)(-w_1) && H = \begin{bmatrix} 0 & -2 \\ -2 & 0 \end{bmatrix}, \lambda_1 = 2, \lambda_2 = -2 \\ &= 0 \end{aligned}$$

Saddle point

$$\frac{\partial^2 L}{\partial w_1^2} = 2(-w_2)(-w_2) = 0 \quad \frac{\partial^2 L}{\partial w_1 \partial w_2} = -2 + 4w_1 w_2 = -2$$

$$\frac{\partial^2 L}{\partial w_2 \partial w_1} = -2 + 4w_1 w_2 = -2 \quad \frac{\partial^2 L}{\partial w_2^2} = 2(-w_1)(-w_1) = 0$$



我们需要 update  $\theta$  !

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

$\therefore U$  is an eigen vector of  $H$

$\lambda$  is the eigen value of  $H$

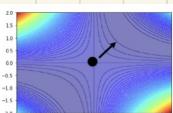
$$\Rightarrow \mathbf{U}^T \mathbf{H} \mathbf{U} = \mathbf{U}^T (\lambda \mathbf{I}) \mathbf{U} = \lambda \|\mathbf{U}\|^2$$

当  $\lambda < 0$  时  $\mathbf{L}(\theta) < \mathbf{L}(\theta')$   $\theta - \theta' = \mathbf{U}$  而  $\theta = \theta' + \mathbf{U}$  Decrease  $\mathbf{L}$

Example:

$$\lambda_2 = -2 \quad \mathbf{U} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

往左走



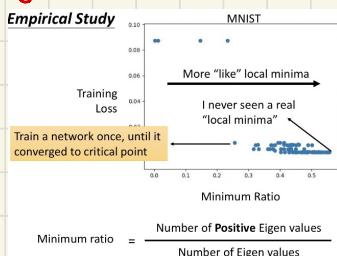
You can escape the saddle point and decrease the loss.

## Saddle Point v.s. Local Minima

在特征空间中是 local Minima, 但现实中数据维度很大, 而 Saddle point 符合 local minima

例:

## Empirical Study 実证研究



## ② Batch and Momentum

### a. Batch

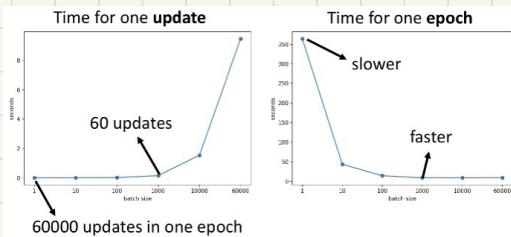
Review: Optimization with Batch

- $\theta^* = \arg \min \mathcal{L}$
  - (Randomly) Pick initial values  $\theta^0$
  - Compute gradient  $\mathbf{g}^0 = \nabla \mathcal{L}(\theta^0)$
  - update  $\theta^1 \leftarrow \theta^0 - \eta \mathbf{g}^0$
  - Compute gradient  $\mathbf{g}^1 = \nabla \mathcal{L}(\theta^1)$
  - update  $\theta^2 \leftarrow \theta^1 - \eta \mathbf{g}^1$
  - Compute gradient  $\mathbf{g}^2 = \nabla \mathcal{L}(\theta^2)$
  - update  $\theta^3 \leftarrow \theta^2 - \eta \mathbf{g}^2$
- 1 epoch = see all the batches once → Shuffle after each epoch

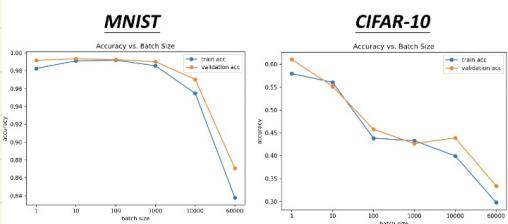
每 shuffle 之后 batch 不一样

### Small Batch vs. Large Batch

- large batch size does not require longer time to compute gradient
- smaller batch requires longer time for one epoch

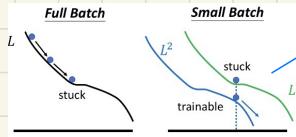


所以小的 batch 在时间上是不吃亏的.



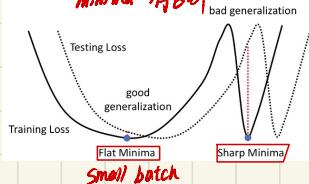
- > Smaller batch size has better performance
- > What's wrong with large batch size? Optimization Issue

- Small batch size has better performance
- "Noisy" update is better for training



Name	Training Accuracy		Testing Accuracy	
	SB	LB	SB	LB
$F_1$	$99.66\% \pm 0.05\%$	$99.92\% \pm 0.01\%$	$98.03\% \pm 0.07\%$	$97.81\% \pm 0.07\%$
$F_2$	$99.99\% \pm 0.03\%$	$98.35\% \pm 2.08\%$	$64.02\% \pm 0.2\%$	$59.45\% \pm 1.05\%$
$C_1$	$99.89\% \pm 0.02\%$	$99.66\% \pm 0.2\%$	$80.04\% \pm 0.12\%$	$77.26\% \pm 0.42\%$
$C_2$	$99.99\% \pm 0.04\%$	$99.99\% \pm 0.01\%$	$89.24\% \pm 0.12\%$	$87.26\% \pm 0.07\%$
$C_3$	$99.56\% \pm 0.44\%$	$99.88\% \pm 0.30\%$	$49.58\% \pm 0.39\%$	$46.45\% \pm 0.43\%$
$C_4$	$99.10\% \pm 1.23\%$	$99.57\% \pm 1.84\%$	$63.08\% \pm 0.5\%$	$57.81\% \pm 0.17\%$

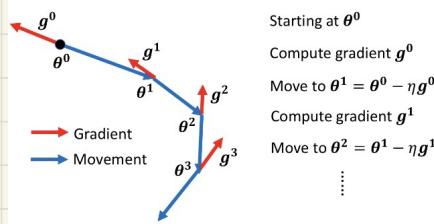
"Noisy" update is better for generalization  
minima 好坏 bad generalization



	Small	Large
Speed for one update (no parallel)	Faster	Slower
Speed for one update (with parallel)	Same	Same (not too large)
Time for one epoch	Slower	Faster
Gradient	Noisy	Stable
Optimization	Better	Worse
Generalization	Better	Worse

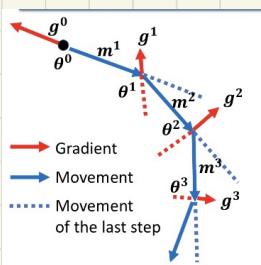
## (3) momentum

Recall (Vanilla) Gradient Descent



## Gradient Descent + Momentum

Movement: movement of last step minus gradient at present



① Starting at  $\theta^0$

$$\text{Movement } m^0 = 0$$

Compute gradient  $g^0$

$$\text{Movement } m^1 = \lambda m^0 - \eta g^0$$

$$\text{Move to } \theta^1 = \theta^0 + m^1$$

② Compute gradients  $g^1$

$$\text{movement } m^2 = \lambda m^1 - \eta g^1$$

$$\text{move to } \theta^2 = \theta^1 + m^2$$

$$\begin{aligned} m^0 &= 0 \\ m^1 &= -\eta g^0 \\ m^2 &= -\lambda \eta g^0 - \eta g^1 \\ &\vdots \end{aligned}$$

Movement not just based on gradient, but previous movement.

周記

