# Tips

**NN Practical Tips**
Nov 10th, 2016

# Applied Deep Learning
YUN-NUNG (VIVIAN) CHEN   WWW.CSIE.NTU.EDU.TW/~YVCHEN/F105-ADL

National Taiwan University

Slide credit from Hung-Yi Lee & Richard Socher

1

# Outline

Data Preprocessing

Activation Function

Loss Function

Optimization
◦ Adagrad
◦ Momentum

Generalization
◦ Early Stopping
◦ Regularization
◦ Dropout

# Outline

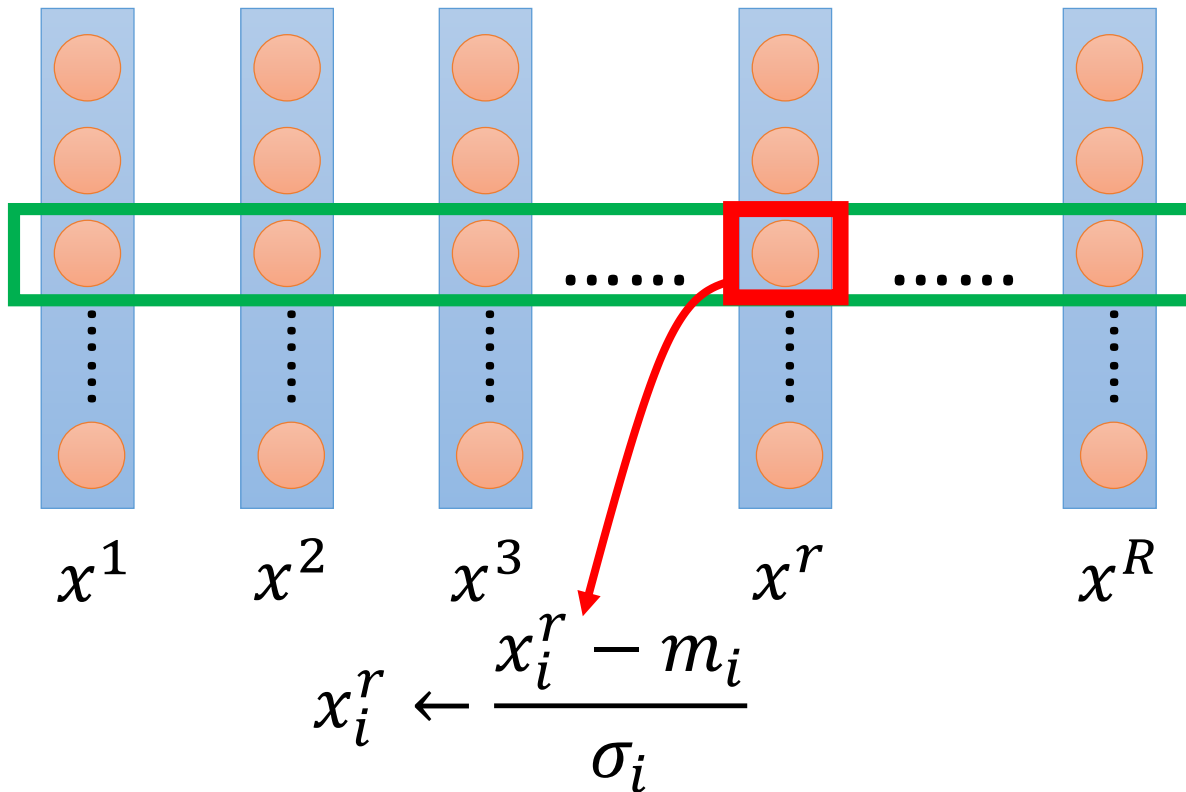**Data Preprocessing**

Activation Function

Loss Function

Optimization
◦ Adagrad
◦ Momentum

Generalization
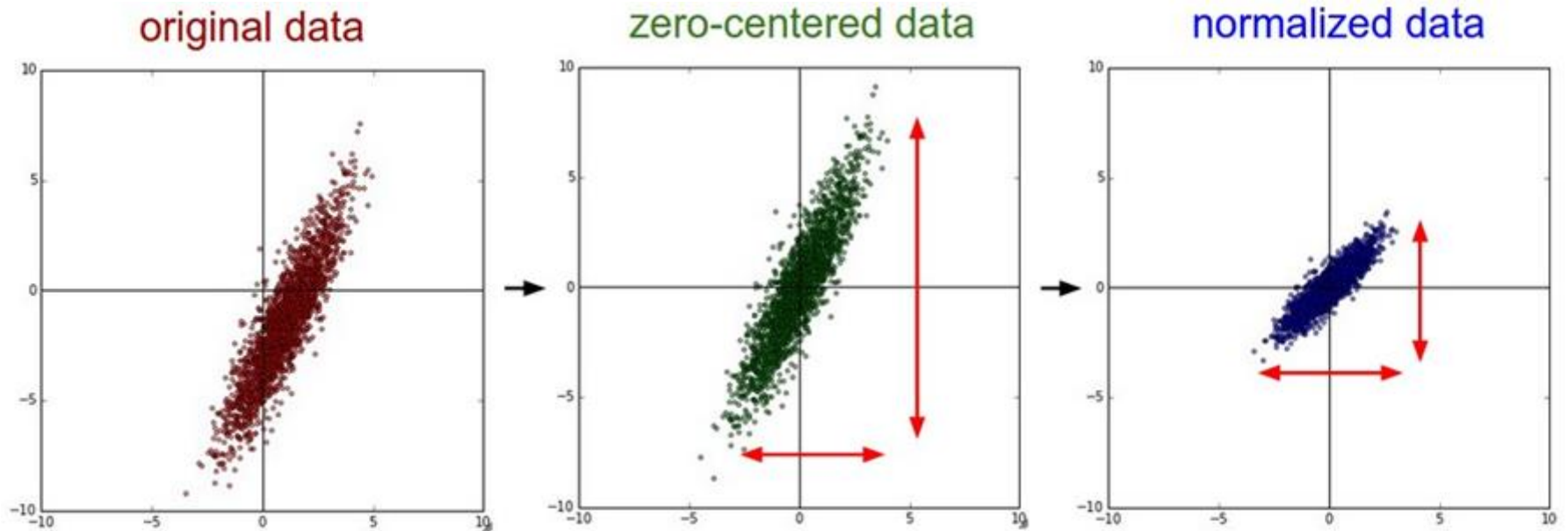◦ Early Stopping
◦ Regularization
◦ Dropout

# Input Normalization



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

# Input Normalization



original data      zero-centered data      normalized data

Normalizing training and testing data in the same way

# Outline

Data Preprocessing

**Activation Function**

Loss Function

Optimization
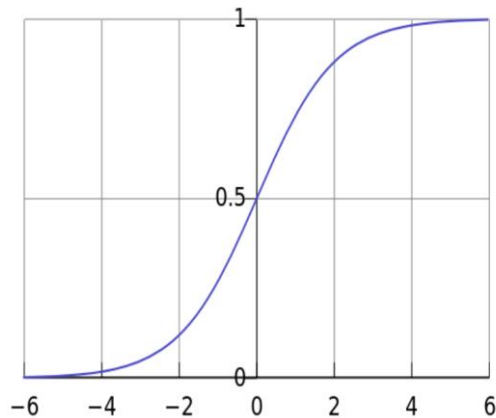 ◦ Adagrad
 ◦ Momentum

Generalization
 ◦ Early Stopping
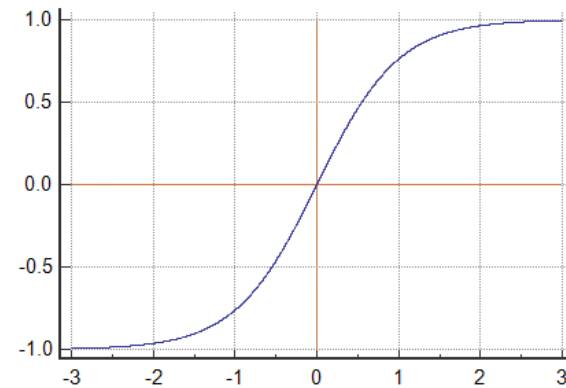 ◦ Regularization
 ◦ Dropout

