

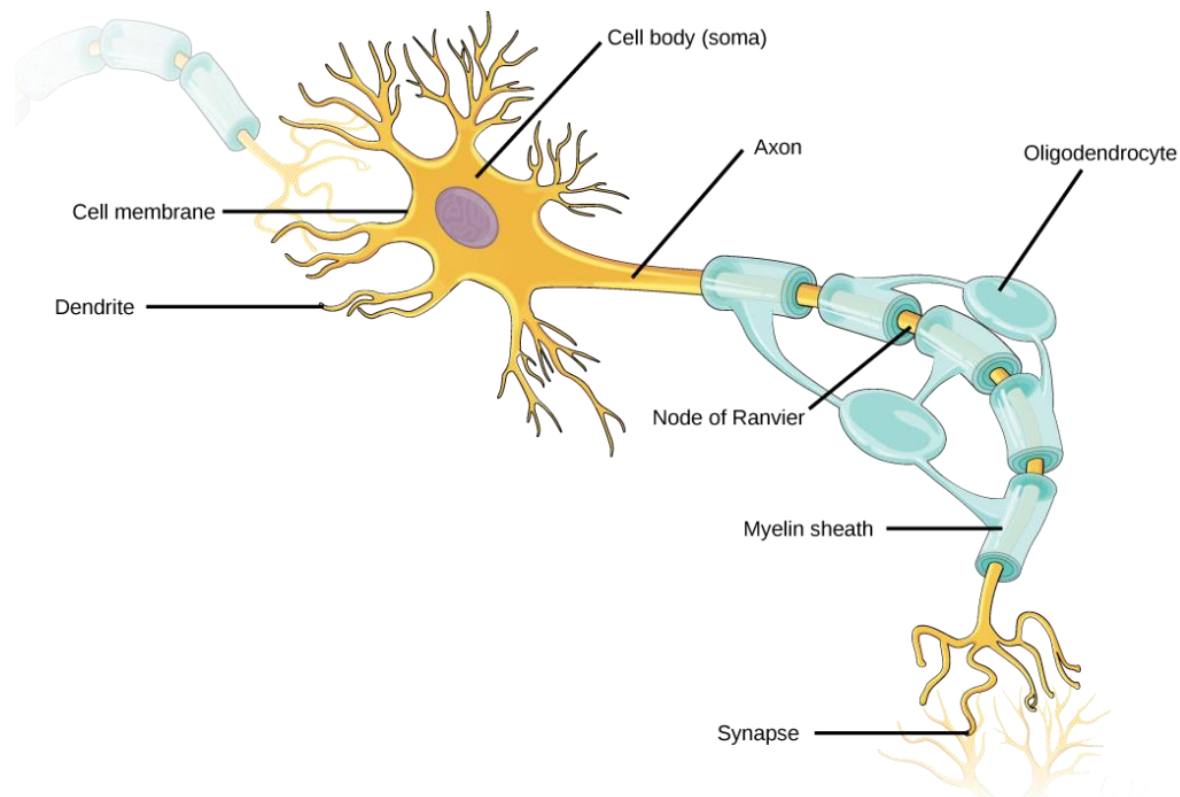
Artificial Neural Network for Classification (1)

