

# A Brief Introduction to Deep Learning

Baojian Zhou, University at Albany, SUNY

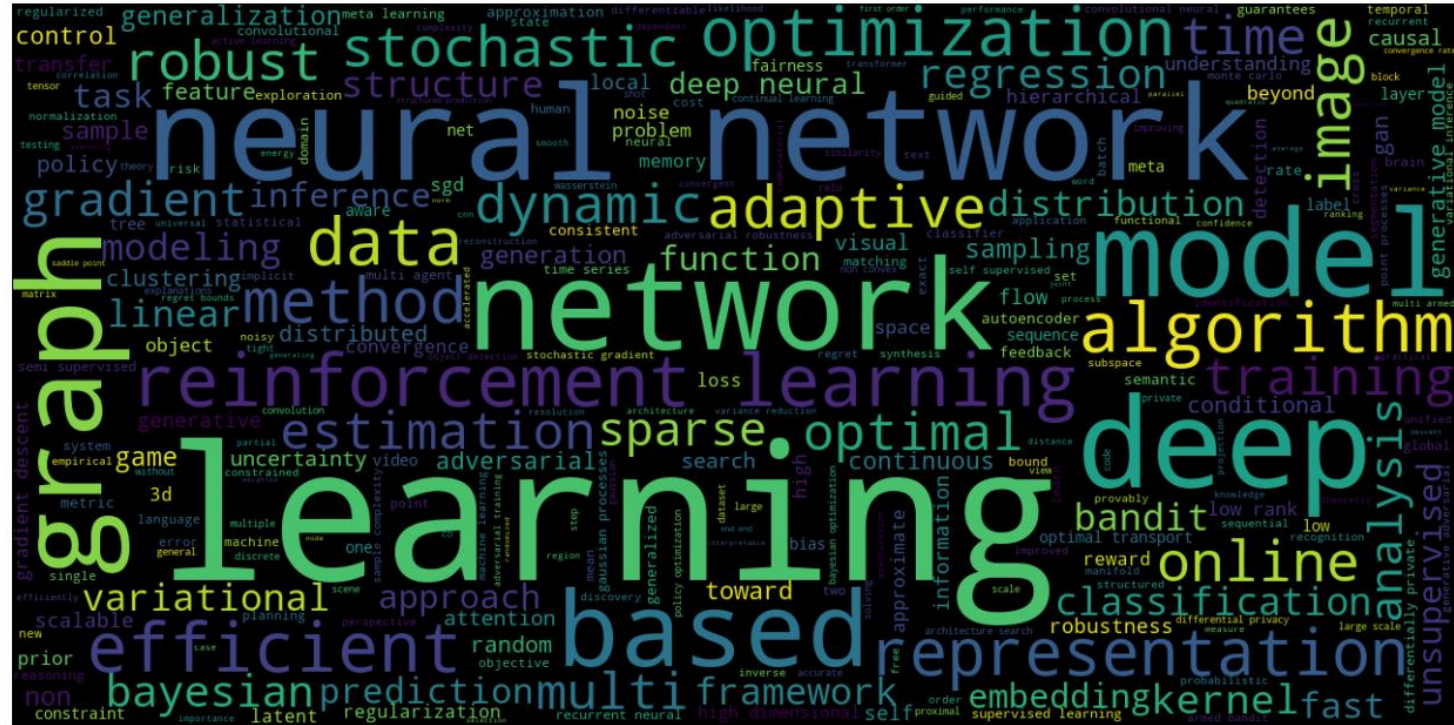
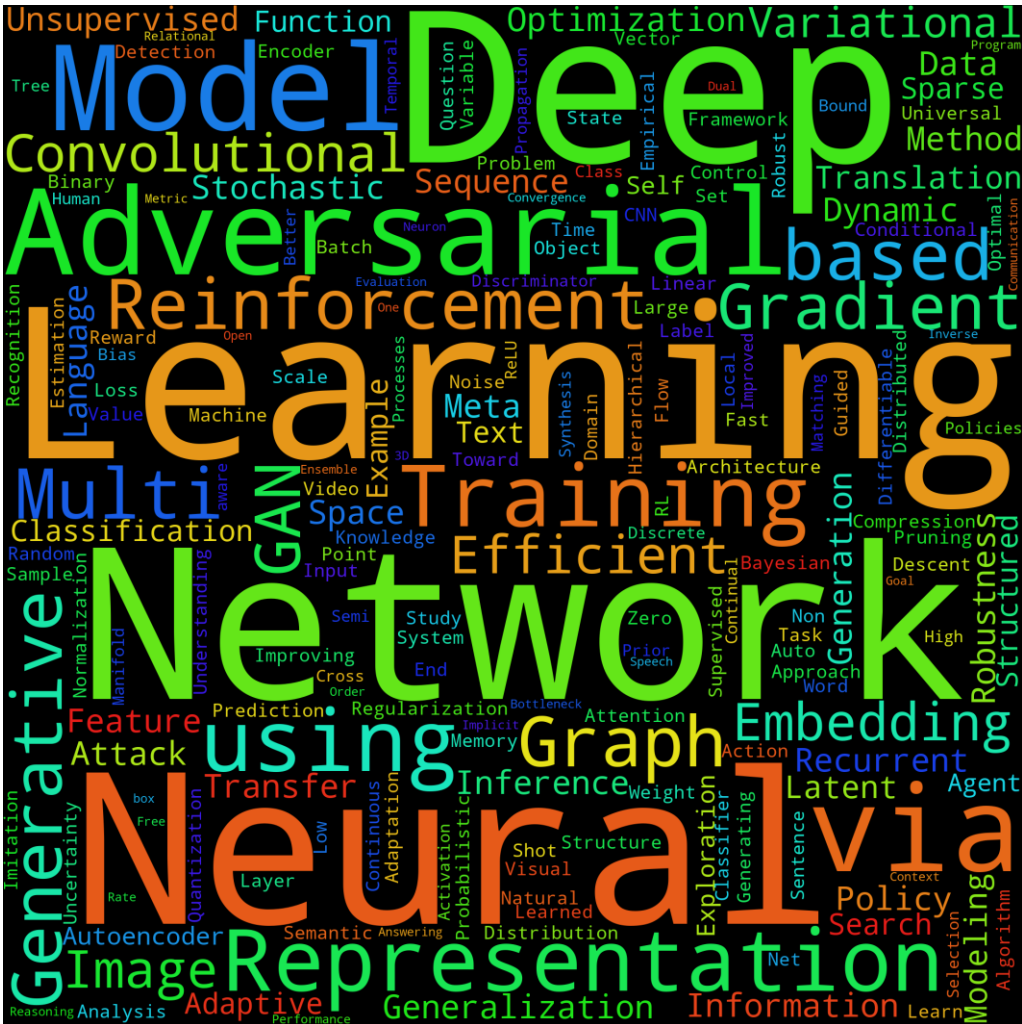
# Outline

- **Successful stories of DL (2 minutes)**
- **DL Basics (28 minutes)**
- **Build two simple neural networks (45 minutes)**
- **TensorFlow and Others (5 minutes)**

# Outline

- **Successful stories of DL (2 minutes)**
- DL Basics (28 minutes)
- Build two simple neural networks (45 minutes)
- TensorFlow and Others (5 minutes)

# AI, machine learning, deep learning, ...

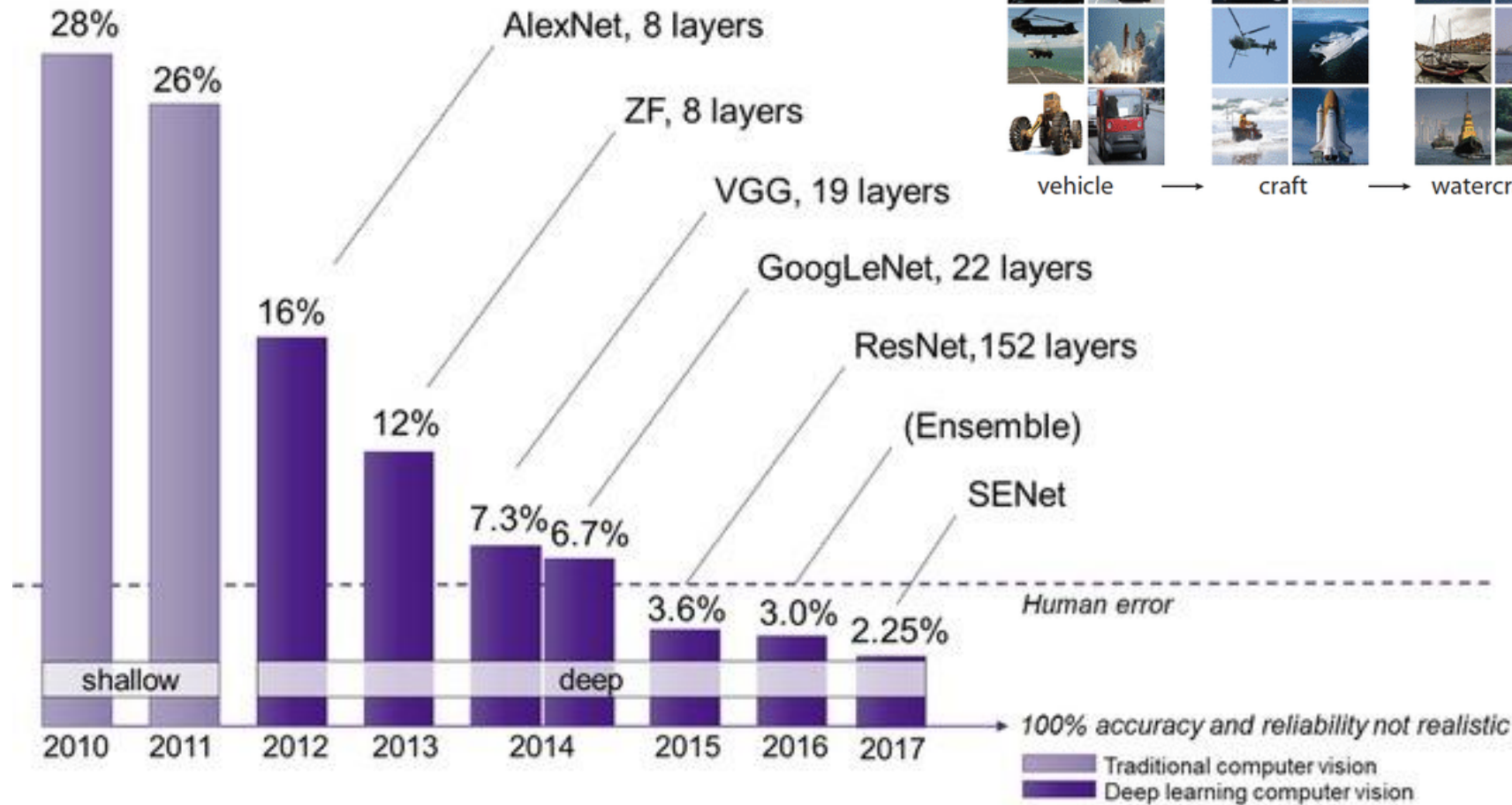


# Word Cloud - NeurIPS, 2019

# Word Cloud - ICLR, 2019

# Successful stories of DL

## ImageNet Classification





# Successful stories of DL



AlphaGo



# Successful stories of DL

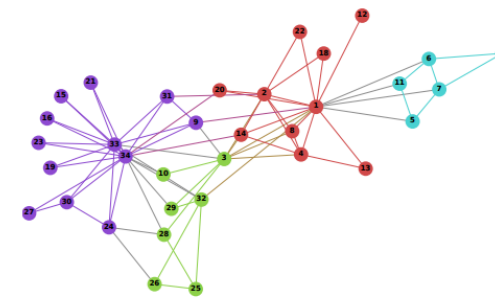
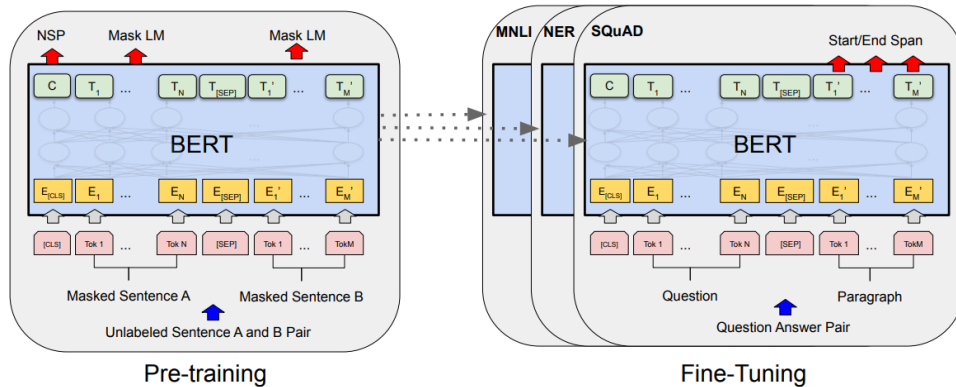


## Self-driving

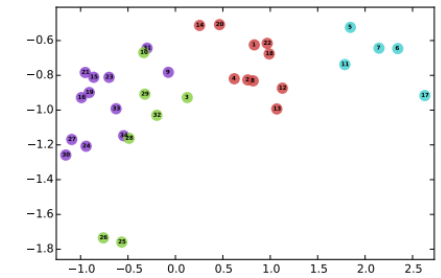


# Successful stories of DL

Generative adversarial network, graph embedding, deep reinforcement learning, and many others ...



(a) Input: Karate Graph



(b) Output: Representation

Word2vec and BERT in NLP

DeepWalk

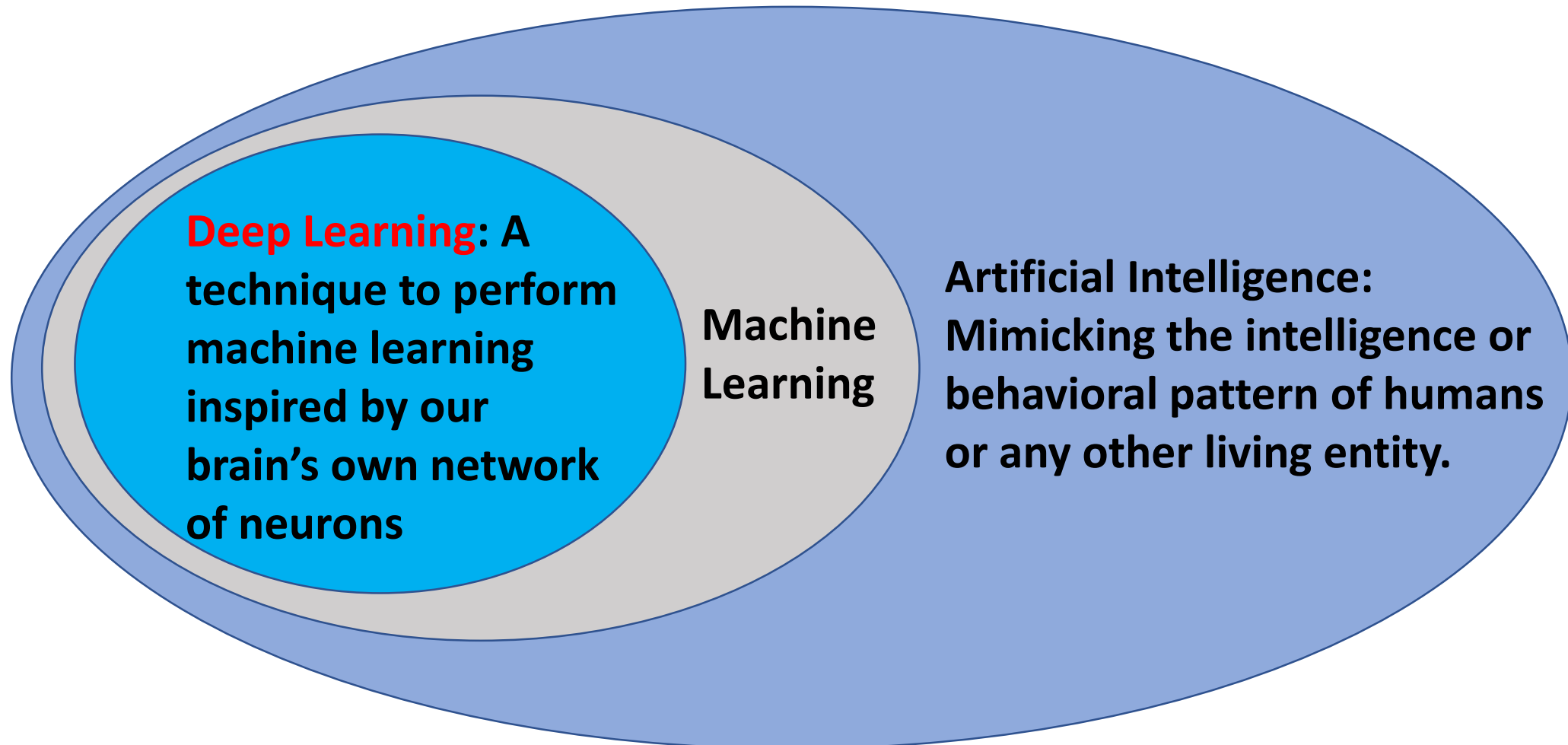


# Outline

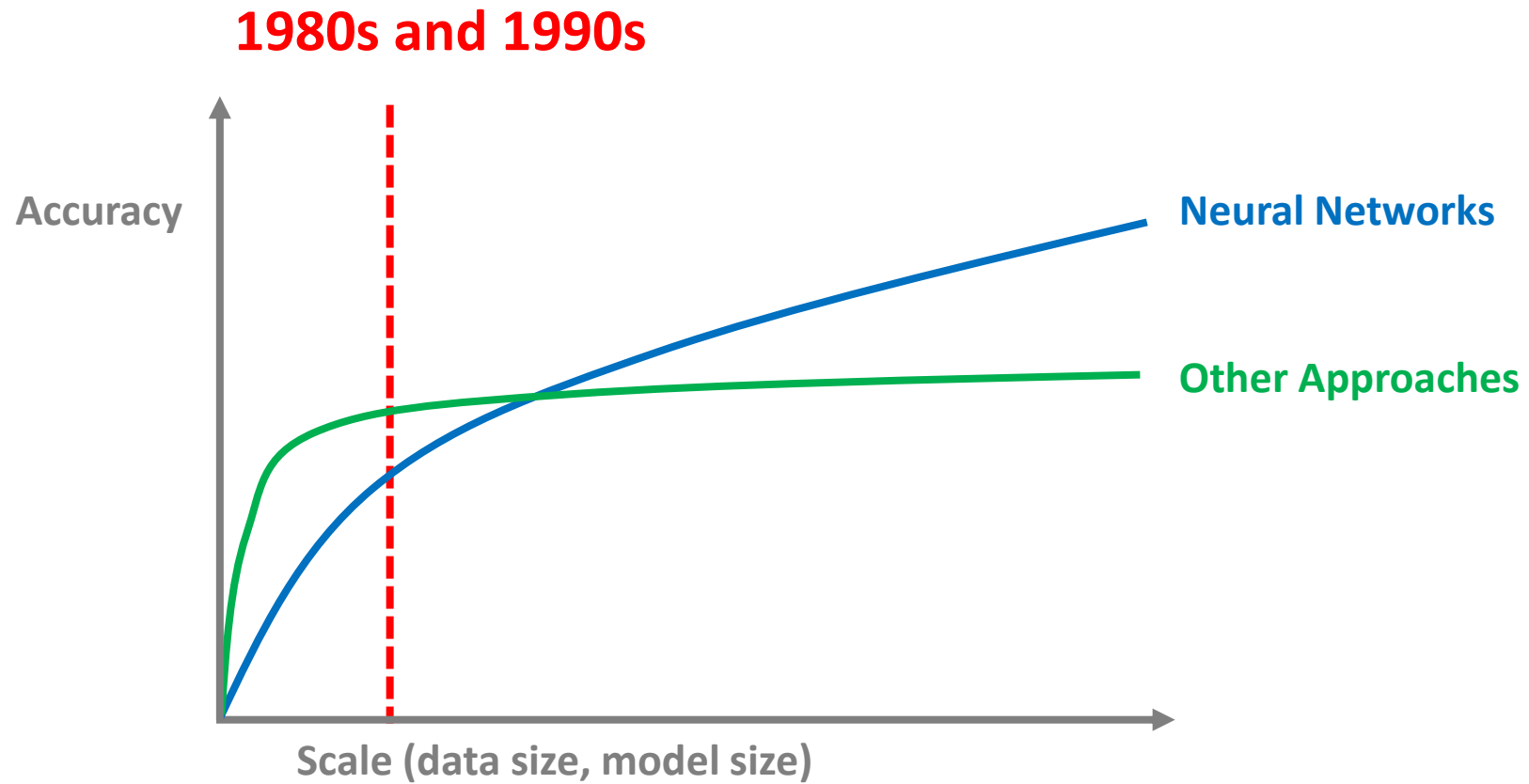
- Successful stories of DL (2 minutes)
- **DL Basics (28 minutes)**
- Build two simple neural networks (45 minutes)
- TensorFlow and Others (5 minutes)

# Deep Learning – Definition

Deep learning is part of a broader family of machine learning methods based on **artificial neural networks**. Learning can be supervised, semi-supervised or unsupervised.

