

# **CSE463: Neural Networks**

## **Deep Neural Networks**

**by:**

**Hossam Abd El Munim**

**Computer & Systems Engineering Dept.,**

**Ain Shams University,**

**1 El-Sarayut Street, Abbassia, Cairo 11517**

# Deep learning attracts lots of attention.

- Google Trends

Deep learning obtains many exciting results.

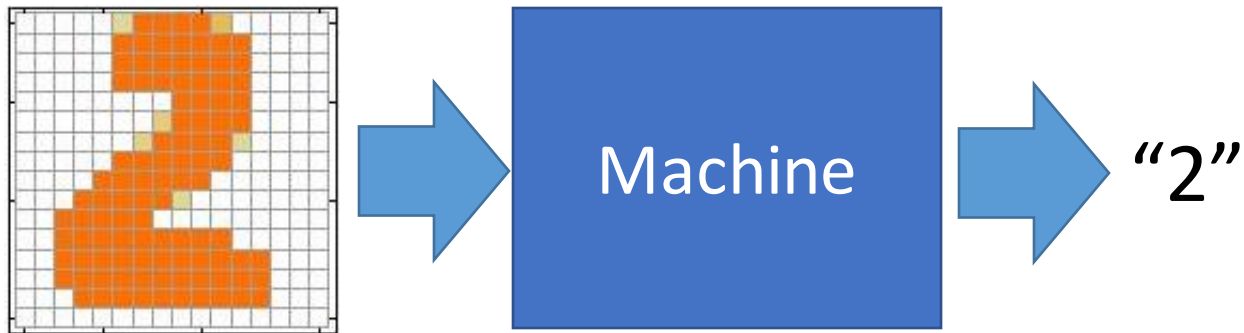


Interest over time 



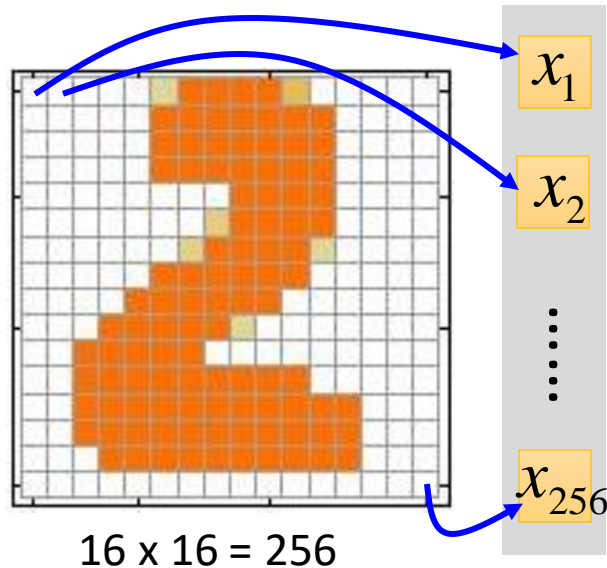
# Example Application

- Handwriting Digit Recognition



# Handwriting Digit Recognition

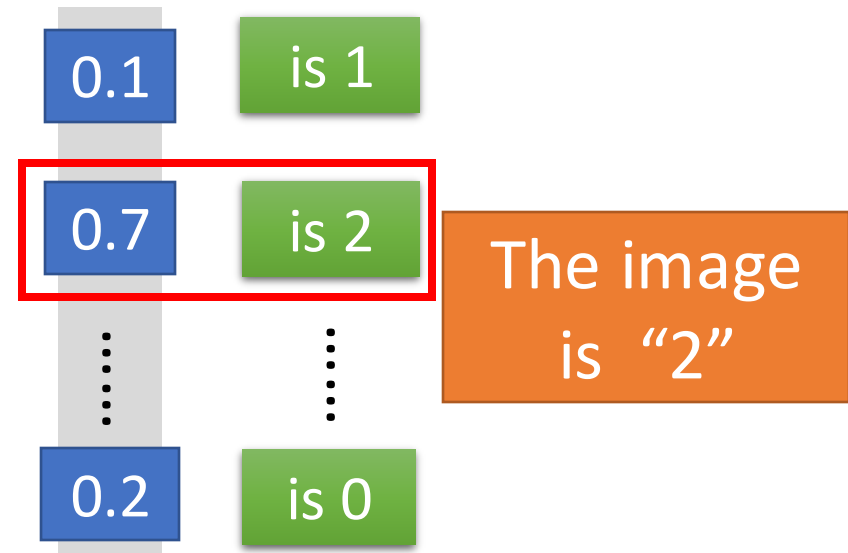
## Input



Ink  $\rightarrow$  1

No ink  $\rightarrow$  0

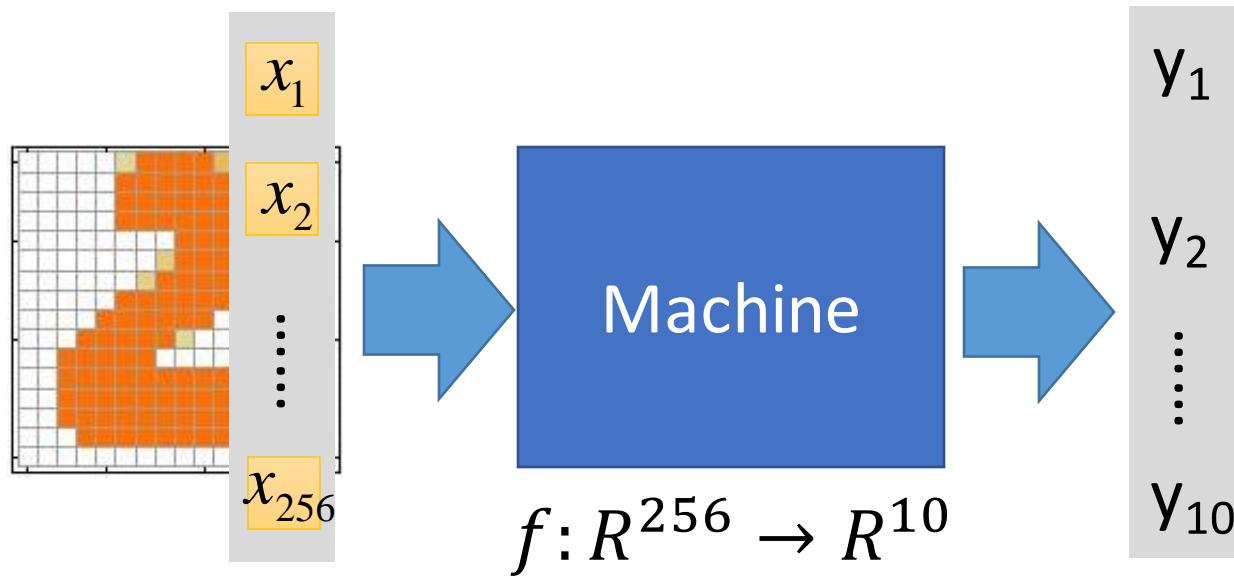
## Output



Each dimension represents the confidence of a digit.

