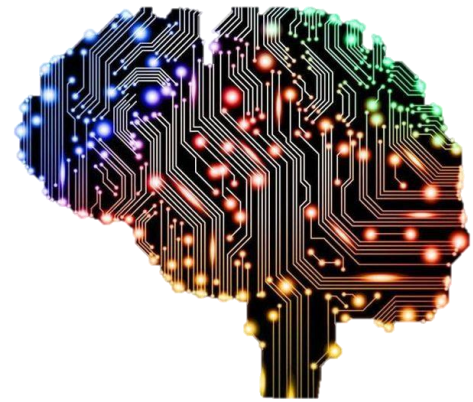


智慧化企業整合 Intelligent Integration of Enterprise

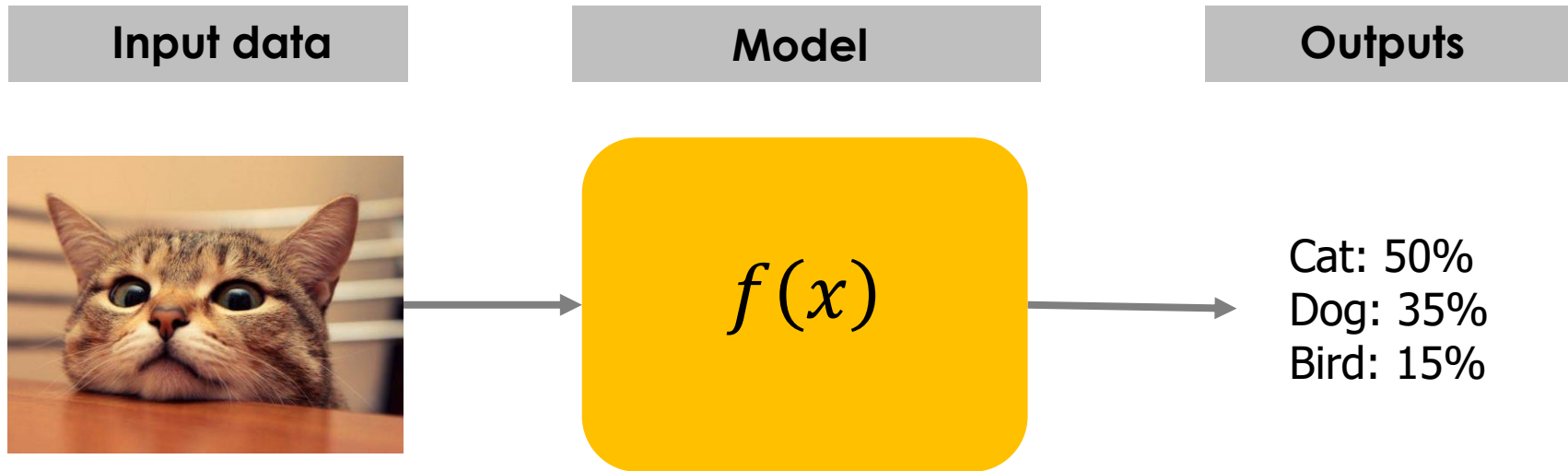
Introduction of Deep Learning

Outline

- Neural Network
- Loss function
- Gradient Descent



System



Steps of Define Functions

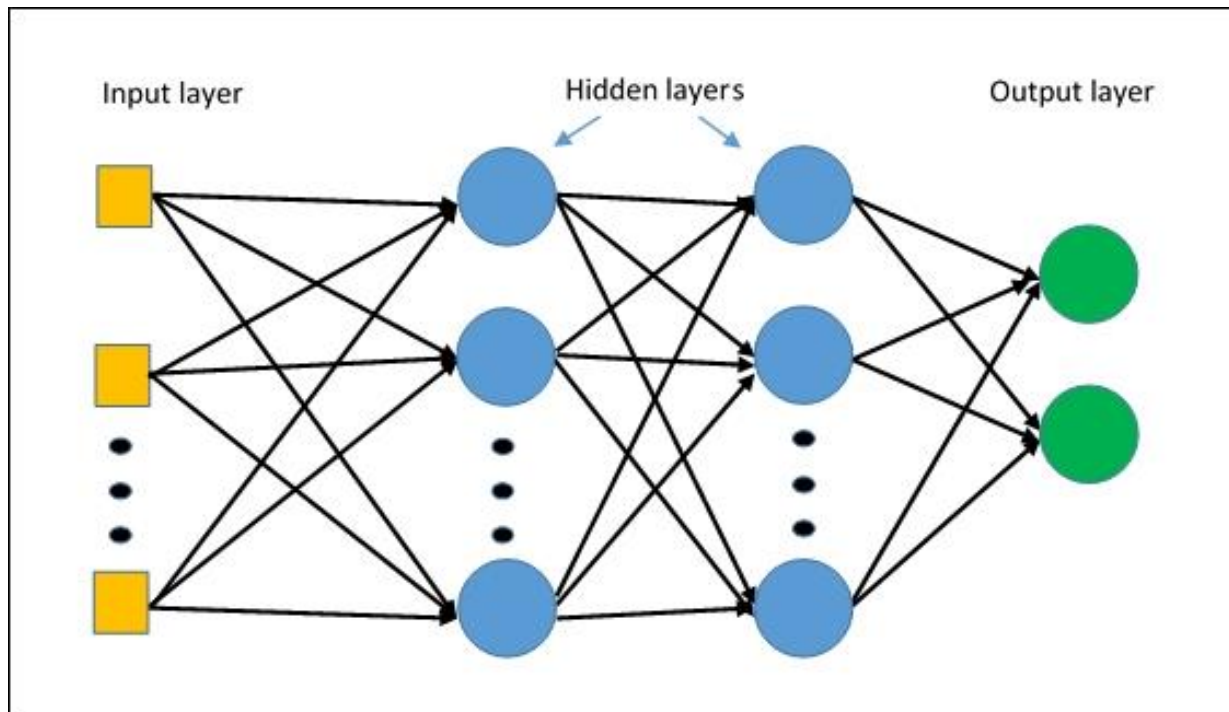
Step1: Define your network (function) structure