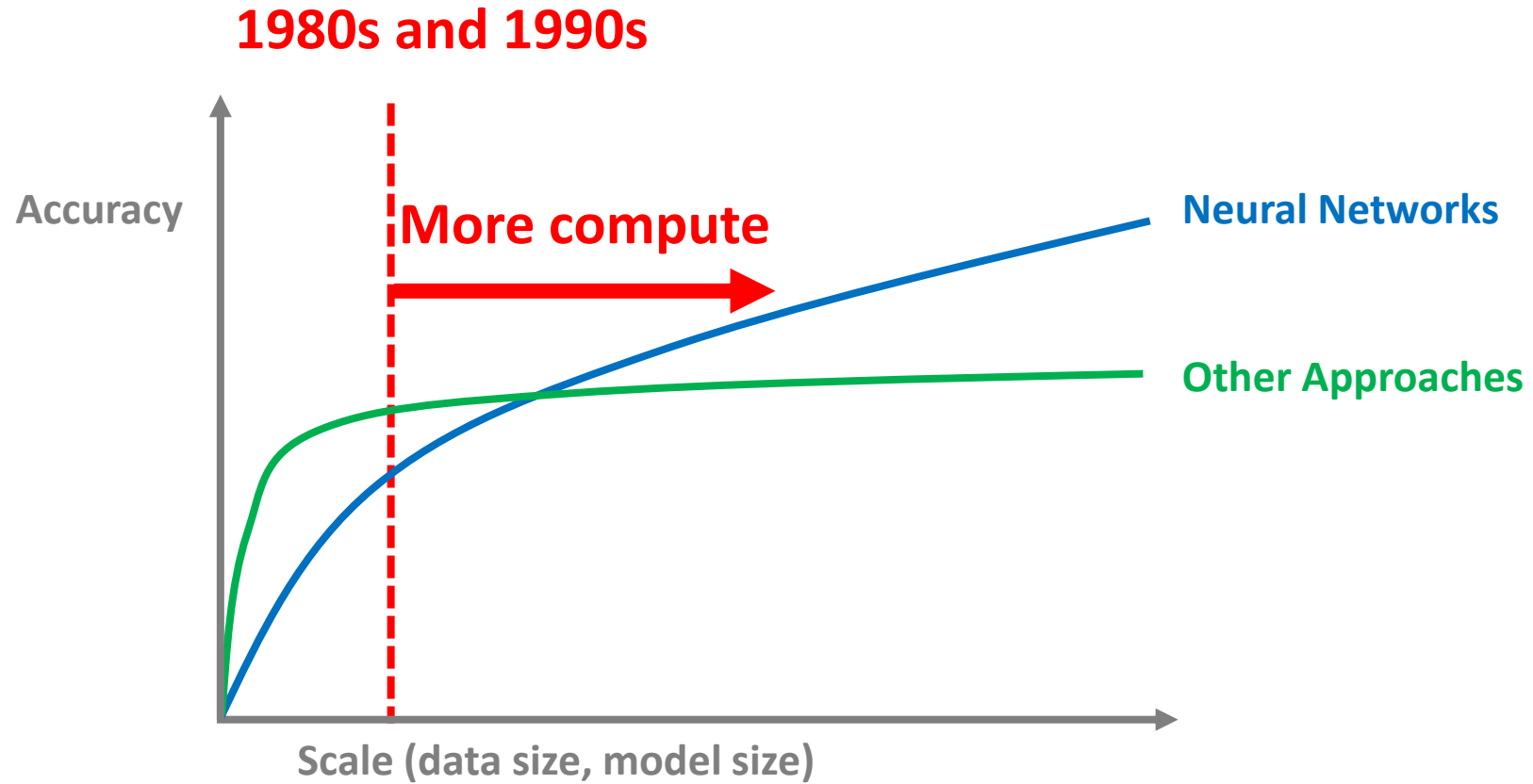


# Why deep learning ?



Jeff Dean's Lecture for YC AI

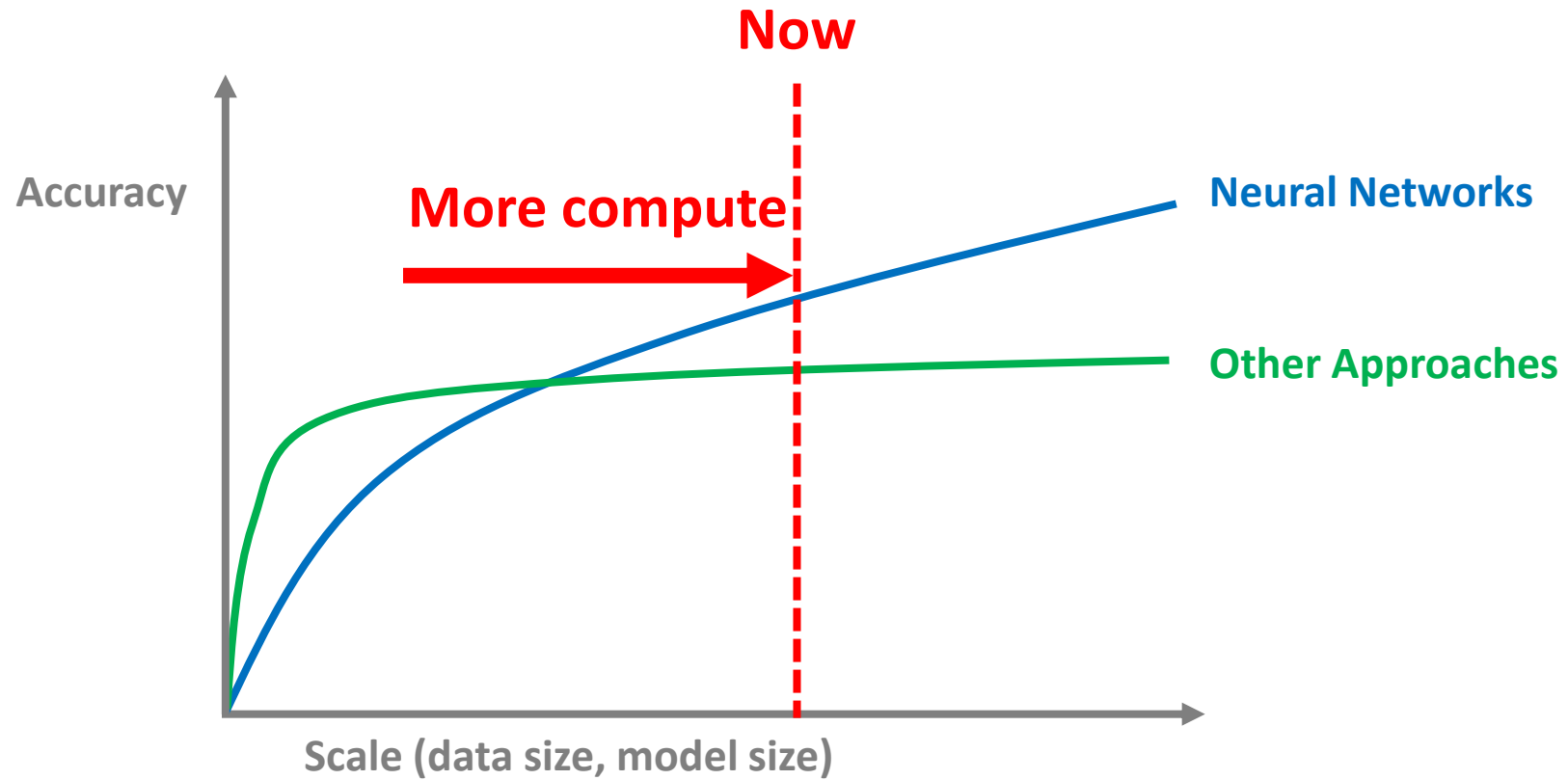
# Why deep learning ?



Jeff Dean's Lecture for YC AI



# Why deep learning ?

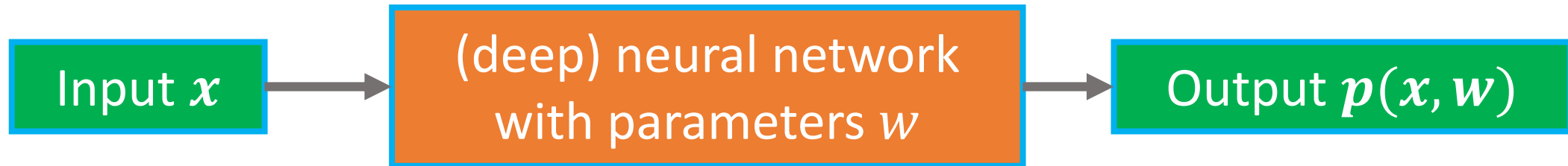


Jeff Dean's Lecture for YC AI

# Artificial Neural Networks – Basics

What is a neural network

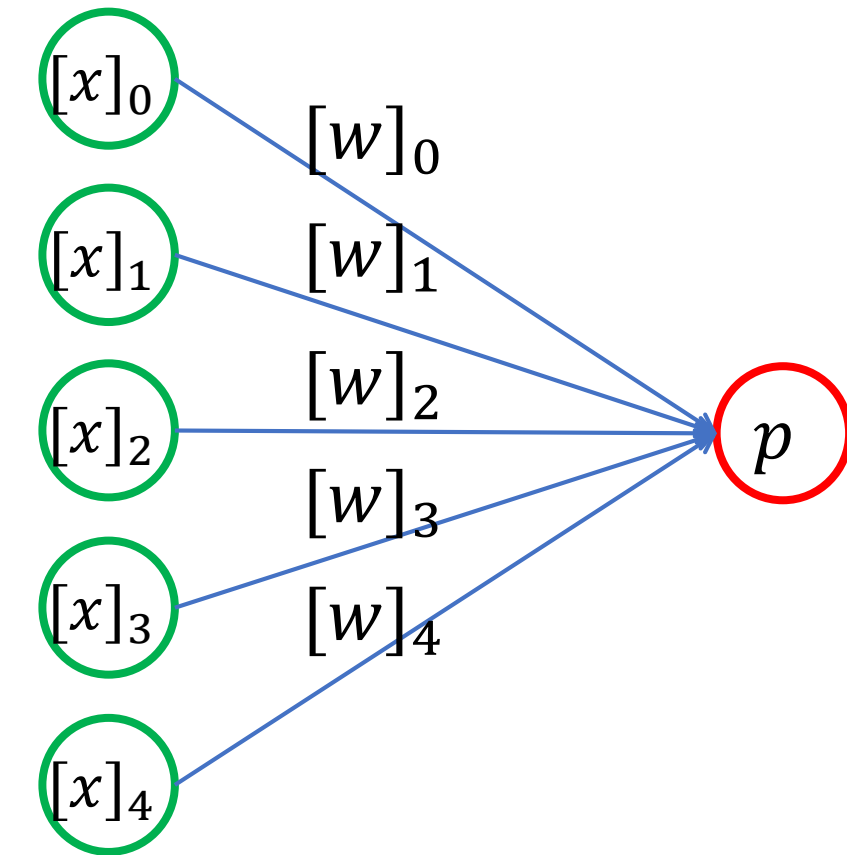
- A brain ?
- A computational graph ?
- A function ?



**We can train the neural network, i.e., make it learn, so that, given any input  $x$ , the output  $p(x; w)$  is what we want to see.**

# Artificial Neural Networks – Basics

## A simple neural network



Input Layer

Output Layer

**Nodes** (Neurons) to represent input/output values.  
**Edges** to represent multiplication by a weight.

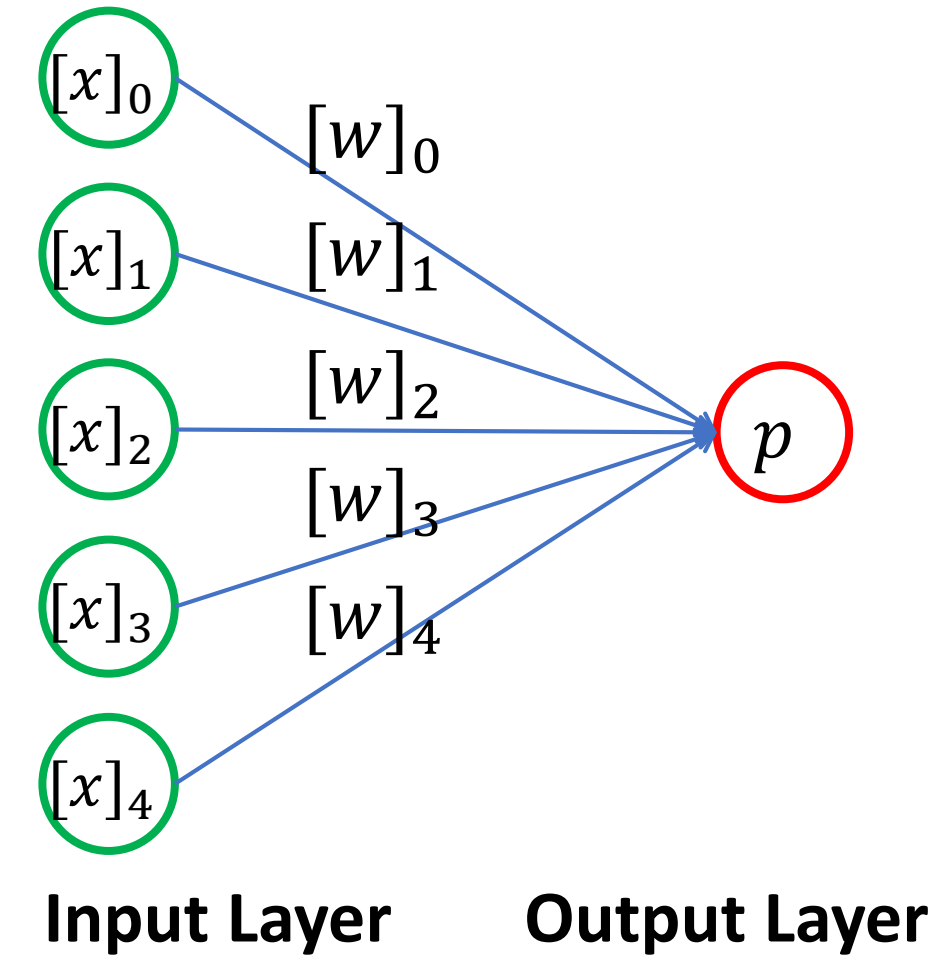
The output value is given by

$$p(w; x) = x_1 w_1 + x_2 w_2 + x_3 w_3 + x_4 w_4 + x_5 w_5$$

$$p(w; x) = \sum_{i=1}^5 x_i w_i = \mathbf{x}^T \mathbf{w}$$

# Artificial Neural Networks – Basics

## A simple neural network for **regression**



Suppose we have a set of training data  $\{\mathbf{x}_i, y_i\}_{i=1}^N$  and suppose we want the weights so that

$$\mathbf{x}_i^\top \mathbf{w} \approx y_i \quad \forall i \in \{1, \dots, n\}.$$

Then, we our goal is

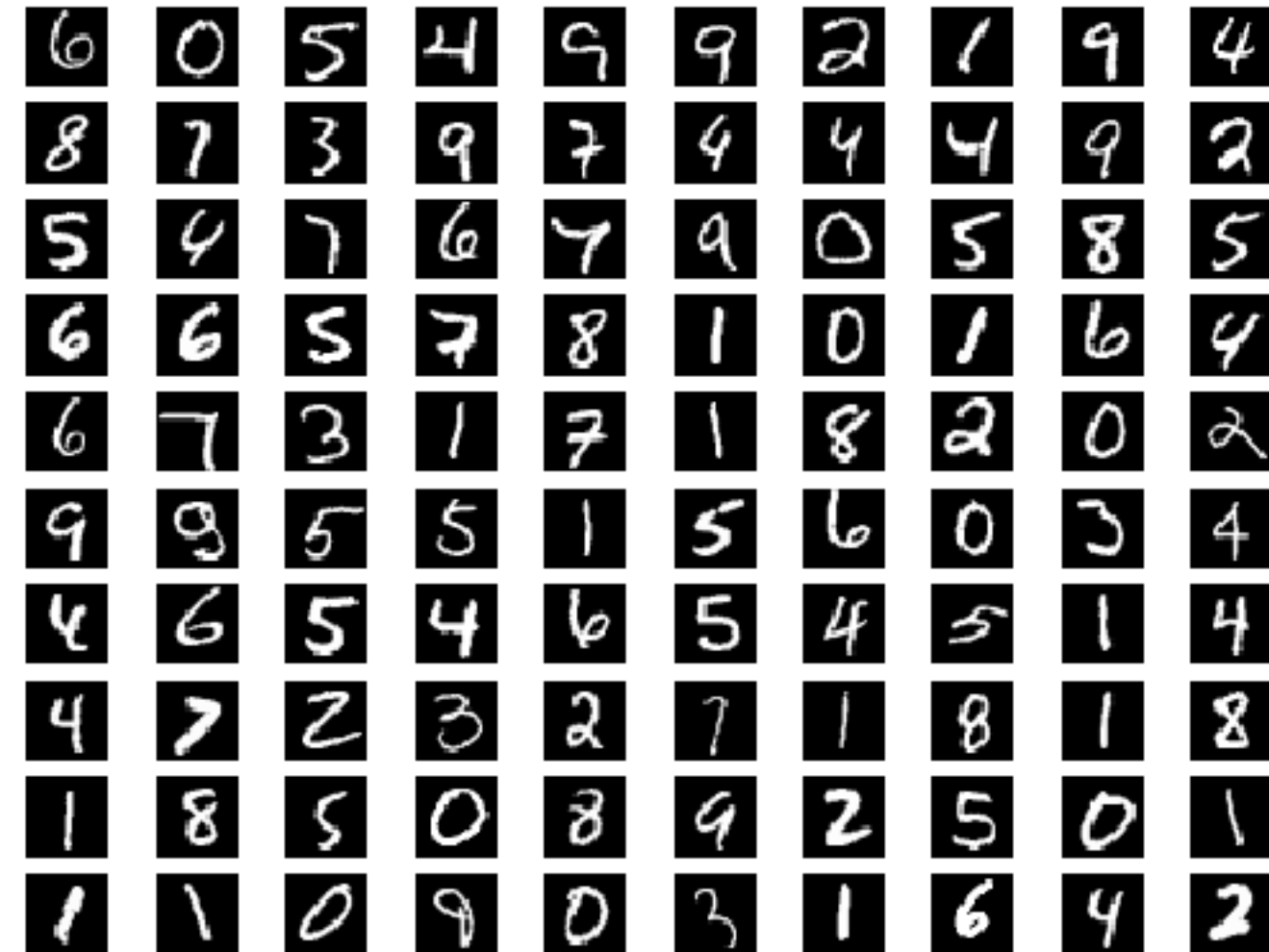
$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i^\top \mathbf{w}, y_i),$$

for some loss  $\ell$ . This is one way to describe linear regression using a network.



