

# Lecture 12

# Machine Learning

## 12.1 Machine Learning

Machine Learning is a field of study focusing on how to make computer system learn from a set of examples. To be more precise, Machine Learning focuses on *automatically constructing models from the collected dataset without requiring explicit programming*.

Several types of machine learning techniques have been developed until now, e.g. supervised learning, unsupervised learning, reinforcement learning, etc. In this course, we focus only on *supervised machine learning*.

### 12.1.1 Supervised Machine Learning

Given a set of input-output pairs

$$\mathcal{D} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N$$

where,

$\mathcal{D}$  is a training set,

$N$  is the number of examples,

$\mathbf{x}_i \in \mathcal{X}$  is an input vector, and

$t_i \in \mathcal{T}$  is a target output,

find a function  $\hat{f} : \mathcal{X} \rightarrow \mathcal{T}$  that predicts the value of  $t$  from an unknown  $\mathbf{x}$ .

Each input vector  $\mathbf{x}_i$  is usually represented as a  $D$ -dimensional feature vector, i.e.

$$\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{iD}]^T \in \mathbb{R}^D$$

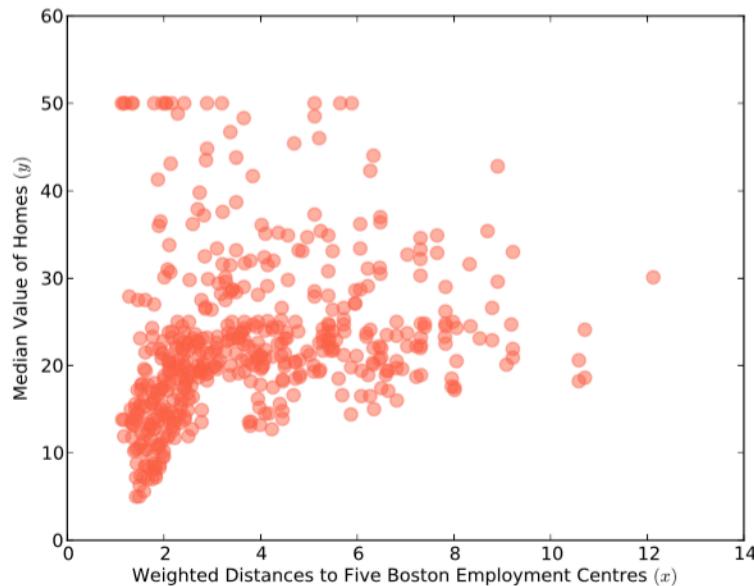
It can be represented as a point in the  $D$ -dimensional space.

It is assumed that each pair  $(\mathbf{x}_i, t_i)$  is drawn i.i.d. (independently and identically distributed) from an unknown distribution on  $\mathcal{X}$ .

### 12.1.2 Classification and Regression

When  $\mathcal{T}$  is a finite set or  $t_i$  is categorical, the task of the supervised learning is called *classification*. The task is called *regression* when  $t_i$  is real-valued i.e.  $\mathcal{T} = \mathbb{R}$  or  $\mathcal{T} = \mathbb{R}^D$ . In this course, we mainly focus on the classification task where the *binary classification* is the most frequently studied problem. It is the classification task when  $\mathcal{T} = \{-1, +1\}$ .

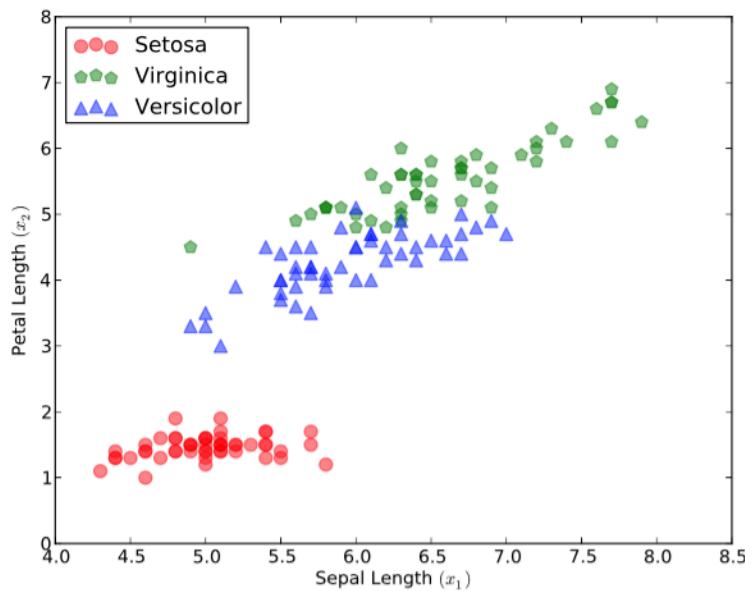
**Example 12.1** Predict the housing values in suburbs of Boston.<sup>1</sup>