Step2: Measure the goodness of your function

Step3: Find the best function

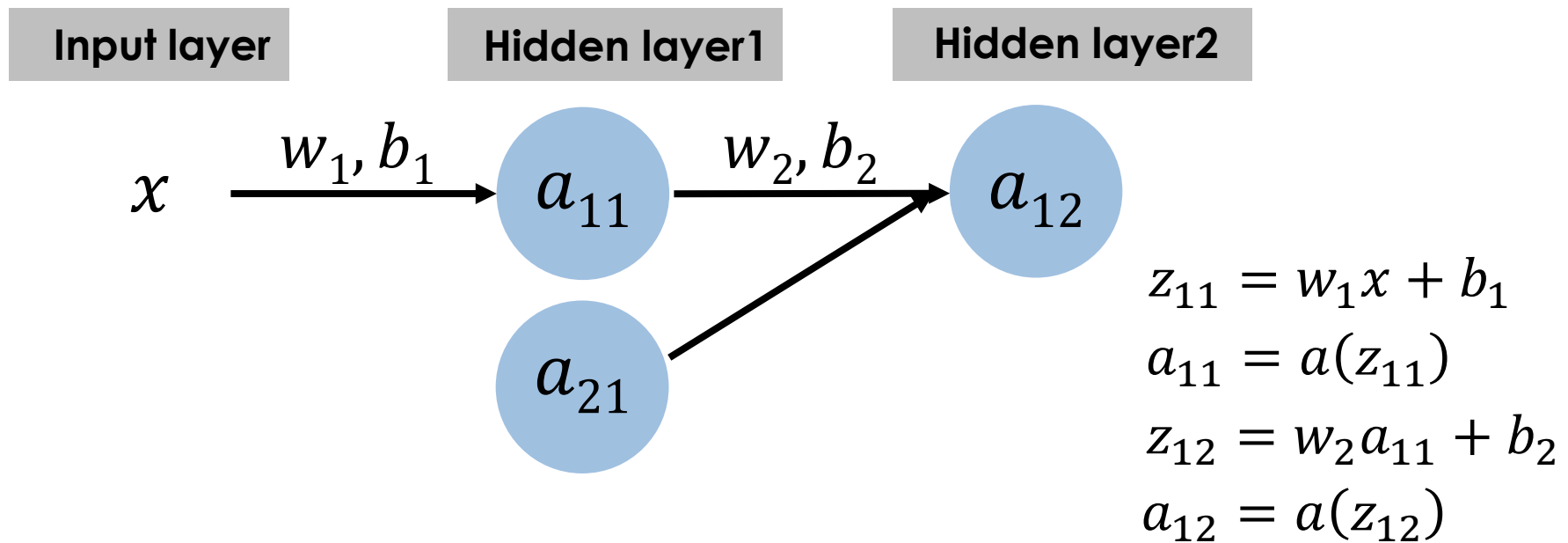
Step 1: Function Structure

Neural Network



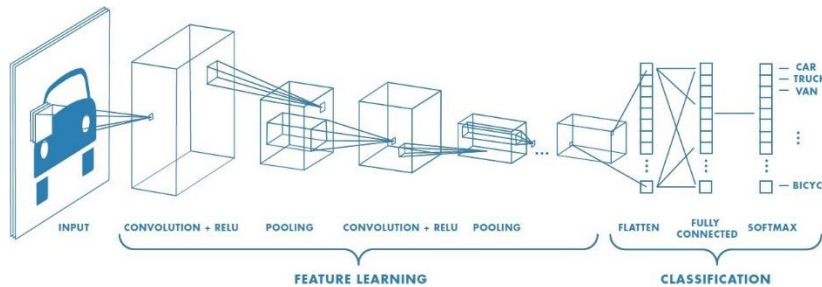
- Fully connected feedforward neural network
- A large function set that consists of lots of variables

Neurons

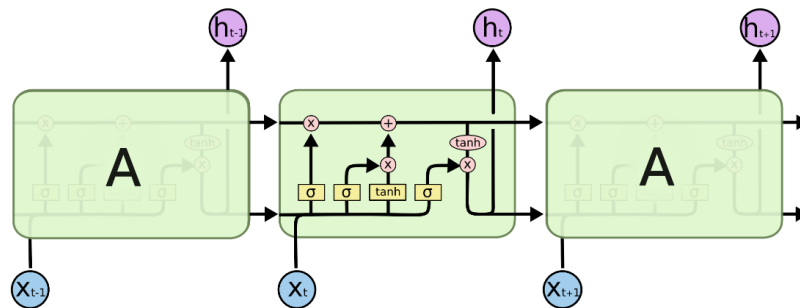


