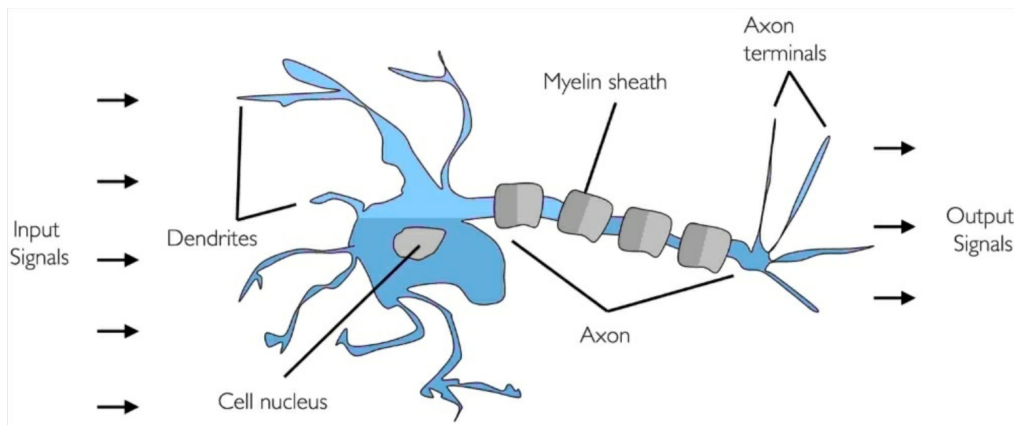
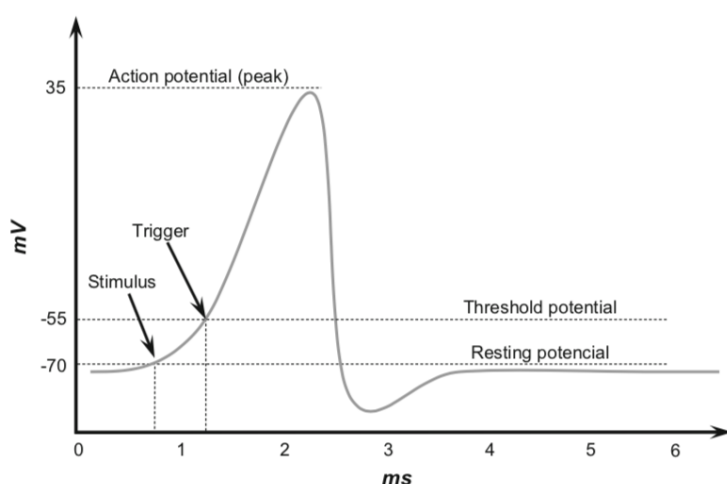


## Perceptron

The Perceptron network is one of the simplest ANN architectures, invented in 1957 by Frank Rosenblatt. The Perceptron network is inspired by the biological neuron shown below. Neurons are interconnected nerve cells in the brain that are involved in the processing and transmitting of chemical and electrical signals, which is illustrated in the following figure:

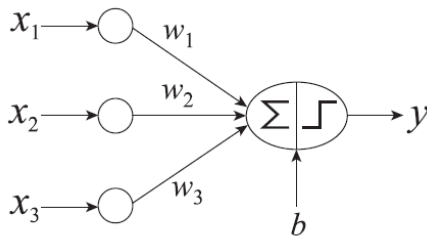


When the nervous cell is stimulated with an impulse higher than its activation threshold ( $-55$  mV), caused by the variation of internal concentrations of sodium ( $\text{Na}^+$ ) and potassium ( $\text{K}^+$ ) ions, it triggers an electrical impulse which will propagate throughout its axon with a maximum amplitude of  $35$  mV. The stages related to variations of the action voltage within a neuron during its excitation are shown in the Figure below.



## 2. Perceptron: The Basics

The Perceptron network consists of one artificial neuron. Figure below illustrates a Perceptron network composed of 3 input signals, representing the problem being analyzed, and just one output.



In mathematical notation, the inner processing performed by the Perceptron can be described by the following expressions:

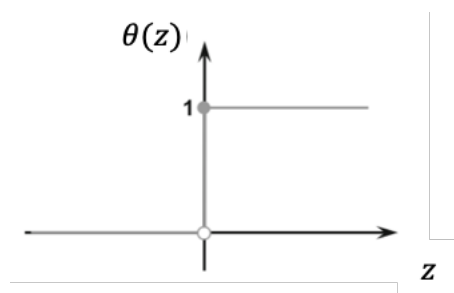
$$\begin{cases} z = w_1x_1 + w_2x_2 + w_3x_3 + b \\ y = \theta(z) \end{cases}$$

where  $x_i$  is a network input,  $w_i$  is the weight associated with the  $i^{th}$  input,  $b$  is the activation threshold (or bias),  $\theta(\cdot)$  is the activation function and  $z$  is the activation potential or sometimes called input signal.