# Activation Function

Sigmoid $f(x) = \dfrac{1}{1 + e^{-x}}$

Tanh $f(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$
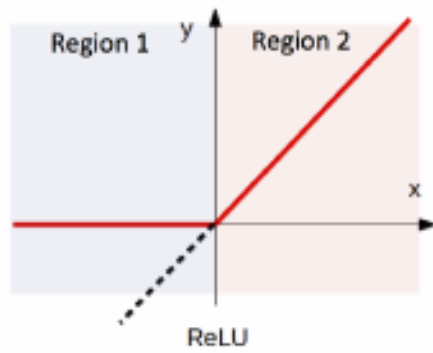
$$f'(x) = f(x)(1 - f(x))$$

$$f'(x) = 1 - f(x)^2$$

tanh is just a rescaled and shifted sigmoid, but better for many models
- Initialization: values close to 0
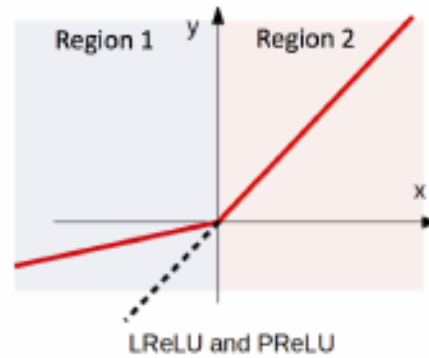- Convergence: faster in practice
- Nice derivative (similar to sigmoid)

# Variants
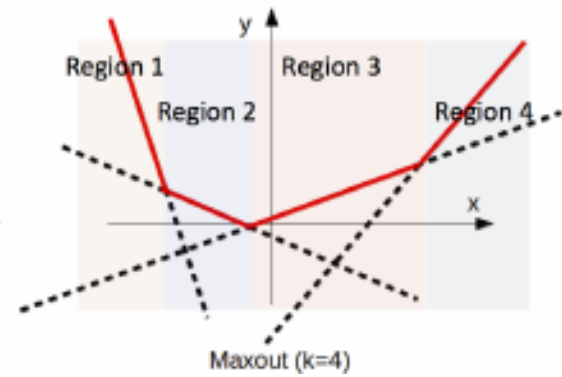
ReLU      LReLU & PReLU     Maxout

# Rectified Linear Unit (ReLU)

$\sigma(z)$

$$a$$

$a = z$

$a = 0$

$z$

$\sigma'(z)$

$$a$$

$1$

$0$

$z$

# Rectified Linear Unit (ReLU)

$\sigma(z)$

Reason

1. Fast to compute
2. Biological reason
3. Infinite sigmoid with different biases
4. Solution for vanishing gradient

$a$

$a = z$

$a = 0$

$z$

# Backpropagation

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \frac{\partial C(\theta)}{\partial z_i^l} \frac{\partial z_i^l}{\partial w_{ij}^l}$$

$$\delta_i^l$$

Error signal

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

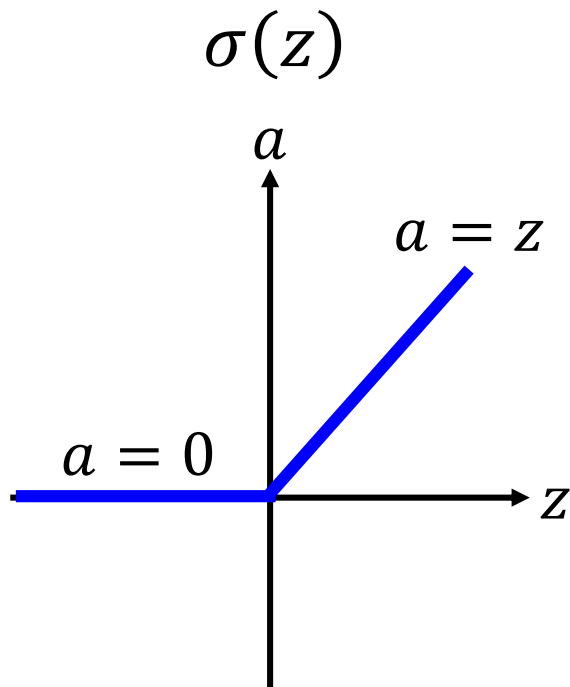Layer $l-1$    Layer $l$

$w_{ij}^l$

**_Backward Pass_**

$$\delta^L = \sigma'(z^L) \odot \nabla C(y)$$
$$\delta^{L-1} = \sigma'(z^{L-1}) \odot (W^L)^T \delta^L$$
$$\vdots$$
$$\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$$
$$\vdots$$

**_Forward Pass_**

$$z^1 = W^1 x + b^1$$
$$a^1 = \sigma(z^1)$$
$$\vdots$$
$$z^l = W^l a^{l-1} + b^l$$
$$a^l = \sigma(z^l)$$
$$\vdots$$

# Backpropagation

$$\frac{\partial C(\theta)}{\partial w_{ij}^l} = \boxed{\frac{\partial C(\theta)}{\partial z_i^l}}\frac{\partial z_i^l}{\partial w_{ij}^l}$$
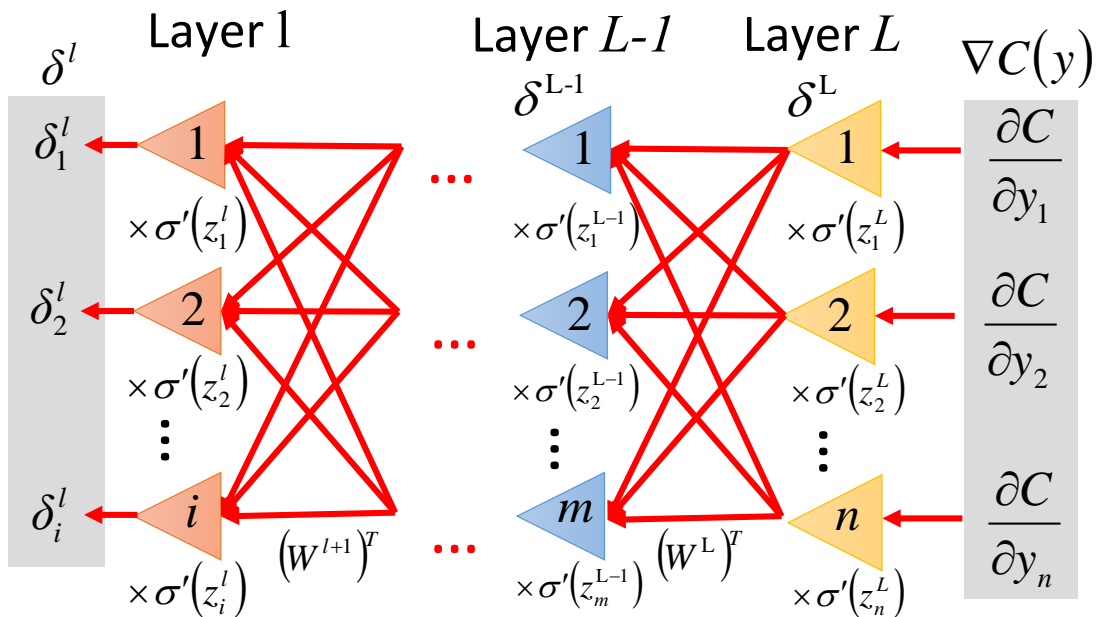
$$\boxed{\delta_i^l}$$ Error signal

**Backward Pass**

$$\delta^L = \sigma'(z^L) \odot \nabla C(y)$$
$$\delta^{L-1} = \sigma'(z^{L-1}) \odot (W^L)^T \delta^L$$
$$\vdots$$
$$\delta^l = \sigma'(z^l) \odot (W^{l+1})^T \delta^{l+1}$$
$$\vdots$$



