

Deep Neural Network Foundation

Hao Dong

2019, Peking University

Deep Neural Network Foundation

- Single Neuron
- Activation Functions
- Multi-layer Perceptron
- Loss Functions
- Optimization
- Regularization
- Implementation

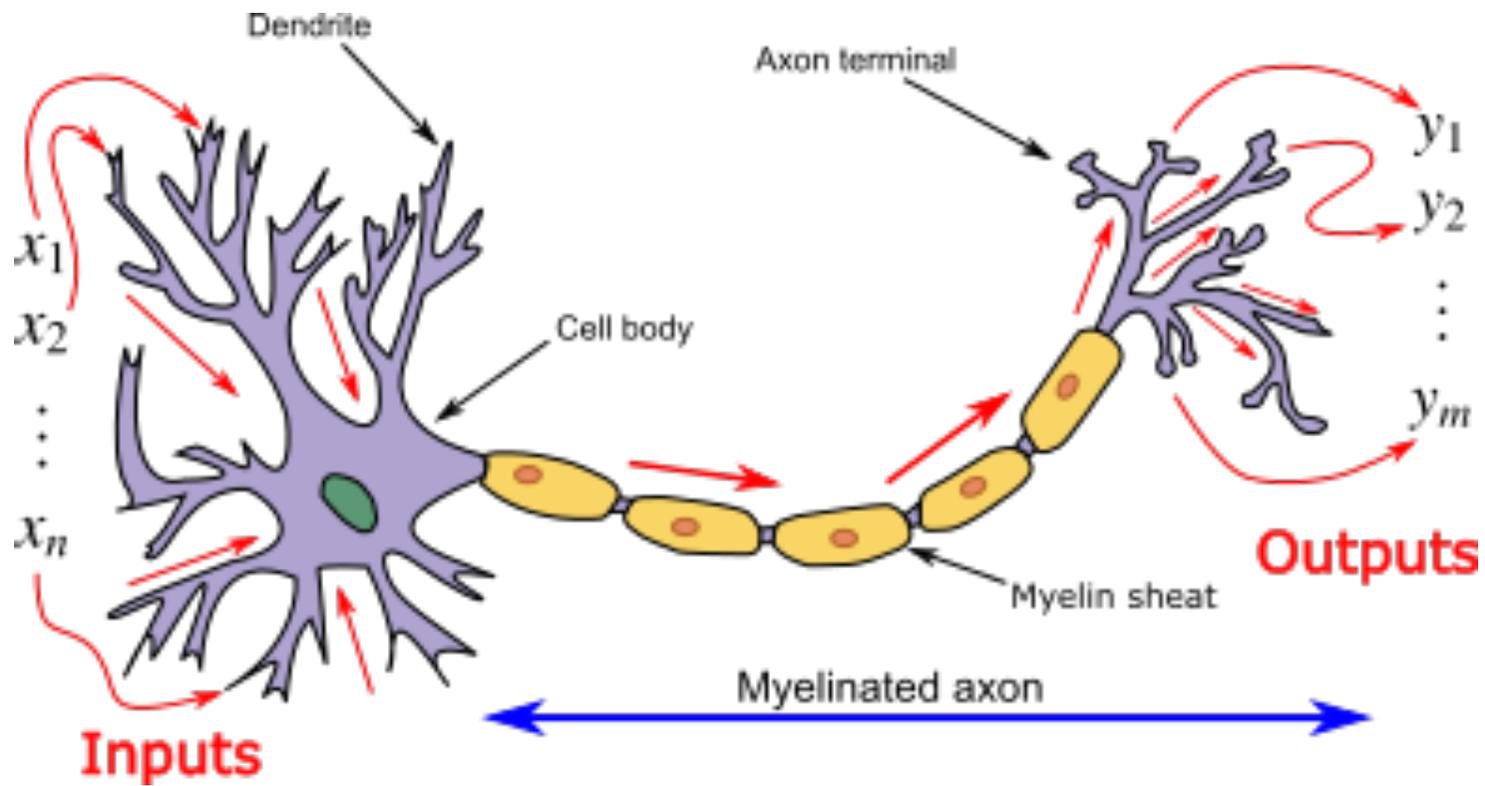
Goal: Understand the basic knowledge of deep neural networks

Single Neuron



Single Neuron

- Motivation



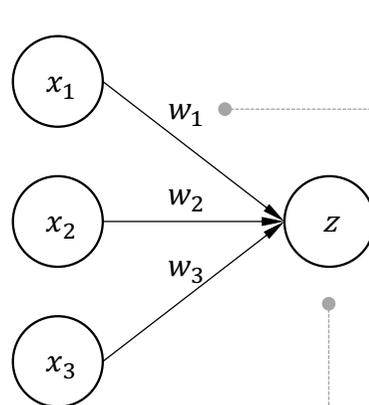
Single Neuron

- 3 inputs and 1 output

The output is a linear combination of the inputs

$$z = x_1w_1 + x_2w_2 + x_3w_3$$

input layer output layer



A single neuron

- A larger absolute value of the **weight** means the output more sensitive to the specific input

For example, z may be a score determining if we are to play football. If z is large, then we play. To determine this score, x_1 represents the weather, x_2 is the expense of the football field rental, and x_3 is the distance to the field. These inputs are considered the features w.r.t the output. If the weather is the most critical factor, then we can set w_1 to a large positive value and set w_2 and w_3 to smaller positive values. If any w is set to zero, then the corresponding input feature is ignored.

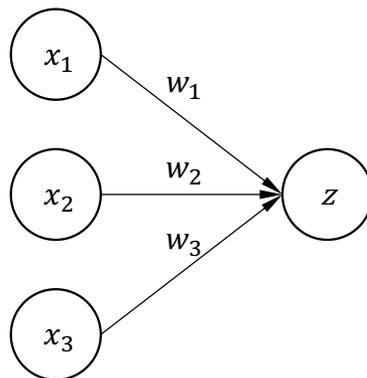
Single Neuron

- 3 inputs and 1 output

The output is a linear combination of the inputs

$$z = x_1w_1 + x_2w_2 + x_3w_3$$

input layer *output layer*



- A single neuron is a network that has only one (output) layer and one output
- As the output connected to all inputs, this layer is called “fully connected layer” or “dense layer”

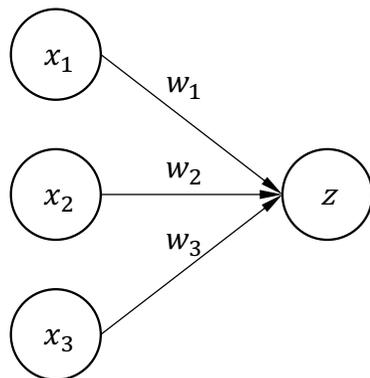
Single Neuron

- 3 inputs and 1 output

This neuron can be represented by a matrix multiplication

$$z = x_1w_1 + x_2w_2 + x_3w_3$$

input layer output layer



column format

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$z = \mathbf{w}^T \mathbf{x}$$

$$z = [w_1 \quad w_2 \quad w_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

row format

$$\mathbf{x} = [x_1 \quad x_2 \quad x_3] \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$z = \mathbf{xw}$$

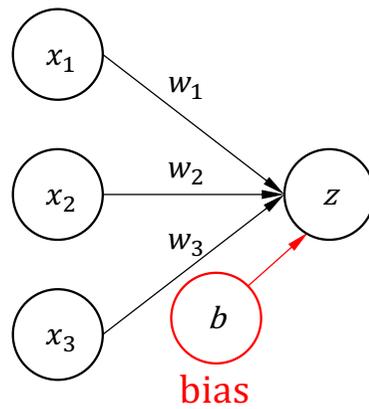
$$z = [x_1 \quad x_2 \quad x_3] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

Single Neuron

- Bias

A **bias** value allows the output value to be shifted higher or lower to better fit the input data.

input layer output layer



$$z = x_1w_1 + x_2w_2 + x_3w_3 + b$$

column format

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$z = \mathbf{w}^T \mathbf{x} + b$$

$$z = [w_1 \quad w_2 \quad w_3] \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + b$$

row format

$$\mathbf{x} = [x_1 \quad x_2 \quad x_3] \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$z = \mathbf{xw} + b$$

$$z = [x_1 \quad x_2 \quad x_3] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + b$$

Single Neuron

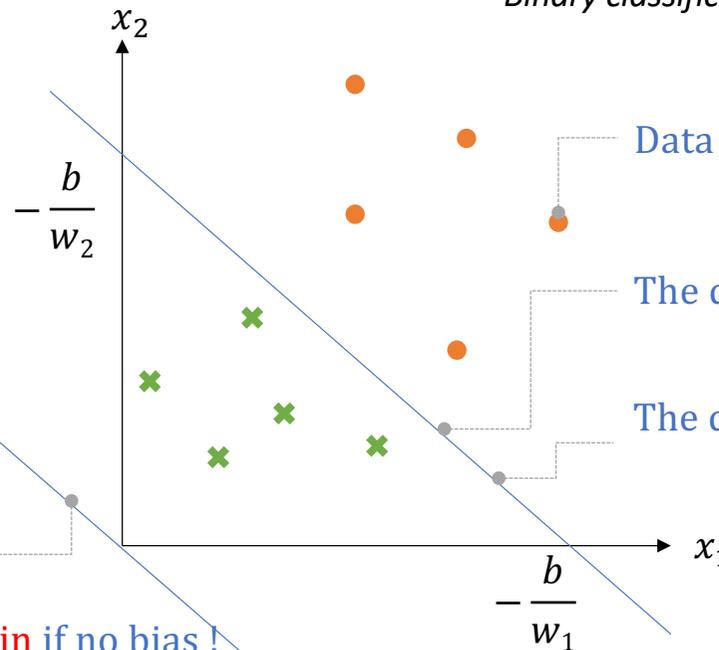
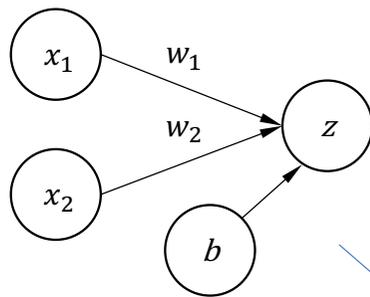
- Classification

A bias value allows the output value to be shifted higher or lower to better fit the input data.

$$z = x_1w_1 + x_2w_2 + b$$

$$\text{Binary classification: } a = \begin{cases} 0, & \text{if } z \leq 0 \\ 1, & \text{if } z > 0 \end{cases}$$

input layer output layer



Data samples with **two** features (x_1, x_2)

The decision boundary is a **line** for $z = 0$

The decision boundary can be shifted left or right via the bias

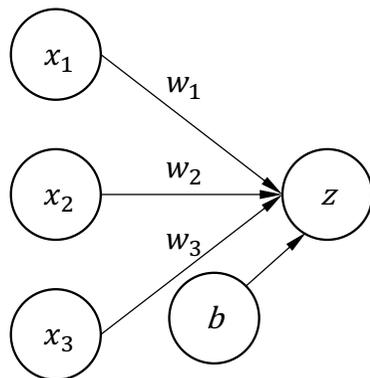
The decision boundary must cross the **origin** if no bias !

Single Neuron

- Classification

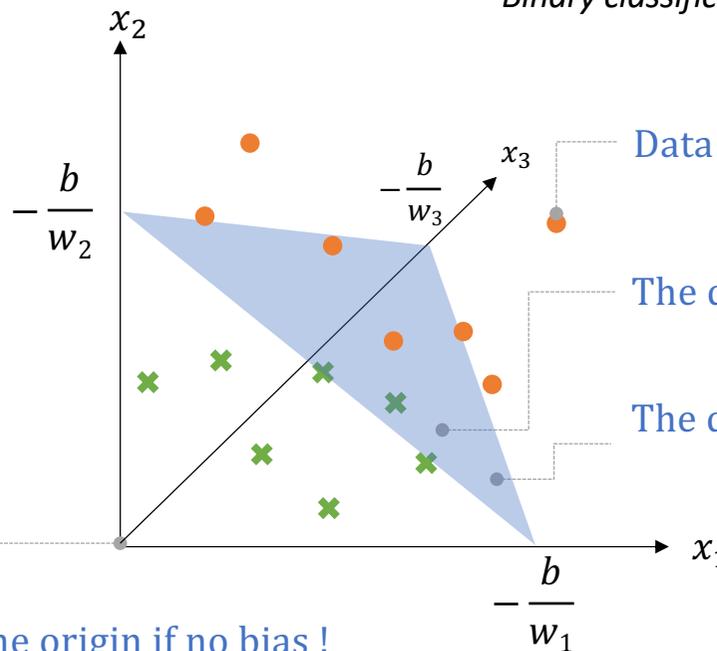
A bias value allows the output value to be shifted higher or lower to better fit the input data.

input layer output layer



$$z = x_1w_1 + x_2w_2 + x_3w_3 + b$$

$$\text{Binary classification: } a = \begin{cases} 0, & \text{if } z \leq 0 \\ 1, & \text{if } z > 0 \end{cases}$$



Data samples with **three** features (x_1, x_2, x_3)

The decision boundary is a **surface** for $z = 0$

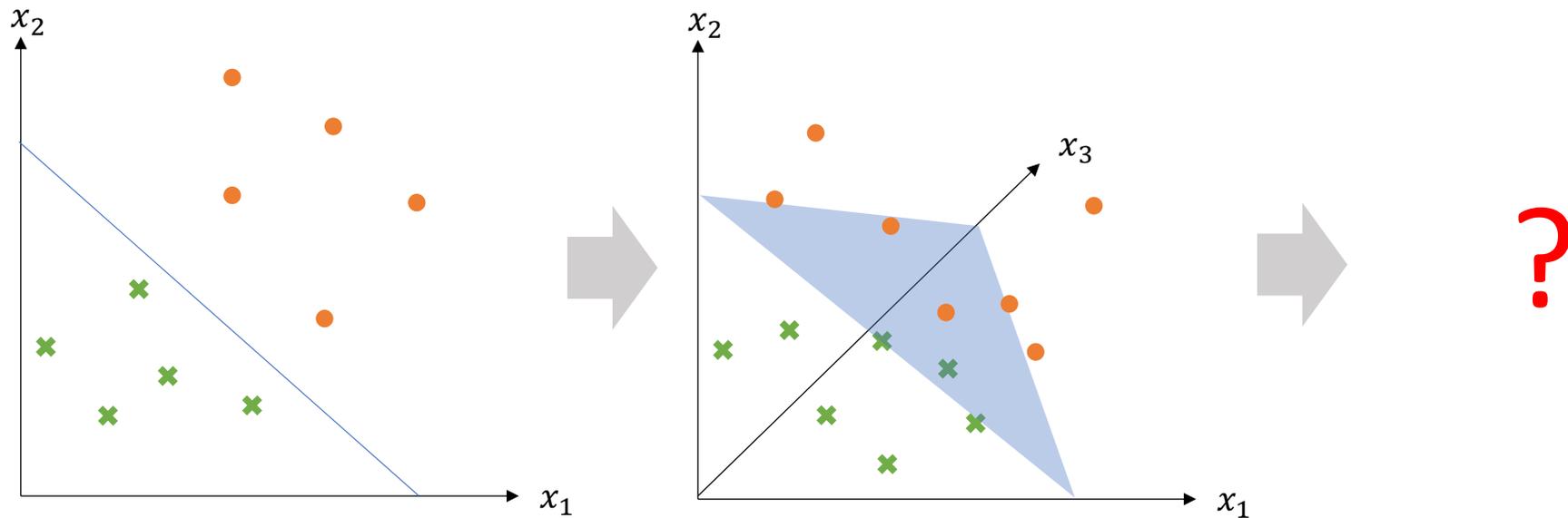
The decision boundary can be shifted left or right via the bias

The decision boundary must cross the origin if no bias !

Single Neuron

- Classification

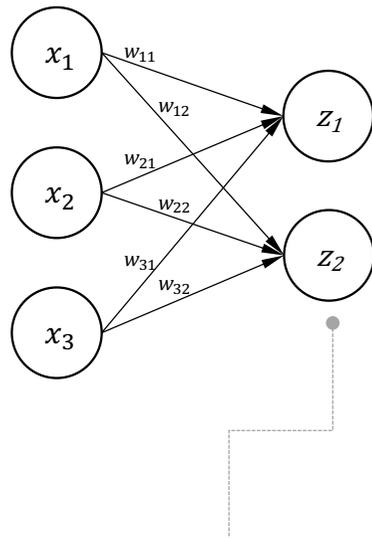
- Two input features: Decision boundary is a line
- Three input features: Decision boundary is a surface
- Many input features: Decision boundary is a **hyperplane** or hypersurface



Single Neuron

- 3 inputs and **2** outputs

input layer output layer



Multiple neurons

Expanding from a single neuron, a network can have multiple outputs

(The biases are not drawn to simplify the explanation)

$$z_1 = x_1w_{11} + x_2w_{21} + x_3w_{31} + b_1$$

$$z_2 = x_1w_{12} + x_2w_{22} + x_3w_{32} + b_2$$

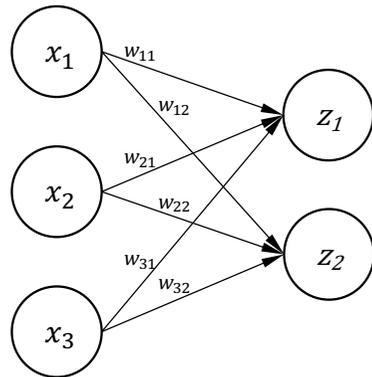
By using multiple neurons, we can obtain multiple outputs. For example, the outputs can represent the scores if we should play football or basketball.

They are all linear!

