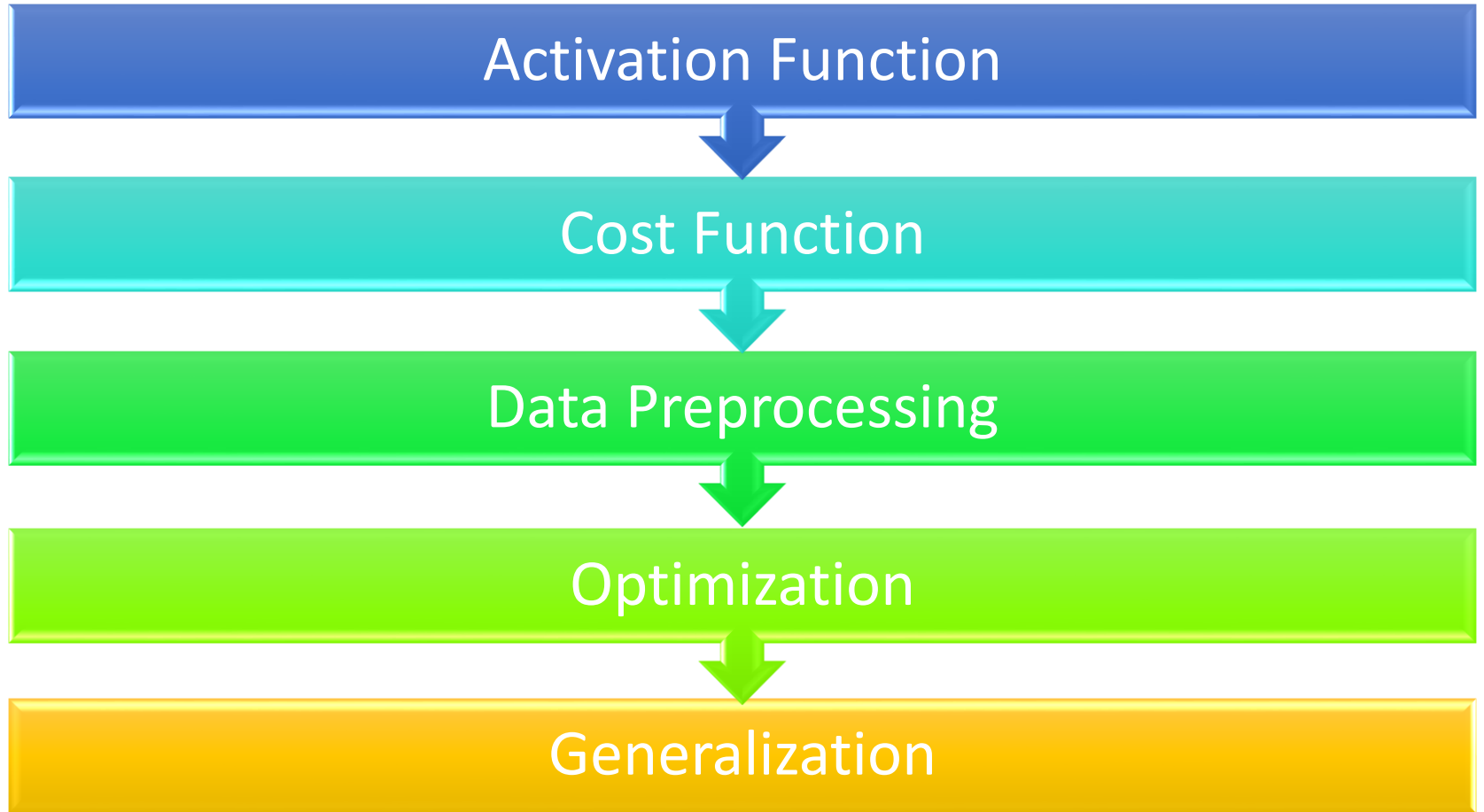# Tips for Training Deep Neural Network
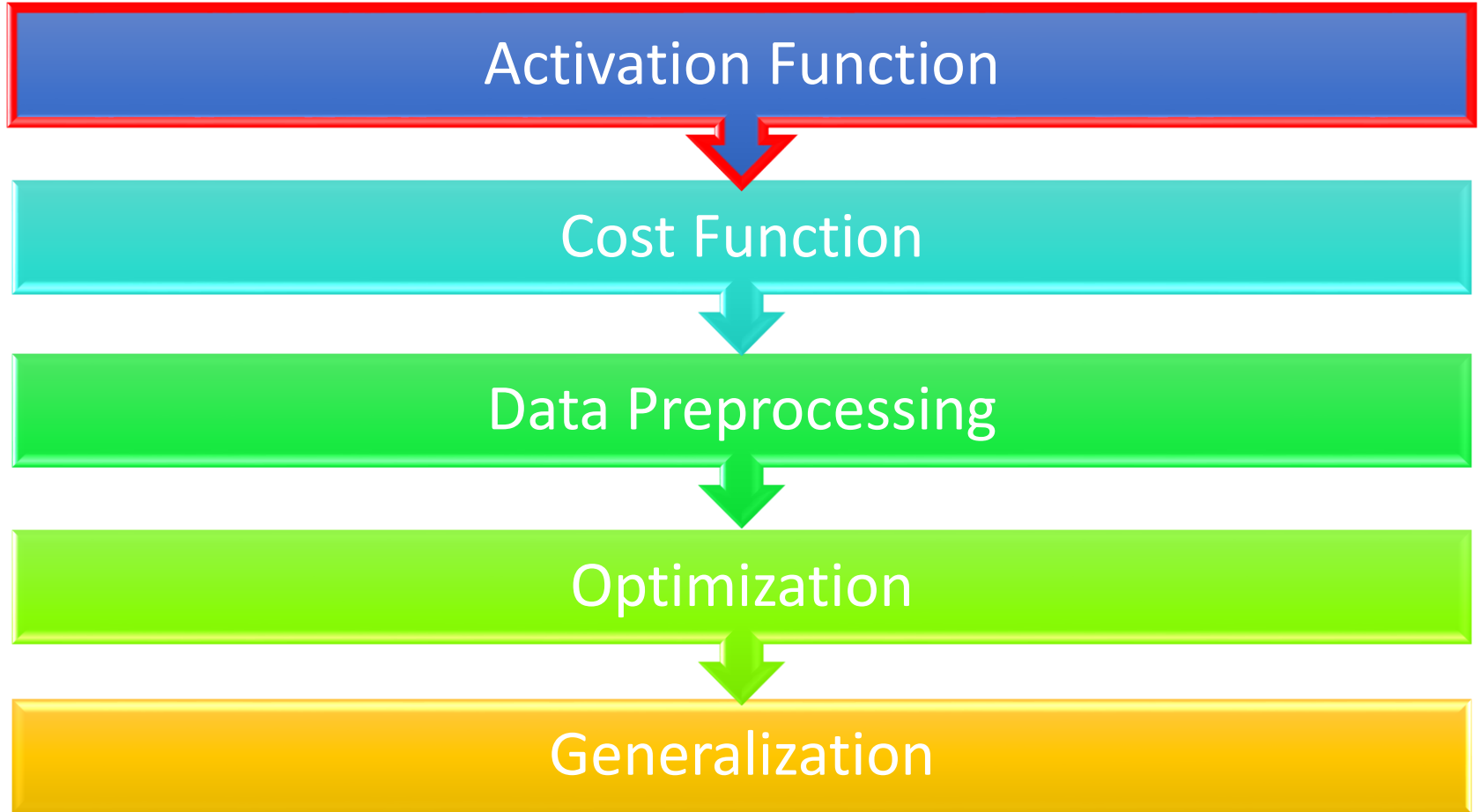
## Hung-yi Lee

# Announcement

- 分組
  - 請確認在 ceiba 上的分組是否正確
- HW1
  - 截止日期: 10/23 2:00 p.m. (下週五上課前)
- HW2
  - 公告日期: 10/23
  - 截止日期: 11/13 2:00 p.m.
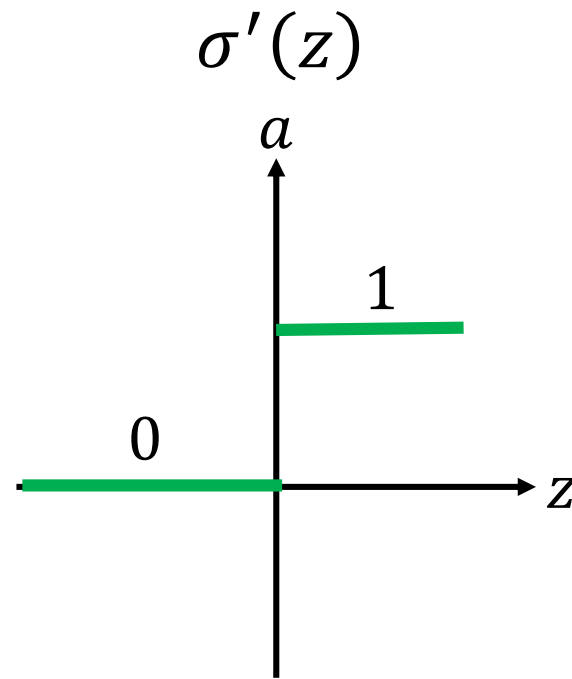    - 比第一堂課公告的提早一週截止

# Outline

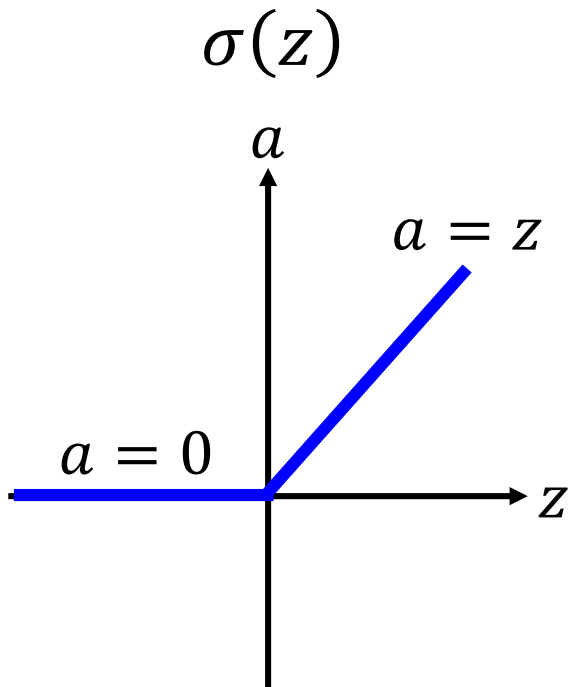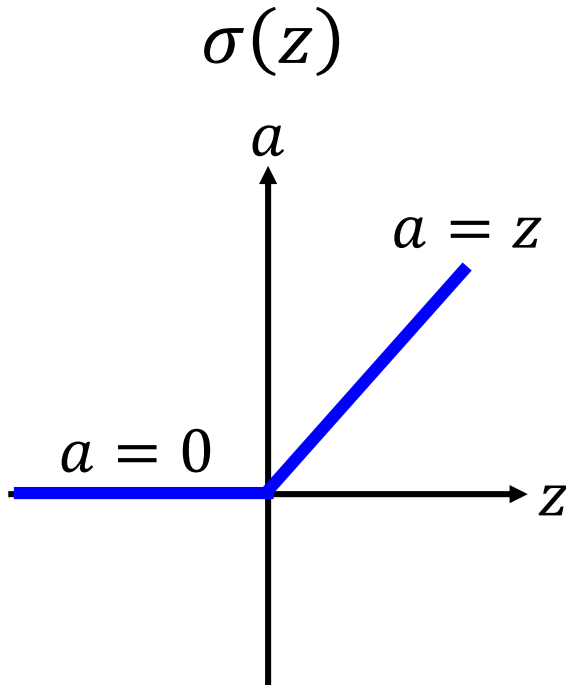Activation Function

Cost Function

Data Preprocessing

Optimization

Generalization

# Outline

Activation Function

Cost Function

Data Preprocessing

Optimization

Generalization

# ReLU

- Rectified Linear Unit (ReLU)

$$\sigma(z)$$

$$\sigma'(z)$$

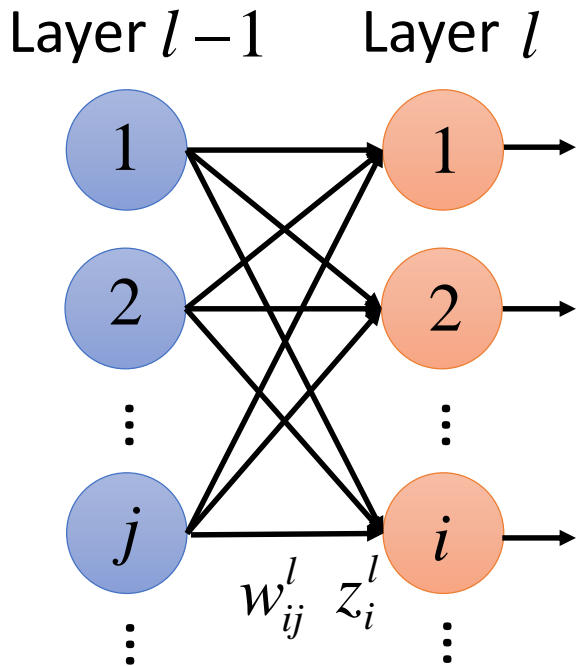# ReLU

- Rectified Linear Unit (ReLU)

$\sigma(z)$



**Reason:**

1. Fast to compute

2. Biological reason

3. Infinite sigmoid with different biases

4. Vanishing gradient problem

# Review: Backpropagation

$$\frac{\partial C_x}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \frac{\partial C_x}{\partial z_i^l}$$

Error signal

$$\delta_i^l$$

$$\begin{cases} a_j^{l-1} & l > 1 \\ x_j & l = 1 \end{cases}$$

Layer $l-1$    Layer $l$



$w_{ij}^l$   $z_i^l$

### **Forward Pass**

$$z^1 = W^1 x + b^1$$

$$a^1 = \sigma(z^1)$$

......

$$z^{l-1} = W^{l-1} a^{l-2} + b^{l-1}$$

$$a^{l-1} = \sigma(z^{l-1})$$

### **Backward Pass**

$$\delta^L = \sigma'(z^L) \bullet \nabla C_x(y)$$

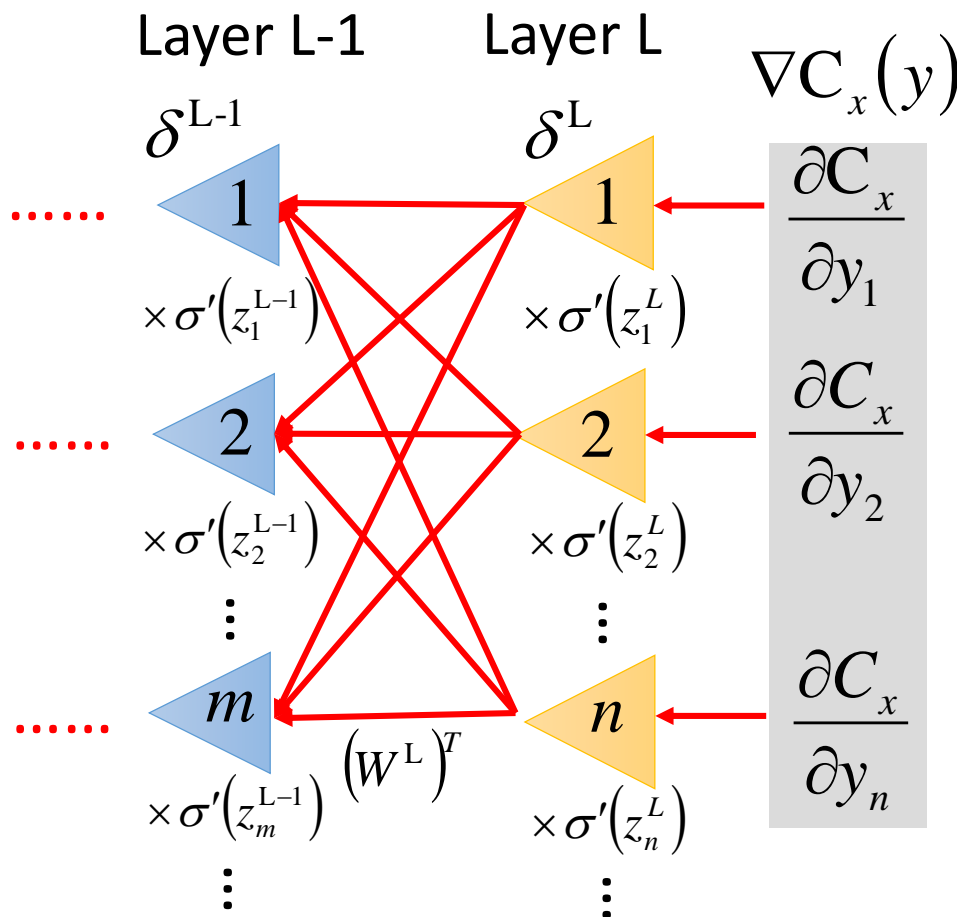$$\delta^{L-1} = \sigma'(z^{L-1}) \bullet (W^L)^T \delta^L$$

......

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^T \delta^{l+1}$$

......

# Review: Backpropagation

$$\frac{\partial C_x}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \frac{\partial C_x}{\partial z_i^l}$$

Error signal

$\delta_i^l$



Layer L-1    Layer L    $\nabla C_x(y)$

$\delta^{L-1}$    $\delta^{L}$

1    1    $\dfrac{\partial C_x}{\partial y_1}$

$\times \sigma'(z_1^{L-1})$    $\times \sigma'(z_1^{L})$

2    2    $\dfrac{\partial C_x}{\partial y_2}$

$\times \sigma'(z_2^{L-1})$    $\times \sigma'(z_2^{L})$

$m$    $n$    $\dfrac{\partial C_x}{\partial y_n}$

$(W^L)^T$

$\times \sigma'(z_m^{L-1})$    $\times \sigma'(z_n^{L})$

## Backward Pass

$$\delta^L = \sigma'(z^L) \bullet \nabla C_x(y)$$

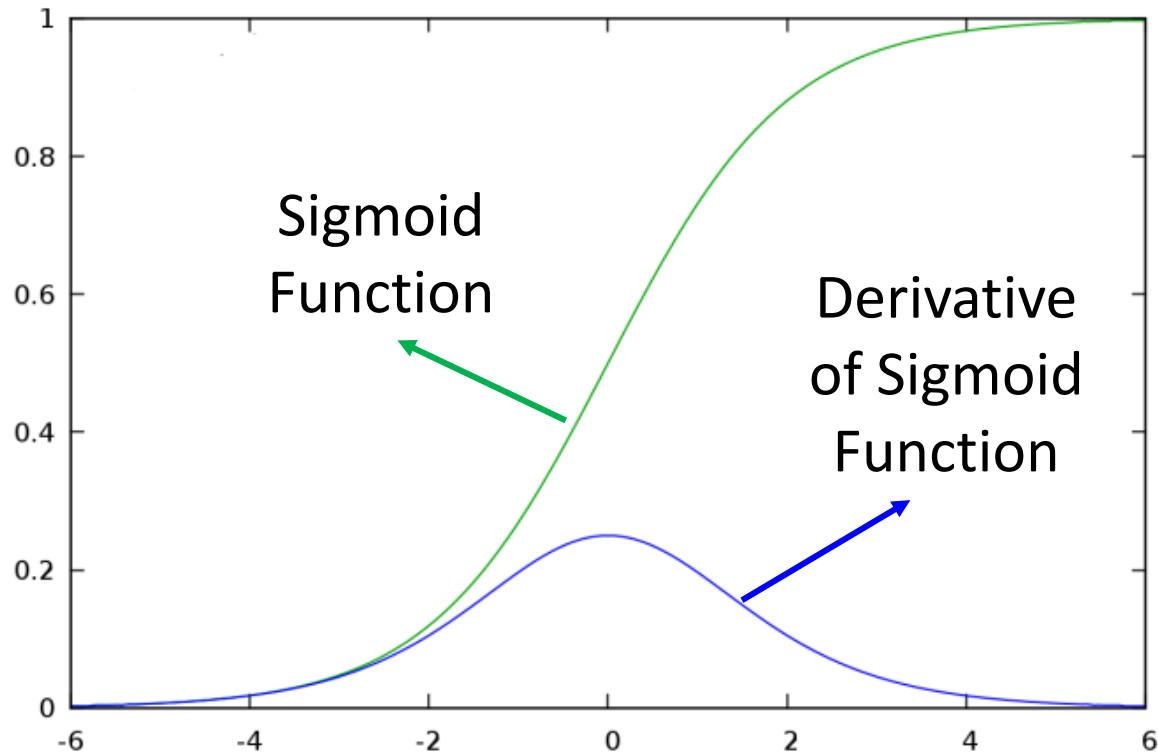$$\delta^{L-1} = \sigma'(z^{L-1}) \bullet (W^L)^T \delta^L$$

......

$$\delta^l = \sigma'(z^l) \bullet (W^{l+1})^T \delta^{l+1}$$
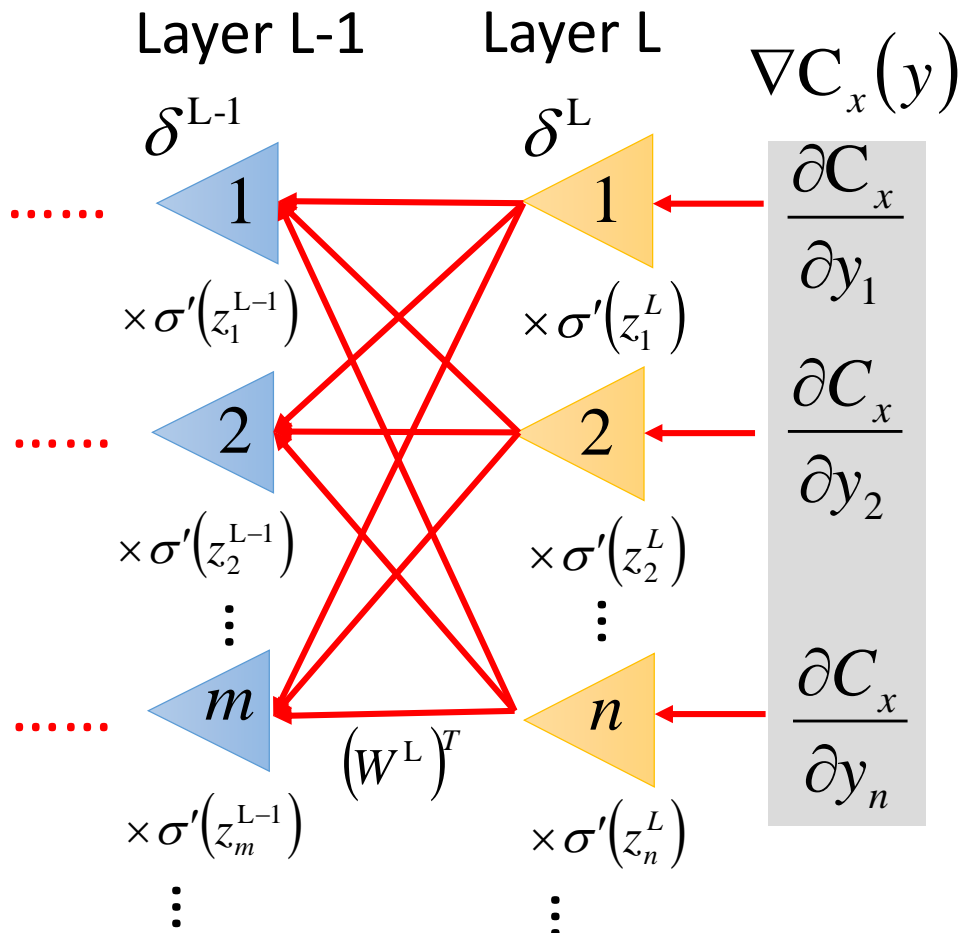
......