- $a(z)$ is an activation function
- $a(z)$ could be different in different layers
- # of neurons in different layers could be different

Common Architectures



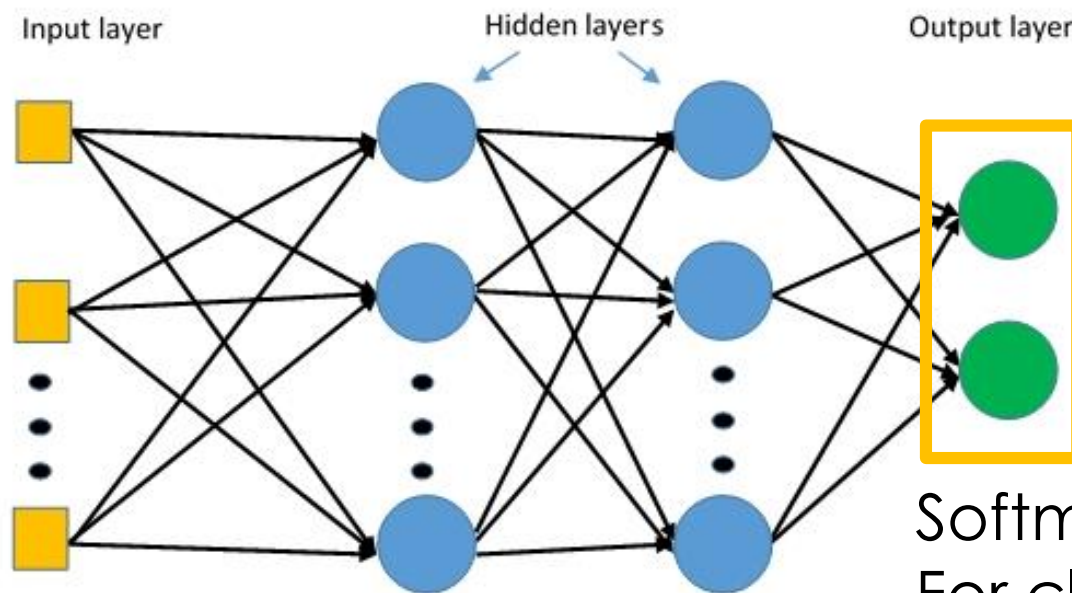
Convolutional neural network (CNN)



Recurrent neural network (RNN/LSTM cell)

The repeating module in an LSTM contains four interacting layers.