# Example Application

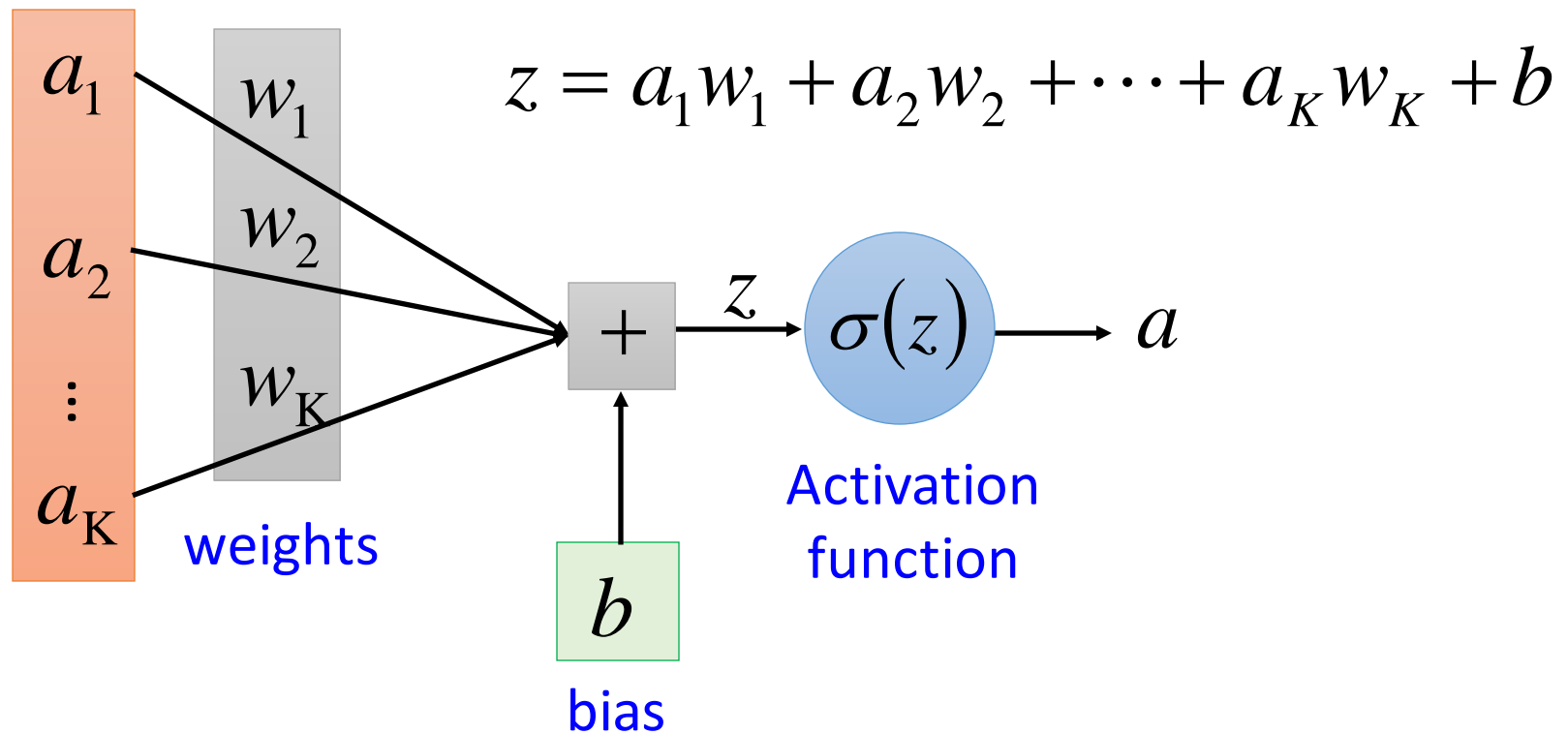
- Handwriting Digit Recognition



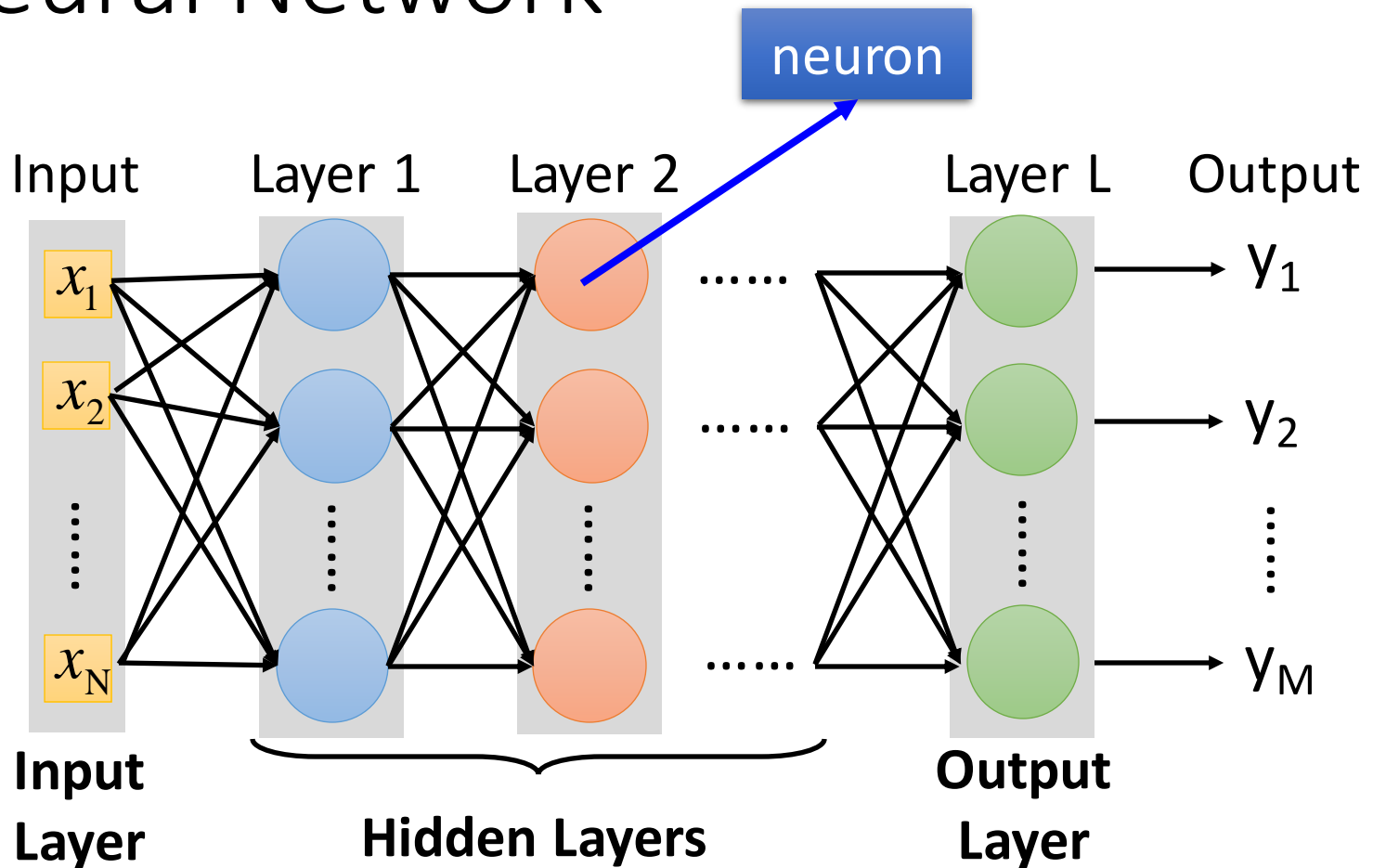
In deep learning, the function  $f$  is represented by neural network

# Element of Neural Network

**Neuron**  $f: R^K \rightarrow R$

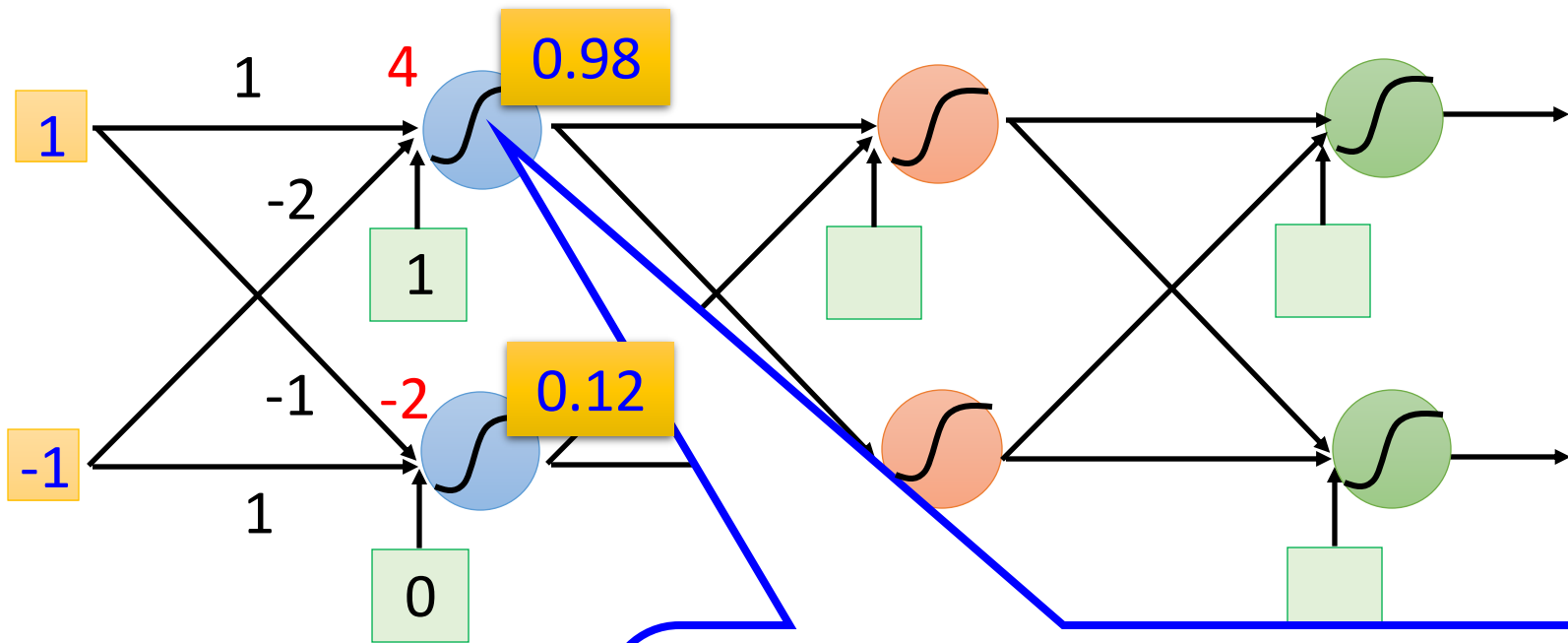


# Neural Network



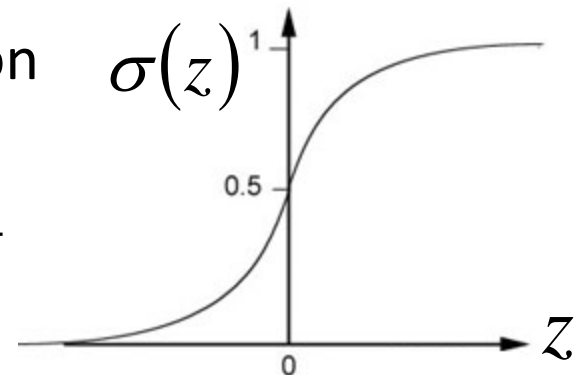
Deep means many hidden layers

# Example of Neural Network



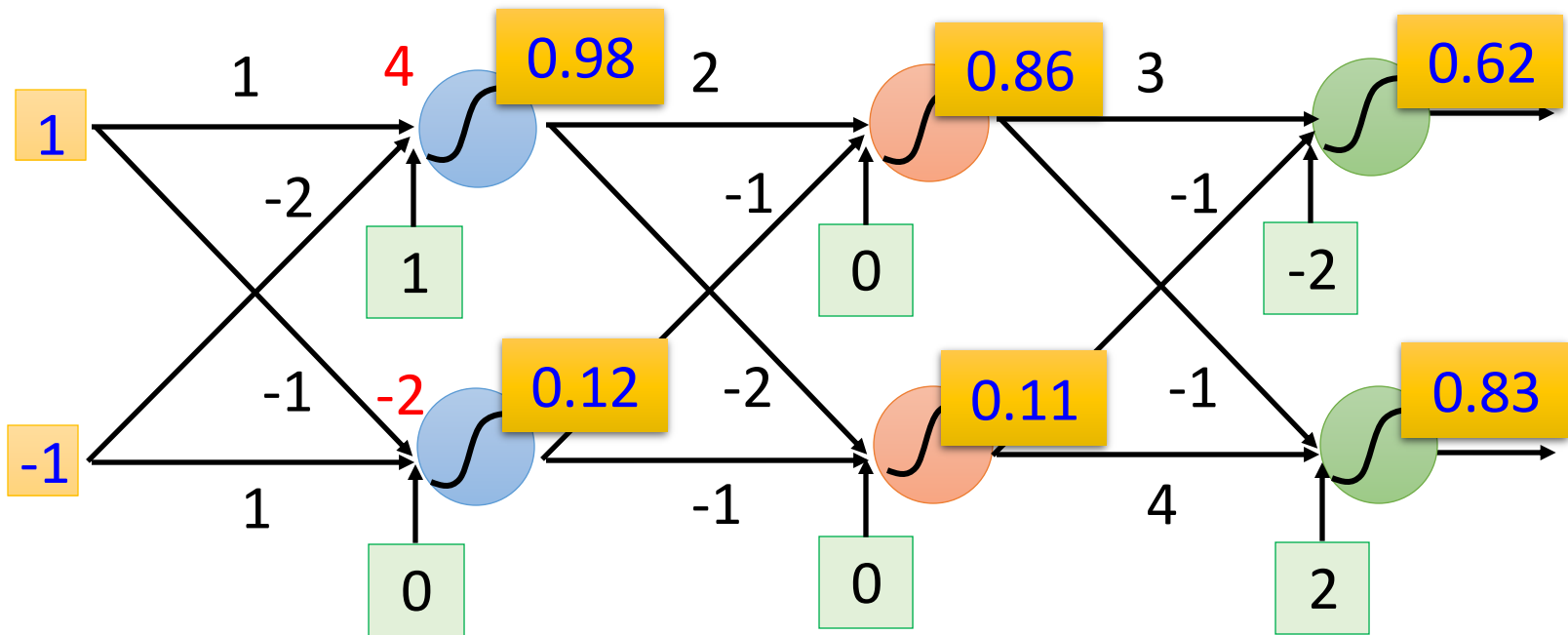
Sigmoid Function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