$\delta^l$   Layer 1   Layer $L$-1   Layer $L$   $\nabla C(y)$

$\delta_1^l$   1   $\times \sigma'(z_1^l)$

$\delta_2^l$   2   $\times \sigma'(z_2^l)$

$\delta_i^l$   $i$   $\times \sigma'(z_i^l)$   $(W^{l+1})^T$

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

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

$m$   $\times \sigma'(z_m^{L-1})$   $(W^L)^T$

$\delta^L$   1   $\times \sigma'(z_1^L)$

2   $\times \sigma'(z_2^L)$

$n$   $\times \sigma'(z_n^L)$

$\frac{\partial C}{\partial y_1}$

$\frac{\partial C}{\partial y_2}$

$\frac{\partial C}{\partial y_n}$

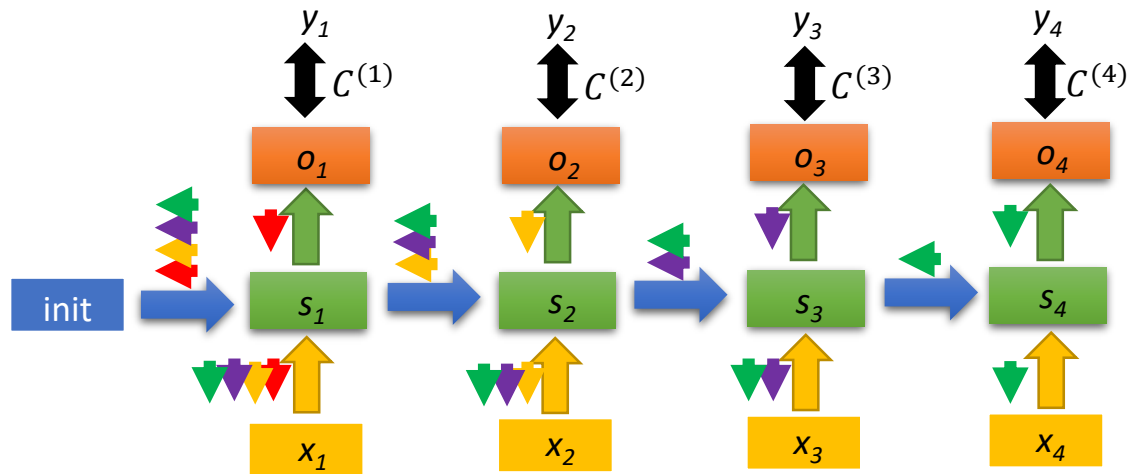# Sigmoid Issue



Sigmoid Function

Derivative of Sigmoid Function

Derivative of the sigmoid function is always smaller than 1

# Vanishing Gradient Problem



$$\delta^l = \boxed{\sigma'(z^l) \odot (W^{l+1})^T} \delta^{l+1}$$
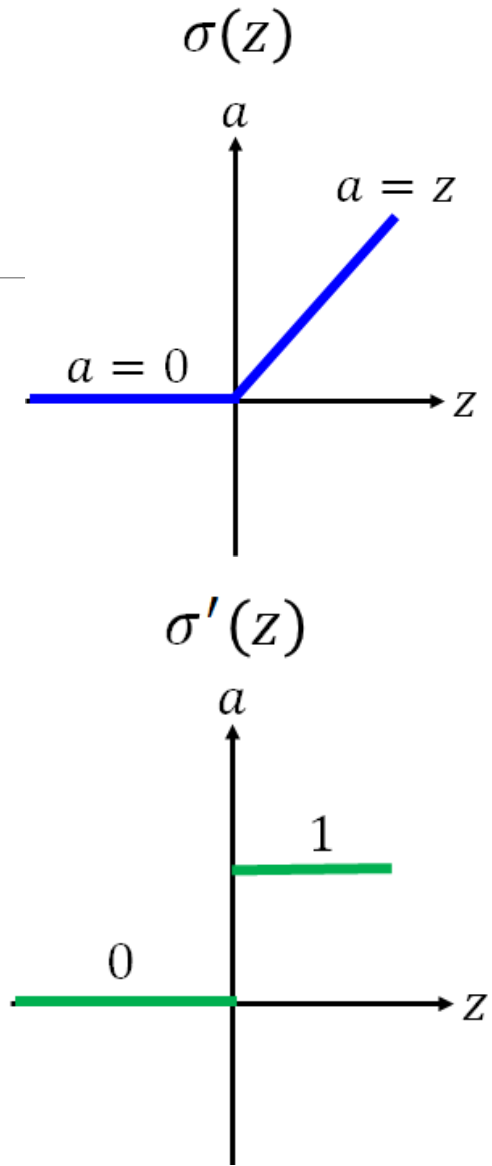
The error signal is getting smaller and smaller due to $\sigma'(z) < 1$
$\rightarrow$ **vanishing gradient**

# Vanishing Gradient Problem



Input      Layer 1      Layer 2      Layer L      Output

vector **x**

$x_1$
$x_2$
$x_3$

$y_1$
$y_2$

vector **y**

Learn very slowly

Still random

Learn faster

Already converge

The weights are converged based on random!?

# ReLU

# ReLU

$$\sigma(z)$$

$a$

$a = z$

$a = 0$

$z$

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

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

$\delta_1^l$    1    1    1    $\dfrac{\partial C}{\partial y_1}$

0    1    1

$\delta_2^l$    2    2    2    $\dfrac{\partial C}{\partial y_2}$

1    0    1

$\delta_i^l$    $i$    $m$    $n$    $\dfrac{\partial C}{\partial y_n}$

1    $(W^{l+1})^T$    $(W^L)^T$    0    0

$$\sigma'(z)$$

$a$

1

0    $z$

# ReLU



A thinner network without any attenuation

# ReLU – Forward Pass



Input

Output

$x_1$

$x_2$

0

0

0

0

$y_1$

$y_2$

# ReLU – Backward Pass



Input

Output

$x_1$

$x_2$

good influence to the output

$y_1$

$y_2$

# Variant ReLU

*Leaky ReLU*

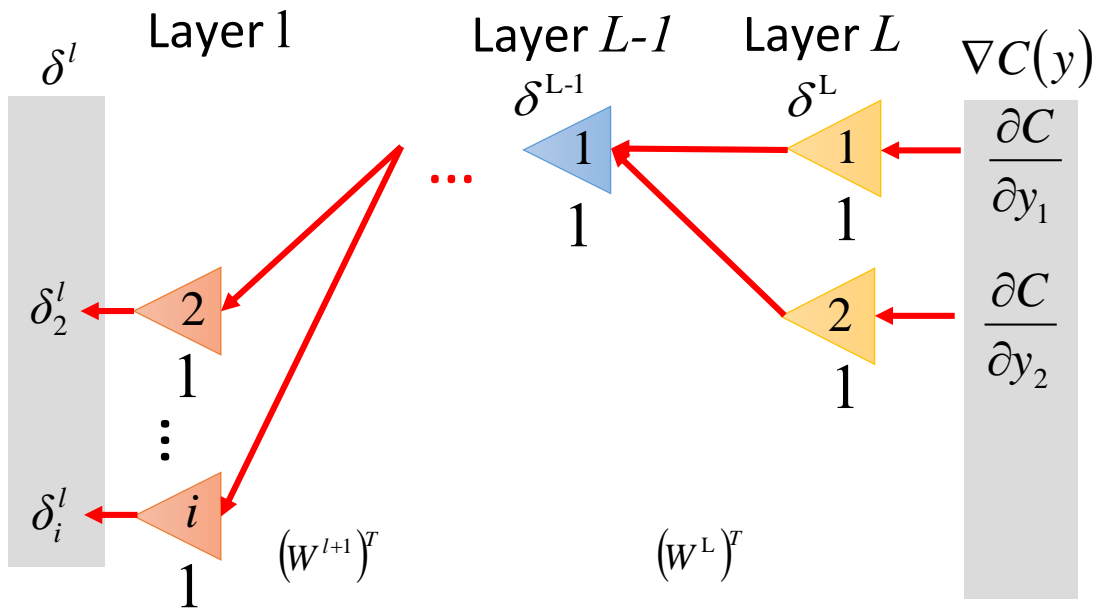$a$

$a = z$

$z$

$a = 0.01z$

*Parametric ReLU*

$a$

$a = z$

$z$

$a = \alpha z$

α is also learned by gradient descent

# Maxout



Input

$x_1$

$x_2$

$x$

$z_1^1$

$z_2^1$

$z_3^1$

$z_4^1$

$z^1$

$max\{z_1^1, z_2^1\}$

Max

Max

$a_1^1$

$a_2^1$

$a^1$

$z_1^2$

$z_2^2$

$z_3^2$

$z_4^2$

$z^2$

Max

Max

$a_1^2$

$a_2^2$

$a^2$

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

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

# Maxout – ReLU is a special case

Input

$x$

$1$

$w$

$b$

$z$

ReLU $\rightarrow a$

$a$

$z = wx + b$

$x$

Input

$x$

$1$

$w$

$b$

$0$

$0$

$+$ $\rightarrow z_1$

$+$ $\rightarrow z_2$

Max $\rightarrow a$

$max\{z_1, z_2\}$

$a$

$z_1 = wx + b$

$z_2 = 0$

$x$

# Maxout – ReLU is a special case

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

$z = wx + b$

Input
$w$
$b$
$x$   $w'$
$b'$
1
+ → $z_1$
+ → $z_2$
Max → $a$
$max\{z_1, z_2\}$

Learnable activation function

$z_1 = wx + b$

$z_2 = w'x + b'$

# Maxout - Training

Given training data $x$, we decide $z$ for maxout



$max\{z_1^1, z_2^1\}$

# Maxout - Training

Given training data $x$, we decide $z$ for maxout



Training this thin and linear network

# Outline

Data Preprocessing

Activation Function

**Loss Function**

Optimization
◦ Adagrad
◦ Momentum

Generalization
◦ Early Stopping
◦ Regularization
◦ Dropout

# Loss Function – Square Error



Layer L-1

Layer L
(Output layer)