**Example 12.2** Predict the class of iris using its shape.<sup>2</sup>

<sup>1</sup>Housing Data Set – <http://archive.ics.uci.edu/ml/datasets/Housing>

<sup>2</sup>Iris Data Set – <http://archive.ics.uci.edu/ml/datasets/Iris>



**Example 12.3** Recognize handwritten digits.<sup>3</sup>

3 6 8 1 7 9 6 6 9 1  
 6 7 5 7 8 6 3 4 8 5  
 2 1 7 9 7 1 2 8 4 6  
 4 8 1 9 0 1 8 8 9 4  
 7 6 1 8 6 4 1 5 6 0  
 7 5 9 2 6 5 8 1 9 7  
 2 2 2 2 2 3 4 4 8 0  
 0 2 3 8 0 7 3 8 5 7  
 0 1 4 6 4 6 0 2 4 3  
 7 1 2 8 7 6 9 8 6 1

Each handwritten digit is preprocessed and transformed into a  $28 \times 28$  greyscale image. Therefore, each input vector composes of 784 features.

---

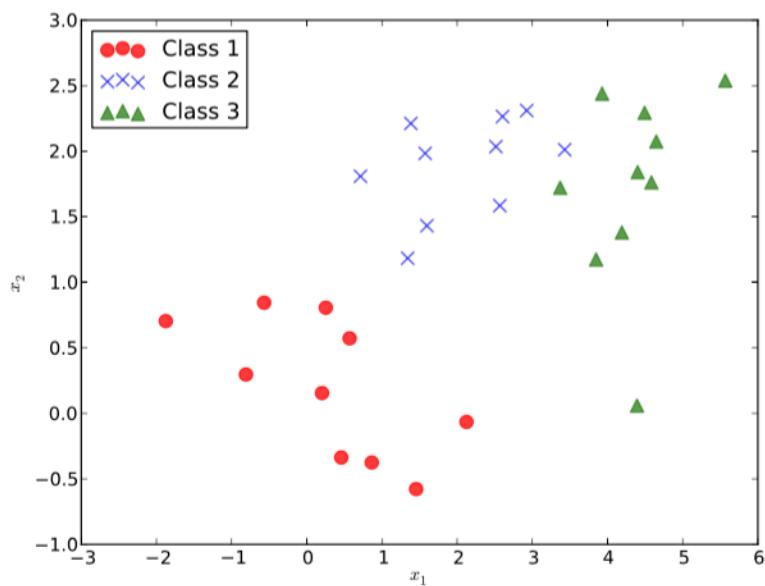
<sup>3</sup>The MNIST Database of Handwritten Digits (<http://yann.lecun.com/exdb/mnist/>); Y. LeCun, L. Bottou, Y. Bengio and P. Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998

**Example 12.4** Classify a set of artificially generated examples

$x_1$	$x_2$	$t$
2.124	-0.065	1
0.253	0.807	1
1.454	-0.578	1
0.569	0.573	1
0.458	-0.337	1
-0.809	0.297	1
0.864	-0.375	1
0.202	0.155	1
-1.875	0.705	1
-0.569	0.845	1

$x_1$	$x_2$	$t$
1.342	1.182	2
2.568	1.583	2
2.515	2.034	2
1.384	2.211	2
2.926	2.310	2
0.714	1.808	2
3.430	2.011	2
1.575	1.983	2
1.597	1.430	2
2.604	2.264	2