Hanwool Jeong

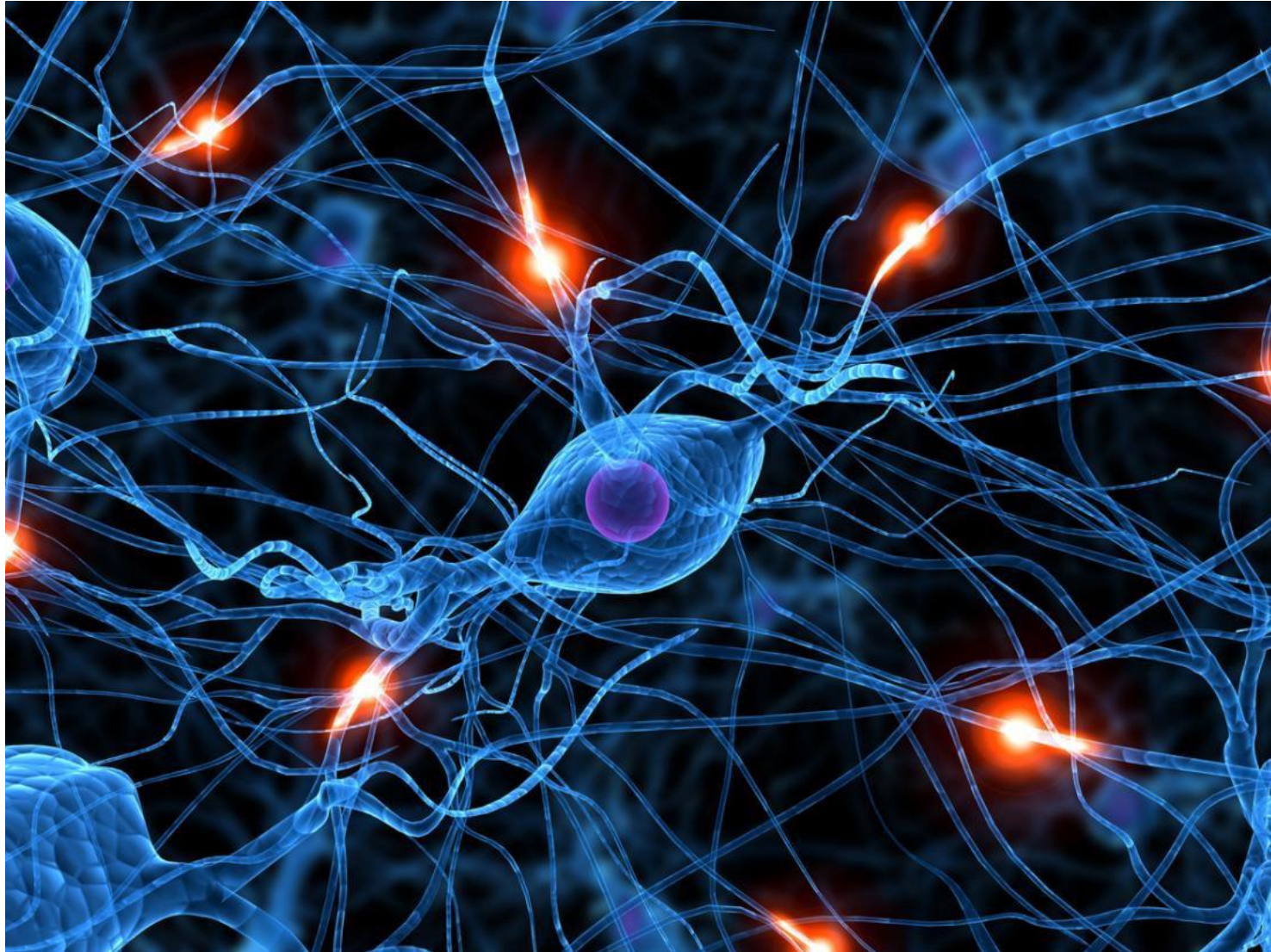
hwjeong@kw.ac.kr

Neuron in Human Brain

- Considered as Information processing unit
 - Dendrite accepts signals from other neurons
 - Axon conducts electrical impulses away from the cell body.

