

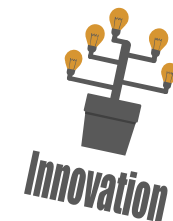
# 训练神经网络的技巧

## Tricks for Training Artificial Neural Networks

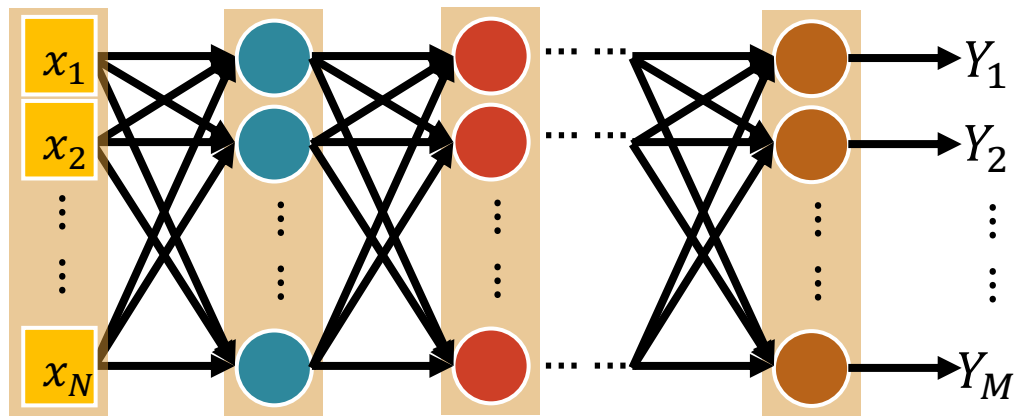
郝 奇  
南京大学 天文与空间科学学院

**Application of  
Machine Learning  
in Astronomy**

**机器学习在天文中的应用**



## Contents



**01** Recipe of ANN

---

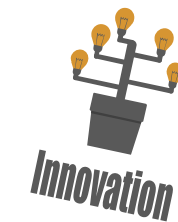
**02** Tricks from training process

---

**03** Tricks from testing process

---

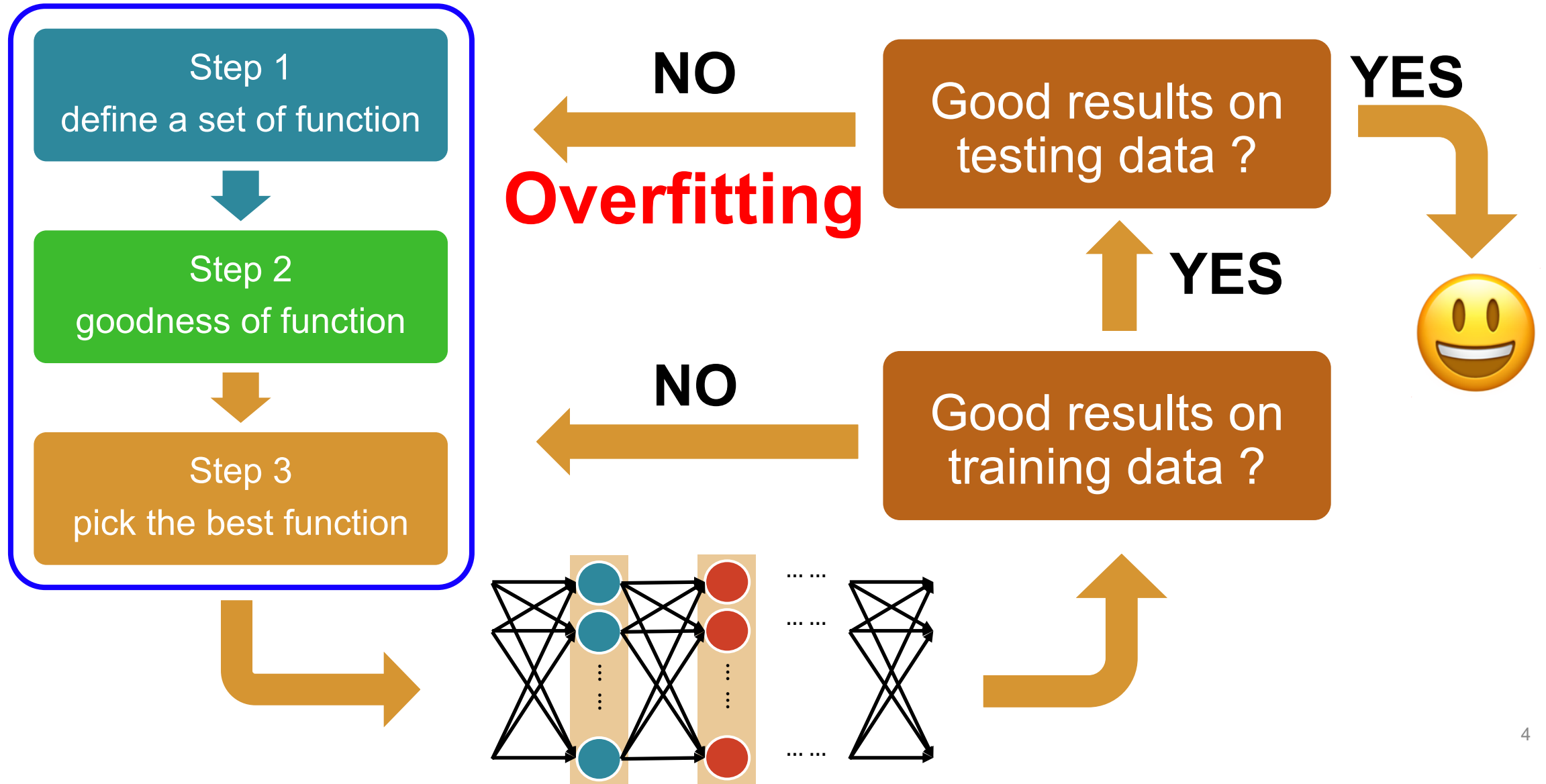
# Recipe of Artificial Neural Networks



01

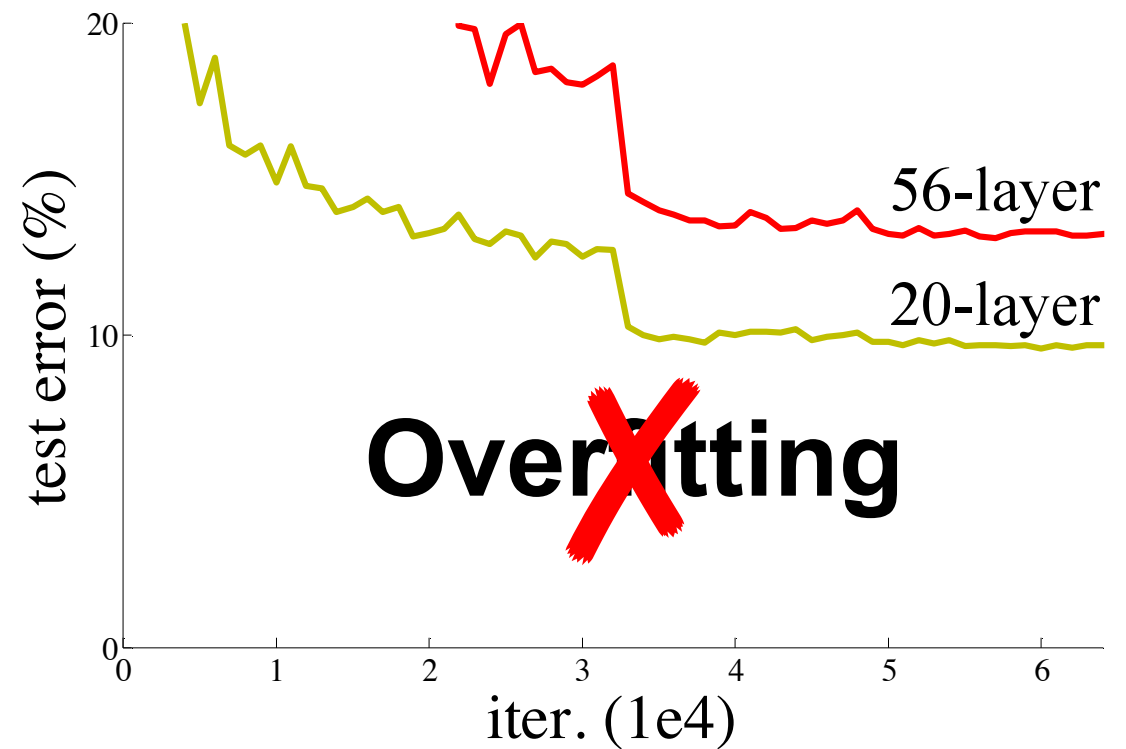
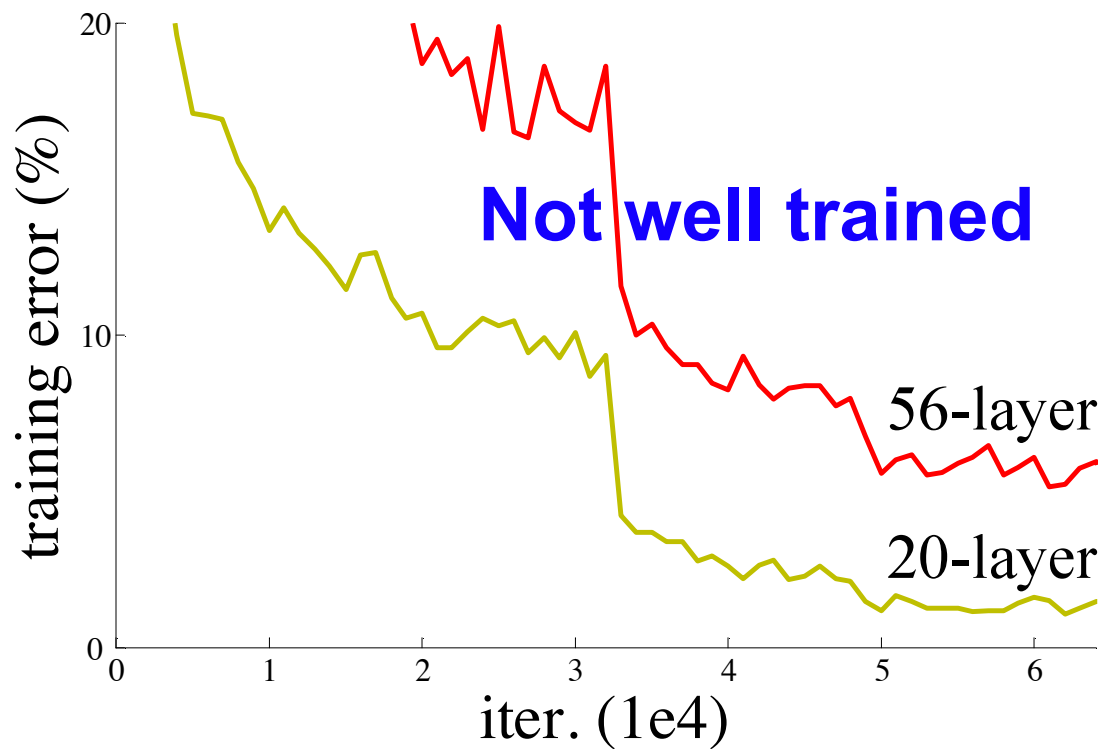


# 1. Recipe of Artificial Neural Networks



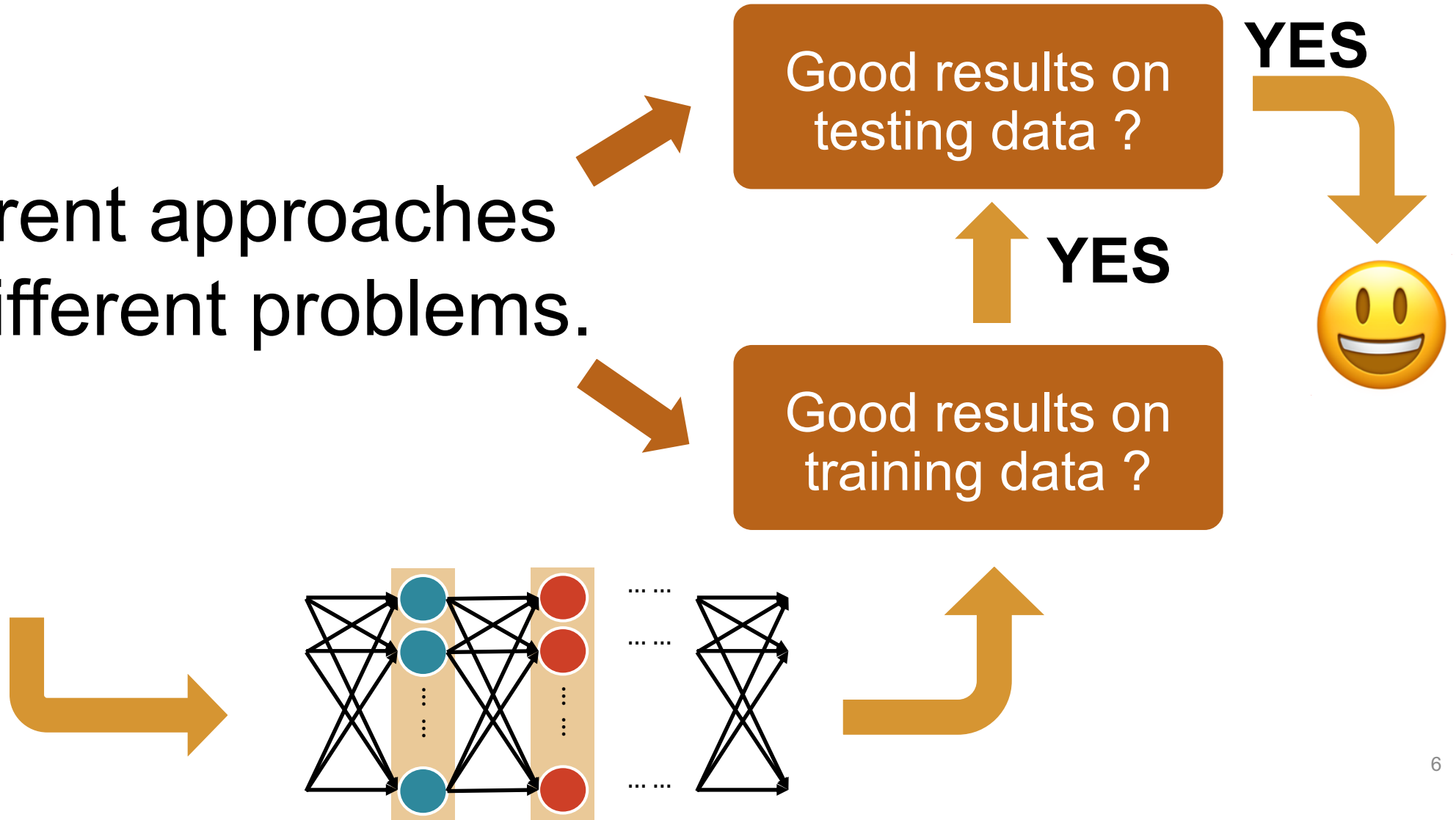
# 1. Recipe of Artificial Neural Networks

## Do not always blame overfitting



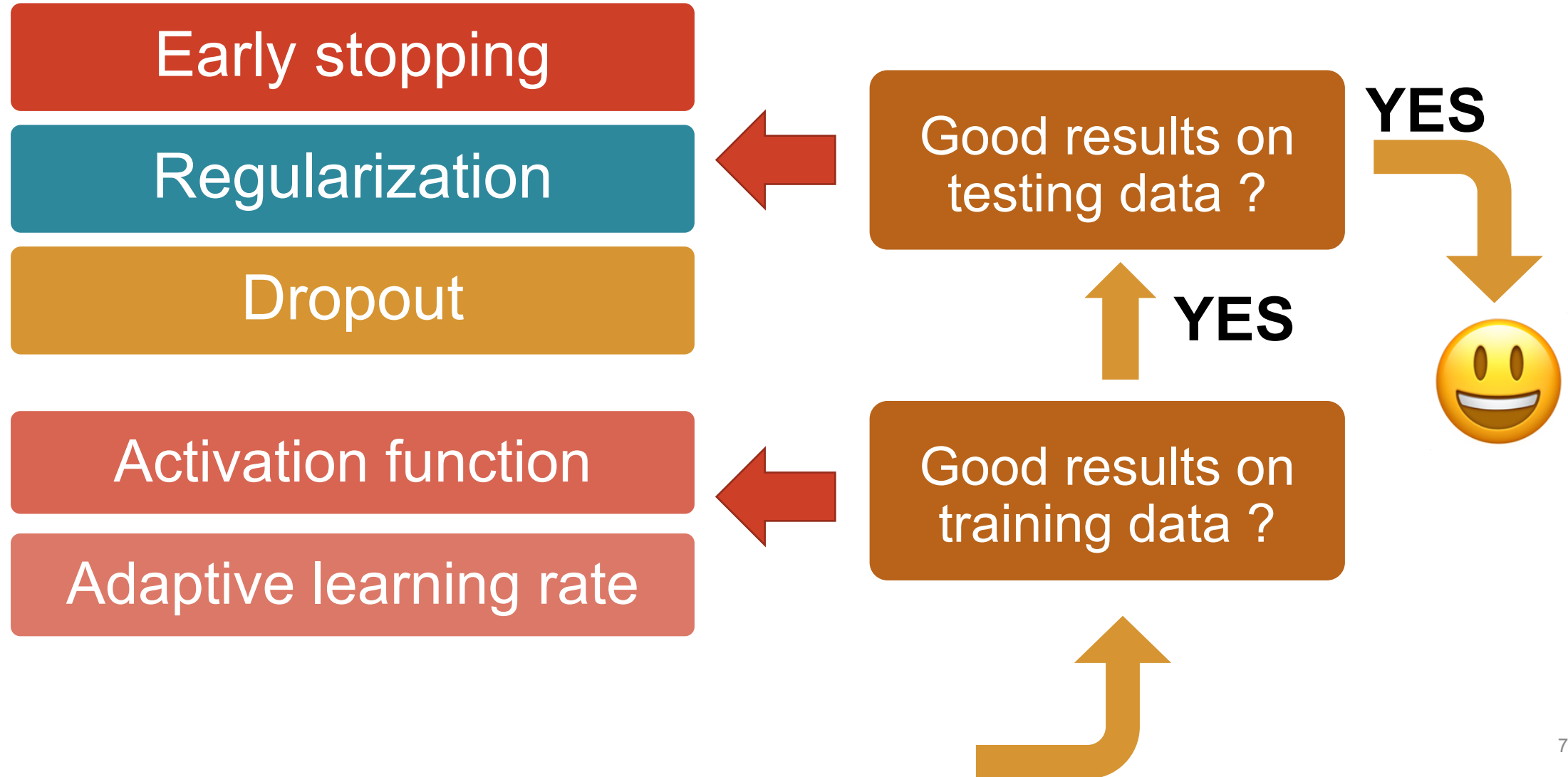
# 1. Recipe of Artificial Neural Networks

Different approaches for different problems.