$x_1$	$x_2$	$t$
4.393	0.059	3
4.584	1.761	3
3.370	1.720	3
4.192	1.379	3
4.649	2.073	3
3.849	1.173	3
4.401	1.839	3
4.494	2.292	3
5.567	2.538	3
3.928	2.438	3



### 12.1.3 Probabilistic Prediction

Using the probability theory, we can find a prediction function from

$$y = \hat{f}(\mathbf{x}) = \operatorname{argmax}_{c=1}^C P(t = c | \mathbf{x}, \mathcal{D})$$

This is known as a **MAP estimate**.<sup>4</sup>

## 12.2 K-Nearest Neighbors

A simple technique to predict the class of an unknown feature vector  $\mathbf{x}$  is to check the classes of the  $K$  nearest vectors in the training set. Then, predicting the class of  $\mathbf{x}$  using the majority class.

The distance between two vectors can be simply calculated by the Euclidean norm of the difference of the two vectors i.e.

$$\begin{aligned} \text{dist}(\mathbf{x}_i, \mathbf{x}_j) &= \|\mathbf{x}_i - \mathbf{x}_j\| \\ &= \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \cdots + (x_{iD} - x_{jD})^2} \end{aligned}$$

Formally, we can compute the prediction probability from

$$P(t = c | \mathbf{x}, \mathcal{D}, K) = \frac{1}{K} \sum_{i \in N_K(\mathbf{x}, \mathcal{D})} \mathbb{I}(t_i = c)$$

where  $N_K(\mathbf{x}, \mathcal{D})$  is the set of  $K$  nearest points to  $\mathbf{x}$  in  $\mathcal{D}$ , and  $\mathbb{I}(e)$  is the indicator function defined as

$$\mathbb{I}(e) = \begin{cases} 1 & \text{if } e \text{ is true} \\ 0 & \text{if } e \text{ is false} \end{cases}$$

---

<sup>4</sup>MAP stands for *maximum a posteriori*

**Exercise 12.1** From the following dataset, use the  $K$ -nearest neighbors technique to predict the classes of the following examples with different values of  $K$ .

$x_1$	$x_2$	$t$	$x_1$	$x_2$	$t$
0.5	0.5	0	2.0	1.5	1
0.2	1.0	0	3.0	1.0	1
1.0	0.8	0	4.0	3.2	1

1.  $[0.0, 0.0]^\top$  using  $K = 3$

2.  $[1.0, 1.0]^\top$  using  $K = 5$

## 12.3 Model Selection

Since we can run KNN using different value of  $K$ , which value of  $K$  should we use?

A basic way to evaluate a classifier is to use the *misclassification rate* on the training set. It can be defined as:

$$\text{err}(f, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(f(\mathbf{x}_i) \neq t_i)$$

However, the error on the training set should not represent the *generalization error*.

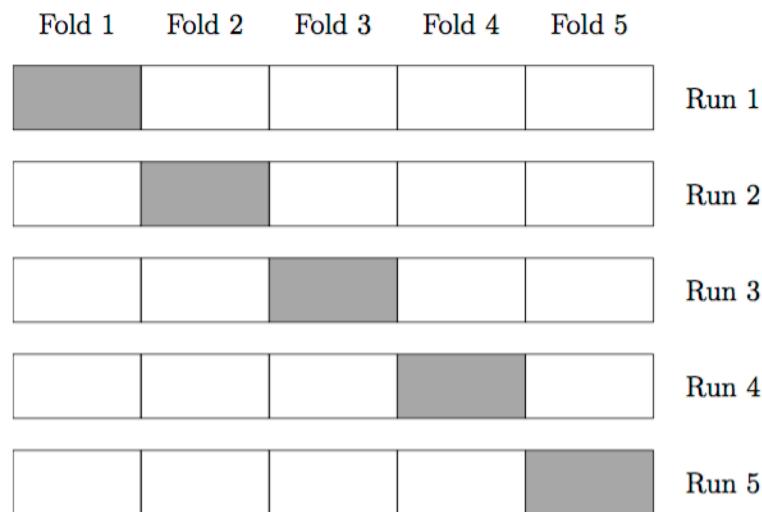
### 12.3.1 Holdout Method

The collected dataset is split into two subsets, i.e. *training* and *validation* sets. Generally, about 80% of the data are for the training set, and 20% of the data are for the validation set.