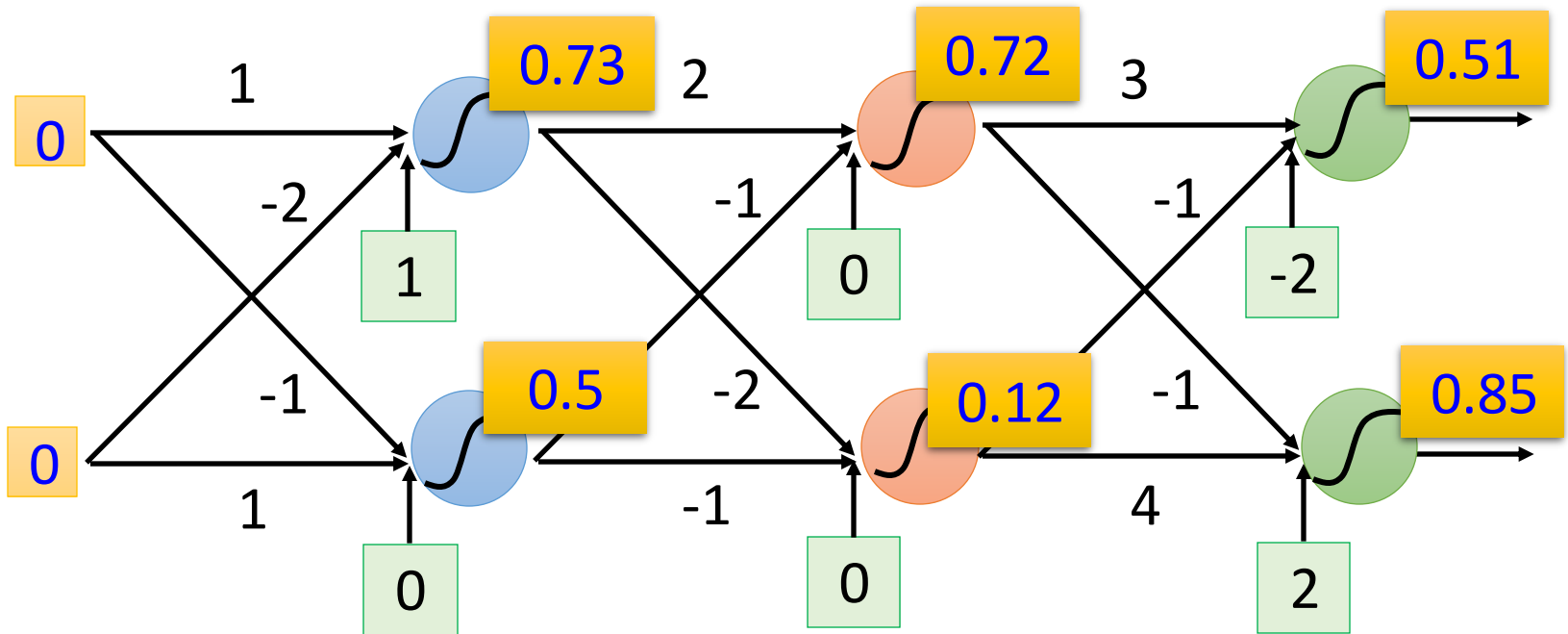




# Example of Neural Network



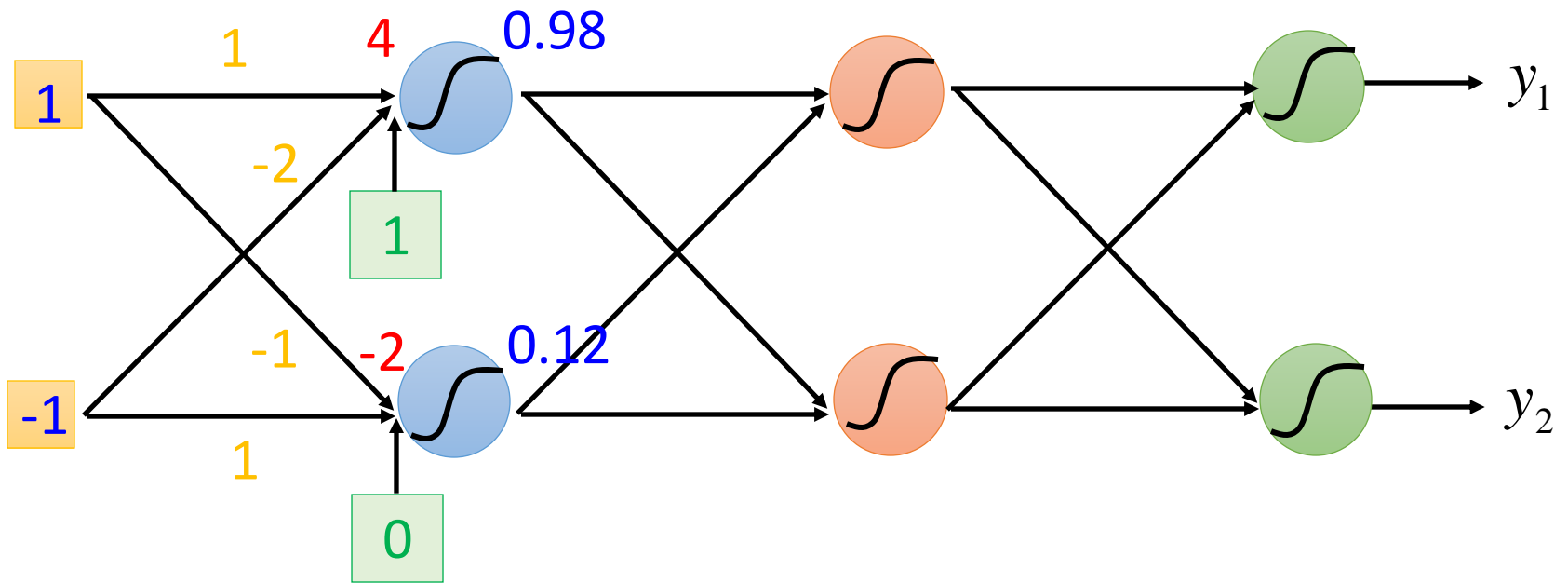
# Example of Neural Network



$$f: \mathbb{R}^2 \rightarrow \mathbb{R}^2 \quad f\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 0.62 \\ 0.83 \end{bmatrix} \quad f\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} 0.51 \\ 0.85 \end{bmatrix}$$

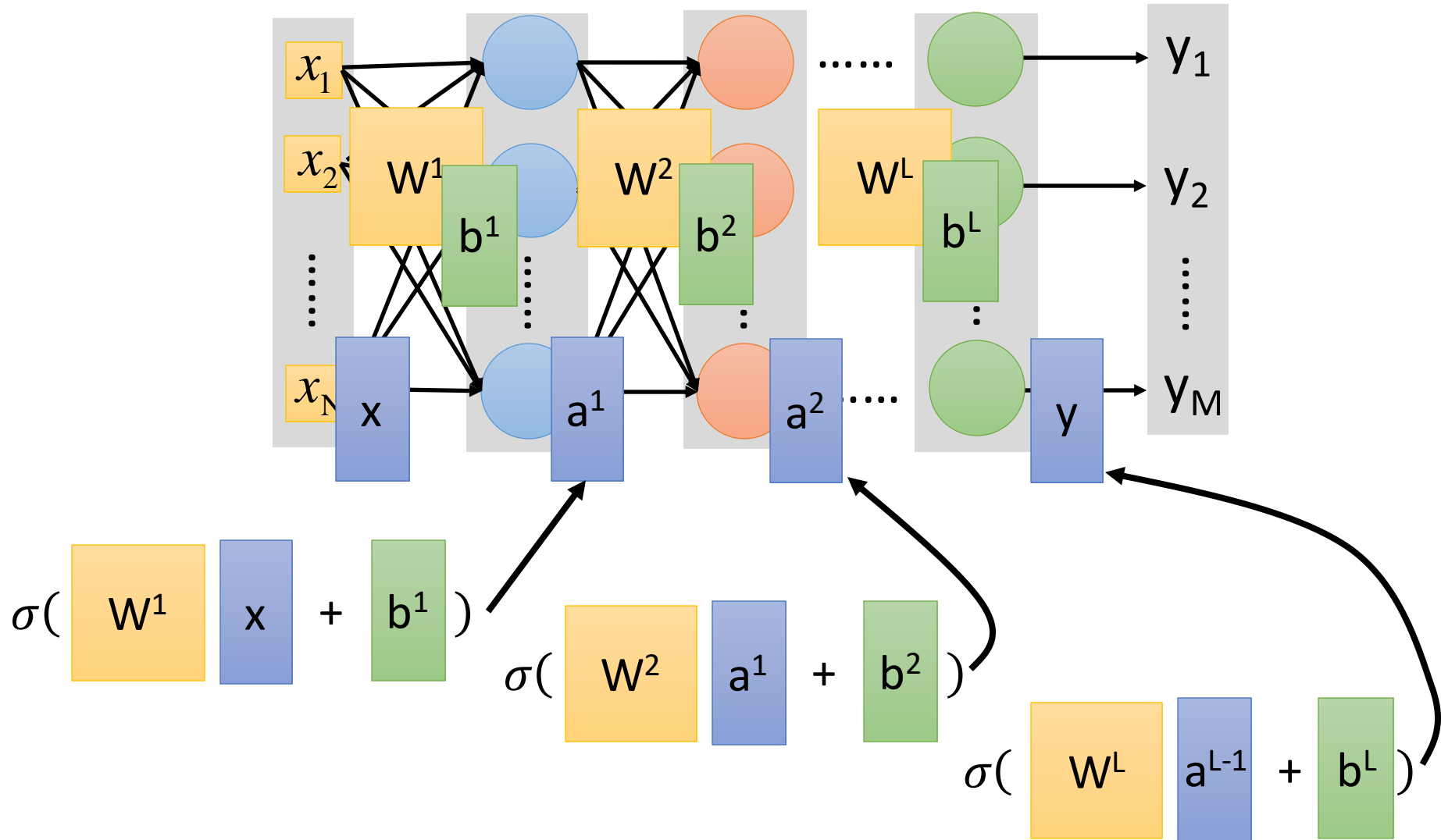
Different parameters define different function

# Matrix Operation

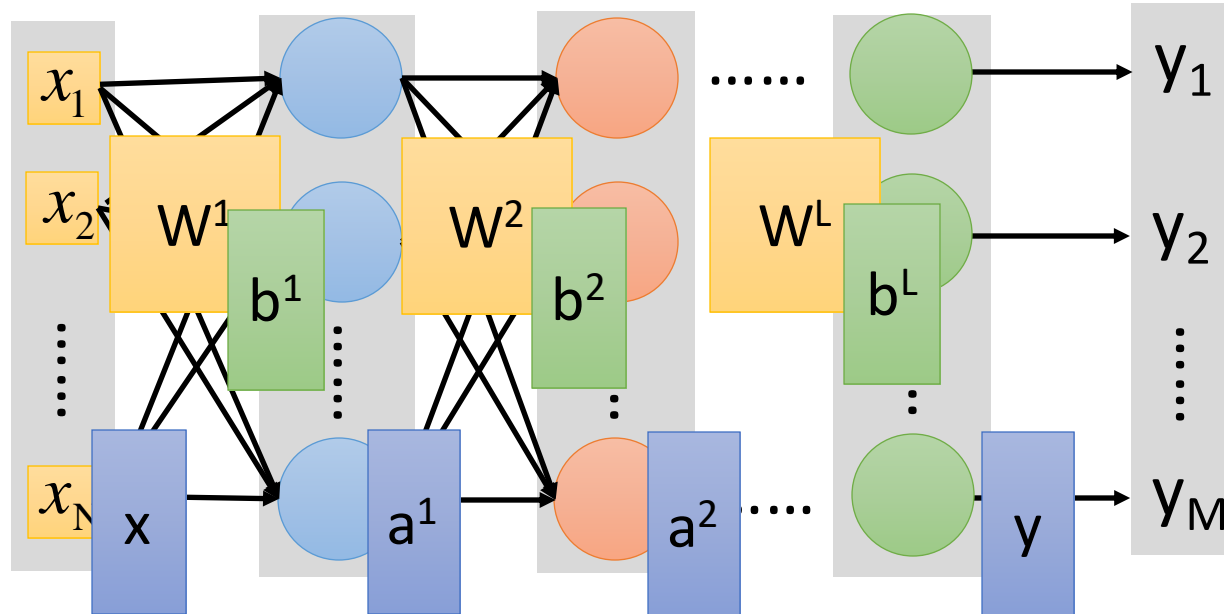


$$\sigma\left( \underbrace{\begin{bmatrix} 1 & -2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}}_{\begin{bmatrix} 4 \\ -2 \end{bmatrix}} \right) = \begin{bmatrix} 0.98 \\ 0.12 \end{bmatrix}$$

# Neural Network



# Neural Network



$$\mathbf{y} = f(\mathbf{x})$$

Using parallel computing techniques  
to speed up matrix operation

$$= \sigma(\mathbf{W}^L \dots \sigma(\mathbf{W}^2 \sigma(\mathbf{W}^1 \mathbf{x} + \mathbf{b}^1) + \mathbf{b}^2) \dots + \mathbf{b}^L)$$

# Softmax

- Softmax layer as the output layer

## Ordinary Layer

$$z_1 \longrightarrow \sigma \longrightarrow y_1 = \sigma(z_1)$$

$$z_2 \longrightarrow \sigma \longrightarrow y_2 = \sigma(z_2)$$

$$z_3 \longrightarrow \sigma \longrightarrow y_3 = \sigma(z_3)$$

In general, the output of network can be any value.