# Artificial Neural Networks – Basics

## A simple neural network for **classification**

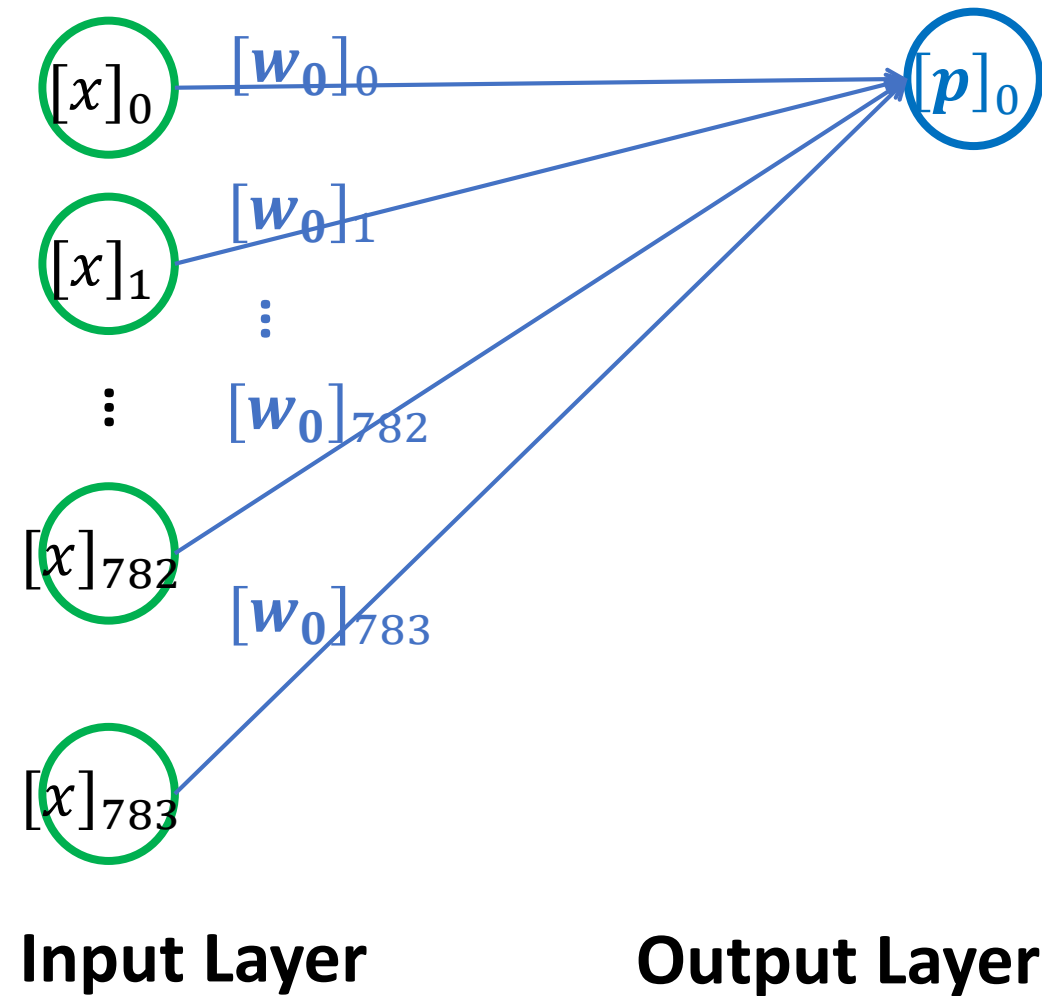


Suppose our inputs are images of digits we want to classify (MNIST dataset). The potential outputs are 0, 1, 2, ..., 9.

Each digit image is  $d = 28 \times 28$

# Artificial Neural Networks – Basics

## A simple neural network for **classification**



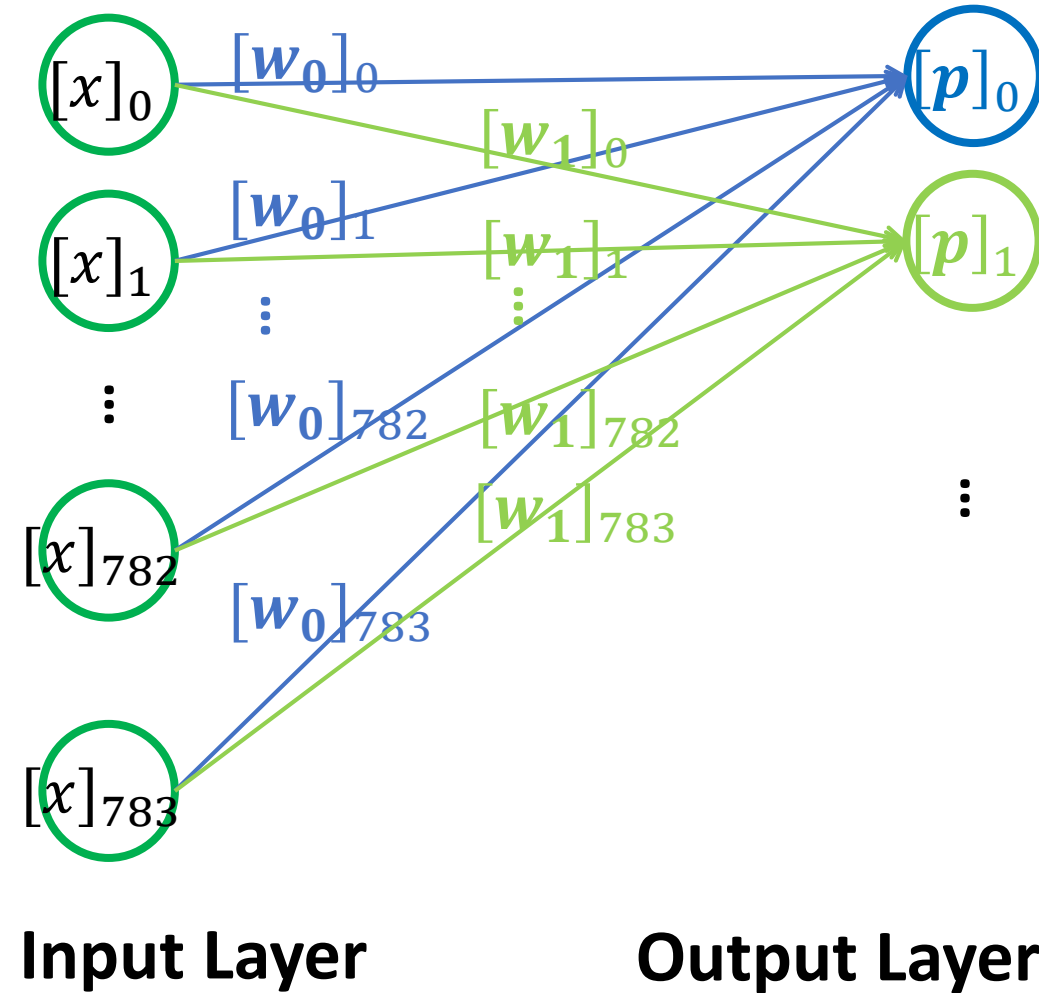
Now, instead of a single vector  $w$ , we want 10 weight vectors (one for each digit). Define the vector  $\mathbf{w}_j$  for  $j \in \{0, 1, 2, \dots, 9\}$ .

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_0^\top \\ \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \vdots \\ \mathbf{w}_9^\top \end{bmatrix}. \text{ Then, } \mathbf{p} = \mathbf{W}\mathbf{x}.$$

We want to choose  $\mathbf{W}$  such that, when the correct digit is  $j \in \{0, 1, 2, \dots, 9\}$ ,  $[\mathbf{p}]_j > [\mathbf{p}]_i$  for  $i \neq j$ .

# Artificial Neural Networks – Basics

## A simple neural network for **classification**



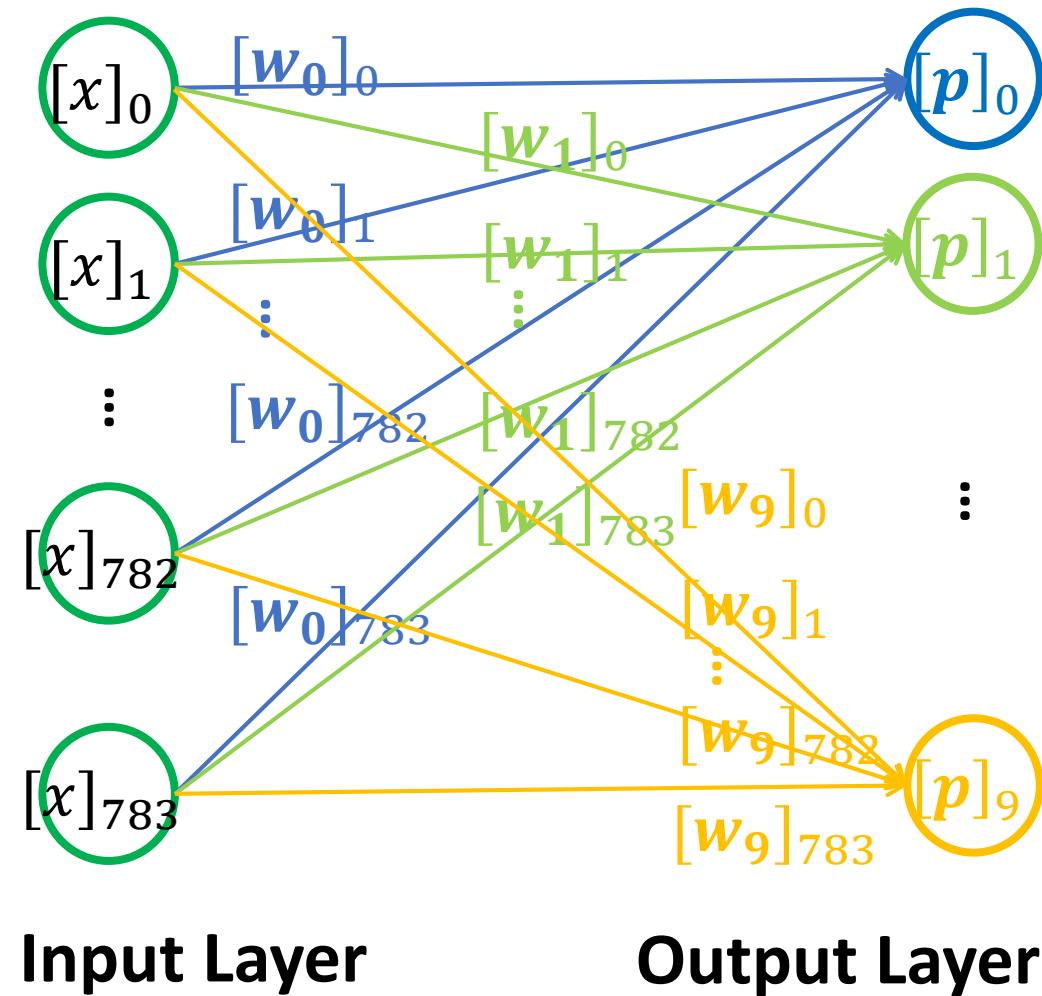
Now, instead of a single vector  $w$ , we want 10 weight vectors (one for each digit). Define the vector  $w_j$  for  $j \in \{0, 1, 2, \dots, 9\}$ .

$$W = \begin{bmatrix} w_0^\top \\ w_1^\top \\ w_2^\top \\ \vdots \\ w_9^\top \end{bmatrix}. \text{ Then, } p = Wx.$$

We want to choose  $W$  such that, when the correct digit is  $j \in \{0, 1, 2, \dots, 9\}$ ,  $[p]_j > [p]_i$  for  $i \neq j$ .

# Artificial Neural Networks – Basics

## A simple neural network for **classification**



Now, instead of a single vector  $w$ , we want 10 weight vectors (one for each digit). Define the vector  $\mathbf{w}_j$  for  $j \in \{0, 1, 2, \dots, 9\}$ .

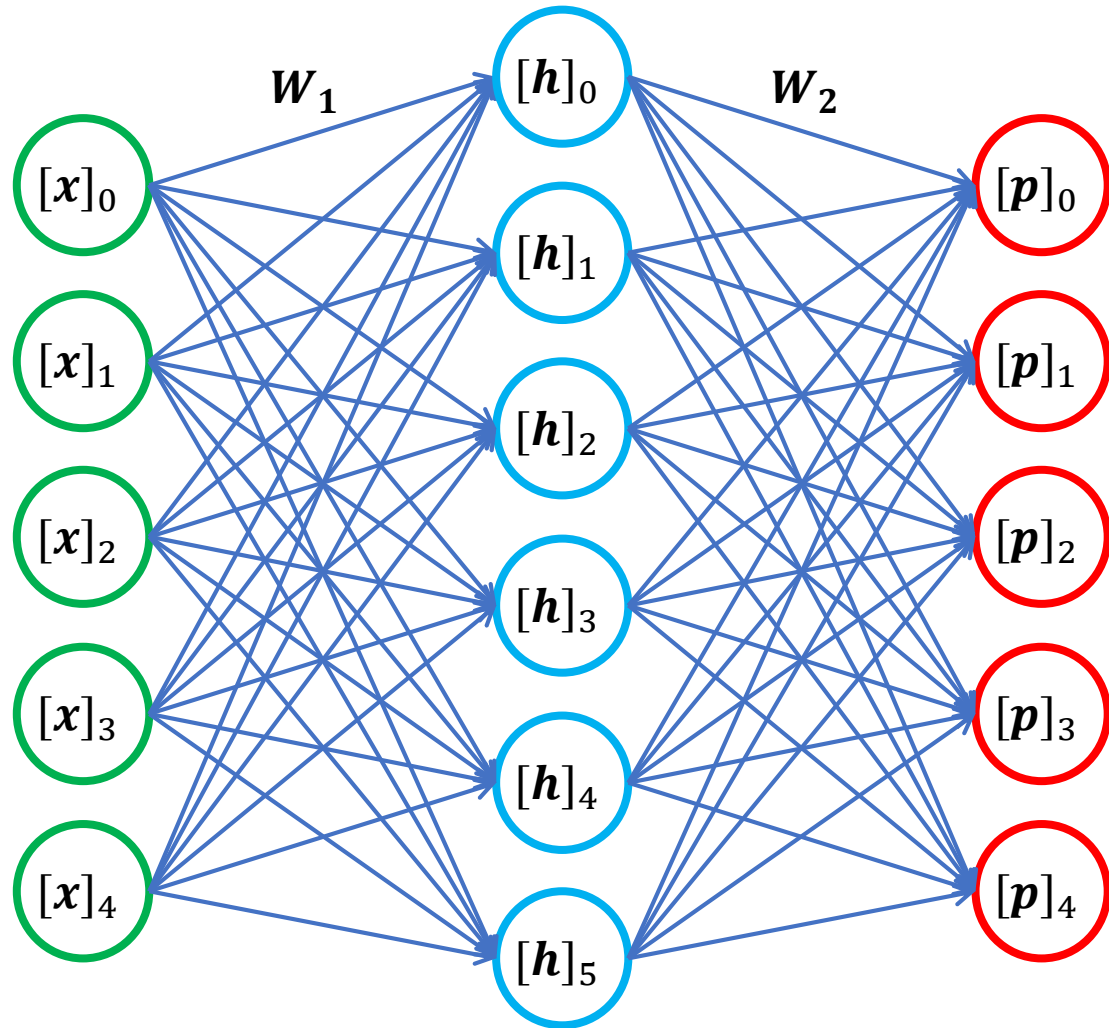
$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_0^\top \\ \mathbf{w}_1^\top \\ \mathbf{w}_2^\top \\ \vdots \\ \mathbf{w}_9^\top \end{bmatrix}. \text{ Then, } \mathbf{p} = \mathbf{W}\mathbf{x}.$$

We want to choose  $\mathbf{W}$  such that, when the correct digit is  $j \in \{0, 1, 2, \dots, 9\}$ ,  $[\mathbf{p}]_j > [\mathbf{p}]_i$  for  $i \neq j$ .



# Artificial Neural Networks – Basics

## Hidden Layers



Input Layer

Hidden Layer

Output Layer

Rather than go directly from inputs to outputs. We can add one or more hidden layers. Now the network performs the computation. Then,

$$\mathbf{p} = \mathbf{W}_2 \mathbf{h}(\mathbf{W}_1 \mathbf{x}).$$

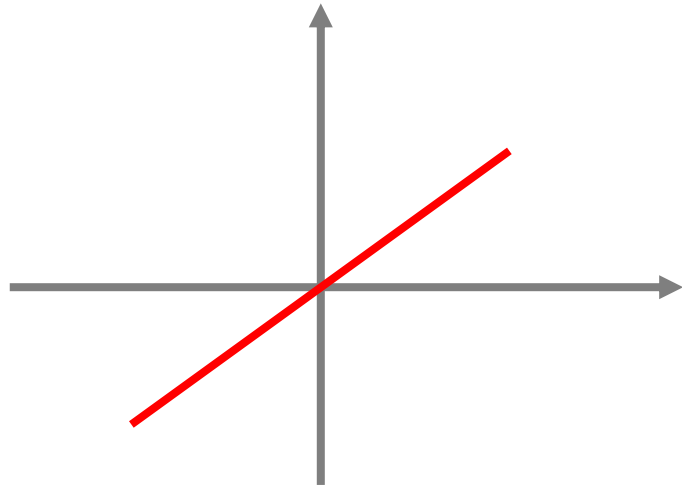
Can also include bias terms such that

$$\mathbf{p} = \mathbf{W}_2 \mathbf{h}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_0) + \mathbf{b}_1.$$

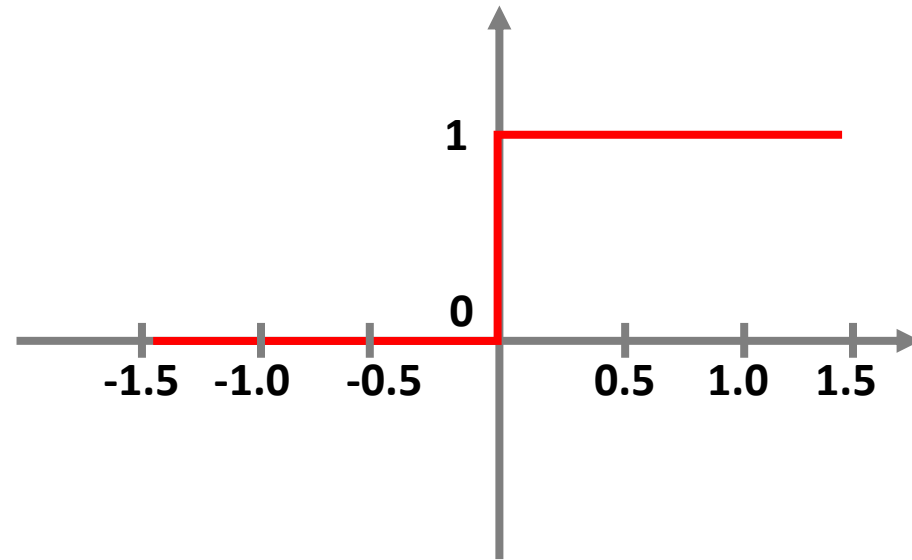
# Artificial Neural Networks – Basics

## Activation functions $h$

Inspired by neuroscience, use activation functions to approximate the “firing or not” of a neuron. We want to approximate a 0-1 step function.