However, the evaluation result depends heavily on the examples in the training set and the validation set. We may get too good or too bad result easily.

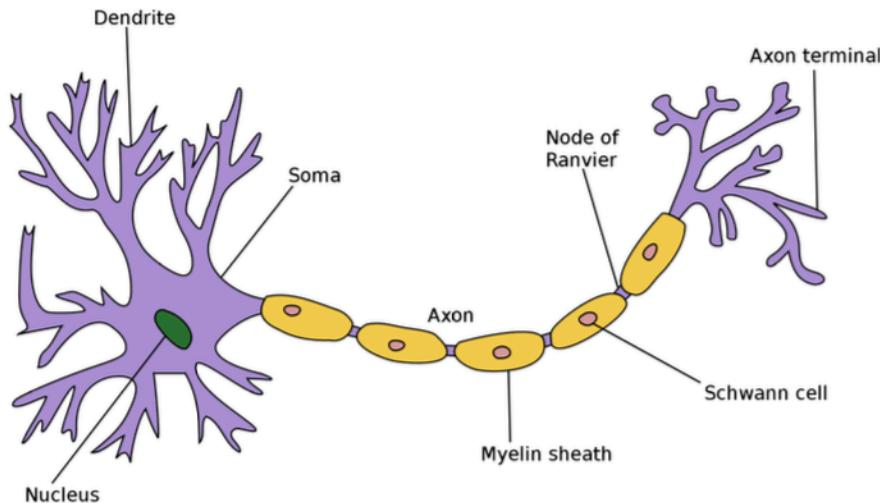
### 12.3.2 Cross Validation Method

The cross validation method splits the dataset into  $k$  folds. Then, it runs the holdout evaluation  $k$  runs. In each run, the  $k$ th fold is used as the validation set, while the combination of the rest folds are used as the training set in that run.