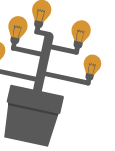
# 1. Recipe of Artificial Neural Networks

---



# Tricks from training process

# 02

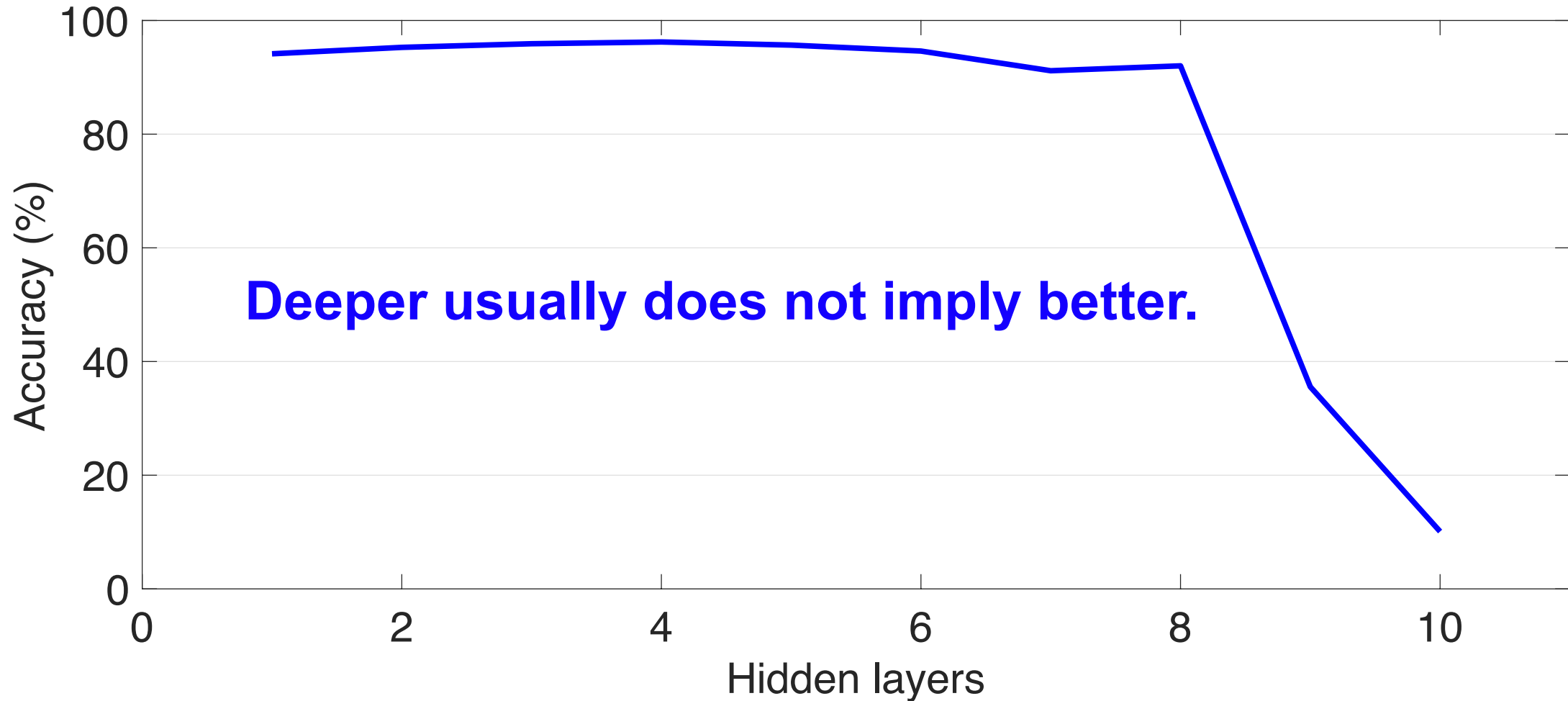




## 2. Tricks from training process

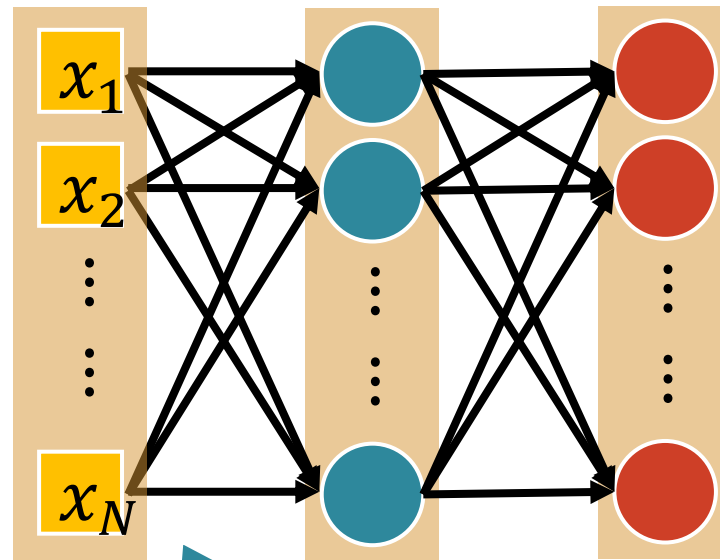
---

### 2.1 Activation function



## 2. Tricks from training process

### 2.1 Activation function

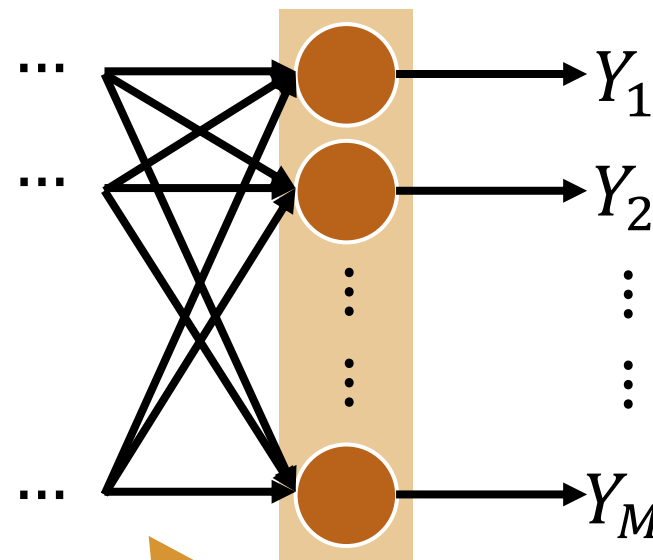


Smaller gradients

Learn very slow

Almost random

### Vanishing gradient problem



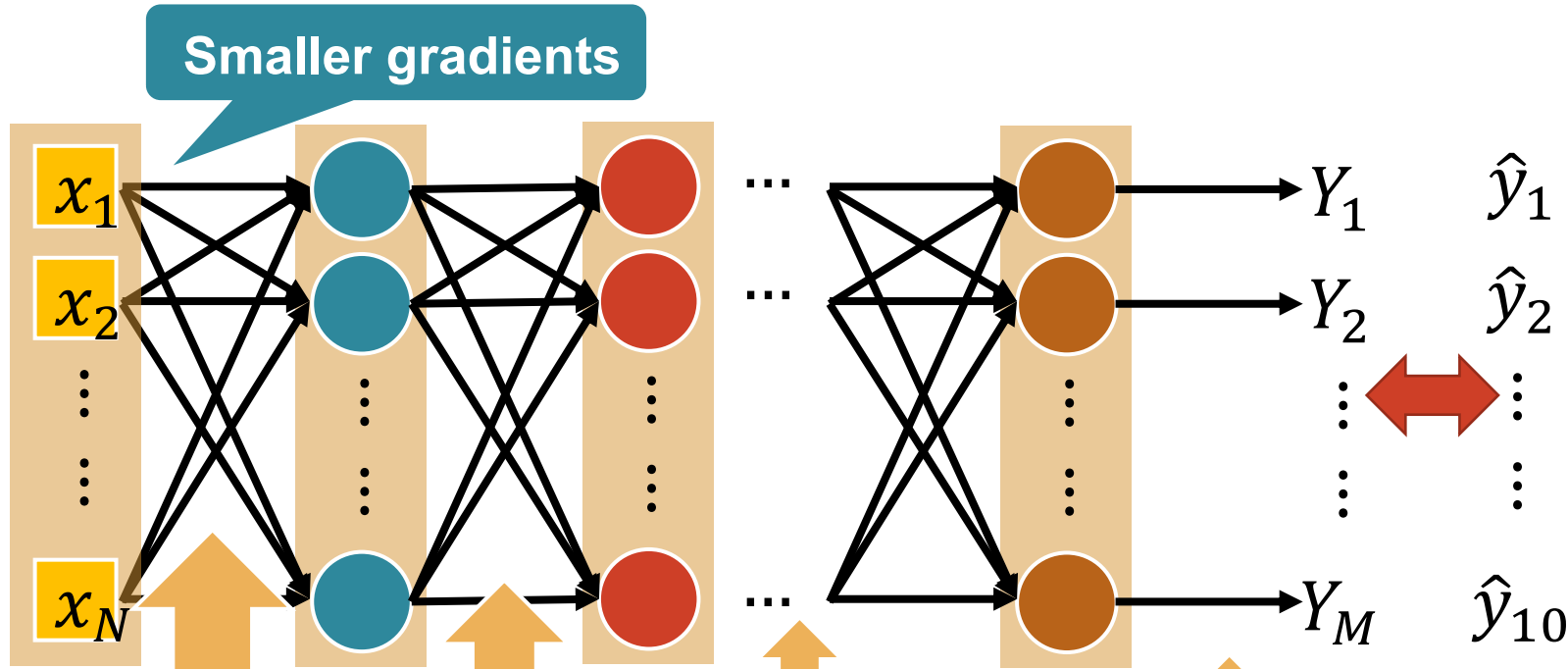
Larger gradients

Learn very fast

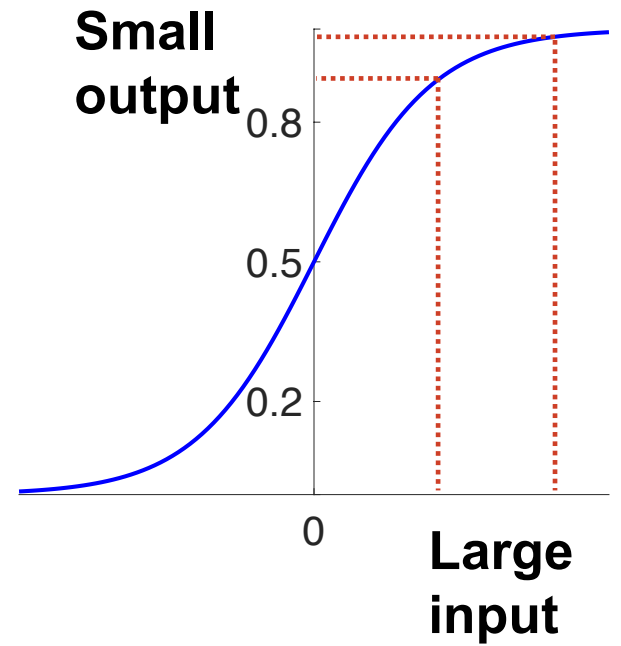
Already converge

## 2. Tricks from training process

### 2.1 Activation function



### Vanishing gradient problem



$$\frac{\partial C}{\partial w} = ? \quad \frac{\Delta C}{\Delta w}$$

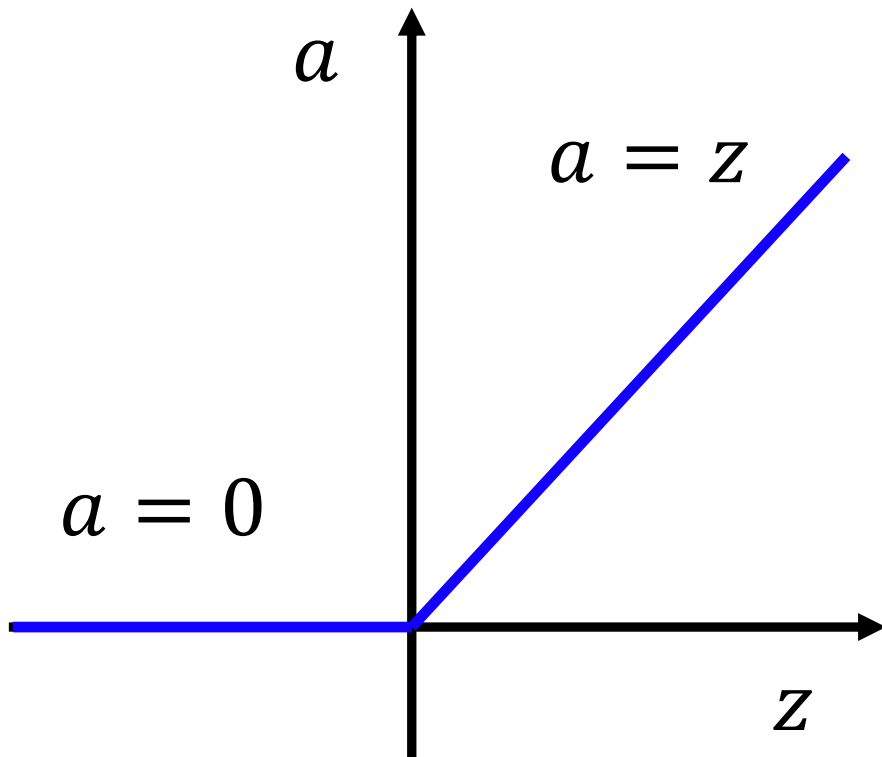
## 2. Tricks from training process

---

### 2.1 Activation function

ReLU

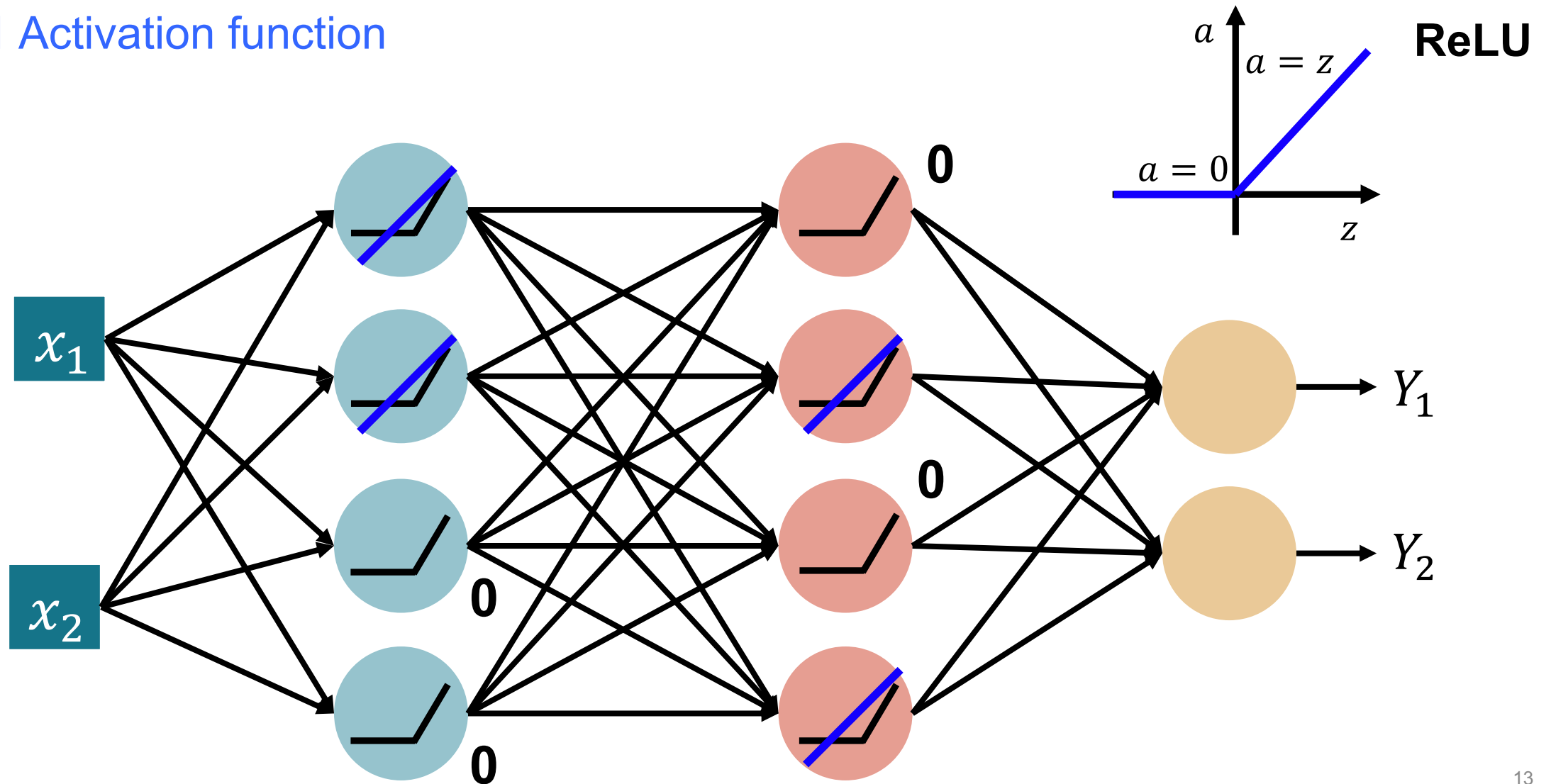
#### Rectified Linear Unit (ReLU)



- Fast to compute;
- Biological reason;
- Infinite sigmoid with different biases;
- Solve vanishing gradient problem.