# Problem of Sigmoid



Derivative of Sigmoid Function is always smaller than 1

# Vanishing Gradient Problem
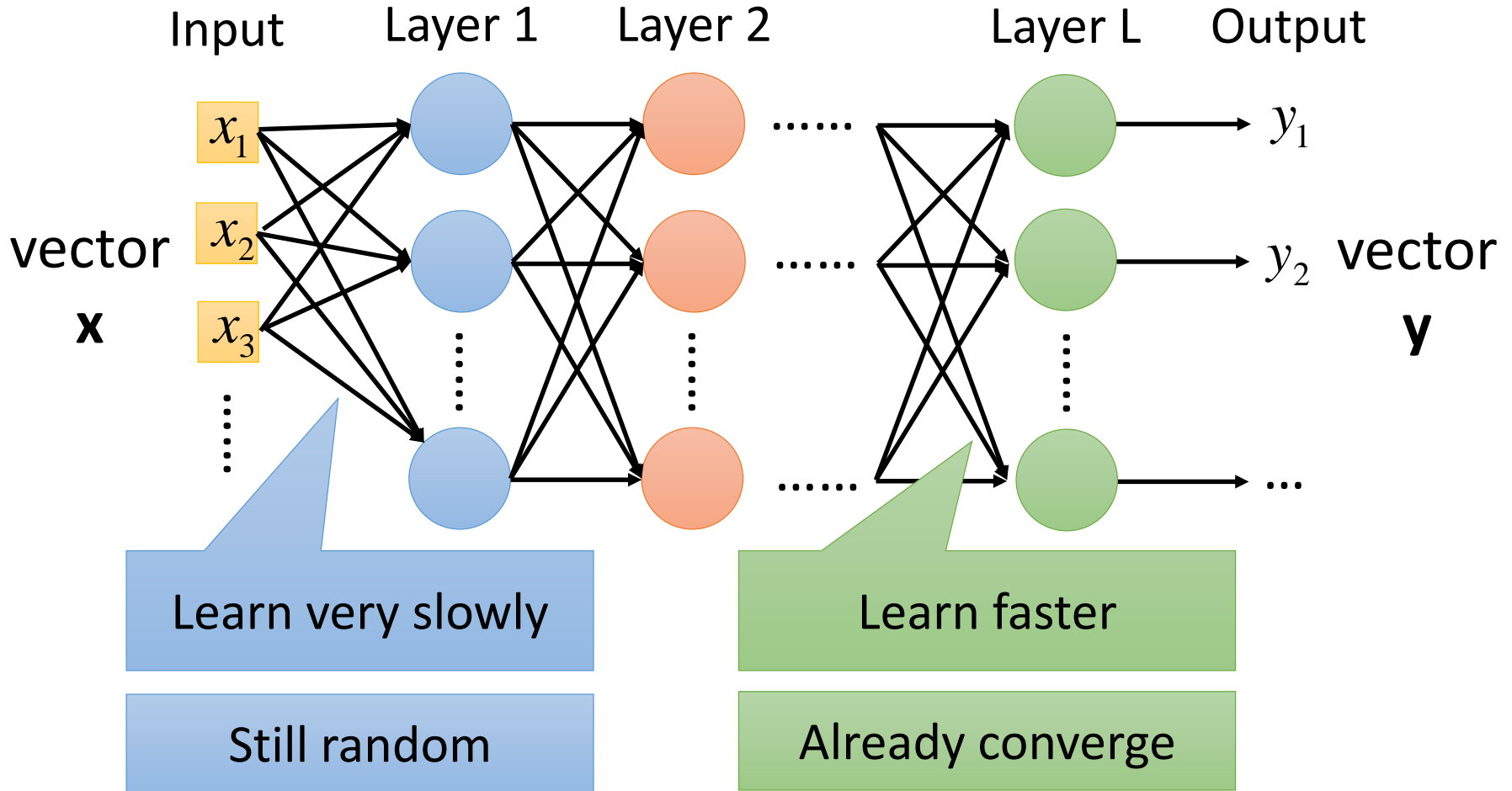
Backward Pass:



- For sigmoid function, $\sigma'(z)$ always smaller than 1

- Error signal is getting smaller and smaller

Gradient is smaller

$$\frac{\partial C_x}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \boxed{\frac{\partial C_x}{\partial z_i^l}} \rightarrow \boxed{\delta_i^l}$$
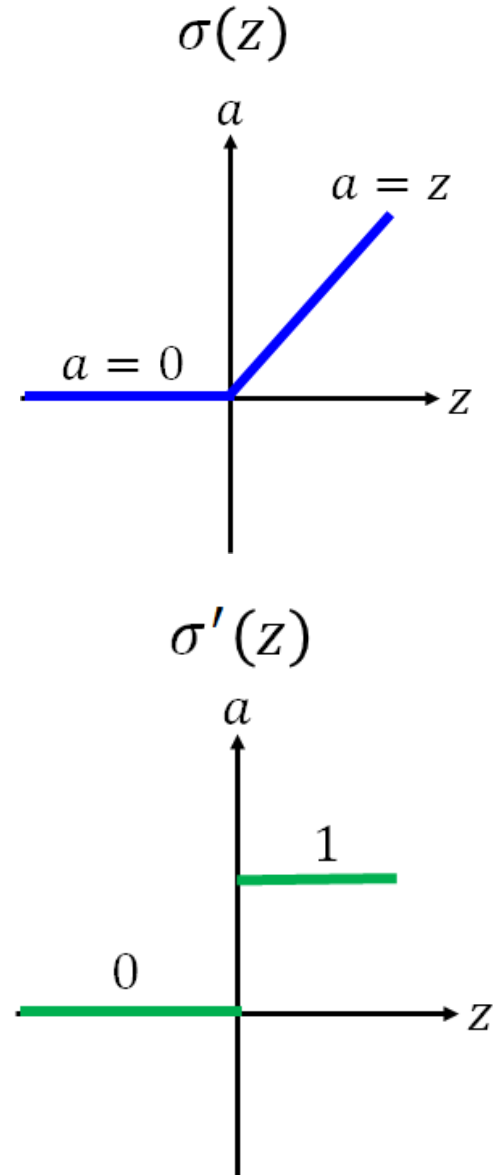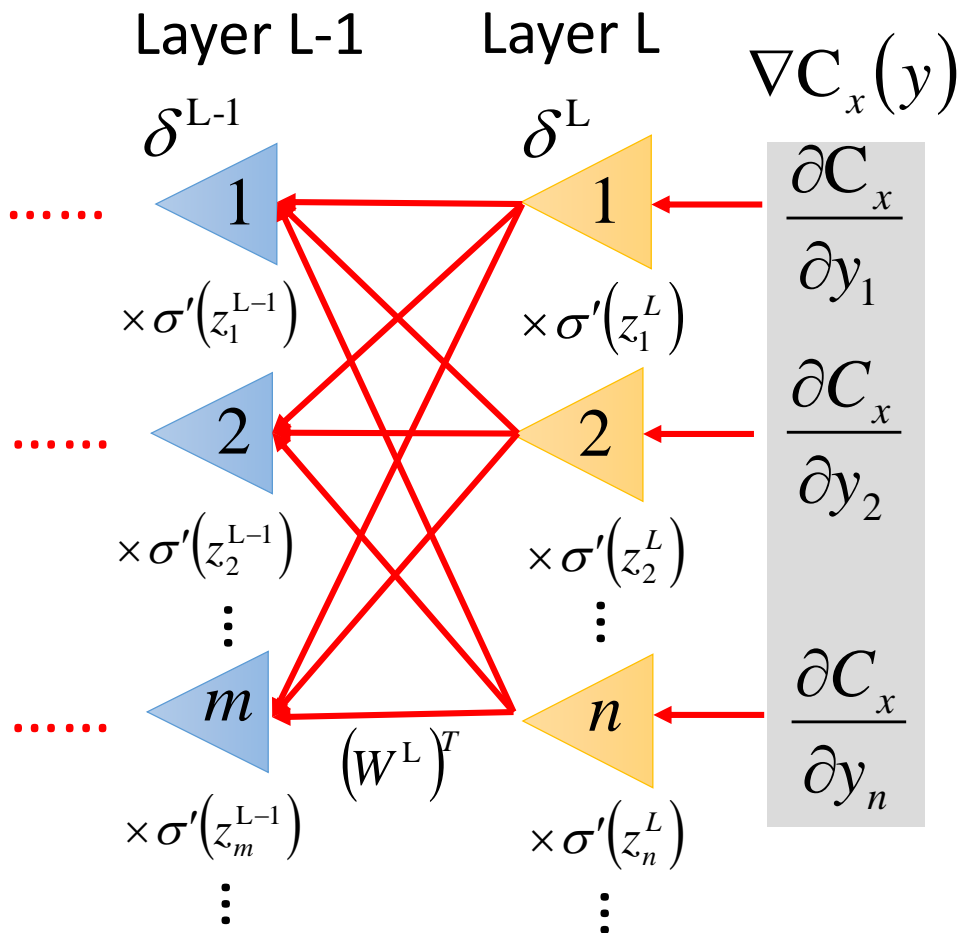
# Vanishing Gradient Problem

# ReLU

## Backward Pass:

Layer L-1    Layer L    $\nabla C_x(y)$

$\delta^{L-1}$    $\delta^{L}$

1    1    $\dfrac{\partial C_x}{\partial y_1}$

$\times \sigma'(z_1^{L-1})$    $\times \sigma'(z_1^{L})$

2    2    $\dfrac{\partial C_x}{\partial y_2}$

$\times \sigma'(z_2^{L-1})$    $\times \sigma'(z_2^{L})$

$m$    $n$    $\dfrac{\partial C_x}{\partial y_n}$

$(W^L)^T$

$\times \sigma'(z_m^{L-1})$    $\times \sigma'(z_n^{L})$
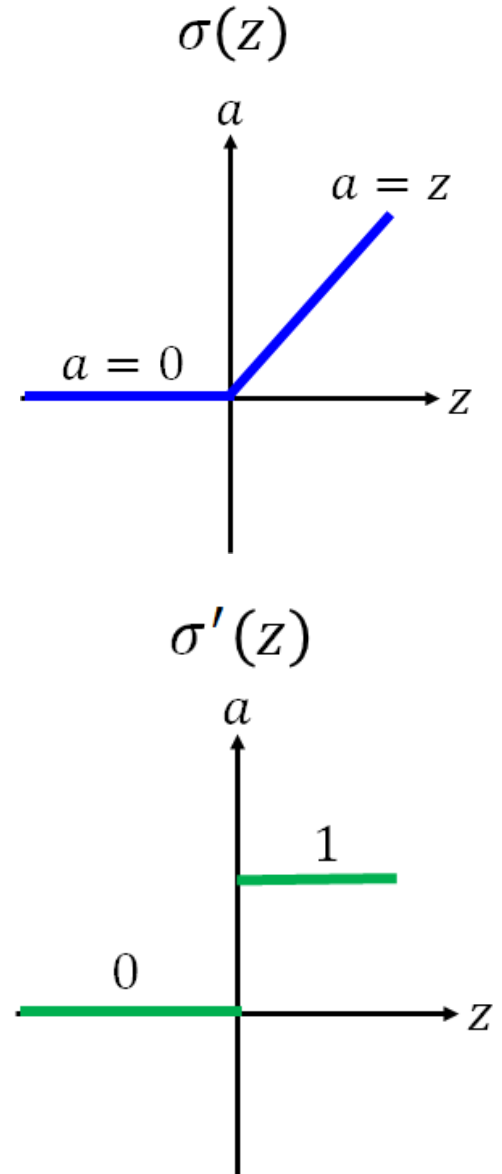
$\sigma(z)$

$a$

$a = z$

$a = 0$

$z$

$\sigma'(z)$

$a$

1

0

$z$

# ReLU

## Backward Pass:

# ReLU

## Backward Pass:



Layer L-1     Layer L     $\nabla C_x(y)$

$\delta^{L-1}$     $\delta^{L}$

$\frac{\partial C_x}{\partial y_1}$

$\frac{\partial C_x}{\partial y_2}$

A thinner network without any attenuation

$\sigma(z)$

$a = z$

$a = 0$

$\sigma'(z)$
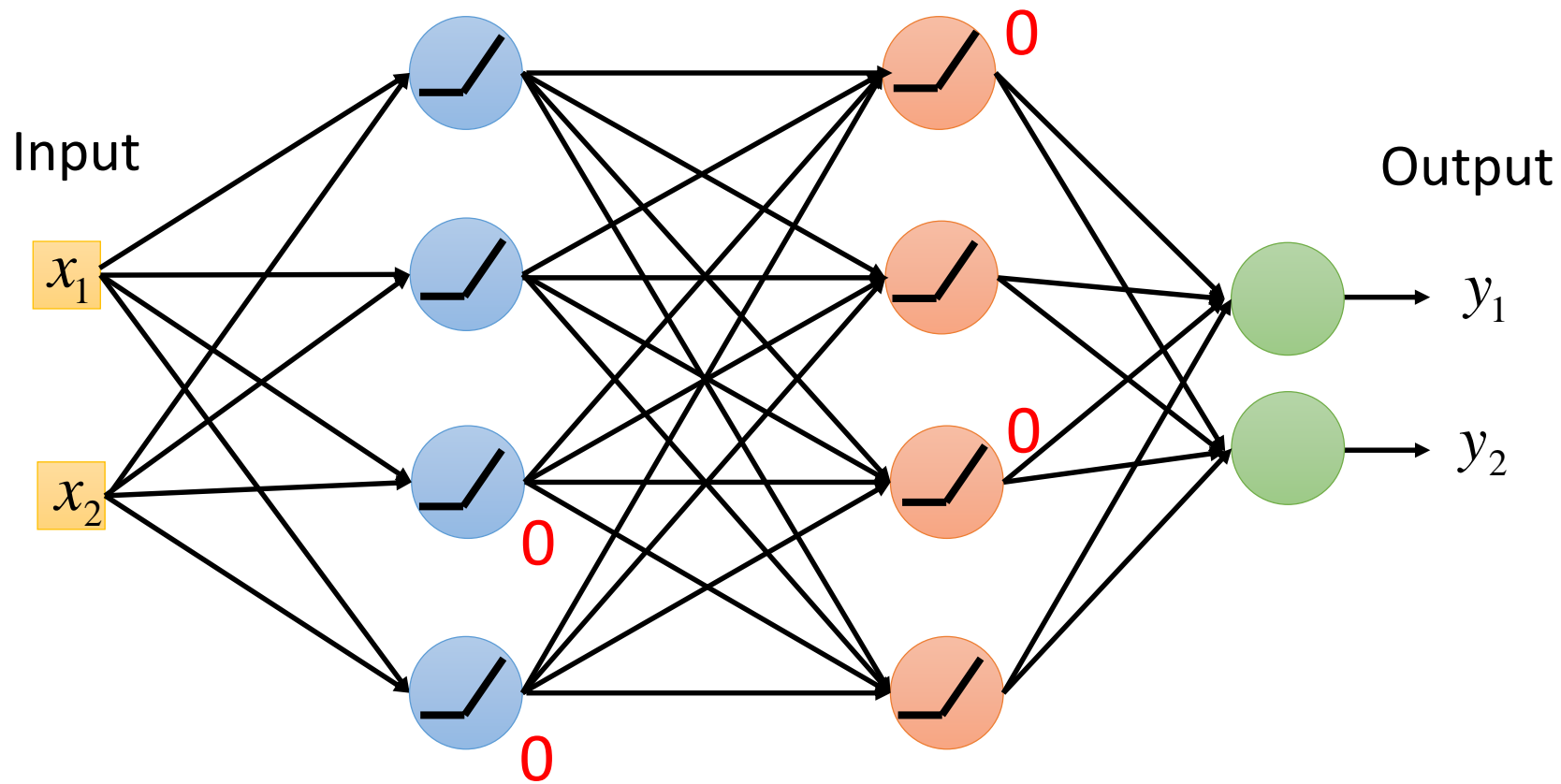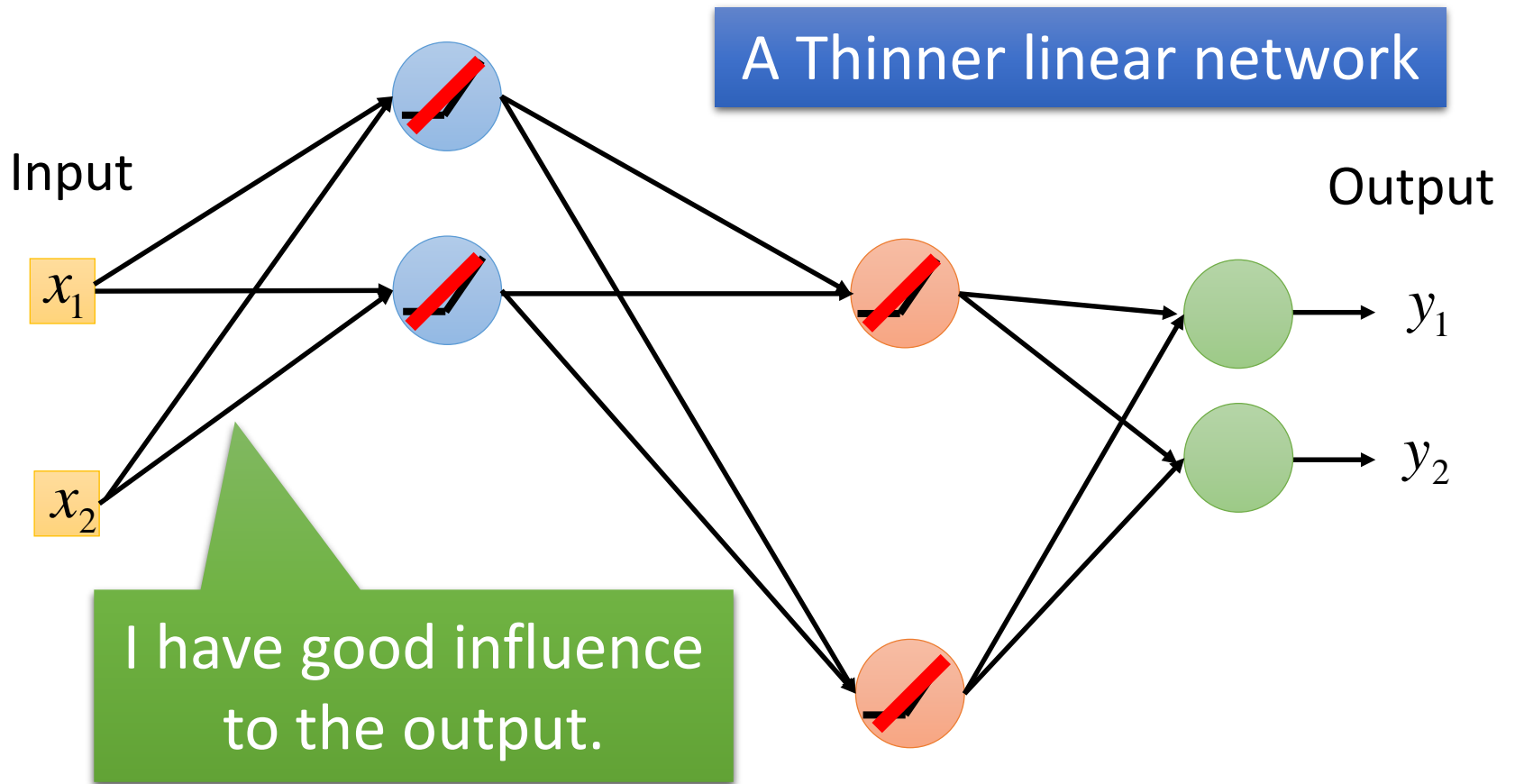
# ReLU