Single Neuron

- 3 inputs and 2 outputs

input layer output layer



$$z_1 = x_1 w_{11} + x_2 w_{21} + x_3 w_{31}$$

$$z_2 = x_1 w_{12} + x_2 w_{22} + x_3 w_{32}$$

column format

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

$$\mathbf{z} = \mathbf{W}^T \mathbf{x}$$

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

row format

$$\mathbf{z} = [z_1 \quad z_2] \quad \mathbf{x} = [x_1 \quad x_2 \quad x_3]$$

$$\mathbf{W} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

$$\mathbf{z} = \mathbf{x} \mathbf{W}$$

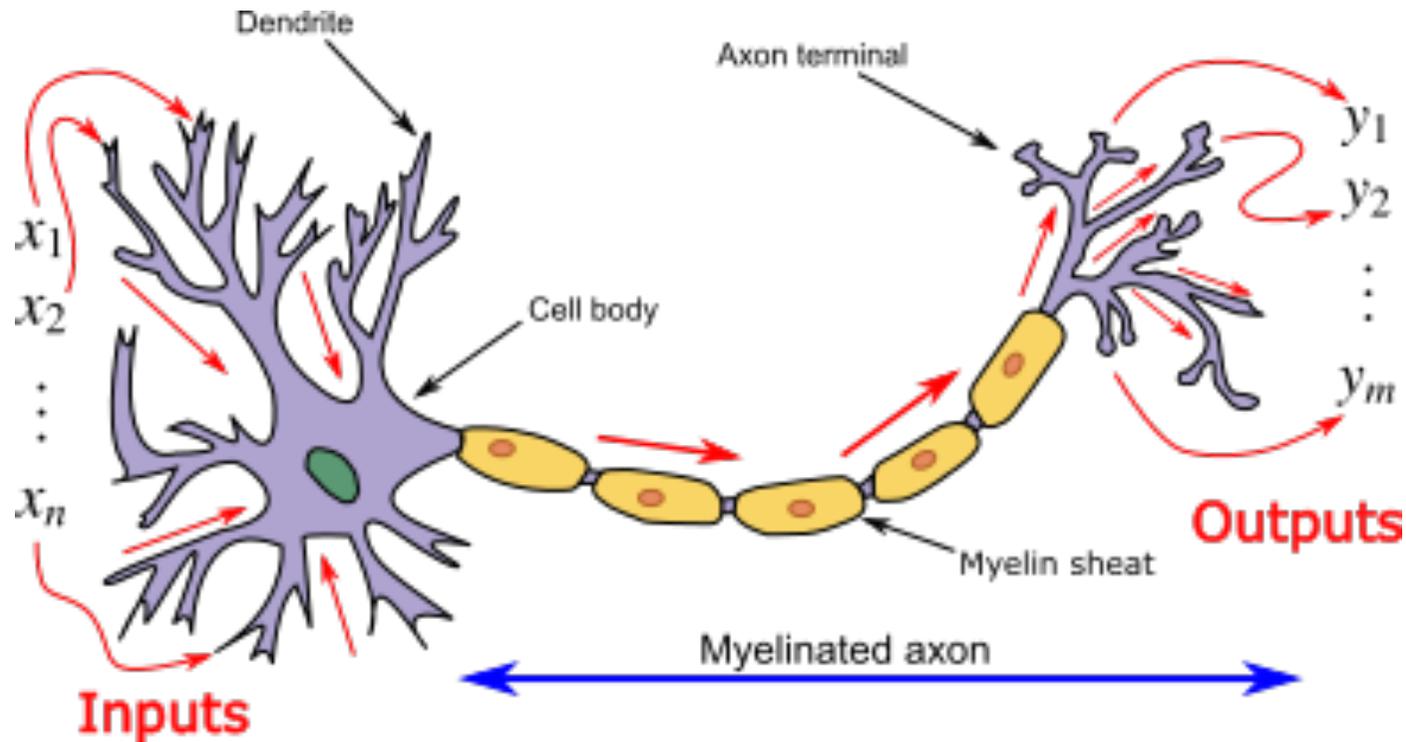
$$[z_1 \quad z_2] = [x_1 \quad x_2 \quad x_3] \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix}$$

Activation

Activation

- Motivation

Activation functions provide the **non-linearity** on the layers outputs, and their design remains an active researched area.



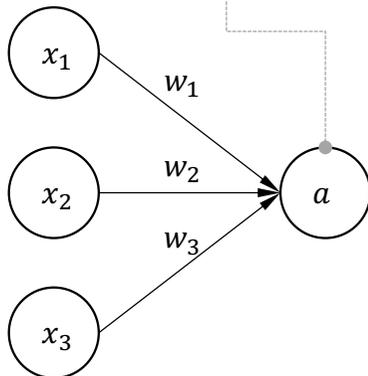
Activation

- Sigmoid or logistic function

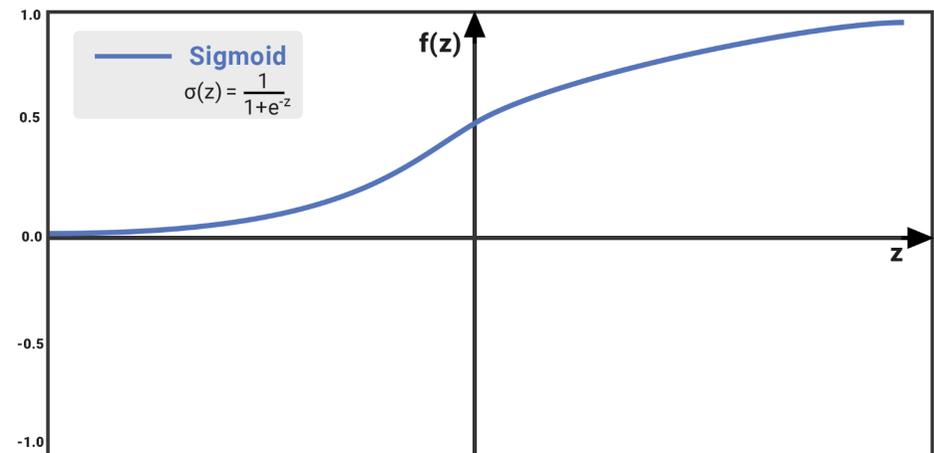
Continuing with our example, for a given neural network, the output can represent specific scores, such as the probability of playing football. To represent the probability from 0% to 100%, it is of common practice to apply a function to scale the output to a value between 0 to 1.

$$a = f(z) = f(x_1w_1 + x_2w_2 + b)$$

activation value



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

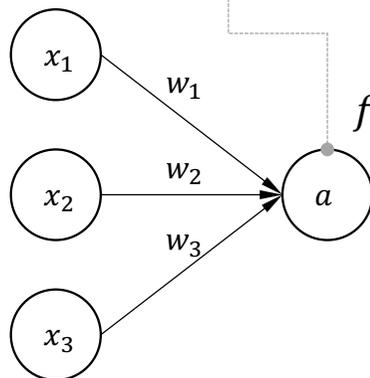


Activation

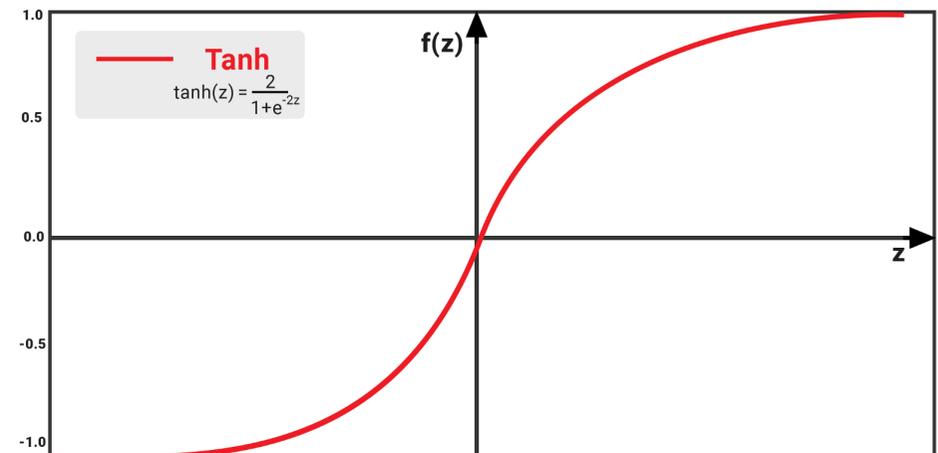
- Tanh

Similar to the sigmoid, the hyperbolic tangent (tanh) also scales the output layer to a limited range of values. With an output range between -1 to 1, this function is often used for regression, such as for an output image with pixel values between -1 to 1.

$$a = f(z) = f(x_1w_1 + x_2w_2 + b)$$



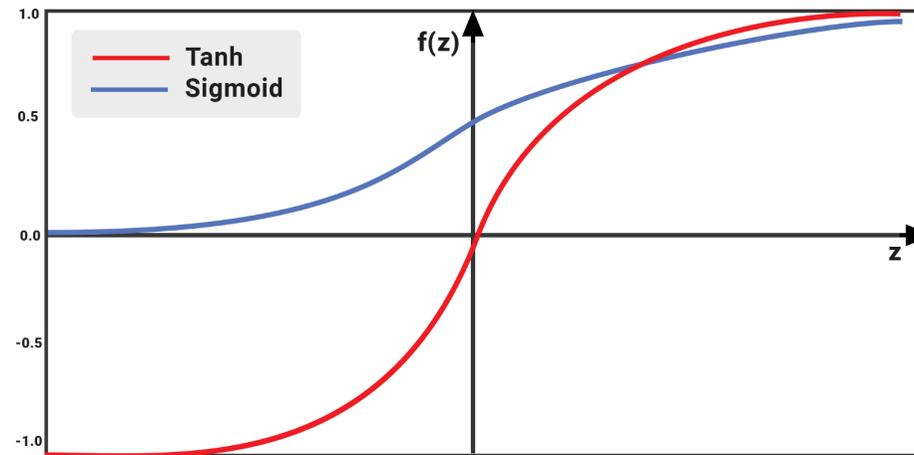
$$f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{2}{1 + e^{-2z}}$$



Activation

- Sigmoid vs Tanh

The sigmoid and hyperbolic tangent functions are used to provide non-linearity to the network.

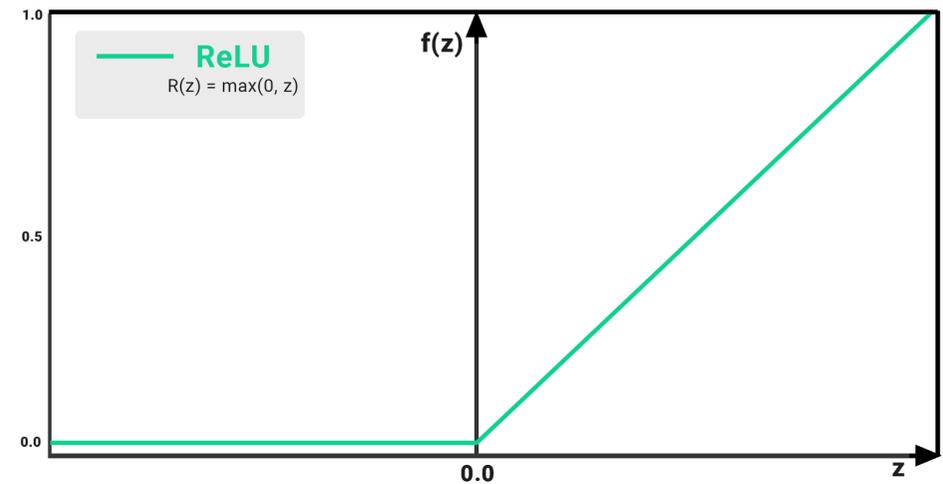
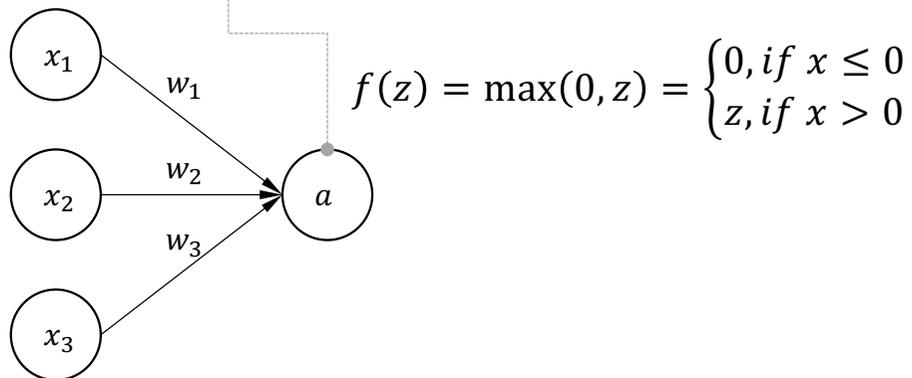


Activation

- Rectifier, Rectified linear unit (ReLU)

A function that sets the negative values to zero for the purpose of feature selection and provide a simple way to compute the derivative which will be used in the optimization.

$$a = f(z) = f(x_1w_1 + x_2w_2 + b)$$

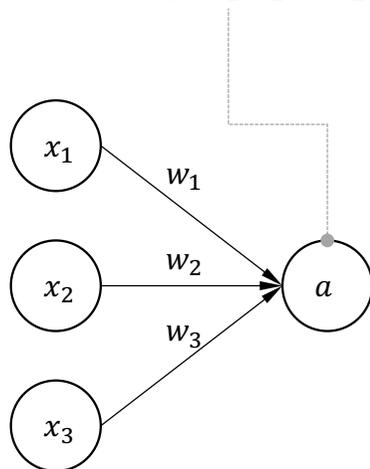


Activation

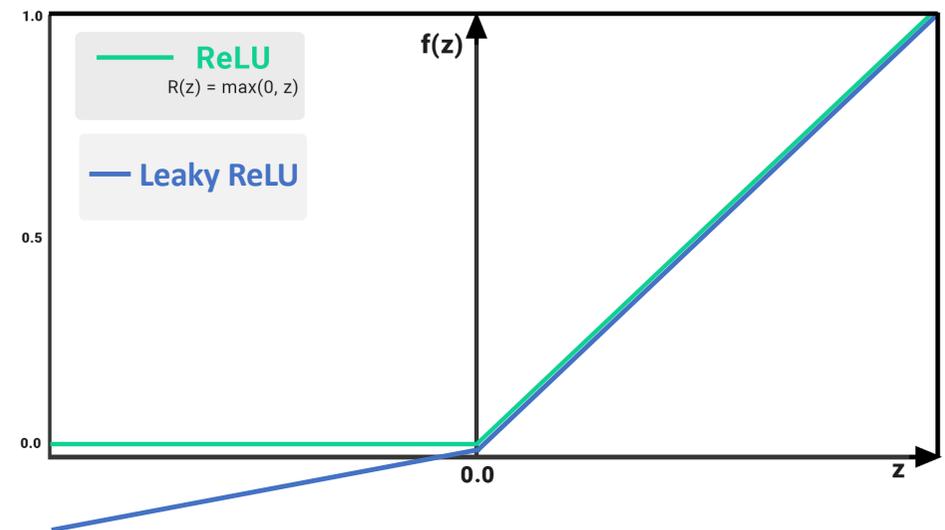
- Leaky ReLU and Parametric Leaky ReLU

However, merely setting negative values to zero will lead to information loss. A solution was proposed with the leaky ReLU where α is a small positive value to control the slope (e.g., 0.1 and 0.2) so that the information from the negative values can pass through to the output. In addition, the parametric ReLU (PReLU) was also proposed to consider α as a network parameter.

$$a = f(z) = f(x_1w_1 + x_2w_2 + b)$$



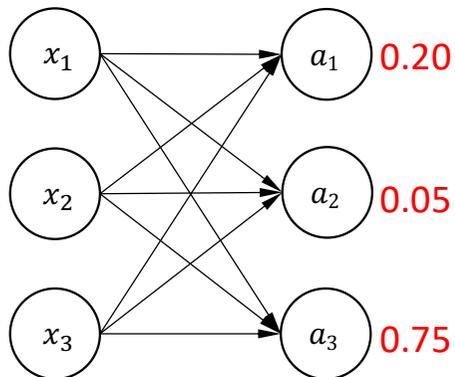
$$f(z) = \begin{cases} \alpha z, & \text{if } x \leq 0 \\ z, & \text{if } x > 0 \end{cases}$$



Activation

- Softmax

An network with three outputs and three inputs, with these multiple outputs, multi-class classification can be performed, i.e., classify the input into one of three or more classes. The Softmax function is designed for the output layer of a multi-class classifier to not only limit all outputs to 0 to 1, as with the sigmoid, but also ensure the sum of each output equals 1, i.e., the sum of all probabilities must be 100%.

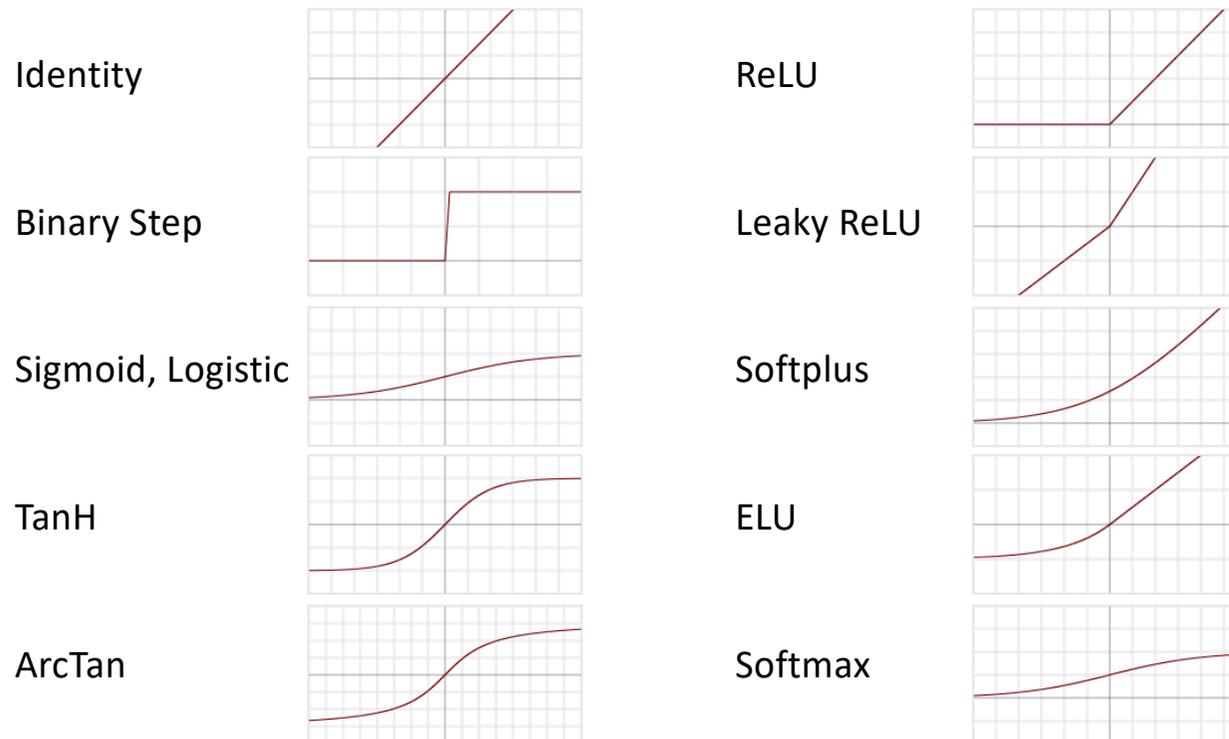


Give a output vector $\mathbf{z} = [z_1, z_2, z_3]$ then $K = 3$ is the vector length, we obtain an activation vector $\mathbf{a} = [a_1, a_2, a_3]$ as follow

$$a_i = f(z)_i = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$$

Softmax first applies an exponential function to each output and then normalizes each by dividing it by the sum of all outputs.