## 2. Tricks from training process

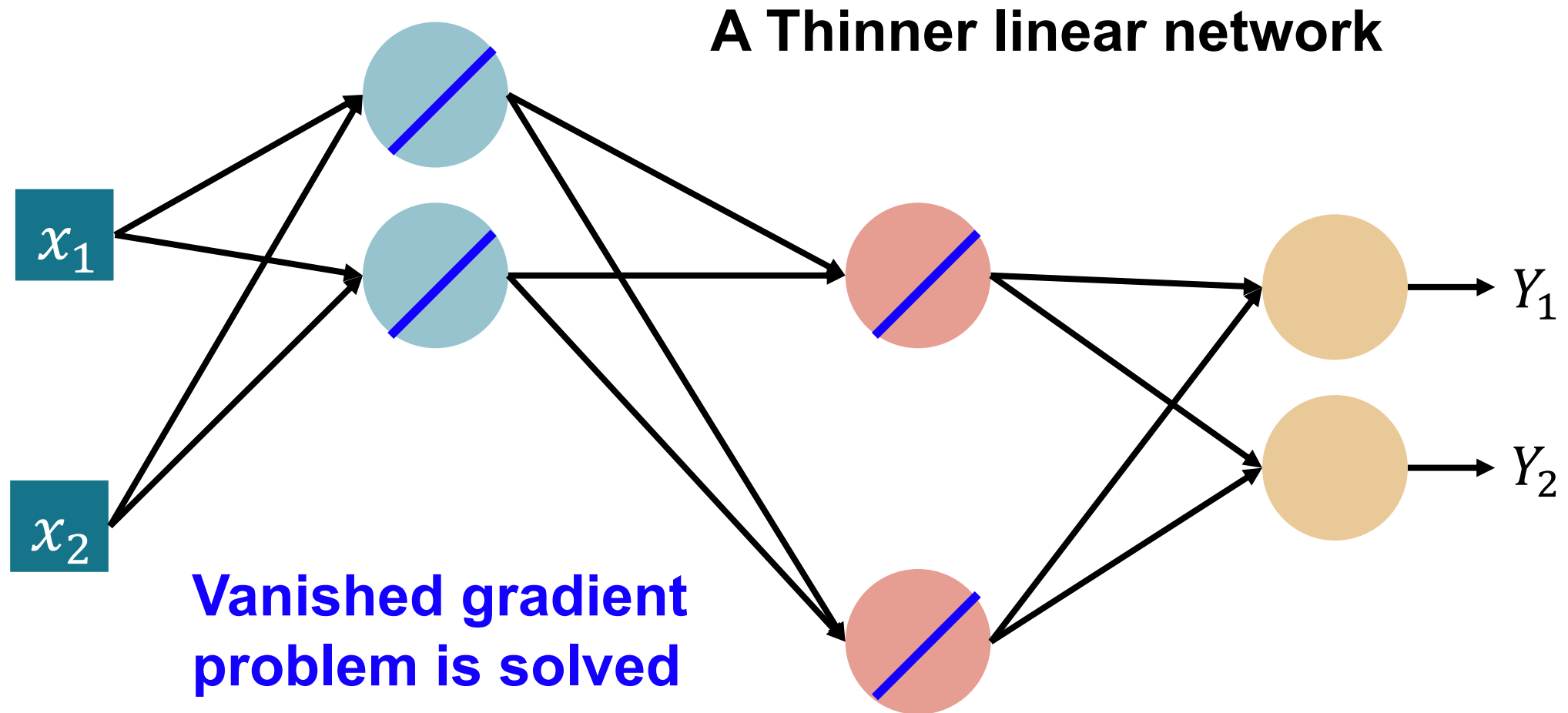
### 2.1 Activation function



## 2. Tricks from training process

### 2.1 Activation function

ReLU



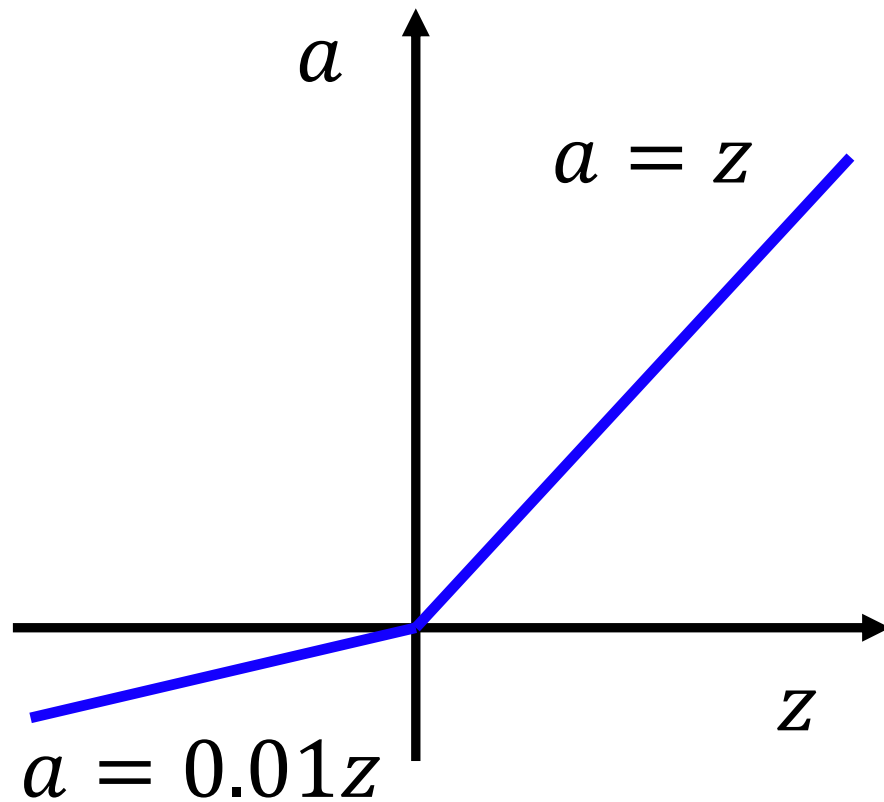
## 2. Tricks from training process

---

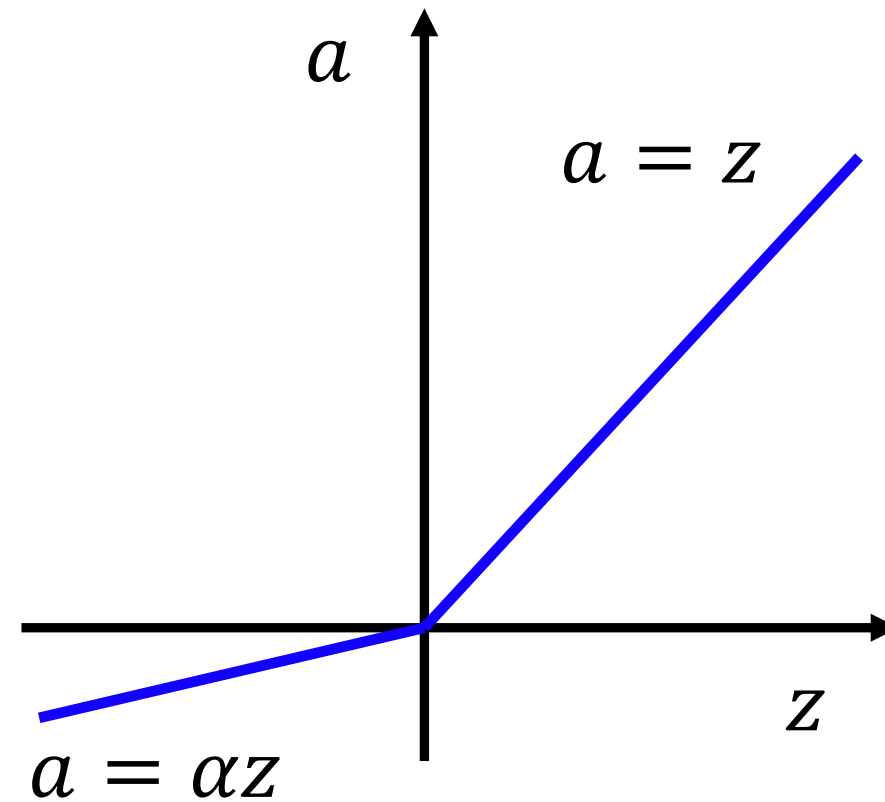
### 2.1 Activation function

ReLU variant

Leaky ReLU



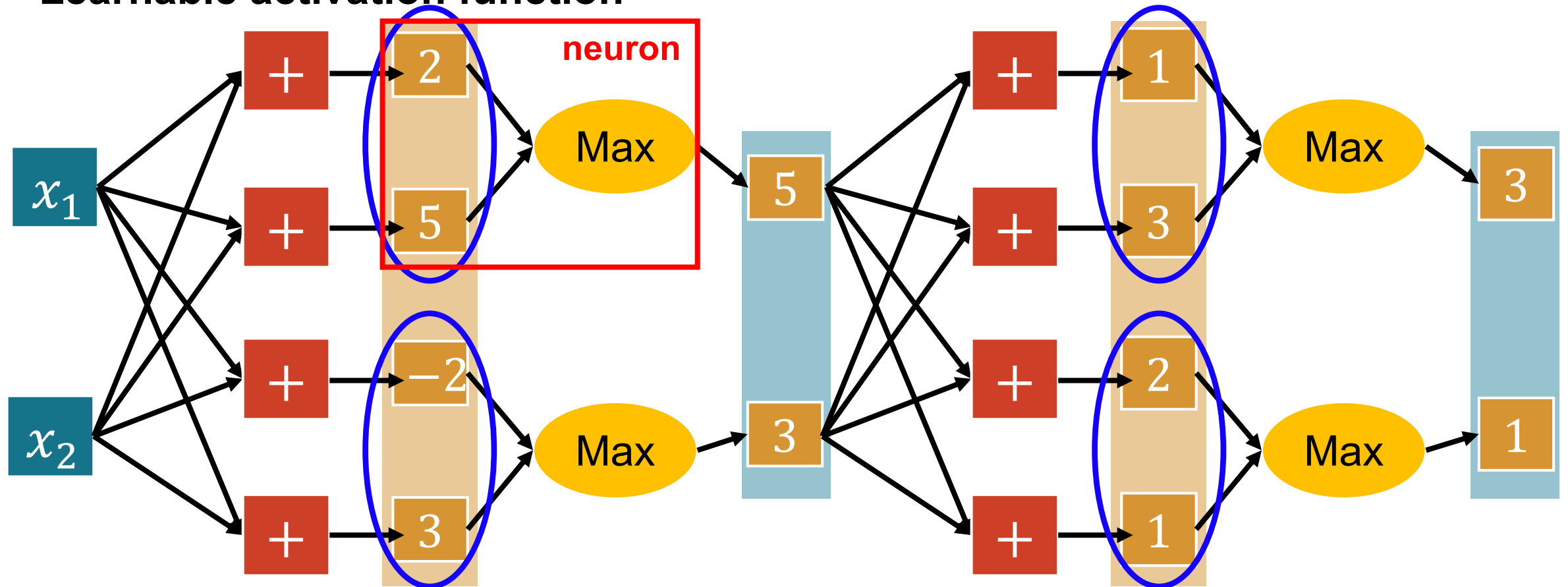
Parametric ReLU



## 2. Tricks from training process

### 2.1 Activation function

#### Learnable activation function



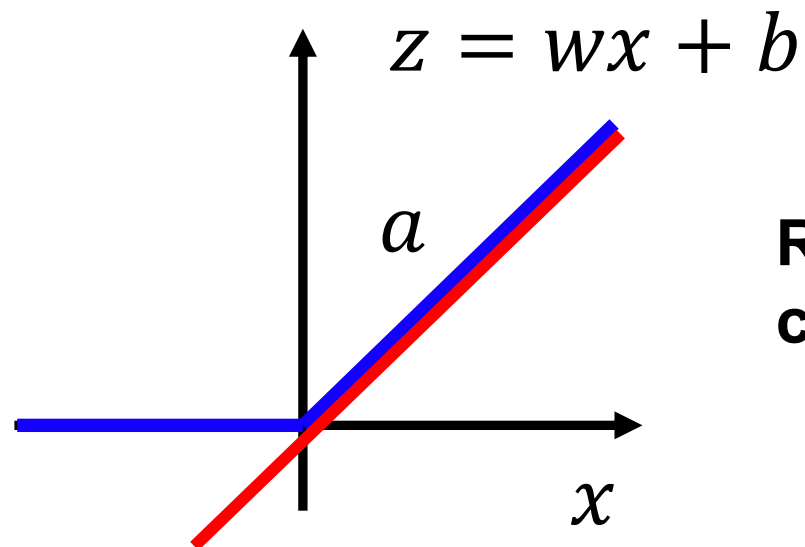
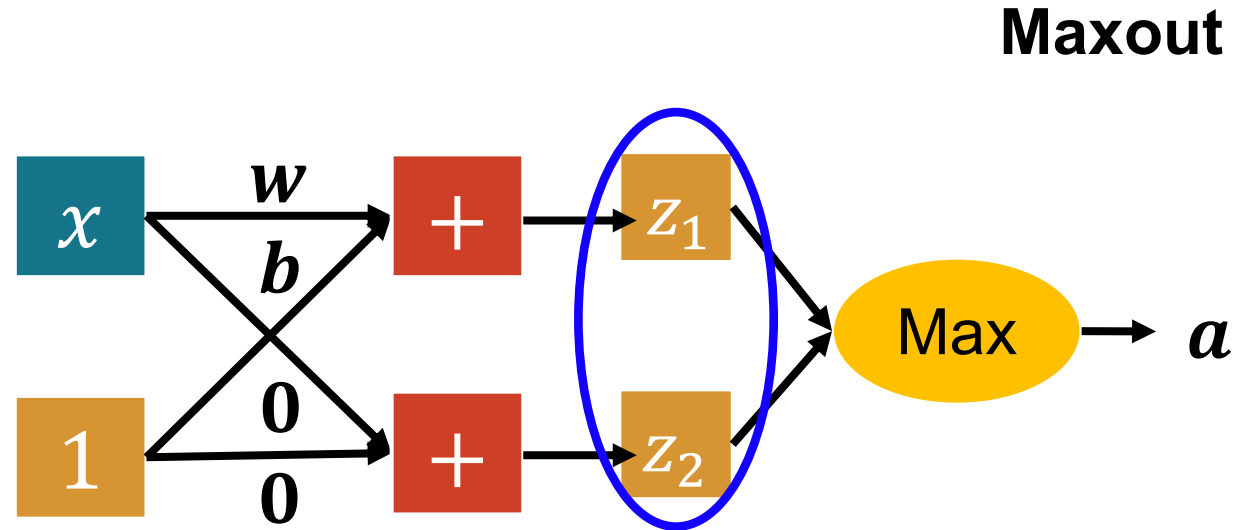
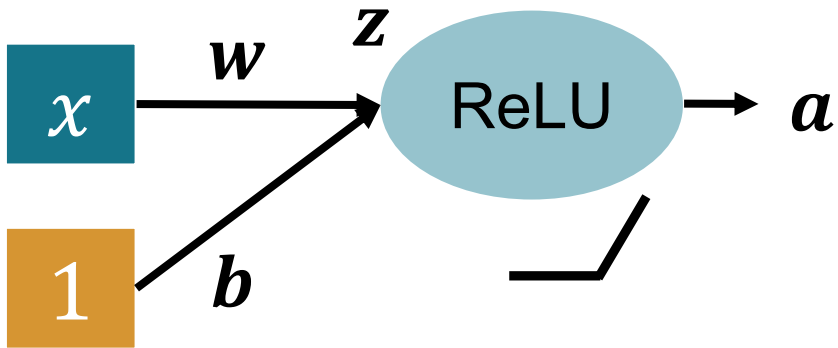
You can have more than 2 elements in a group