Output layer



Softmax Function
For classification
task

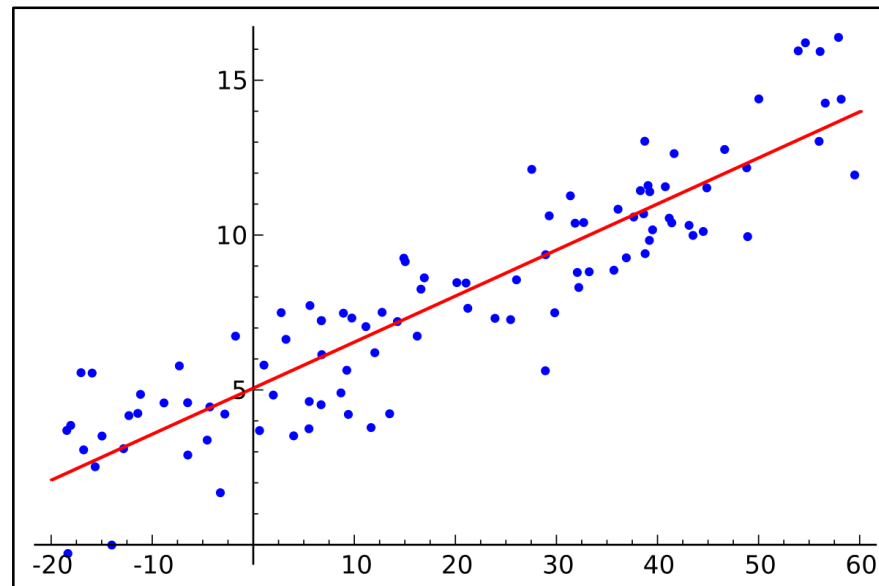
Softmax Function

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$

- Output : [0,1]
- Sum of outputs is 1
- Simulation of probability
- Usually used in classification problems

Step2: Goodness of Function

Loss Function



- A measure of goodness
- Distance (error rate) between predictions and true labels
- Common loss functions: MSE, MAE, Cross entropy...





Loss Function

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

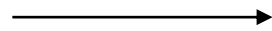
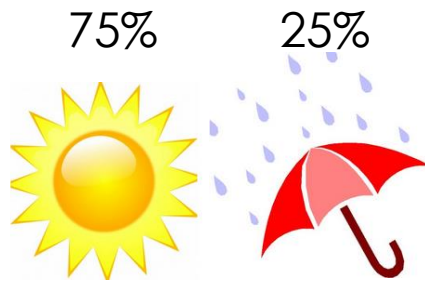
- Distance (error rate)
- Minimize loss function (Step3)

Shannon's Entropy

50%	50%		
		→	$-0.5 \times \log_2 0.5 - 0.5 \times \log_2 0.5 = 1$
75%	25%		
		→	$-0.75 \times \log_2 0.75 - 0.25 \times \log_2 0.25$ $= 0.31 + 0.5 = 0.81$