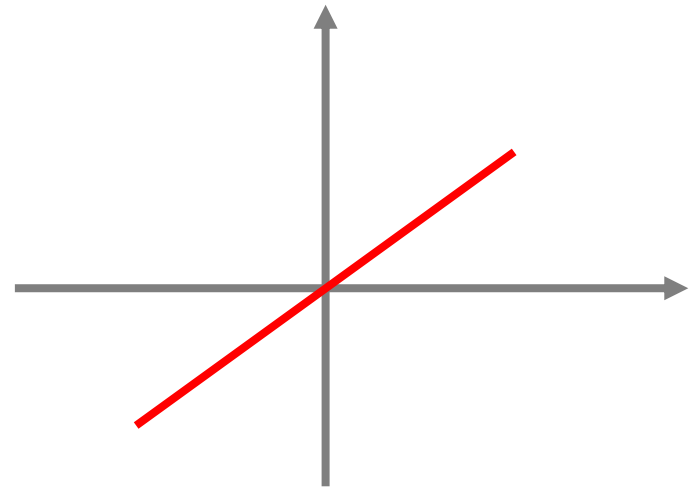
**Identity (a.k.a Linear)**



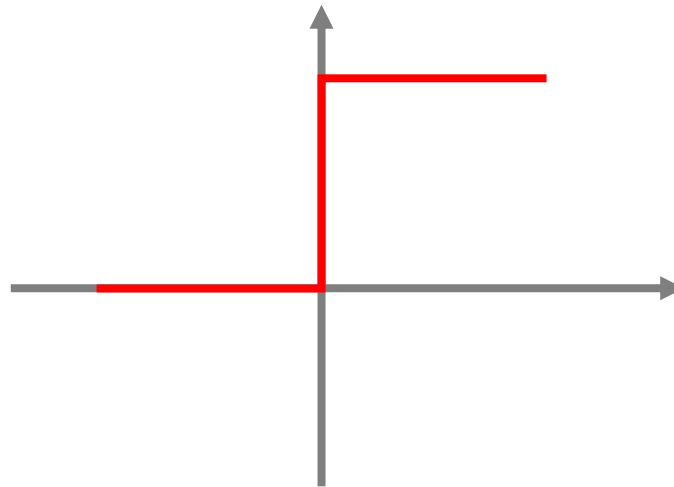
**Binary Step**

# Artificial Neural Networks – Basics

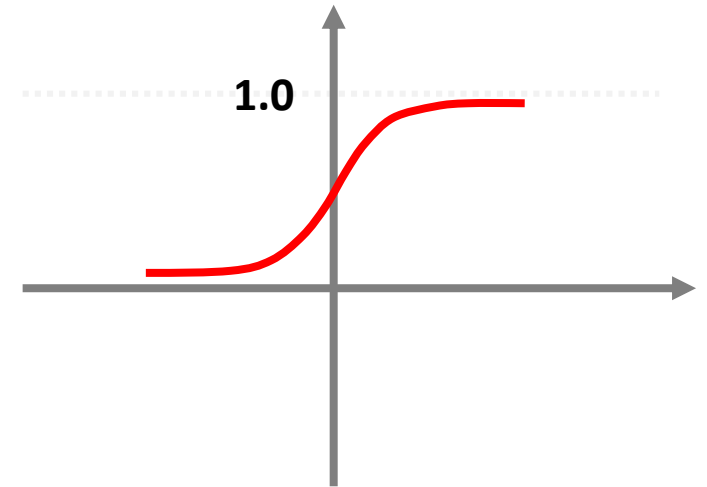
## Activation functions $h$



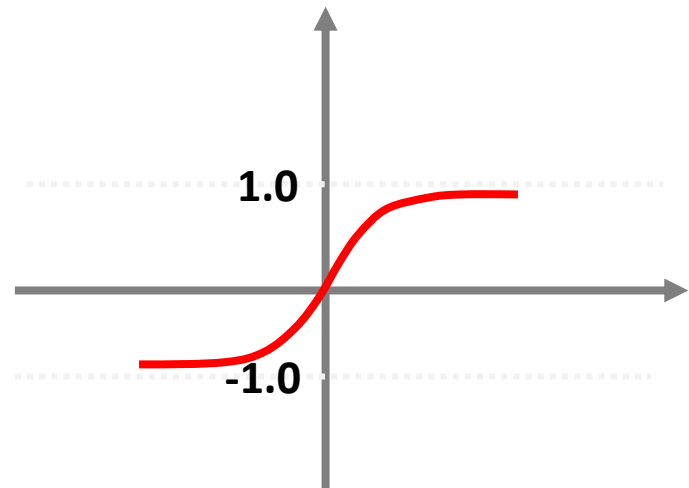
**Identity (a.k.a Linear)**



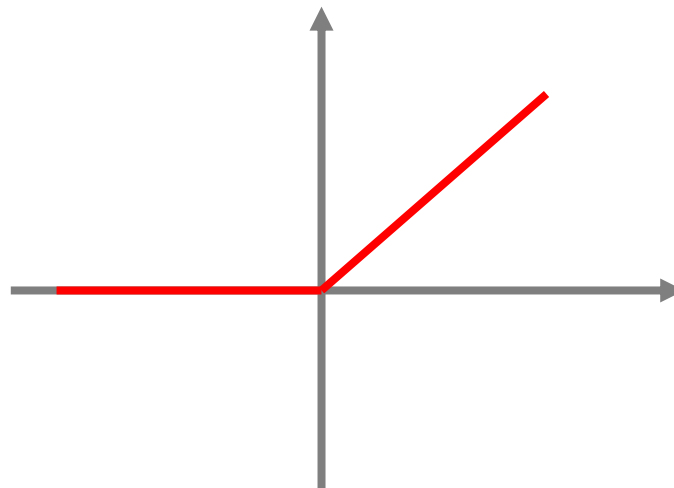
**Binary Step**



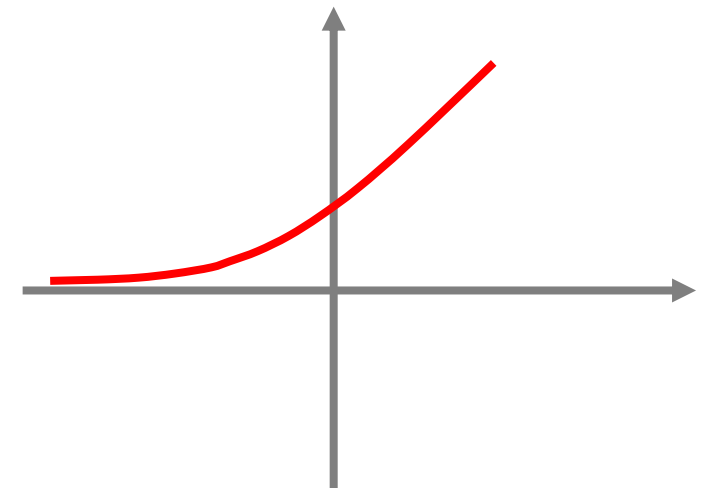
**Logistic (a.k.a Sigmoid)**



**TanH**



**Rectified Linear Unit (ReLU)**



**SoftPlus**

# Artificial Neural Networks – Basics

To design a neural net, you need to

Even for a basic Neural Network, there are many design decisions to make:

- **Number of hidden layers (depth)**
- **Number of units per hidden layer (width)**
- **Type of activation function (nonlinearity)**
- **Form of objective function**
- **Recurrent or feed-forward?**

These parameters will affect the overall performance such as computation cost and accuracy.



# Artificial Neural Networks – Basics

## Computation cost

For example, I want to design a neural network for classifying image digits (MNIST). Let us consider the following neural network:

- $28 \times 28 = 784$  neurons in the input layer
- 512 neurons in the first hidden layer (fully connected)
- 256 neurons in the second hidden layer (fully connected)
- 10 neurons in the output layer (fully connected)

**Question: How many optimization variables are there?**

# Artificial Neural Networks – Basics

## Computation cost

For example, I want to design a neural network for classifying image digits (MNIST). Let us consider the following neural network:

- $28 \times 28 = 784$  neurons in the input layer
- 512 neurons in the first hidden layer (fully connected)
- 256 neurons in the second hidden layer (fully connected)
- 10 neurons in the output layer (fully connected)

**Question: How many optimization variables are there?**

In the first layer:  $784 \times 512$

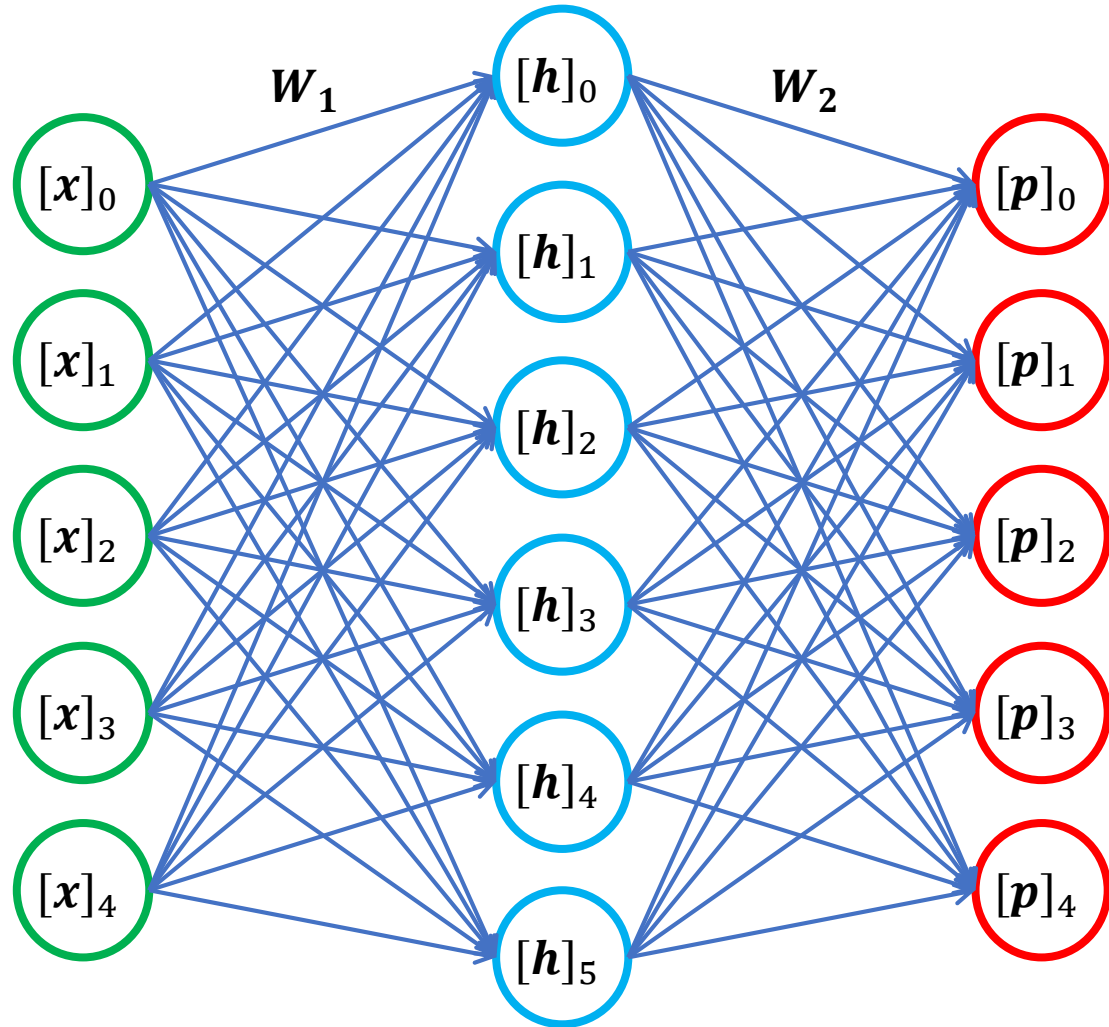
In the second layer:  $512 \times 256$

In the third layer:  $256 \times 10$

$784 \times 512 + 512 \times 256 + 256 \times 10 = \mathbf{535040}.$

# Artificial Neural Networks – Basics

How to update the weights  $W_1, W_2$ ?



Input Layer

Hidden Layer

Output Layer

In general, we minimize a loss

$$\min_{W_1, W_2} \frac{1}{n} \sum_{i=1}^n \ell(W_1, W_2; x_i, y_i)$$

Now the question is how to optimize the loss w.r.t  $W_1, W_2$ .

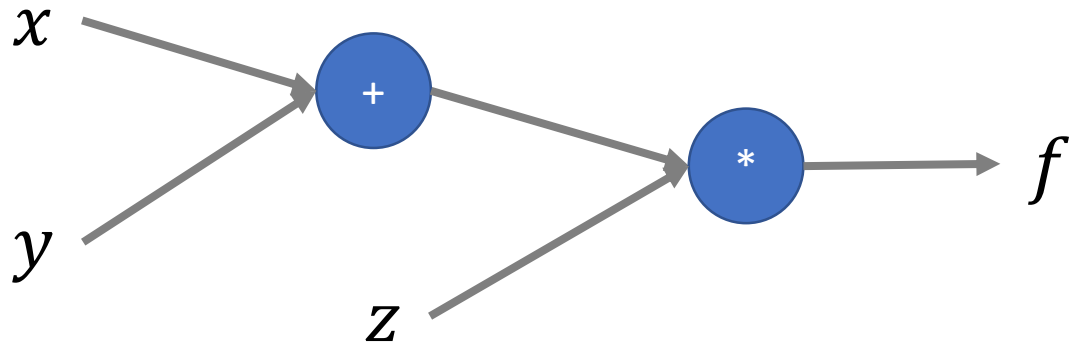
**Stochastic Gradient Descent:**

$$[W_1]_{ij}^{t+1} = [W_1]_{ij}^t - \eta_t \frac{d\ell}{d[W_1]_{ij}^t}$$

$$[W_2]_{ij}^{t+1} = [W_2]_{ij}^t - \eta_t \frac{d\ell}{d[W_2]_{ij}^t}$$

# Artificial Neural Networks – Basics

## Backpropagation: a simple example



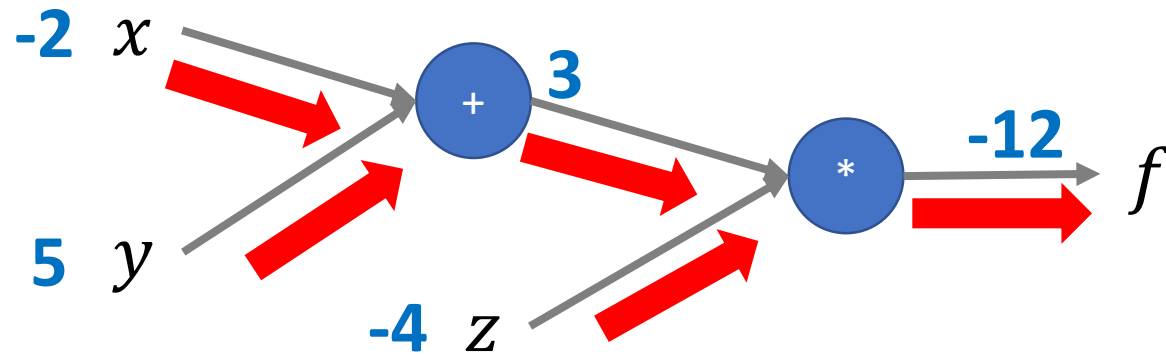
$$f(x, y, z) = (x + y) * z$$

We want to calculate:

$$\frac{\partial f}{\partial x}, \quad \frac{\partial f}{\partial y}, \quad \frac{\partial f}{\partial z}$$

# Artificial Neural Networks – Basics

## Backpropagation: a simple example



Let  $x = -2, y = 5, z = -4$ ,  
then  $f(x, y, z) = -12$ .

**Feed Forward**

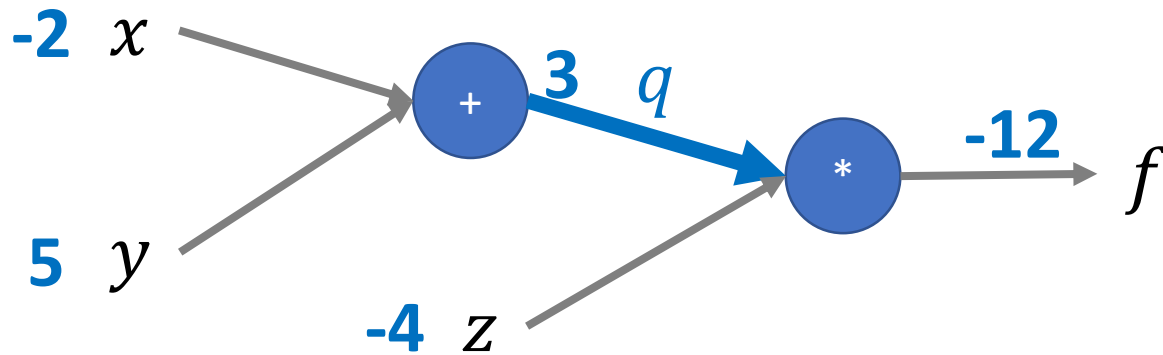
$$f(x, y, z) = (x + y) * z$$

We want to calculate:

$$\frac{\partial f}{\partial x}, \quad \frac{\partial f}{\partial y}, \quad \frac{\partial f}{\partial z}$$

# Artificial Neural Networks – Basics

## Backpropagation: a simple example



$$f(x, y, z) = (x + y) * z$$

We want to calculate:

$$\frac{\partial f}{\partial x}, \quad \frac{\partial f}{\partial y}, \quad \frac{\partial f}{\partial z}$$

Let  $x = -2, y = 5, z = -4$ ,  
then  $f(x, y, z) = -12$ .

We let the intermediate result  
be a function

$$q(x, y) = (x + y).$$

$$\text{Then } f(x, y, z) = q(x, y)z.$$

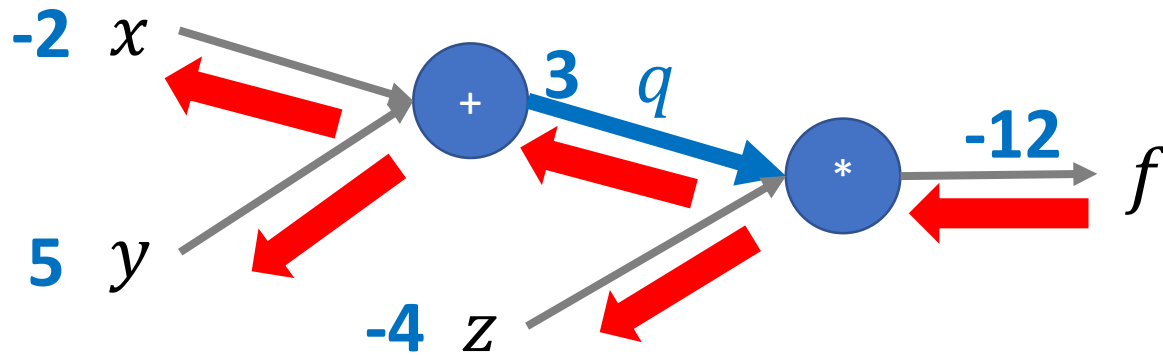


$$\begin{array}{ll} \frac{\partial q}{\partial x} = 1, & \frac{\partial q}{\partial y} = 1 \\ \frac{\partial f}{\partial q} = z, & \frac{\partial f}{\partial z} = q \end{array}$$



# Artificial Neural Networks – Basics

## Backpropagation: a simple example



$$f(x, y, z) = (x + y) * z$$

We want to calculate:

$$\frac{\partial f}{\partial x}, \quad \frac{\partial f}{\partial y}, \quad \frac{\partial f}{\partial z}$$

**Solution: By using the chain Rule**

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial x}, \quad \frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \cdot \frac{\partial q}{\partial y}, \quad \frac{\partial f}{\partial z} = q$$

Let  $x = -2, y = 5, z = -4$ ,  
then  $f(x, y, z) = -12$ .