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

$$= \frac{1}{2}\sum_{n}\left(y_n - \hat{y}_n\right)^2$$

# 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

**_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** $y_3 = e^{z_3^L} \Big/ \sum_{j=1}^{3} e^{z_j^L}$

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

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

Training labels indicate positive and negative samples in stead of the actual values

# Outline

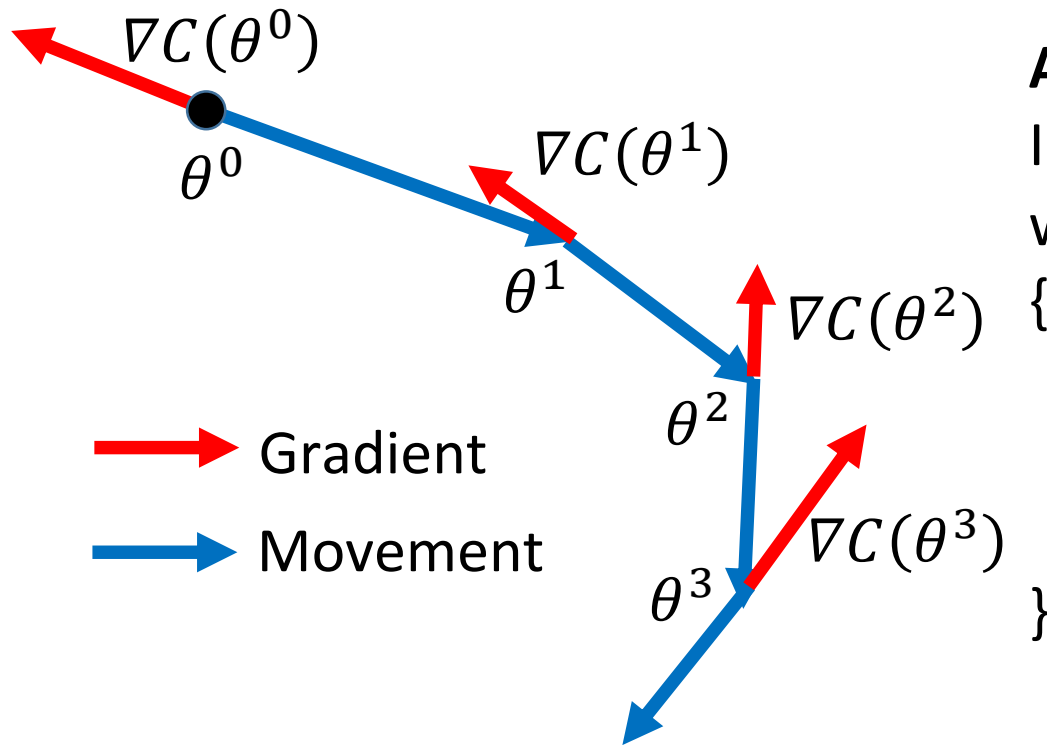Data Preprocessing

Activation Function

Loss Function

**Optimization**
- Adagrad
- Momentum

Generalization
- Early Stopping
- Regularization
- Dropout

# Gradient Descent for Optimization



$\nabla C(\theta^0)$

$\theta^0$

$\nabla C(\theta^1)$

$\theta^1$

$\nabla C(\theta^2)$

$\theta^2$

$\nabla C(\theta^3)$

$\theta^3$

→ Gradient

→ Movement

**Algorithm**

Initialization: start at $\theta^0$

while($\theta^{(i+1)} \neq \theta^i$)

{

    compute gradient at $\theta^i$

    update parameters

      $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta C(\theta^i)$

}

1) How to determine the learning rates → learning rate
2) How to avoid stucking at local minima or saddle points → learning direction

# Outline

Data Preprocessing

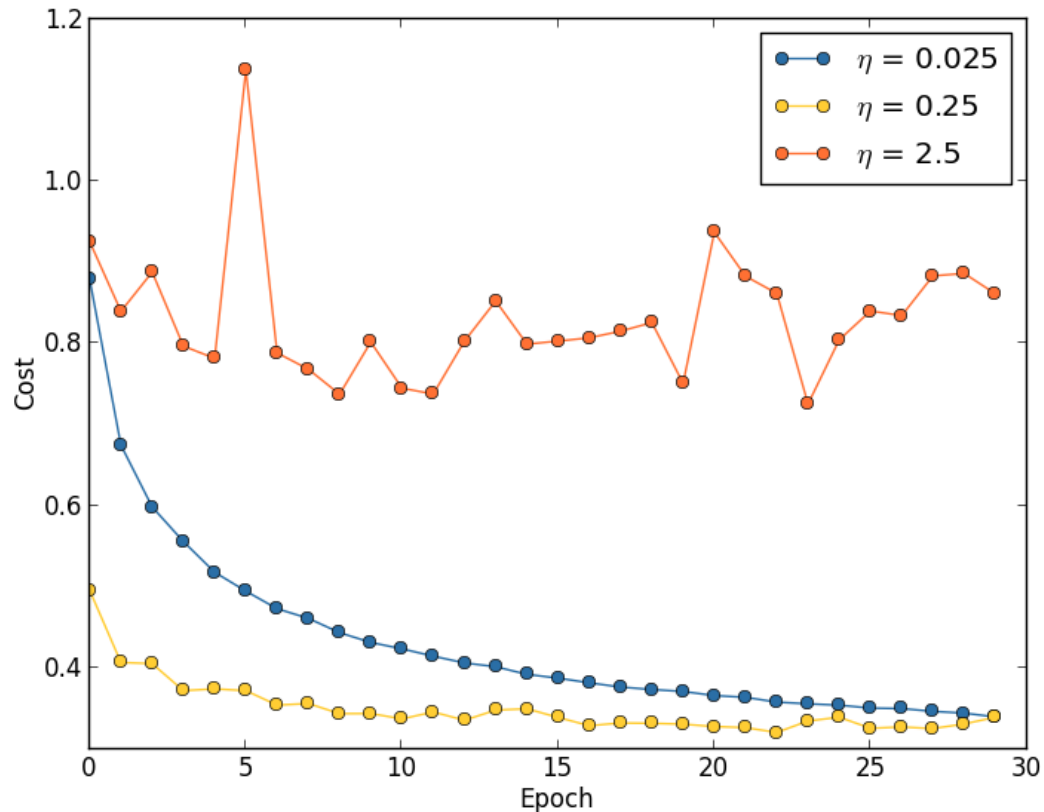Activation Function

Loss Function

Optimization
◦ **Adagrad**
◦ Momentum

Generalization
◦ Early Stopping
◦ Regularization
◦ Dropout

# Learning Rate



The proper learning rate is important to find the optimal point

# Learning Rate

Idea: reduce the learning rate every few epochs

- At the beginning, we are far from the destination, so we use a larger learning rate
- After several epochs, we are close to the destination, so we reduce the learning rate

## Manually set learning rate

1) Reduce by 0.5 when validation error stops improving
2) 1/t decay: $\eta^t = \eta / \sqrt{t+1}$ due to theoretical convergence guarantees