# ReLU

# ReLU

$$\frac{\partial C_x}{\partial w_{nj}^L} = \frac{\partial z_n^L}{\partial w_{nj}^L} \frac{\partial C_x}{\partial z_n^L}$$
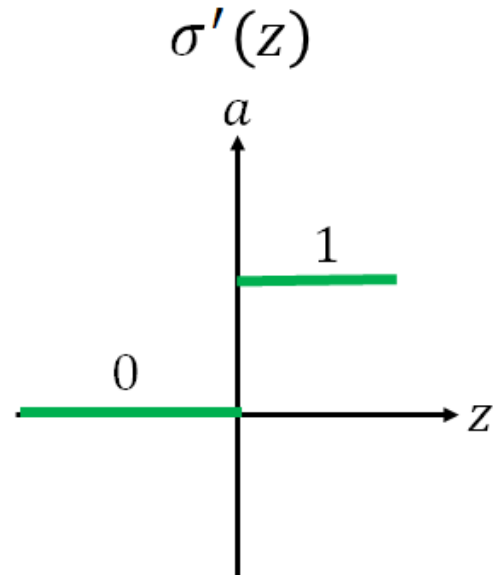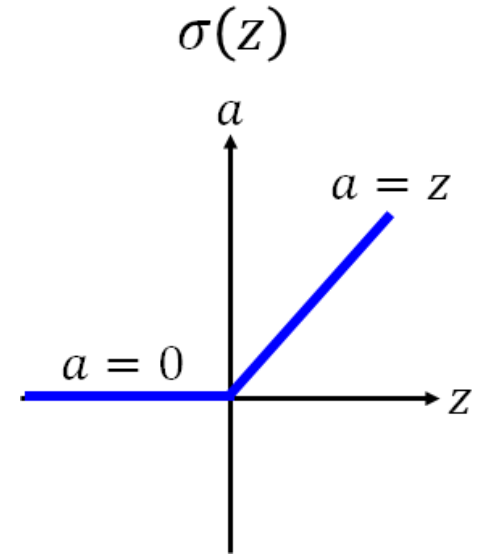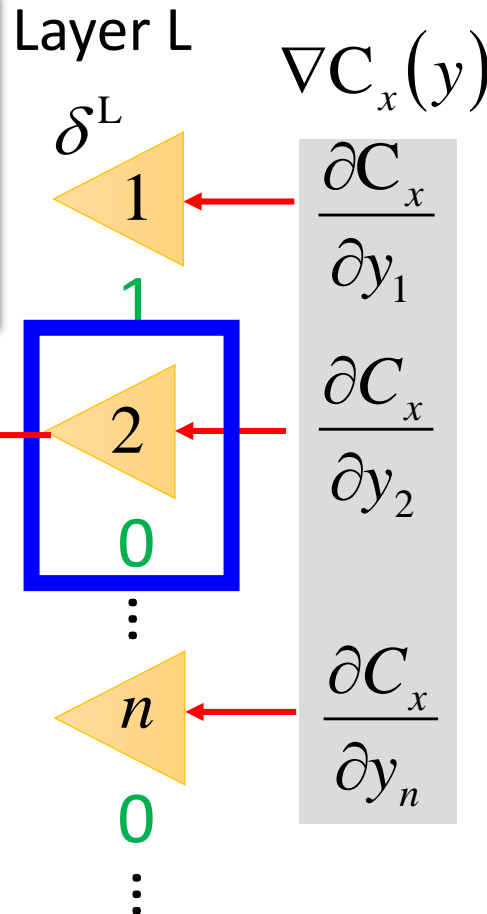
## Backward Pass:

All the weights connected to this neuron will not update.

Layer L

$\delta^L$

$\nabla C_x(y)$

$$\delta_n^L = \frac{\partial C_x}{\partial z_n^L} = 0$$

**_Possible solution:_**

1. softplus

2. Initialize with large bias



$\sigma(z)$

$\sigma'(z)$

# ReLU - variant

*Leaky ReLU*



$a$

$a = z$

$z$

$a = 0.01z$

*Parametric ReLU*



$a$

$a = z$

$z$

$a = \alpha z$

α also learned by gradient descent

# Maxout

- All ReLU variants are just special cases of Maxout



$$z^1 = W^1 x$$

$$z^2 = W^2 a^1$$

# Maxout – ReLU is special case



Input $x$, $1$ → $w$, $z$ → ReLU → $a$

Input $x$, $1$ → $w$, $b$, $0$, $0$ → $+$ $z_1$, $+$ $z_2$ → Max → $a$

$max\{z_1, z_2\}$

$z = wx + b$

$z_1 = wx + b$

$z_2 = 0$

# Maxout – ReLU is special case
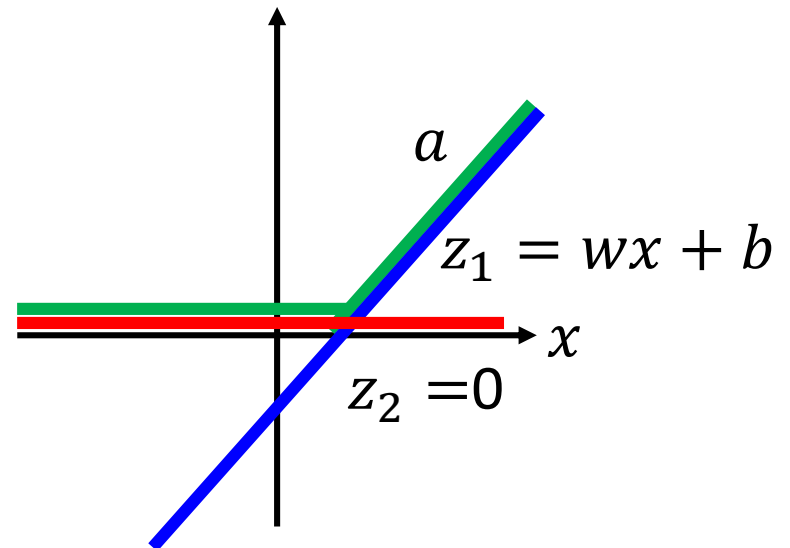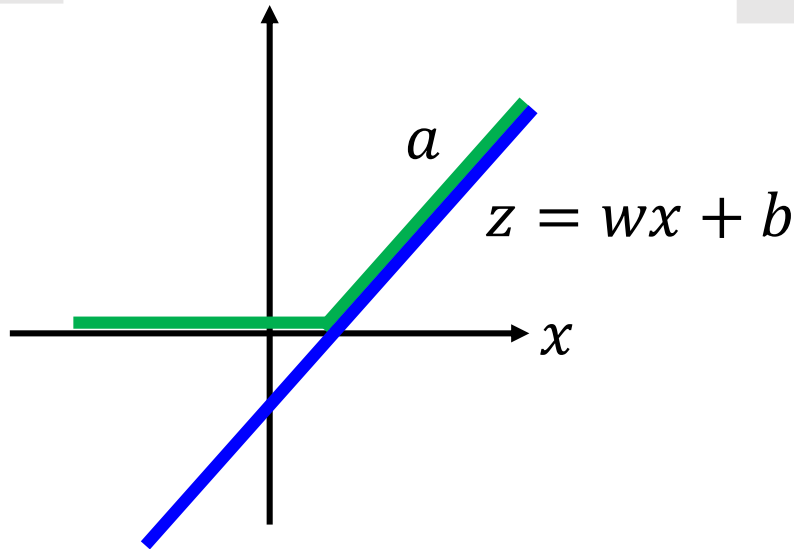


Input

$z$

$w$ ReLU $\rightarrow a$

$x$

$b$

$1$

$z = wx + b$

$a$

$x$

Input

$w$

$b$ $+ \rightarrow z_1$

$x$

$w'$ $+ \rightarrow z_2$

$b'$

$1$

Max $\rightarrow a$

$max\{z_1, z_2\}$

Learnable Activation Function

$a$

$z_1 = wx + b$

$x$

$z_2 = w'x + b'$

# Maxout - Training

- Given a training data x, we know which z would be the max

# Maxout - Training

- Given a training data x, we know which z would be the max



- Train this thin and linear netowrk

# Outline

# Cost Function



Layer L-1

Layer L
(Output layer)

$$C = \frac{1}{2}\|y - \hat{y}\|^2$$

$$= \frac{1}{2}\sum_n (y_n - \hat{y}_n)^2$$
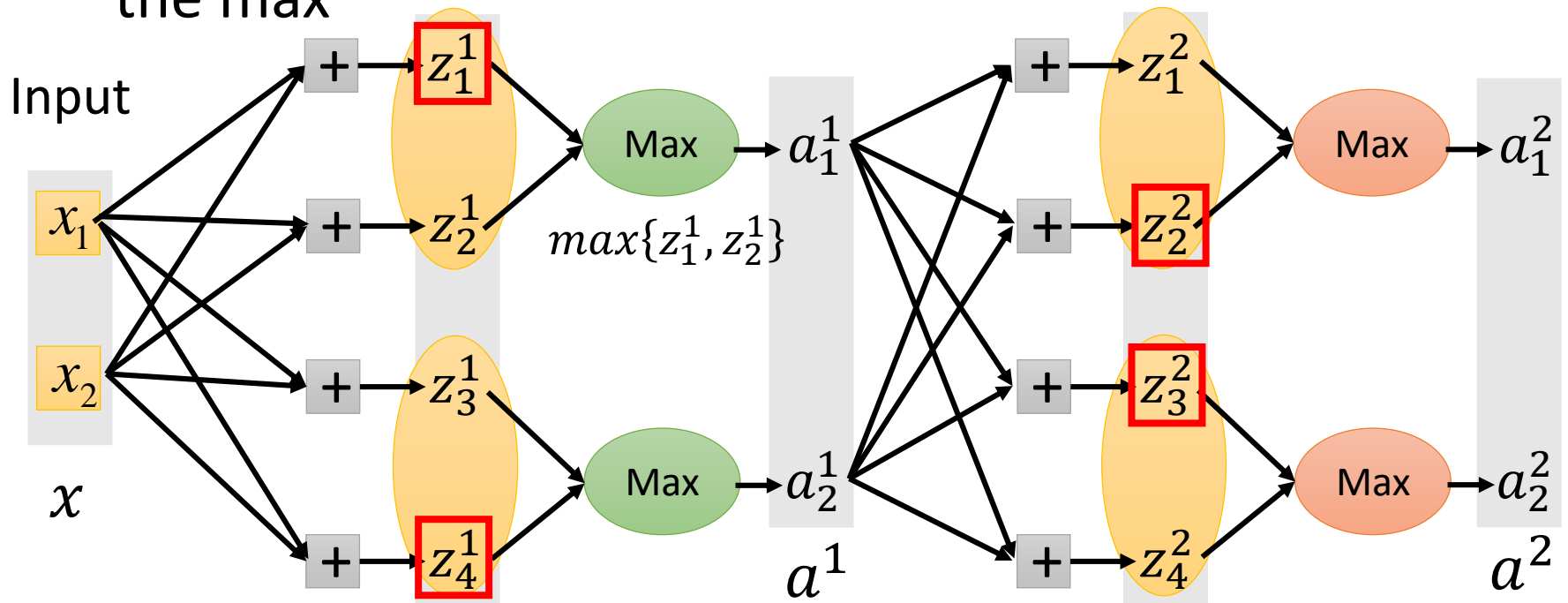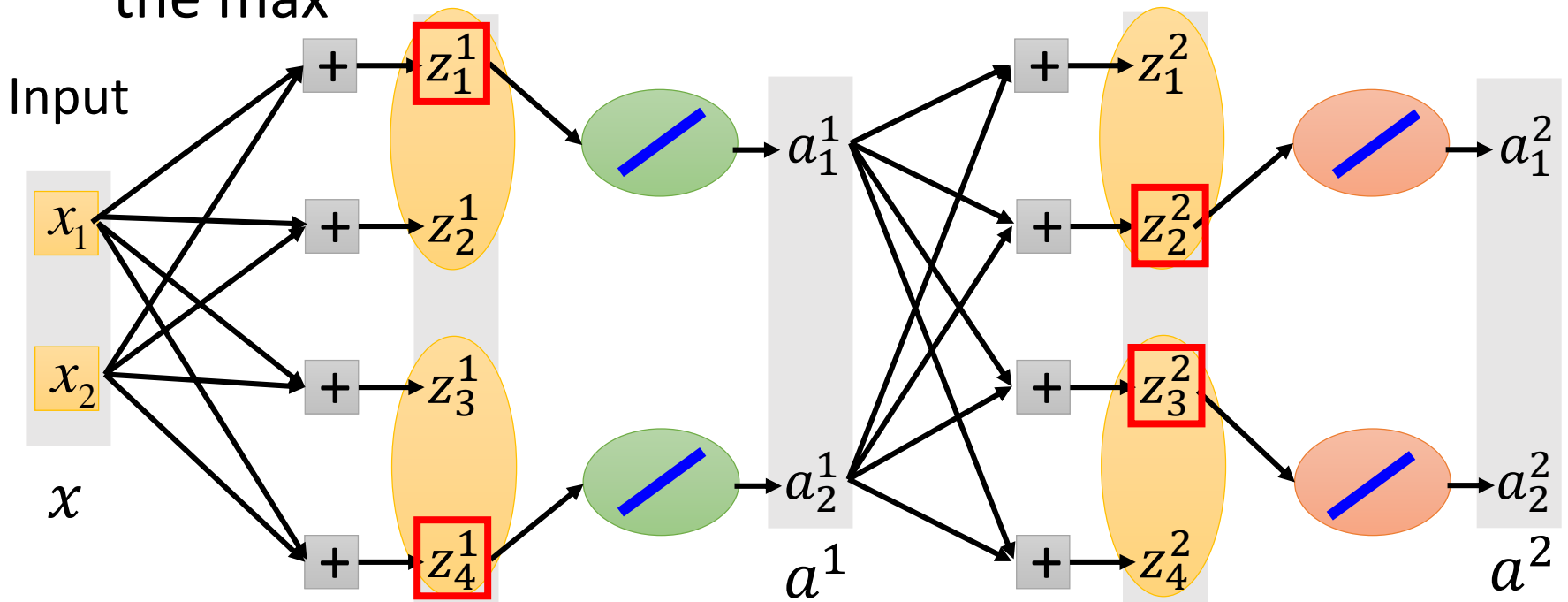
# Output Layer

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

ReLU $\quad y = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1.2 \\ \vdots \end{bmatrix}$

More similar?

ReLU $\quad y = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 2 \\ \vdots \end{bmatrix}$

➤ Larger output means larger confidence

Better?

**_Classification Task:_**

Only one dimension is 1, and others are all 0

It is better to let the output bounded.

# Softmax

- Softmax layer as the output layer

***Ordinary Output layer***

$z_1^L \longrightarrow \sigma \longrightarrow y_1 = \sigma\left(z_1^L\right)$

$z_2^L \longrightarrow \sigma \longrightarrow y_2 = \sigma\left(z_2^L\right)$

$z_3^L \longrightarrow \sigma \longrightarrow y_3 = \sigma\left(z_3^L\right)$

# Softmax

- Softmax layer as the output layer

**_Probability_**:
- $1 > y_i > 0$
- $\sum_i y_i = 1$

**_Softmax Layer_**



$z_1^L$ **3**  $\rightarrow$  $e$  $\rightarrow$  $e^{z_1^L}$ **20**  $\rightarrow$  $\div$  $\rightarrow$ **0.88** $y_1 = e^{z_1^L} \Big/ \sum_{j=1}^{3} e^{z_j^L}$

$z_2^L$ **1**  $\rightarrow$  $e$  $\rightarrow$  $e^{z_2^L}$ **2.7**  $\rightarrow$  $\div$  $\rightarrow$ **0.12** $y_2 = e^{z_2^L} \Big/ \sum_{j=1}^{3} e^{z_j^L}$

$z_3^L$ **-3**  $\rightarrow$  $e$  $\rightarrow$  $e^{z_2^L}$ **0.05**  $\rightarrow$  $\div$  $\rightarrow$ **≈0** $y_3 = e^{z_3^L} \Big/ \sum_{j=1}^{3} e^{z_j^L}$

$+$  $\sum_{j=1}^{3} e^{z_j^L}$

# Softmax

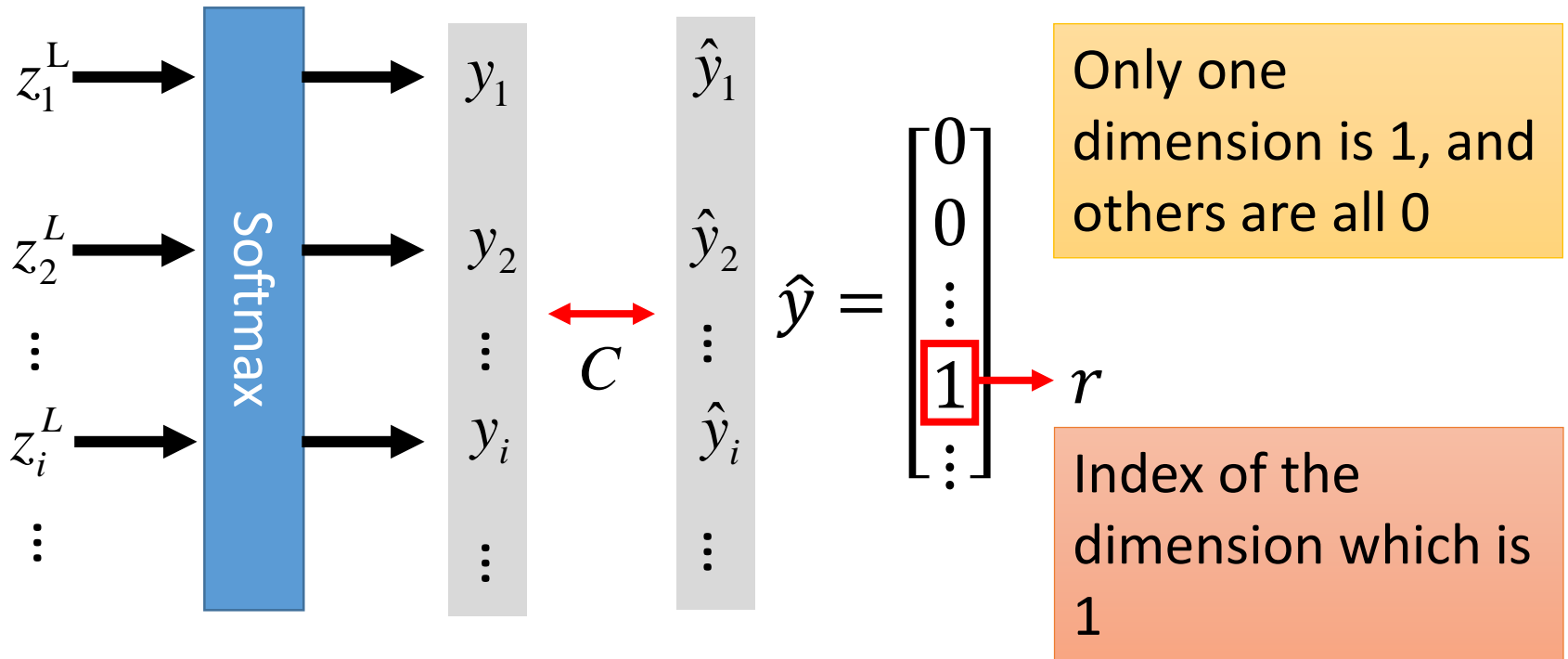- What kind of cost function should we used for softmax layer output?