Activation

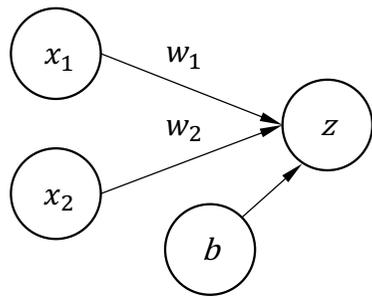


Reference: https://en.wikipedia.org/wiki/Activation_function

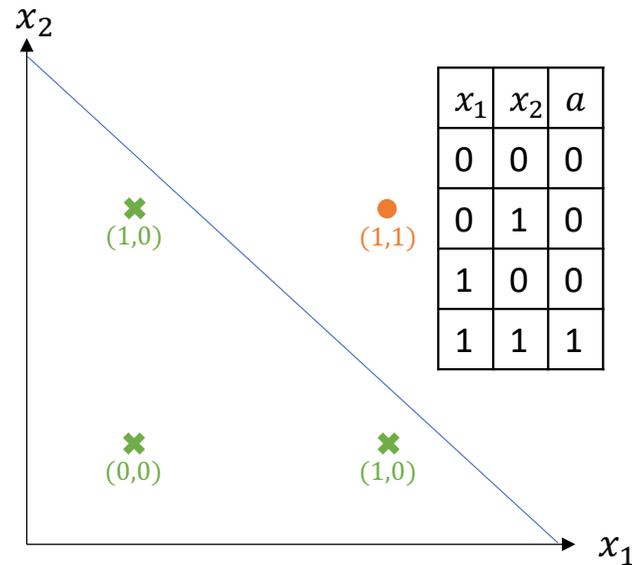
Multi-layer Perceptron

Multi-layer Perceptron

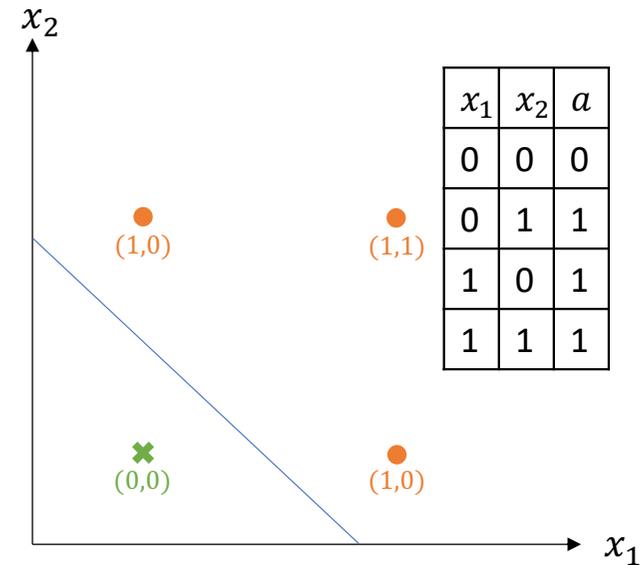
- Motivation : XOR Classification Problem



$$z = x_1w_1 + x_2w_2 + b$$



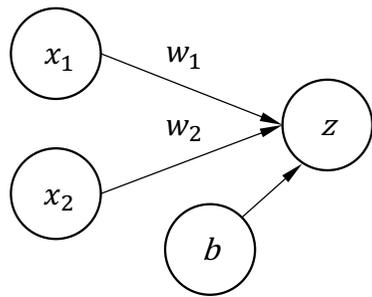
AND



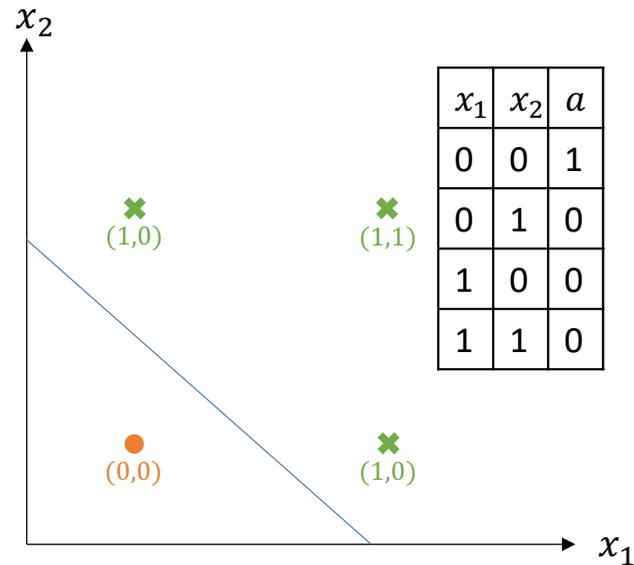
OR

Multi-layer Perceptron

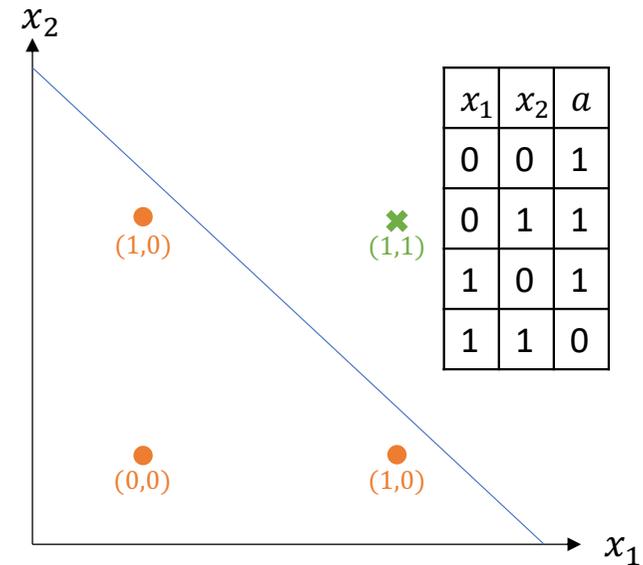
- Motivation : XOR Classification Problem



$$z = x_1 w_1 + x_2 w_2 + b$$



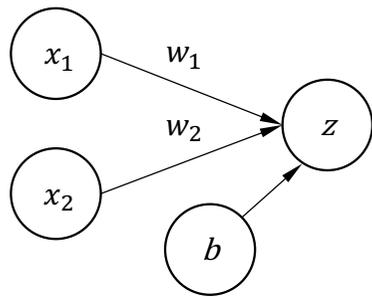
NOR



NAND

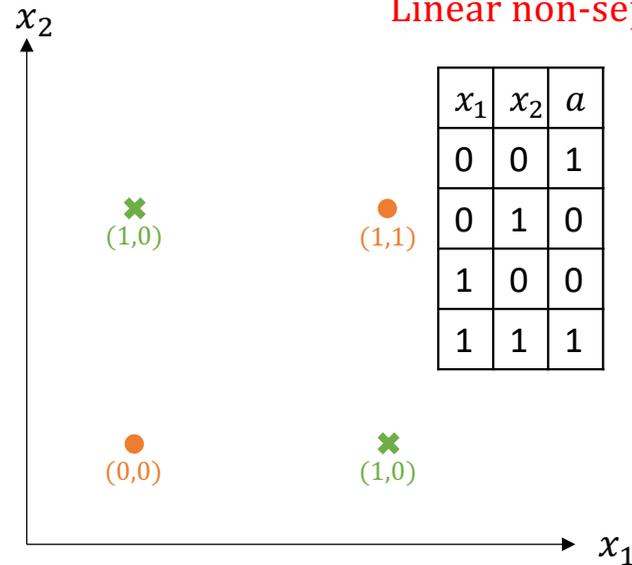
Multi-layer Perceptron

- Motivation : XOR Classification Problem

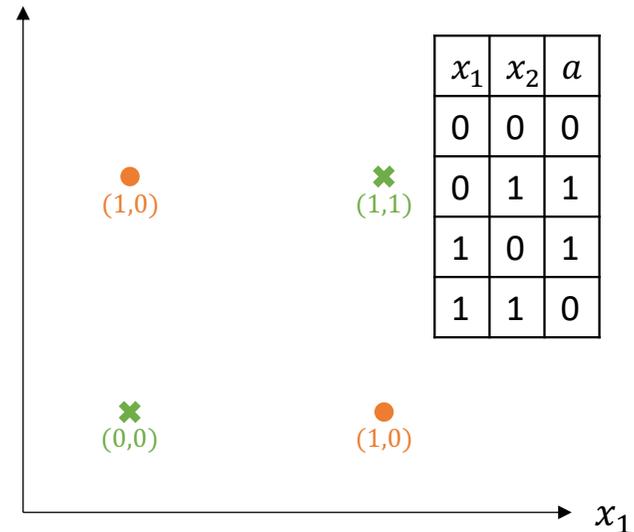


$$z = x_1w_1 + x_2w_2 + b$$

Cannot find a line to segment the data points
Linear non-separable problems



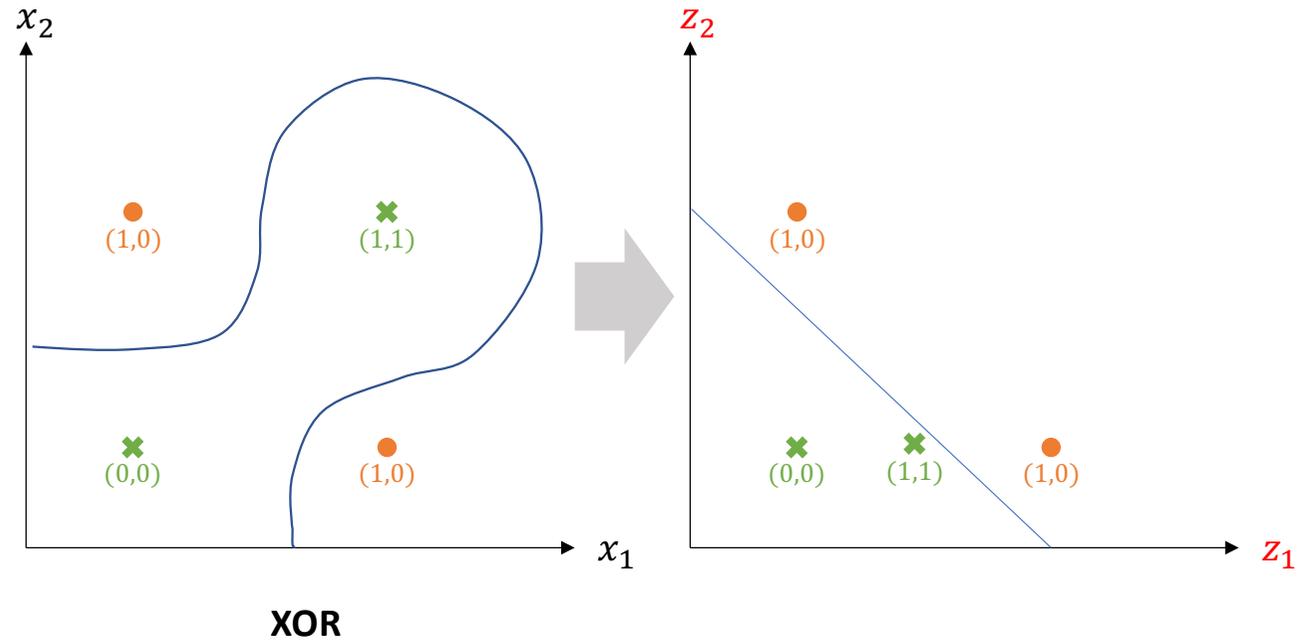
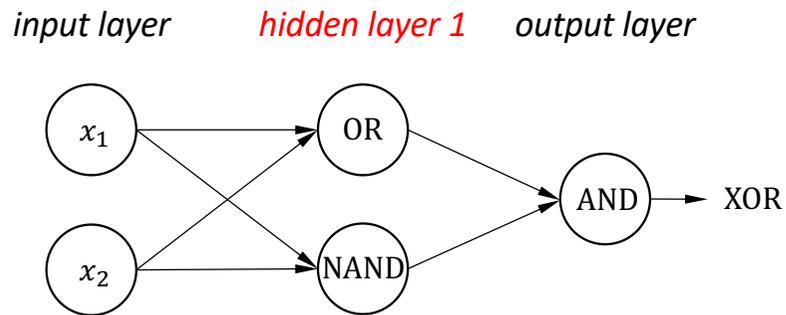
XNOR



XOR

Multi-layer Perceptron

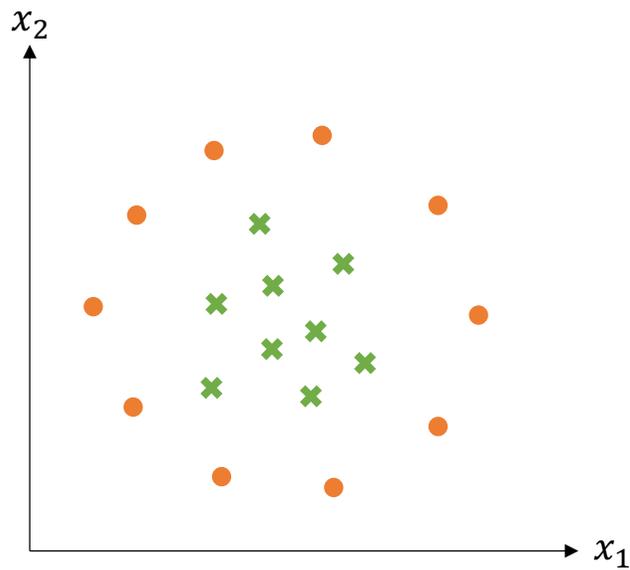
- XOR Classification Problem



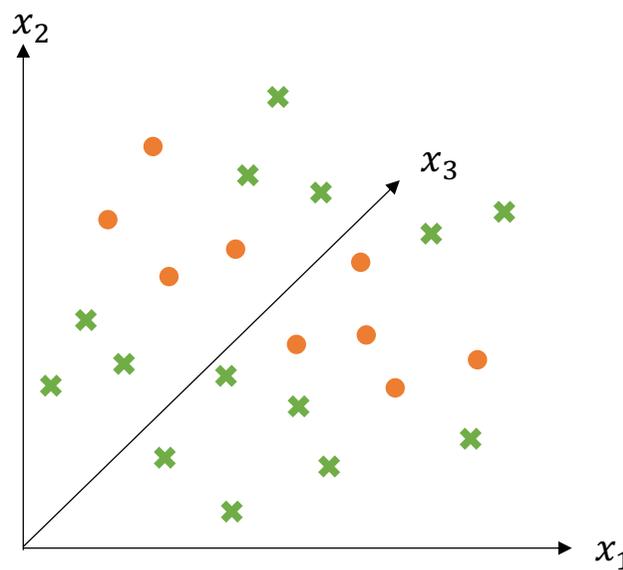
Multi-layer Perceptron

- More Complex Problems

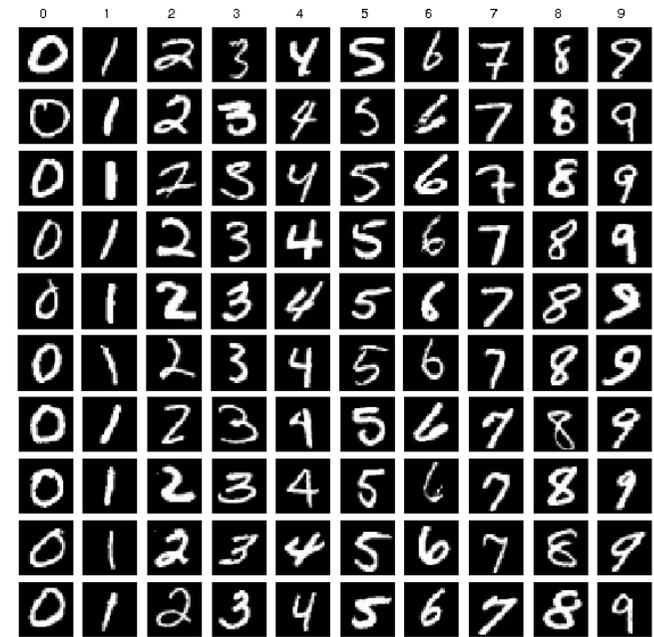
We need a network model with higher capacity



$$\mathbf{x} = [x_1, x_2]$$



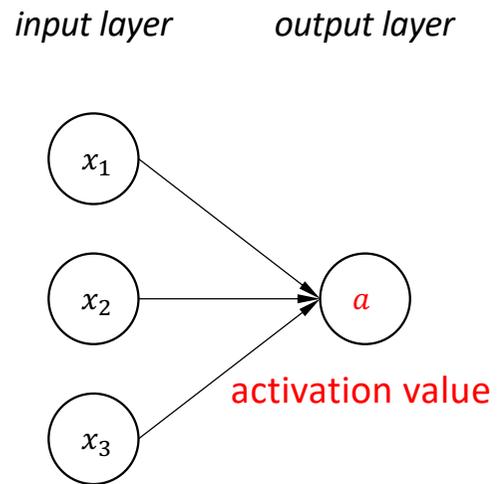
$$\mathbf{x} = [x_1, x_2, x_3]$$



$$\mathbf{x} = [x_1, x_2, x_3, \dots, x_{784}]$$

Multi-layer Perceptron

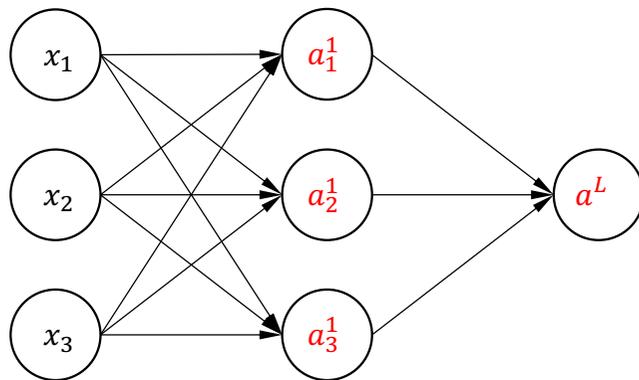
- Representation Capacity



Multi-layer Perceptron

- Representation Capacity

input layer hidden layer 1 output layer



A multi-layer perceptron (MLP) extends from a single, fully connected layer, and consists of at least two fully connected layers.