We let the intermediate result  
be a function

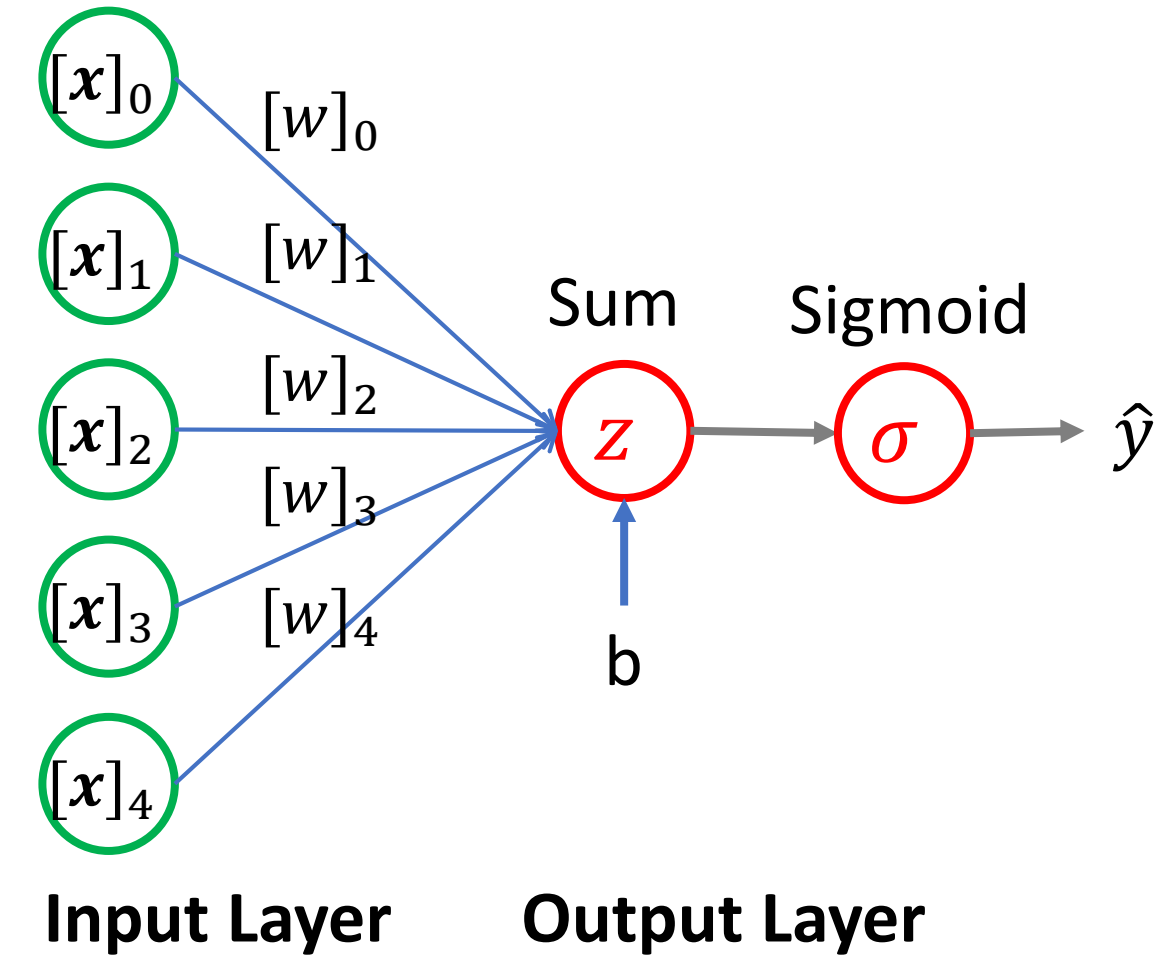
$$q(x, y) = (x + y).$$

$$\text{Then } f(x, y, z) = q(x, y)z.$$

$$\begin{aligned} \frac{\partial q}{\partial x} &= 1, & \frac{\partial q}{\partial y} &= 1 \\ \frac{\partial f}{\partial q} &= z, & \frac{\partial f}{\partial z} &= q \end{aligned}$$

# Artificial Neural Networks – Basics

## Backpropagation



Given a training sample  $(\mathbf{x}, y)$ ,  
compute the loss

$$z = \mathbf{x}^\top \mathbf{w} + b, \hat{y} = \sigma(z),$$

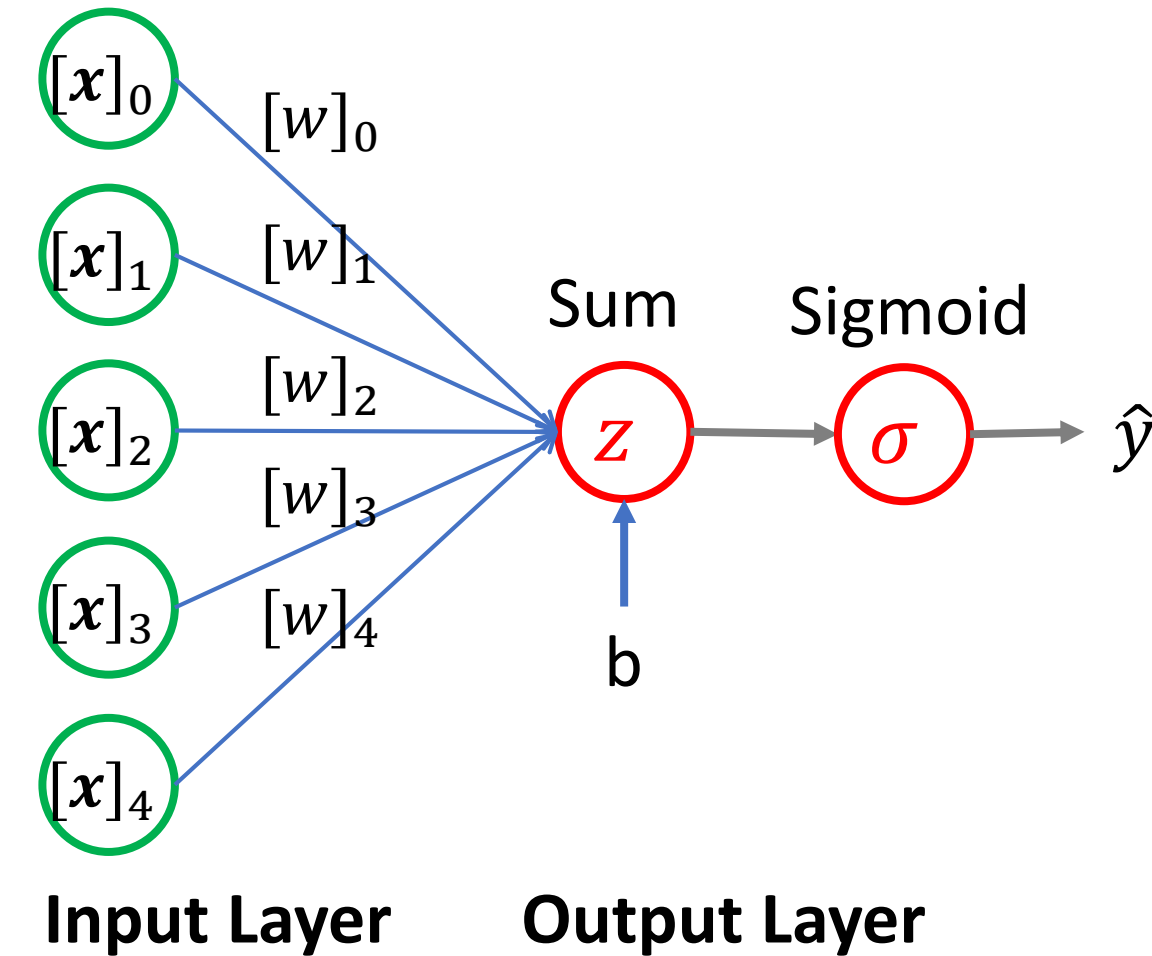
$$\ell = \frac{1}{2} (\hat{y} - y)^2$$

**Q: calculate**

$$\frac{\partial \ell}{\partial w}, \frac{\partial \ell}{\partial b}$$

# Artificial Neural Networks – Basics

## Backpropagation



Given a training sample  $(\mathbf{x}, y)$ ,  
compute the loss

$$z = \mathbf{x}^\top \mathbf{w} + b, \hat{y} = \sigma(z),$$
$$\ell = \frac{1}{2} (\hat{y} - y)^2$$

Compute the derivatives

$$\frac{d\ell}{d\hat{y}} = (\hat{y} - y),$$

$$\frac{d\ell}{dz} = \frac{d\ell}{d\hat{y}} \cdot \frac{d\hat{y}}{dz} = \frac{d\ell}{d\hat{y}} \cdot \sigma(z) \cdot (1 - \sigma(z)),$$

$$\frac{\partial \ell}{\partial w} = \frac{d\ell}{dz} \cdot \frac{dz}{\partial w} = \frac{d\ell}{dz} \cdot \mathbf{x},$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial \ell}{\partial b} \cdot \frac{dz}{\partial b} = \frac{d\ell}{dz} \cdot 1.$$

# Outline

- Successful stories of DL (2 minutes)
- DL Basics (28 minutes)
- **Build two simple neural networks (45 minutes)**
- TensorFlow and Others (5 minutes)

# Build our first neural network

## Here is our task

Assume that we have one training example where  
 $(\mathbf{x} = [0.05, 0.1]^\top, \mathbf{y} = [0.01, 0.99]^\top)$ .

Given  $\mathbf{x}$ , we want to design a neural network such that it outputs  
 $[\mathbf{p}]_0 \approx 0.01, [\mathbf{p}]_1 \approx 0.99$ .

Obviously, this is a regression problem, we want to minimize the least square loss

$$\min \ell = \frac{1}{2} \|\mathbf{p} - \mathbf{y}\|^2 = \frac{1}{2} \sum_{i=0}^1 ([p]_i - [y]_i)^2$$

# Build our first neural network

## Here is our task

Assume that we have one training example where  
 $(\boldsymbol{x} = [0.05, 0.1]^\top, \boldsymbol{y} = [0.01, 0.99]^\top)$ .

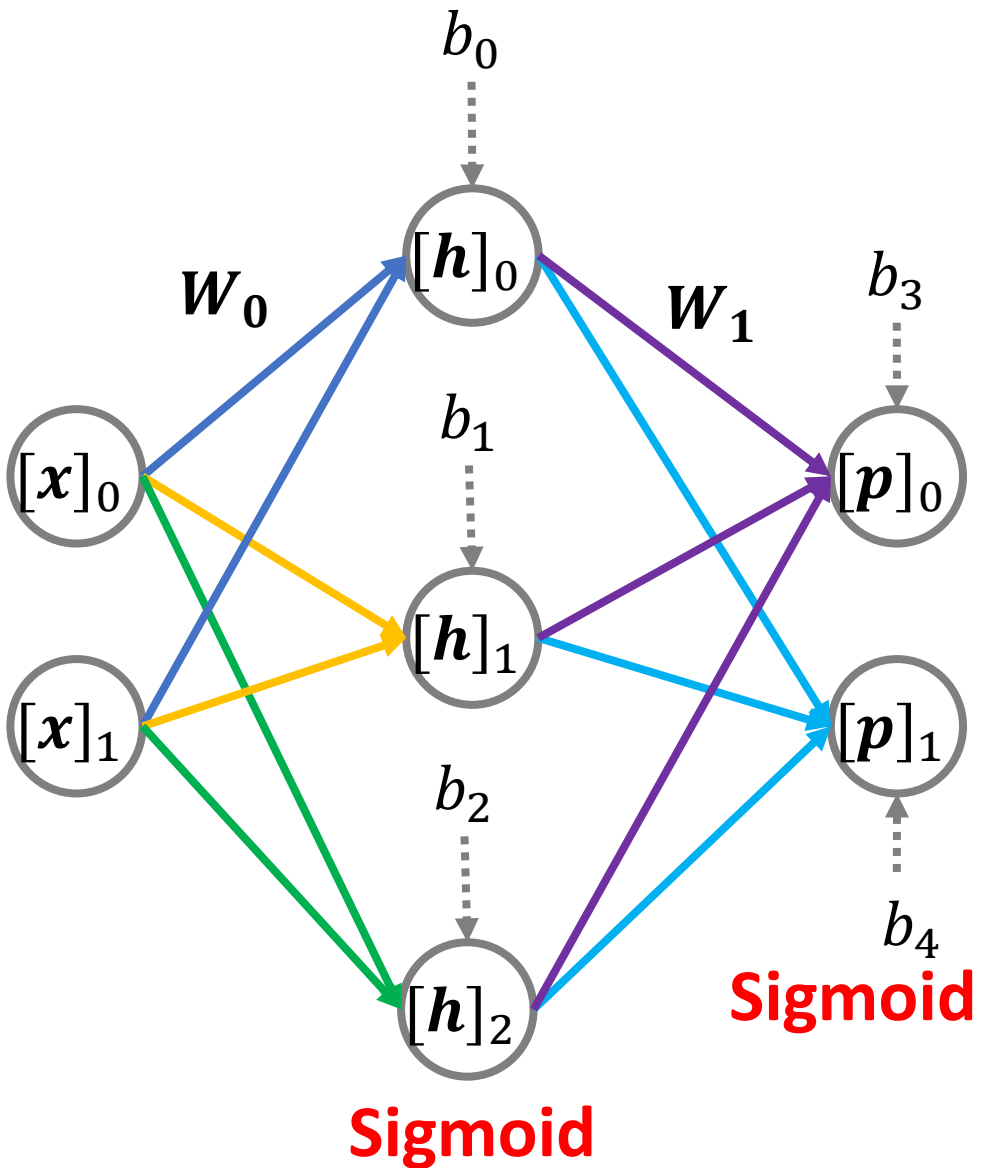
Given  $\boldsymbol{x}$ , we want to design a neural network such that it outputs  
 $[\boldsymbol{p}]_0 \approx 0.01, [\boldsymbol{p}]_1 \approx 0.99$ .

**To finish this task, we need to**

- 1. design a neural net**
- 2. implement this neural net**
- 3. load the dataset and train the neural net**
- 4. check the training losses**

# Build our first neural network

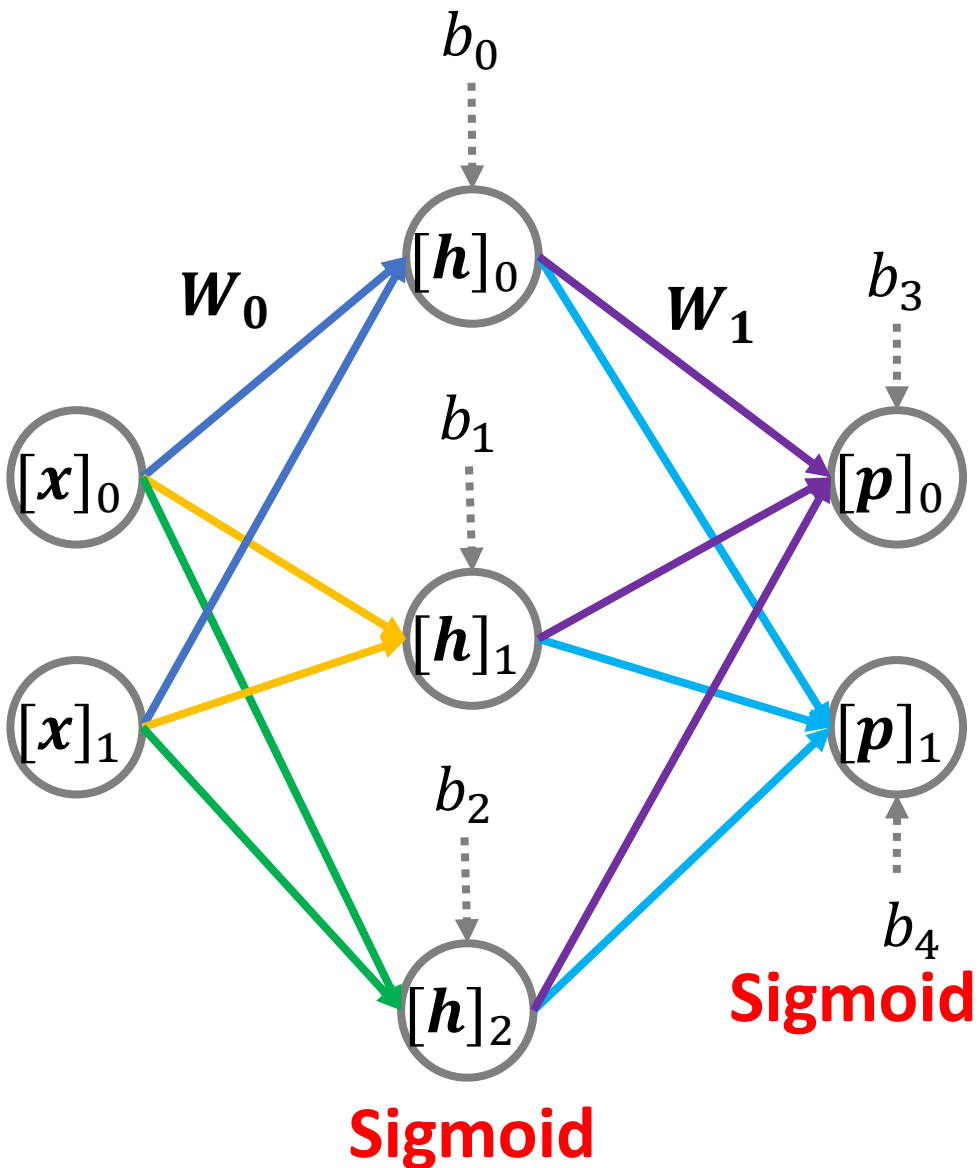
## 1. design a neural network





# Build our first neural network

## 1. design a neural network



More details:

$$W_0 = \begin{bmatrix} \mathbf{w}_0^T \\ \mathbf{w}_1^T \\ \mathbf{w}_2^T \end{bmatrix}, \quad W_1 = \begin{bmatrix} \mathbf{w}_3^T \\ \mathbf{w}_4^T \end{bmatrix}.$$

$$z_0 = [x]_0 [\mathbf{w}_0]_0 + [x]_1 [\mathbf{w}_0]_1 + b_0 = \mathbf{x}^T \mathbf{w}_0 + b_0,$$

$$z_1 = [x]_0 [\mathbf{w}_1]_0 + [x]_1 [\mathbf{w}_1]_1 + b_1 = \mathbf{x}^T \mathbf{w}_1 + b_1,$$

$$z_2 = [x]_0 [\mathbf{w}_2]_0 + [x]_1 [\mathbf{w}_2]_1 + b_2 = \mathbf{x}^T \mathbf{w}_2 + b_2,$$

$$[h]_0 = \sigma(z_0) = \frac{1}{1+e^{-z_0}},$$

$$[h]_1 = \sigma(z_1) = \frac{1}{1+e^{-z_1}},$$

$$[h]_2 = \sigma(z_2) = \frac{1}{1+e^{-z_2}},$$

$$z_3 = [h]_0 [\mathbf{w}_3]_0 + \dots + [h]_2 [\mathbf{w}_3]_2 + b_3 = \mathbf{h}^T \mathbf{w}_3 + b_3,$$

$$z_4 = [h]_0 [\mathbf{w}_4]_0 + \dots + [h]_2 [\mathbf{w}_4]_2 + b_4 = \mathbf{h}^T \mathbf{w}_4 + b_4,$$

$$[p]_0 = \sigma(z_3) = \frac{1}{1+e^{-z_3}},$$

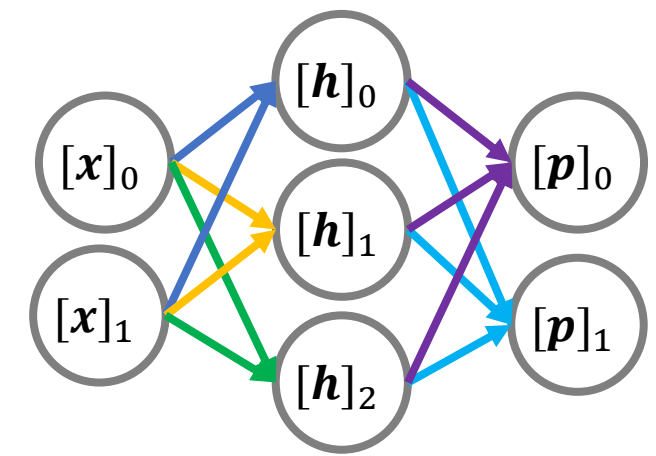
$$[p]_1 = \sigma(z_4) = \frac{1}{1+e^{-z_4}}.$$

# Build our first neural network

## 2. implement this neural network

We have the following information about **our network**:

1. The number of inputs ( i.e. dimension of the inputs).
2. The number of neurons in the **hidden** layer.
3. The number of neurons in the **output** layer.
4. Hidden layer and output layer use **Sigmoid** as the activation.
5. The loss function used is the **least square loss**.