May not be easy to interpret

# Softmax

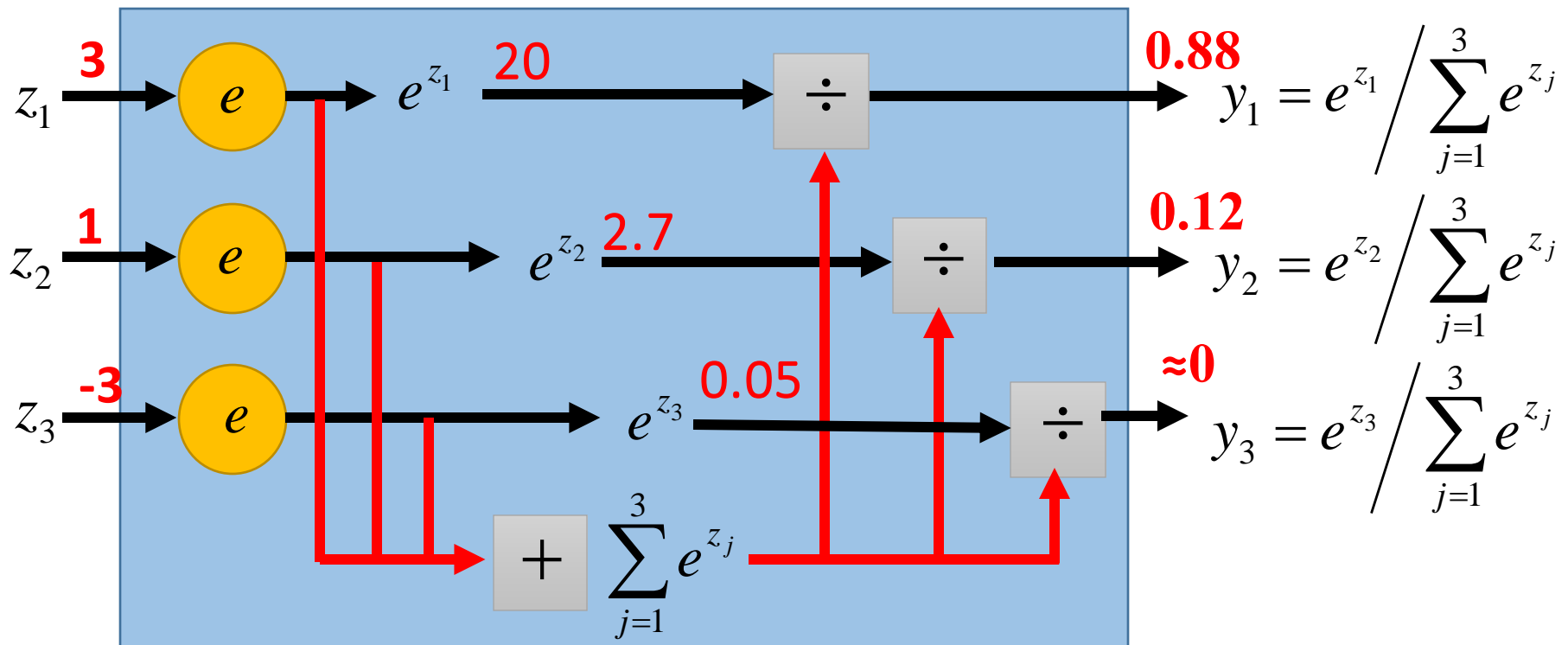
- Softmax layer as the output layer

**Probability:**

■  $1 > y_i > 0$

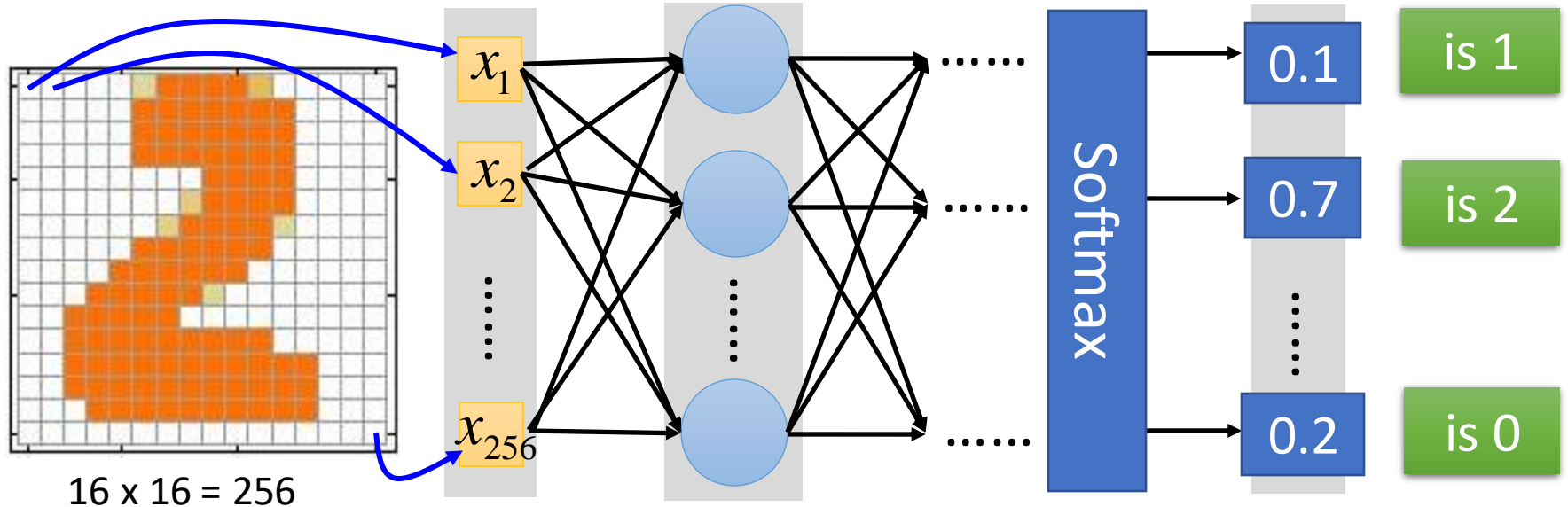
■  $\sum_i y_i = 1$

## Softmax Layer



# How to set network parameters

$$\theta = \{W^1, b^1, W^2, b^2, \dots, W^L, b^L\}$$



$16 \times 16 = 256$

Ink  $\rightarrow 1$

No ink  $\rightarrow 0$

Set the network parameters  $\theta$  such that .....

Input:  ... value

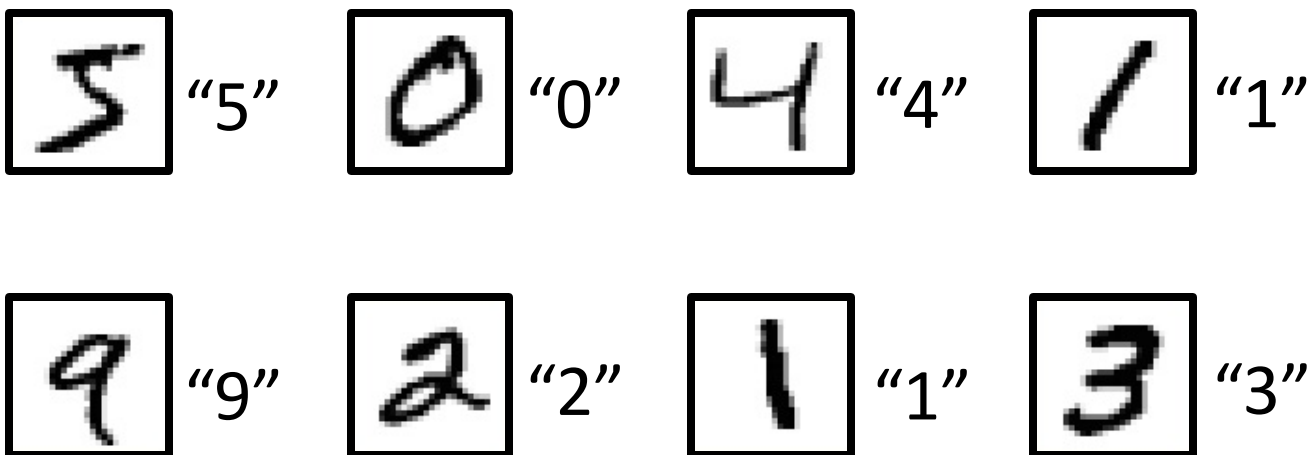
Input:   $y_2$  has the maximum value

How to let the neural  
network achieve this



# Training Data

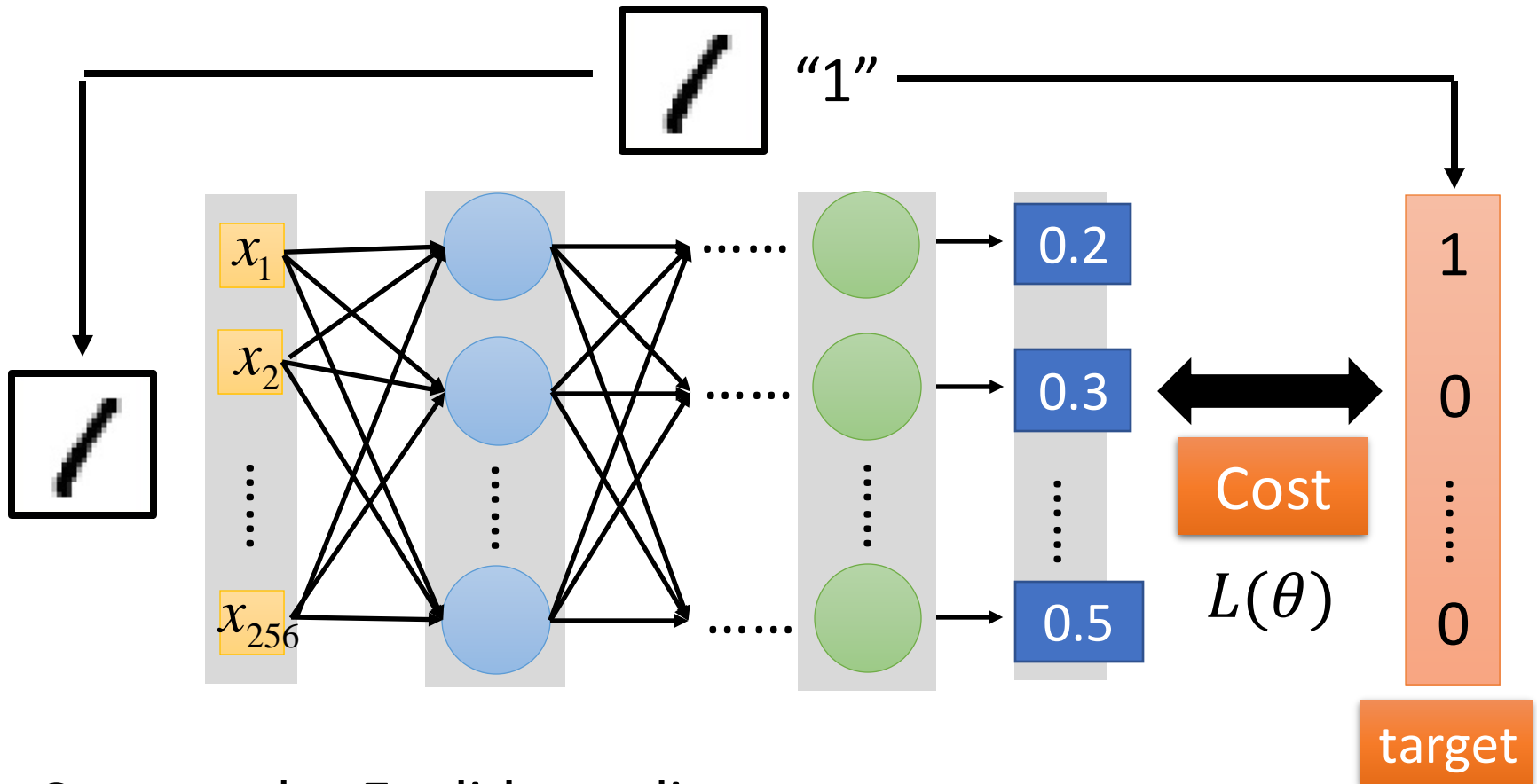
- Preparing training data: images and their labels



Using the training data to find  
the network parameters.