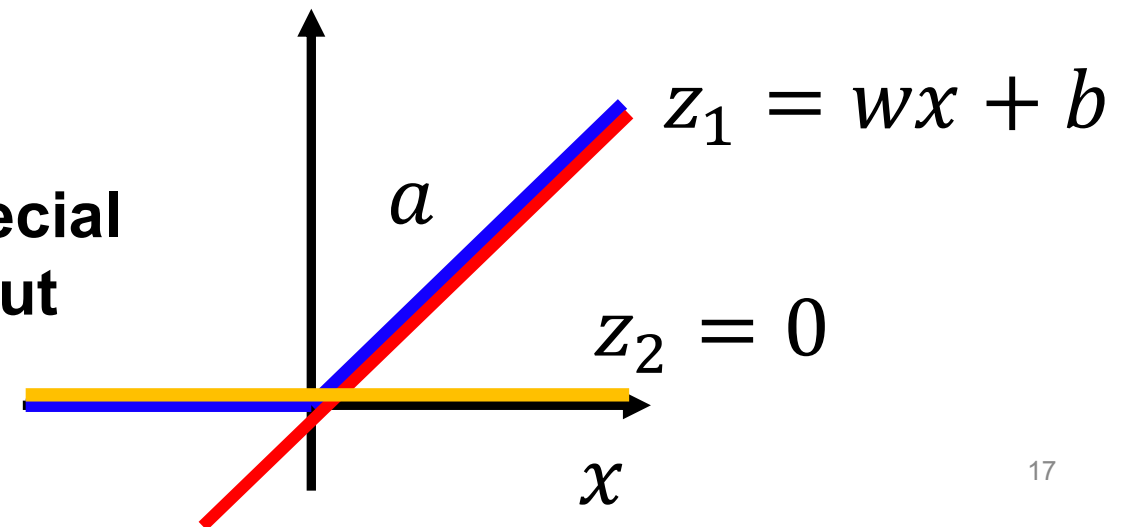
## 2. Tricks from training process

### 2.1 Activation function

Learnable activation function



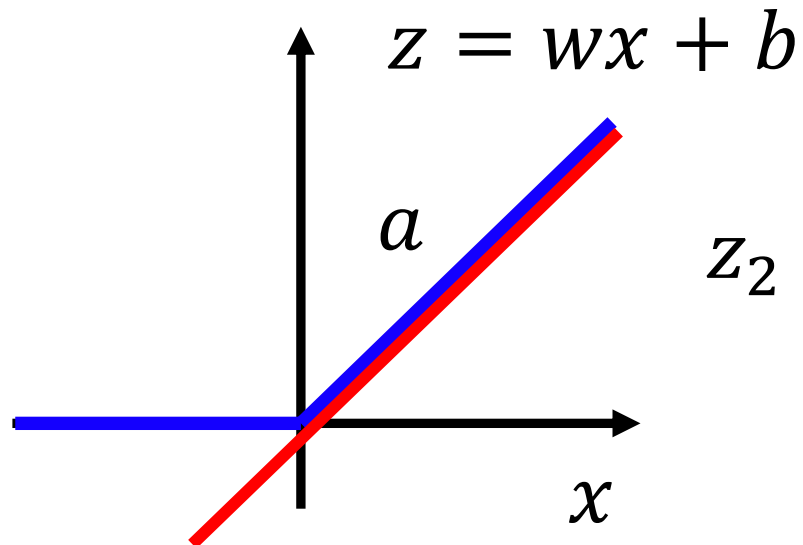
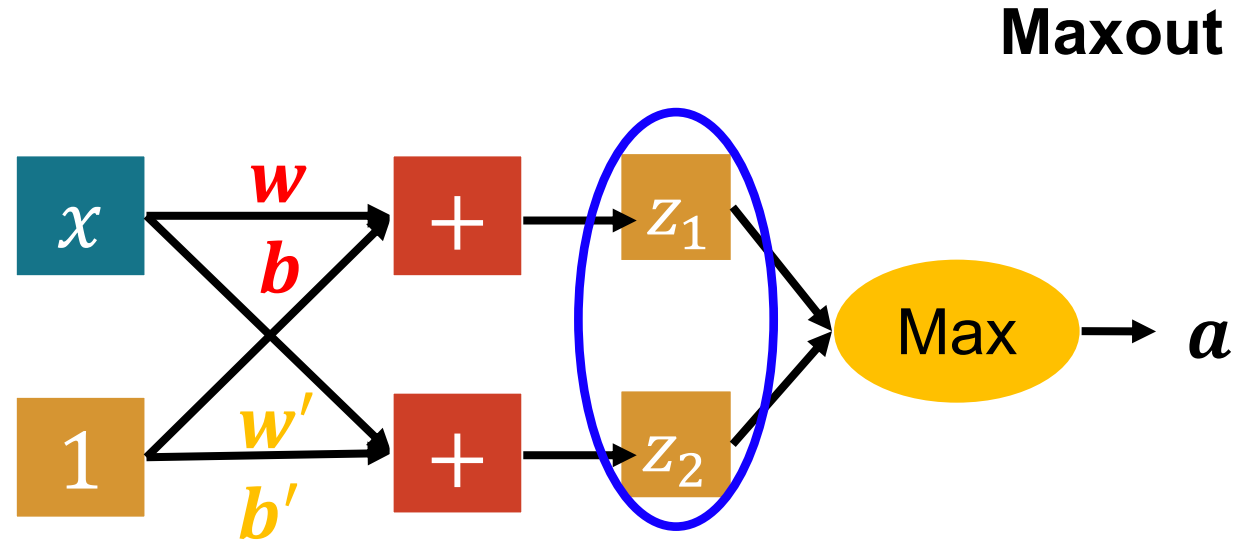
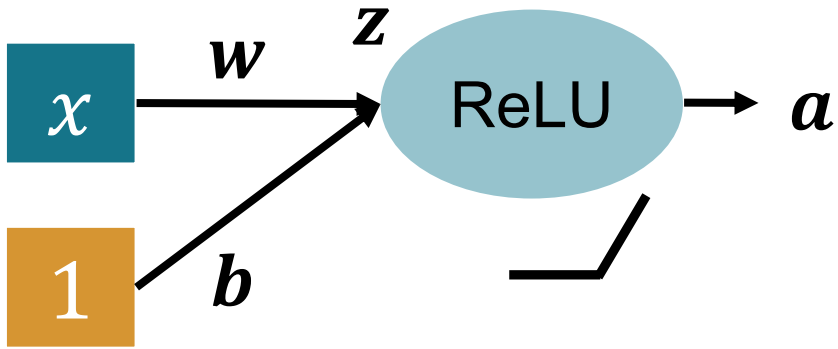
ReLU is a special case of Maxout



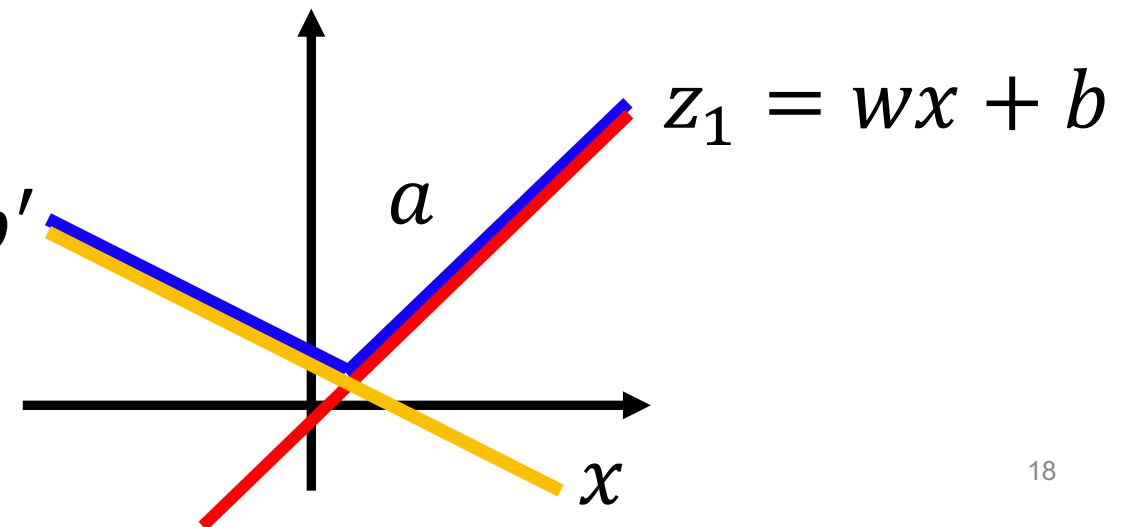
## 2. Tricks from training process

### 2.1 Activation function

#### Learnable activation function



$$z_2 = w'x + b'$$



## 2. Tricks from training process

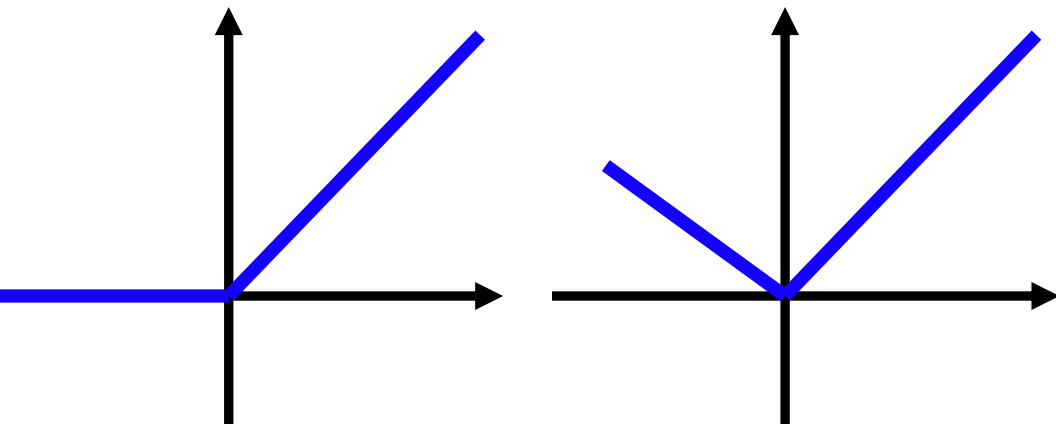
### 2.1 Activation function

Maxout

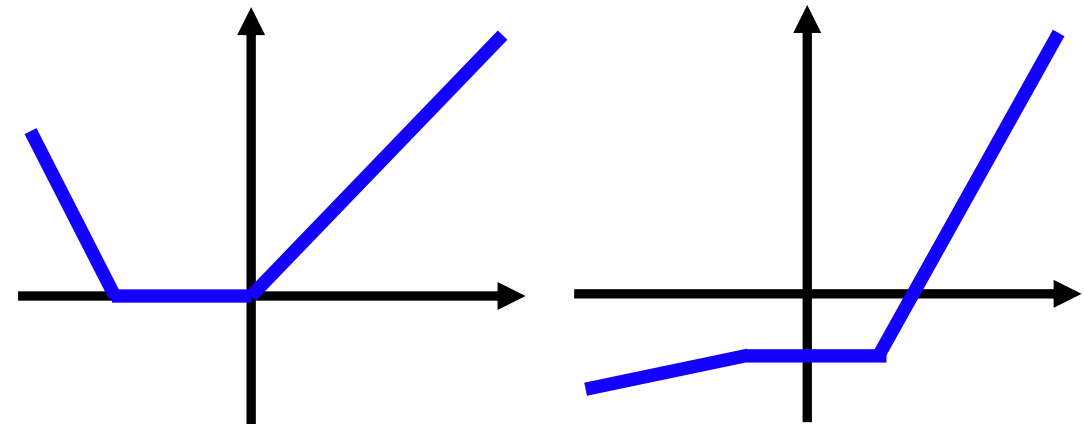
#### Learnable activation function

- Activation function in maxout network can be any piecewise linear convex function;
- How many pieces depending on how many elements in a group.

2 element in a group



3 element in a group



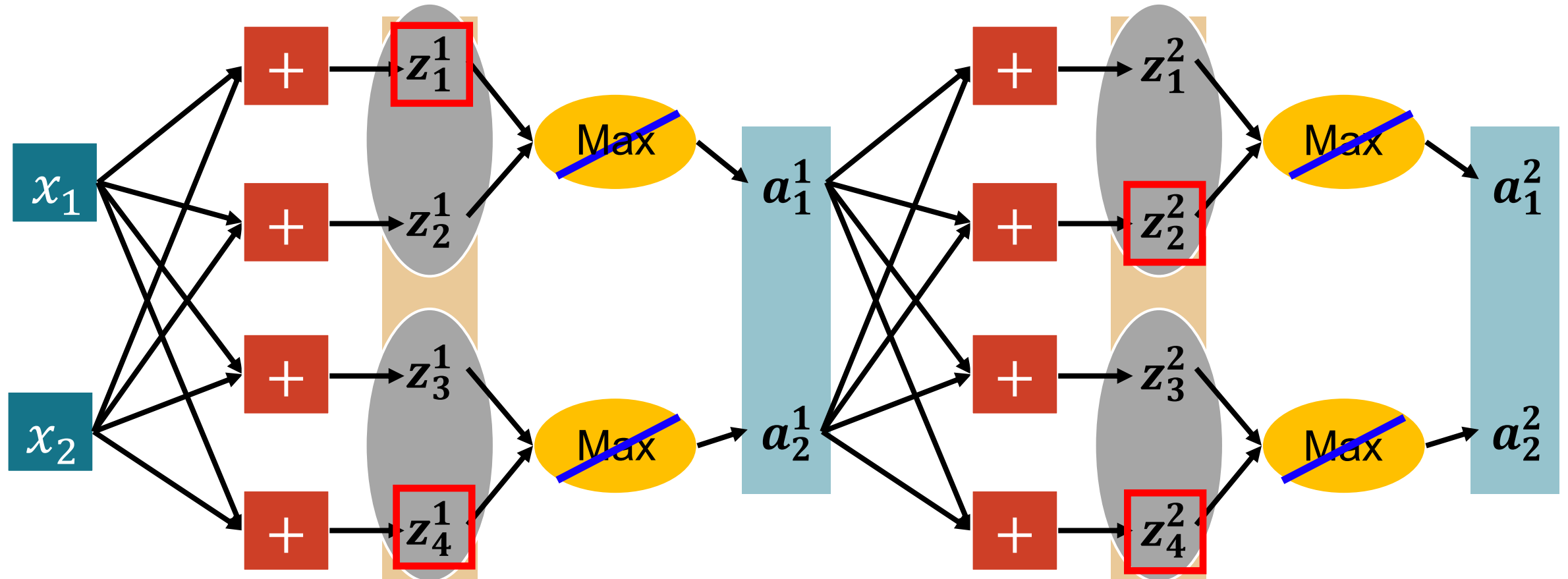
## 2. Tricks from training process

### 2.1 Activation function

Maxout

Training

Given a training data  $x$ , the max of  $z$  is determined

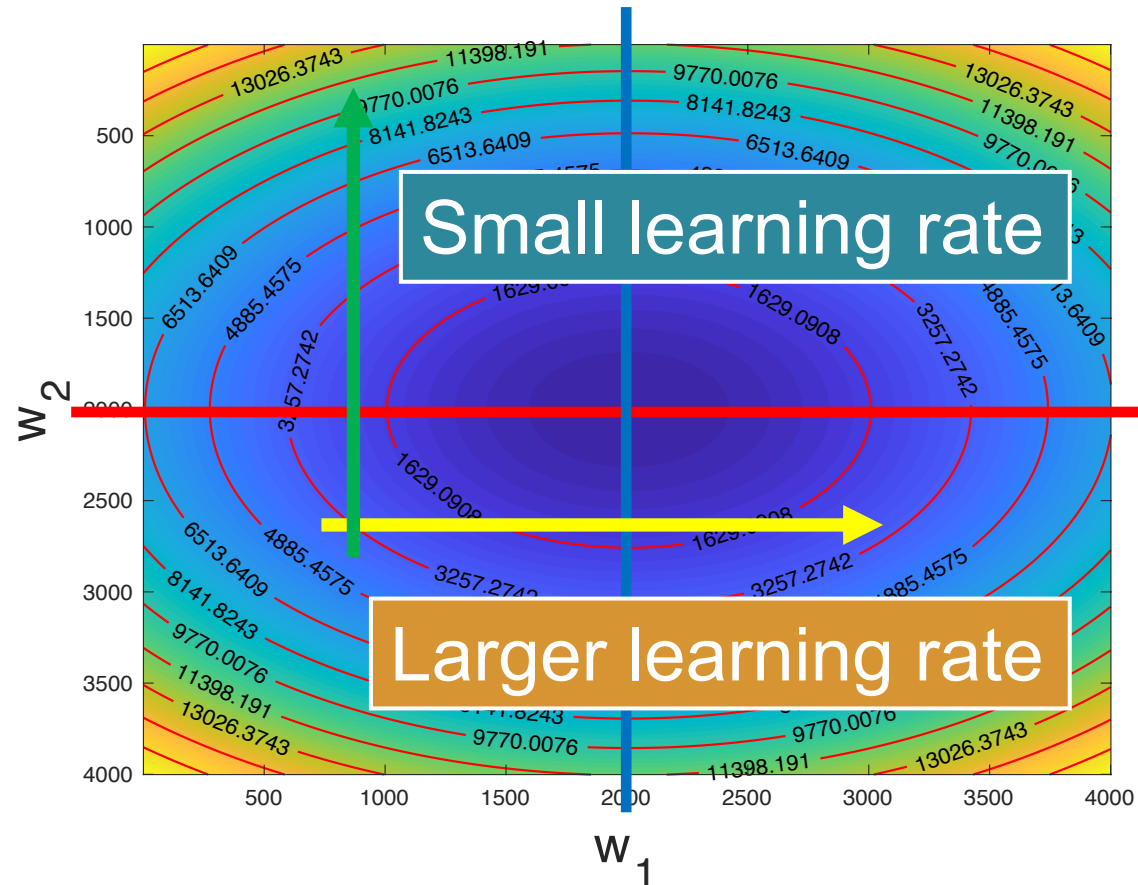


Train this thin and linear network

## 2. Tricks from training process

### 2.2 Adaptive learning rate

#### Adagrad



Previous Lecture

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

|First derivative|

Second derivative

## 2. Tricks from training process

---

### 2.2 Adaptive learning rate

#### RMSProp

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

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

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

⋮

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

**Root mean square of the gradients with previous gradients being decayed.**

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

$$\sigma^1 = \sqrt{\alpha(\sigma^0)^2 + (1 - \alpha)(g^1)^2}$$

$$\sigma^2 = \sqrt{\alpha(\sigma^1)^2 + (1 - \alpha)(g^2)^2}$$

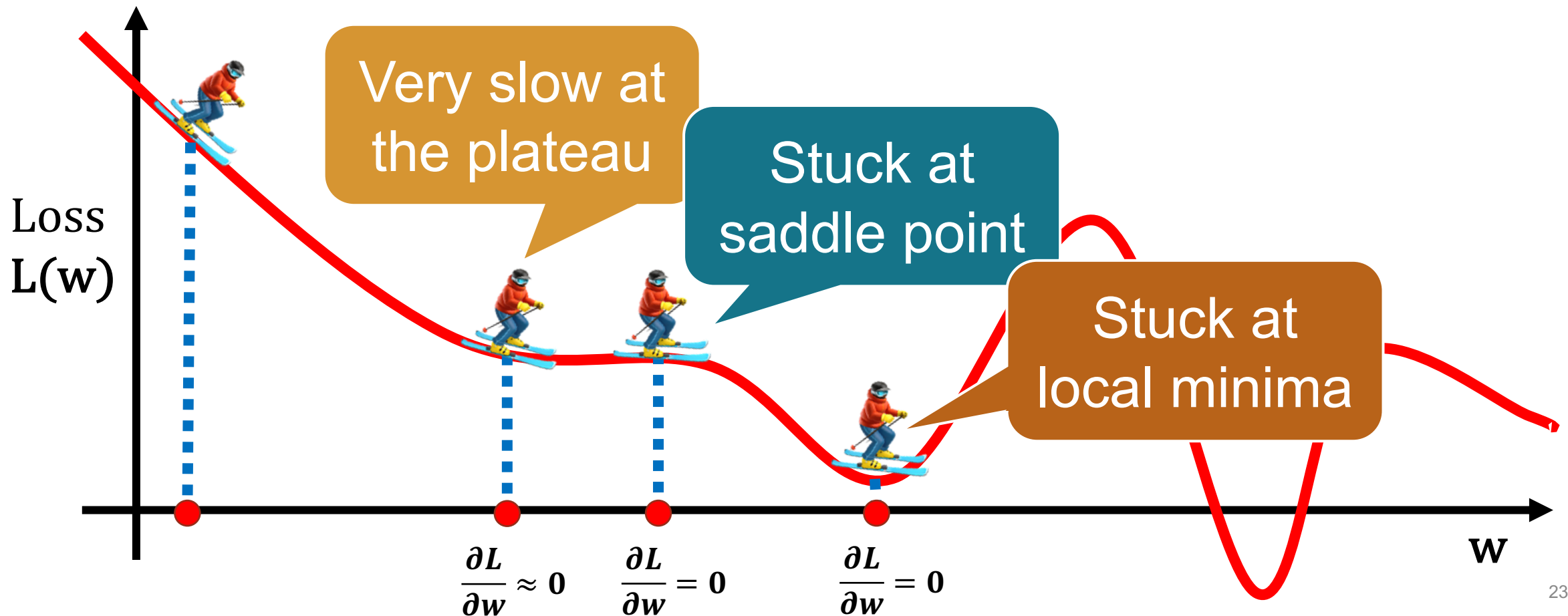
$$\sigma^t = \sqrt{\alpha(\sigma^{t-1})^2 + (1 - \alpha)(g^t)^2}$$