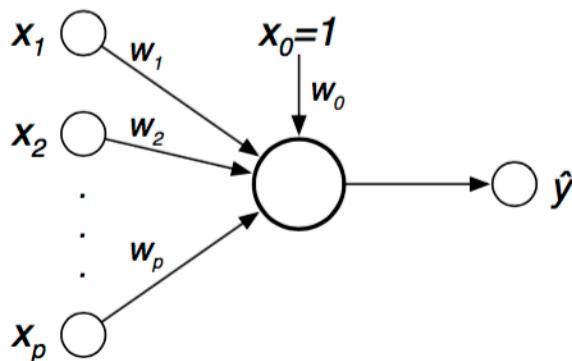
## 12.4 Artificial Neural Networks

**Artificial Neural Network (ANN)** is a supervised learning technique inspired by neurons in the nervous system.



Source: <http://en.wikipedia.org/wiki/Neuron>

**Perceptron** is a type of ANN working as a kind of binary classification based on a linear model.



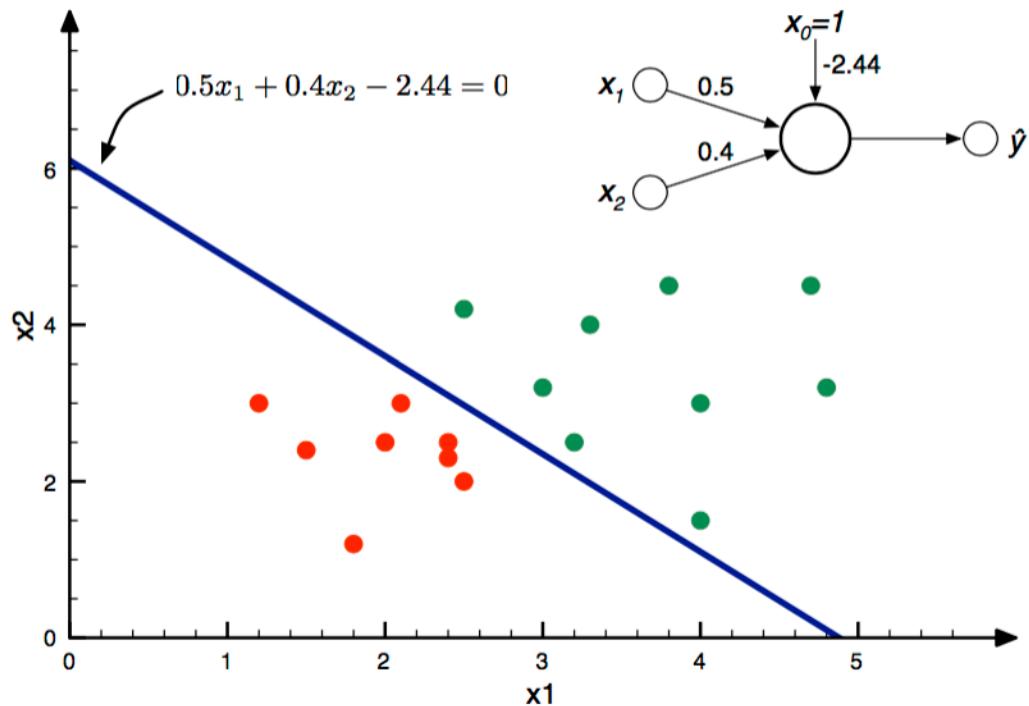
Each perceptron has several inputs expressing instance features denoting  $(x_0, x_1, x_2, \dots, x_p)^\top$ . A weight is attached to each input as a vector  $(w_0, w_1, w_2, \dots, w_p)^\top$ . These weights are the important part denoting the main characteristic of each perceptron. An output of perceptron denotes the prediction result.

A perceptron works as a **linear classifier** getting features as inputs and estimating the class. The prediction is done by

$$\begin{aligned}
 y &= \operatorname{sgn}\left(\sum_{i=0}^p w_i x_i\right) \\
 &= \operatorname{sgn}(w_0 + w_1 x_1 + \cdots + w_p x_p)
 \end{aligned}$$

where,  $\operatorname{sgn}(x) = \begin{cases} 1 & \text{if } x > 0, \\ 0 & \text{otherwise} \end{cases}$

This classifier deals with a two-class problem (1 and 0). We have  $w_0 + w_1 x_1 + \cdots + w_p x_p = 0$  as a linear hyperplane.



**Exercise 12.2** Use the perceptron in the previous figure to predict the following examples.

$$1. [0.0, 0.0]^T$$

$$2. [4.0, 4.0]^T$$

$$3. [3.0, 2.0]^T$$

**Exercise 12.3** Answer the questions from the following dataset.

$x_1$	$x_2$	t
0	0	0
0	1	0
1	0	0
1	1	1

1. Draw the structure of a perceptron for learning from this dataset.
2. Find the weights that makes the perceptron correctly classify the training examples.

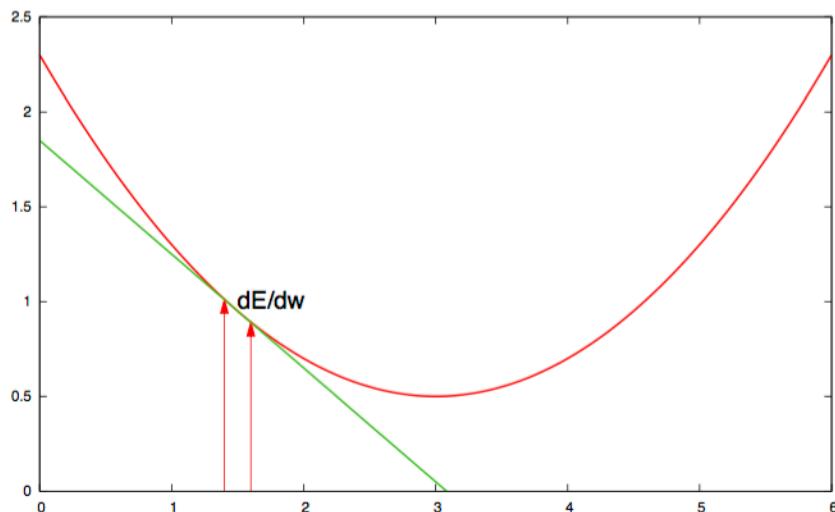
### 12.4.1 Perceptron Training Algorithm

Weights are the main component of each perceptron. If each weight is set to a suitable value, the perceptron will be able to classify unlabeled instances. We need a technique to assign or adjust weights until we obtain the most suitable values.