# Softmax

Define cost: $C = -\log y_r$

$$y_i = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

**Do we have to consider other dimensions?**



$z_1^L \rightarrow$ Softmax $\rightarrow y_1$ ⟷ $\hat{y}_1$

$z_2^L \rightarrow y_2$ ⟷ $\hat{y}_2$

$\vdots$

$z_i^L \rightarrow y_i$ ⟷ $\hat{y}_i$

$C$

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix} \rightarrow r$$

Only one dimension is 1, and others are all 0
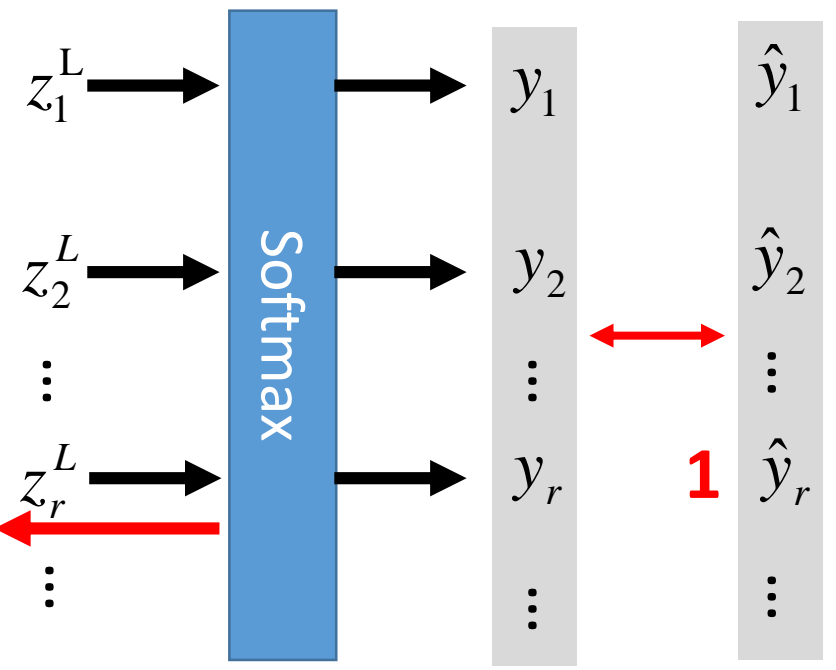
Index of the dimension which is 1

$$y_i = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

$$C = -\log y_r$$



$$\delta_r^{\mathrm{L}} = \frac{\partial C}{\partial z_r^L}$$

$$y_r - 1$$

$$\frac{\partial C}{\partial y_r} \frac{\partial y_r}{\partial z_r^L}$$

$$\delta_r^{\mathrm{L}} = \frac{\partial C}{\partial z_r^L} = -\frac{1}{y_r} \frac{\partial y_r}{\partial z_r^L} = -\frac{1}{y_r}\left(y_r - y_r^2\right) = y_r - 1$$

$$y_r = \frac{e^{z_r^L}}{\sum_j e^{z_j^L}}$$

$z_r^L$ appears in both numerator and denominator

The absolute value of $\delta_r^L$ is larger when $y_r$ is far from 1

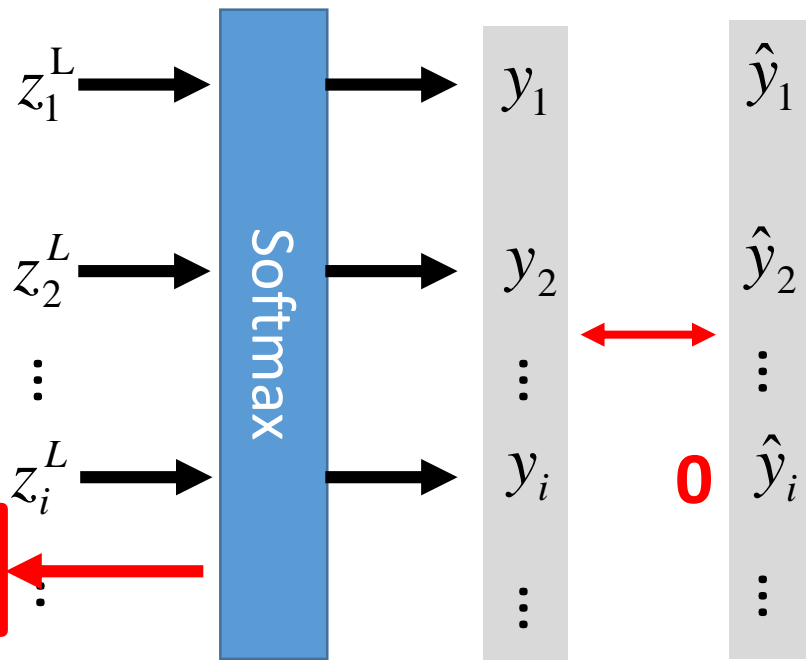$$y_i = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

$$C = -\log y_r$$

$$i \neq r$$

$$\delta_i^L = \frac{\partial C}{\partial z_i^L} = -\frac{1}{y_r}\frac{\partial y_r}{\partial z_i^L} = -\frac{1}{y_r}\left(-y_r y_i\right) = \underline{y_i}$$

$$y_r = \frac{e^{z_r^L}}{\sum_j e^{z_j^L}}$$

$$\frac{\partial C}{\partial y_r}\frac{\partial y_r}{\partial z_i^L}$$

$$\delta_i^L = \frac{\partial C}{\partial z_i^L}$$

$y_i - 0$

Softmax

$z_1^L$ → → $y_1$    $\hat{y}_1$

$z_2^L$ → → $y_2$    $\hat{y}_2$

$z_i^L$ → → $y_i$    **0** $\hat{y}_i$

$z_i^L$ appears only in denominator

The absolute value of $\delta_i^L$ is larger when $y_i$ is larger

# Outline

Activation Function

Cost Function

Data Preprocessing
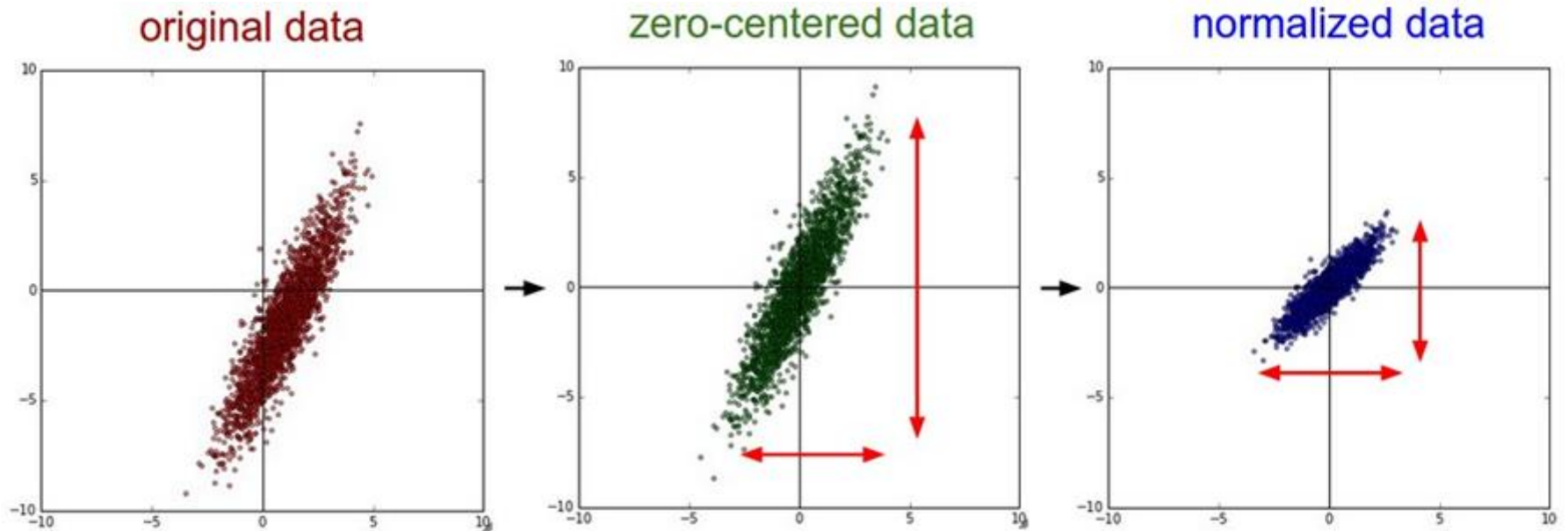
Optimization

Generalization

# Normalizing Input



For each dimension i:

mean: $m_i$

standard deviation: $\sigma_i$

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

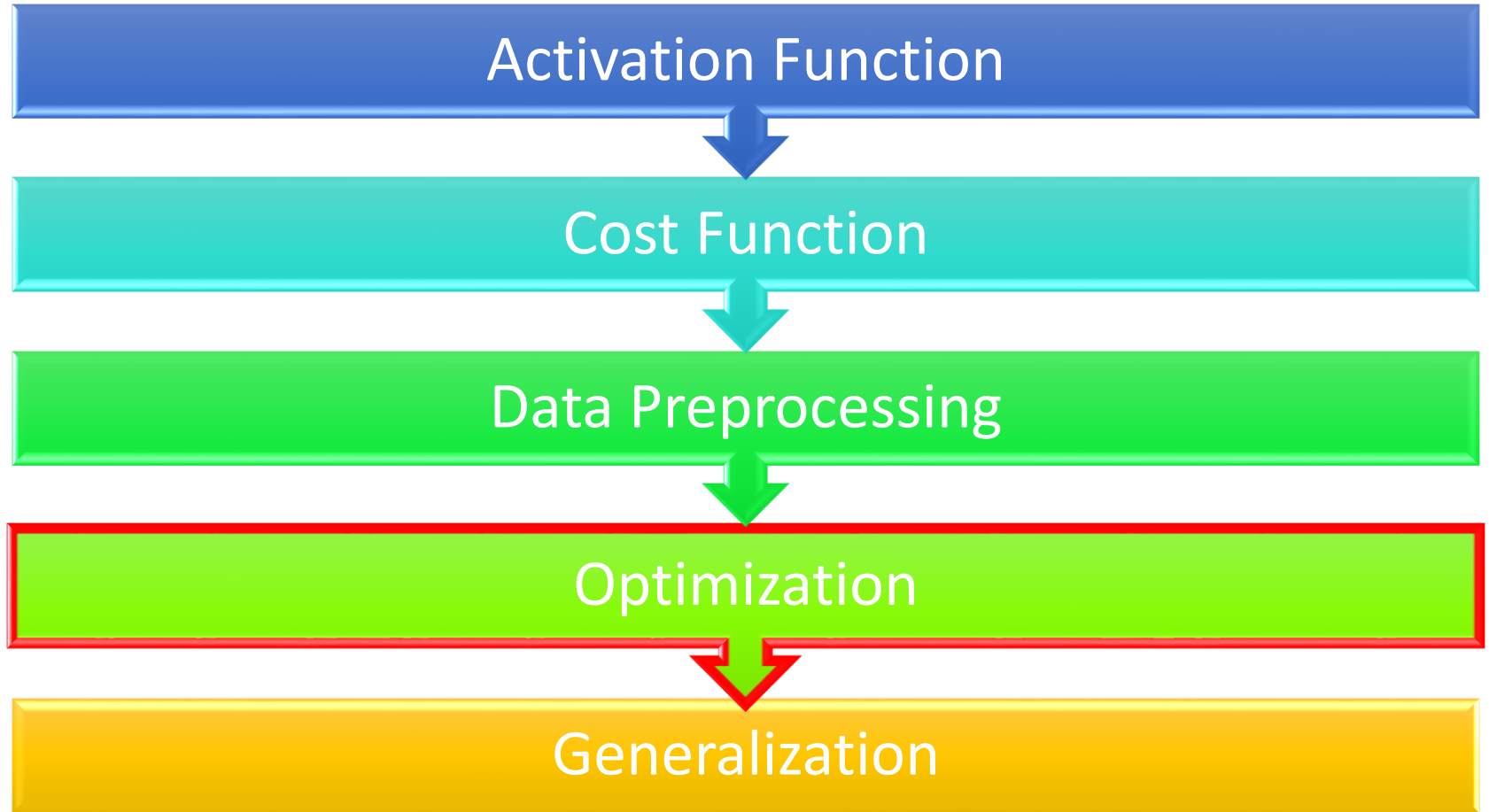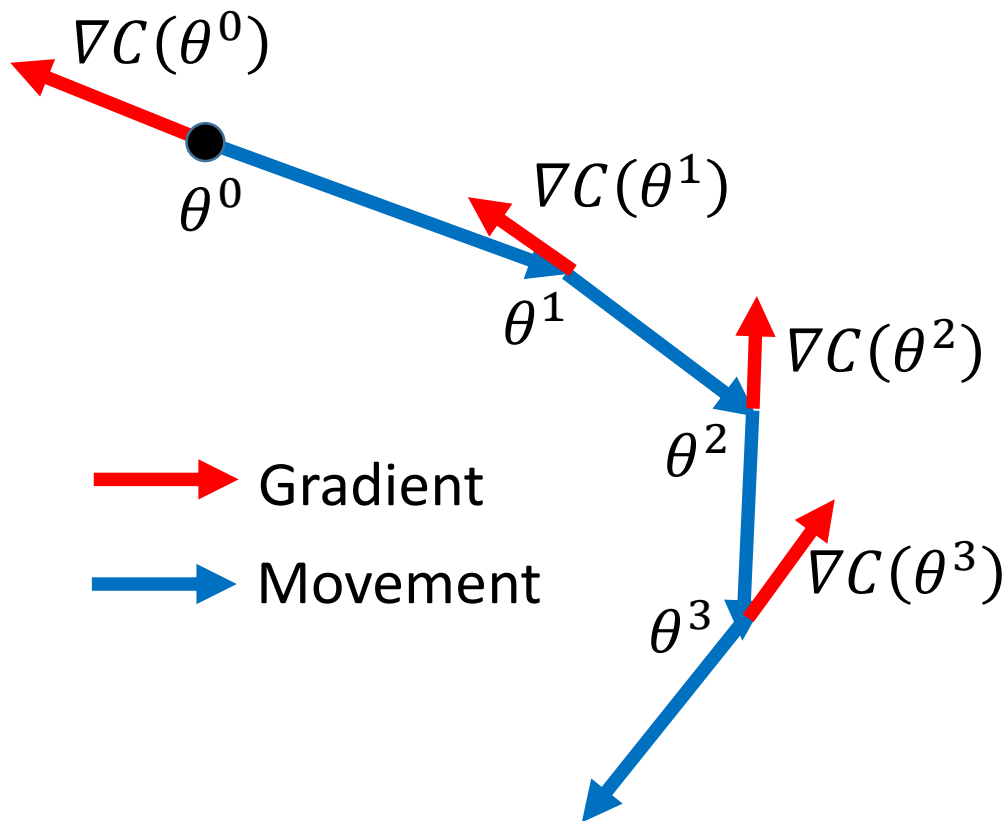The means of all dimensions are 0, and the variances are all 1

$x^1 \quad x^2 \quad x^3 \quad x^r \quad x^R$

# Normalizing Input



Source of figure: http://cs231n.github.io/neural-networks-2/

Normalizing your training and testing data in the same way.

# Outline

# Vanilla Gradient Descent



$\nabla C(\theta^0)$

$\theta^0$

$\nabla C(\theta^1)$

$\theta^1$

$\nabla C(\theta^2)$

$\theta^2$

$\nabla C(\theta^3)$

$\theta^3$

Gradient

Movement

Start at position $\theta^0$

Compute gradient at $\theta^0$

Move to $\theta^1 = \theta^0 - \eta \nabla C(\theta^0)$
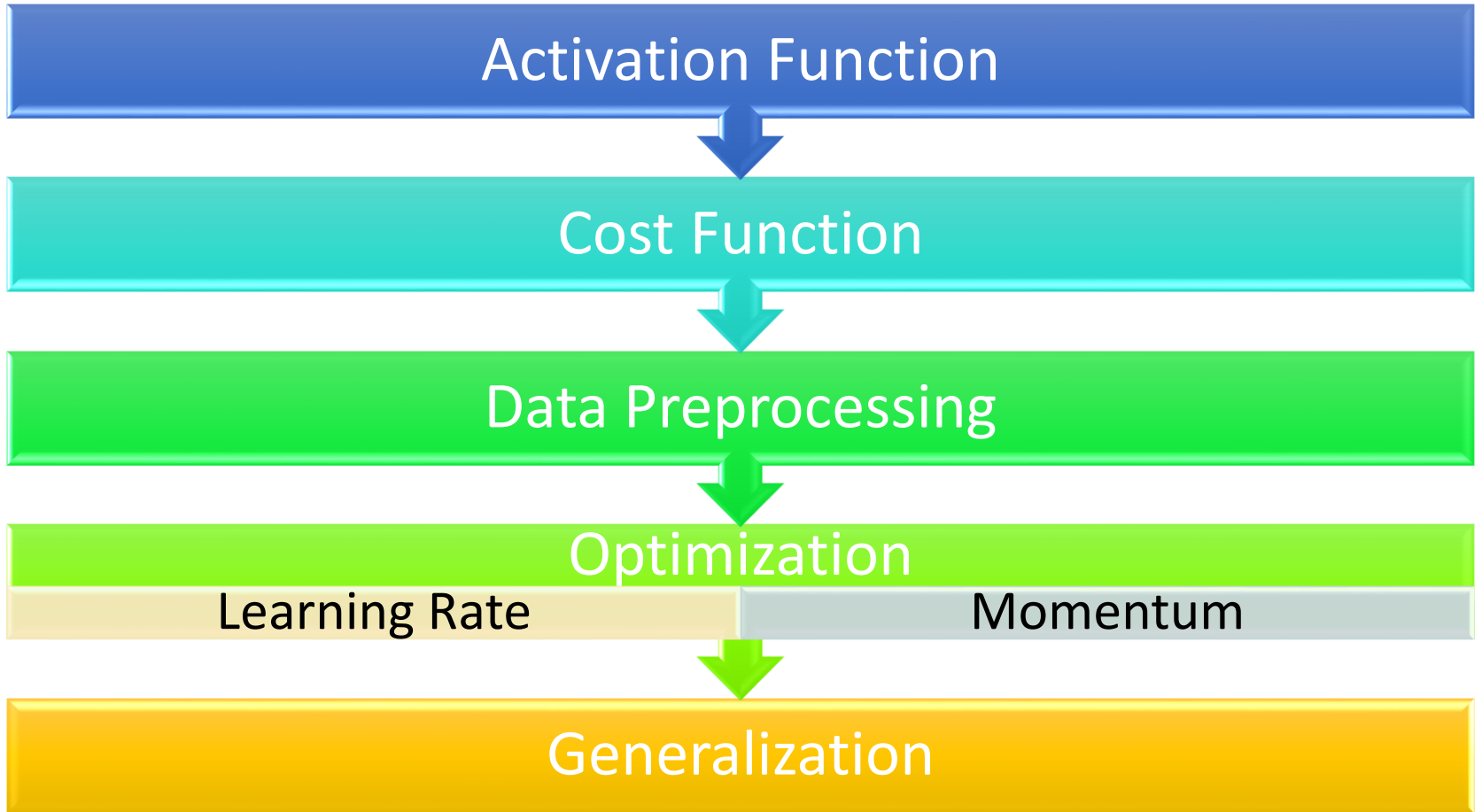
Compute gradient at $\theta^1$