The layers between the inputs and outputs are called “hidden” because they cannot be directly accessed from outside the network.

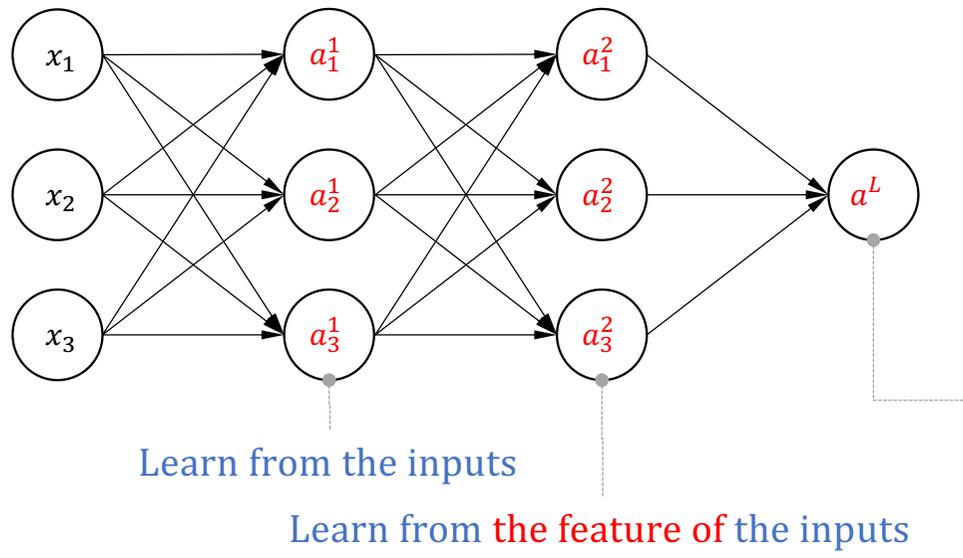
By stacking a new layer on top of an existing layer, the new layer is considered to use the output of the previously existing layer as its input features. Therefore, compared with a single fully connected layer, MLP can fit more complex input data. In other words, MLP can have more representational capability than a single layer.

The biases are not drawn to simplify the figure.

Multi-layer Perceptron

- Representation Capacity

input layer hidden layer 1 hidden layer 2 output layer

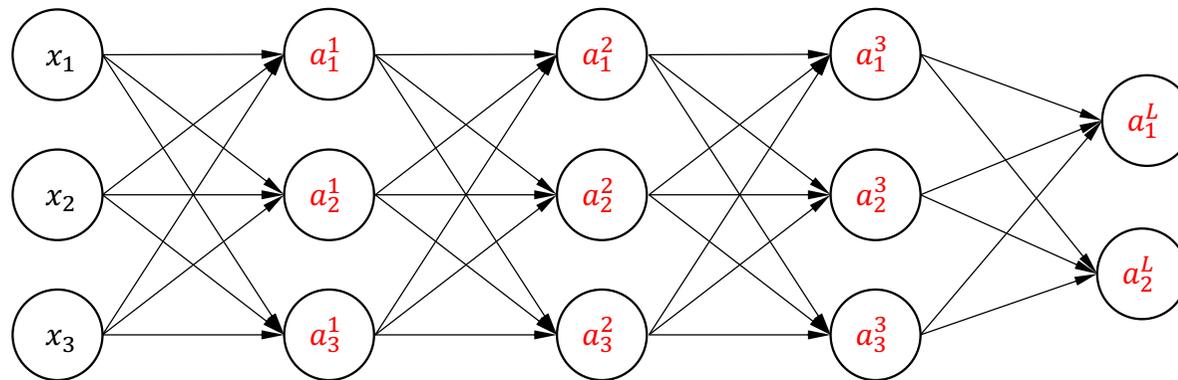


Learn from the feature of the feature of the inputs

Multi-layer Perceptron

- Representation Capacity: Hierarchical Representation

input layer hidden layer 1 hidden layer 2 hidden layer 3 output layer



a_k^l

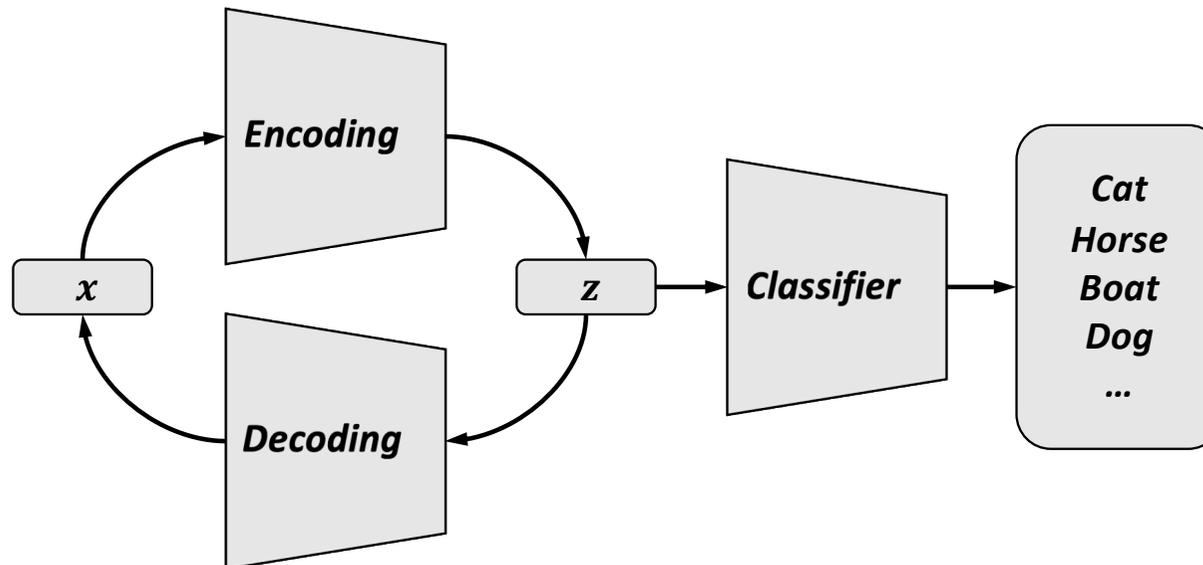
Layer index $l = 1 \dots L$ from the first hidden layer to the output layer.
The input layer can be written as $\mathbf{x} = \mathbf{a}^0$.

Output index $k = 1 \dots K$ where K is the number of units of this layer.

Activation output $\mathbf{a}^l = f(\mathbf{z}^l)$

Multi-layer Perceptron

- Representation Capacity: Encoding and Decoding



The successful of deep learning is the approximate capacity of the neural network that finds a good latent representation z for the visible input data x . In deep learning, the process of transforming visible input data, such as an image, text or video, into a latent representation, alternatively known as embedding or hidden representation, is referred to as “encoding”. The inverted process is referred to as “decoding”.

Loss Functions

Loss Functions

- Motivation

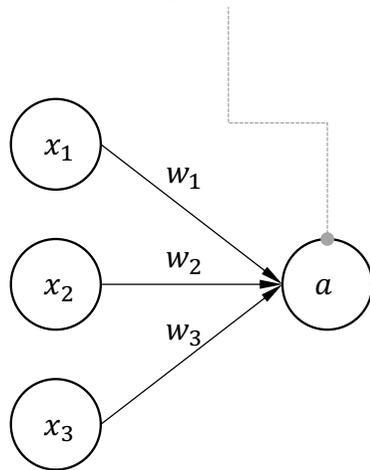
Loss functions are defined to quantify an error, known as the loss value, between the predicted and targeted (i.e., ground truth) outputs. The loss value is used as the goal for optimizing the neural network parameters, such as the weights and biases. Specifically, optimizing a given neural network minimizes the defined loss value by updating the network parameters. Gradient descent is commonly used to update the parameter by computing the partial derivatives of the loss w.r.t the network parameters. Details about gradient descent and neural network training are included in the next Section..

We don't set the values of parameters manually, we define a loss function and use data to optimize the parameters

Loss Functions

- Logistic Regression Loss

$$a = f(z) = f(x_1w_1 + x_2w_2 + b)$$



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

The model only has one output

$$\mathcal{L} = y \log(a) + (1 - y)\log(1 - a)$$

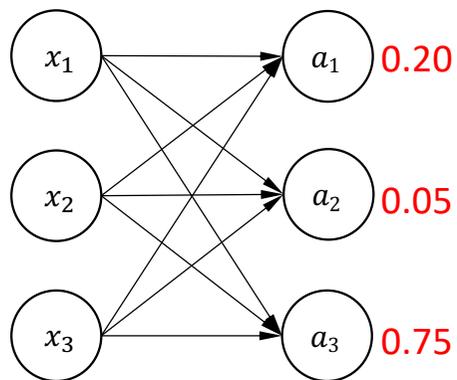
Given M data samples

$$\mathcal{L} = \sum_{m=1}^M (y^m \log(a^m) + (1 - y^m)\log(1 - a^m))$$

Loss Functions

- Cross-Entropy Loss

Output a vector
 $\mathbf{a} = \text{softmax}(\mathbf{z})$



The model has multiple outputs

$$\mathcal{L} = \sum_{k=1}^K \mathbf{y}_k \log(\mathbf{a}_k)$$

Given M data samples

$$\mathcal{L} = \sum_{m=1}^M \mathbf{y}_k^m \log(\mathbf{a}_k^m)$$

Loss Functions

- \mathcal{L}_p norm

Measure the scale of a vector

$$\|\mathbf{x}\|_p = \left(\sum_{k=1}^K |\mathbf{x}_k|^p\right)^{\frac{1}{p}}$$
$$\|\mathbf{x}\|_p^p = \sum_{k=1}^K |\mathbf{x}_k|^p$$

Measure the difference between two vectors

$$\mathcal{L}_p = \|\mathbf{y} - \mathbf{a}\|_p^p = \sum_{k=1}^K |\mathbf{y}_k - \mathbf{a}_k|^p$$

Loss Functions

- Mean Squared Error (MSE)

MSE is the \mathcal{L}_2 norm over samples and is used for regression problems in which the output of the neural networks contains continuous values, such as the pixels of an image or a scalar value

$$\mathcal{L}_{MSE} = \|\mathbf{y} - \mathbf{a}\|_2^2$$

Given M data samples

$$\mathcal{L}_{MSE} = \frac{1}{M} \sum_{m=1}^M \|\mathbf{y}^m - \mathbf{a}^m\|_2^2$$

Loss Functions

- Mean Absolute Error (MAE)

MAE is the \mathcal{L}_1 norm over samples, known as the least square error, MAE is also used for regression problems and is expressed as follows.

$$\mathcal{L}_{MAE} = \|\mathbf{y} - \mathbf{a}\|$$

Given M data samples

$$\mathcal{L}_{MAE} = \frac{1}{M} \sum_{m=1}^M |\mathbf{y}^m - \mathbf{a}^m|$$

Loss Functions

- Many many more



Optimization

Optimization

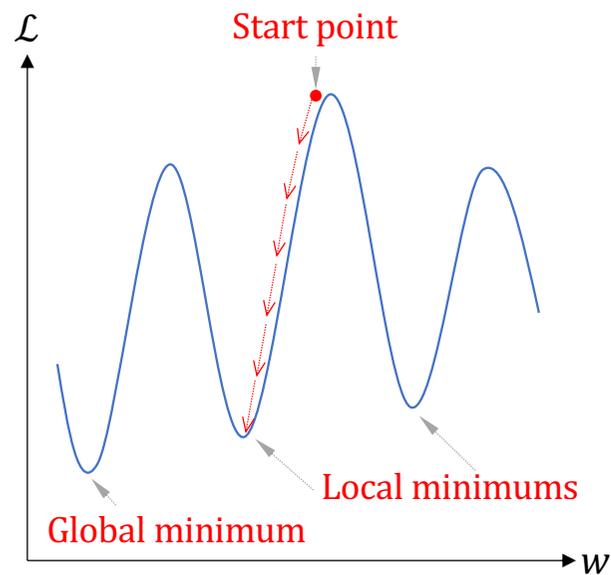
- Motivation

Given a network $f(x; \theta)$ and a loss function \mathcal{L} , training the network is the process to minimise the loss value \mathcal{L} by updating the network parameters θ , such as the weights and biases. Gradient descent is one method to update the network parameters. Even though there are some other optimisation methods exist, such as limited memory BFGS (L-BFGS) and conjugate gradient (CG), due to the drawbacks related to larger computation requirements, they are not often applied.

Given $f(x; \theta)$ and \mathcal{L} find a good θ

Optimization

- Gradient Descent

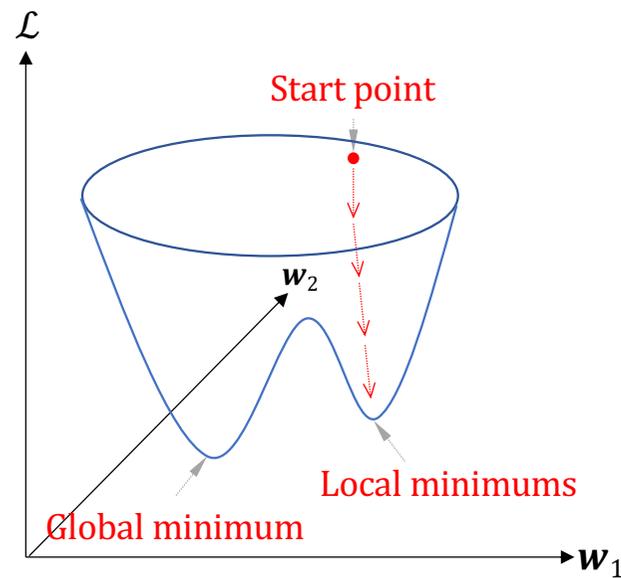


Assuming one parameter only

$$w := w - \alpha \frac{d\mathcal{L}}{dw}$$

Optimization

- Gradient Descent



Assuming two parameters

$$w_j := w_j - \alpha \frac{\partial \mathcal{L}}{\partial w_j} \quad \mathbf{w} = [w_1, w_2]$$

Optimization

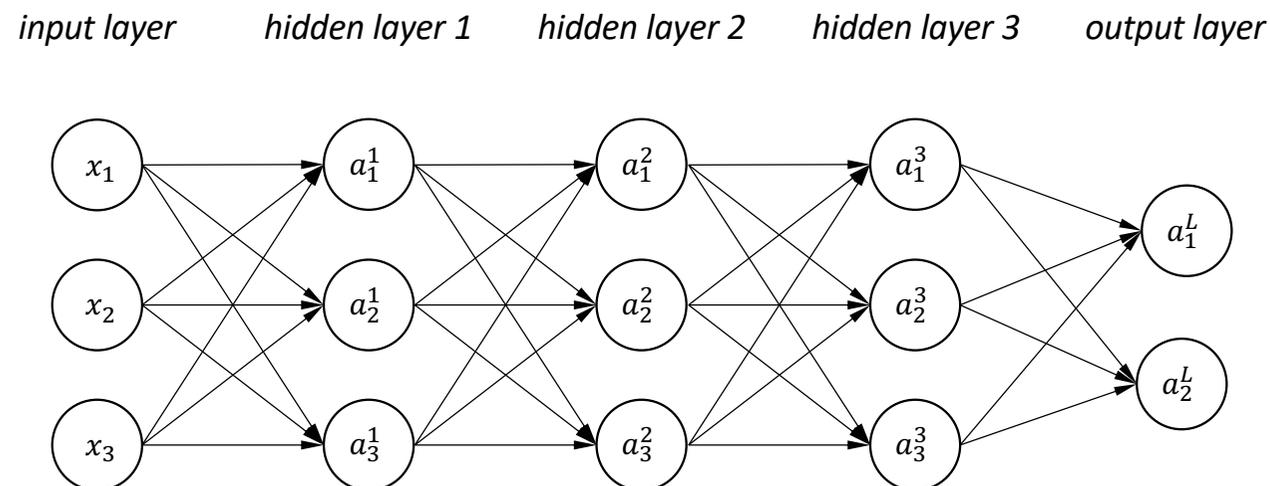
- Gradient Descent

All we need: $\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}$

Optimization

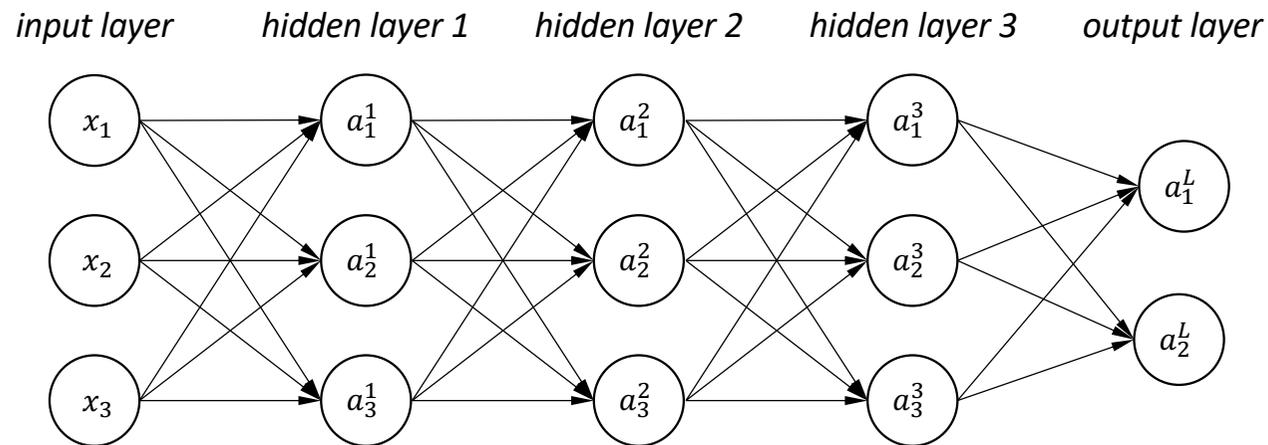
- Error Back-Propagation

The process of error back-propagation computes the **gradients** $\frac{\partial \mathcal{L}}{\partial \theta}$ for every parameters in the network. When computing the gradient, an **intermediate result** $\delta = \frac{\partial \mathcal{L}}{\partial z}$ is introduced, which is the ∂z partial derivative of the loss \mathcal{L} w.r.t the layer's output z . Based on this intermediate result, the process next computes the partial derivative of the loss \mathcal{L} w.r.t every parameters $\frac{\partial \mathcal{L}}{\partial \theta}$ which is then used to update the parameter values.