# Cost

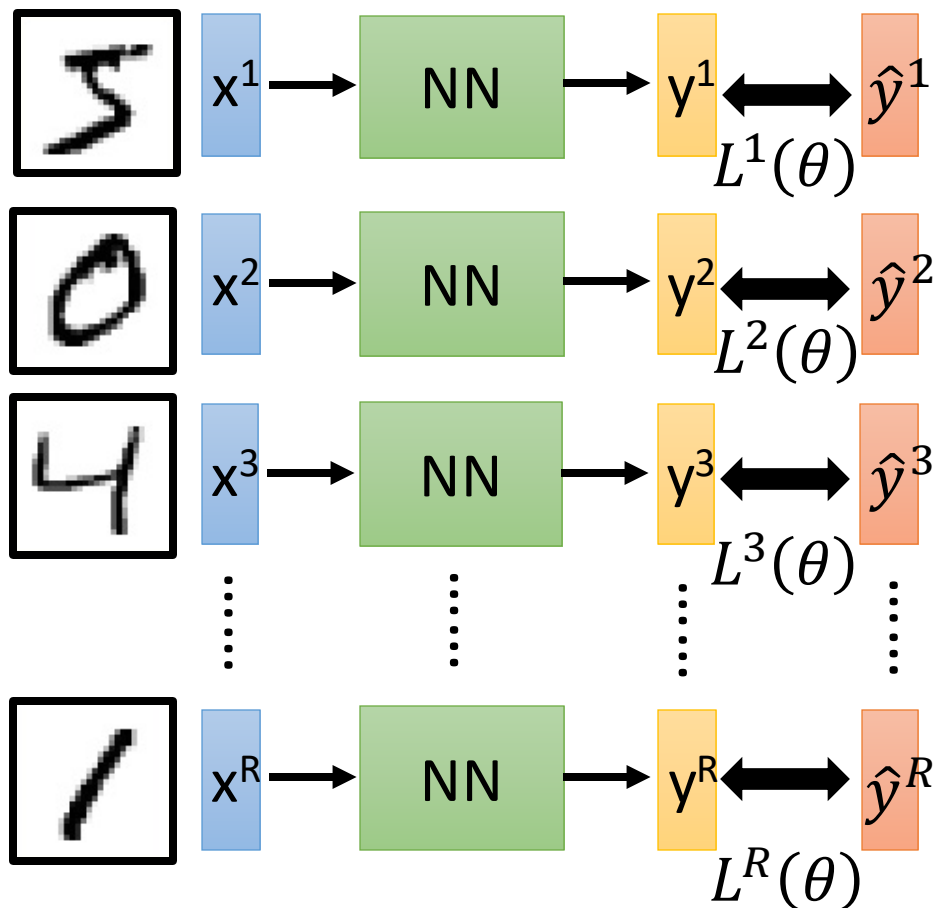
Given a set of network parameters  $\theta$ , each example has a cost value.



Cost can be Euclidean distance or cross entropy of the network output and target

# Total Cost

For all training data ...



Total Cost:

$$C(\theta) = \sum_{r=1}^R L^r(\theta)$$

How bad the network parameters  $\theta$  is on this task

Find the network parameters  $\theta^*$  that minimize this value

# Backpropagation

- A network can have millions of parameters.
  - Backpropagation is the way to compute the gradients efficiently (not today)
  - Ref:  
[http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS\\_2015\\_2/Lecture/DNN%20backprop.ecm.mp4/index.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/index.html)
- Many toolkits can compute the gradients automatically

theano



Ref:

[http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS\\_2015\\_2/Lecture/Theano%20DNN.ecm.mp4/index.html](http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/Theano%20DNN.ecm.mp4/index.html)

Why Deep?

# Deeper is Better?

Layer X Size	Word Error Rate (%)
1 X 2k	24.2
2 X 2k	20.4
3 X 2k	18.4
4 X 2k	17.8
5 X 2k	17.2
7 X 2k	17.1

Not surprised, more parameters, better performance

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

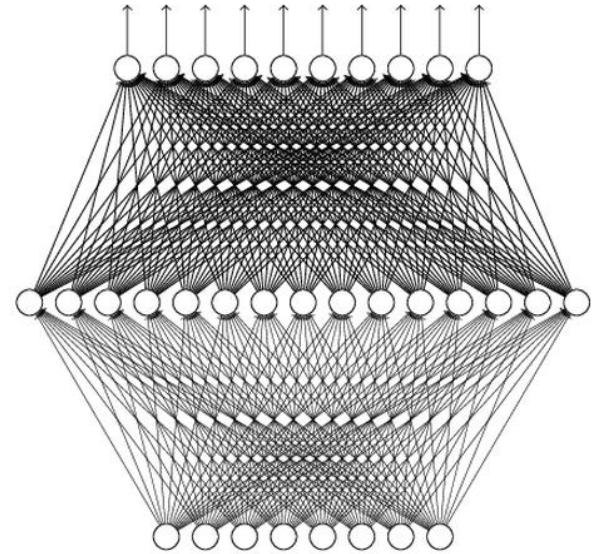
# Universality Theorem

Any continuous function  $f$

$$f : R^N \rightarrow R^M$$

Can be realized by a network  
with one hidden layer

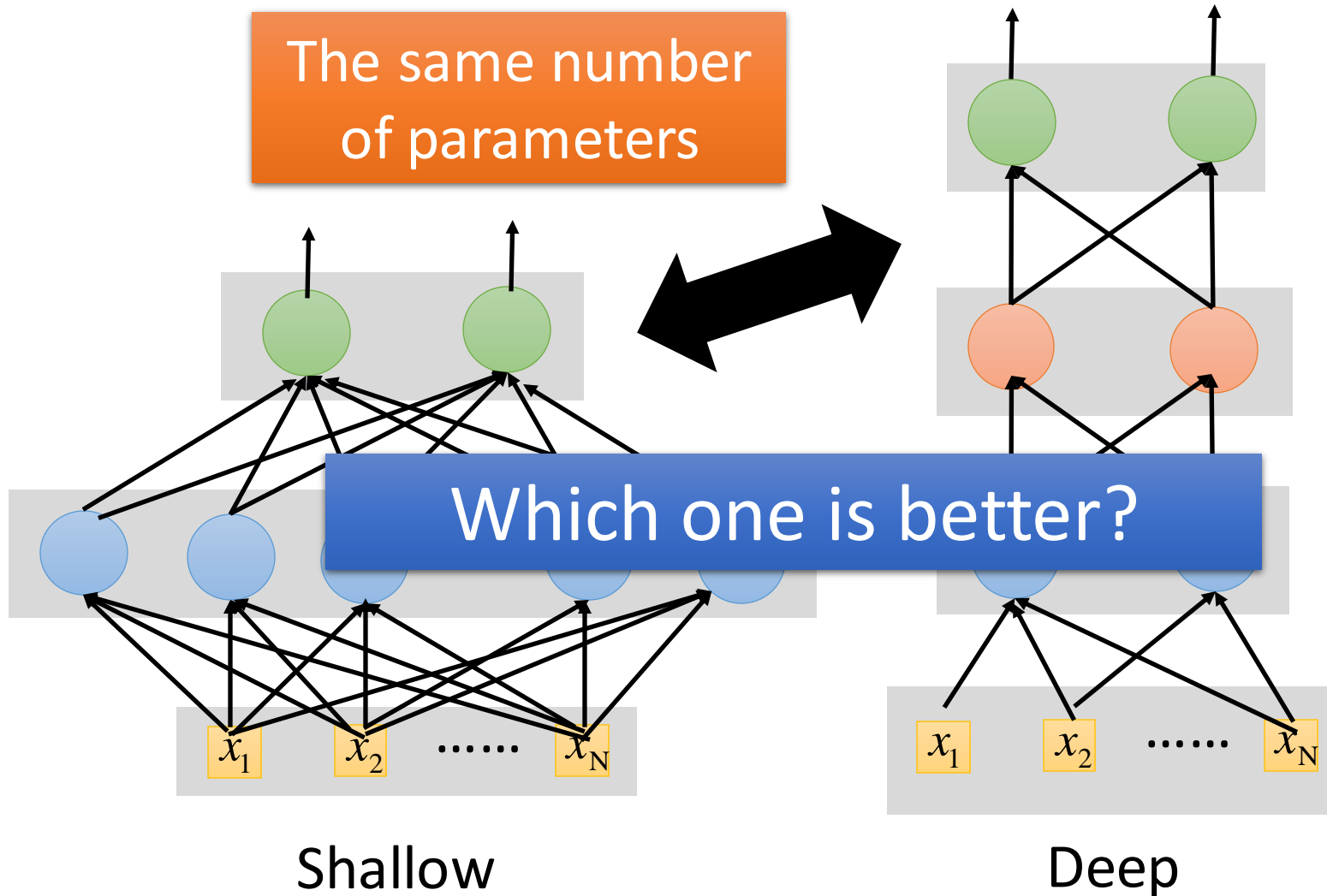
(given **enough** hidden  
neurons)



Reference for the reason:  
<http://neuralnetworksanddeeplearning.com/chap4.html>

Why “Deep” neural network not “Fat” neural network?

# Fat + Short v.s. Thin + Tall





# Fat + Short v.s. Thin + Tall

Layer X Size	Word Error Rate (%)	Layer X Size	Word Error Rate (%)
1 X 2k	24.2		
2 X 2k	20.4		
3 X 2k	18.4		
4 X 2k	17.8		
5 X 2k	17.2	1 X 3772	22.5
7 X 2k	17.1	1 X 4634	22.6
		1 X 16k	22.1

Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

# Why Deep?

## Shallow Neural Network

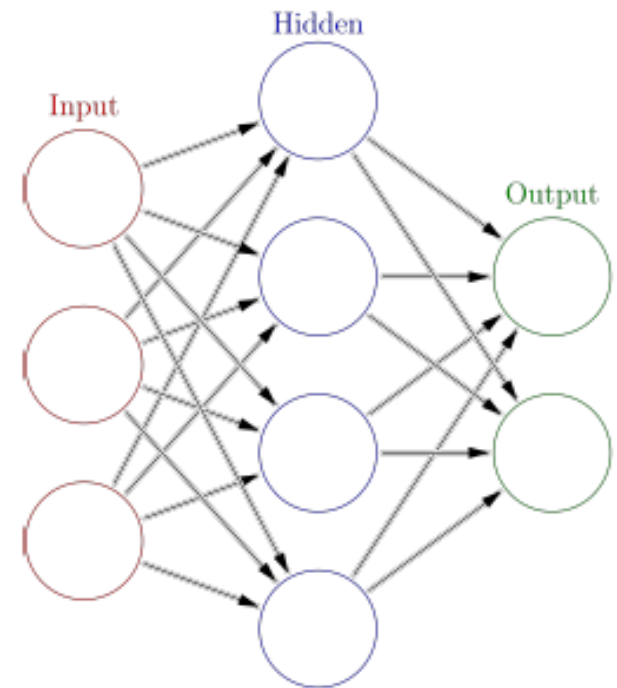
Shallow Neural Network is similar to most of the conventional supervised models

### Pros:

1. Easy to train and test
2. Able to approach any continuous function

### Cons:

1. Performance depends on well-designed features
2. Difficult to generalize the prediction



# Why Deep?

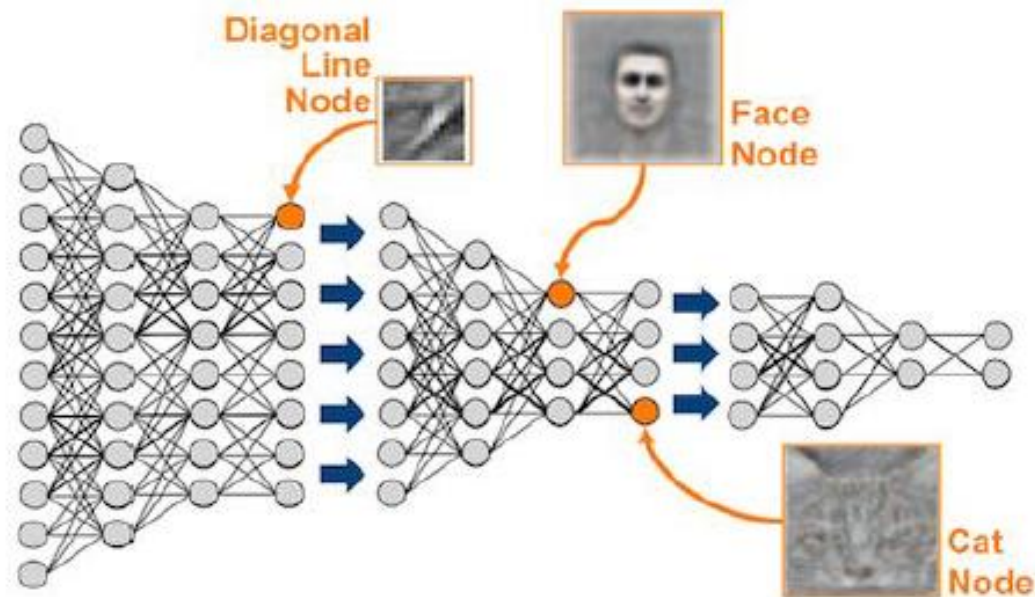
## Deep Neural Network

### Pros:

1. Automatically learn the High-level features representation
2. Modularity: DNN can be composed like LEGO bricks
3. Able to do **Transfer Learning**