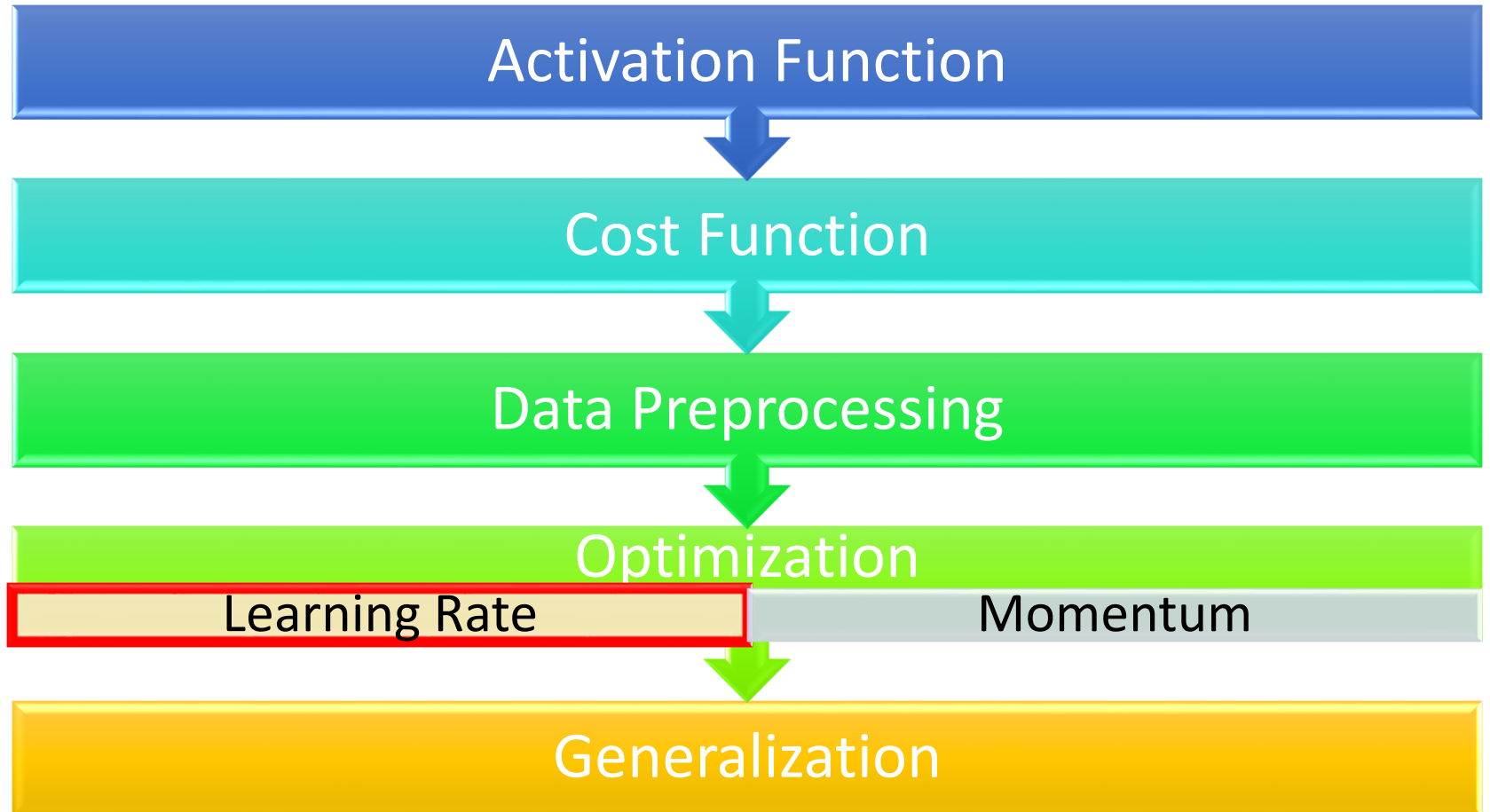Move to $\theta^2 = \theta^1 - \eta \nabla C(\theta^1)$

Stop until $\nabla C(\theta^t) \approx 0$

1. How to determine the learning rates
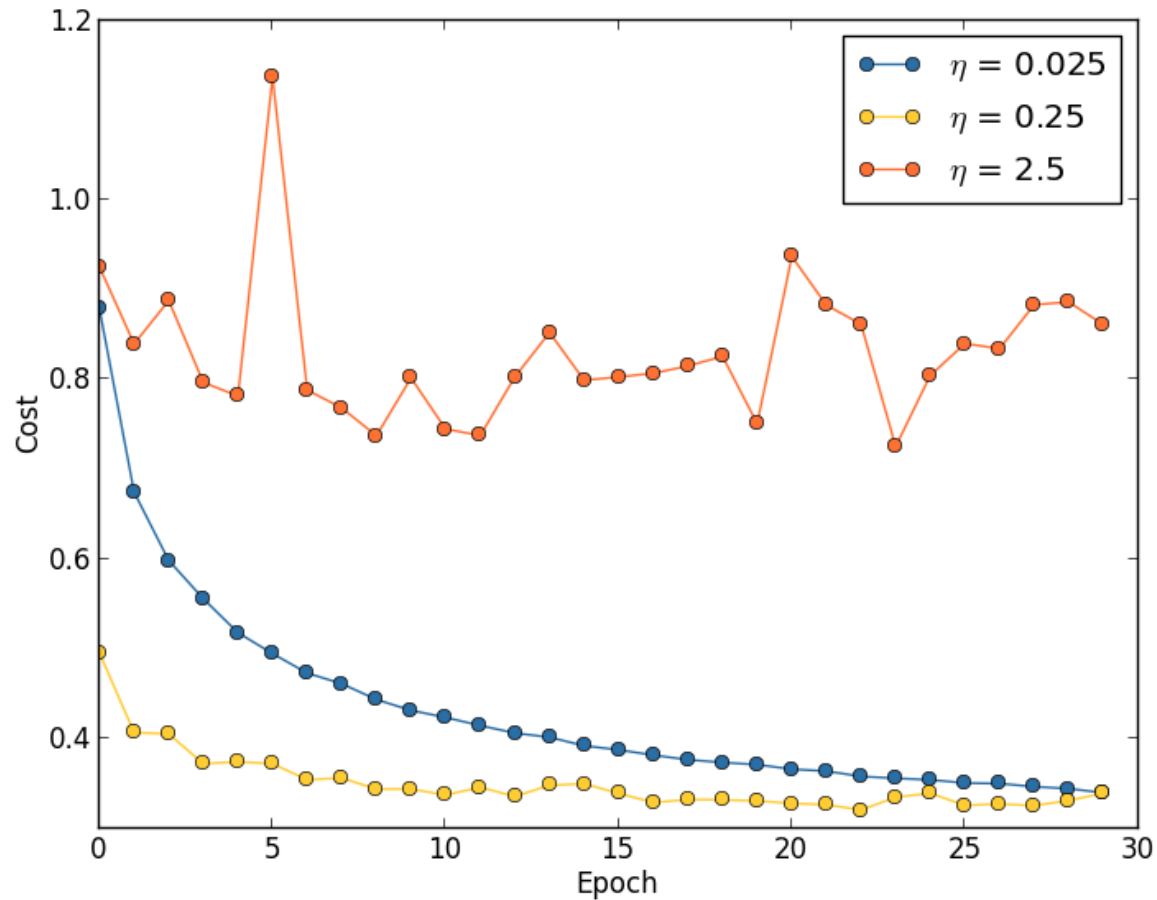2. Stuck at local minima or saddle points

# Outline

# Outline

Activation Function

Cost Function

Data Preprocessing

Optimization

Learning Rate | Momentum

Generalization

# Learning Rates



Source:
http://neuralnetworksanddeeplearning.com/chap3.html

# Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
  - At the beginning, we are far from the destination, so we use larger learning rate
  - After several epochs, we are close to the destination, so we reduce the learning rate
  - E.g. 1/t decay: $\eta^t = \eta/\sqrt{t+1}$
- Learning rate cannot be one-size-fits-all
  - Give different parameters different learning rates

# Adagrad

$$g^t = \frac{\partial C(\theta^t)}{\partial w} \qquad \eta^t = \frac{\eta}{\sqrt{t+1}}$$

- Divide the learning rate of each parameter by the ***root mean square of its previous derivatives***

**_Vanilla Gradient descent_**

$$w^{t+1} \leftarrow w^t - \eta^t g^t$$

w is one parameters

**_Adagrad_**

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$\sigma^t$: ***root mean square*** of the previous derivatives of parameter w

Parameter dependent

# Adagrad

$$w^1 \leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0$$

$$\sigma^0 = g^0$$

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1$$

$$\sigma^1 = \sqrt{\frac{1}{2}[(g^0)^2 + (g^1)^2]}$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2$$

$$\sigma^2 = \sqrt{\frac{1}{3}[(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

$$\vdots$$

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

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

# Adagrad

- Divide the learning rate of each parameter by the *root mean square of its previous derivatives*

$$\eta^t = \frac{\eta}{\sqrt{t+1}}$$   <span style="color:red">1/t decay</span>

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

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

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^{t} (g^i)^2}} g^t$$

# Contradiction?

$$g^t = \frac{\partial C(\theta^t)}{\partial w} \quad \eta^t = \frac{\eta}{\sqrt{t+1}}$$

## ***Vanilla Gradient descent***

$$w^{t+1} \leftarrow w^t - \eta^t \underline{g^t}$$

Larger gradient, larger step

## ***Adagrad***

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}} \underline{g^t}$$

Larger gradient, larger step
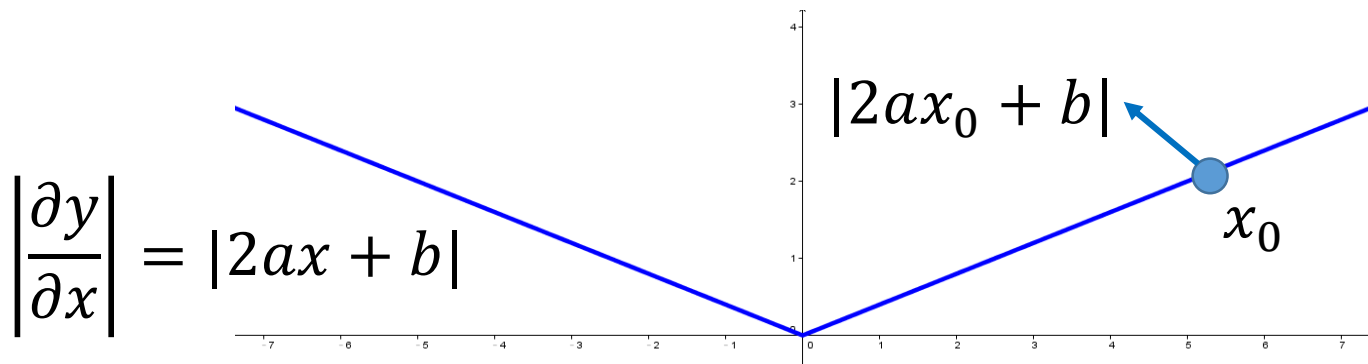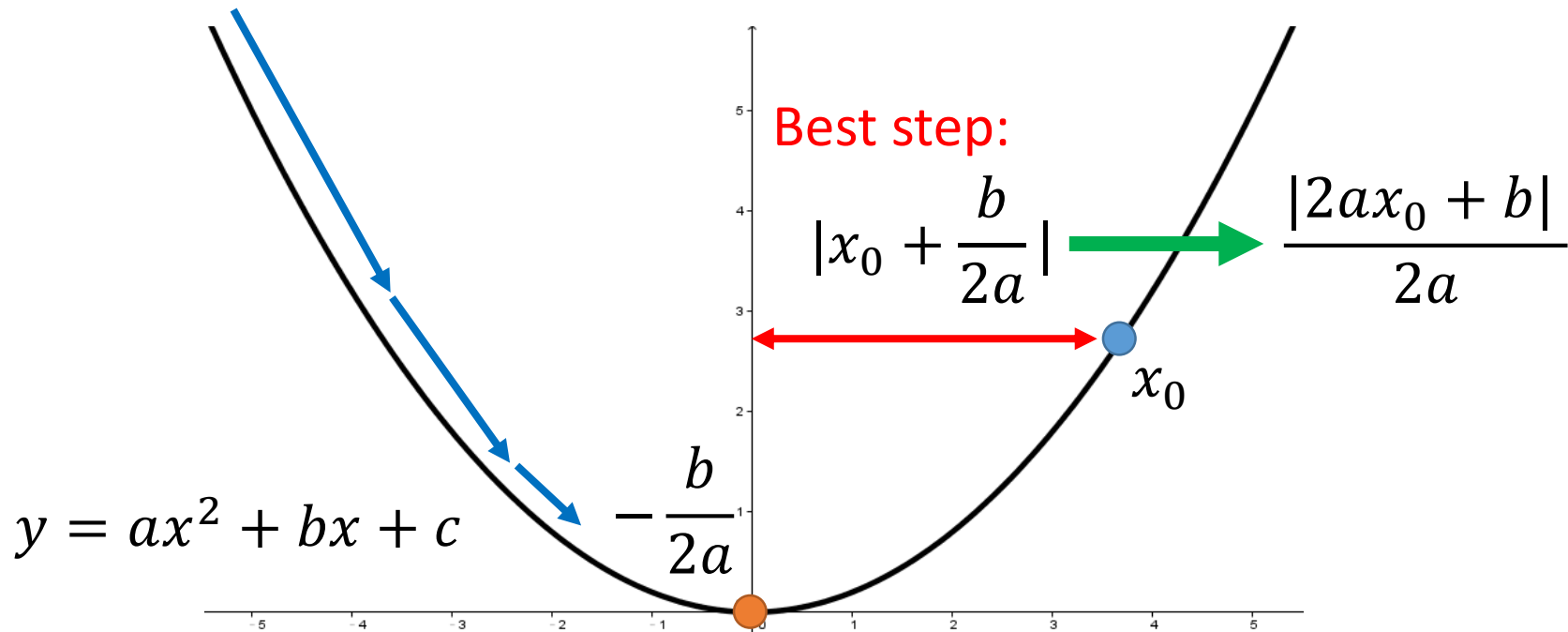
Larger gradient, smaller step

# Intuitive Reason

- *反差*

| $g^0$ | $g^1$ | $g^2$ | $g^3$ | $g^4$ | ...... |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 0.001 | 0.001 | 0.003 | 0.002 | 0.1 | ...... |

| $g^0$ | $g^1$ | $g^2$ | $g^3$ | $g^4$ | ...... |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 10.8 | 20.9 | 31.7 | 12.1 | 0.1 | ...... |

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^{t}(g^i)^2}} g^t$$

造成反差的效果

# Larger gradient, larger steps?

$$y = ax^2 + bx + c$$

$$-\frac{b}{2a}$$

Best step:

$$\left| x_0 + \frac{b}{2a} \right|$$

$$\frac{|2ax_0 + b|}{2a}$$

$x_0$

$$\left| \frac{\partial y}{\partial x} \right| = |2ax + b|$$

$$|2ax_0 + b|$$

$x_0$

# Second Derivative

Best step:

$$|x_0 + \frac{b}{2a}| \quad \xrightarrow{\hspace{2cm}} \quad \frac{|2ax_0 + b|}{\boxed{2a}}$$

$$y = ax^2 + bx + c$$

$$-\frac{b}{2a}$$

$x_0$

$$\left|\frac{\partial y}{\partial x}\right| = |2ax + b|$$

$$|2ax_0 + b|$$
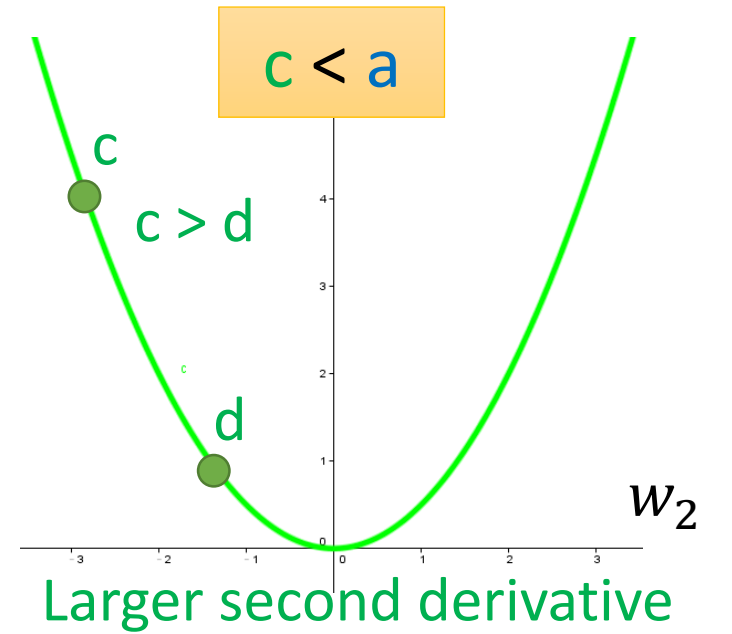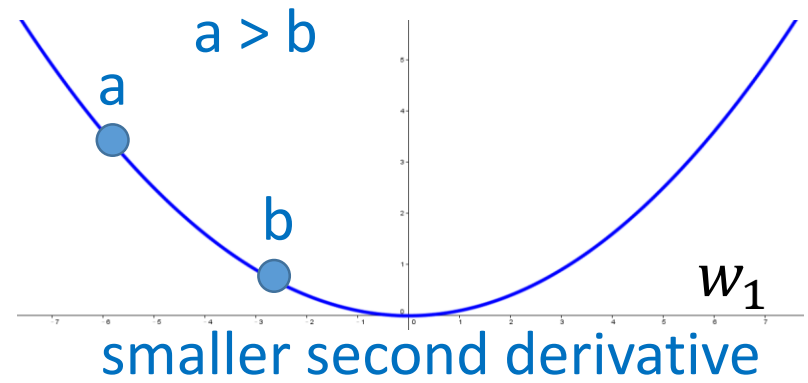
$x_0$

$$\frac{\partial^2 y}{\partial x^2} = 2a$$
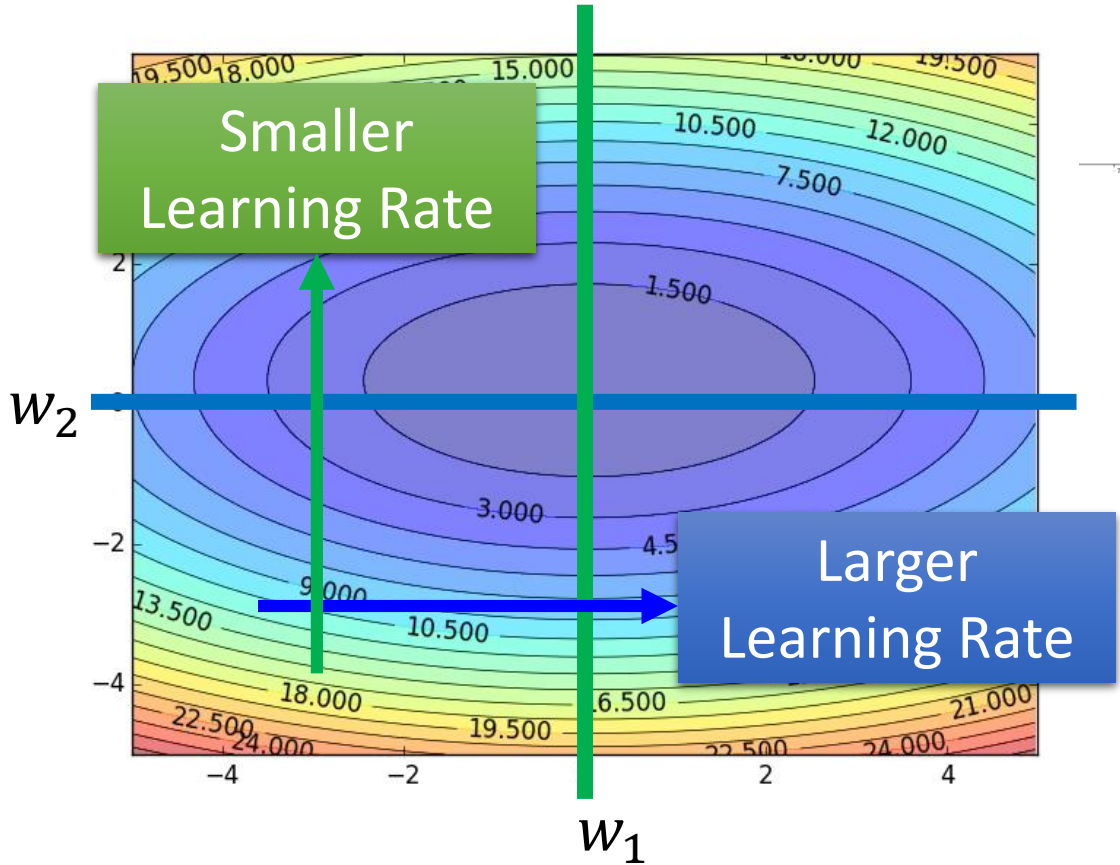
The best step is $\dfrac{|\text{First derivative}|}{\text{Second derivative}}$