Learning rate cannot be one-size-fits-all
→ different parameters have different learning rates

# Adagrad

Idea: adaptive learning rates for each parameter

Approach: divide the learning rate of each parameter by the ***root mean square of its previous derivatives***

**_Vanilla Gradient Descent_**

$w$ is a parameter

1/t decay

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

**_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

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

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1 \qquad \sigma^1 = \sqrt{\frac{1}{2}[(g^0)^2 + (g^1)^2]}$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2 \qquad \sigma^2 = \sqrt{\frac{1}{3}[(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

$$\vdots$$

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

# Adagrad

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

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

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad \text{1/t decay}$$

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

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

Learning rate is adapting differently for each parameter and rare parameters get larger updates than frequently occurring parameters

Duchi et al., "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization," in COLT, 2010. [link]

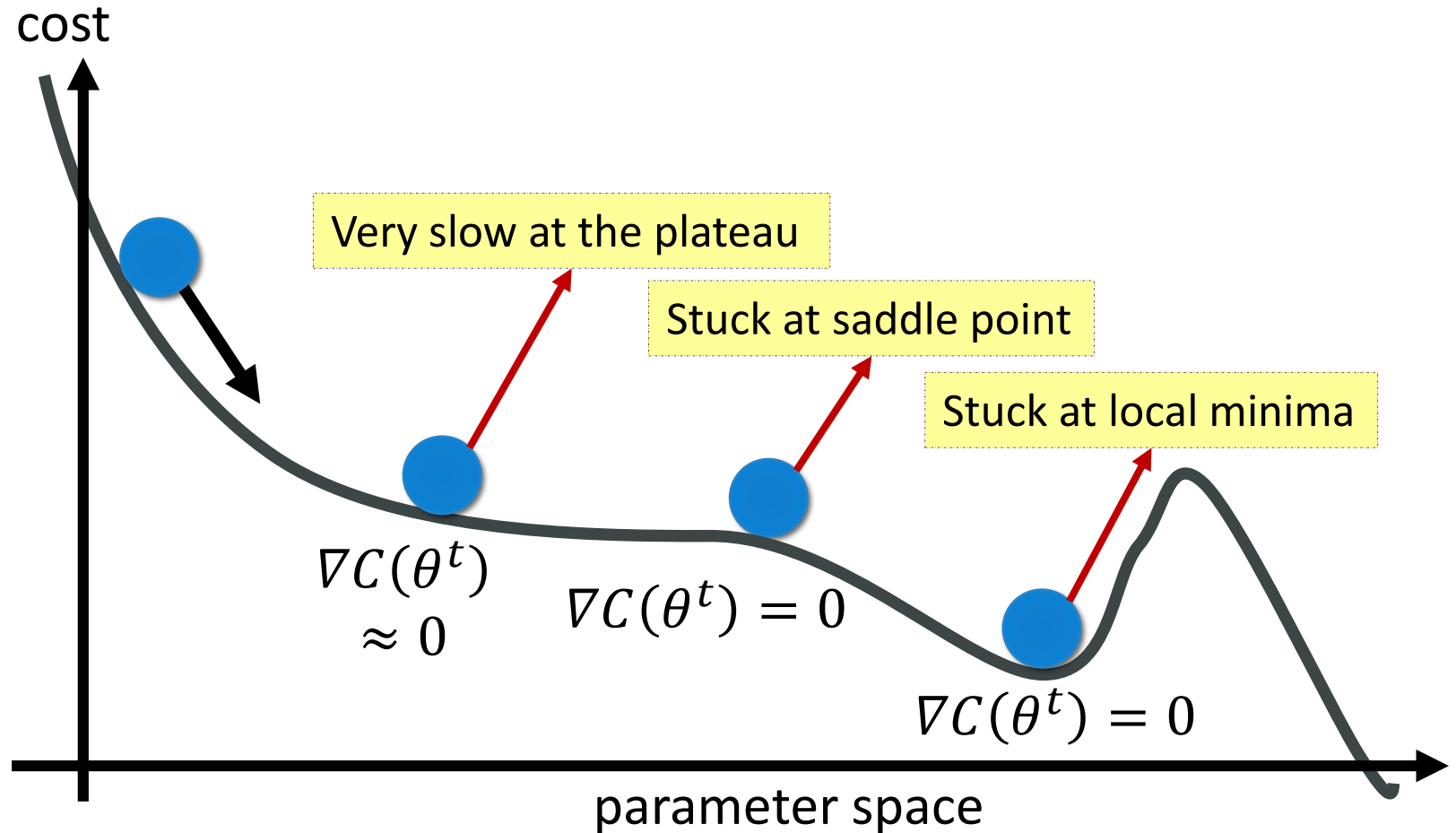# Outline

Data Preprocessing

Activation Function

Loss Function

**Optimization**
- Adagrad
- **Momentum**

Generalization
- Early Stopping
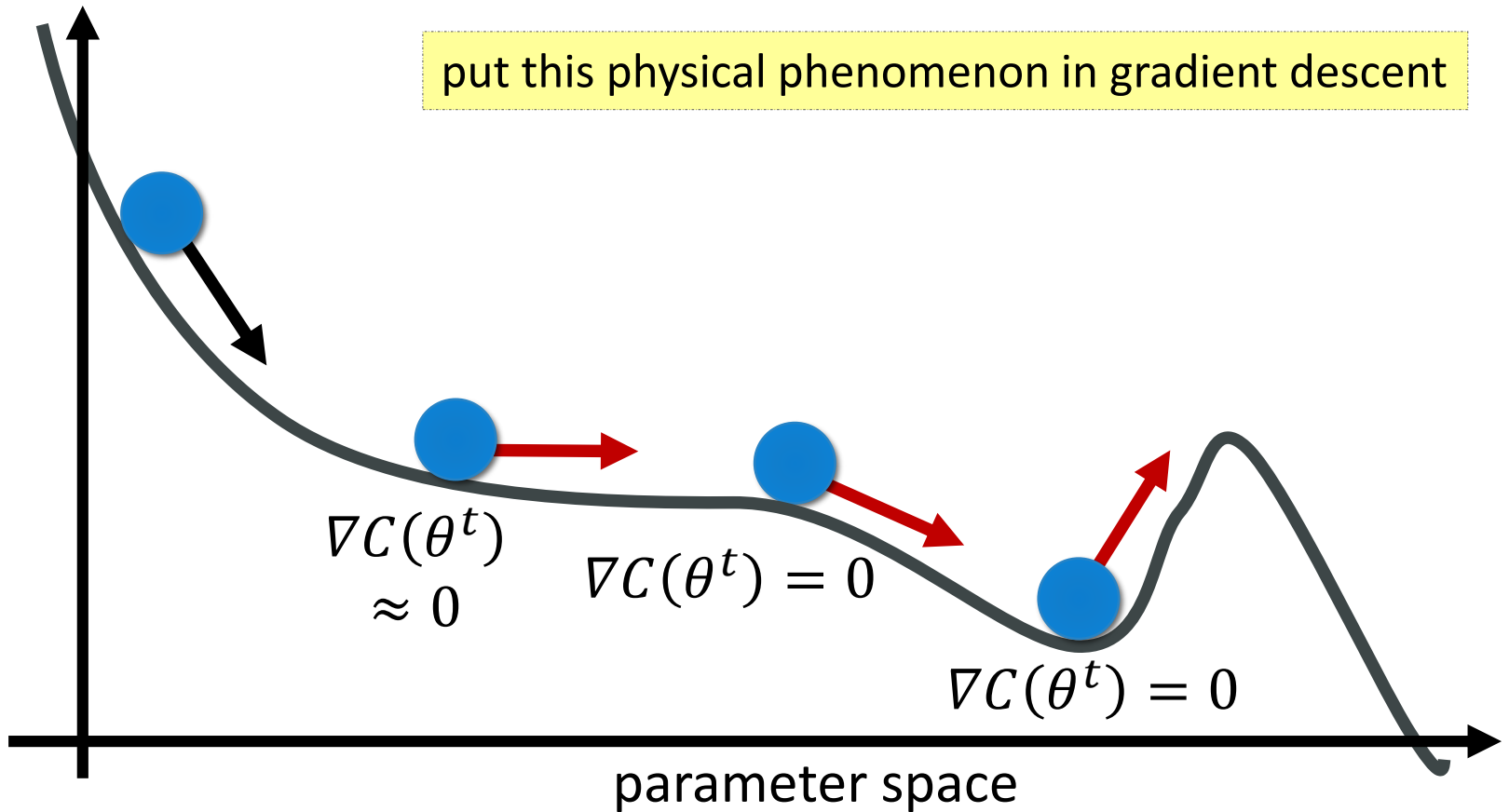- Regularization
- Dropout