$$S = - \sum_i P_i \log P_i.[1]$$

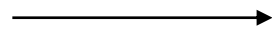
Cross Entropy



$$-0.75 \times \log_2 0.75 - 0.25 \times \log_2 0.25$$

$$= 0.31 + 0.5 = 0.81$$

25% 75%



$$-0.75 \times \log_2 0.25 - 0.25 \times \log_2 0.75$$

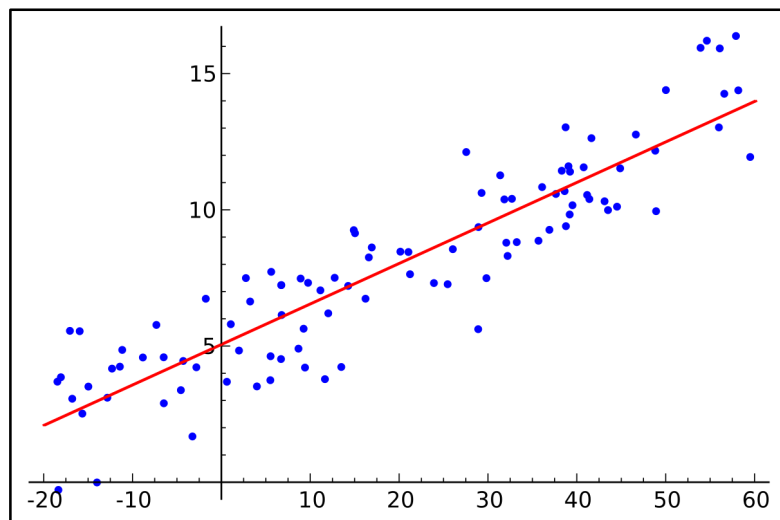
$$= 1.5 + 0.1 = 1.6$$

$$H(p, q) = - \sum_x p(x) \log q(x).$$

- Usually used in classification task
- Predicted probability distribution: $q(x)$
- Real probability distribution: $p(x)$
- Goal: minimize cross entropy

Step3: The best Function

Goal



- Assume $y = wx + b$
- w, b are coefficients we want to learn
- Minimize loss function (EX: MSE)

$$L(w, b) = \sum_{n=1}^k (\hat{y}^n - (b + w \times x^n))^2$$

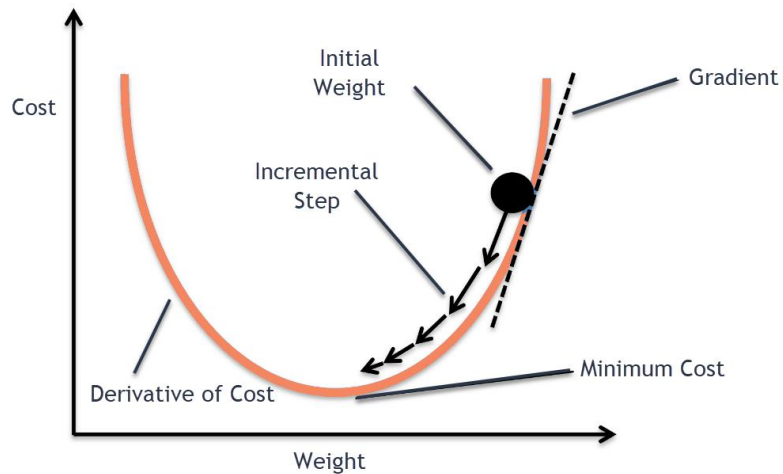
Gradient

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \\ \vdots \end{bmatrix}$$

- Partial derivative of all variable
- The direction you should travel to **increase** the value of f most rapidly