## 2. Tricks from training process

### 2.2 Adaptive learning rate

Previous Lecture

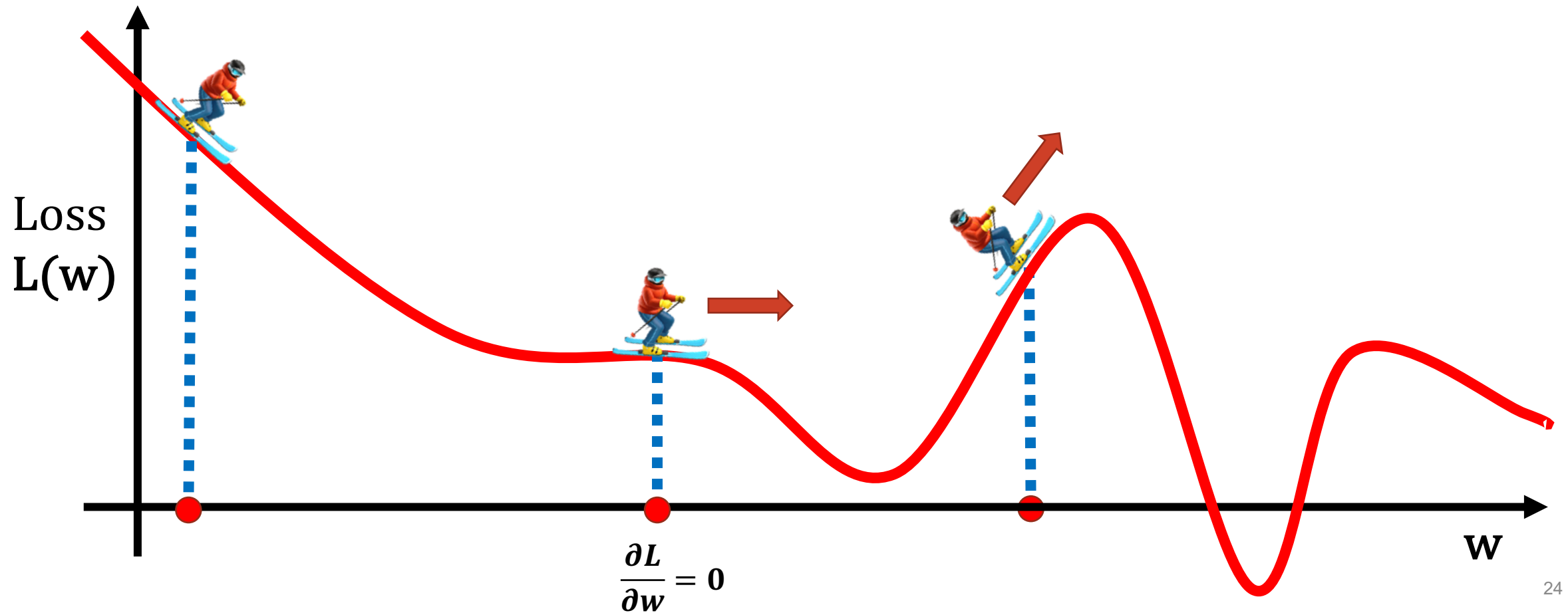
#### Momentum



## 2. Tricks from training process

### 2.2 Adaptive learning rate

#### Momentum

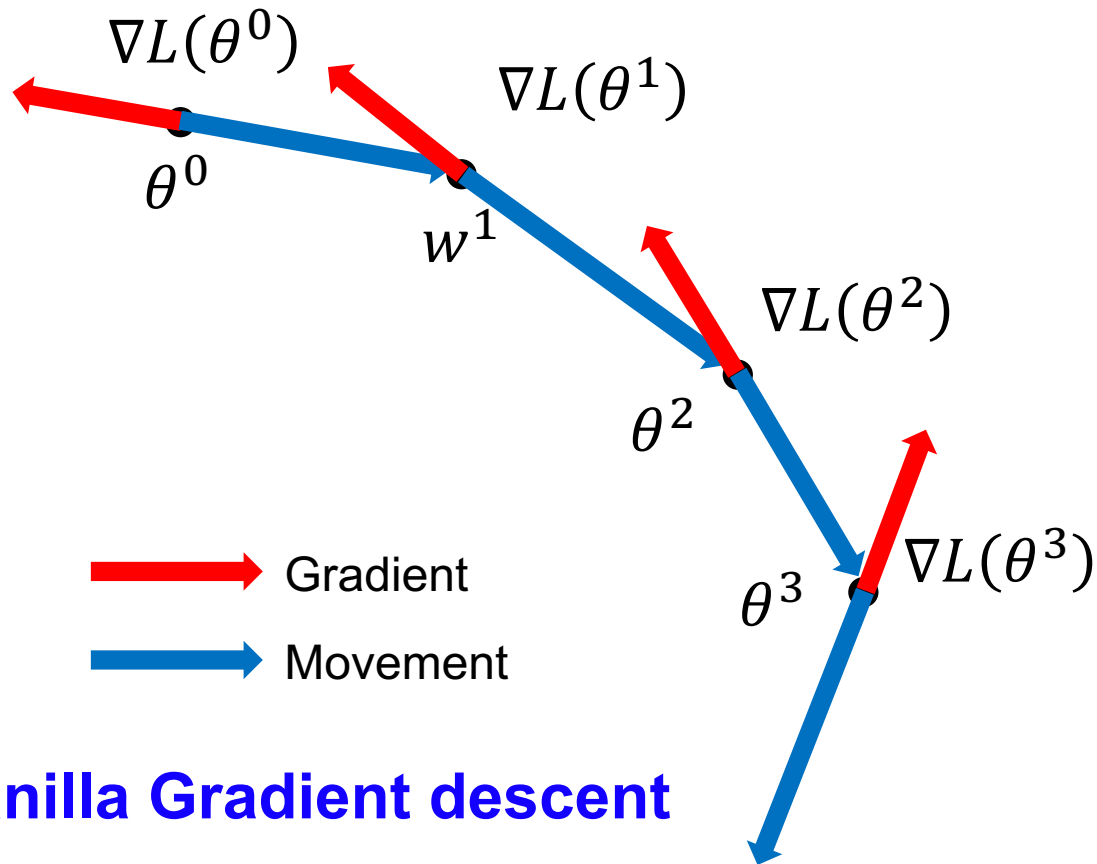




## 2. Tricks from training process

### 2.2 Adaptive learning rate

#### Momentum



#### Vanilla Gradient descent

### Previous Lecture

Start at position  $\theta^0$

Calculate gradient at  $\theta^0$ :  $\nabla L(\theta^0)$

Move to  $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Calculate gradient at  $\theta^1$ :  $\nabla L(\theta^1)$

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

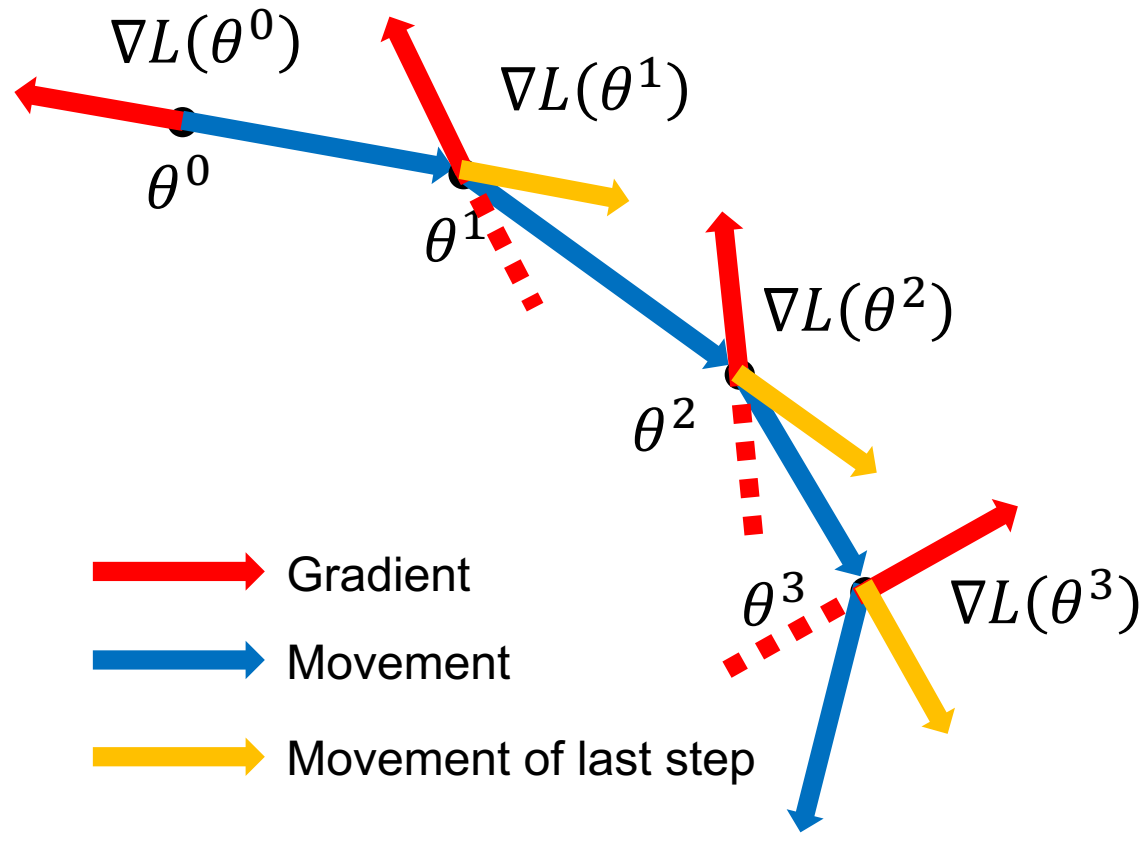
⋮  
⋮

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

## 2. Tricks from training process

### 2.2 Adaptive learning rate

**Momentum**      **Movement:** movement of last step minus gradient at present.



Start at position  $\theta^0$

Movement  $v^0 = 0$

Calculate gradient at  $\theta^0$ :  $\nabla L(\theta^0)$

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

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

Calculate gradient at  $\theta^1$ :  $\nabla L(\theta^1)$

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

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

⋮

**Movement is not only based on gradient, but also consider the previous movement.**

## 2. Tricks from training process

---

### 2.2 Adaptive learning rate

**Momentum**      **Movement:** movement of last step minus gradient at present.

$v^i$  is the weighted sum of all the previous gradient:

$$\nabla L(\theta^0), \nabla L(\theta^1), \dots \nabla L(\theta^{i-1})$$

$$v^0 = 0$$

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

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

Start at position  $\theta^0$

Movement  $v^0 = 0$

Calculate gradient at  $\theta^0$ :  $\nabla L(\theta^0)$

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

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

Calculate gradient at  $\theta^1$ :  $\nabla L(\theta^1)$

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

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

⋮

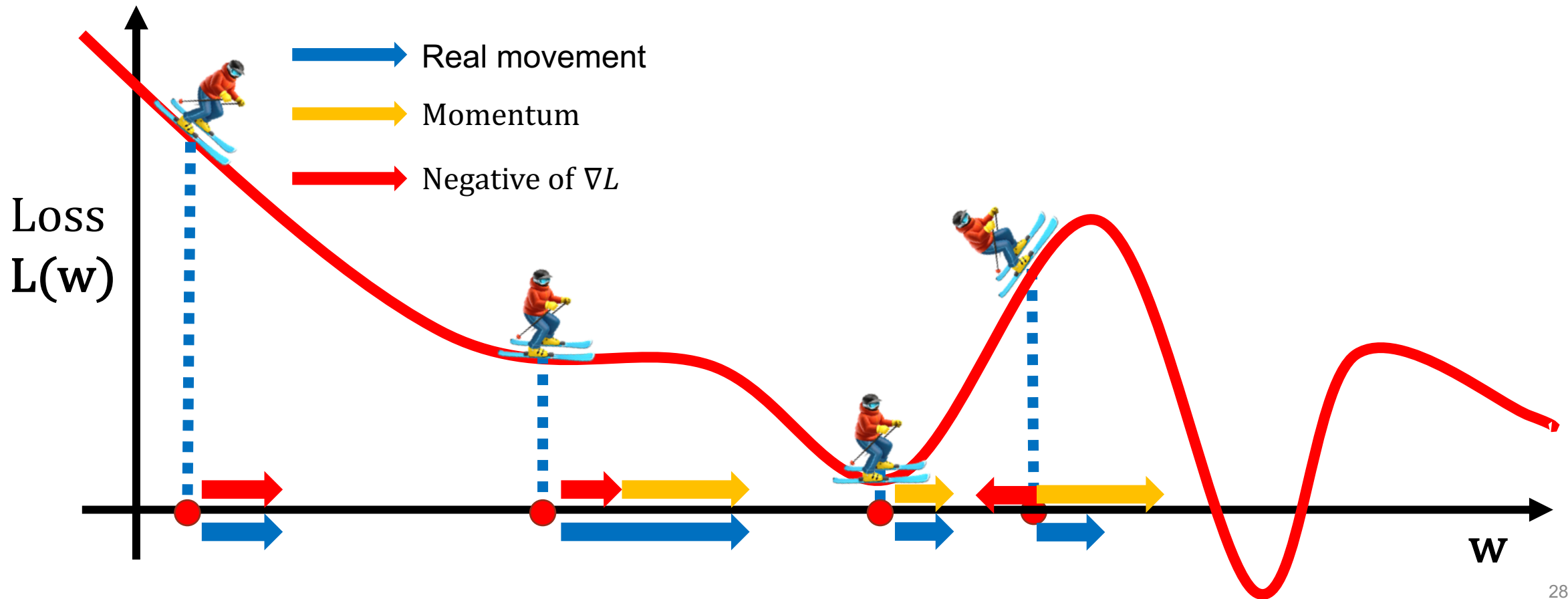
**Movement is not only based on gradient, but also consider the previous movement.**