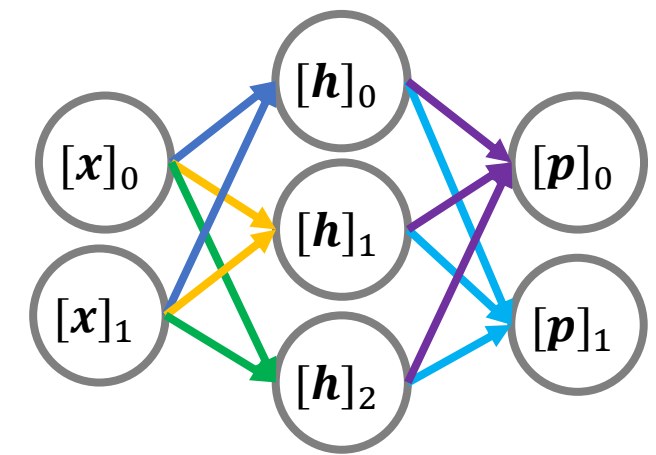
```
def neural_network_create(num_inputs, num_hidden, num_outputs,
                          hidden_layer_weights=None, hidden_layer_bias=None,
                          output_layer_weights=None, output_layer_bias=None):
    hidden_layer = {'num_neurons': num_hidden,
                    'neurons': [neuron_create(num_inputs, 'sigmoid', hidden_layer_bias, 1, _)
                                for _ in range(num_hidden)]}
    output_layer = {'num_neurons': num_outputs,
                   'neurons': [neuron_create(num_hidden, 'sigmoid', output_layer_bias, 2, _)
                                for _ in range(num_outputs)]}
```

# Build our first neural network

## 2. implement this neural network

We have the following information about **our neuron**:

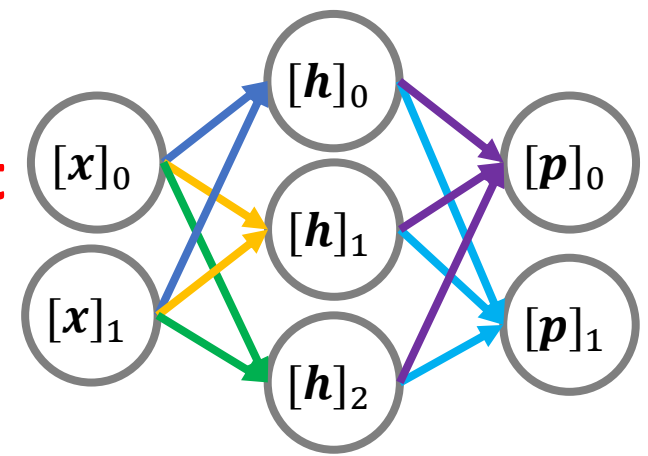
1. The number of inputs: 2.
2. Weight vector
3. Input vector ( dimension is 2 for hidden layer)
4. Bias
5. Output
6. Activation: Sigmoid



```
def neuron_create(num_inputs, activation, bias, layer_id, neuron_id):  
    return {'weights': np.zeros(num_inputs),  
            'inputs': np.zeros(num_inputs),  
            'bias': random.random() if bias is None else bias,  
            'output': 0.0,  
            'activation': activation,  
            'layer_id': layer_id,  
            'neuron_id': neuron_id}
```

# Build our first neural network

## 2. implement this neural network – **Constructed Neural Net**



```
neural_network = {'num_inputs': num_inputs,  
                  'num_hidden': num_hidden,  
                  'num_outputs': num_outputs,  
                  'hidden_layer': hidden_layer,  
                  'output_layer': output_layer}  
print('created a new neural networks and finished the initialization ...')  
return neural_network
```

After we construct the neural net, we need to have the initial weights ...

# Build our first neural network

## 2. implement this neural network - **Initialization**

initialize weights from inputs to hidden layer neurons

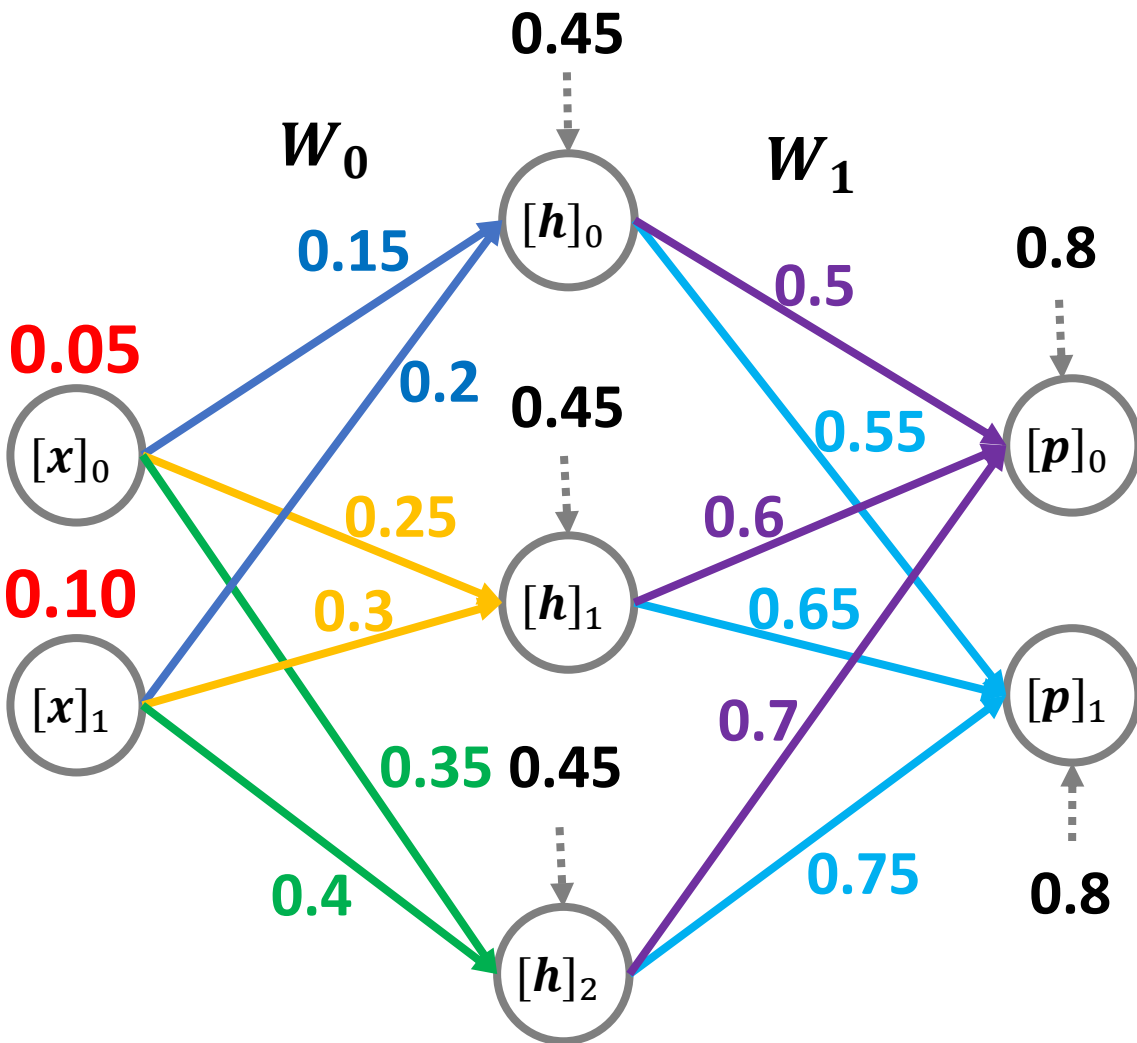
```
weight_num = 0
for h in range(hidden_layer['num_neurons']):
    for i in range(num_inputs):
        if not hidden_layer_weights:
            hidden_layer['neurons'][h]['weights'][i] = random.random()
        else:
            hidden_layer['neurons'][h]['weights'][i] = hidden_layer_weights[weight_num]
        weight_num += 1
```

initialize weights from hidden layer neurons to output layer neurons

```
weight_num = 0
for o in range(output_layer['num_neurons']):
    for h in range(hidden_layer['num_neurons']):
        if not output_layer_weights:
            output_layer['neurons'][o]['weights'][h] = random.random()
        else:
            output_layer['neurons'][o]['weights'][h] = output_layer_weights[weight_num]
        weight_num += 1
```

# Build our first neural network

## 2. implement this neural network – Calculate the square loss



Proposed neural net

$$W_0 = \begin{bmatrix} 0.15 & 0.2 \\ 0.25 & 0.3 \\ 0.35 & 0.4 \end{bmatrix}, \quad W_1 = \begin{bmatrix} 0.5 & 0.6 & 0.7 \\ 0.55 & 0.65 & 0.75 \end{bmatrix}.$$

$$z_0 = 0.05 * 0.15 + 0.10 * 0.2 + 0.45 = 0.4775,$$

$$z_1 = 0.05 * 0.25 + 0.10 * 0.3 + 0.45 = 0.4925,$$

$$z_2 = 0.05 * 0.35 + 0.10 * 0.4 + 0.45 = 0.5057,$$

$$[h]_0 = \frac{1}{1+e^{-z_0}} \approx ?,$$

$$[h]_1 = \frac{1}{1+e^{-z_1}} \approx ?,$$

$$[h]_2 = \frac{1}{1+e^{-z_2}} \approx ?,$$

$$z_3 \approx ?, z_4 \approx ?,$$

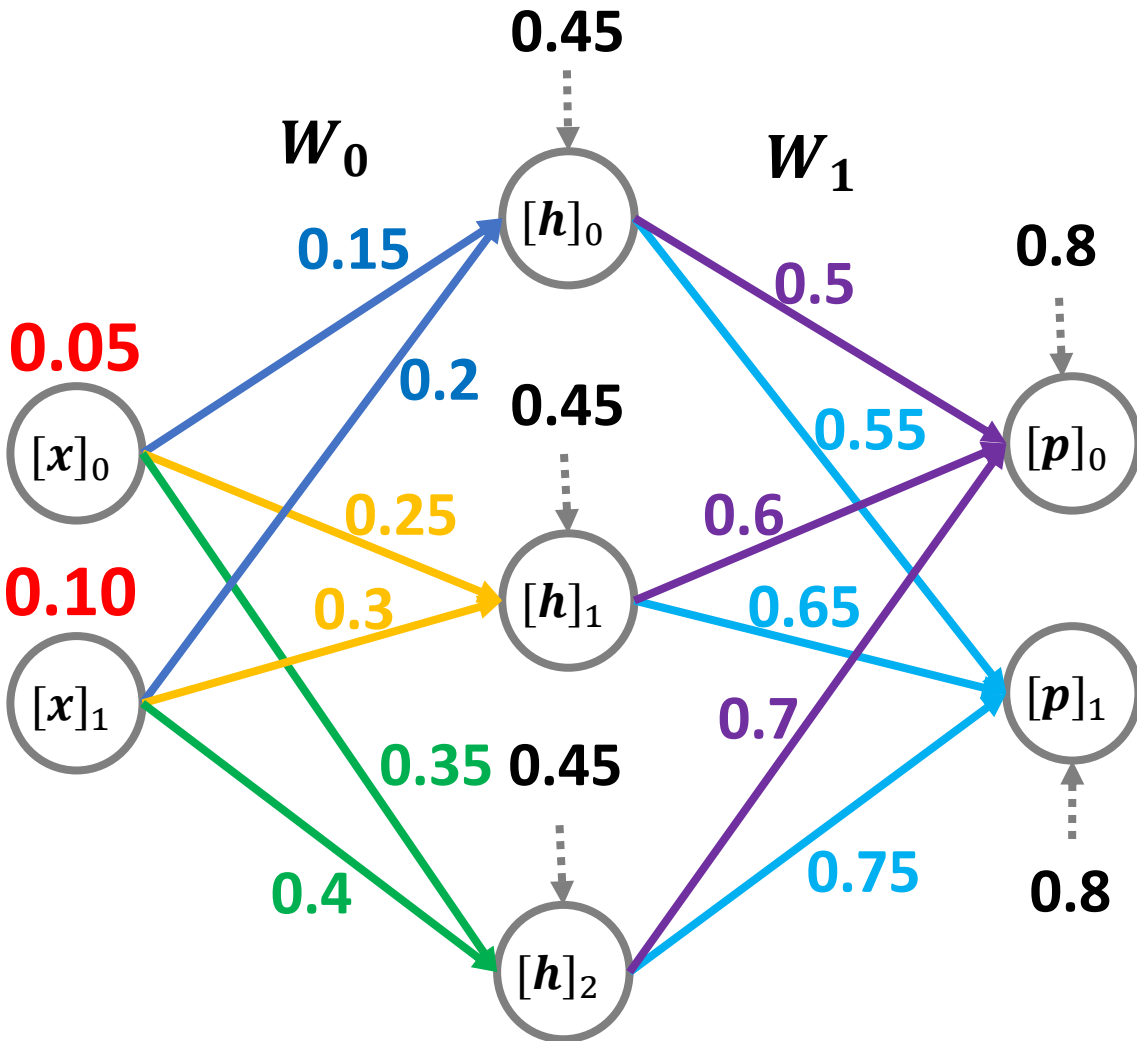
$$[p]_0 = \frac{1}{1+e^{-z_3}} \approx ?, [p]_1 = \frac{1}{1+e^{-z_4}} \approx ?$$

Quiz!

The loss  $\ell$  ?

# Build our first neural network

## 2. implement this neural network – Calculate the square loss



Proposed neural net

$$W_0 = \begin{bmatrix} 0.15 & 0.2 \\ 0.25 & 0.3 \\ 0.35 & 0.4 \end{bmatrix}, \quad W_1 = \begin{bmatrix} 0.5 & 0.6 & 0.7 \\ 0.55 & 0.65 & 0.75 \end{bmatrix}.$$

$$z_0 = 0.05 * 0.15 + 0.10 * 0.2 + 0.45 = 0.4775,$$

$$z_1 = 0.05 * 0.25 + 0.10 * 0.3 + 0.45 = 0.4925,$$

$$z_2 = 0.05 * 0.35 + 0.10 * 0.4 + 0.45 = 0.5057,$$

$$[h]_0 = \frac{1}{1+e^{-z_0}} \approx 0.6172,$$

$$[h]_1 = \frac{1}{1+e^{-z_1}} \approx 0.6207,$$

$$[h]_2 = \frac{1}{1+e^{-z_2}} \approx 0.6242,$$

$$z_3 \approx 1.8245, \quad z_4 \approx 2.1038,$$

$$[p]_0 = \frac{1}{1+e^{-z_3}} \approx 0.8611, \quad [p]_1 = \frac{1}{1+e^{-z_4}} \approx 0.8913.$$

The loss  $\ell \approx 0.3671$



# Build our first neural network

2. implement this neural network – **Calculate the square loss**

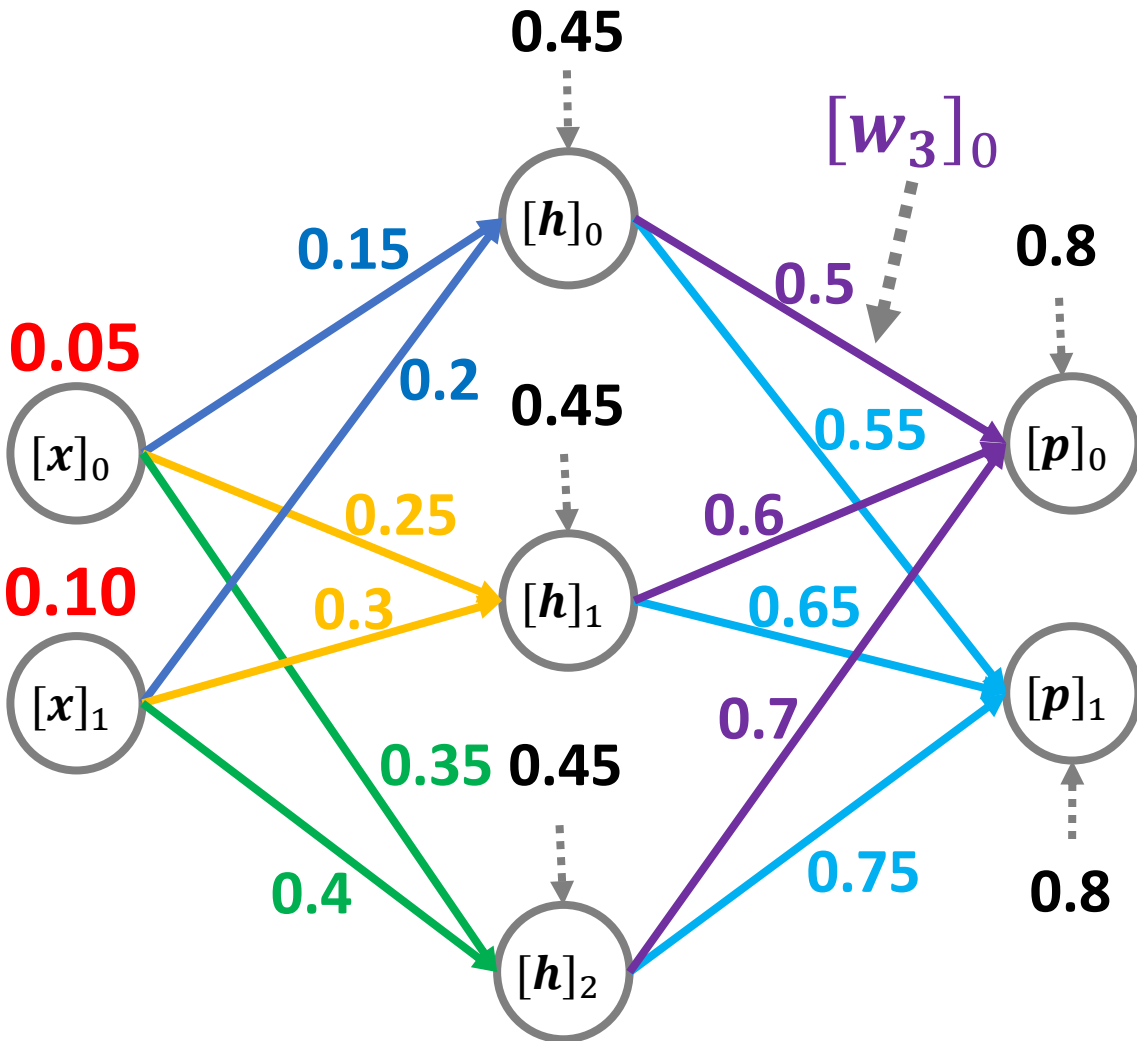
## Feed Forward

```
def neural_network_feed_forward(neural_network, train_xi):  
    hidden_layer = neural_network['hidden_layer']  
    output_layer = neural_network['output_layer']  
    # hidden layer feeds forward activation.  
    hidden_layer_outputs = np.zeros(hidden_layer['num_neurons'])  
    for index, neuron in enumerate(hidden_layer['neurons']):  
        hidden_layer_outputs[index] = neuron_cal_output(neuron, train_xi)  
    # output layer feeds forward activation.  
    outputs = np.zeros(output_layer['num_neurons'])  
    for index, neuron in enumerate(output_layer['neurons']):  
        outputs[index] = neuron_cal_output(neuron, hidden_layer_outputs)  
    # final prediction  
    return outputs
```

Given the squared loss, we can how  
calculate the gradient w.r.t  $w$  ?

# Build our first neural network

## 2. implement this neural network – Calculate the gradient



Proposed neural net

Given the loss, how can we calculate the gradient w.r.t  $W_1, W_2, b$  ?