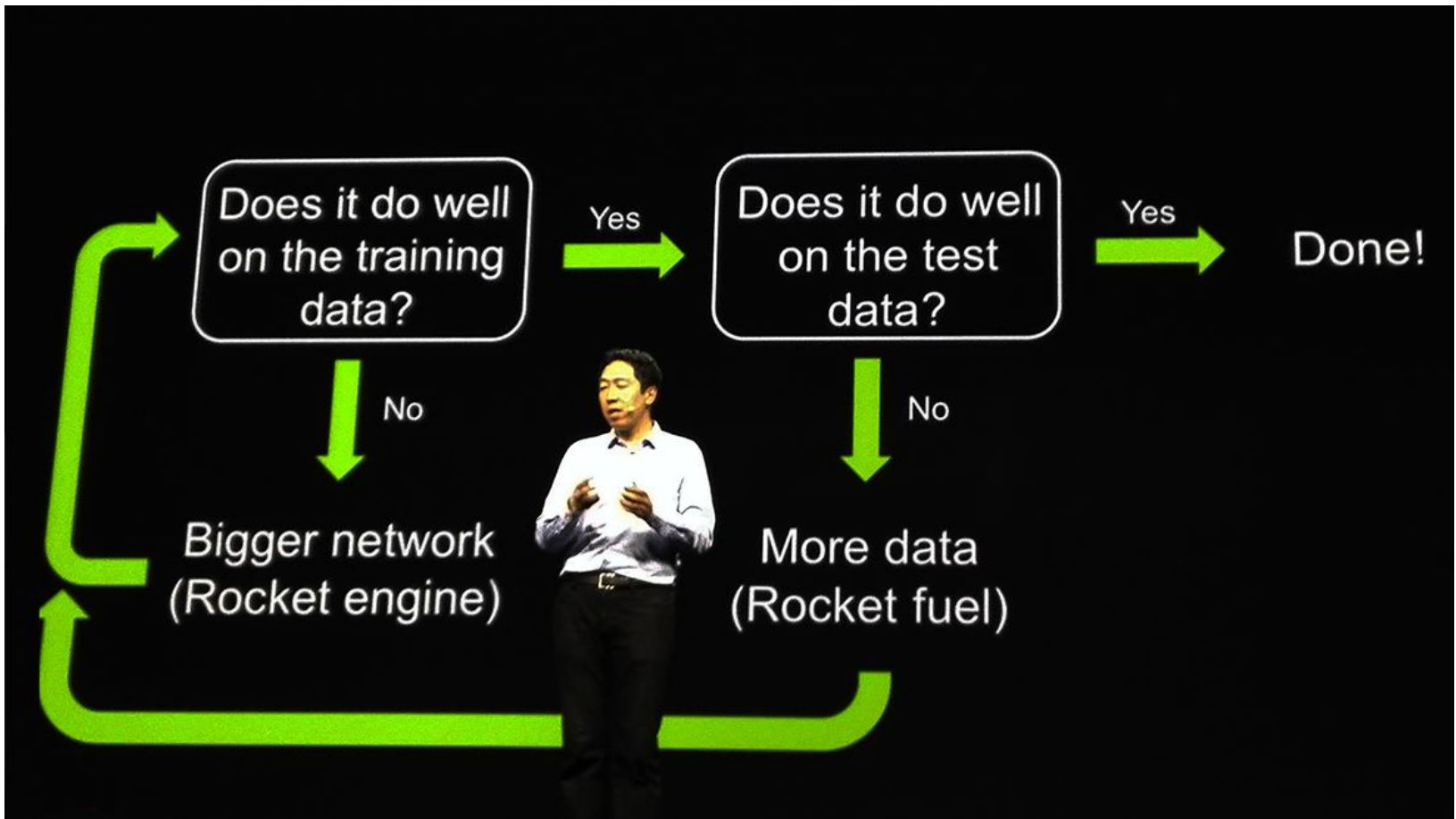
### Cons:

1. Requires tons of data for training
2. Expensive computation power for training and testing (no CV)



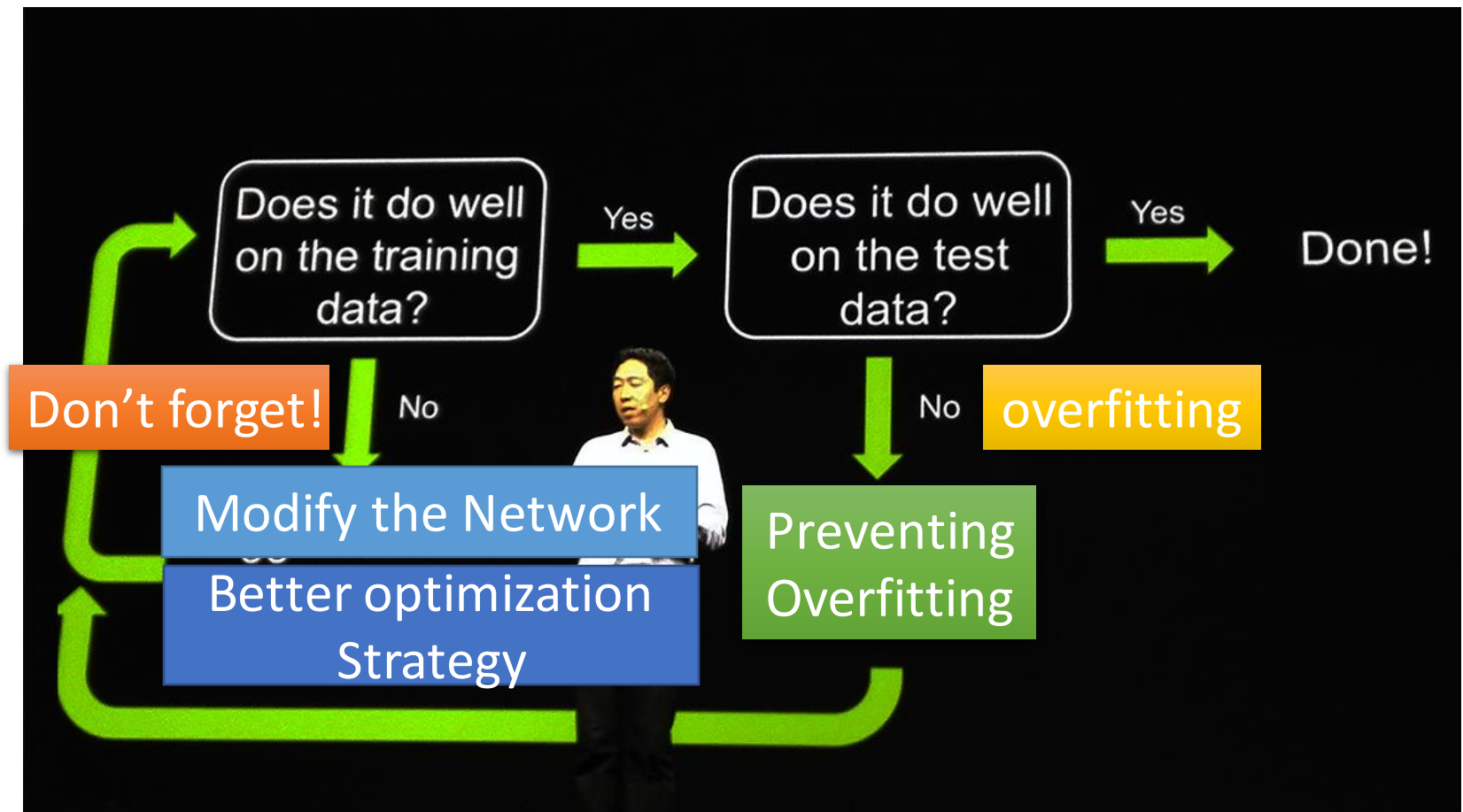
# Tips for Training DNN

# Recipe for Learning



<http://www.gizmodo.com.au/2015/04/the-basic-recipe-for-machine-learning-explained-in-a-single-powerpoint-slide/>

# Recipe for Learning



<http://www.gizmodo.com.au/2015/04/the-basic-recipe-for-machine-learning-explained-in-a-single-powerpoint-slide/>

# Recipe for Learning

## Modify the Network

- New activation functions, for example, ReLU or Maxout

## Better optimization Strategy

- Adaptive learning rates

## Prevent Overfitting

- Dropout

Only use this approach when you already obtained good results on the training data.

# Tips for Training DNN

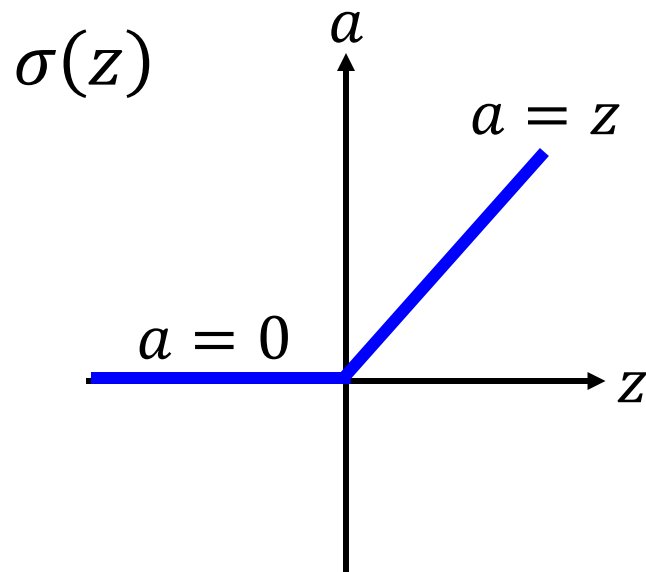
New Activation Function



# ReLU

- Rectified Linear Unit (ReLU)