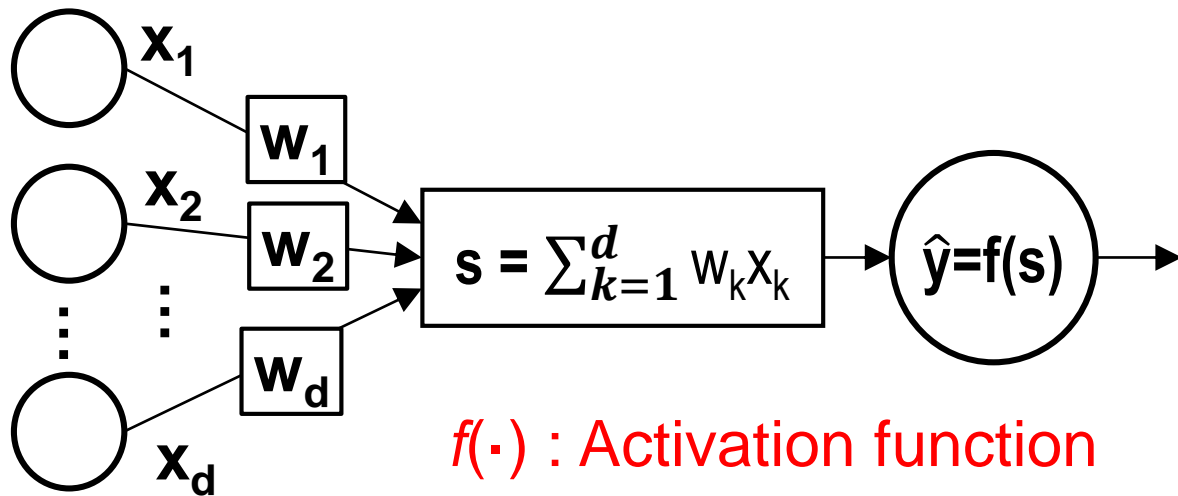


Neural Network

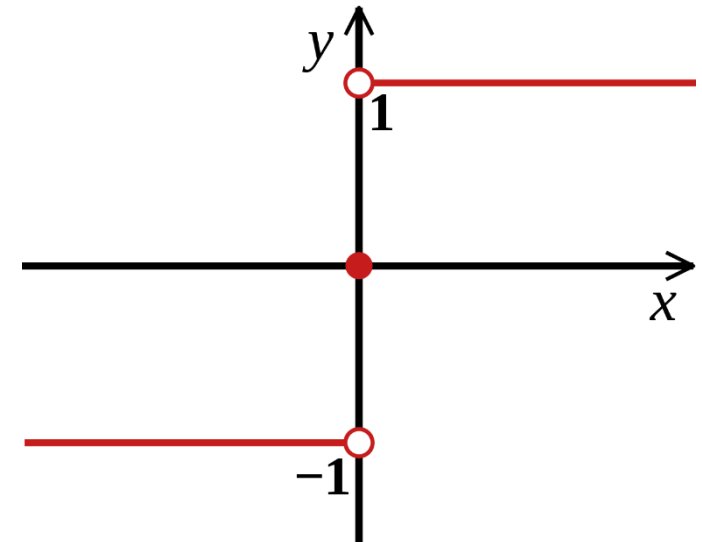


Perceptron

- Structure and operation of perceptron



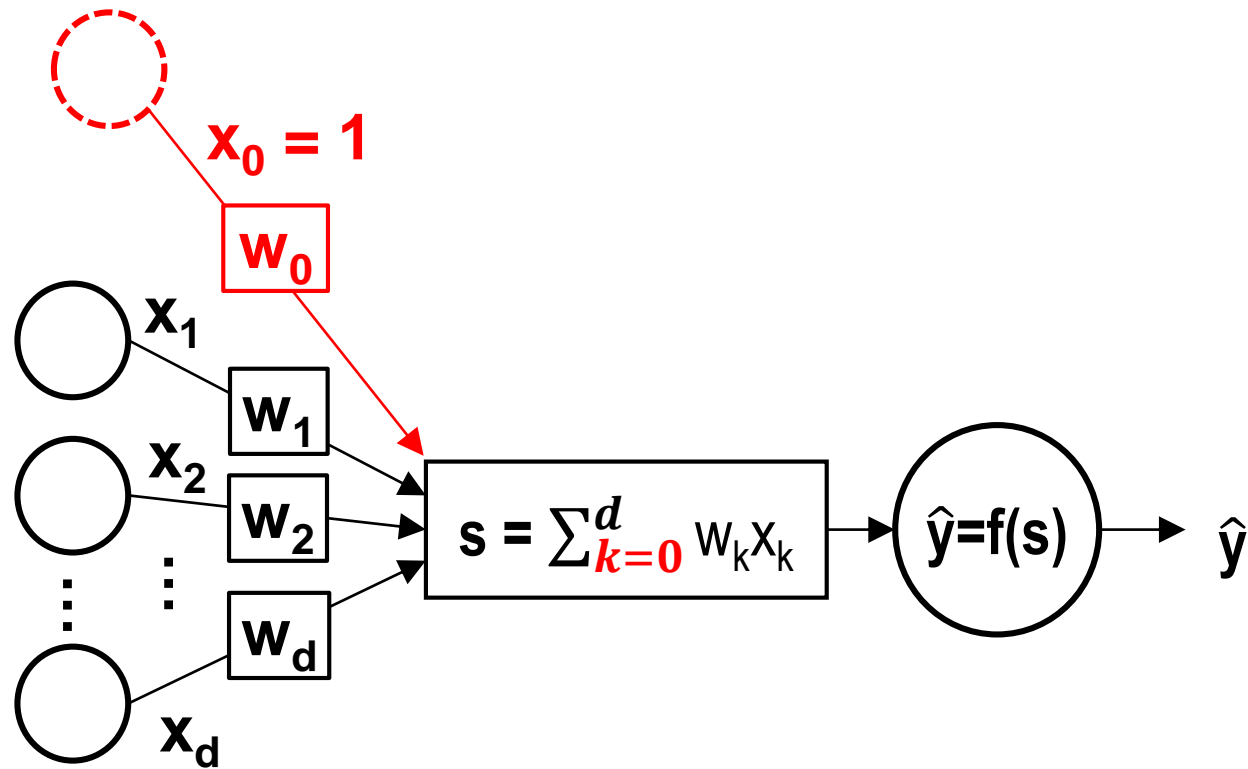
Widely used Activation Function
Sign function or $\text{sgn}(x)$