$$L(w, b) = \sum_{n=1}^k (\hat{y}^n - (b + w \times x^n))^2$$

Gradient Descent



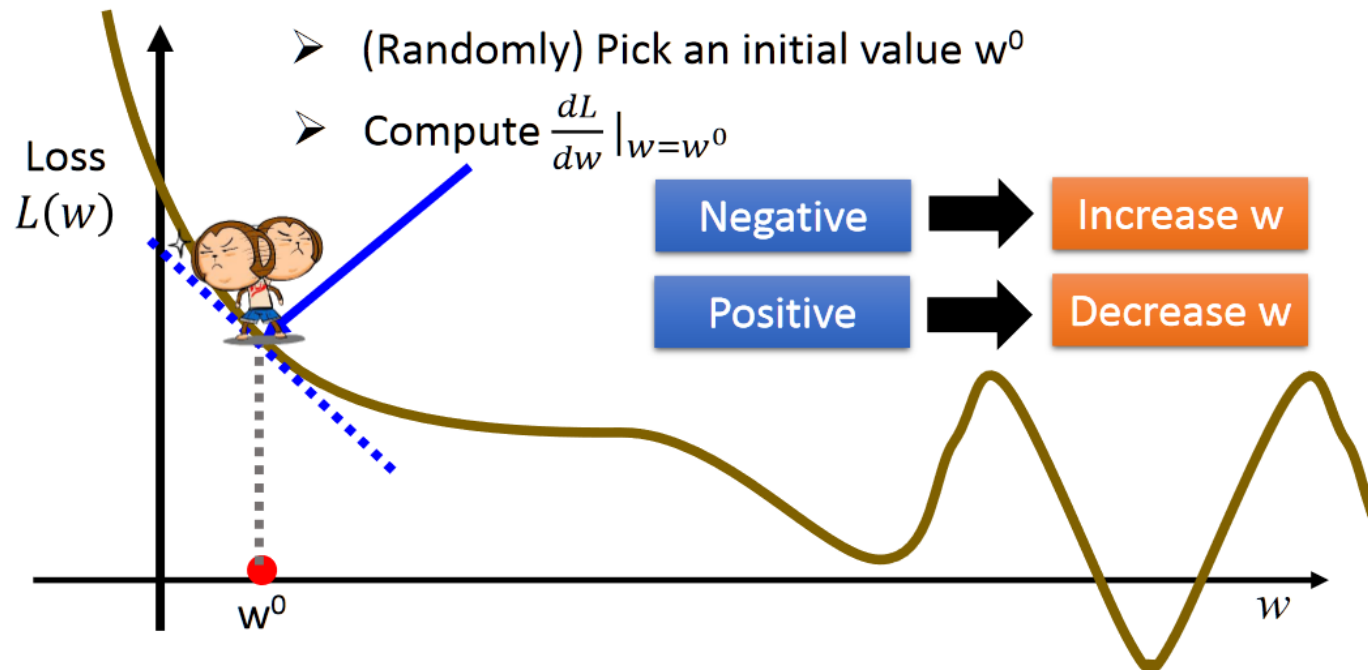
$$L(w, b) = \sum_{n=1}^k (\hat{y}^n - (b + w \times x^n))^2$$

- w, b are coefficients we want to learn
- We want to minimize $L(w, b)$
 1. Choose an initial weight
 2. Calculate the gradient
 3. Move in the opposite direction of the gradient

Gradient Descent

$$w^* = \underset{w}{\operatorname{arg\,min}} L(w)$$

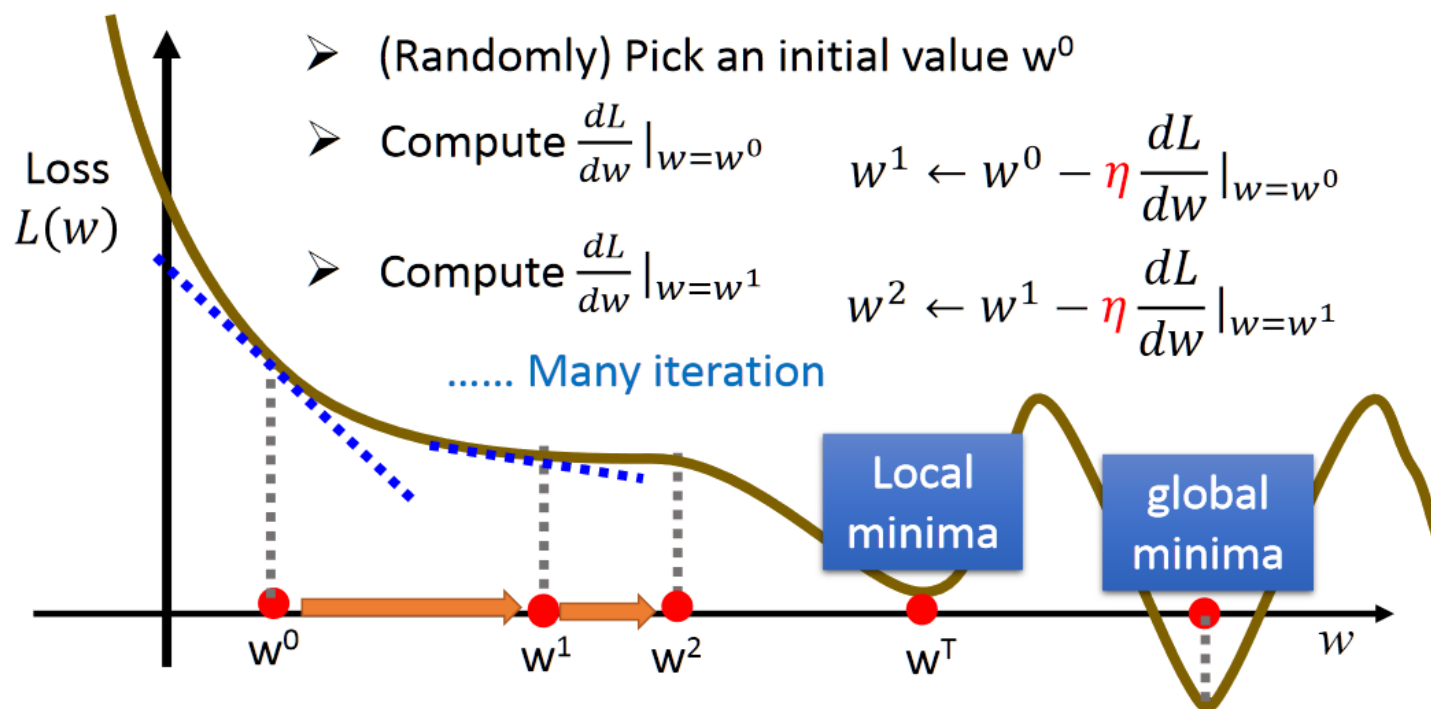
- Consider loss function $L(w)$ with one parameter w :