# More than one parameters

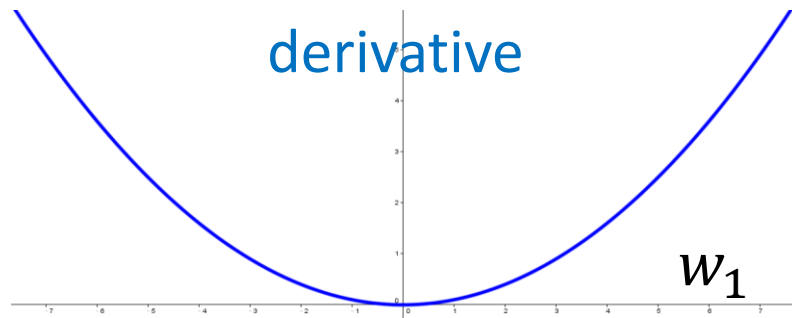The best step is $\dfrac{|\text{First derivative}|}{\text{Second derivative}}$



a > b

a

b

smaller second derivative

c < a

c

c > d

d

Larger second derivative

Smaller Learning Rate

Larger Learning Rate

$w_2$

$w_1$

$w_1$

$w_2$

# What to do with Adagrad?

The best step is

$$\frac{|\text{First derivative}|}{\text{Second derivative}}$$
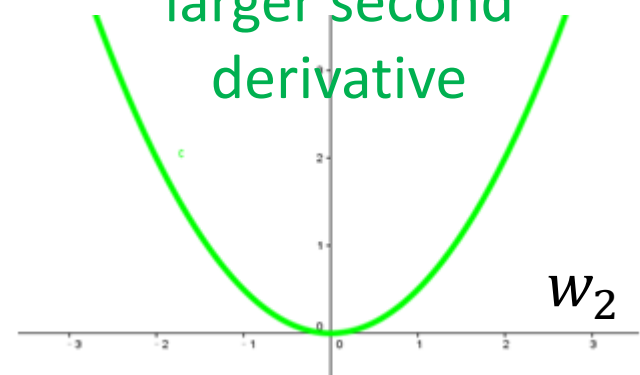
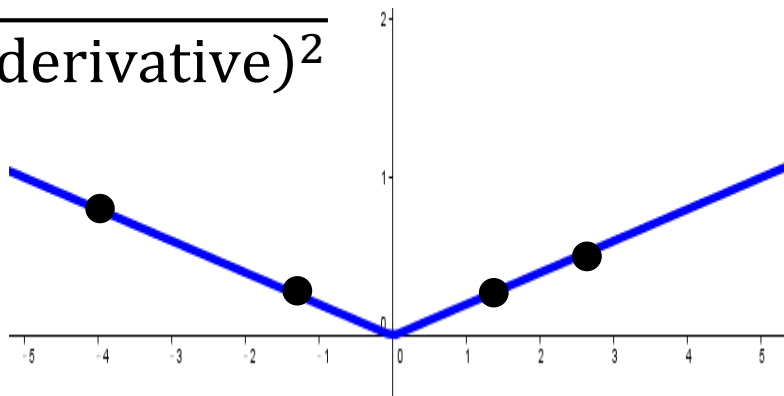Use *first derivative* to estimate *second derivative*

smaller second derivative

$w_1$

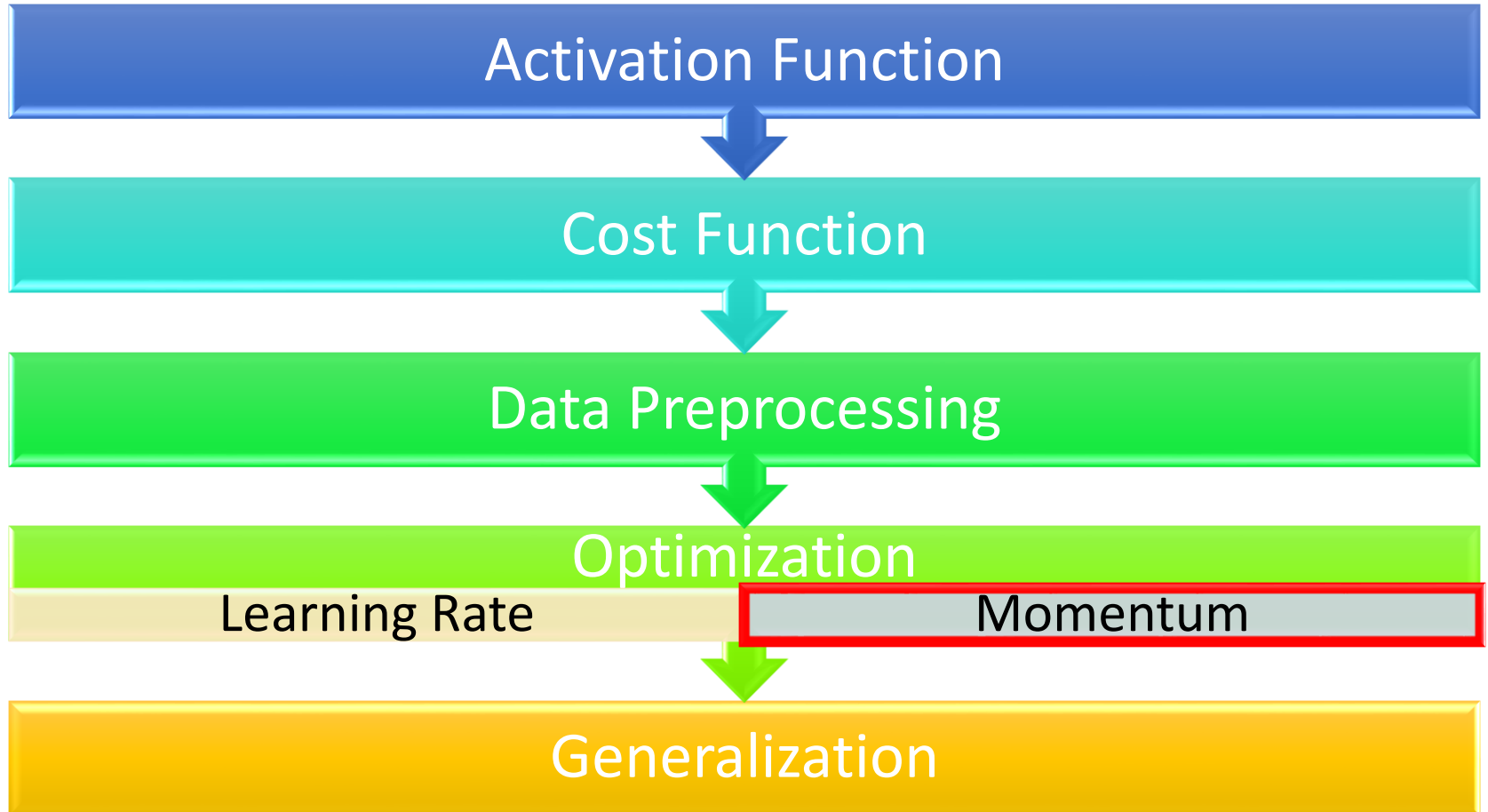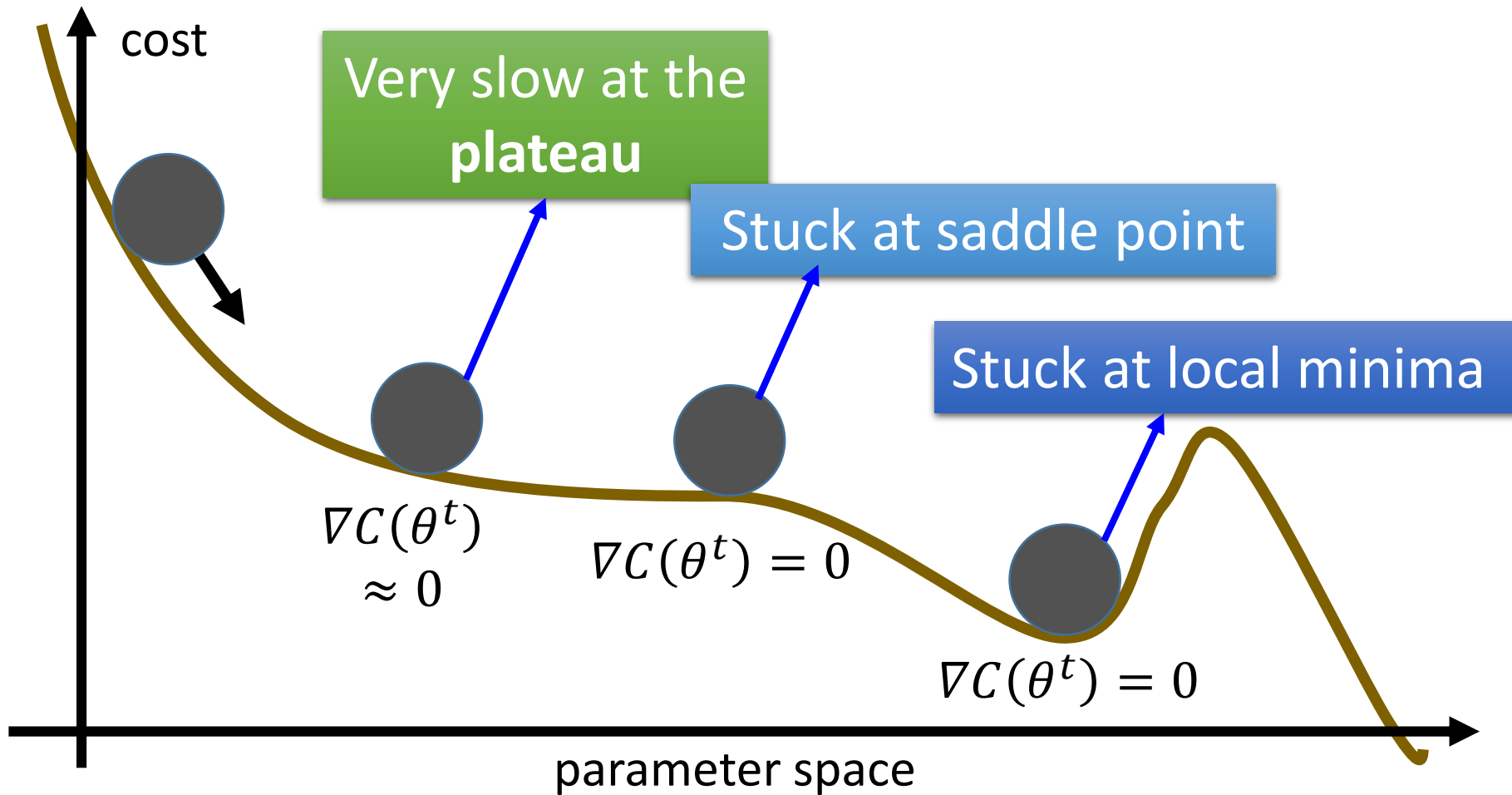larger second derivative

$w_2$

$$\sqrt{(\text{first derivative})^2}$$

# Outline

Activation Function

Cost Function

Data Preprocessing

Optimization

Learning Rate

Momentum

Generalization

# Easy to stuck



cost

Very slow at the **plateau**

Stuck at saddle point

Stuck at local minima

$\nabla C(\theta^t) \approx 0$

$\nabla C(\theta^t) = 0$

$\nabla C(\theta^t) = 0$

parameter space

# In physical world ……

- Momentum



How about put this phenomenon in gradient descent?

# Momentum

Movement: movement of last step minus gradient at present

$\nabla C(\theta^0)$

$\nabla C(\theta^1)$

$\theta^0$

$\theta^1$

$\nabla C(\theta^2)$

$\theta^2$

$\theta^3$

$\nabla C(\theta^3)$

→ Gradient

→ Movement

⋯ Movement of last step

Start at point $\theta^0$

Movement $v^0$=0

Compute gradient at $\theta^0$

Movement $v^1 = \lambda v^0 - \eta \nabla C(\theta^0)$

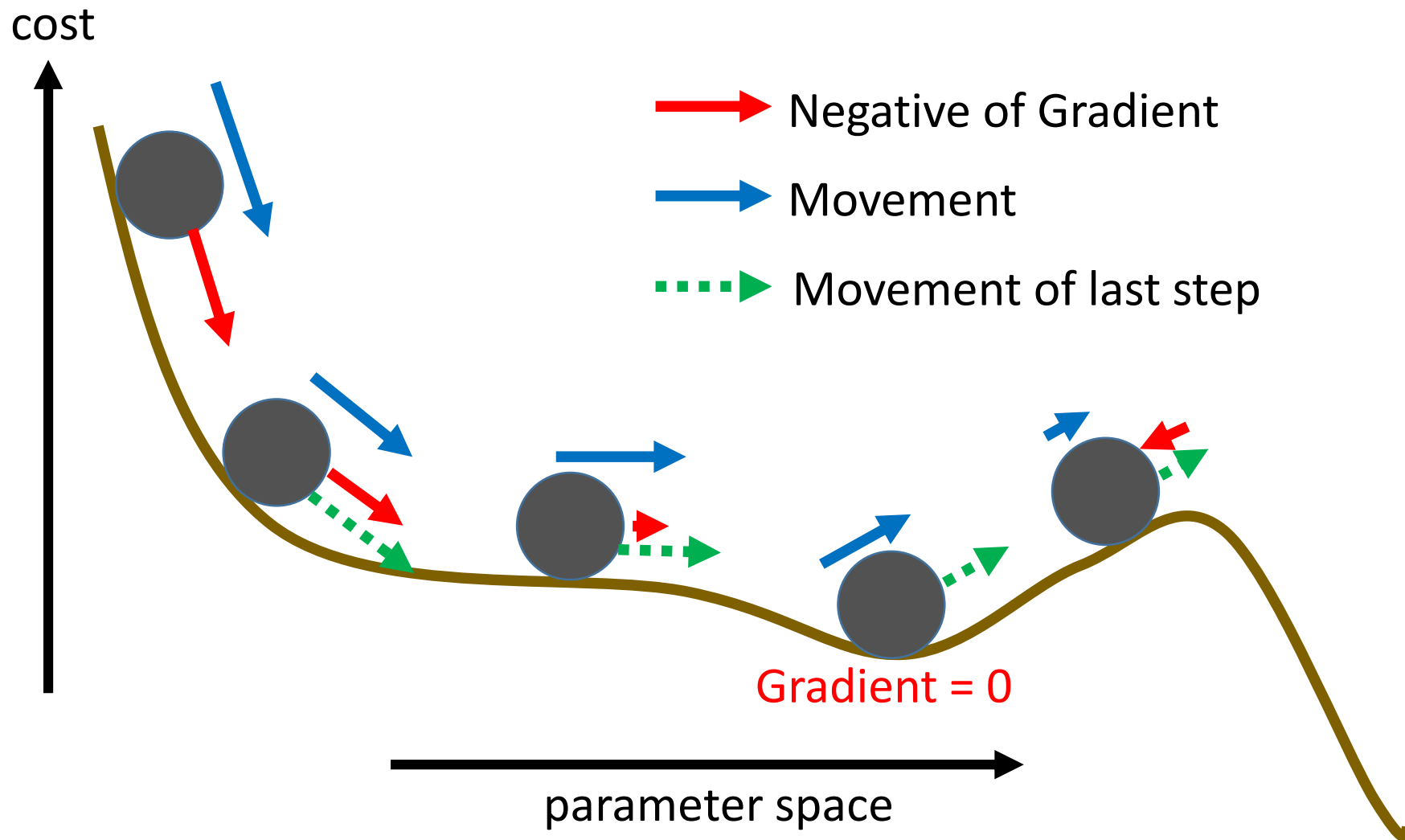Move to $\theta^1 = \theta^0 + v^1$

Compute gradient at $\theta^1$

Movement $v^2 = \lambda v^1 - \eta \nabla C(\theta^1)$

Move to $\theta^2 = \theta^1 + v^2$

Movement not just based on gradient, but previous movement.

# Momentum

Movement: movement of last step minus gradient at present

$v^i$ is actually the weighted sum of all the previous gradient:
$\nabla C(\theta^0), \nabla C(\theta^1), \dots \nabla C(\theta^{i-1})$

$v^0 = 0$

$v^1 = -\eta \nabla C(\theta^0)$

$v^2 = -\lambda \eta \nabla C(\theta^0) - \eta \nabla C(\theta^1)$

⋮

Start at point $\theta^0$

Movement $v^0 = 0$

Compute gradient at $\theta^0$

Movement $v^1 = \lambda v^0 - \eta \nabla C(\theta^0)$

Move to $\theta^1 = \theta^0 + v^1$

Compute gradient at $\theta^1$

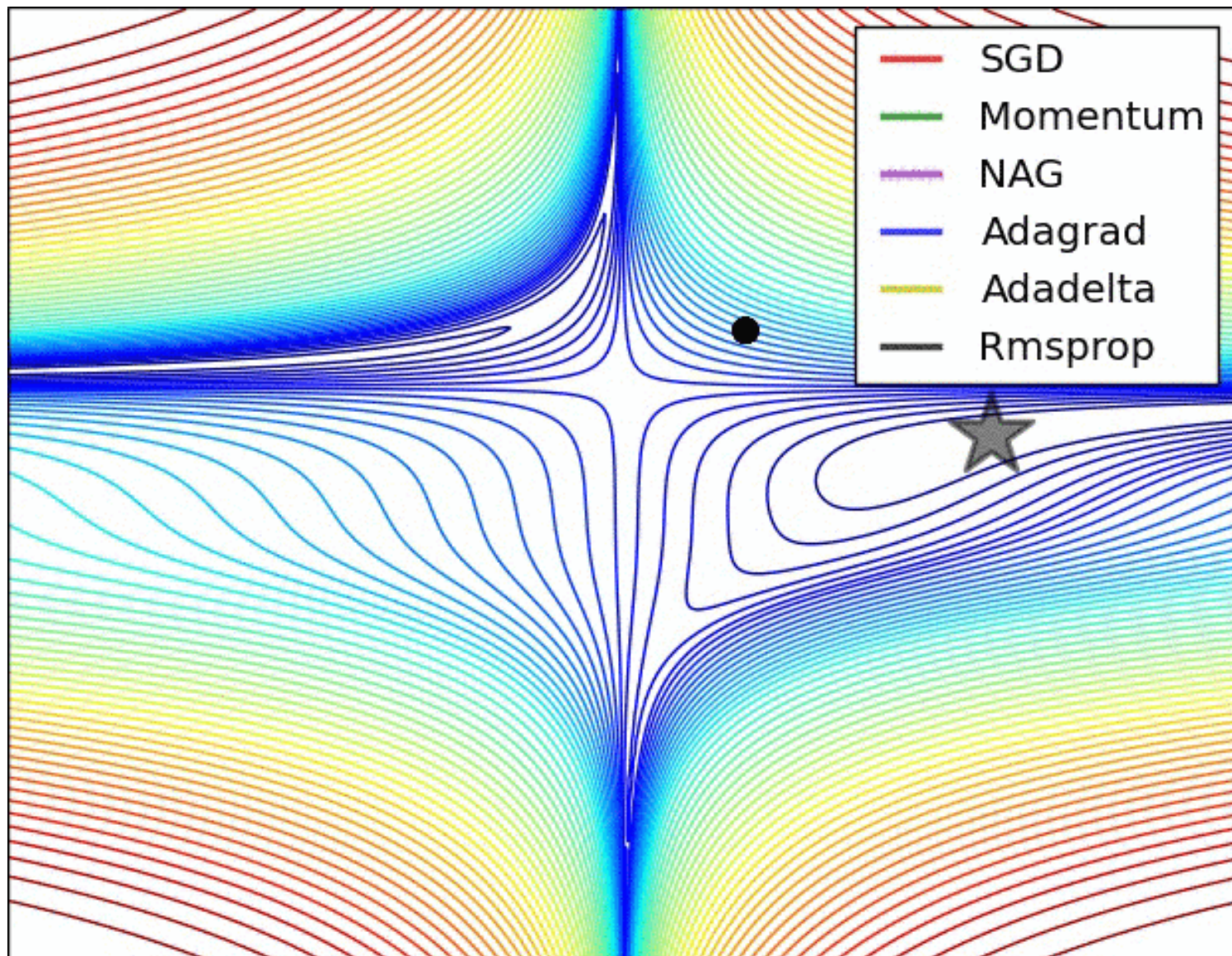Movement $v^2 = \lambda v^1 - \eta \nabla C(\theta^1)$