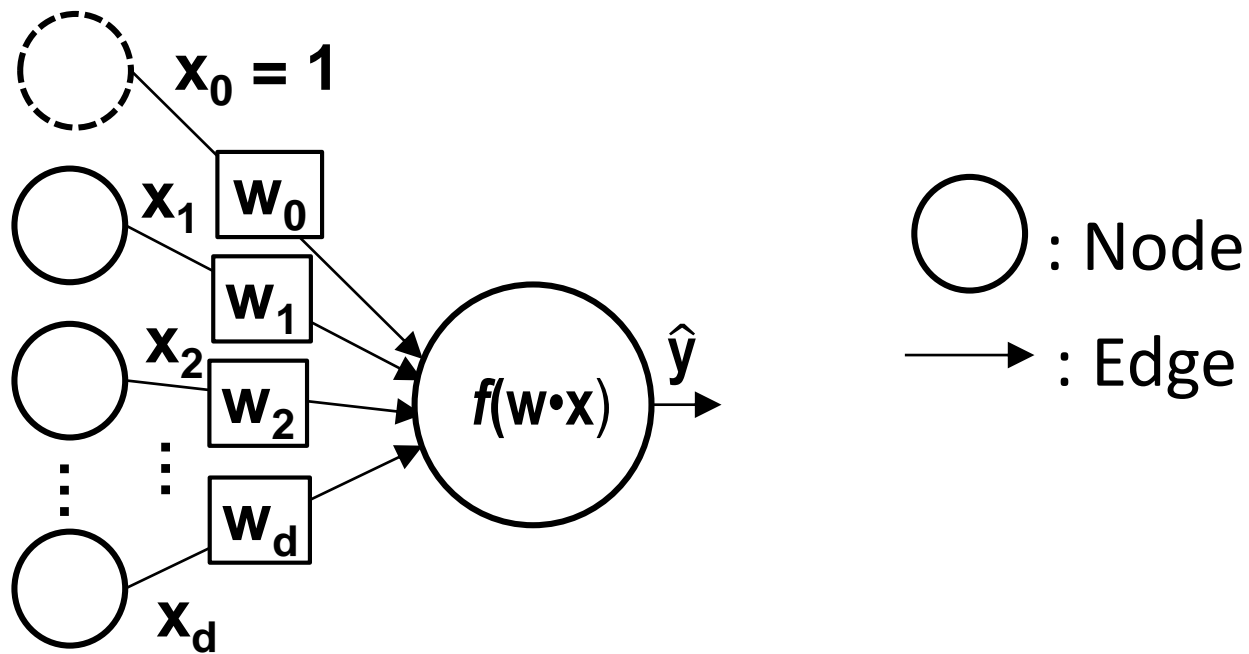
- We can represent $\sum_{k=1}^d w_k x_k$ with inner product of \mathbf{w} and \mathbf{x} .