## 2. Tricks from training process

### 2.2 Adaptive learning rate

#### Momentum

Movement = Negative of  $\nabla L$  + Momentum



## 2. Tricks from training process

---

### 2.2 Adaptive learning rate

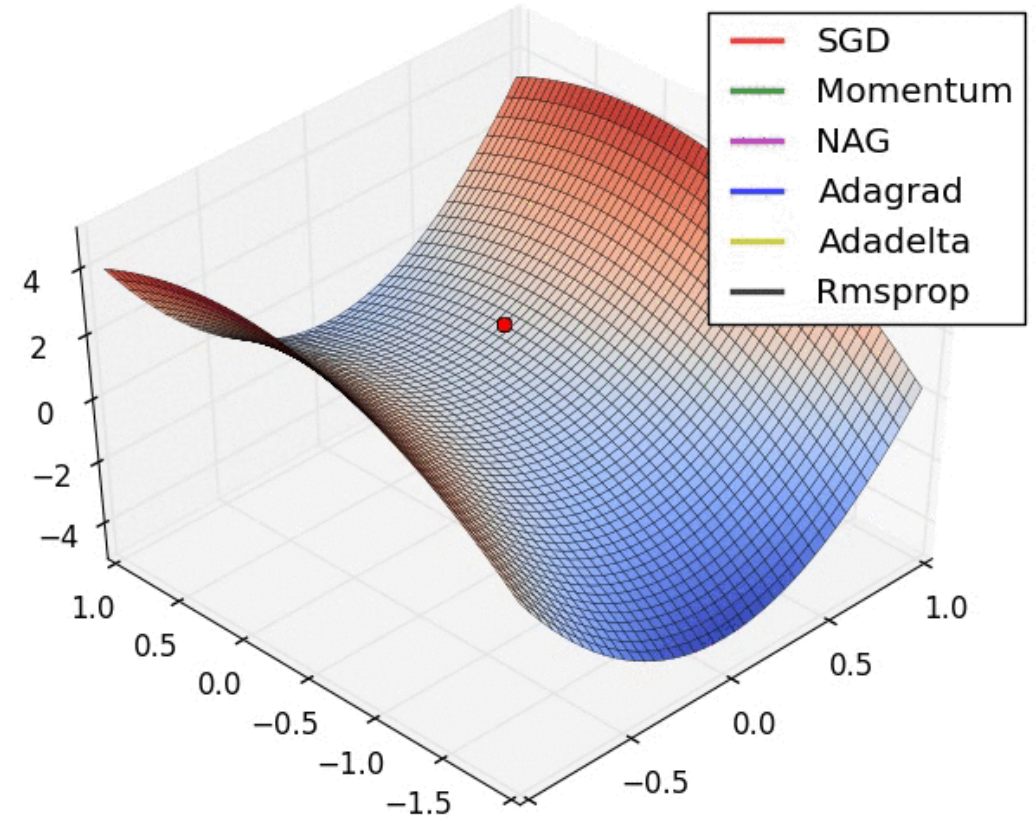
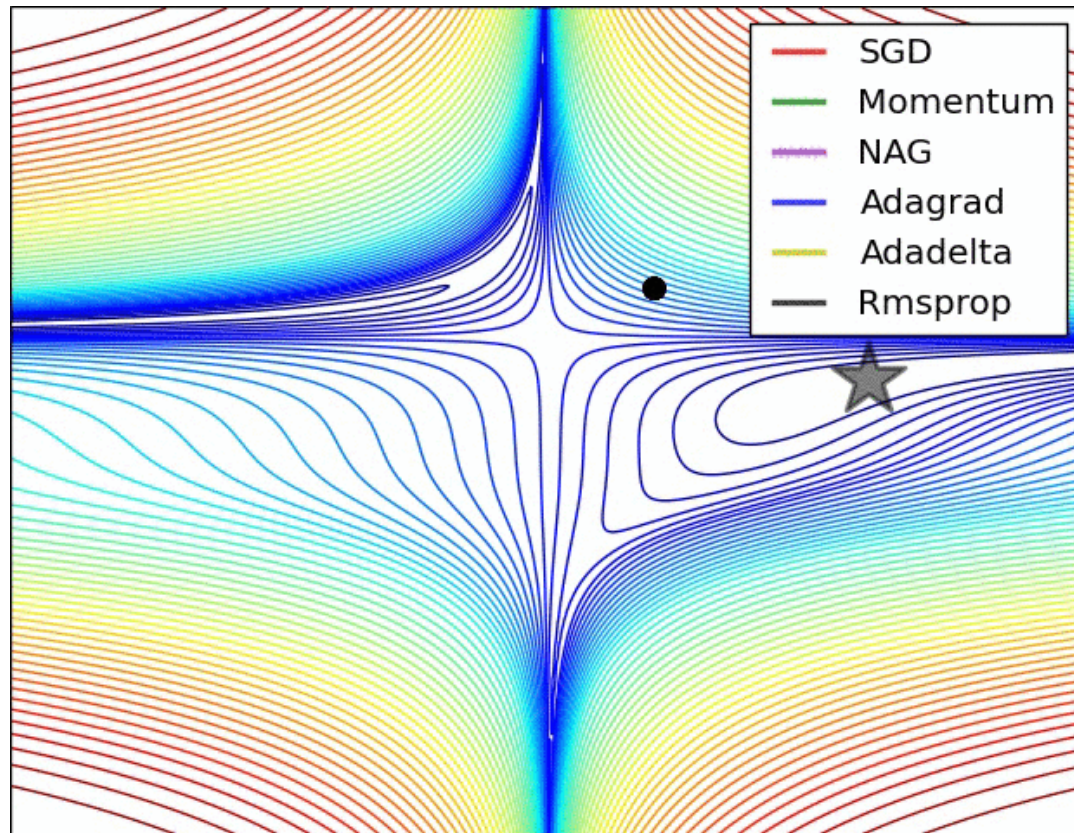
**Momentum**

**Adam =**

**RMSProp + Momentum**

## 2. Tricks from training process

### 2.2 Adaptive learning rate



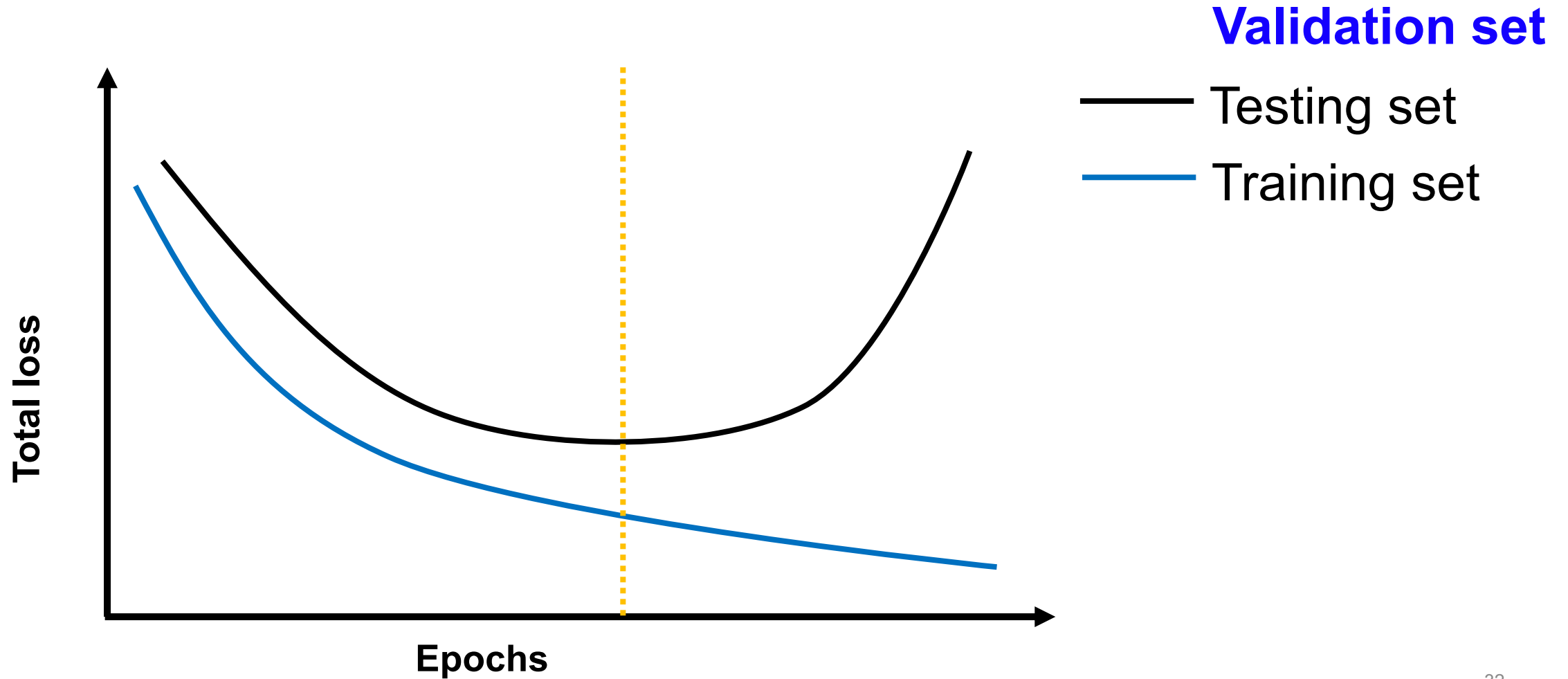
# Tricks from testing process

# 03



# 3. Tricks from testing process

## 3.1 Early stopping





### 3. Tricks from testing process

---

#### 3.2 Regularization

Previous Lecture

**Regularization** — Adding a penalty on model complexity

$$L = \sum_n \left( \hat{y}^n - \left( b + \sum w_i x_i \right) \right)^2 + \lambda \sum (w_i)^2$$

**L2-norm**

**Ridge regression**

$$L = \sum_n \left( \hat{y}^n - \left( b + \sum w_i x_i \right) \right)^2 + \lambda \sum |w_i|$$

**L1-norm**

**Lasso regression**

Least absolute shrinkage and selection operator

### 3. Tricks from testing process

---

#### 3.2 Regularization

New loss function to be minimized: not only minimizing original cost but also let it be close to zero

$$L'(\theta) = L(\theta) + \lambda \frac{1}{2} \|\theta\|^2 \leftarrow \text{Regularization term}$$

**Original loss**

e.g. minimize square error,  
cross entropy ...

$$\theta = \{w_1, w_2, \dots \dots \}$$

**L2 regularization:**  $\|\theta\|^2 = (w_1)^2 + (w_2)^2 + \dots$

### 3. Tricks from testing process

---

#### 3.2 Regularization

#### L2 regularization:

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

Update:

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

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

**Close to zero**

**Weight decay**

### 3. Tricks from testing process

---

#### 3.2 Regularization

**L1 regularization:**  $\|\theta\|^1 = |w_1| + |w_2| + \dots$

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

Update:

$$\begin{aligned} w^{t+1} &= w^t - \eta \frac{\partial L'}{\partial w} = w^t - \eta \left( \frac{\partial L}{\partial w} + \lambda \text{sign}(w^t) \right) \\ &= w^t - \eta \frac{\partial L}{\partial w} - \eta \lambda \text{sign}(w^t) \quad \text{Always delete} \\ &\quad (1 - \eta \lambda) w^t - \eta \frac{\partial L}{\partial w} \end{aligned}$$

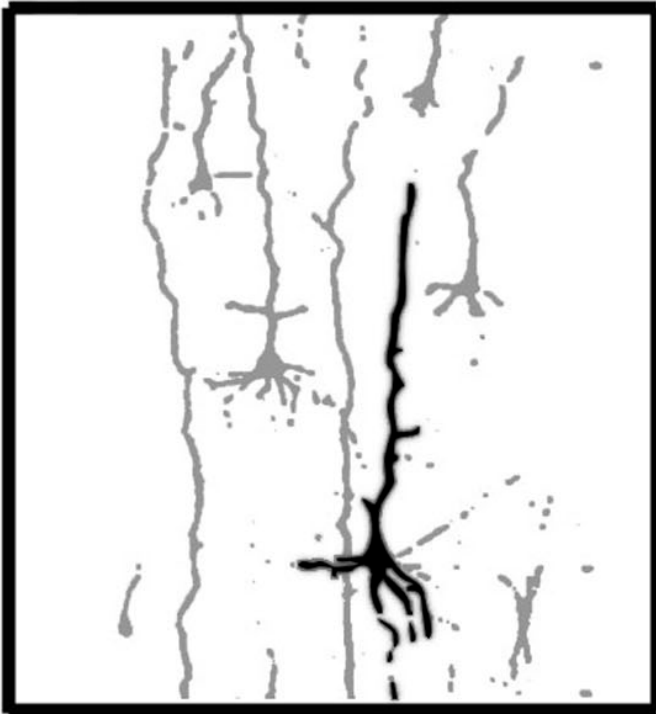
### 3. Tricks from testing process

---

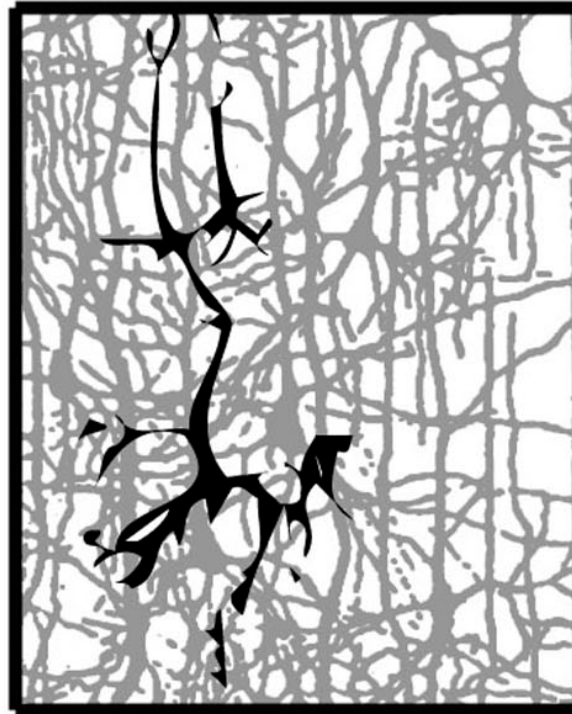
#### 3.2 Regularization

#### Weight decay

**Birth**



**7 Years**



**Adult**

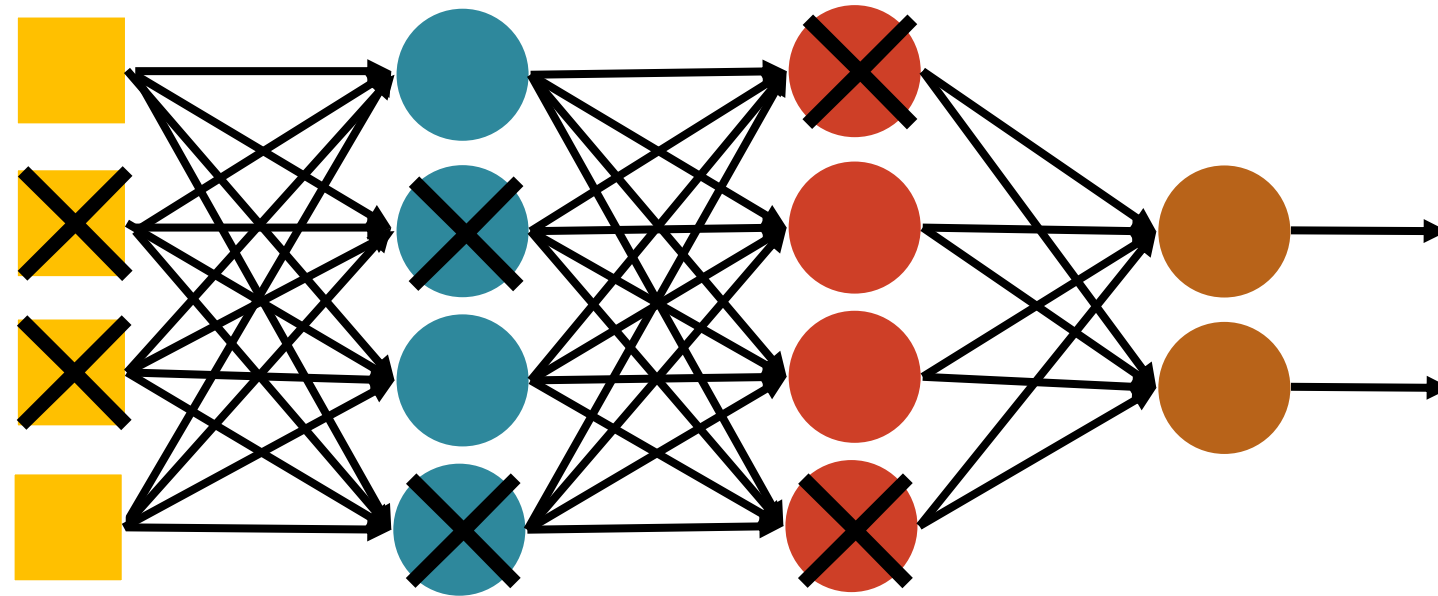


## 3. Tricks from testing process

### 3.3 Dropout

For each mini-batch, we resample the dropout neurons

Training

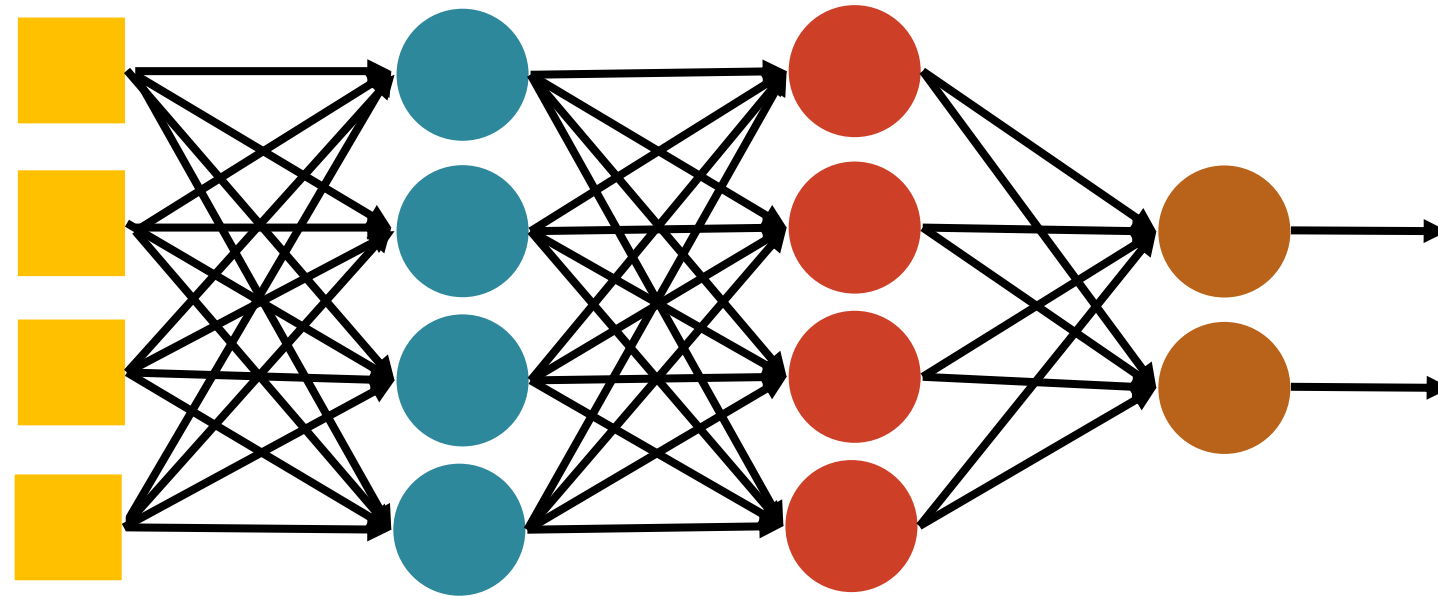


- Each time before updating the parameters :
  - Each neuron has probability  $p$  to dropout, then the structure of the network is changed;
  - Using the new network for training.

## 3. Tricks from testing process

### 3.3 Dropout

Testing



No dropout

- For each neuron has probability  $p$  to dropout in training process, the weight times  $(1 - p)$ ;
- E.g., Assume the dropout rate is 40%, if a weight  $w = 1$  by training, set  $w = 0.6$  for testing.