# Gradient Descent Stuck Issue

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

# Momentum

cost

put this physical phenomenon in gradient descent

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

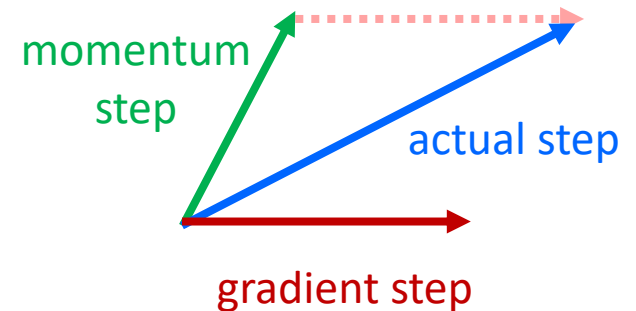$$\nabla C(\theta^t) = 0$$

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

parameter space

# Momentum

Parameters build up velocity in direction of consistent gradient

$$v^{i+1} = \lambda v^i - \eta \nabla C_\theta(\theta^i)$$
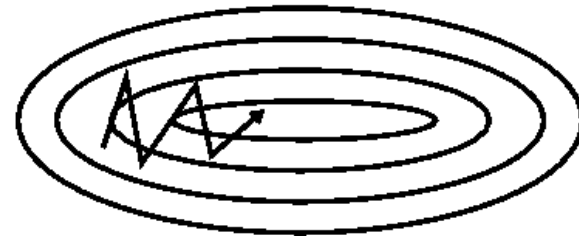$$\theta^{i+1} = \theta^i + v^{i+1}$$

momentum step

actual step

gradient step

E.g. convex function optimization dynamics

without momentum

with momentum

# Momentum

$\nabla C(\theta^0)$

$\nabla C(\theta^1)$

$\theta^0$

$\theta^1$

$\nabla C(\theta^2)$

$\theta^2$

$\theta^3$

$\nabla C(\theta^3)$

→ Gradient

→ Movement

⋯⋯ Movement of last step

**Algorithm**
Initialization: model parameters start at $\theta^0$, movement $v^0 = 0$
while($\theta^{(i+1)} \neq \theta^i$)
{

    compute gradient at $\theta^i$
    update parameters
$$v^{i+1} = \lambda v^i - \eta \nabla C_\theta(\theta^i)$$
$$\theta^{i+1} = \theta^i + v^{i+1}$$

}

Movement is not only based on gradient, but also previous movement

# Momentum

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

$$v^0 = 0$$

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

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

**Algorithm**

Initialization: model parameters start at $\theta^0$, movement $v^0 = 0$

while($\theta^{(i+1)} \neq \theta^i$)

{

    compute gradient at $\theta^i$

    update parameters

$$v^{i+1} = \lambda v^i - \eta \nabla C_\theta(\theta^i)$$

$$\theta^{i+1} = \theta^i + v^{i+1}$$

}

Movement is not only based on gradient, but also previous movement

# Outline

# Generalization

Practical tricks

1) find the right network structure and implement and optimize it properly

2) make the model overfit on training data

3) prevent overfitting
   ◦ Reduce the model size by lowering the number of units and layers/hyperparameters
   ◦ Early stopping
   ◦ Standard L1 or L2 regularization
   ◦ Sparsity constraint

# Outline

# Early Stopping



Check performance on validation set to prevent training too many iterations

# Outline

# Regularization

Idea: the parameters closer to zero are preferred

Training data:

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

Testing data:

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

$$x' = x + \varepsilon$$

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

$$z = w \cdot x$$

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

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

# Regularization

Idea: optimize a new cost function to find a set of weight that 1) minimizes original cost and 2) is close to zero

$$C'(\theta) = \underline{C(\theta)} + \lambda \frac{1}{2} \underline{\|\theta\|^2} \longrightarrow \text{regularization term}$$

original cost

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

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

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

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

→ not consider biases

# Regularization

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

Idea: optimize a new cost function to find a set of weight that 1) minimizes original cost and 2) is close to zero

$$C'(\theta) = C(\theta) + \lambda \frac{1}{2}\|\theta\|^2$$

Gradient: $\dfrac{\partial C'}{\partial w} = \dfrac{\partial C}{\partial w} + \lambda w$

Update:

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

$$= (1 - \eta\lambda)w^t - \eta \frac{\partial C}{\partial w^t}$$

Smaller and smaller

# Outline

# Dropout

For each iteration of training,
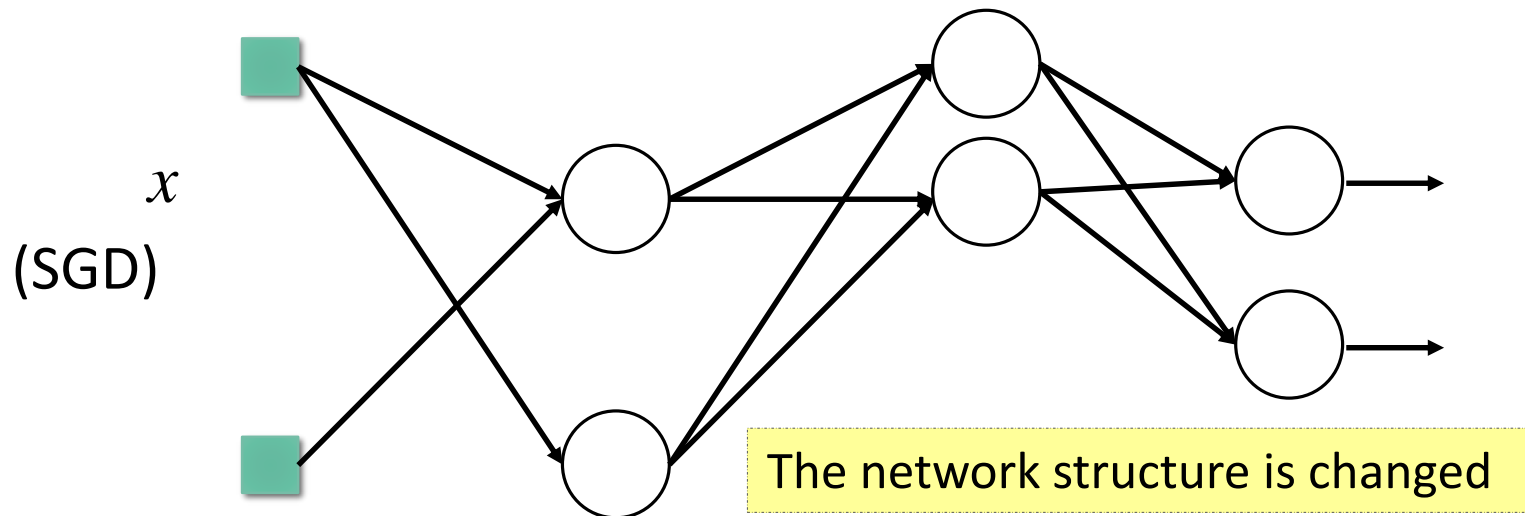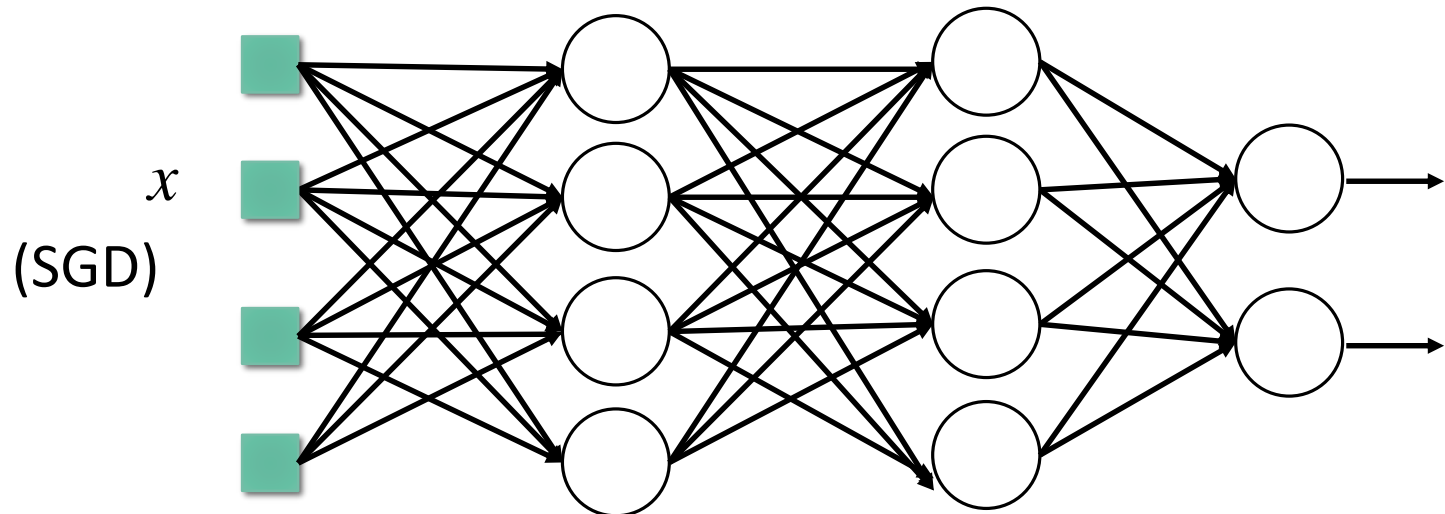◦ each neuron has $p$% to dropout