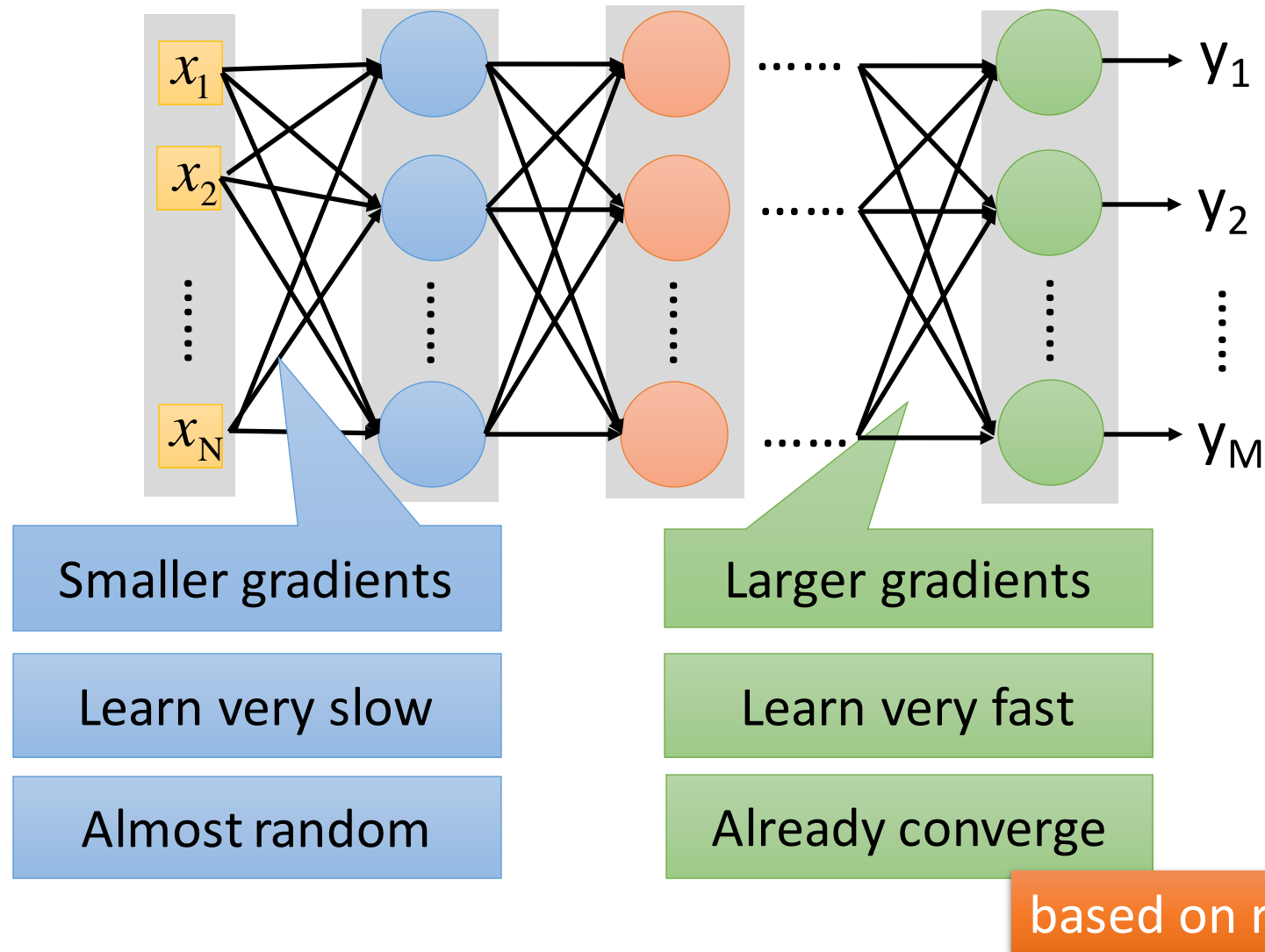
**Reason:**



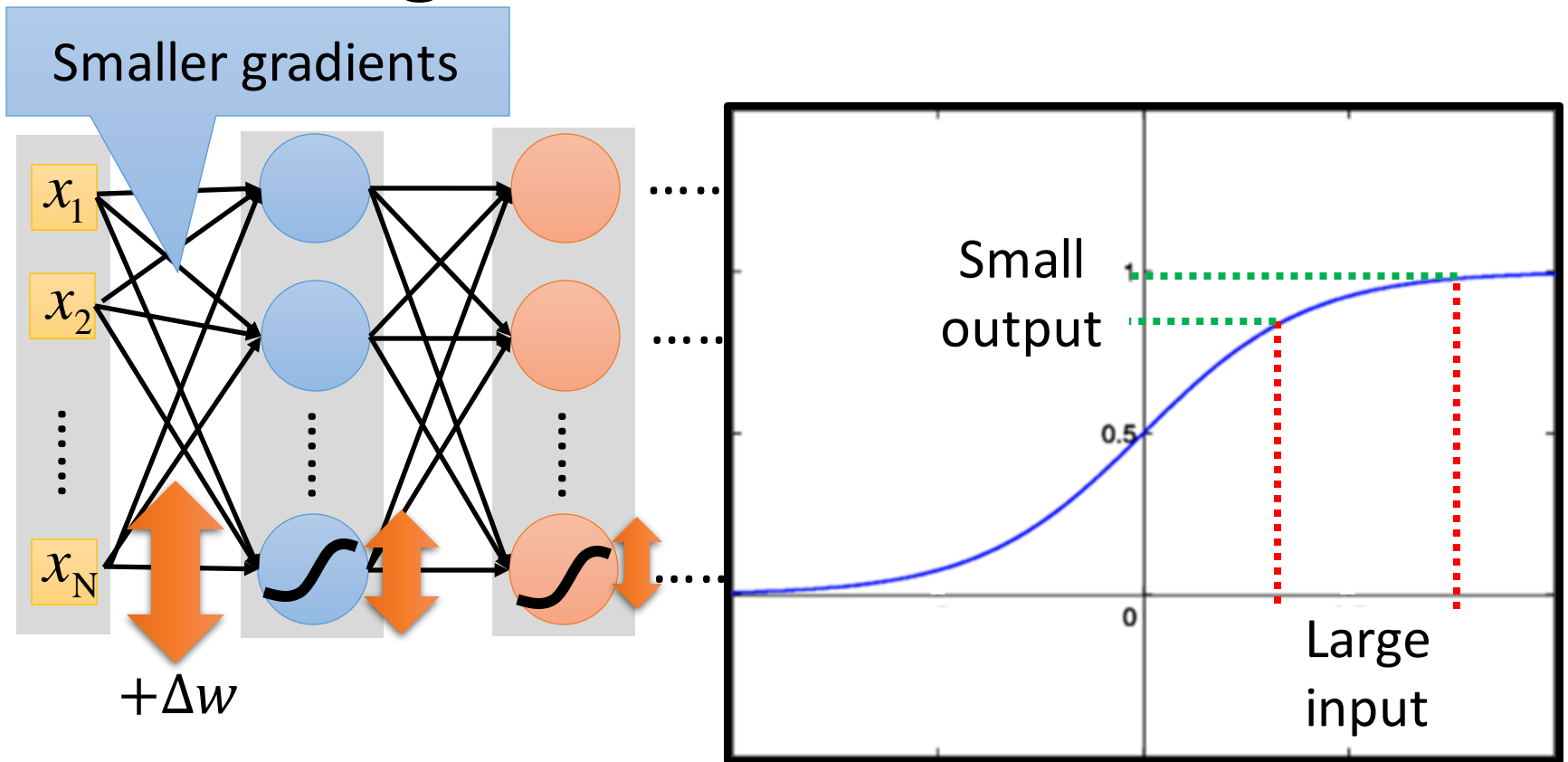
Vanishing gradient  
problem

[Xavier Glorot, AISTATS'11]  
[Andrew L. Maas, ICML'13]  
[Kaiming He, arXiv'15]

# Vanishing Gradient Problem



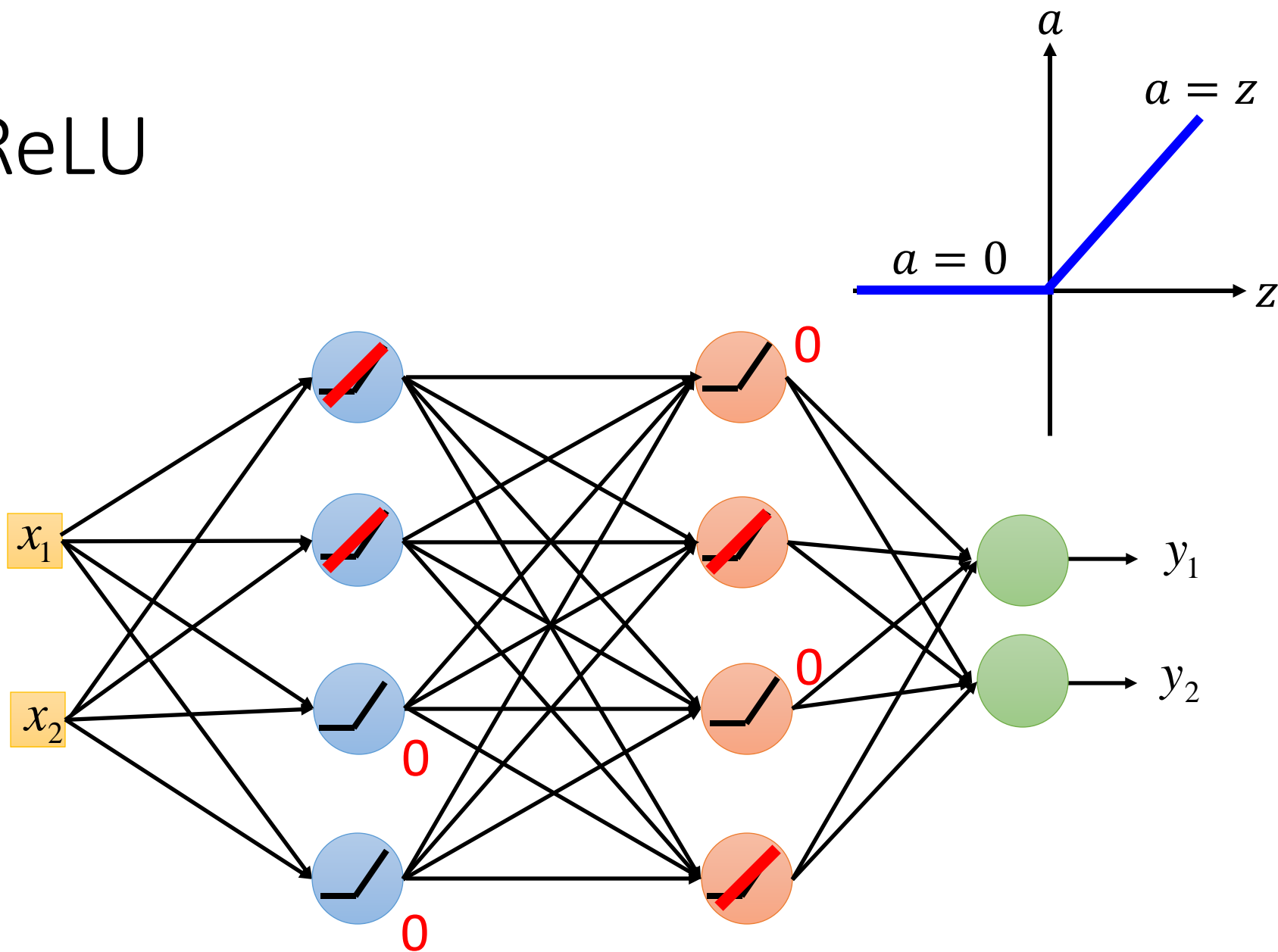
# Vanishing Gradient Problem



Intuitive way to compute the gradient ...

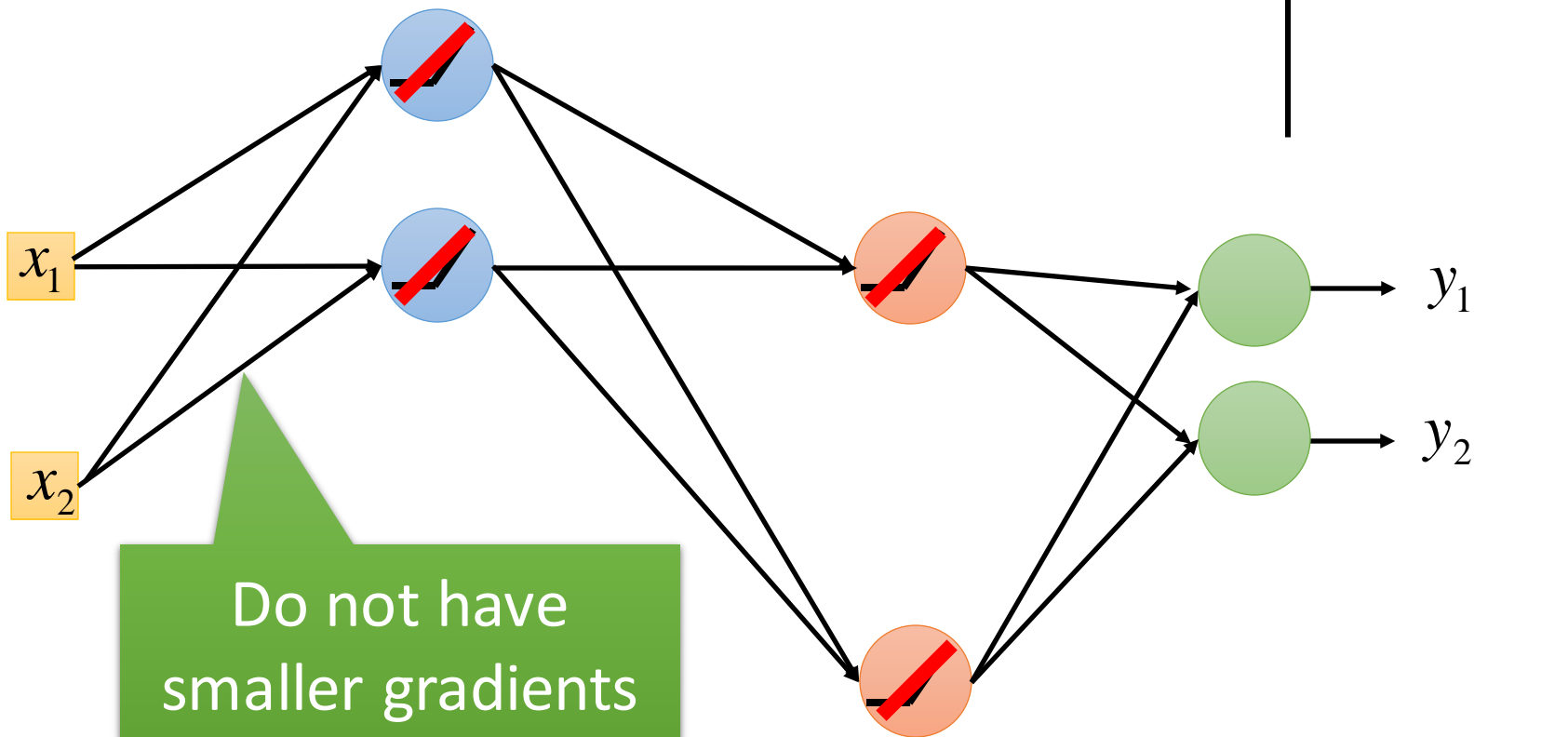
$$\frac{\partial C}{\partial w} = ? \quad \frac{\Delta C}{\Delta w}$$