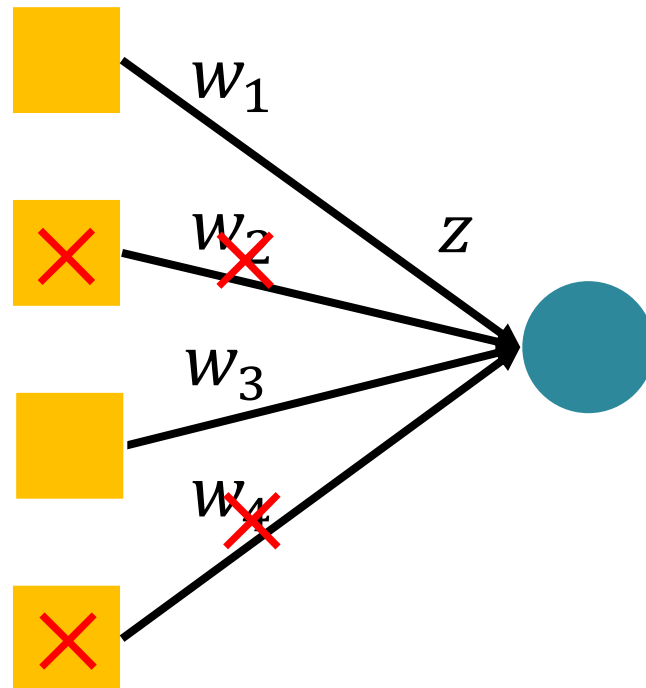
# 3. Tricks from testing process

## 3.3 Dropout

### Why need the weights multiply $(1 - p)$ when testing?

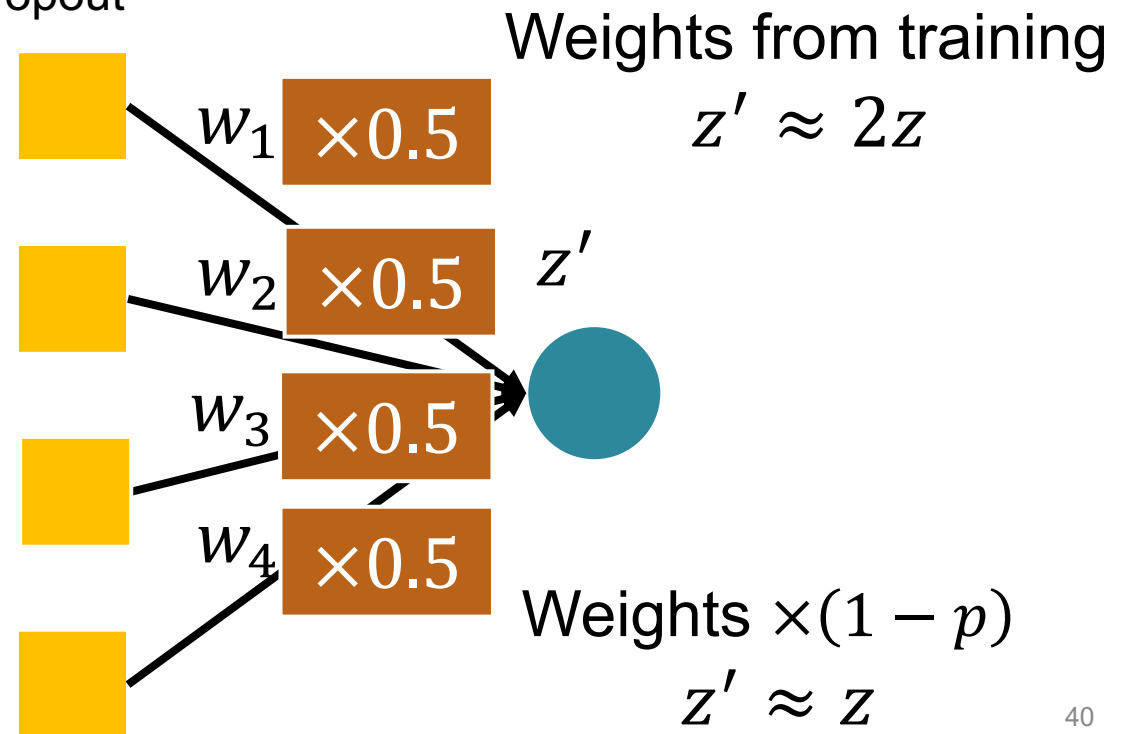
#### Training processing

Assume the dropout rate is 50%



#### Testing processing

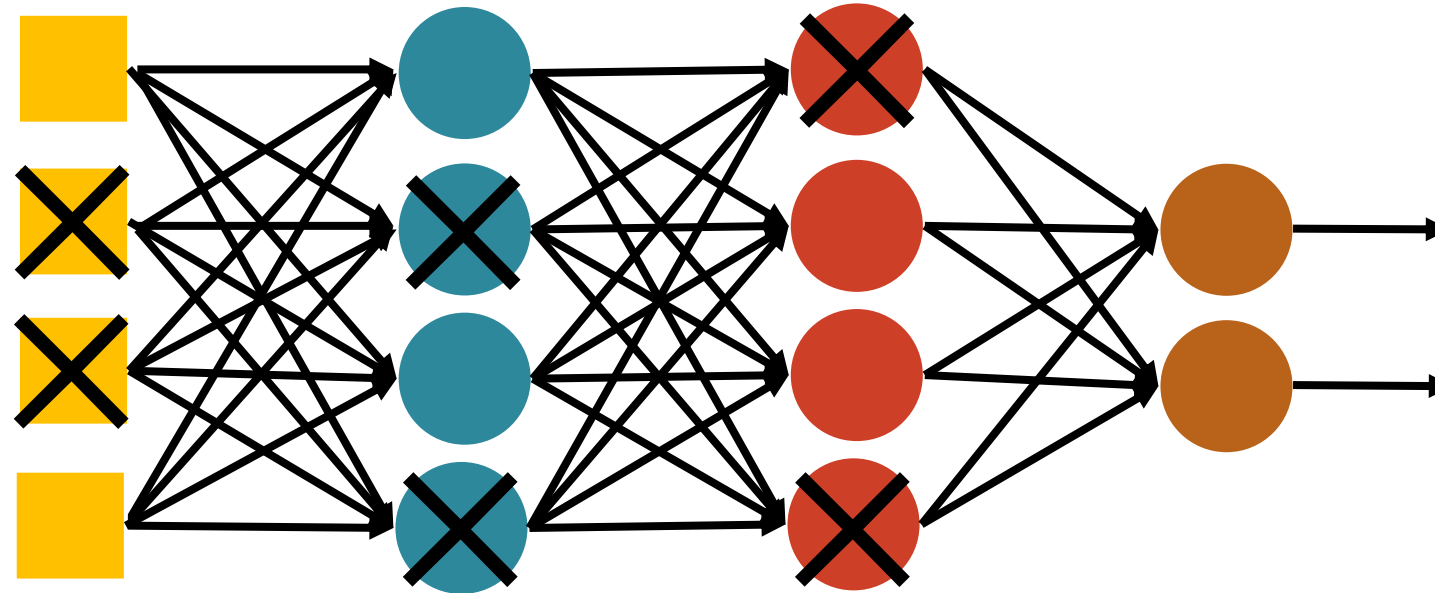
No dropout





## 3. Tricks from testing process

### 3.3 Dropout



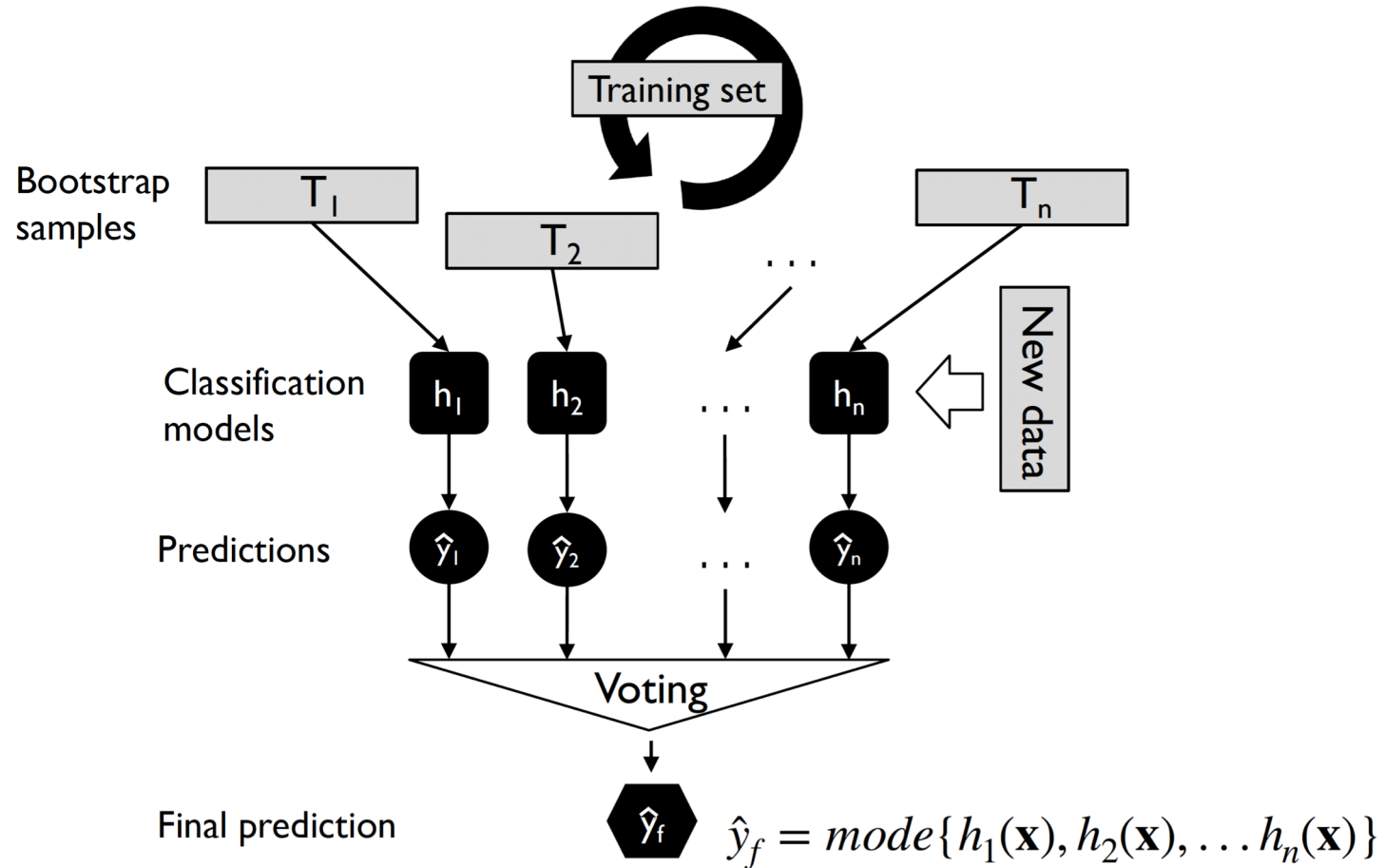
#### ■ Considering a teamwork

- If everyone expect the partner will do the work but do nothing himself, nothing will be done finally;
- However, if everyone know some partner will dropout, you need to work harder to finish the task;
- When submitting (testing), no one dropout actually, the results would be good eventually.