### Obtaining the Weight Updating Rule

*Gradient Descent* technique can be used to find the most suitable weights. It is a simple technique to solve optimization problem on a continuous environment.

1. Define an objective function. In this case, it should be a function of  $\mathbf{w}$  i.e.  $E(\mathbf{w})$ .
2. Randomly select initial values for  $\mathbf{w}$ .
3. Update weights using  $\mathbf{w} - \eta \frac{dE}{d\mathbf{w}}$  where  $\eta$  is a positive small number.
4. Iteratively update until  $\frac{dE}{d\mathbf{w}} \approx 0$



We can define an objective function using weights  $\mathbf{w}$  as a parameter. The function should return 0 if the weights are adjusted to the best values. In the simplest way, we can define an objective function as

$$\begin{aligned} E(\mathbf{w}) &\stackrel{d}{=} \frac{1}{2} \sum_{d=1}^N (t_d - y_d)^2 \\ &= \frac{1}{2} \sum_{d=1}^N \left( t_d - \sum_{i=0}^p w_i x_{di} \right) \end{aligned}$$

We compute gradient of the objective function which is denoted as

$$\nabla E(\mathbf{w}) \stackrel{d}{=} \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_p} \right]$$

Then, we find the partial differentiate of  $E$  by  $w_i$ :

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \left( \frac{1}{2} \sum_{d=1}^N (t_d - y_d)^2 \right) \\ &= \frac{1}{2} \sum_{d=1}^N \frac{\partial}{\partial w_i} ((t_d - y_d)^2) \\ &= \frac{1}{2} \sum_{d=1}^N 2(t_d - y_d) \frac{\partial}{\partial w_i} (t_d - y_d) \\ &= \sum_{d=1}^N (t_d - y_d) \frac{\partial}{\partial w_i} \left( t_d - \sum_{i=0}^p w_i x_{di} \right) \\ &= \sum_{d=1}^N (t_d - y_d)(-x_{di}) \end{aligned}$$

From the gradient of  $E$ , we have

$$w_i^* = w_i - \eta \nabla E(w_i)$$

We then update  $w_i$  by

$$w_i \leftarrow w_i + \eta(t - y)(x_i)$$

A training algorithm for perceptron is as below:

1. Start with random weights.
2. Apply each training instance. If the perceptron misclassifies, then adjust all the weights. Each weight  $w_i$  for an input  $x_i$  is adjusted by

$$w_i \leftarrow w_i + \eta(t - y)x_i$$

where  $\eta$  is a positive constant called the *learning rate*.

3. Iteratively, perform step 2 through the training set until the perceptron classifies all training instances correctly.

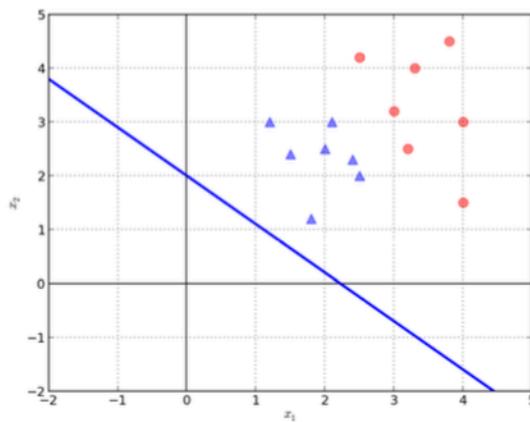
**Example 12.5** Find the most suitable weights from the following examples:

No	$x_1$	$x_2$	$t$	No	$x_1$	$x_2$	$t$
1	2.5	2.0	0	8	4.0	3.0	1
2	1.2	3.0	0	9	3.8	4.5	1
3	2.1	3.0	0	10	3.2	2.5	1
4	2.4	2.3	0	11	3.3	4.0	1
5	2.0	2.5	0	12	2.5	4.2	1
6	1.5	2.4	0	13	4.0	1.5	1
7	1.8	1.2	0	14	3.0	3.2	1

We start from random weights and the learning rate:

$$w_0 = -0.40; \quad w_1 = 0.18; \quad w_2 = 0.20; \quad \eta = 0.01$$

From this setting, we have a linear classifier as the following figure.