$x$

(SGD)

# Dropout

For each iteration of training,
  ◦ each neuron has $p$% to dropout
  ◦ training using the new network

**Training:**  $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta C(\theta^i)$

$x$

(SGD)

The network structure is changed

# Dropout

For each iteration of training,
◦ each neuron has $p$% to dropout
◦ training using the new network

**Training:** $\theta^{i+1} \leftarrow \theta^i - \eta \nabla_\theta C(\theta^i)$

For testing,
◦ no dropout
◦ if the dropout rate at training is p%, all the weights times (1-p)%
◦ e.g. the dropout rate is 50%, $w_{ij}^l = 1$ from training → $w_{ij}^l = 0.5$ for testing

$x$

(SGD)
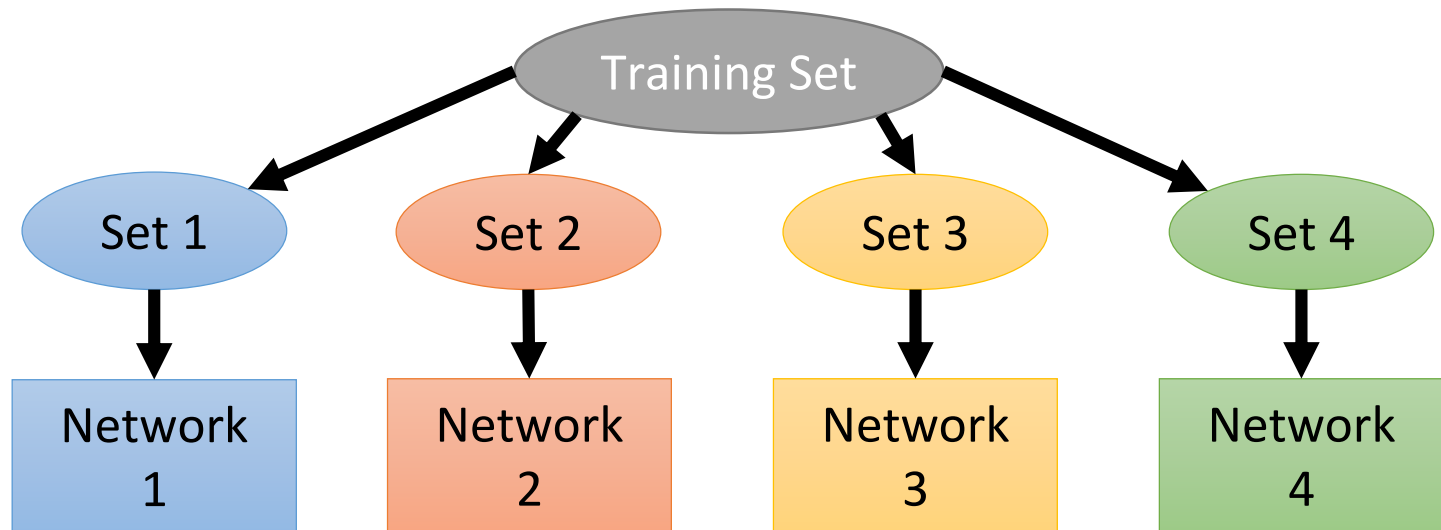
# Dropout – Intuitive Reason

**Training: dropout**
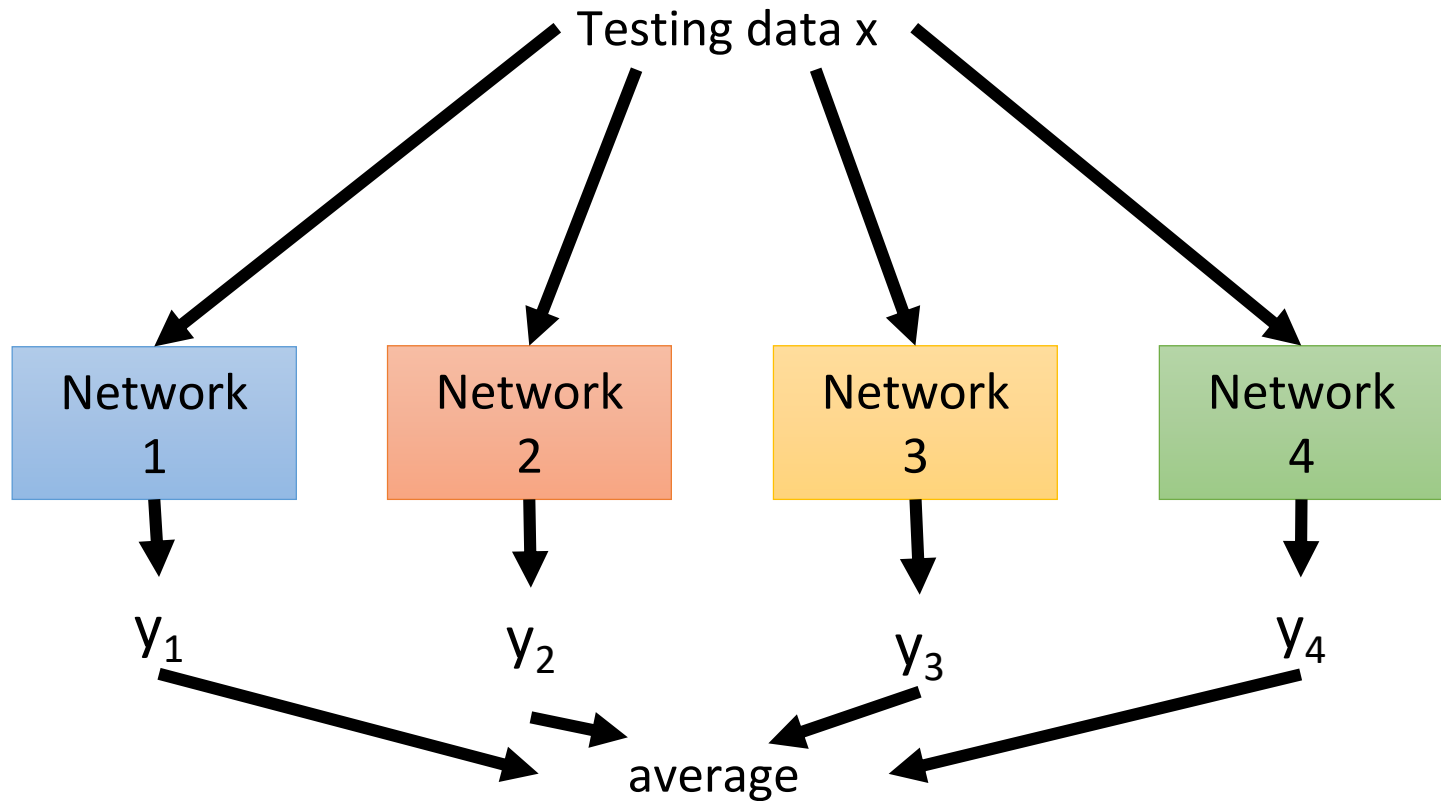
**Testing: no dropout**

# Dropout

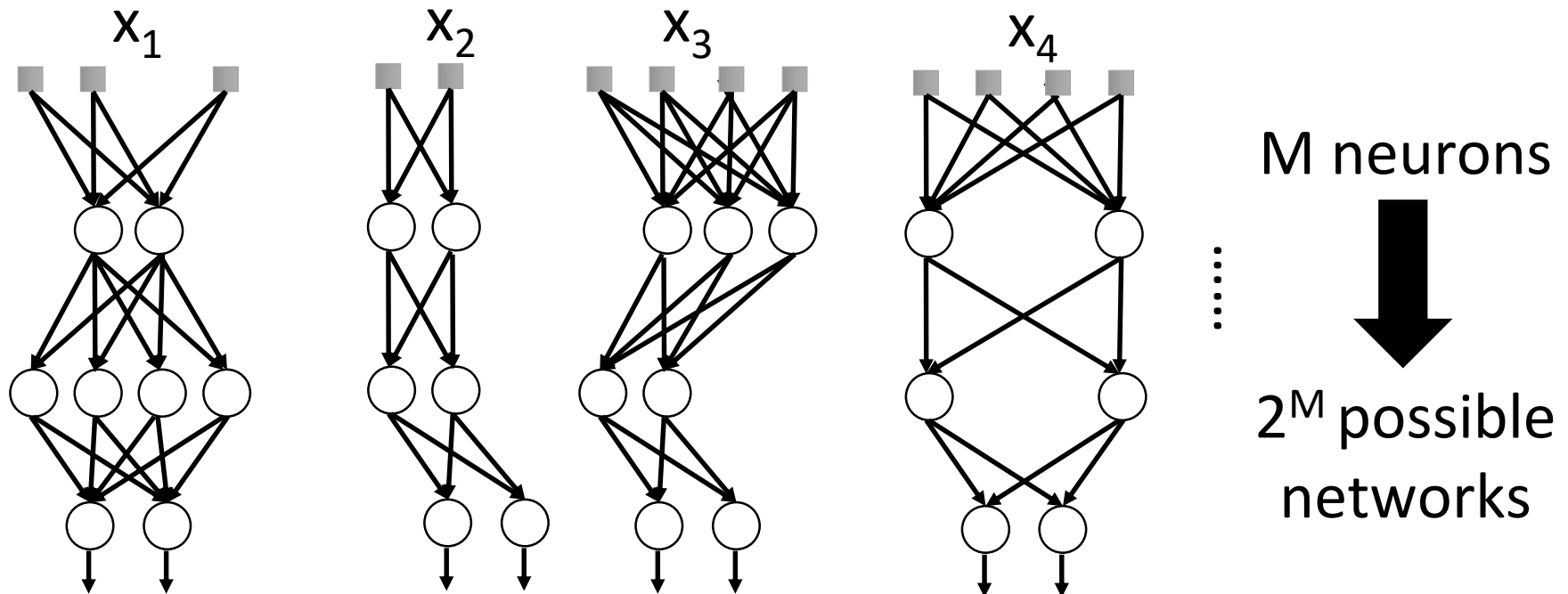Train a bunch of networks with different structures
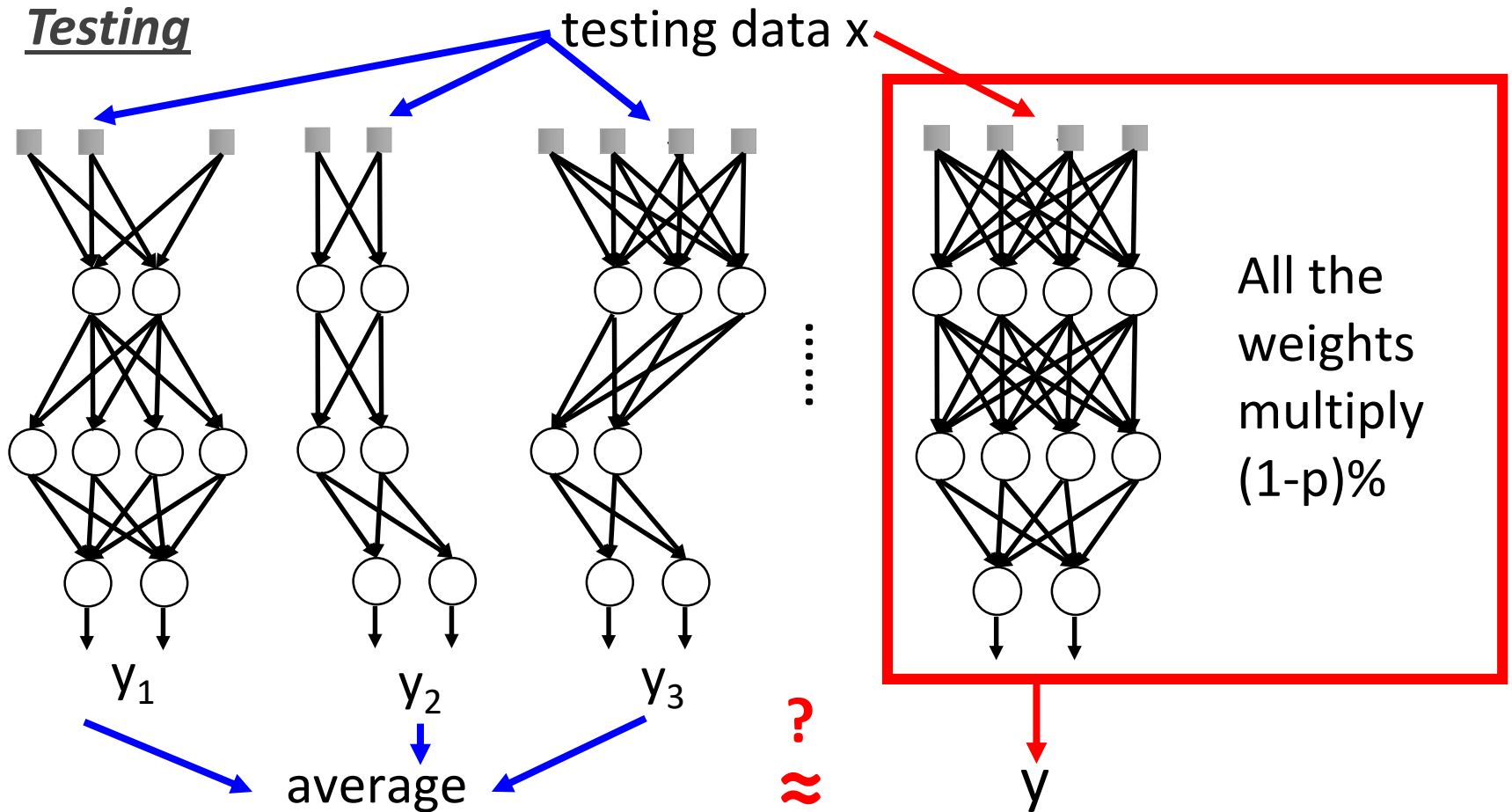
# Dropout – Ensemble

**_Ensemble_**

# Dropout – Ensemble

***Training***

◦ Using one data to train one network

◦ Some parameters in the network are shared



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

M neurons

$2^M$ possible networks

# Dropout – Ensemble

**_Testing_**

testing data x

$y_1$     $y_2$     $y_3$

average

**?**
**≈**

All the
weights
multiply
(1-p)%

y

# Dropout Tips

Larger network

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

Longer training time

Higher learning rate

Larger momentum

# Concluding Remarks

Data Preprocessing: Input Normalization

Activation Function: ReLU, Maxout

Loss Function: Softmax

Optimization
◦ Adagrad: Learning Rate Adaptation
◦ Momentum: Learning Direction Adaptation

Generalization
◦ Early Stopping: avoid too many iterations from overfitting
◦ Regularization: minimize the effect of noise
◦ Dropout: leverage the benefit of ensemble