Move to $\theta^2 = \theta^1 + v^2$

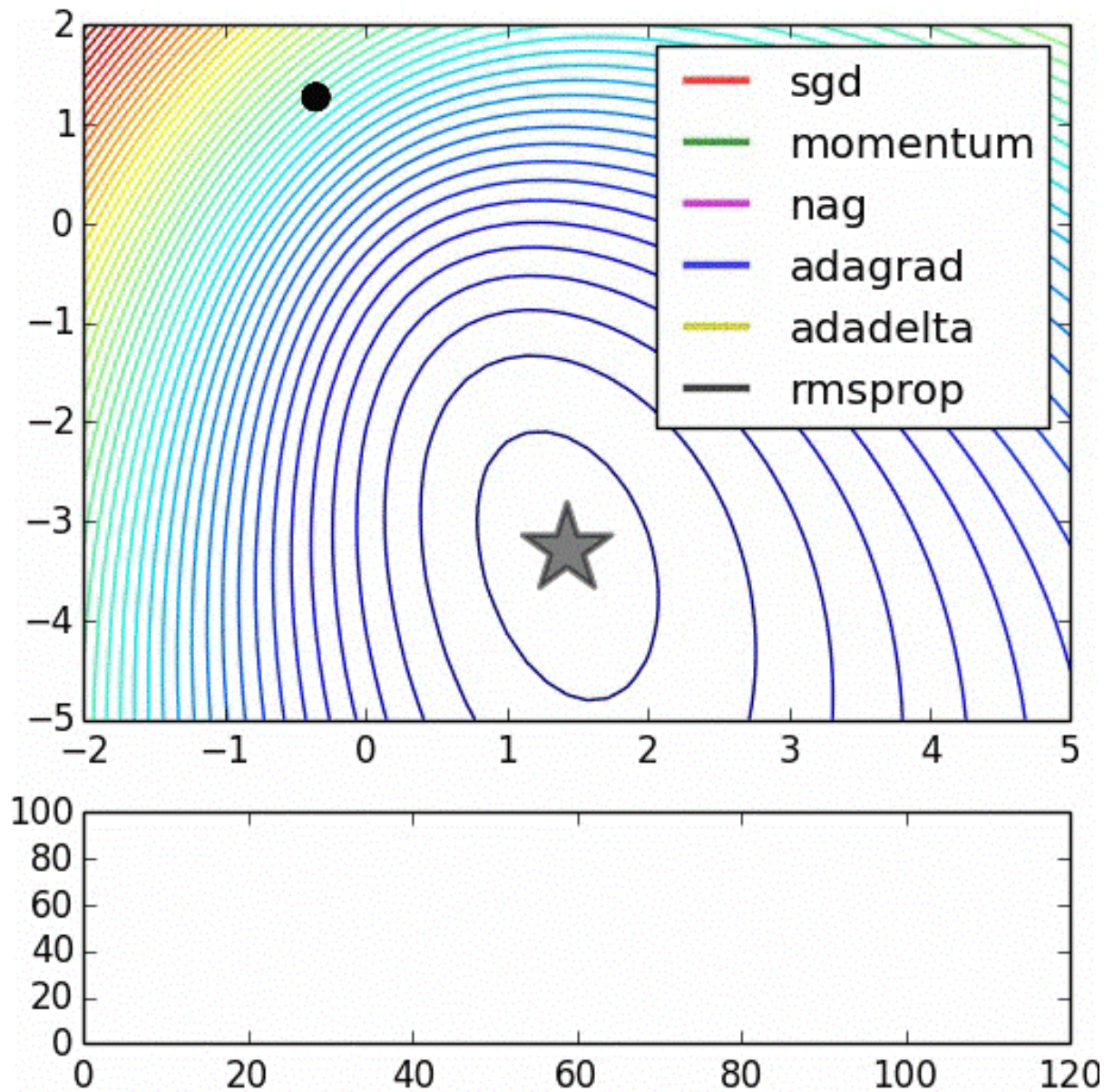Movement not just based on gradient, but previous movement
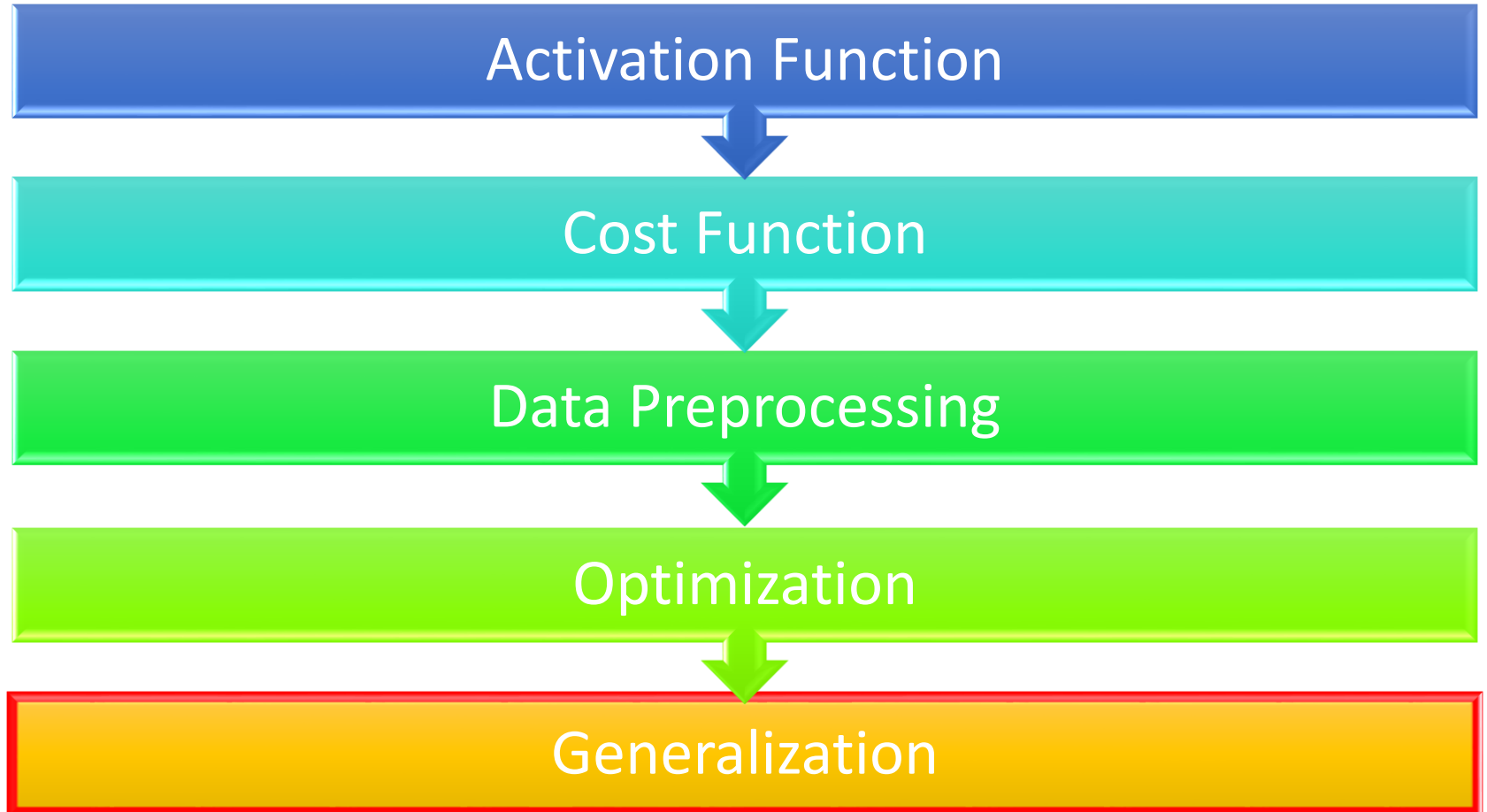
SGD
Momentum
NAG
Adagrad
Adadelta
Rmsprop

http://www.reddit.com/r/MachineLearning/comments/2gopfa/visualizing_gradient_optimization_techniques/cklhott **(By Alec Radford)**

# Outline

Activation Function

↓

Cost Function

↓
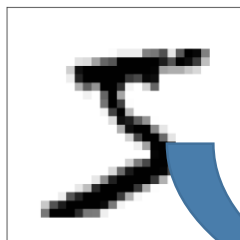
Data Preprocessing
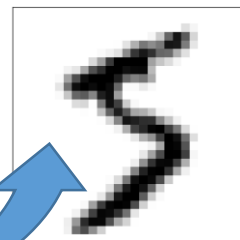
↓

Optimization

↓

Generalization

# Panacea

- Have more training data
- *Create* more training data (?)
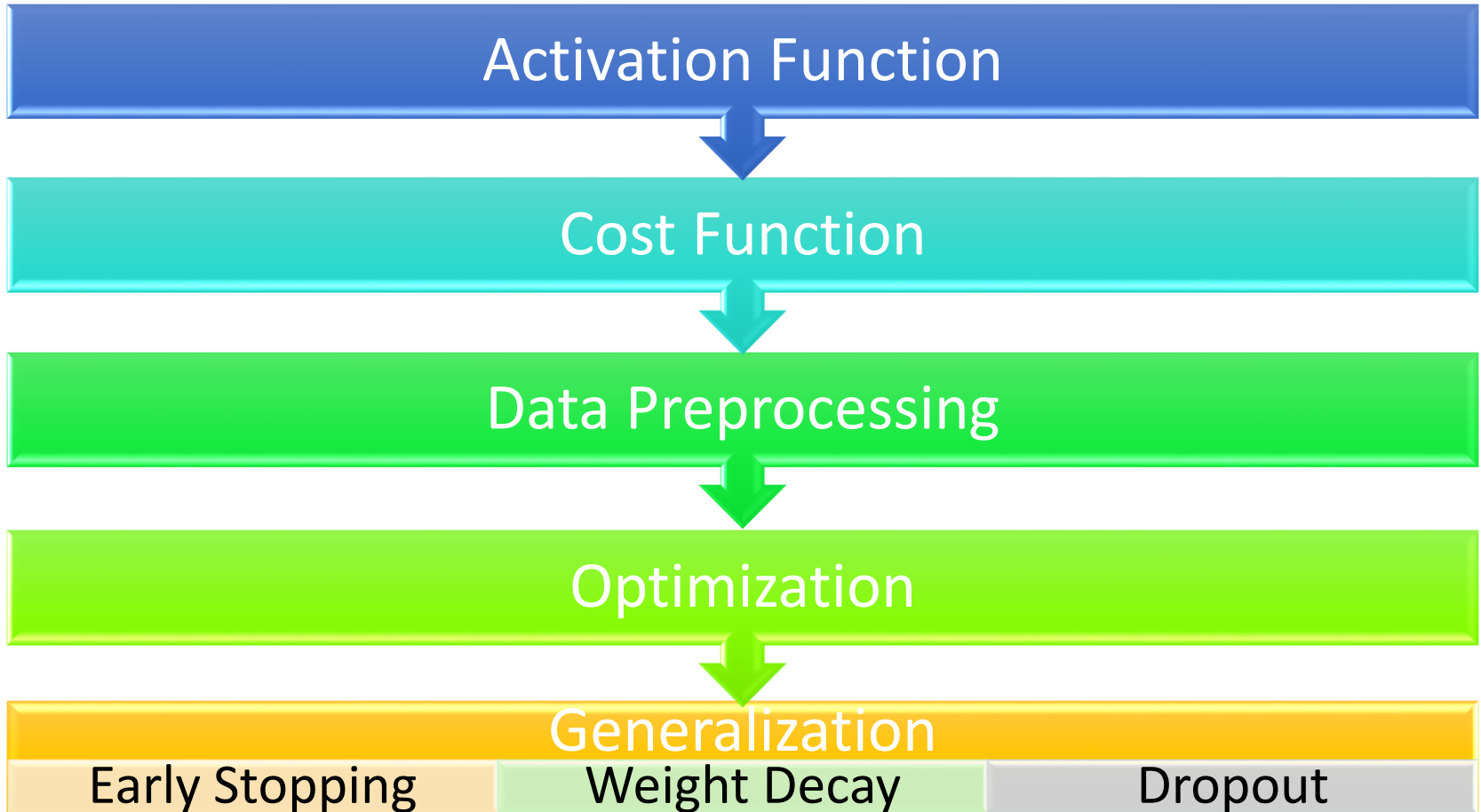
Handwriting recognition:

Original
Training Data:

Created
Training Data:

Shift 15°

# Outline

Activation Function

↓

Cost Function

↓

Data Preprocessing

↓

Optimization

↓

Generalization

| Early Stopping | Weight Decay | Dropout |

# Outline

Activation Function

↓

Cost Function

↓

Data Preprocessing

↓

Optimization

↓

Generalization

| Early Stopping | Weight Decay | Dropout |

# Early Stopping



How many parameter updates do we need?

# Outline

Activation Function

Cost Function

Data Preprocessing

Optimization

Generalization

| Early Stopping | Weight Decay | Dropout |

# Weight Decay

- The parameters closer to zero is preferred.

Training data:

$$\{(x, \hat{y}), \dots\}$$

$$x \begin{cases} x_1 \ w_1 \\ x_2 \ w_2 \\ \vdots \quad w_i \\ x_i \\ \vdots \end{cases} \xrightarrow{\quad} \sigma(z) \longrightarrow \hat{y}$$

$$z = w \cdot x$$

Testing data:

$$\{(x', \hat{y}), \dots\}$$

$$x' = x + \varepsilon$$

$$z' = w \cdot (x + \varepsilon)$$
$$= w \cdot x + w \cdot \varepsilon$$
$$= z + w \cdot \varepsilon$$

To minimize the effect of noise, we want w close to zero.

# Weight Decay

- New cost function to be minimized
  - Find a set of weight not only minimizing original cost but also close to zero

$$C'(\theta) = \underline{C(\theta)} + \lambda \frac{1}{2} \underline{\|\theta\|^2}$$ → Regularization term:

$$\theta = \{\mathbf{W}^1, \mathbf{W}^2, \ldots\}$$

Original cost (e.g. minimize square error, cross entropy …)

$$\|\theta\|^2 = \left(w_{11}^1\right)^2 + \left(w_{12}^1\right)^2 + \ldots$$

$$+ \left(w_{11}^2\right)^2 + \left(w_{12}^2\right)^2 + \ldots$$

(not consider biases. why?)

# Weight Decay

$$\|\theta\|^2 = \left(w_{11}^1\right)^2 + \left(w_{12}^1\right)^2 + \ldots$$
$$+ \left(w_{11}^2\right)^2 + \left(w_{12}^2\right)^2 + \ldots$$
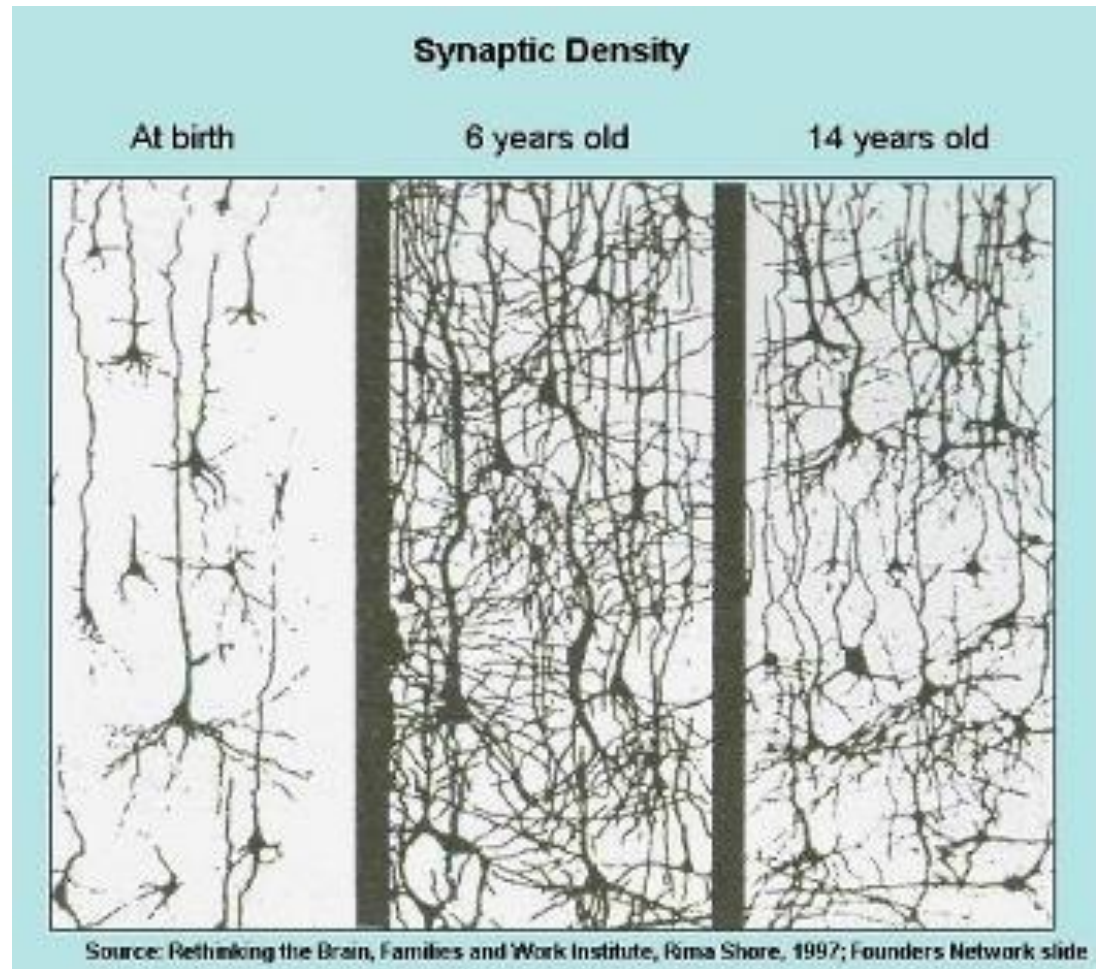
- New cost function to be minimized

$$C'(\theta) = C(\theta) + \lambda \frac{1}{2}\|\theta\|^2 \quad \text{Gradient:} \quad \frac{\partial C'}{\partial w} = \frac{\partial C}{\partial w} + \lambda w$$

Update:
$$w^{t+1} \to w^t - \eta \frac{\partial C'}{\partial w} = w^t - \eta \left(\frac{\partial C}{\partial w} + \lambda w^t\right)$$
$$= (1 - \eta\lambda)w^t - \eta \frac{\partial C}{\partial w}$$

Smaller and smaller

# Weight Decay

- Our Brain



**Synaptic Density**

At birth     6 years old     14 years old

Source: Rethinking the Brain, Families and Work Institute, Rima Shore, 1997; Founders Network slide

# Outline

Activation Function

Cost Function

Data Preprocessing

Optimization

Generalization

| Early Stopping | Weight Decay | Dropout |

# Dropout

$$\theta^t \leftarrow \theta^{t-1} - \eta \nabla C_x(\theta^{t-1})$$
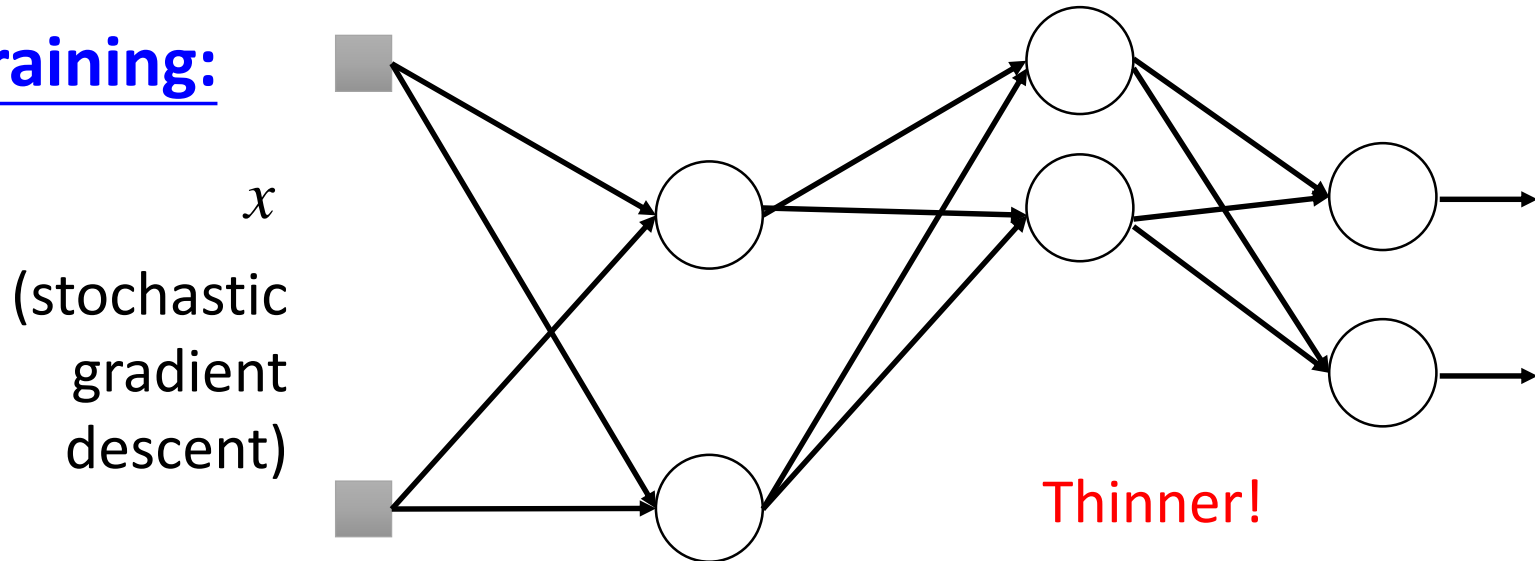
**Training:**

$x$

(stochastic
gradient
descent)



➢ In each *iteration*

● Each neuron has p% to dropout

# Dropout

$$\theta^t \leftarrow \theta^{t-1} - \eta \nabla C_x(\theta^{t-1})$$

**Training:**

$x$

(stochastic gradient descent)

Thinner!