Optimization

- Error Back-Propagation



Layer index $l = 1 \dots L$ from the first hidden layer to the output layer.
The input layer can be written as $\mathbf{x} = \mathbf{a}^0$.

$$\mathbf{a}^l = f(\mathbf{z}^l) = \frac{1}{1 + e^{-\mathbf{z}^l}}$$

$$\mathbf{z}^l = \mathbf{W}^{lT} \mathbf{a}^{l-1} + \mathbf{b}^l$$

$$\mathcal{L} = \frac{1}{2} (\mathbf{y} - \mathbf{a}^L)^2$$

We use this model and loss as an example:

Optimization

• Error Back-Propagation: **Vectors in Column Format**

1. Given:

- $\mathbf{a}^l = f(\mathbf{z}^l) = \frac{1}{1+e^{-z^l}}$
- $\mathbf{z}^l = \mathbf{W}^{lT} \mathbf{a}^{l-1} + \mathbf{b}^l$
- $\mathcal{L} = \frac{1}{2} (\mathbf{y} - \mathbf{a}^L)^2$

2. Then we can have the following derivatives:

- $\frac{\partial \mathbf{a}^l}{\partial \mathbf{z}^l} = f^{-1}(\mathbf{z}^l) = \mathbf{a}^l \circ (1 - \mathbf{a}^l)$
- $\frac{\partial \mathcal{L}}{\partial \mathbf{a}^l} = (\mathbf{a}^L - \mathbf{y})$ Hadamard (element-wise) product
- $\frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} = \mathbf{a}^{l-1}$ and $\frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} = \mathbf{1}$

3. Then the error of the output layer

$$\delta^L = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^L} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^L} \frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L} = (\mathbf{a}^L - \mathbf{y}) \circ (\mathbf{a}^L \circ (1 - \mathbf{a}^L))$$

Chain rule !

4. Then the error of the other layers $l = 1 \dots L - 1$:

- $\delta^l = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{l+1}} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \delta^{l+1} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l}$
- $\mathbf{z}^{l+1} = \mathbf{W}^{l+1T} \mathbf{a}^l + \mathbf{b}^{l+1}$
- $\frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \mathbf{W}^{l+1T} f^{-1}(\mathbf{z}^l) = \mathbf{W}^{l+1T} \circ (\mathbf{a}^l \circ (1 - \mathbf{a}^l))$
- $\delta^l = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{l+1}} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \mathbf{W}^{l+1T} \delta^{l+1} \circ (\mathbf{a}^l \circ (1 - \mathbf{a}^l))$

5. The the gradients are:

- $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} = \delta^l \frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} = \delta^l \mathbf{a}^{l-1T}$
- $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} = \delta^l \frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} = \delta^l$

6. Update parameters iteratively

$$\mathbf{W}^l := \mathbf{W}^l - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} \quad \mathbf{b}^l := \mathbf{b}^l - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{b}^l}$$

Optimization

• Error Back-Propagation: **Vectors in Raw Format for Python**

1. Given:

- $\mathbf{a}^l = f(\mathbf{z}^l) = \frac{1}{1+e^{-z}}$
- $\mathbf{z}^l = \mathbf{a}^{l-1}\mathbf{W}^l + \mathbf{b}^l$
- $\mathcal{L} = \frac{1}{2}(\mathbf{y} - \mathbf{a}^L)^2$

2. Then we can have the following derivatives:

- $\frac{\partial \mathbf{a}^l}{\partial \mathbf{z}^l} = f^{-1}(\mathbf{z}^l) = \mathbf{a}^l \circ (1 - \mathbf{a}^l)$
- $\frac{\partial \mathcal{L}}{\partial \mathbf{a}^L} = (\mathbf{a}^L - \mathbf{y})$
- $\frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} = \mathbf{a}^{l-1}$ and $\frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} = \mathbf{1}$

3. Then the error of the output layer

- $\delta^L = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^L} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^L} \frac{\partial \mathbf{a}^L}{\partial \mathbf{z}^L} = (\mathbf{a}^L - \mathbf{y}) \circ (\mathbf{a}^L \circ (1 - \mathbf{a}^L))$

4. Then the error of the other layers $l = 1 \dots L - 1$:

- $\delta^l = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{l+1}} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \delta^{l+1} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l}$
- $\mathbf{z}^{l+1} = \mathbf{a}^l \mathbf{W}^{l+1} + \mathbf{b}^{l+1}$
- $\frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \mathbf{W}^{l+1} \circ f^{-1}(\mathbf{z}^l) = \mathbf{W}^{l+1} \circ (\mathbf{a}^l \circ (1 - \mathbf{a}^l))$
- $\delta^l = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{l+1}} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \delta^{l+1} \mathbf{W}^{l+1 T} \circ (\mathbf{a}^l \circ (1 - \mathbf{a}^l))$

5. Then the gradients are:

- $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} = \delta^l \frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l} = \mathbf{a}^{l-1 T} \delta^l$
- $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} = \delta^l \frac{\partial \mathbf{z}^l}{\partial \mathbf{b}^l} = \delta^l$

6. Update parameters iteratively

$$\mathbf{W}^l := \mathbf{W}^l - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{W}^l} \quad \mathbf{b}^l := \mathbf{b}^l - \alpha \frac{\partial \mathcal{L}}{\partial \mathbf{b}^l}$$

Optimization

- Error Back-Propagation: Gradient Vanish

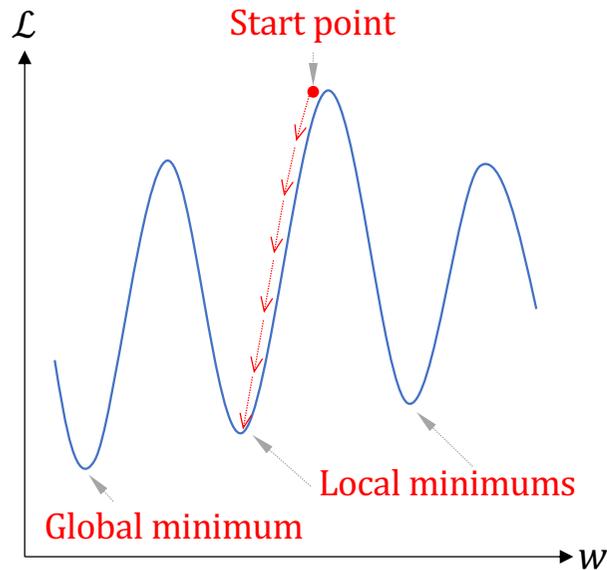
- $$\delta^l = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{l+1}} \frac{\partial \mathbf{z}^{l+1}}{\partial \mathbf{z}^l} = \delta^{l+1} \mathbf{W}^{l+1T} \circ (\mathbf{a}^l \circ (1 - \mathbf{a}^l))$$

Following the previous example, δ has a term $(\mathbf{a}^l \circ (1 - \mathbf{a}^l))$ in which if the activation output \mathbf{a} is close to 0 or 1, the δ^l is small. In chain rule, $\delta^l = \delta^{l+1} \dots$ so as the error propagate from the output layer to the input layer, the δ will become more and more small, which will leads to difficult for updating the parameters of the layers on the front (near the input).

- Solution 1: use ReLU to replace Sigmoid function. (commonly used in practice)
- Solution 2: pre-trained the parameters layer-by-layer ... Deep Belief Net
- ...

Optimization

- Stochastic Gradient Descent (SGD)



Motivation

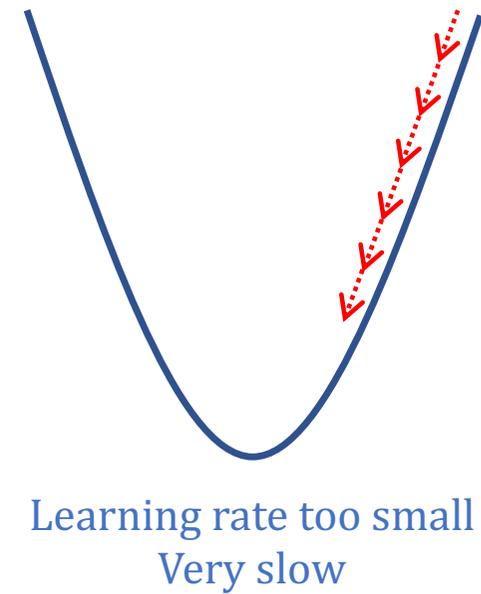
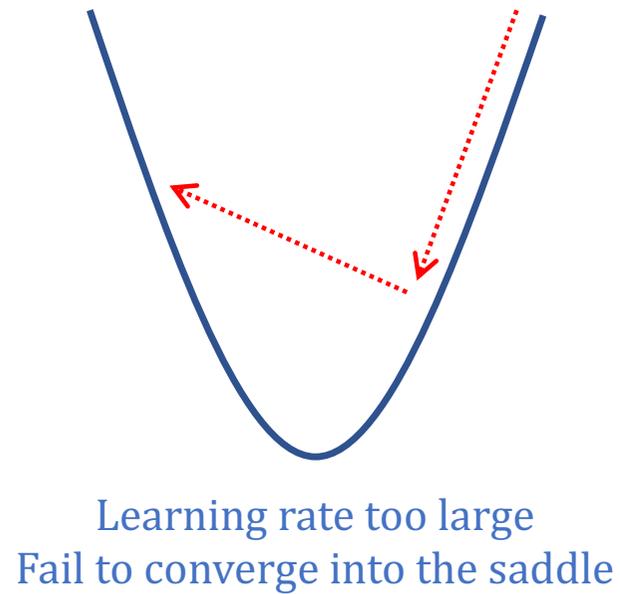
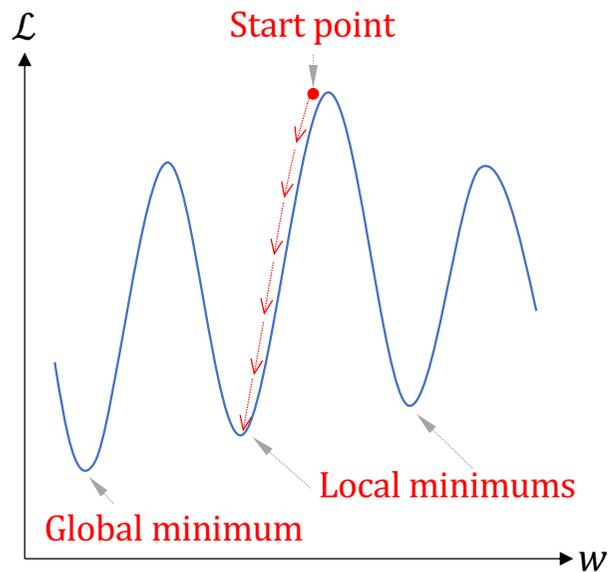
- Error back-propagation updates the parameters iteratively.
- Compute the loss \mathcal{L} using all data samples every time would be SLOW!

Method

- Randomly select a “batch” of samples from the training set to compute the \mathcal{L} .
- These data are called a “mini-batch”, and the quantity of data is called the “batch size”
- By updating the parameter multiple times, the mini-batches will cover the entire training set. Epoch represents the mini-batch has looped over the entire training dataset.

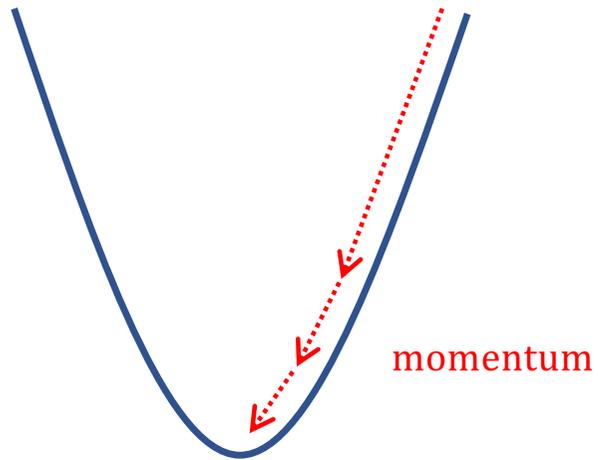
Optimization

- Adaptive Learning Rate



Optimization

- Adaptive Learning Rate

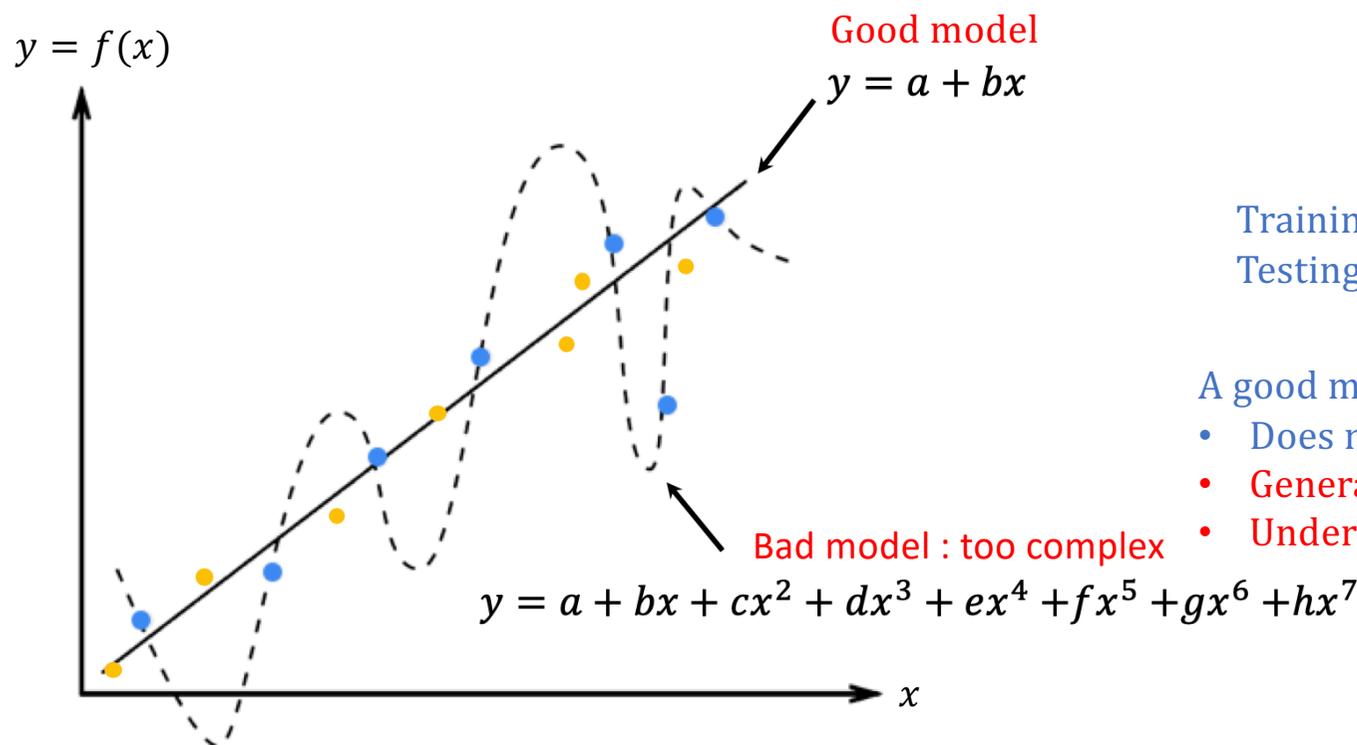


RMSProp, Adagrad, Adam, AMSGrad, AdaBound (PKU Undergraduate)

Optimization

- Hyper-Parameter Selection

Training Data vs Testing Data



- Training data
- Testing data

Training data is used to train the model
Testing data is used to test the model

A good model $f(x; \theta)$ (algorithm):

- Does not **overfit** to training data
- **Generalizes** well to testing data
- **Underfitting** can be solved by a deeper model

Optimization

- Hyper-Parameter Selection

Training Data vs Validation Data vs Testing Data

Hyper-parameters are the settings of the models, such as

- The number of layers
- The number of units
- The activation functions
- The loss function
- The batch size
- The number of epochs
- ...



