

## A. Perceptron

(Rosenblatt, 1957)

Firing of neuron is binary classification task. Use step function as activation function  $\phi(\tilde{z})$ , i.e.

$$\phi(\tilde{z}) = \begin{cases} 1 & \tilde{z} \geq 0 \\ -1 & \tilde{z} < 0 \end{cases}$$

where

$$\tilde{z} = \sum_{j=1}^m w_j x_j$$

 $x_1, \dots, x_m$  : input data $w_1, \dots, w_m$  : weightsand suppose there are  $N$  data samples ( $i = 1, \dots, N$ )

Define

$$\begin{cases} w_0 = -\theta \\ x_0 = 1 \end{cases} \quad (\text{bias})$$

then

$$z = \sum_{j=0}^m w_j x_j = -\theta + \tilde{z}$$

$$\phi(z) = \begin{cases} 1 & z \geq 0 \\ -1 & z < 0 \end{cases}$$

Perceptron algorithm :

- initialise weights randomly  $\epsilon (0, 1)$
- for each training sample  $(x_1^i, \dots, x_m^i)$  with true output  $y_i \in \{-1, 1\}$

1. compute :  $\hat{y}^i = \phi(z^i)$

2. update :  $\hat{w}_j = w_j + \Delta w_j$

$$\Delta w_j = \eta \left( \hat{y}^i - y^i \right) x_j \quad j=1, \dots, m$$

$\eta$  : learning rate

Note : all weights updated simultaneously for each  $i$

Cases : 1.  $\hat{y}^i - y^i = 0$  correct prediction of output

$$\rightarrow \Delta w_j = 0$$

2.  $\hat{y}^i - y^i \neq 0$  incorrect

a.  $\hat{y}^i = 1 - y^i = -1$

$$\rightarrow \Delta w_j = -2\eta x_j \rightarrow \text{note } w_-$$

b.  $\hat{y}^i = -1 \rightarrow y^i = 1$

$$\rightarrow \Delta w_j = 2\eta x_j \rightarrow \text{note } w_+$$

2b. assume  $\eta = 1$  &  $x_j^i = \frac{1}{2}$   
 [ sample misclassified with  
 $y_j^i = -1$  instead of 1.]

$$\rightarrow \Delta w_j = 1$$

pushed such that  $z$  is  
 larger, so that it is more  
 likely that  $z^i \geq 0$  and  
 $y_j^i = f(z^i) = +1$

Remarks :- the output is also  
 called a class label  

- classes should be  
 separable by a  
 linear decision  
 boundary.

## B. Adaptive linear neuron

In the perceptron, the weight were updated based on the step function. In the Adaline, the weight are updated using the linear activation function

$$\phi(z) = z$$

where  $z = \sum_{j=0}^m w_j x_j$

After this, a quantizer (e.g. a step function) can be used

to predict the class labels,

i.e.  $q(z) = \begin{cases} 1 & z \geq 0 \\ -1 & z < 0 \end{cases}$

where again the bias is taken into account through  $w_0$  and  $x_0$ .

To determine the weights, we define a cost function

$$J(\underline{w}) = \frac{1}{2} \sum_i \left( y^i - \phi(z^i) \right)^2$$

where the sum is over the samples.

The mathematical problem to solve is to determine  $\underline{w}_*$  as

$$\underline{w}_* = \min_{\underline{w}} J(\underline{w})$$

7

## Method 1. Gradient descend

iteration index  $k$ ,  $\underline{w}_0 = \underline{w}_{\text{initial}}$

then

$$\left[ \begin{array}{l} \underline{w}_{k+1} = \underline{w}_k + \Delta \underline{w}_{k+1} \\ \Delta \underline{w}_{k+1} = -\eta \nabla \mathcal{E}(\underline{w}_k) \end{array} \right]$$

↑ learning rate  
 ↓ down gradient

Now  $(\nabla \mathcal{E})_j$  can be computed from

$$\begin{aligned} (\nabla \mathcal{E})_j &= \frac{\partial}{\partial w_j} \left( \frac{1}{2} \sum_i (y^i - \phi(z^i))^2 \right) \\ &= \sum_i (y^i - \phi(z^i)) \frac{\partial}{\partial w_j} (y^i - \phi(z^i)) \\ \text{use } \phi(z^i) &= \sum_{j=0}^m x_j^i w_j \\ &= \sum_i (y^i - \phi(z^i)) (-x_j^i) \end{aligned}$$

Hence :

$$(\Delta \underline{w}_{k+1})_j = \gamma \equiv (y^i - \phi(z^i)) x_j^i$$

Remarks : - convergence can be slow when small values of  $\gamma$  are needed. In this case, it often helps to rescale the input

data (with  $x_j = (x_j^1, \dots, x_j^n)$ )

$$\underline{x}_j^s = \frac{\underline{x}_j - \mu_j}{\sigma_j}$$

where  $\mu_j = \frac{1}{N} \sum_{i=1}^N x_j^i$  (mean)

and  $\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_j^i - \mu_j)^2$

Method 2 : Stochastic gradient descend

Update of weights :

$$(\Delta w_{k+1})_j = \eta (y^i - \phi(z^i)) x_j^i$$

Advantages : - cheaper  
 - more easily to  
 hop over 'local'  
 minima of  $E$ .

Variants : - use varying  $\eta$   
 - shuffle training  
 data for every  
 iteration (epoch)