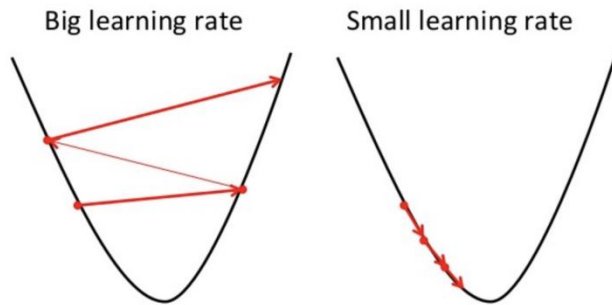
Gradient Descent

$$w^* = \underset{w}{\operatorname{arg\,min}} L(w)$$

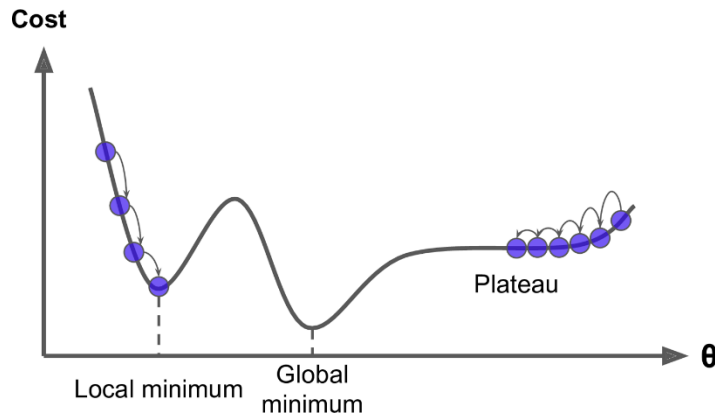
- Consider loss function $L(w)$ with one parameter w :



Problems of Gradient Descent



- Small learning rate \rightarrow slow training process
- Big learning rate \rightarrow difficult to converge



- Training process may be hindered by local minimum or plateau
- We need to adjust learning rate to improve this process (Adagrad, Momentum, Adam...)

References

- [Regression and gradient descent by Hung-Yi Lee](#)
- [Introduction of entropy](#)
- [Introduction of loss function](#)

Class assignment

- Please search , read and learn the definition of the following terms.
- Write down your own introduction of these terms in your report.
- Turn in your report with the format of pdf
- (maximum: 3 pages).
 - 1. Stochastic gradient descent (SGD)**
 - 2. Momentum (in machine learning)**
 - 3. Adagrad**
 - 4. RMS prop**

Homework

- Please refer to the 'gradient descent.ipynb' and see how this algorithm updates parameters
- Please modify some hyper parameters of this notebook and illustrate of those modifications with corresponding effects in this notebook
- Turn in this homework with the format of ipynb