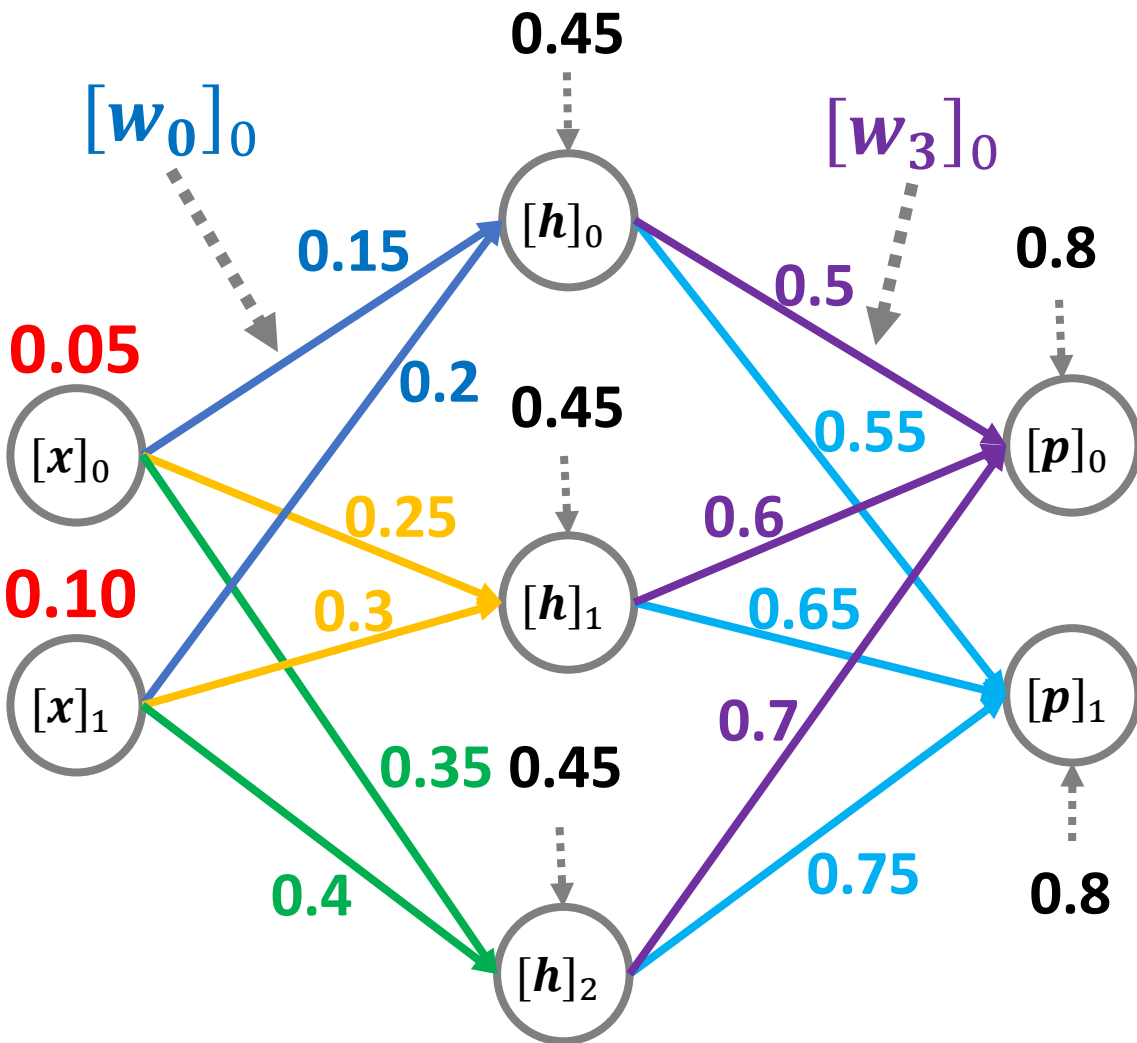
### Backpropagation!

For example, we to calculate the gradient of  $\ell$  w.r.t.  $[w_3]_0$ .

$$\begin{aligned}\frac{\partial \ell}{\partial [w_3]_0} &= \frac{\partial \ell}{\partial [p]_0} \cdot \frac{\partial [p]_0}{\partial z_3} \cdot \frac{\partial z_3}{\partial [w_3]_0} \\ &= ([p]_0 - [y]_0) \cdot (\sigma(z_3)(1 - \sigma(z_3))) [h]_0.\end{aligned}$$

# Build our first neural network

## 2. implement this neural network – Calculate the gradient



Proposed neural net

Given the loss, how can we calculate the gradient w.r.t  $W_1, W_2, b$  ?

### Backpropagation!

For example, we to calculate the gradient of  $\ell$  w.r.t.  $[w_3]_0$ .

$$\begin{aligned}\frac{\partial \ell}{\partial [w_3]_0} &= \frac{\partial \ell}{\partial [p]_0} \cdot \frac{\partial [p]_0}{\partial z_3} \cdot \frac{\partial z_3}{\partial [w_3]_0} \\ &= ([p]_0 - [y]_0) \cdot \sigma(z_3)(1 - \sigma(z_3)) [h]_0 \\ &= ([p]_0 - [y]_0) \cdot [p]_0(1 - [p]_0) [h]_0\end{aligned}$$

Use chain rule, to calculate

$$\frac{\partial \ell}{\partial [w_0]_0} = 0.0006$$

Quiz!

# Build our first neural network

## 2. implement this neural network – Update the weights

Update all weights by using stochastic gradient descent

$$[w_i]_j^{t+1} = [w_i]_j^t - \eta \frac{\partial \ell}{\partial [w_i]_j^t}, \text{ where } \eta = 0.5$$

```
# 3. Update output neuron weights
for o in range(output_layer['num_neurons']):
    for w_ho in range(len(output_layer['neurons'][o]['weights'])):
        #  $\partial E_j / \partial w_{ij} = \partial E / \partial z_j * \partial z_j / \partial w_{ij}$ 
        pd = neuron_cal_pd_total_net_input_wrt_weight(output_layer['neurons'][o], w_ho)
        pd_error_wrt_weight = pd_errors_wrt_output_neuron_total_net_input[o] * pd
        #  $\Delta w = \alpha * \partial E_j / \partial w_{ij}$ 
        output_layer['neurons'][o]['weights'][w_ho] -= learning_rate * pd_error_wrt_weight
```

```
# 4. Update hidden neuron weights
for h in range(hidden_layer['num_neurons']):
    for w_ih in range(len(hidden_layer['neurons'][h]['weights'])):
        #  $\partial E_j / \partial w_{ij} = \partial E / \partial z_j * \partial z_j / \partial w_{ij}$ 
        pd = neuron_cal_pd_total_net_input_wrt_weight(hidden_layer['neurons'][h], w_ih)
        pd_error_wrt_weight = pd_errors_wrt_hidden_neuron_total_net_input[h] * pd
        #  $\Delta w = \alpha * \partial E_j / \partial w_{ij}$ 
        hidden_layer['neurons'][h]['weights'][w_ih] -= learning_rate * pd_error_wrt_weight
```

# Build our first neural network

## 3. load the dataset and train the neural network

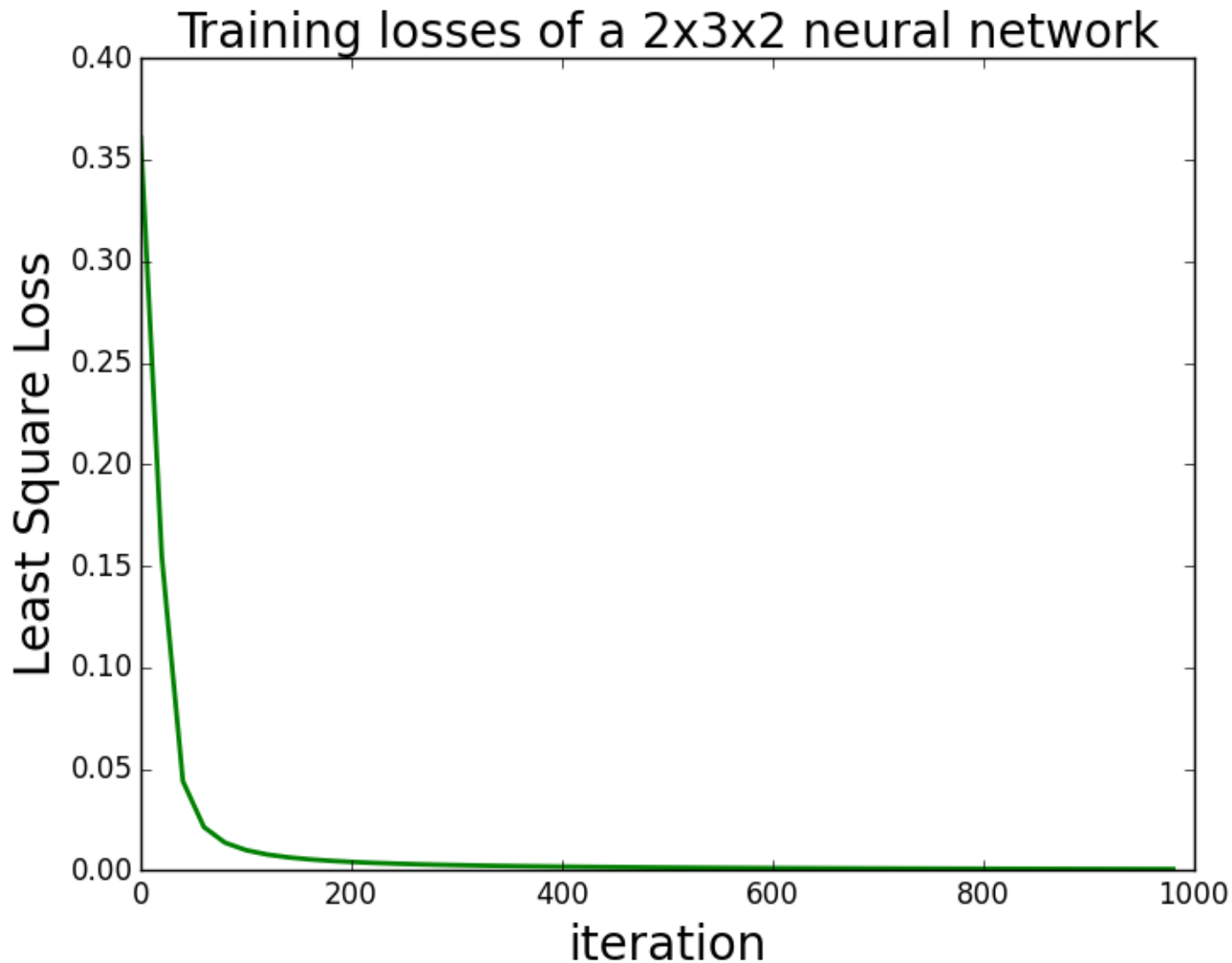
```
train_x = [[0.05, 0.1]]
train_y = [[0.01, 0.99]]
for i in range(1000):
    train_xi, train_yi = train_x[0], train_y[0]
    neural_network_train(nn, train_x=train_xi, train_y=train_yi, learning_rate=.5)
    if i % 5 == 0:
        cur_loss = neural_network_total_error(nn, train_x=train_x, train_y=train_y)
        train_losses.append(cur_loss)
    print(i, train_losses[-1])
```

How can we decide whether this neural network is a good one? There is no a standard criteria, but we can at least take a look at the losses during the training process.

**Important: the total training error will decrease in general.**

# Build our first neural network

## 4. check the training losses



As we can see, during the learning process, the total loss is always decreasing.

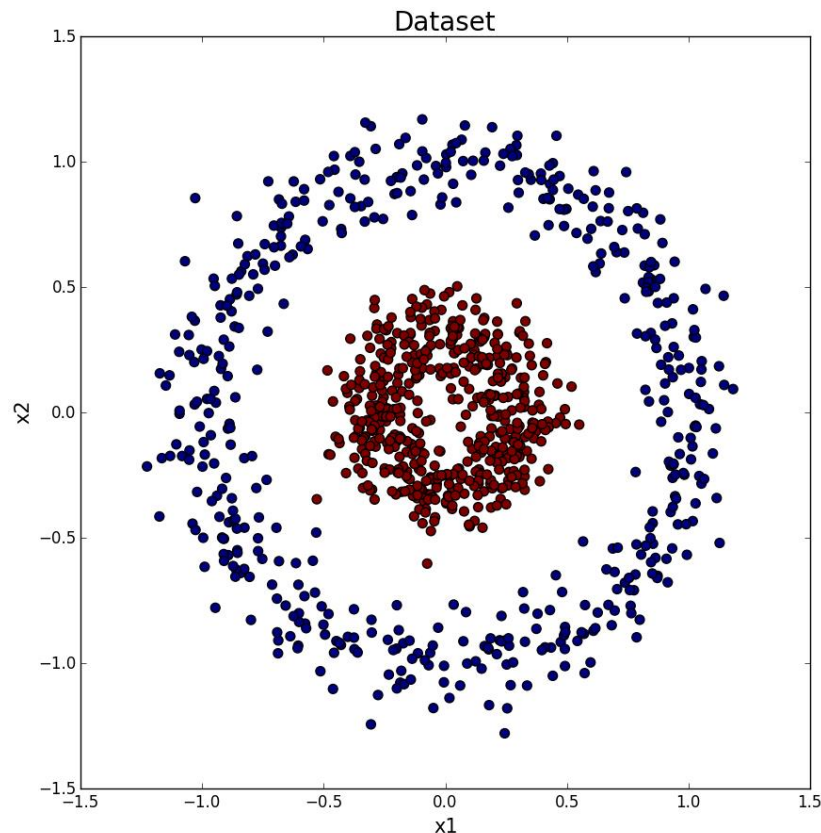
**Quiz!**

Try to run the program!

# Build our second neural network

## Here is our task:

Suppose, we have a classification problem.  
Each  $(x_i, y_i)$  is a training sample from the  
following distribution



We want to minimize the cross-entropy loss

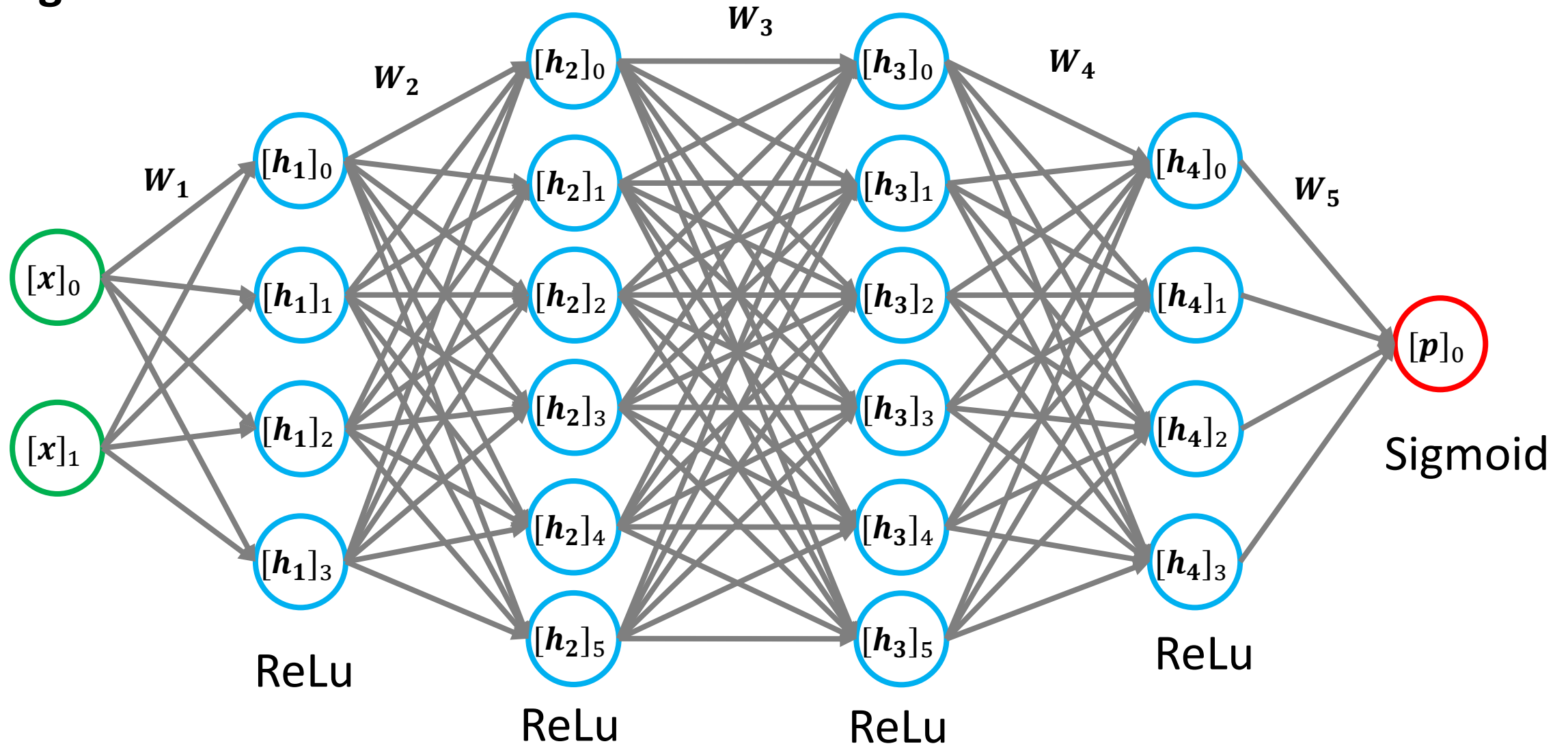
$$\ell = \frac{1}{n} \sum_{i=1}^n -y_i \log p_i$$

To finish this task, we need to

1. design a neural network
2. implement this neural network
3. load the dataset and train the neural net
4. check the training losses

# Build our second neural network

## 1. design a neural network





# Build our second neural network

## 2. implement this neural network

```
def neural_network_create(input_dim, list_node_sizes, output_dim=2):  
    model = dict()  
    # first hidden layer  
    model['W1'] = np.random.randn(input_dim, list_node_sizes[0])  
    model['b1'] = np.zeros((1, list_node_sizes[0]))  
    # second hidden layer  
    model['W2'] = np.random.randn(list_node_sizes[0], list_node_sizes[1])  
    model['b2'] = np.zeros((1, list_node_sizes[1]))  
    # third hidden layer  
    model['W3'] = np.random.randn(list_node_sizes[1], list_node_sizes[2])  
    model['b3'] = np.zeros((1, list_node_sizes[2]))  
    # fourth hidden layer  
    model['W4'] = np.random.randn(list_node_sizes[2], list_node_sizes[3])  
    model['b4'] = np.zeros((1, list_node_sizes[3]))  
    # output layer  
    model['W5'] = np.random.randn(list_node_sizes[3], output_dim)  
    model['b5'] = np.zeros((1, output_dim))  
    return model
```

# Build our second neural network

## 3. load the dataset and train the neural network

```
def test_batch(x_tr, y_tr, x_te, y_te):
    model = neural_network_create(input_dim=x_tr.shape[1], list_node_sizes=[4, 6, 6, 4],
    output_dim=2)
    train_batch(model, x_tr, y_tr, num_passes=50, learning_rate=0.001)
    output = neural_network_feed_forward(model, x_te)
    success = 0
    for ind, item in enumerate(output[-1]):
        if y_te[ind] == np.argmax(item):
            success += 1
    print(success / float(len(y_te)))

def test_stochastic(x_tr, y_tr, x_te, y_te):
    model = neural_network_create(input_dim=x_tr.shape[1], list_node_sizes=[4, 6, 6, 4],
    output_dim=2)
    train_stochastic(model=model, x_tr=x_tr, y_tr=y_tr, num_passes=50, batch_size=20,
    learning_rate=0.001)
    output = neural_network_feed_forward(model, x_te)
    success = 0
    for ind, item in enumerate(output[-1]):
        if y_te[ind] == np.argmax(item):
            success += 1
    print('test accuracy: %.4f' % (success / float(len(y_te))))
```

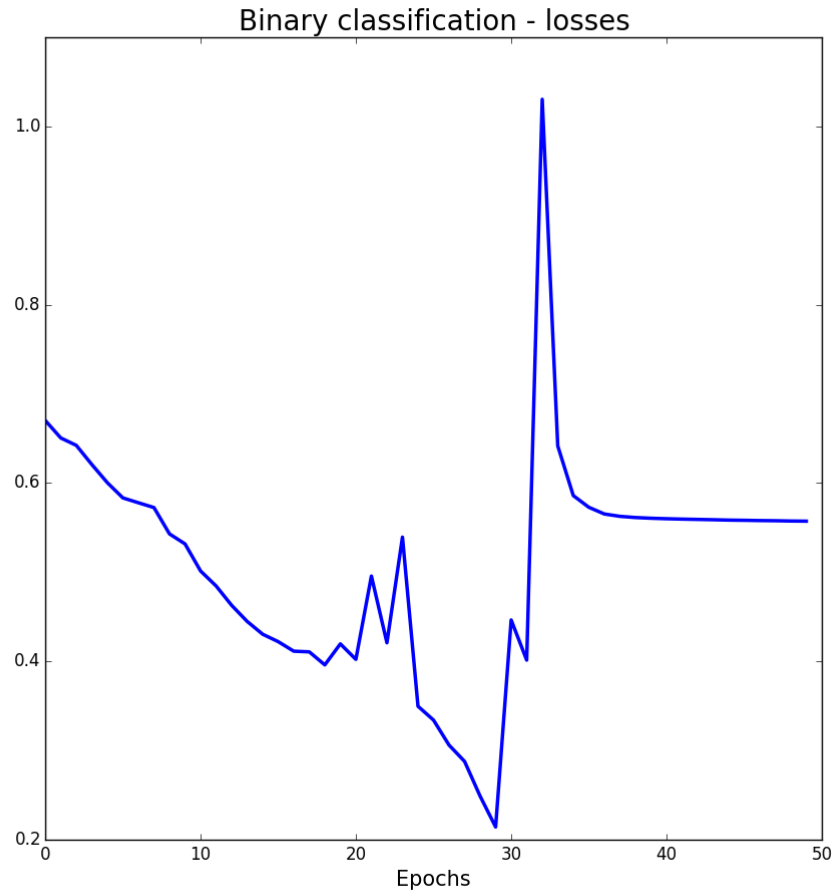
Try to use two different ways:

1. Batch training
2. Stochastic training

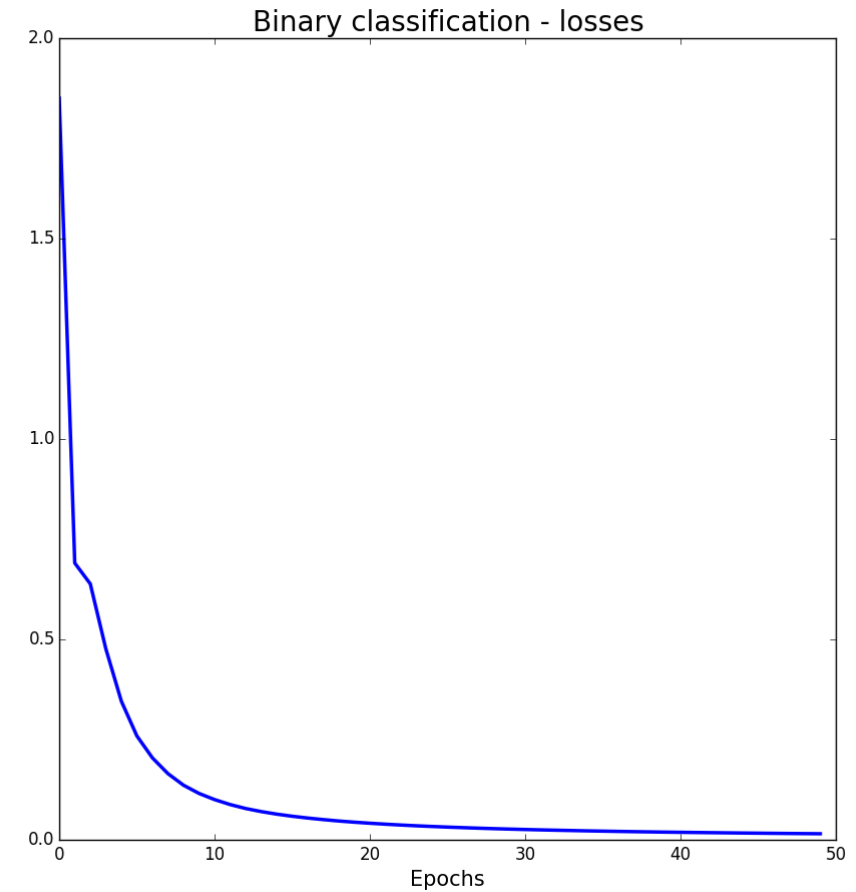
# Build our second neural network

## 4. check the training losses

### Batch training



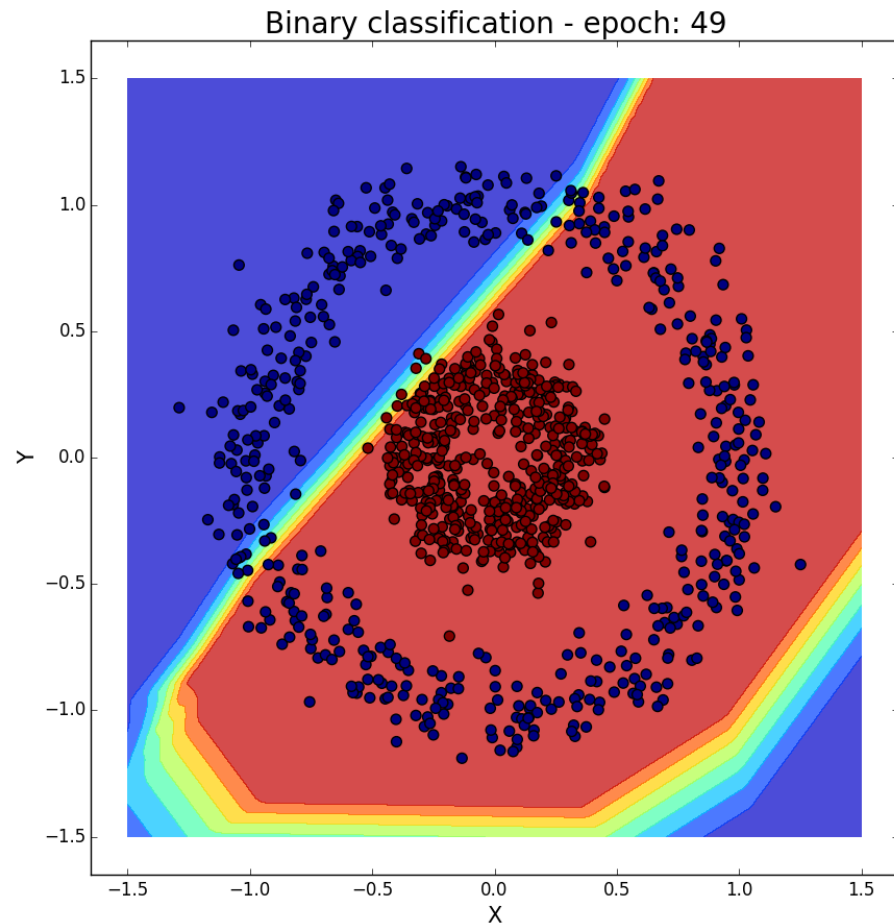
### Stochastic training



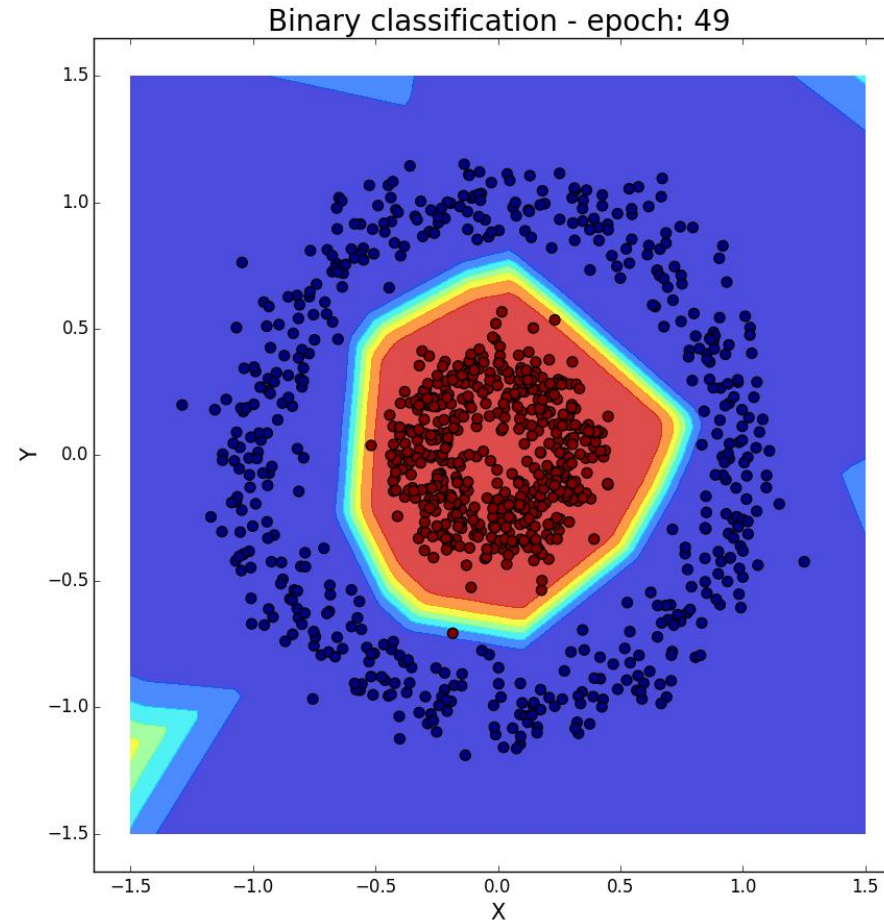
# Build our second neural network

## 4. check the model

### Batch training



### Stochastic training



**Quiz!**

**Try to run the code  
several times to find  
a potential pattern.**

# Outline

- Successful stories of DL (2 minutes)
- DL Basics (28 minutes)
- Build two simple neural networks (45 minutes)
- **TensorFlow and Others (5 minutes)**

# TensorFlow – Overview

tensorflow / tensorflow

Used by 53.7k Watch 8.6k Unstar 137k Fork 78.3k

Code Issues 2,780 Pull requests 229 Projects 1 Security Insights

An Open Source Machine Learning Framework for Everyone <https://tensorflow.org>

tensorflow machine-learning python deep-learning deep-neural-networks neural-network ml distributed

71,170 commits 37 branches 97 releases 2,261 contributors Apache-2.0

C++ 60.7% Python 30.6% HTML 3.7% Go 1.3% MLIR 0.9% Java 0.8% Other 2.0%

Branch: master New pull request Create new file Upload files Find file Clone or download

## Keras

keras-team / keras

Used by 35.5k Watch 2.1k Unstar 45.2k Fork 17.2k

Code Issues 2,630 Pull requests 22 Projects 1 Wiki Security Insights

Deep Learning for humans <http://keras.io/>

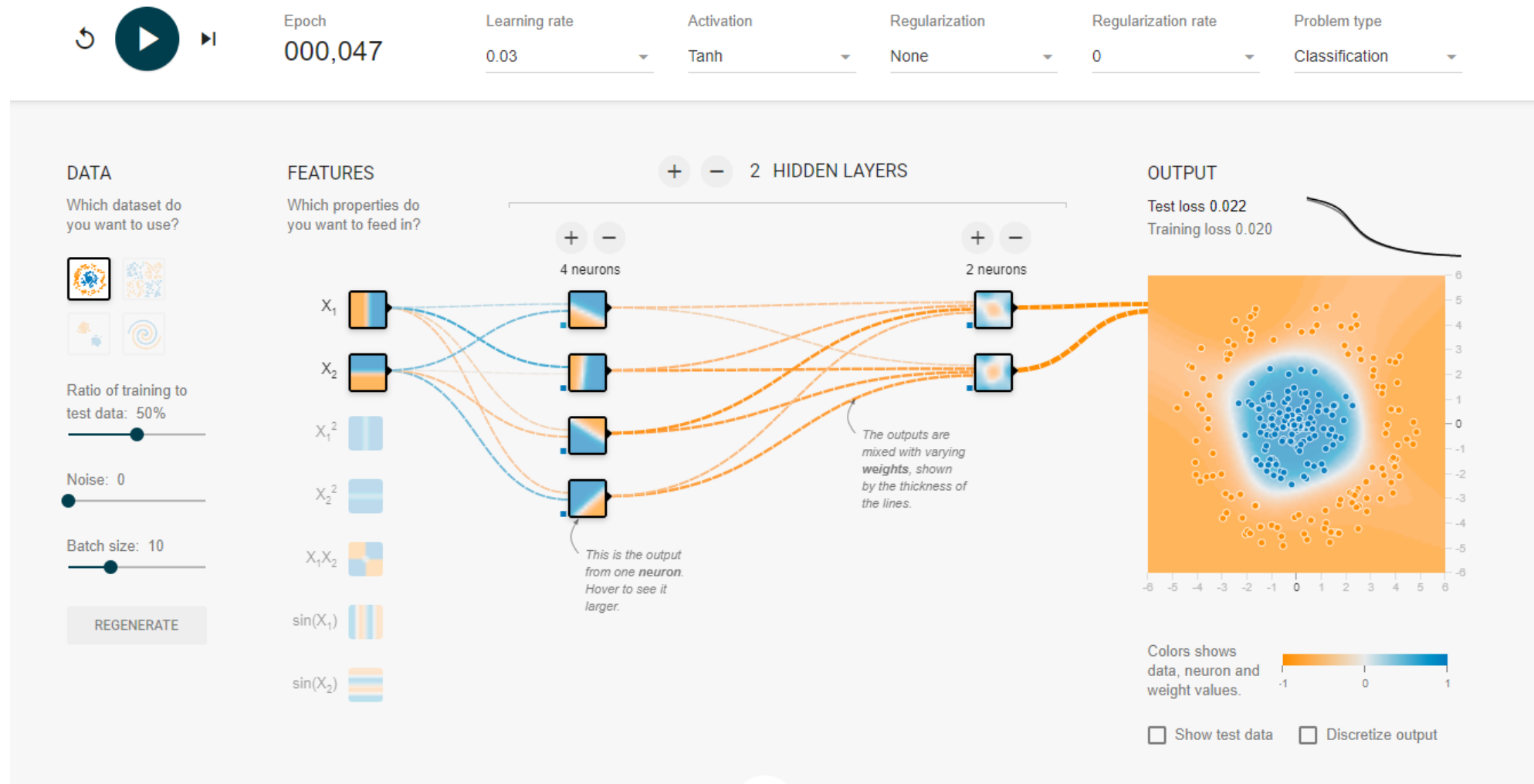
deep-learning tensorflow neural-networks machine-learning data-science python

5,341 commits 6 branches 49 releases 825 contributors View license

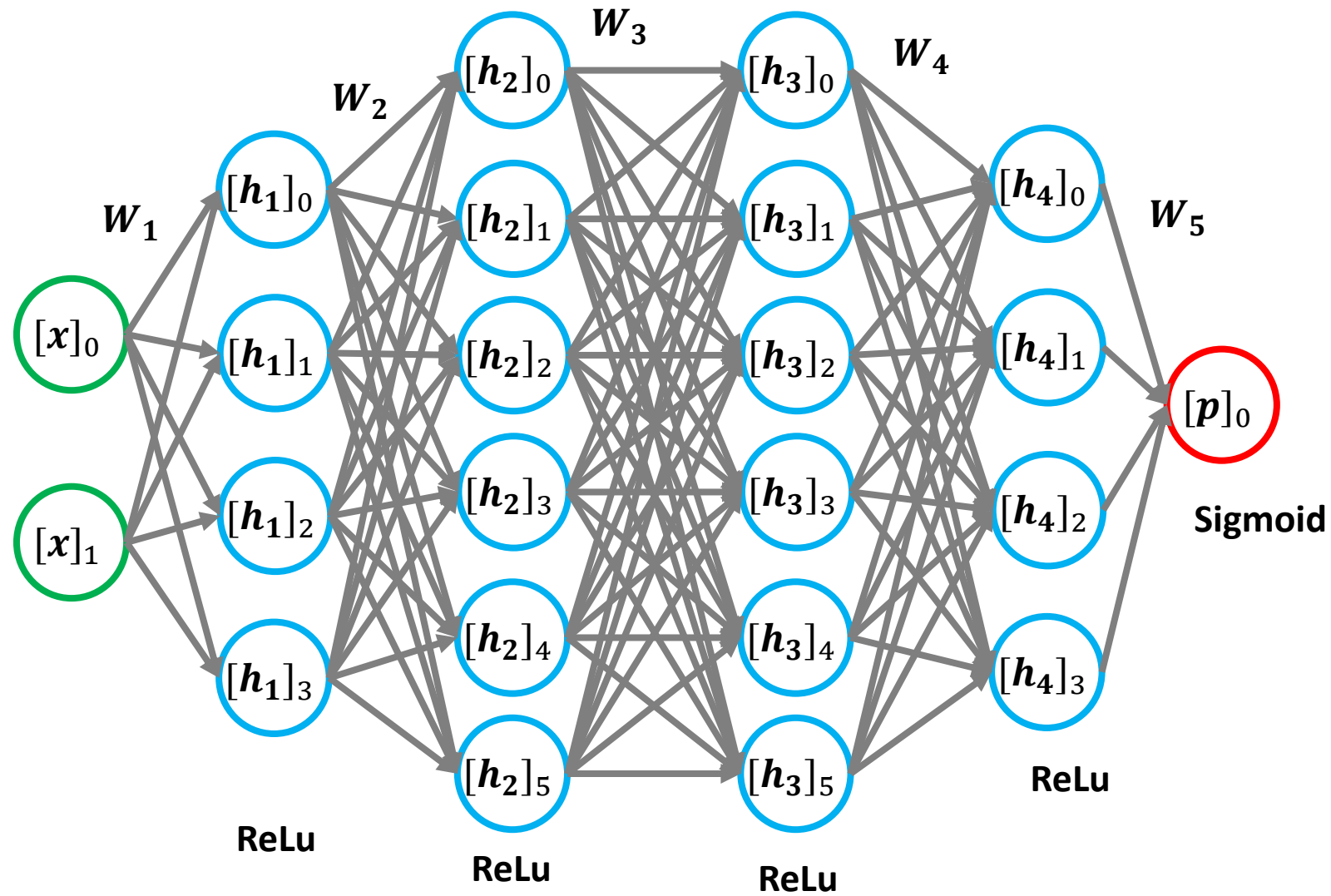
Python 99.9% Other 0.1%

# TensorFlow – Playground

Have fun: <https://playground.tensorflow.org>

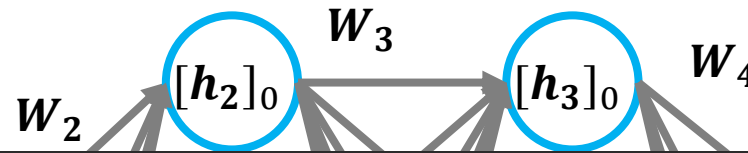


# TensorFlow – Build a deep learning model






# TensorFlow – Build a deep learning model



```
from keras.layers import Dense
from keras.models import Sequential
model = Sequential()
model.add(Dense(4, input_dim=2, activation='relu'))
model.add(Dense(6, activation='relu'))
model.add(Dense(6, activation='relu'))
model.add(Dense(4, activation='relu'))
model.add(Dense(1, activation='sigmoid'))
```



# PyTorch - <https://github.com/pytorch/pytorch>

 **pytorch / pytorch**

Watch 1.4k

Star 33.4k

Fork 8.3k

Code

Issues 3,522

Pull requests 1,030

Actions

Projects 5

Wiki

Security

Insights

Tensors and Dynamic neural networks in Python with strong GPU acceleration <https://pytorch.org>

neural-network

autograd

gpu

numpy

deep-learning

tensor

python

machine-learning

22,068 commits

2,629 branches

29 releases

1,225 contributors

View license

C++ 51.5%

Python 31.6%

Cuda 7.7%

C 5.2%

CMake 1.8%

Objective-C++ 0.7%

Other 1.5%

Branch: master

New pull request

Create new file

Upload files

Find file

Clone or download

# References

- <https://github.com/HarisIqbal88/PlotNeuralNet>
- <http://www.emergentmind.com/neural-network>
- [http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture4.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture4.pdf)
- <https://medium.com/@dcharrezt/neurips-2019-stats-c91346d31c8f> [word cloud, NIPS, 2019]
- [https://www.reddit.com/r/MachineLearning/comments/9jhhvb/d\\_iclr\\_2019\\_submissions\\_are\\_viewable\\_which\\_ones/](https://www.reddit.com/r/MachineLearning/comments/9jhhvb/d_iclr_2019_submissions_are_viewable_which_ones/) [word cloud, ICLR, 2019]
- Jeff Dean's Lecture for YC AI: <https://blog.ycombinator.com/jeff-deans-lecture-for-yc-ai/>
- Page 4, AlphaGo: <https://intellipaat.com/blog/power-of-deep-learning-alphago-vs-lee-sedol-case-study/>

# A Timeline of AI

- 1943: First neuron, Walter Pitts, a logician, and Warren McCulloch, “A logical calculus of the ideas immanent in nervous activity”.
- 1945: ENIAC, in full Electronic Numerical Integrator and Computer, the first programmable general-purpose electronic digital computer.
- 1950: Alan Turing, “Computing Machinery and Intelligence”, Can machines think?
- 1952: **Machine Learning**, Arthur Lee Samuel, “Some Studies in Machine Learning Using the Game of Checkers. II—Recent Progress”, (1959)
- 1957: Frank Rosenblatt, Rosenblatt, a psychologist, submitted a paper entitled “The Perceptron: A Perceiving and Recognizing Automaton” to Cornell Aeronautical Laboratory in 1957.