$$\sum_{k=1}^d w_k x_k = \mathbf{w} \cdot \mathbf{x}$$

Adding Bias

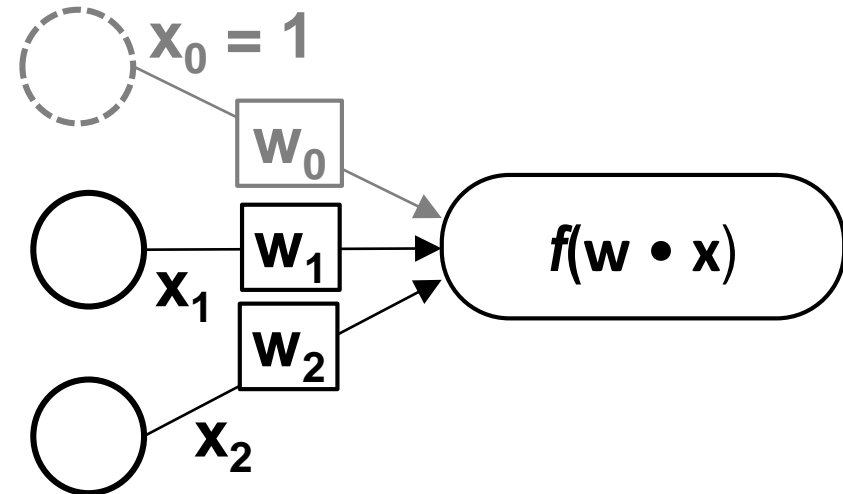
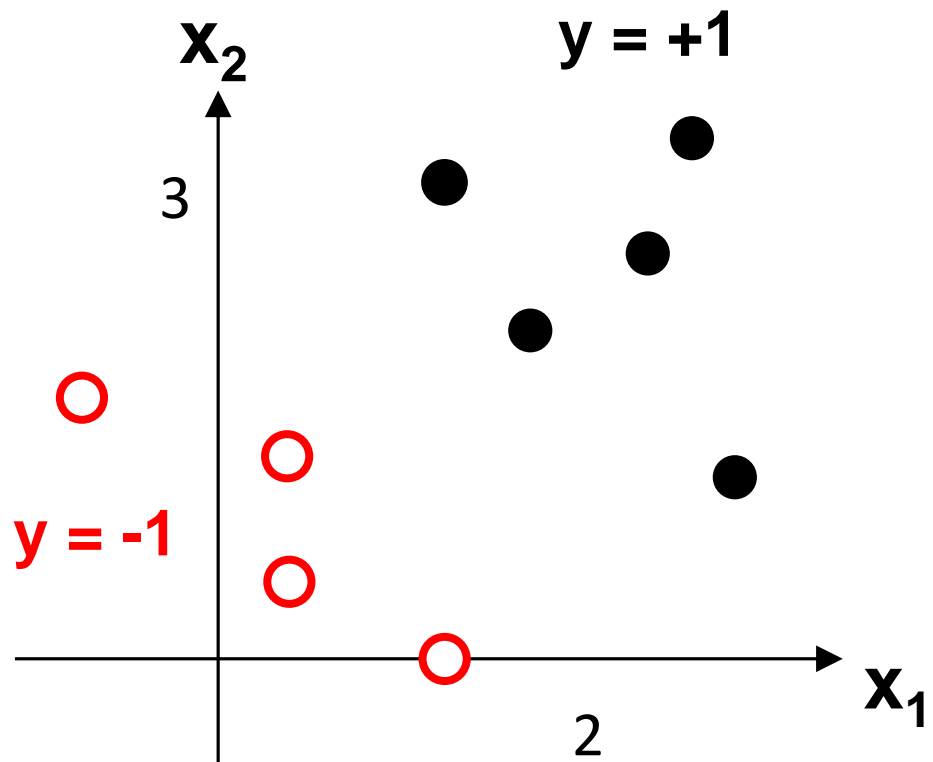


Structure of Perceptron



Simplest Example

- Can you apply perceptron to classification shown below?
→ How should you determine w_0, w_1, w_2 ?
- That is, to determine the parameter set or prediction model.