**DO NOT select the hyper-parameters based on the performance on the testing set
IT IS CHEATING**

Training data is used to train the model

Validation data is used to select the hyper-parameters

Testing data is used to test the model

Optimization

- Hyper-Parameter Selection & Cross Validation

Training Data vs ~~Validation Data~~ vs Testing Data

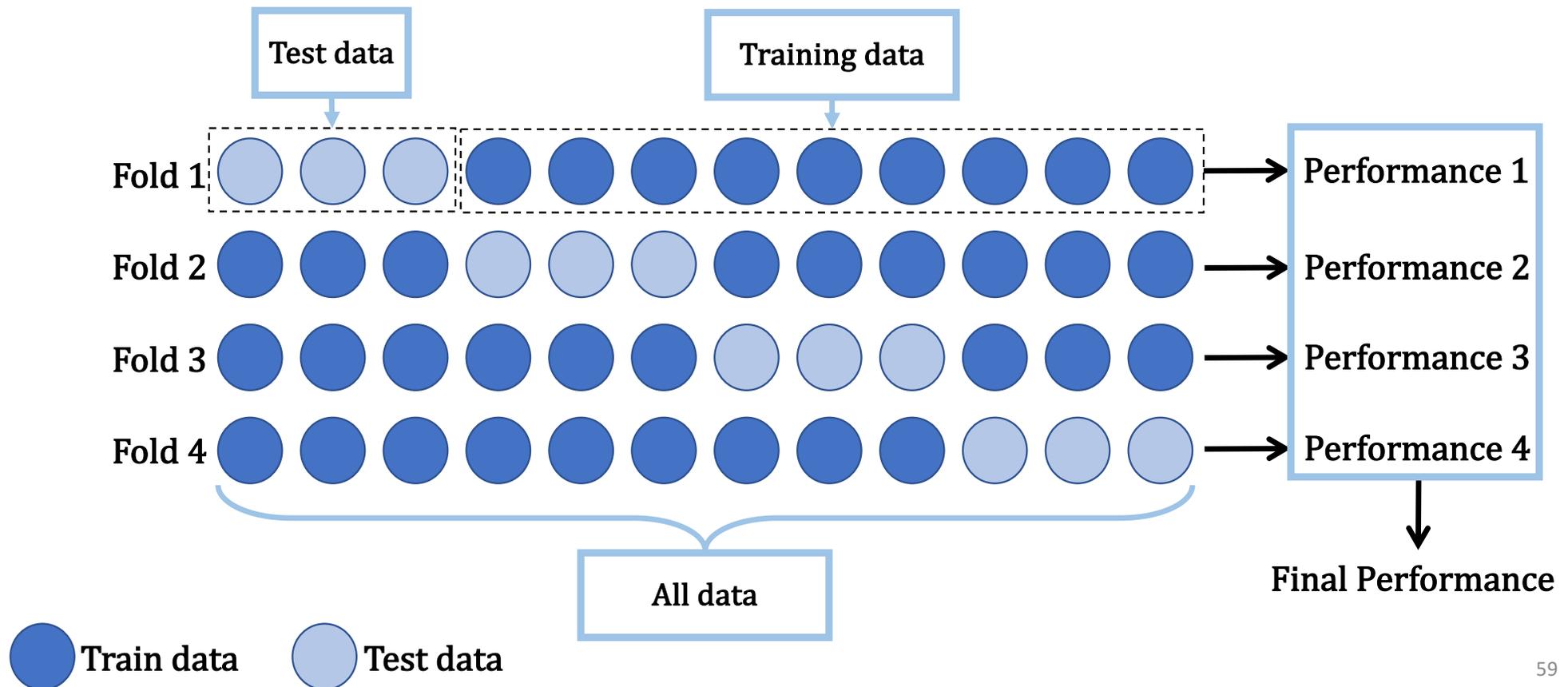
For small datasets, splitting the data into training, validating, and testing sets may be challenging. If the size of the training data is too small, then the performance of the model will be impacted. On the other hand, if the test data is too small, then the evaluation may not adequately reflect the performance.

How to select hyper-parameters without a validation set?

Optimization

- Hyper-Parameter Selection & **K-Fold Cross Validation**

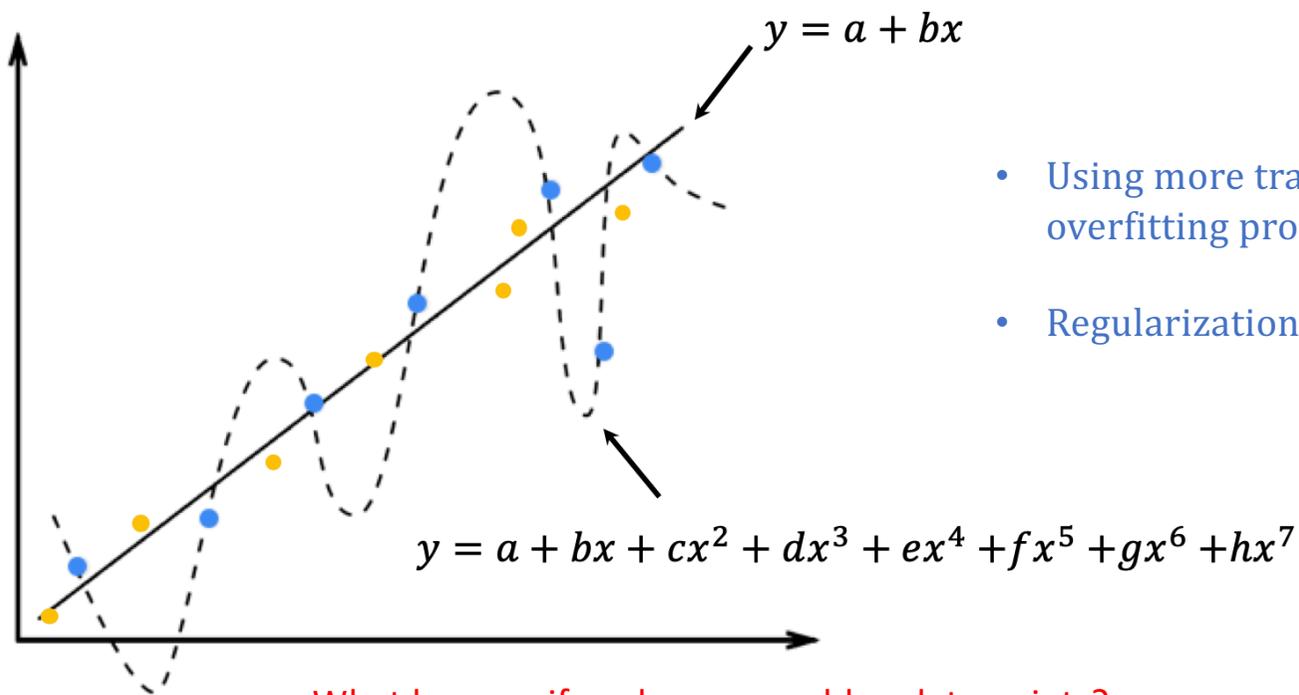
The dataset is separated into K folds of data evenly.



Regularization

Regularization

- Motivation



What happen if we have more blue data points?

- Using more training data can help to alleviate the overfitting problem, but it is expensive.
- Regularization is for alleviating the overfitting problem.

Regularization

- Data Augmentation

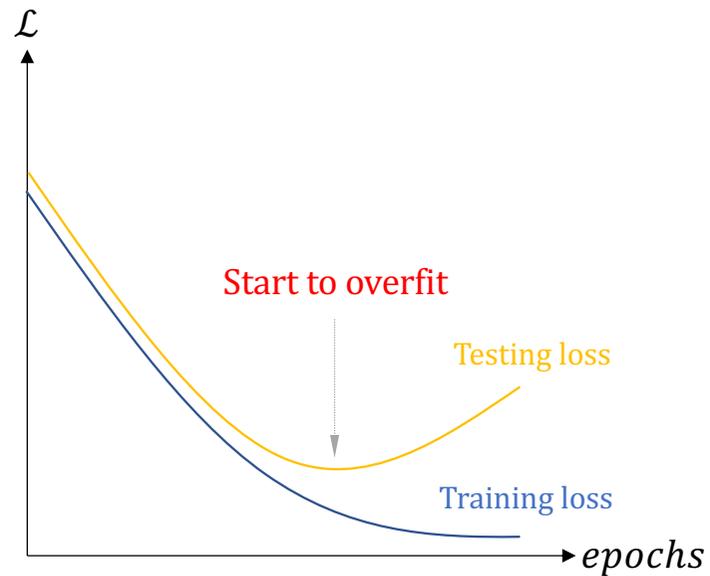


Common image data augmentation methods: horizontal flipping, rotating, shifting, and zooming

Regularization

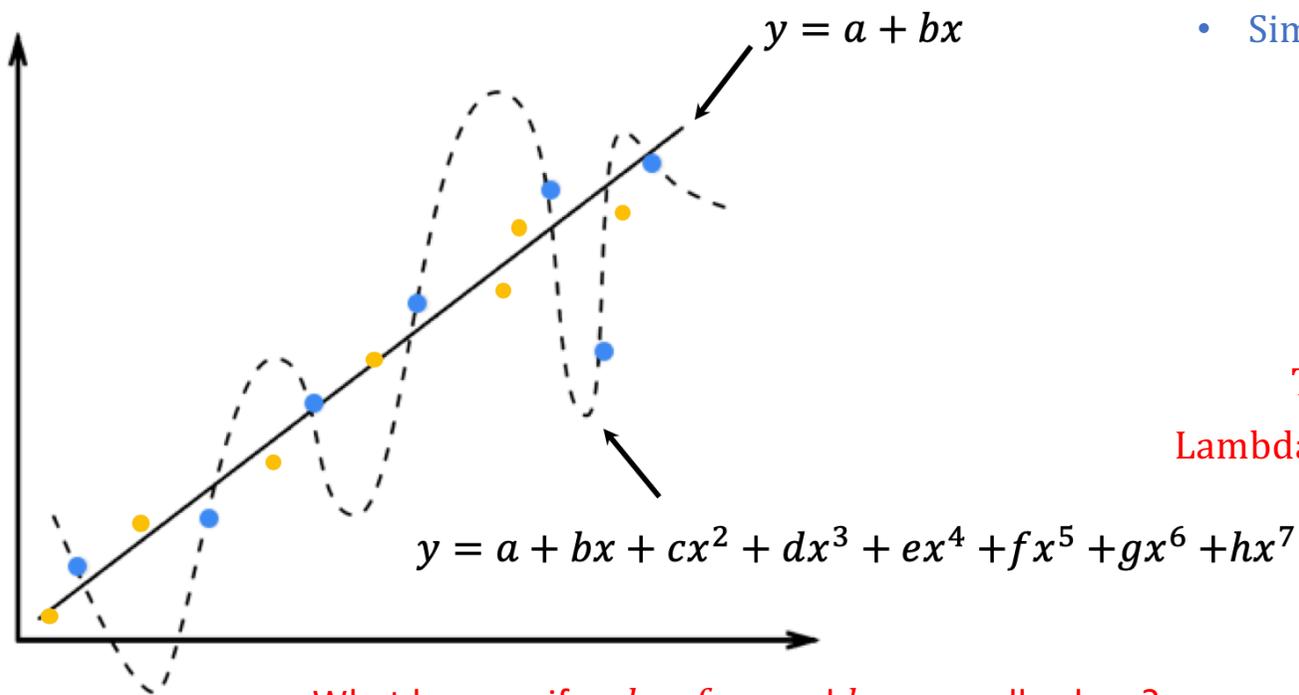
- Early Stopping

- Terminate the training before the model start to overfit the data



Regularization

- Weight Decay



What happen if $c, d, e, f, g,$ and h are small values?

- Simply make the parameters smaller !

$$\mathcal{L}_{total} = \mathcal{L} + \lambda ||\mathbf{W}||$$

The original loss

Regularization term

Lambda: controls the strength of the regularization

Regularization

- \mathcal{L}_1 norm

$$\mathcal{L}_{total} = \mathcal{L} + \lambda \mathcal{L}_1$$

$$\mathcal{L}_1 = \|\mathbf{W}\|$$

- \mathcal{L}_2 norm

$$\mathcal{L}_{total} = \mathcal{L} + \lambda \mathcal{L}_2$$

$$\mathcal{L}_2 = \|\mathbf{W}\|_2^2$$

Note: The weight decay only for the weights not the biases

Regularization

- \mathcal{L}_1 vs \mathcal{L}_2 norm

Given only two weights: w_1 and w_2

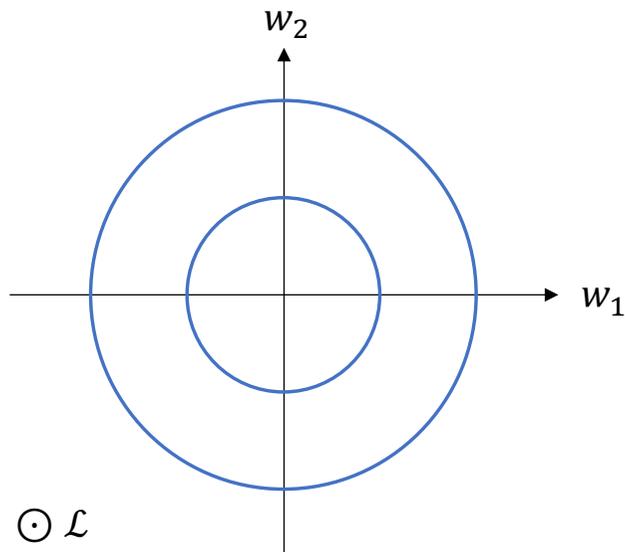
\mathcal{L}_2 penalty become ..

$$\mathcal{L}_2 = \|\mathbf{W}\|_2^2 = w_1^2 + w_2^2$$

Given a specific \mathcal{L}_2 value, it is a circle:

$1 = w_1^2 + w_2^2$ is a circle with a radius of 1

$4 = w_1^2 + w_2^2$ is a circle with a radius of 2



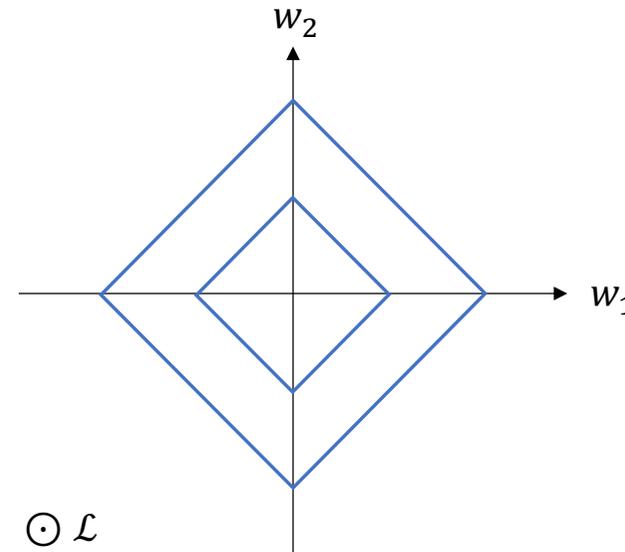
\mathcal{L}_1 penalty become ..

$$\mathcal{L}_1 = |w_1| + |w_2|$$

Given a specific \mathcal{L}_1 value, it is a square:

$1 = |w_1| + |w_2|$ is a square with a diagonal line that has a length of 1

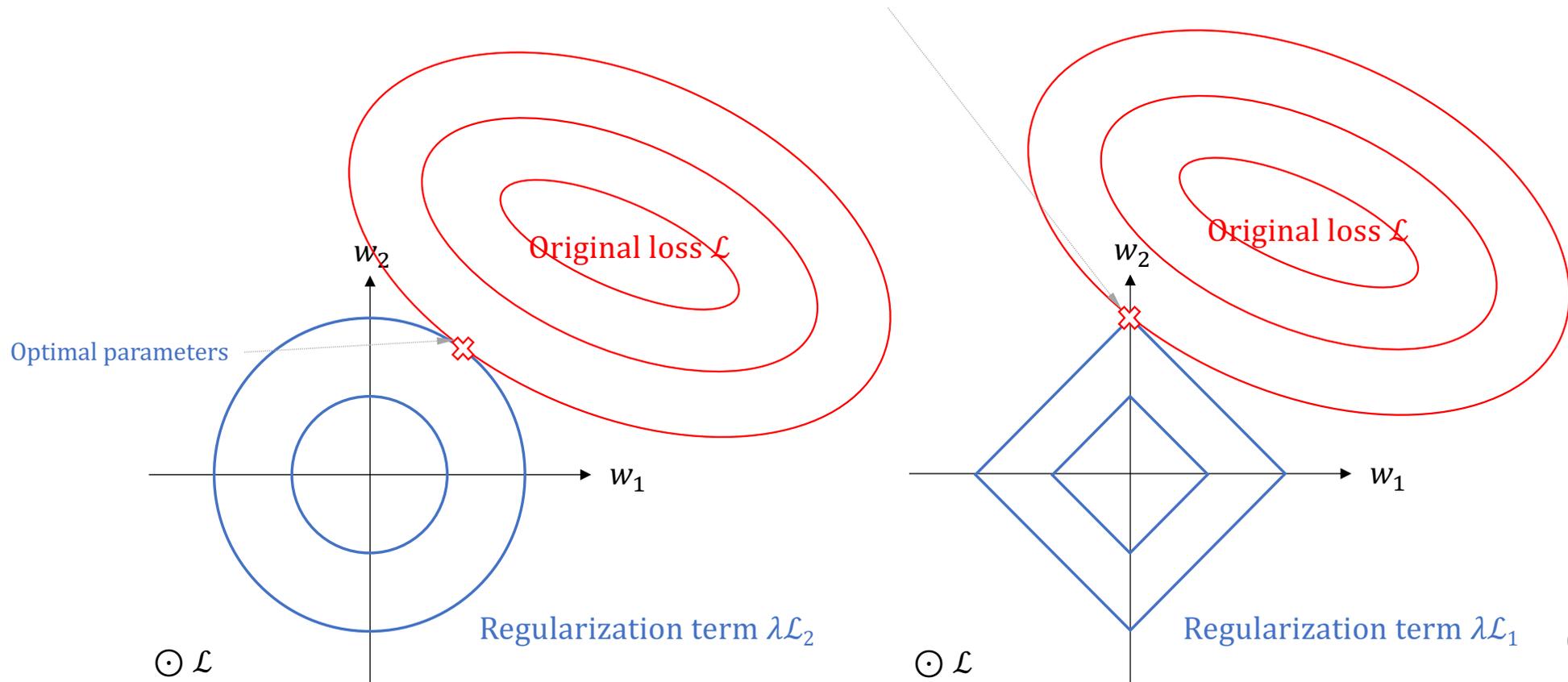
$2 = |w_1| + |w_2|$ is a square with a diagonal line that has a length of 2



Regularization

- \mathcal{L}_1 vs \mathcal{L}_2 norm

\mathcal{L}_1 more likely to has zero weights



Regularization

- \mathcal{L}_1 vs \mathcal{L}_2 norm

The network parameter values are often smaller than 1, so by using \mathcal{L}_2 , a smaller value can result in a much smaller penalty than \mathcal{L}_1 (e.g., $|w| > w^2$ when $w < 1$). In contrast, \mathcal{L}_1 can have a larger penalty than \mathcal{L}_2 for small values.