# ReLU



# ReLU

A Thinner linear network



# Not the whole story .....

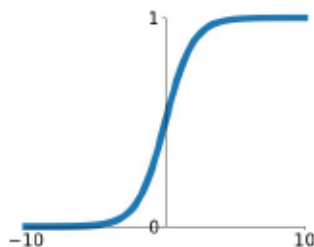
- Adagrad [John Duchi, JMLR'11]
- RMSprop
  - <https://www.youtube.com/watch?v=O3sxAc4hxZU>
- Adadelata [Matthew D. Zeiler, arXiv'12]
- Adam [Diederik P. Kingma, ICLR'15]
- AdaSecant [Caglar Gulcehre, arXiv'14]
- “No more pesky learning rates” [Tom Schaul, arXiv'12]

# Notes:

## Last time: Activation Functions

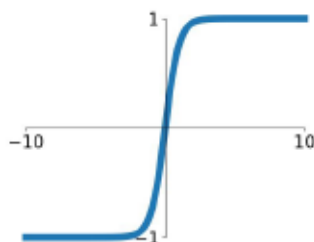
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



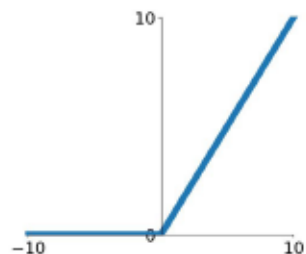
**tanh**

$$\tanh(x)$$



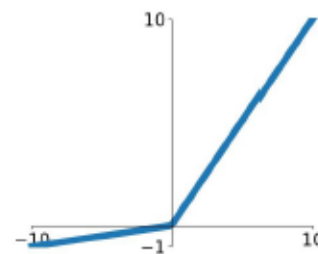
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0.1x, x)$$

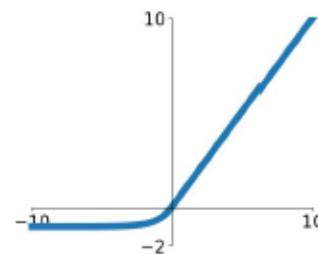


**Maxout**

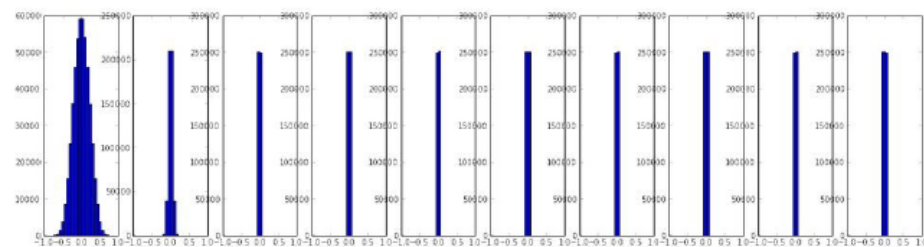
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

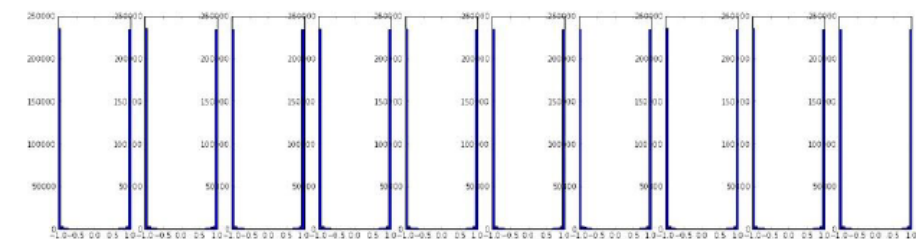


# Last time: Weight Initialization



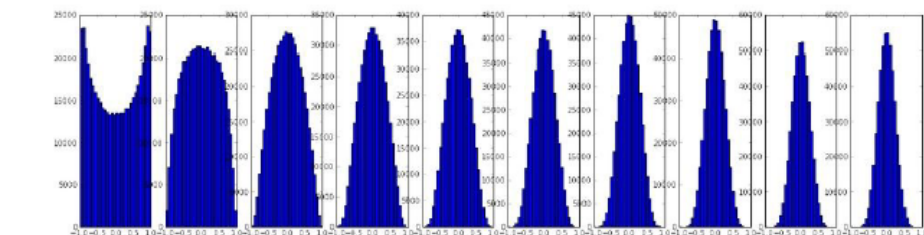
**Initialization too small:**

Activations go to zero, gradients also zero,  
No learning



**Initialization too big:**

Activations saturate (for tanh),  
Gradients zero, no learning

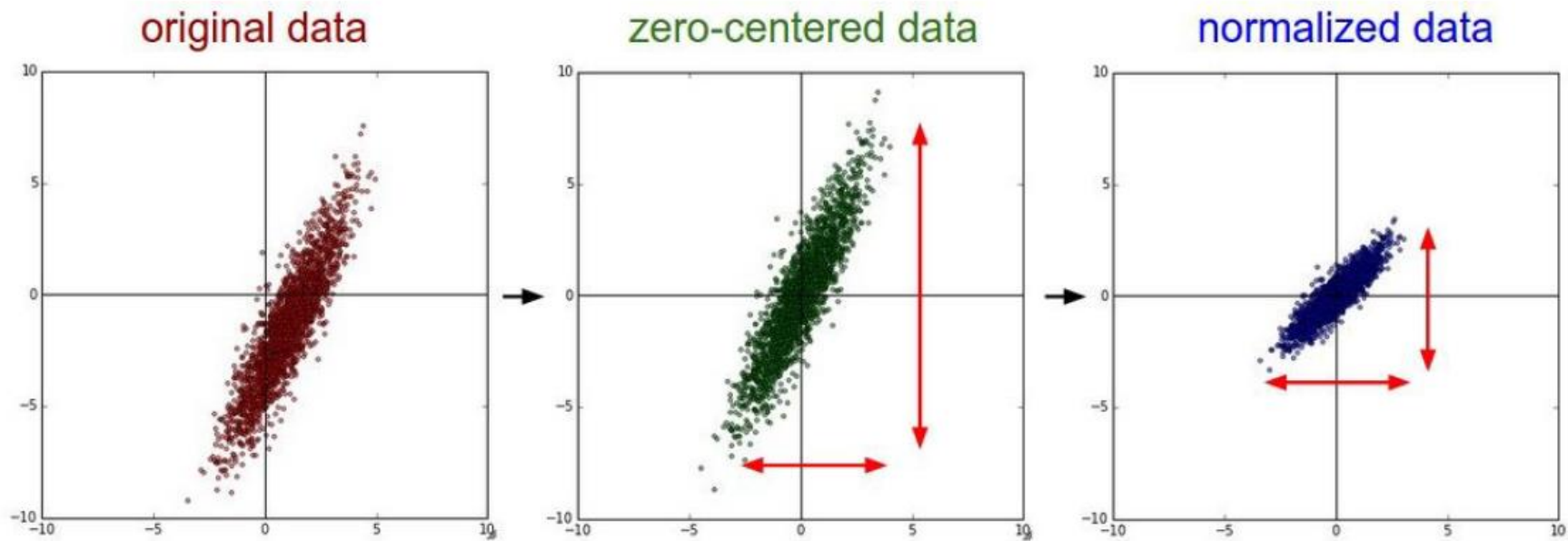


**Initialization just right:**

Nice distribution of activations at all layers,  
Learning proceeds nicely

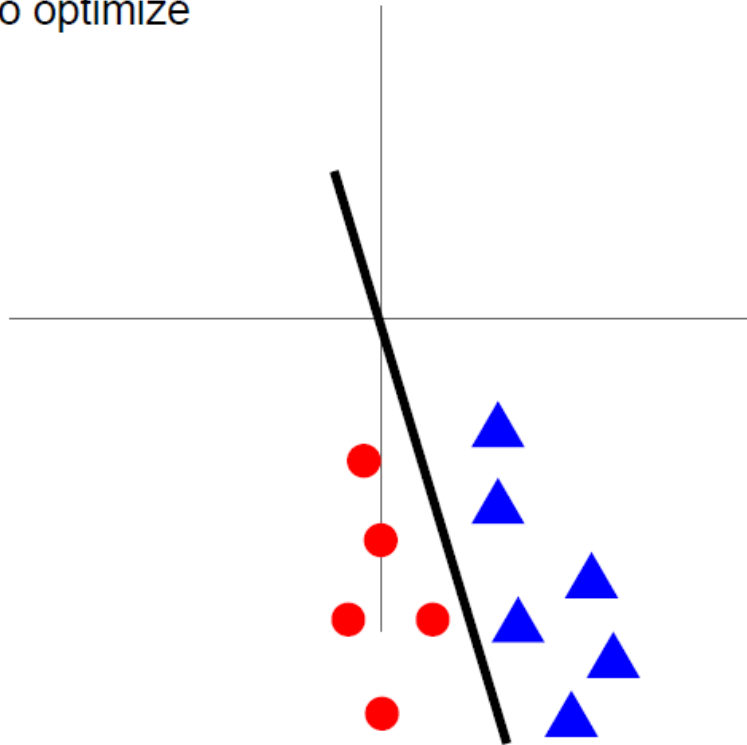


# Last time: Data Preprocessing

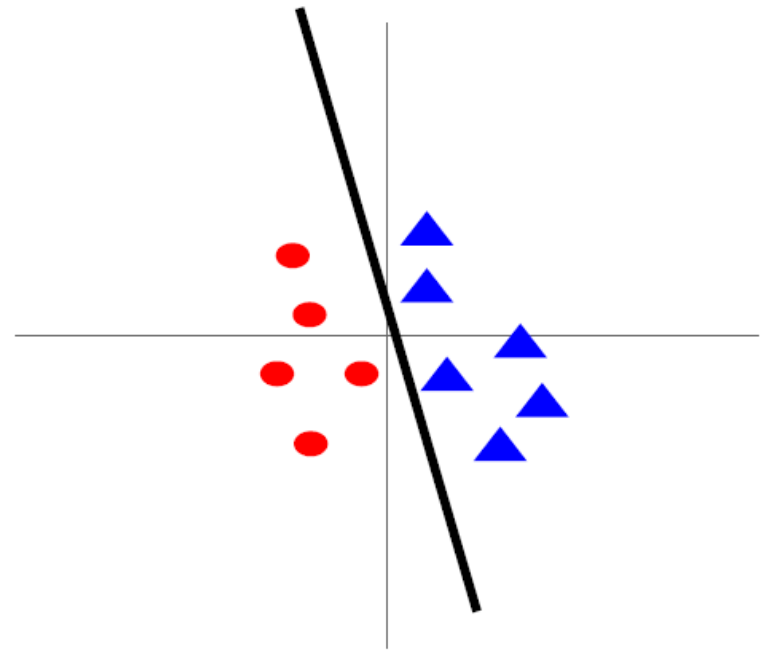


# Last time: Data Preprocessing

**Before normalization:** classification loss  
very sensitive to changes in weight matrix;  
hard to optimize



**After normalization:** less sensitive to small  
changes in weights; easier to optimize



# Last time: Batch Normalization

**Input:**  $x : N \times D$

**Learnable params:**

$$\gamma, \beta : D$$

**Intermediates:**  $\mu, \sigma : D$   
 $\hat{x} : N \times D$

**Output:**  $y : N \times D$

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

# Model Validation