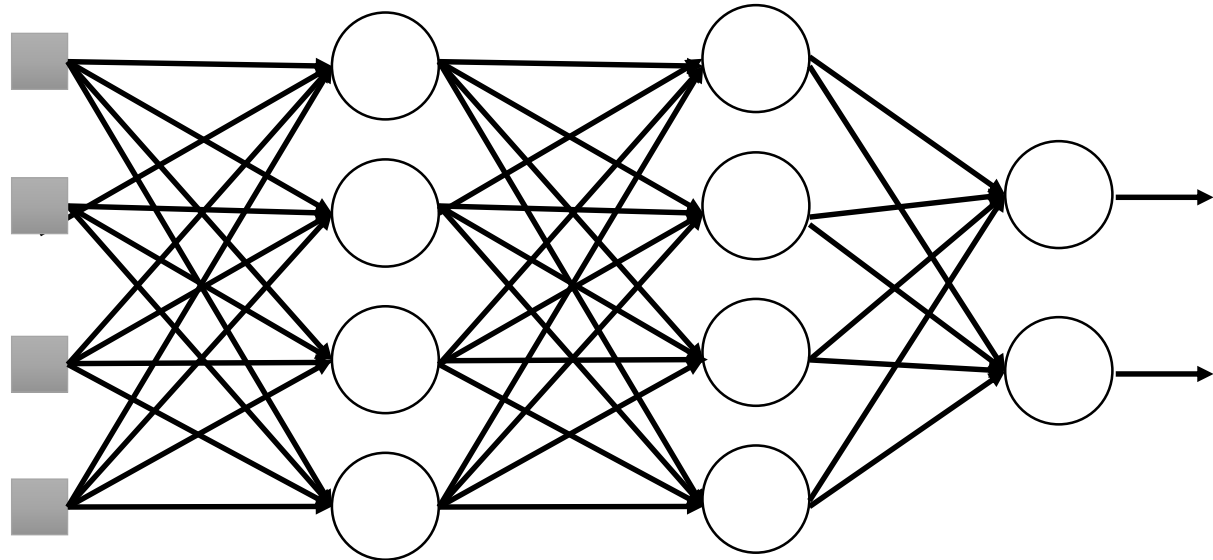➤ In each *iteration*
- Each neuron has p% to dropout

  ➡️ **The structured of the network is changed.**

- Using the new network for training

For each iteration, we resample the dropout neurons

# Dropout

➢ **No dropout**

● If the dropout rate at training is p%,
  all the weights times (1-p)%

● Assume that the dropout rate is 50%.
  If $w_{ij}^l = 1$ from training, set $w_{ij}^l = 0.5$ for testing.

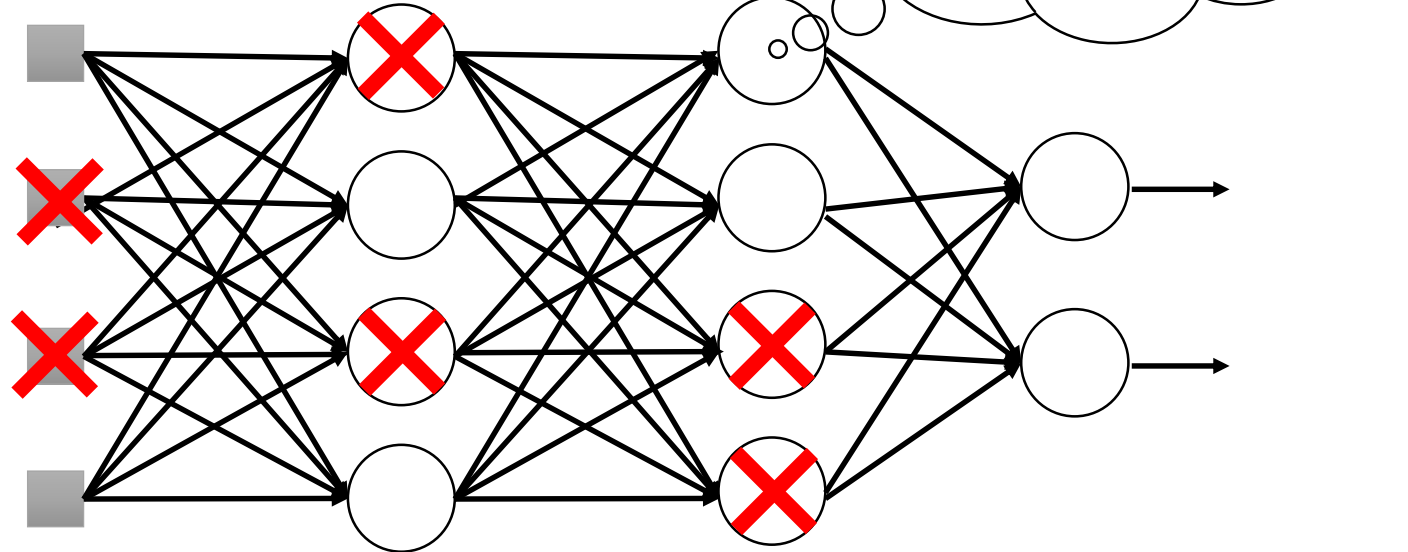# Dropout
# - Intuitive Reason

**_Testing_**

No dropout
(拿下重物後就變很強)

**_Training_**

Dropout (腳上綁重物)

# Dropout
## - Intuitive Reason

我的 partner 會擺爛，所以我要好好做

➤ When teams up, if everyone expect the partner will do the work, nothing will be done finally.

➤ However, if you know your partner will dropout, you will do better.

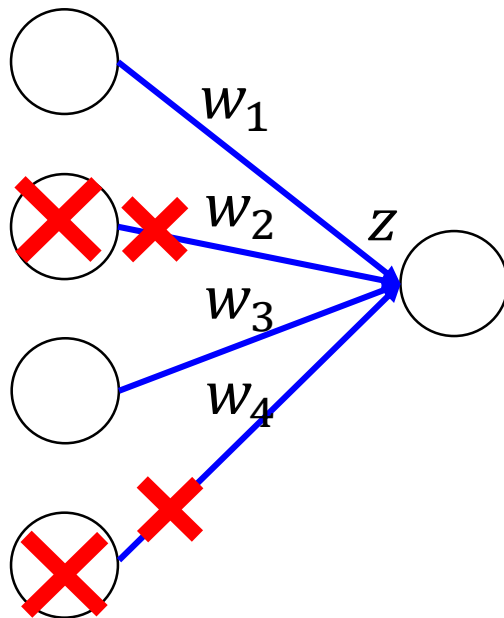➤ When testing, no one dropout actually, so obtaining good results eventually.

# Dropout
# - Intuitive Reason

- Why the weights should multiply (1-p)% (dropout rate) when testing?

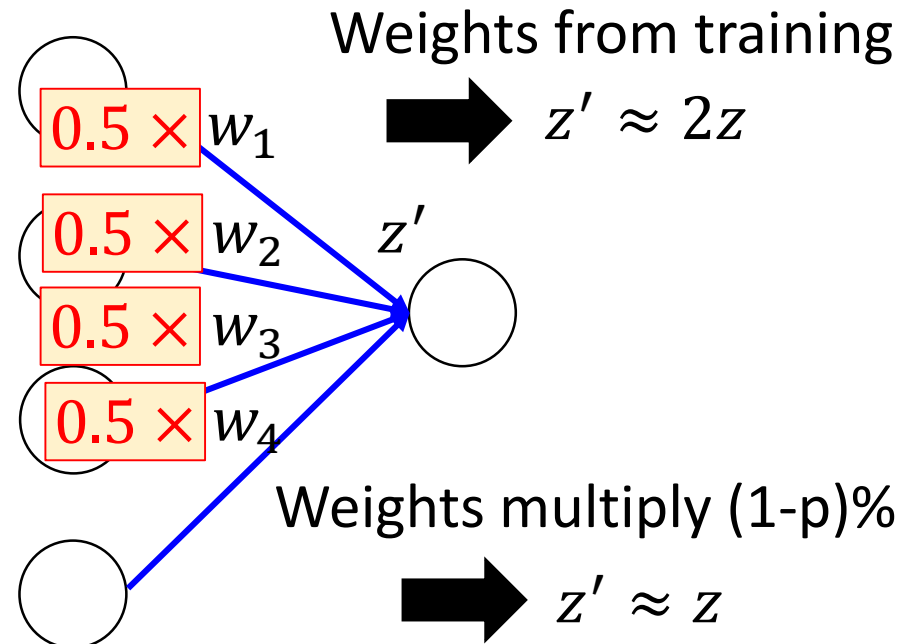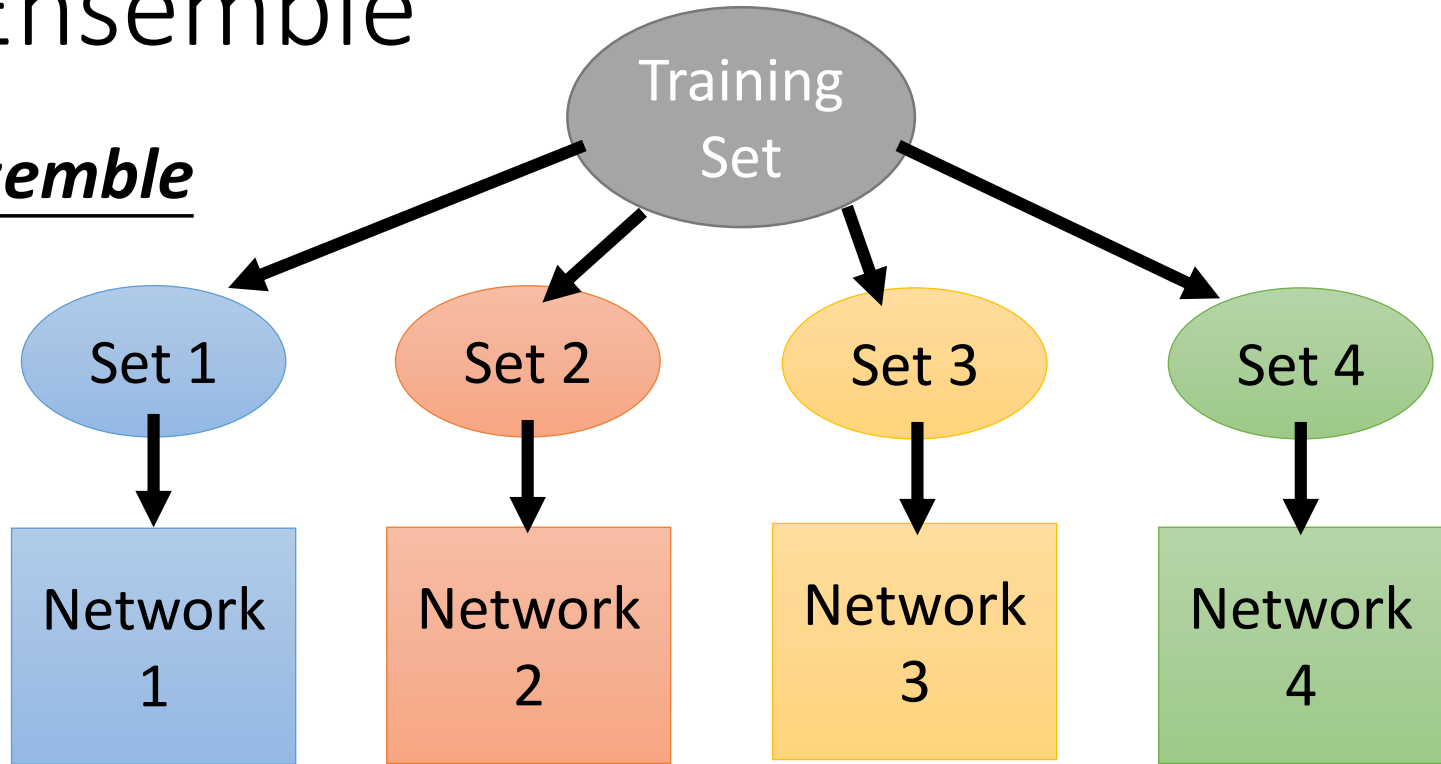***Training of Dropout***

Assume dropout rate is 50%



***Testing of Dropout***

No dropout



Weights from training

$$z' \approx 2z$$

Weights multiply (1-p)%

$$z' \approx z$$

# Dropout
## - Ensemble

***Ensemble***

Train a bunch of networks with different structures

# Dropout
## - Ensemble

### *Ensemble*

Testing data x

Network 1

Network 2

Network 3

Network 4

$y_1$

$y_2$

$y_3$

$y_4$

average

# Dropout
# - Ensemble

**_Dropout ≈ Ensemble_.**

**_Training of Dropout_**

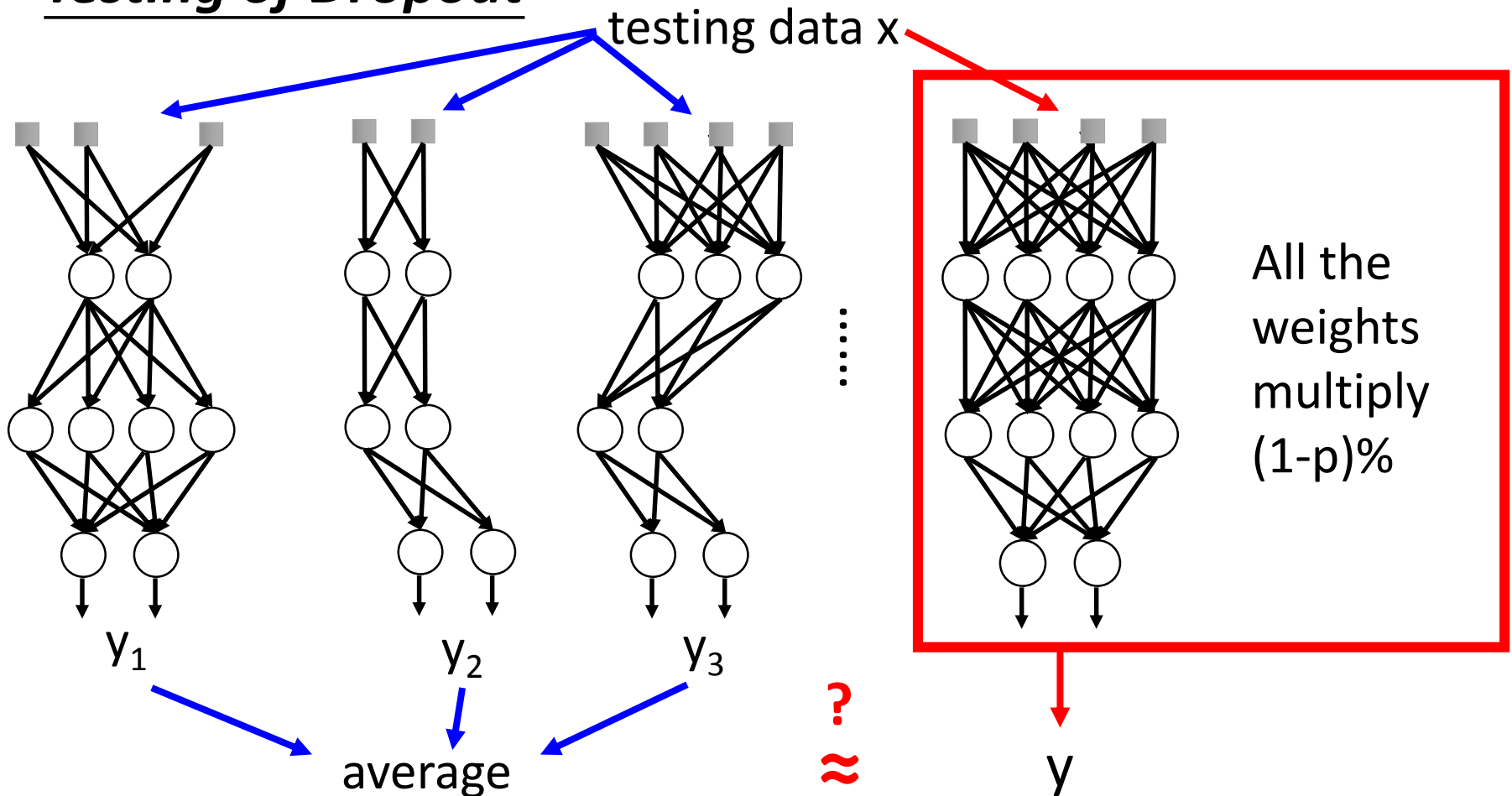$x_1$      $x_2$      $x_3$      $x_4$

M neurons

$2^M$ possible networks

➢ Using one data to train one network

➢ Some parameters in the network are shared

# Dropout - Ensemble

**_Dropout ≈ Ensemble._**

**_Testing of Dropout_**

testing data x



$y_1$

$y_2$

$y_3$

average

All the weights multiply (1-p)%

**?**
**≈**

y

# Dropout
## - Ensemble

- Experiments on hand writing digital classification



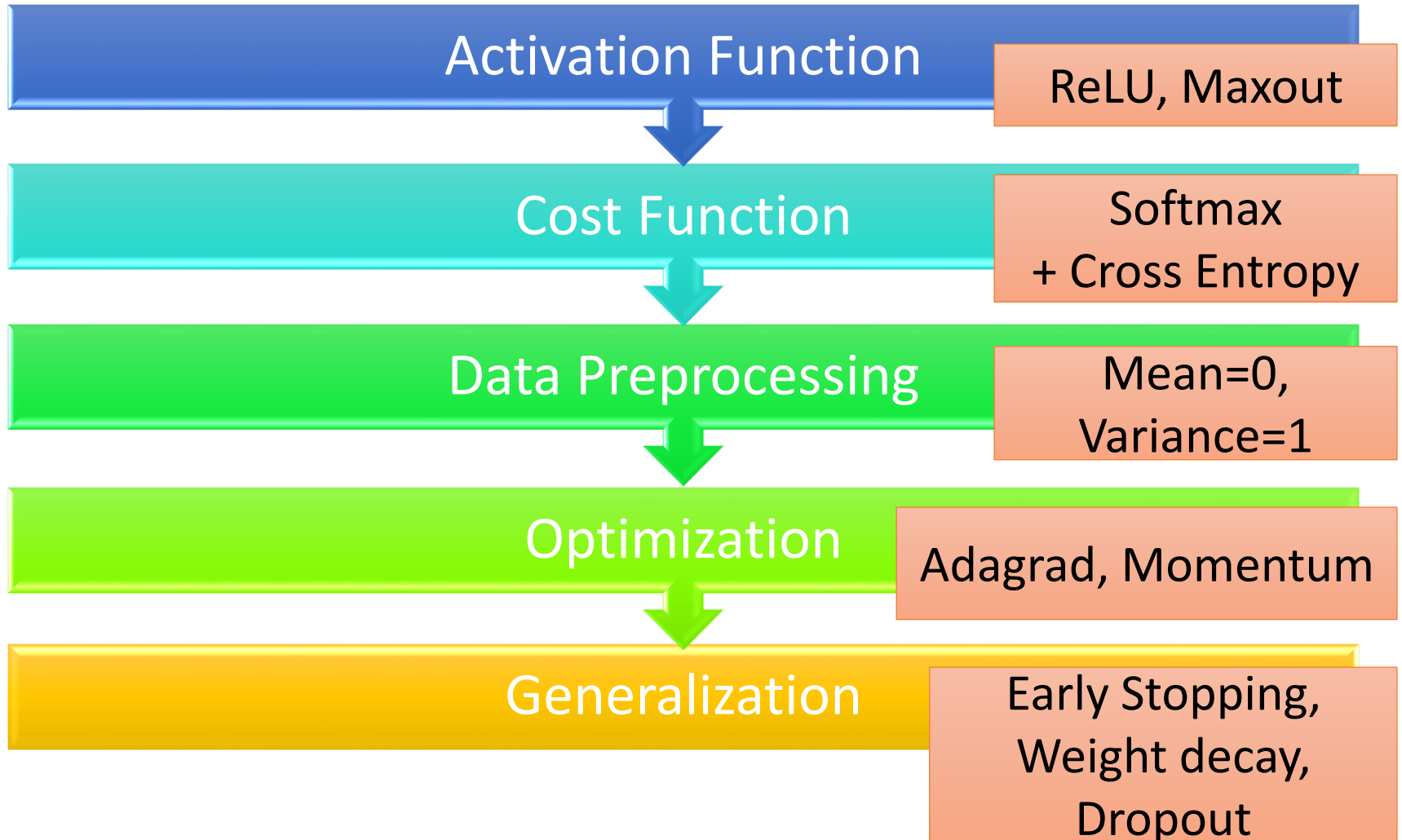Ref: http://arxiv.org/pdf/1302.4389.pdf

# Practical Suggestion for Dropout

- Larger network
  - If you know your task need **$n$** neurons, for dropout rate **$p$**, your network need **$n/(1-p)$** neurons.

- Longer training time

- Higher learning rate

- Larger momentum

# Concluding Remarks

**Activation Function** — ReLU, Maxout

**Cost Function** — Softmax + Cross Entropy

**Data Preprocessing** — Mean=0, Variance=1

**Optimization** — Adagrad, Momentum

**Generalization** — Early Stopping, Weight decay, Dropout

# Acknowledgement

- 感謝 李朋軒 同學糾正投影片上的錯誤
  - 很多地方 p 應該改為 1-p
- 感謝 Ryan Sun 來信指出投影片上的錯誤