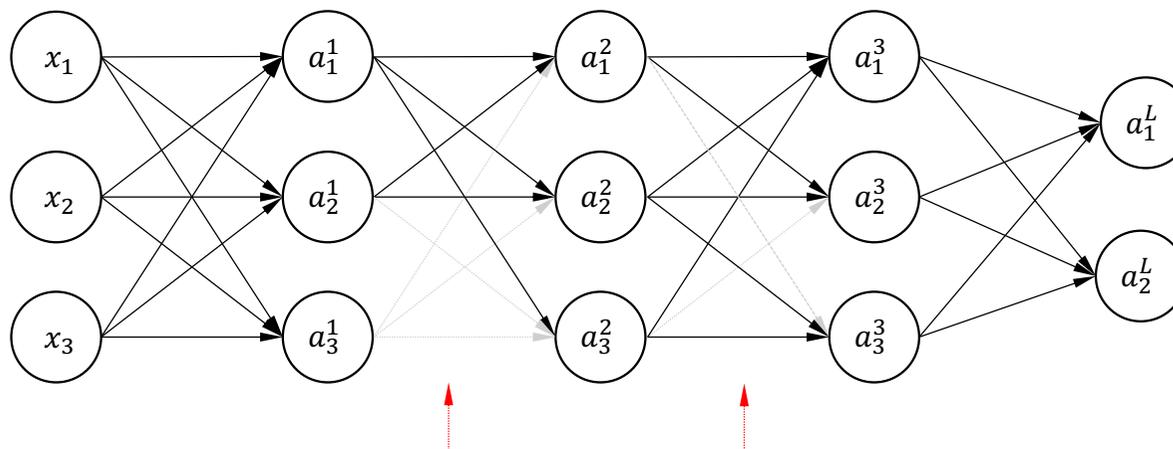
$$|0.5| > 0.5^2$$

Regularization

- \mathcal{L}_1 vs \mathcal{L}_2 norm

\mathcal{L}_1 more likely to have zero parameters, so it has the **sparse property** which enables the networks to perform **“feature selection”**. (ReLU has such property as well)

input layer hidden layer 1 hidden layer 2 hidden layer 3 output layer

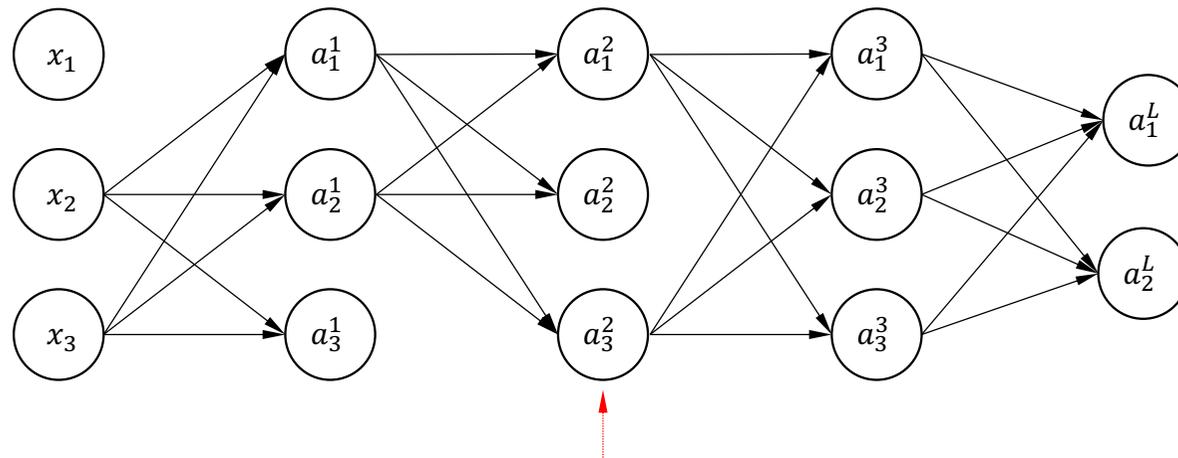


Discarding some input features by setting the corresponding parameters to zero or a very small value

Regularization

- Dropout

input layer hidden layer 1 hidden layer 2 hidden layer 3 output layer



In practice, there will be hundreds and thousands units per layer

- Large neural networks include many parameters making it difficult to deal with the overfitting by combining the predictions from so many parameters.
- Dropout randomly set the hidden outputs to zero, which resembles a random disconnection of the neural units from one layer to the next

Regularization

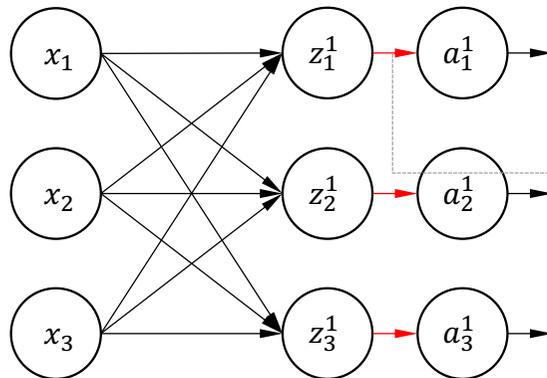
- Dropout

- According to error back-propagation, with a zero-valued output a , the corresponding partial derivative of the loss w.r.t respect to the layer output will be zero. In other words, **only the remaining connected weights will be updated.**
- The dropout method can train many different **sub-networks** while allowing all of them to share the same network parameters.
- **During testing, dropout is disabled**, and no output elements are set to zero. In other words, all sub-networks are used to predict the result represented as using the average from many networks as the final result (i.e., **ensemble learning**)

In practice, only apply Dropout in fully connected layers

Regularization

- Batch Normalization



- Batch normalization is the introduction of a layer that normalizes the inputs to have a mean of 0 and variance of 1 and can improve the performance of a neural network and its training stability.
- During training, the batch normalisation layer estimates the mean and variance of the batch input using a moving average, which updates the moving mean and variance for every iteration to be used for normalizing the batch input.
- During testing, the moving mean and variance are fixed and applied to normalise the input.
- Similar to the dropout process that adds a random factor to the hidden values, the moving mean and variance of batch normalization introduce randomness as they are updated in each iteration according to the random mini-batch. Therefore, the network must learn to be robust enough to deal with the variation.

- Do NOT use bias in the layer before batch norm
- DO NOT use activation operation before batch norm

Implementation

Implementation

Static Model

```

01. import tensorflow as tf
02. from tensorlayer.layers import Input, Dropout, Dense
03. from tensorlayer.models import Model
04.
05. def get_model(inputs_shape):
06.     ni = Input(inputs_shape)
07.     nn = Dropout(keep=0.8)(ni)
08.     nn = Dense(n_units=800, act=tf.nn.relu, name="dense1")(nn)
09.     nn = Dropout(keep=0.8)(nn)
10.     nn = Dense(n_units=800, act=tf.nn.relu)(nn)
11.     nn = Dropout(keep=0.8)(nn)
12.     nn = Dense(n_units=10, act=tf.nn.relu)(nn)
13.     M = Model(inputs=ni, outputs=nn, name="mlp")
14.     return M
15.
16. MLP = get_model([None, 784])
17. MLP.eval()
18. outputs = MLP(data)

```

Lasagne Fashion



<https://tensorlayer.readthedocs.io>

Dynamic Model

```

01. class CustomModel(Model):
02.
03.     def __init__(self):
04.         super(CustomModel, self).__init__()
05.
06.         self.dropout1 = Dropout(keep=0.8)
07.         self.dense1 = Dense(n_units=800, act=tf.nn.relu, in_channels=784)
08.         self.dropout2 = Dropout(keep=0.8)#(self.dense1)
09.         self.dense2 = Dense(n_units=800, act=tf.nn.relu, in_channels=800)
10.         self.dropout3 = Dropout(keep=0.8)#(self.dense2)
11.         self.dense3 = Dense(n_units=10, act=tf.nn.relu, in_channels=800)
12.
13.     def forward(self, x, foo=False):
14.         z = self.dropout1(x)
15.         z = self.dense1(z)
16.         z = self.dropout2(z)
17.         z = self.dense2(z)
18.         z = self.dropout3(z)
19.         out = self.dense3(z)
20.         if foo:
21.             out = tf.nn.relu(out)
22.         return out
23.
24. MLP = CustomModel()
25. MLP.eval()
26. outputs = MLP(data, foo=True) # controls the forward here
27. outputs = MLP(data, foo=False)

```

Chainer Fashion

Implementation

Switching Train/Test Modes

```
01. # method 1: switch before forward
02. Model.train() # enable dropout, batch norm moving avg ...
03. output = Model(train_data)
04. ... # training code here
05. Model.eval() # disable dropout, batch norm moving avg ...
06. output = Model(test_data)
07. ... # testing code here
08.
09. # method 2: switch while forward
10. output = Model(train_data, is_train=True)
11. output = Model(test_data, is_train=False)
```

Implementation

Reuse Weights in Static Model

```
01. def create_base_network(input_shape):
02.     """Base network to be shared (eq. to feature extraction).
03.     """
04.     input = Input(shape=input_shape)
05.     x = Flatten()(input)
06.     x = Dense(128, act=tf.nn.relu)(x)
07.     x = Dropout(0.9)(x)
08.     x = Dense(128, act=tf.nn.relu)(x)
09.     x = Dropout(0.9)(x)
10.     x = Dense(128, act=tf.nn.relu)(x)
11.     return Model(input, x)
12.
13.
14. def get_siamese_network(input_shape):
15.     """Create siamese network with shared base network as layer
16.     """
17.     base_layer = create_base_network(input_shape).as_layer() # convert model as layer
18.
19.     ni_1 = Input(input_shape)
20.     ni_2 = Input(input_shape)
21.     nn_1 = base_layer(ni_1) # call base_layer twice
22.     nn_2 = base_layer(ni_2)
23.     return Model(inputs=[ni_1, ni_2], outputs=[nn_1, nn_2])
24.
25. siamese_net = get_siamese_network([None, 784])
```

Reuse Weights in Dynamic Model

```
01. class MyModel(Model):
02.     def __init__(self):
03.         super(MyModel, self).__init__()
04.         self.dense_shared = Dense(n_units=800, act=tf.nn.relu, in_channels=784)
05.         self.dense1 = Dense(n_units=10, act=tf.nn.relu, in_channels=800)
06.         self.dense2 = Dense(n_units=10, act=tf.nn.relu, in_channels=800)
07.         self.cat = Concat()
08.
09.     def forward(self, x):
10.         x1 = self.dense_shared(x) # call dense_shared twice
11.         x2 = self.dense_shared(x)
12.         x1 = self.dense1(x1)
13.         x2 = self.dense2(x2)
14.         out = self.cat([x1, x2])
15.         return out
16.
17. model = MyModel()
```

Implementation

Print Model Architecture

```

01. import pprint
02.
03. def get_model(inputs_shape):
04.     ni = Input(inputs_shape)
05.     nn = Dropout(keep=0.8)(ni)
06.     nn = Dense(n_units=800, act=tf.nn.relu)(nn)
07.     nn = Dropout(keep=0.8)(nn)
08.     nn = Dense(n_units=800, act=tf.nn.relu)(nn)
09.     nn = Dropout(keep=0.8)(nn)
10.     nn = Dense(n_units=10, act=tf.nn.relu)(nn)
11.     M = Model(inputs=ni, outputs=nn, name="mlp")
12.     return M
13.
14. MLP = get_model([None, 784])
15. pprint.pprint(MLP.config)

```

```

[{'args': {'dtype': tf.float32,
          'layer_type': 'normal',
          'name': '_inputlayer_1',
          'shape': [None, 784]},
  'class': '_InputLayer',
  'prev_layer': None},
 {'args': {'keep': 0.8, 'layer_type': 'normal', 'name': 'dropout_1'},
  'class': 'Dropout',
  'prev_layer': ['_inputlayer_1_node_0']},
 {'args': {'act': 'relu',
          'layer_type': 'normal',
          'n_units': 800,
          'name': 'dense_1'},
  'class': 'Dense',
  'prev_layer': ['dropout_1_node_0']},
 {'args': {'keep': 0.8, 'layer_type': 'normal', 'name': 'dropout_2'},
  'class': 'Dropout',
  'prev_layer': ['dense_1_node_0']},
 {'args': {'act': 'relu',
          'layer_type': 'normal',
          'n_units': 800,
          'name': 'dense_2'},
  'class': 'Dense',
  'prev_layer': ['dropout_2_node_0']},
 {'args': {'keep': 0.8, 'layer_type': 'normal', 'name': 'dropout_3'},
  'class': 'Dropout',
  'prev_layer': ['dense_2_node_0']},
 {'args': {'act': 'relu',
          'layer_type': 'normal',
          'n_units': 10,
          'name': 'dense_3'},
  'class': 'Dense',
  'prev_layer': ['dropout_3_node_0']}

```

Implementation

Print Model Information

```
01. print(MLP) # simply call print function
02.
03. # Model(
04. #   (_inputlayer): Input(shape=[None, 784], name='_inputlayer')
05. #   (dropout): Dropout(keep=0.8, name='dropout')
06. #   (dense): Dense(n_units=800, relu, in_channels='784', name='dense')
07. #   (dropout_1): Dropout(keep=0.8, name='dropout_1')
08. #   (dense_1): Dense(n_units=800, relu, in_channels='800', name='dense_1')
09. #   (dropout_2): Dropout(keep=0.8, name='dropout_2')
10. #   (dense_2): Dense(n_units=10, relu, in_channels='800', name='dense_2')
11. # )
```

Save Weights Only

```
01. MLP.save_weights('./model_weights.h5')
02. MLP.load_weights('./model_weights.h5')
```

Get Specific Weights

```
01. # indexing
02. all_weights = MLP.all_weights
03. some_weights = MLP.all_weights[1:3]
04.
05. # naming
06. some_weights = MLP.get_layer('dense1').all_weights
```

Save Weights + Architecture

```
01. MLP.save('./model.h5', save_weights=True)
02. MLP = Model.load('./model.h5', load_weights=True)
```

Implementation

Customized layers with weights

```
class Dense(Layer):
    """The :class:`Dense` class is a fully connected layer.

    Parameters
    -----
    n_units : int
        The number of units of this layer.
    act : activation function
        The activation function of this layer.
    name : None or str
        A unique layer name. If None, a unique name will be automatically generated.
    """

    def __init__(
        self,
        n_units, # the number of units/channels of this layer
        act=None, # None: no activation, tf.nn.relu: ReLU ...
        name=None, # the name of this layer (optional)
    ):
        super(Dense, self).__init__(name) # auto naming, dense_1, dense_2 ...
        self.n_units = n_units
        self.act = act

    def build(self, inputs_shape): # initialize the model weights here
        shape = [inputs_shape[1], self.n_units]
        self.W = self._get_weights("weights", shape=tuple(shape), init=self.W_init)
        self.b = self._get_weights("biases", shape=(self.n_units, ), init=self.b_init)

    def forward(self, inputs): # call function
        z = tf.matmul(inputs, self.W) + self.b
        if self.act: # is not None
            z = self.act(z)
        return z
```

```
class Dense(Layer):
    """The :class:`Dense` class is a fully connected layer.

    Parameters
    -----
    n_units : int
        The number of units of this layer.
    act : activation function
        The activation function of this layer.
    W_init : initializer
        The initializer for the weight matrix.
    b_init : initializer or None
        The initializer for the bias vector. If None, skip biases.
    in_channels: int
        The number of channels of the previous layer.
        If None, it will be automatically detected when the layer is forwarded for the first time.
    name : None or str
        A unique layer name. If None, a unique name will be automatically generated.
    """

    def __init__(
        self,
        n_units,
        act=None,
        W_init=tl.initializers.truncated_normal(stddev=0.1),
        b_init=tl.initializers.constant(value=0.0),
        in_channels=None, # the number of units/channels of the previous layer
        name=None,
    ):
        # we feed activation function to the base layer, `None` denotes identity function
        # string (e.g., relu, sigmoid) will be converted into function.
        super(Dense, self).__init__(name, act=act)

        self.n_units = n_units
        self.W_init = W_init
        self.b_init = b_init
        self.in_channels = in_channels

        # in dynamic model, the number of input channel is given, we initialize the weights here
        if self.in_channels is not None:
            self.build(self.in_channels)
            self._built = True

        logging.info(
            "Dense %s: %d %s" %
            (self.name, self.n_units, self.act.__name__ if self.act is not None else 'No Activation')
        )

    def build(self, inputs_shape): # initialize the model weights here
        if self.in_channels: # if the number of input channel is given, use it
            shape = [self.in_channels, self.n_units]
        else:
            # otherwise, get it from static model
            self.in_channels = inputs_shape[1]
            shape = [inputs_shape[1], self.n_units]
        self.W = self._get_weights("weights", shape=tuple(shape), init=self.W_init)
        if self.b_init: # if b_init is None, no bias is applied
            self.b = self._get_weights("biases", shape=(self.n_units, ), init=self.b_init)

    def forward(self, inputs):
        z = tf.matmul(inputs, self.W)
        if self.b_init:
            z = tf.add(z, self.b)
        if self.act:
            z = self.act(z)
        return z
```