The most common activation functions used in Perceptron are the step and bipolar step functions:

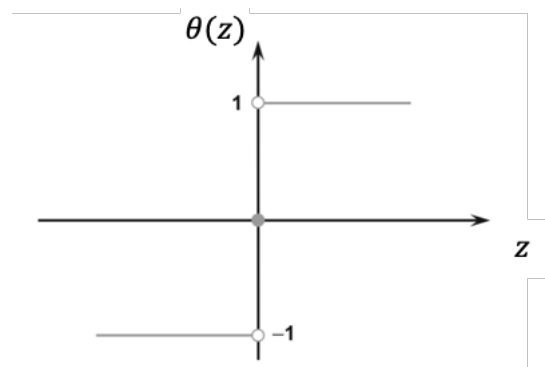
### Step Function

$$\theta(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$



Bipolar Step Function

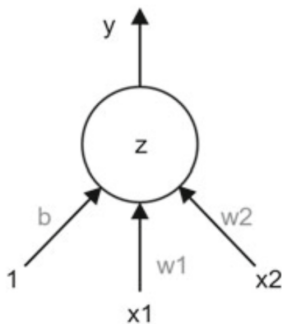
$$\theta(z) = \begin{cases} -1 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$



**Q:** What are the possible output values possibly produced by the Perceptron?

**A:** Only two values can be produced by the network output, that is, 0 or 1 for the step function, and  $-1$  or  $1$  if the bipolar step function is used.

**Remark on Bias Absorption:** We note that it is possible to absorb the bias as one of the weights, so that we only need **a weight update rule**. This is displayed in the figure below: to absorb the bias as a weight, one needs to add an input  $x_0$  with value 1 and **the bias is its weight**.



$$\begin{cases} z = b + \sum_{i=1}^d w_i x_i = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x} \\ y = \theta(z) \end{cases}$$

in which  $\mathbf{w} = [w_0, w_1, \dots, w_d]^T$  and  $\mathbf{x} = [1, x_1, \dots, x_d]^T$  are column vectors.

**Fun Time:** Can perceptron do (1) linear classification (2) linear regression (3) nonlinear classification (4) nonlinear regression?

The Perceptron network can be used for simple linear binary classification. It computes a linear combination of the inputs and if the result exceeds a threshold, it outputs the positive class or else outputs the negative class. Training a Perceptron network means finding the right values for the weights  $\mathbf{w} = [w_0, w_1, \dots, w_d]^T$ .

### 3. Training Process of the Perceptron

The adjustment of Perceptron's weights  $\mathbf{w}$ , in order to classify patterns that belong to one of the two possible classes, is performed by the use of Gradient Decent. For the bipolar step function, the Perceptron can be written in a vector form as:

$$\hat{y} = h(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x})$$

in which  $\mathbf{w} = [w_0, w_1, \dots, w_d]^T$  and  $\mathbf{x} = [1, x_1, \dots, x_d]^T$  are column vectors. Given a training set, we are interested in learning parameters  $\mathbf{w}$  such that  $\hat{y}$  closely resembles the true  $y$  of training instances. This is achieved by using the perceptron learning algorithm given below:

- 1: Let  $D = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\}$  be the set of training examples.
- 2: Initialize the weight vector with random values,  $\mathbf{w}^{(0)}$
- 3: **repeat**
- 4:   **for** each training example  $(\mathbf{x}_i, y_i) \in D$  **do**
- 5:     Compute the predicted output  $\hat{y}_i^{(k)}$
- 6:     **for** each weight  $w_j$  **do**
- 7:       Update the weight,  $w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_{ij}$ .
- 8:     **end for**
- 9:   **end for**
- 10: **until** stopping condition is met

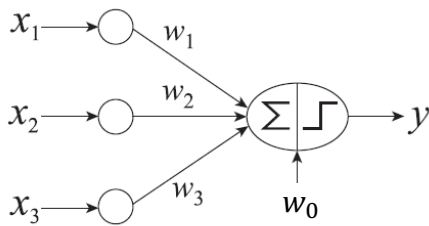
The key computation for this algorithm is the weight update formula given in Step 7 of the algorithm:

$$w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_{ij}$$

where  $w_j^{(k)}$  is the weight parameter associated with the  $j^{th}$  attribute after the  $k^{th}$  iteration,  $\lambda$  is a parameter known as **the learning rate**, and  $x_{ij}$  is the value of the  $j^{th}$  attribute of the training example  $x_i$ .

#### Example

Consider a perceptron with a bipolar step function:



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

The table below shows a dataset containing eight training examples with three Boolean variables  $(x_1, x_2, x_3)$  and an output variable,  $y$ . The output takes on the value  $-1$  if at least two of the three inputs are zero, and  $+1$  if at least two of the inputs are greater than zero.

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 1     | 0     | 0     | -1  |
| 1     | 0     | 1     | 1   |
| 1     | 1     | 0     | 1   |
| 1     | 1     | 1     | 1   |
| 0     | 0     | 1     | -1  |
| 0     | 1     | 0     | -1  |
| 0     | 1     | 1     | 1   |
| 0     | 0     | 0     | -1  |

Let us train the network by sequence with each training example. We will call it an epoch when we go through all the examples. Note that there can be multiple decision boundaries that can separate the two classes, and **the perceptron arbitrarily learns one of these boundaries depending on the random initial values of parameters.** (SVM takes care of the selection of the optimal decision boundary)

Let us start from  $\mathbf{w}^{\text{init}} = [w_0 = 0, w_1 = 0, w_2 = 0, w_3 = 0]^T$  with learning rate  $\lambda = 0.1$

Example 1, Epoch 1

**Fun Time:** What are the updated values for  $w_0$  and  $w_1$ ? (1)  $[0, 0]$  (2)  $[1, 1]$  (3)  $[-0.2, 0.2]$   
 (4)  $[-0.2, -0.2]$  (5)  $[-0.2, 0.0]$ ?