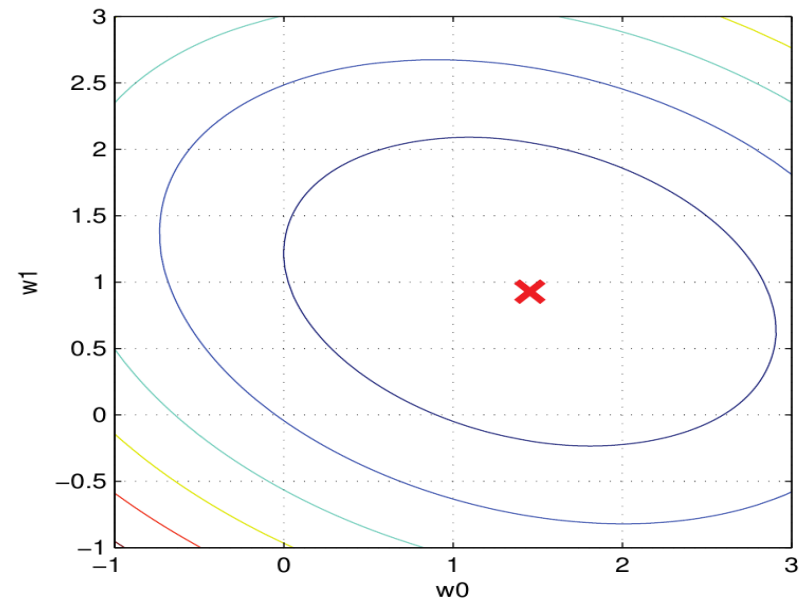
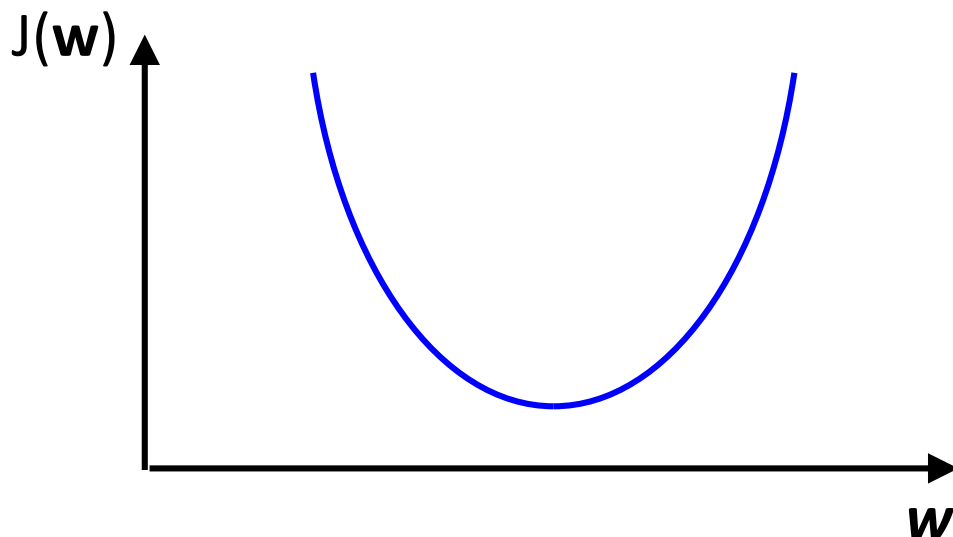
Perceptron as a Linear Classifier

- For 2D feature space, it is line
- How about for 3D feature space? How about larger than 3D?
- We need generalized cost function and minimization process.
→ How do we need to define the cost function?
- Cost function would be dependent to \mathbf{w} and our objective is to find out the \mathbf{w}^* that minimizes the cost.