# 3. Tricks from testing process

## 3.3 Dropout -- A kind of ensemble

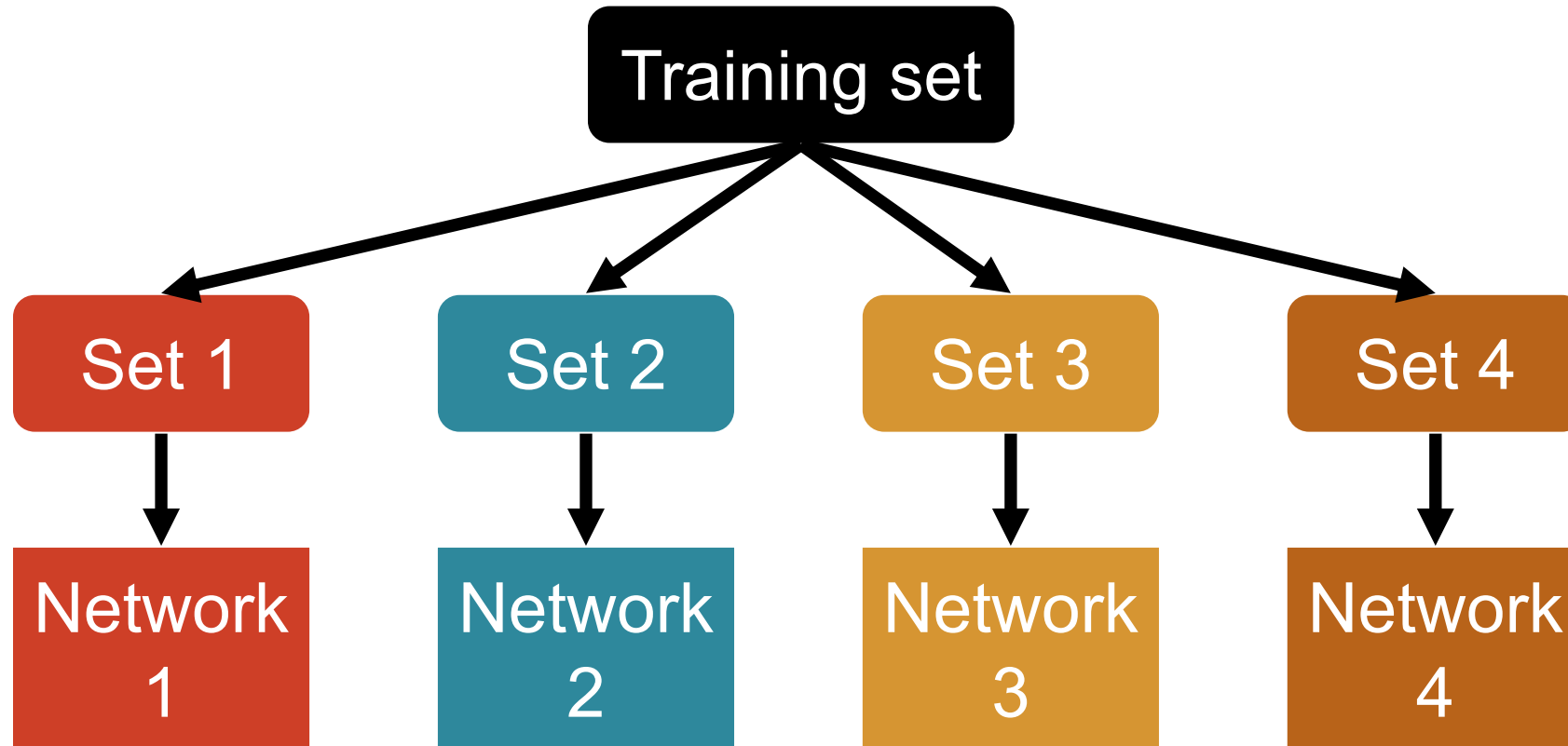
Previous Lecture



## 3. Tricks from testing process

---

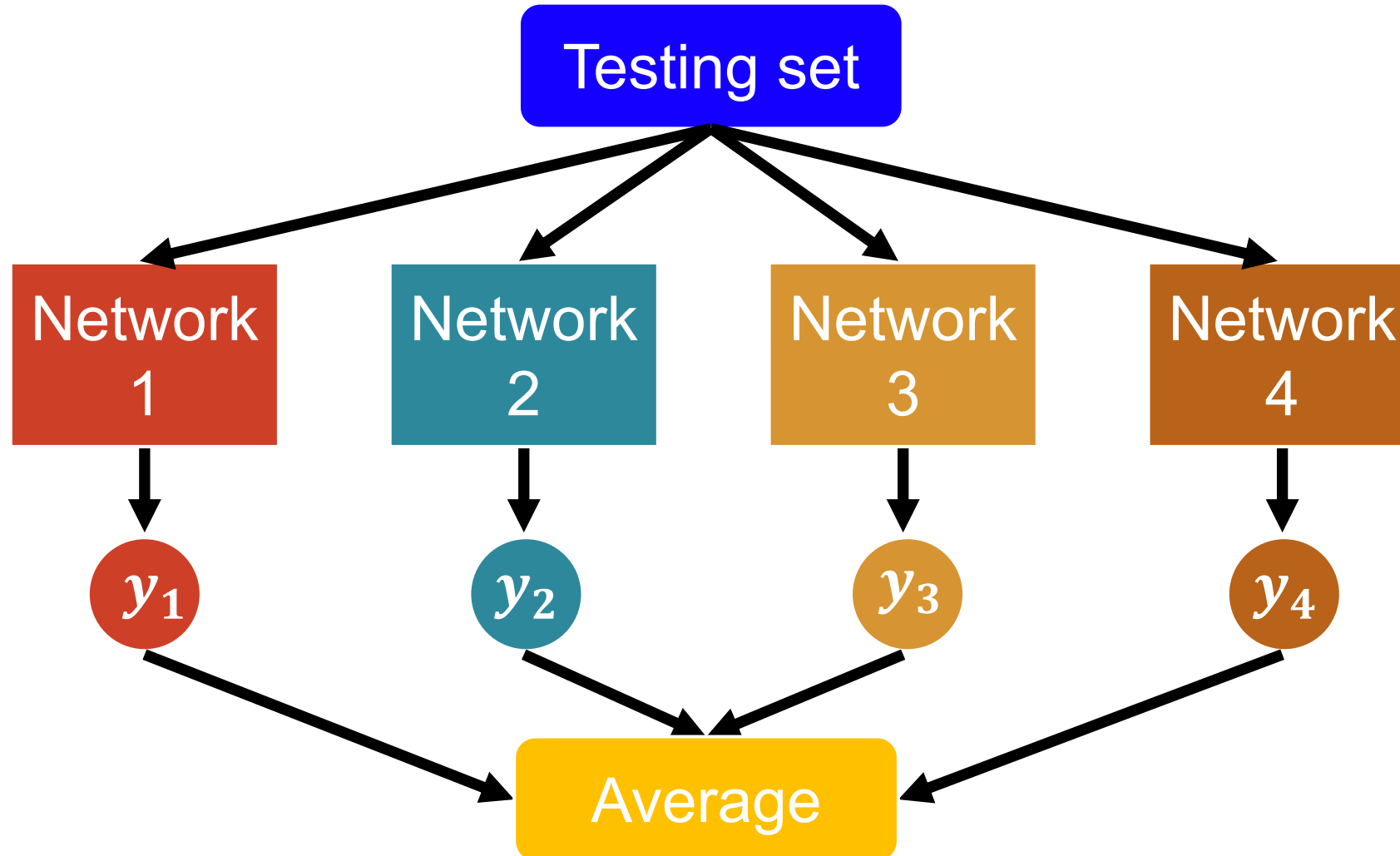
### 3.3 Dropout -- A kind of ensemble



**Train a bunch of networks with different structures**

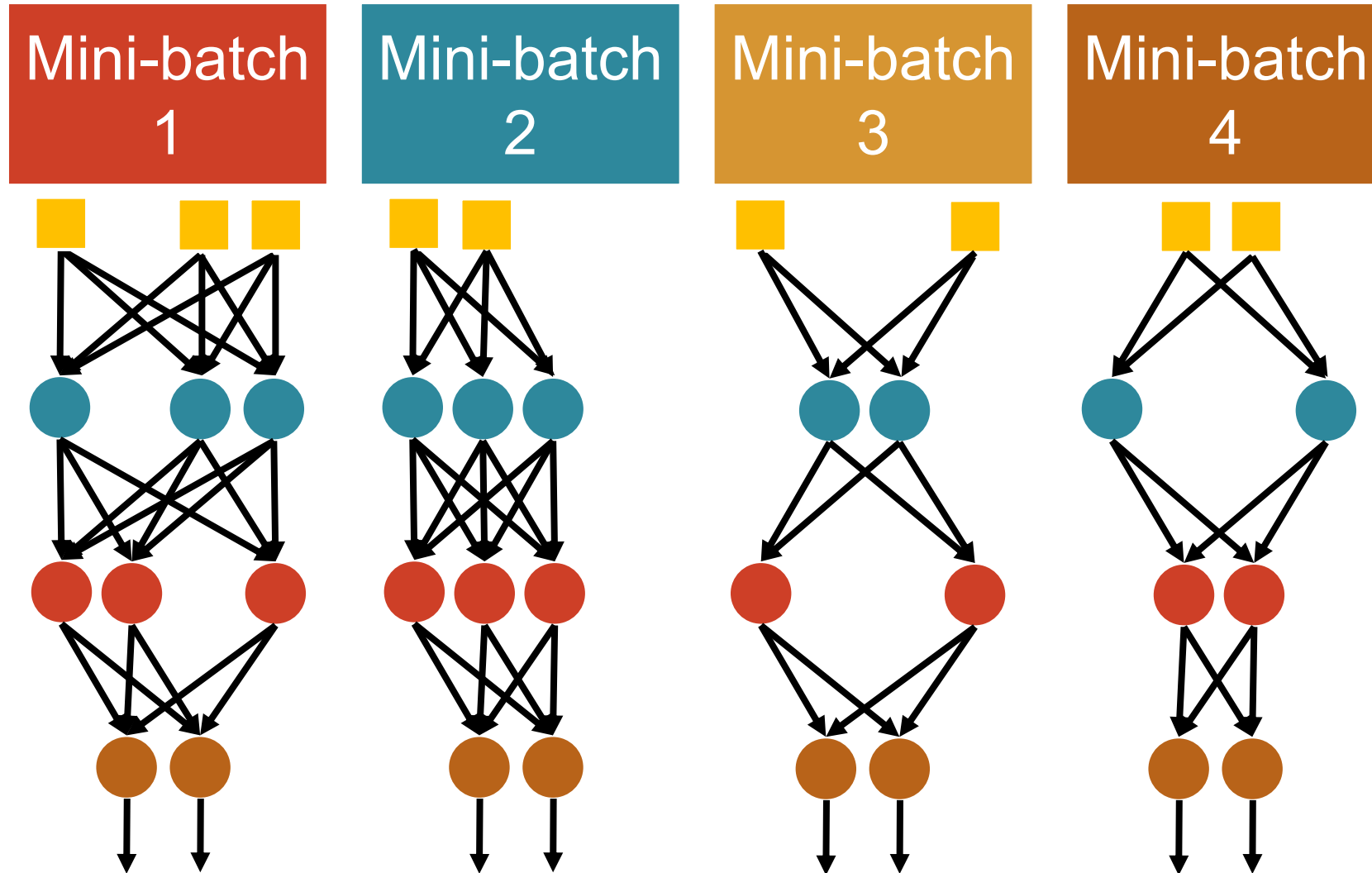
## 3. Tricks from testing process

### 3.3 Dropout -- A kind of ensemble



### 3. Tricks from testing process

#### 3.3 Dropout -- A kind of ensemble

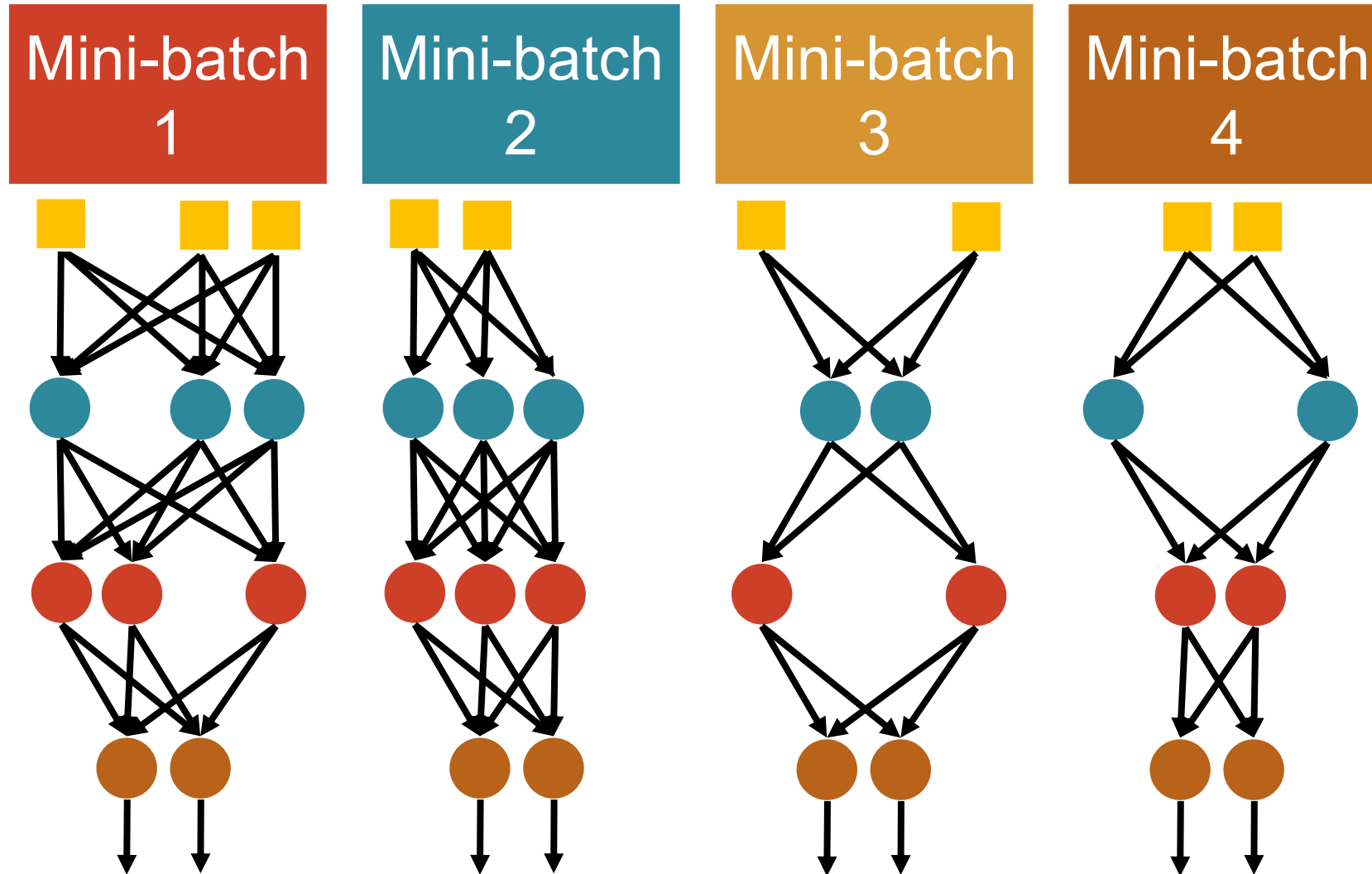


#### Training of dropout

Assume  $M$  neurons;  
Can dropout or not;  
 $2^M$  possible networks.

# 3. Tricks from testing process

## 3.3 Dropout -- A kind of ensemble

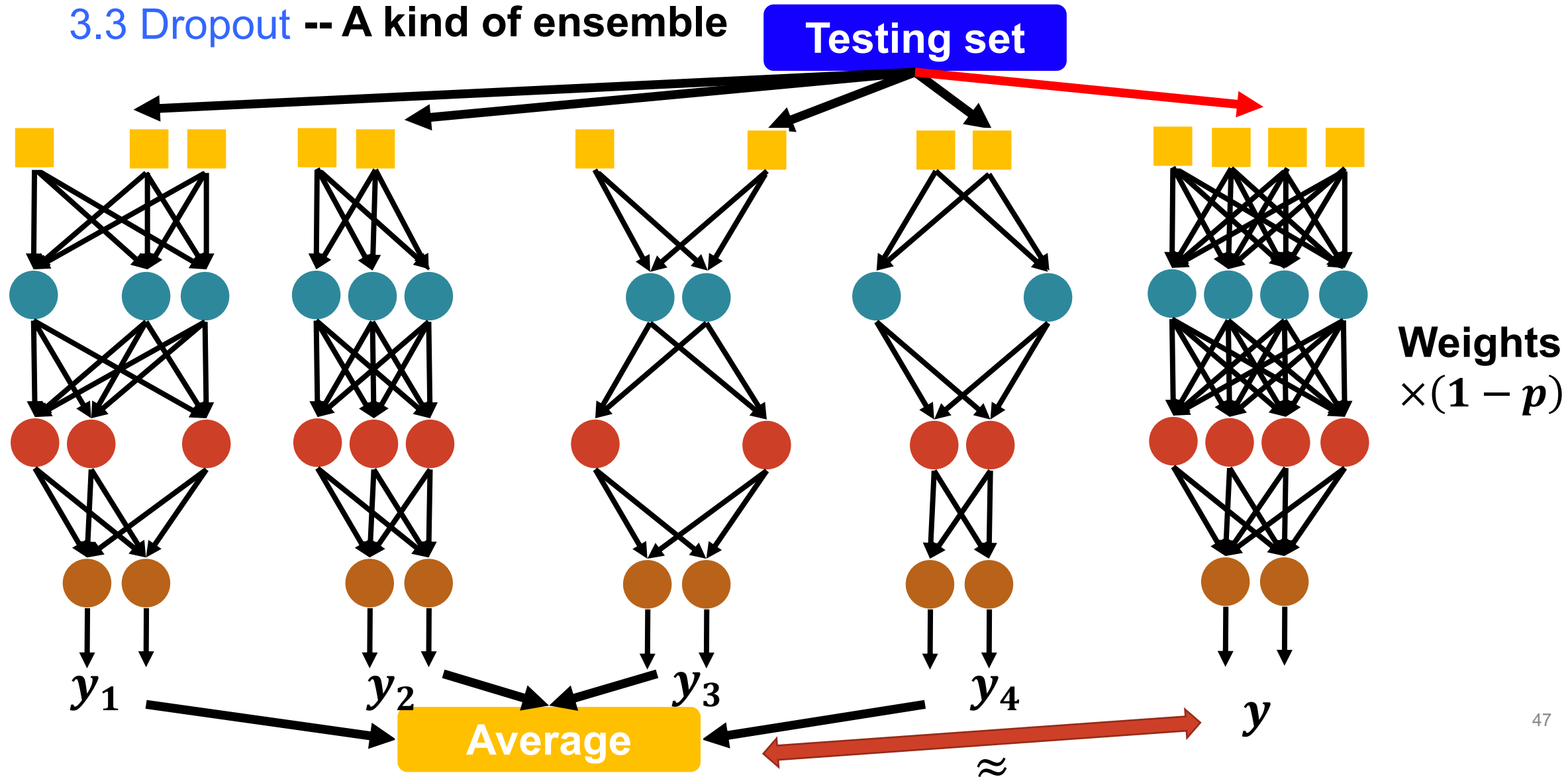


### Training of dropout

- Using one mini-batch to train one network;
- Some parameters in the network are shared.

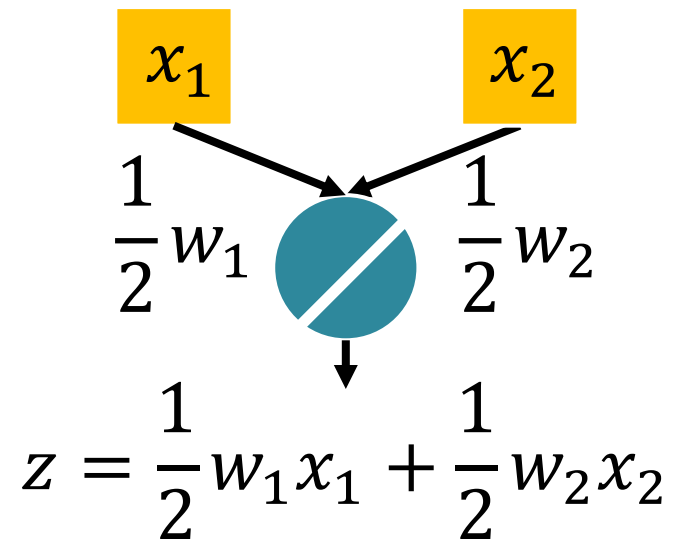
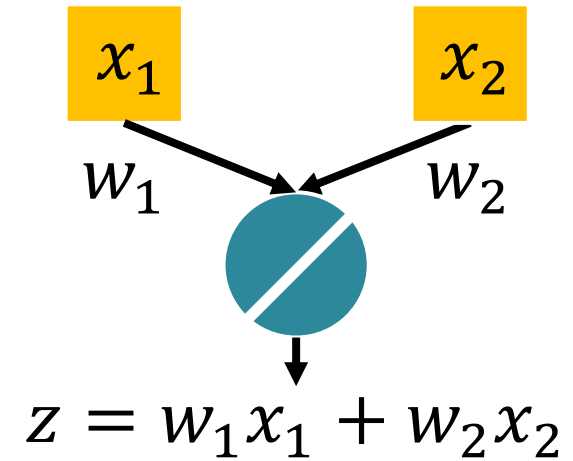
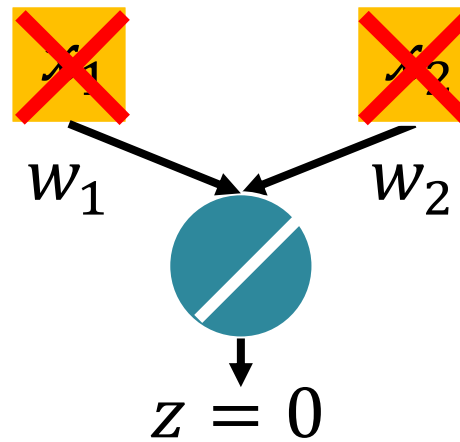
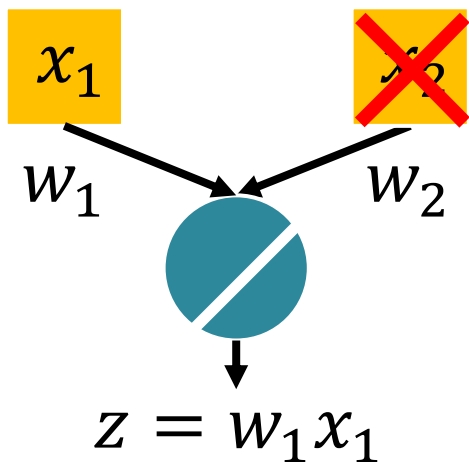
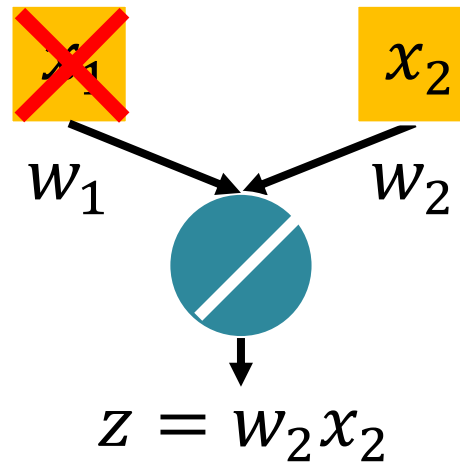
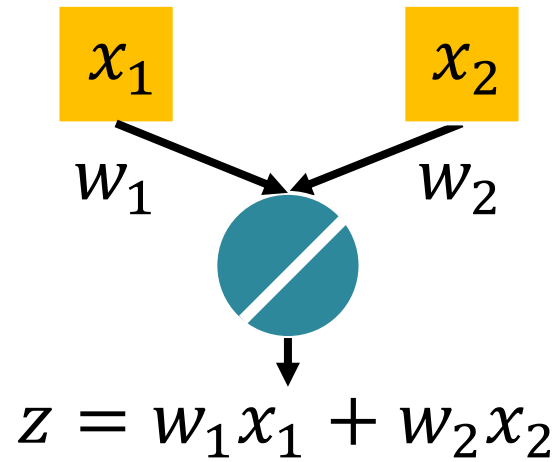
### 3. Tricks from testing process

#### 3.3 Dropout -- A kind of ensemble



### 3. Tricks from testing process

#### 3.3 Dropout -- A kind of ensemble





# Tricks for Training Artificial Neural Networks

Hao, Qi

School of Astronomy and Space Science

# THANKS