Implementation

Customized layers with train/test modes

```
class Dropout(Layer):
    """
    The :class:`Dropout` class is a noise layer which randomly set some
    activations to zero according to a keeping probability.
    Parameters
    -----
    keep : float
        The keeping probability.
        The lower the probability it is, the more activations are set to zero.
    name : None or str
        A unique layer name.
    """

    def __init__(self, keep, name=None):
        super(Dropout, self).__init__(name)
        self.keep = keep

        self.build()
        self._built = True

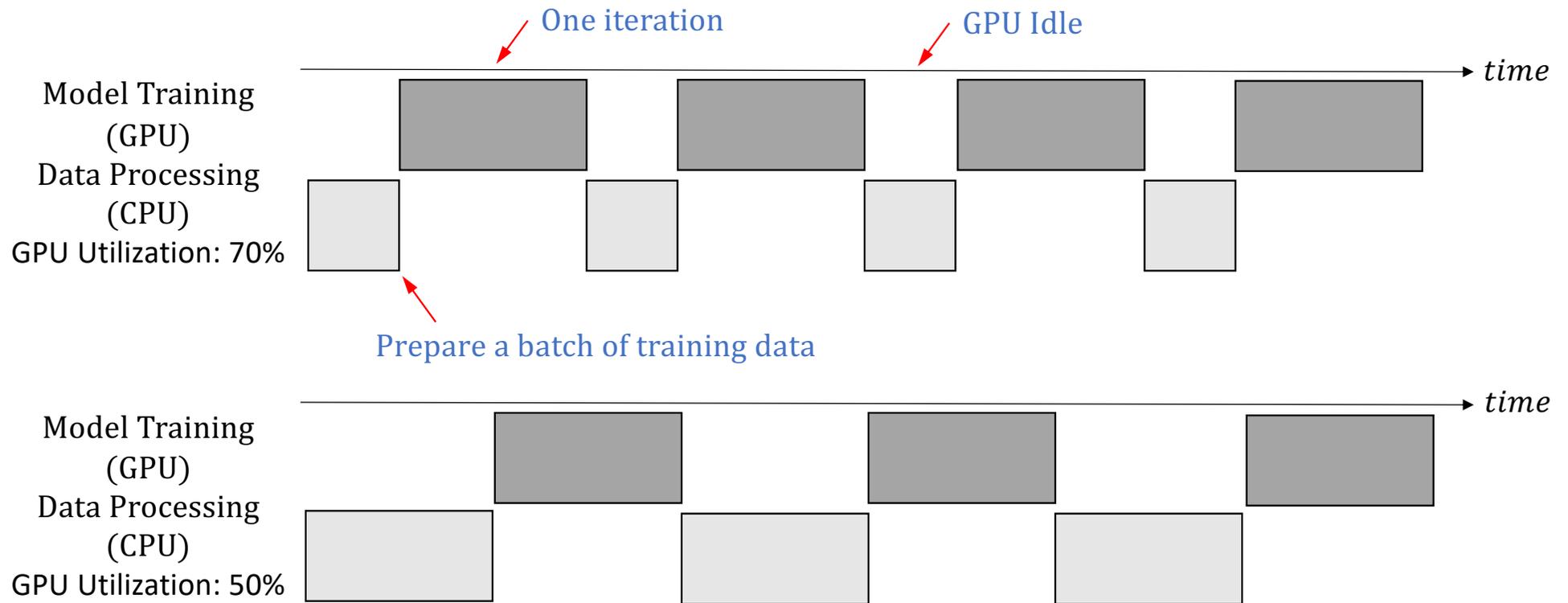
        logging.info("Dropout %s: keep: %f " % (self.name, self.keep))

    def build(self, inputs_shape=None):
        pass # no weights in dropout layer

    def forward(self, inputs):
        if self.is_train: # this attribute is changed by Model.train() and Model.eval()
            outputs = tf.nn.dropout(inputs, rate=1 - (self.keep), name=self.name)
        else:
            outputs = inputs
        return outputs
```

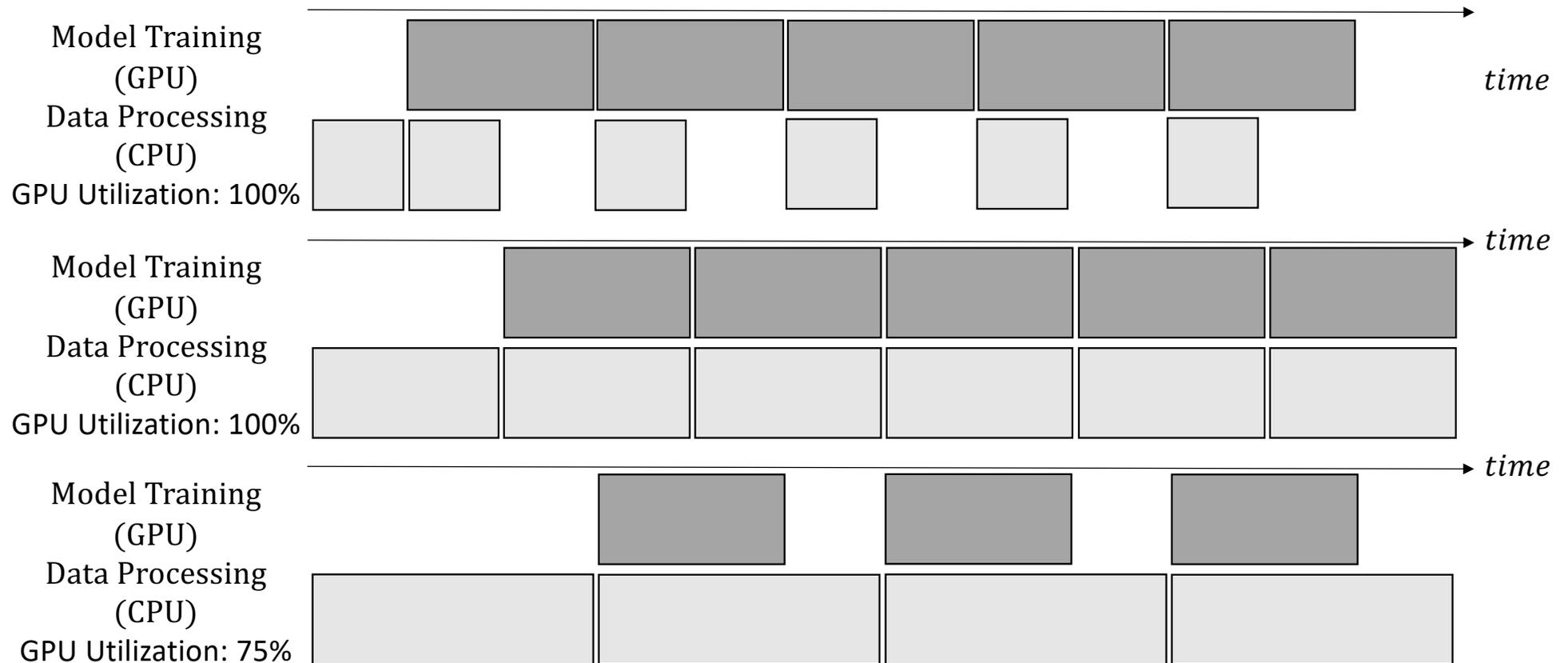
Implementation

- Training without Dataflow



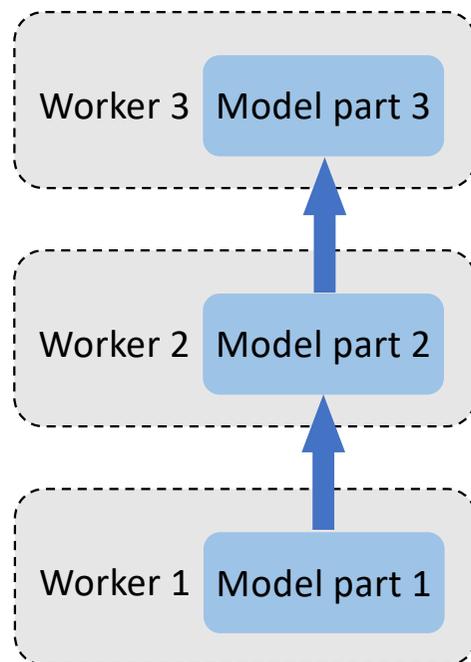
Implementation

- Training with Dataflow

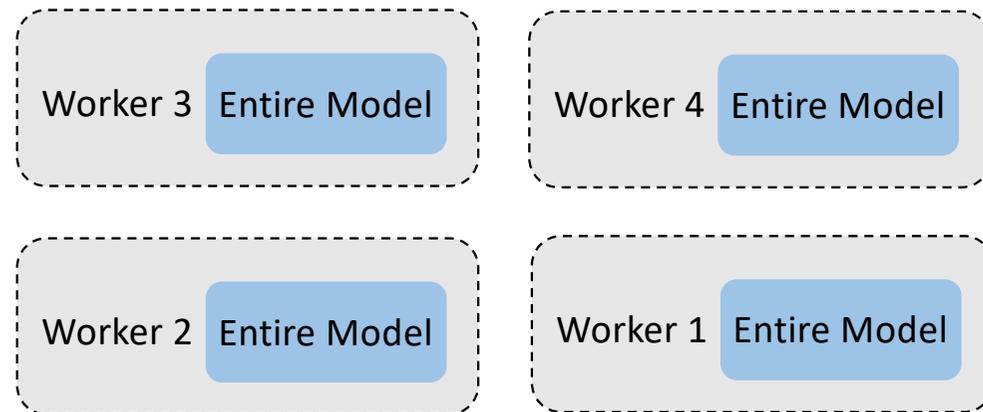


Implementation

- Distributed Training



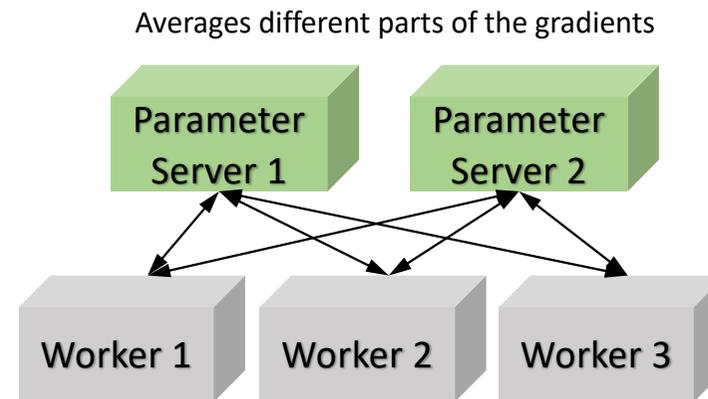
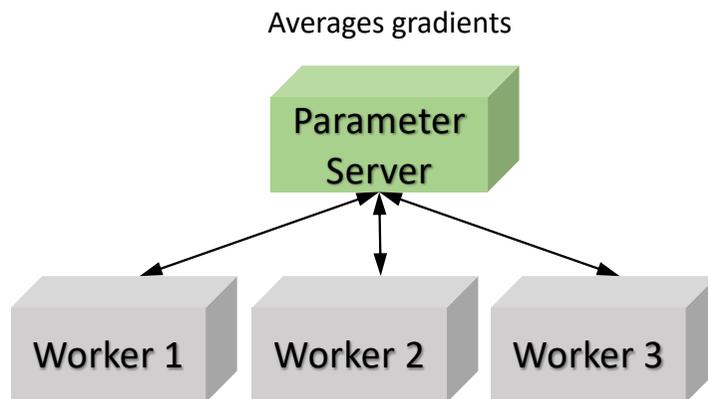
Model Parallelism



Data Parallelism

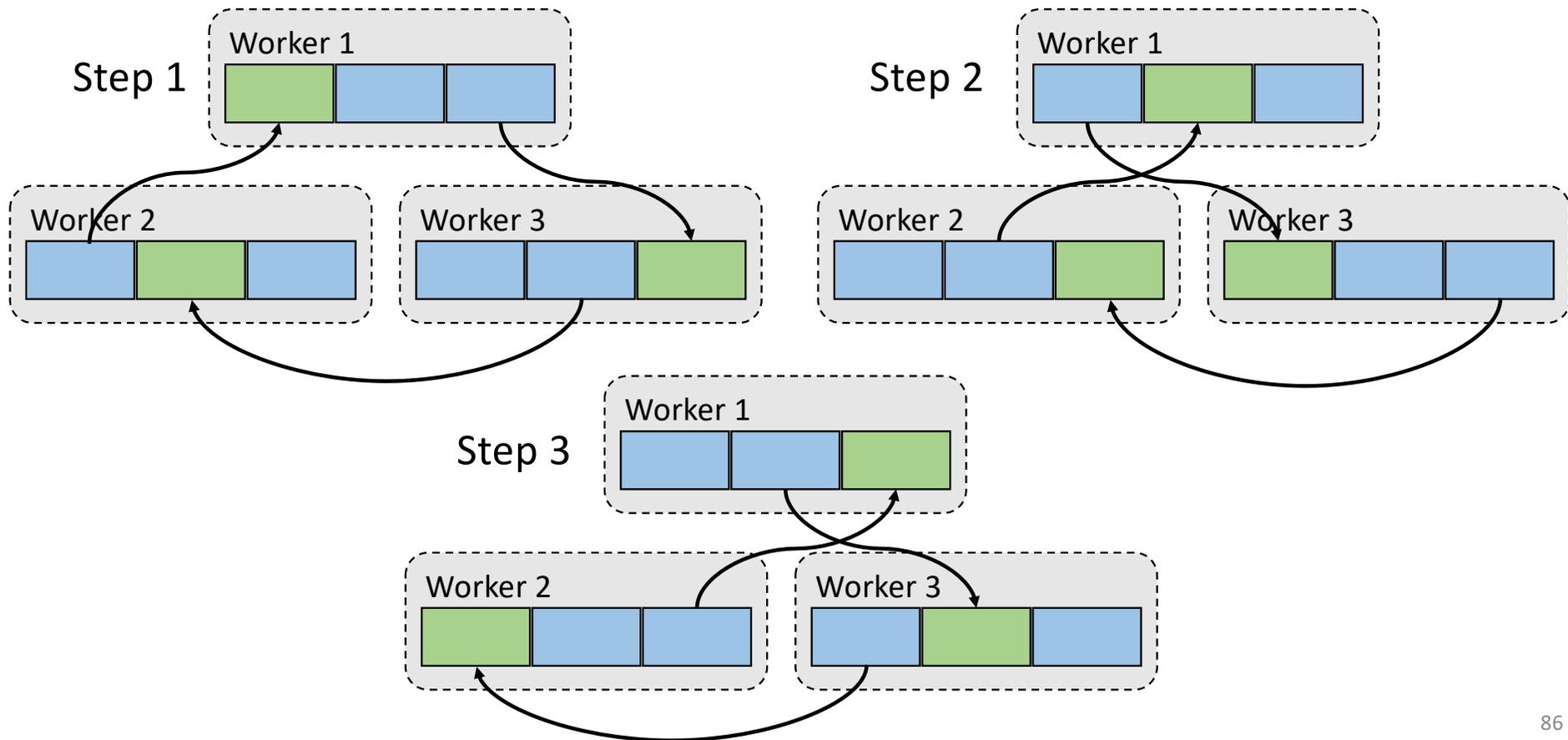
Implementation

- Distributed Training: Data Parallelism - Parameter Server



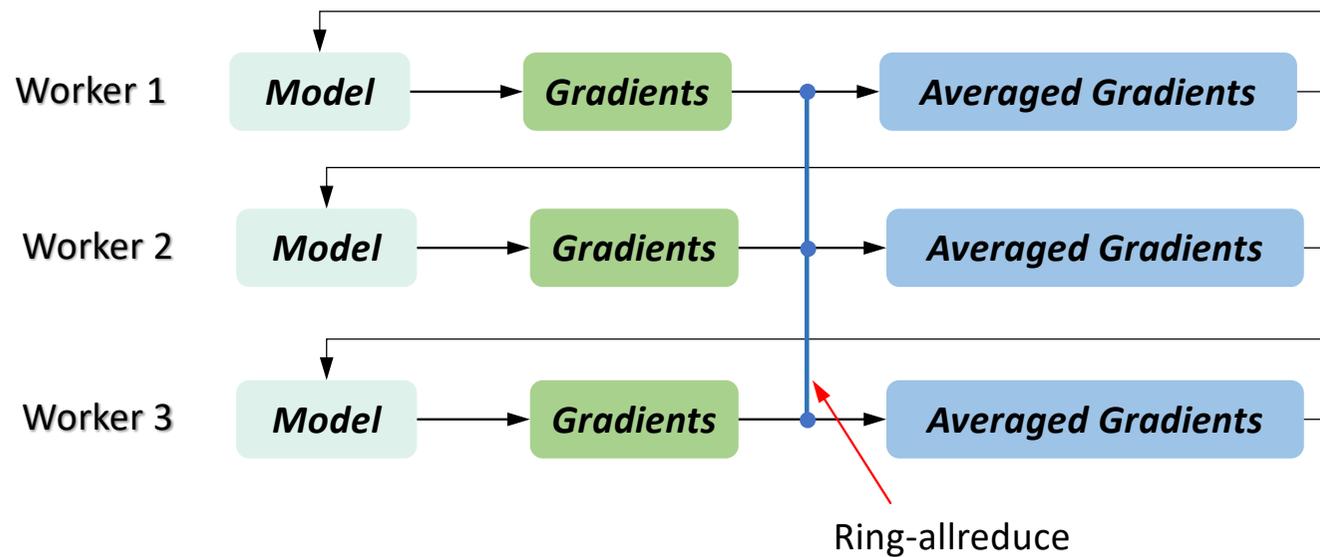
Implementation

- Distributed Training: Data Parallelism - Horovod - All ringreduce



Implementation

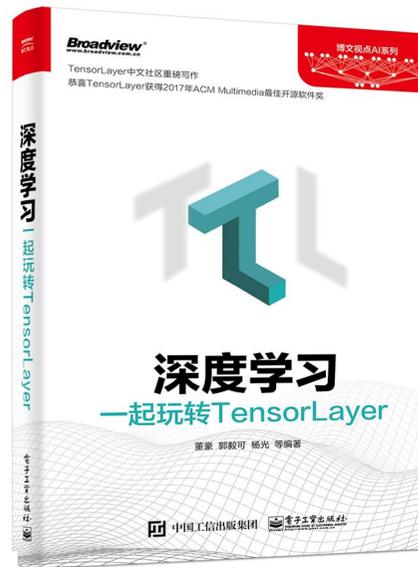
- Distributed Training: Data Parallelism - Horovod - All ringreduce



Implementation

- TensorLayer 2.0

5000 Stars ~ 1000 Forks ~ 180k Downloads



Summary

Deep Neural Network Foundation

- **Single Neuron**
 - The neuron can be represented by the matrix multiplication.
 - A bias value allows the output value to be shifted higher or lower to better fit the input data.
- **Activation Functions**
 - Provide non-linearity to the network.
- **Multi-layer Perceptron**
 - Higher representation capacity.
- **Loss Functions**
 - The goal of the optimization.
- **Optimization**
 - Gradient descent → Error back-propagation → SGD → Adaptive learning rate
- **Regularization**
 - Weight decay, dropout, batch norm ...
- **Implementation**
 - Define static/dynamic models, dataflow, distributed training ...

Deep Neural Network Foundation

Link: <https://github.com/zsdonghao/deep-learning-note>

- Exercise 1:
 - Implement the error back-propagation with the provided Numpy template.
Keywords: regression, classification, Sigmoid, MSE, logistic regression, fully connected layer
- Exercise 2:
 - Use TensorLayer to classify MNIST handwritten digit dataset including training, validating, and testing.
Keywords: classification, cross-entropy, fully connected layer, evaluation
- Exercise 3:
 - Classify the MNIST dataset by modifying the code from exercise 1.
Keywords: classification, Softmax, ReLU, cross-entropy, fully connected layer
- Exercise 4: (Optional)
 - Following exercise 2, implement a dataflow for data augmentation
Keywords: classification, Softmax, ReLU, cross-entropy, fully connected layer

Questions?