Revisit Gradient Descent

- Given $J(\mathbf{w})$, $\frac{\partial J}{\partial \mathbf{w}} = 0$ or $\nabla_{\mathbf{w}} J = 0$ is the solution for minimum/maximum.
- However, it is highly complex or impossible to find solution analytically for many cases.
- Gradient descent method can find the above numerically by repeatedly doing:

$$\mathbf{w}_{\text{next}} = \mathbf{w}_{\text{present}} - \eta \nabla J(\mathbf{w}) \quad \leftarrow \eta : \text{learning rate}$$



Training for Perceptron

- Cost function is defined as

$$J(\mathbf{w}) = -\sum_{x_i \in A} y_i (\mathbf{w} \cdot \mathbf{x}_i)$$

where A is a set of samples that wrongly classified given \mathbf{w}

➔ Why do we use $\mathbf{w} \cdot \mathbf{x}$ instead of $f(\mathbf{w} \cdot \mathbf{x})$?

$$\begin{aligned} \frac{\partial J}{\partial w_k} &= -\frac{\partial}{\partial w_k} \sum_{x_i \in A} y_i (\mathbf{w} \cdot \mathbf{x}_i) = -\frac{\partial}{\partial w_k} \sum_{x_i \in A} y_i \left(\sum_{k=0}^d w_k x_{ik} \right) \\ &= -\sum_{x_i \in A} y_i x_{ik} \end{aligned}$$

- We can find k th component of \mathbf{w}^* by

$$w_{k,\text{next}} = w_{k,\text{present}} + \eta \sum_{x_i \in A} y_i x_{ik} \quad \leftarrow \eta : \text{learning rate}$$

- In vector representation,

$$\mathbf{w}_{\text{next}} = \mathbf{w}_{\text{present}} + \eta \sum_{x_i \in A} y_i \mathbf{x}_i$$

Training Algorithm

Initialize \mathbf{w} randomly

while ($A \neq \Phi$)

$A = \Phi$

 for $i = (1, \dots, N)$:

$$\hat{y}_i = f(\mathbf{w} \cdot \mathbf{x}_i)$$

 if ($\hat{y}_i \neq y_i$):

$A \cup \{\mathbf{x}_i\}$

 If ($A \neq \Phi$):

$$\mathbf{w}_{\text{next}} = \mathbf{w}_{\text{present}} + \eta \sum_{x_i \in A} y_i \mathbf{x}_i$$

$\mathbf{w}^* = \mathbf{w}$

← While ($A \neq \Phi$) can be inefficient

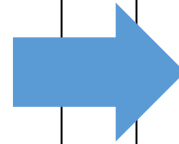
Soft decision can be implemented using
“No change in A”

Batch vs. Stochastic Training

- **Batch training** is training on all samples.
- **Stochastic training** is performing training on one randomly selected sample at a time.
- **Mini-batch training** is training on a part of the overall samples
- Batch training makes the algorithm too slow. Thus, stochastic or mini-batch training is used.

Applying Stochastic Training

```
Initialize  $\mathbf{w}$  randomly
while ( $A \neq \Phi$ )
   $A = \Phi$ 
  for  $i = (1, \dots, N)$ :
     $\hat{y}_i = f(\mathbf{w} \cdot \mathbf{x}_i)$ 
    if ( $\hat{y}_i \neq y_i$ ):
       $A \cup \{\mathbf{x}_i\}$ 
  If ( $A \neq \Phi$ ):
     $\mathbf{w}_{\text{next}} = \mathbf{w}_{\text{present}} + \eta \sum_{\mathbf{x}_i \in A} y_i \mathbf{x}_i$ 
 $\mathbf{w}^* = \mathbf{w}$ 
```



```
Initialize  $\mathbf{w}$  randomly
while ( $A \neq \Phi$ )
  Shuffle  $\mathbf{X}$ 
  for  $i = (1, \dots, N)$ :
     $\hat{y}_i = f(\mathbf{w} \cdot \mathbf{x}_i)$ 
    if ( $\hat{y}_i \neq y_i$ ):
       $\mathbf{w}_{\text{next}} = \mathbf{w}_{\text{present}} + \eta \sum_{\mathbf{x}_i \in A} y_i \mathbf{x}_i$ 
 $\mathbf{w}^* = \mathbf{w}$ 
```

Randomly select ONE sample by
1) Shuffle 2) then pick first error sample

Checkpoints

- ✓ Perceptron is formed mimicking human neural network.
- ✓ Classification for linearly separable data can be done by Linearly boundary (inner product) + Non-linear (sign func.)
- ✓ Gradient descent (batch/stochastic/mini-batch) for training $\theta(=w)$.
- ✓ Is there any way to apply the perceptron to linearly inseparable data?