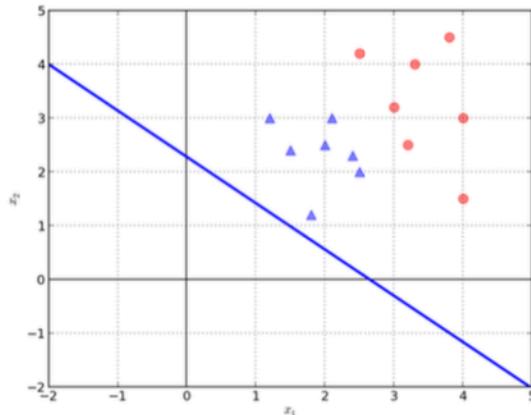


We start from applying an instance  $([2.5, 2.0]^\top, 0)$  to the perceptron, then

$$\begin{aligned} y &= \text{sgn}(0.18 \times 2.5 + 0.2 \times 2.0 - 0.4) \\ &= \text{sgn}(0.45) \\ &= 1 \end{aligned}$$

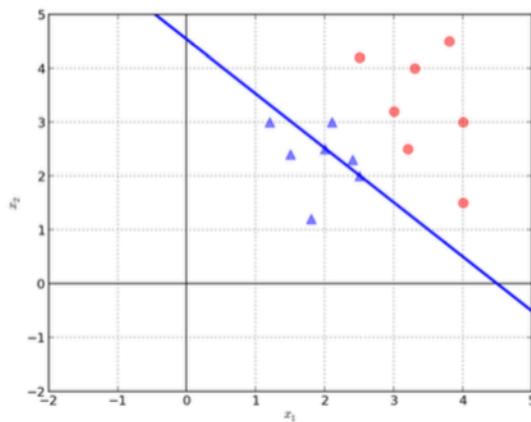
Since  $y \neq t$ , then we adjust weight as follow:

$$\begin{aligned} w_0 &\leftarrow -0.40 + ((0.01)(0 - 1)(1.0)) = -0.410 \\ w_1 &\leftarrow +0.18 + ((0.01)(0 - 1)(2.5)) = +0.155 \\ w_2 &\leftarrow +0.20 + ((0.01)(0 - 1)(2.0)) = +0.180 \end{aligned}$$



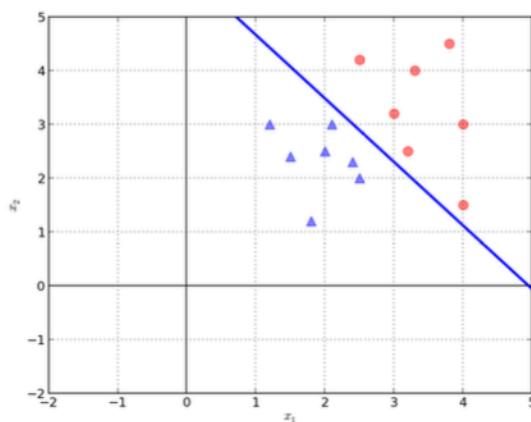
After applying all the remaining 13 examples, we will get weights as:

$$w_0 = -0.440; \quad w_1 = 0.098; \quad w_2 = 0.097$$



After updating the weights for five rounds, the weights become:

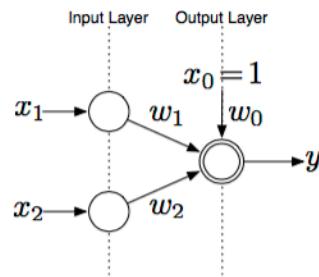
$$w_0 = -0.450; \quad w_1 = 0.091; \quad w_2 = 0.077$$



### 12.4.2 Threshold Function

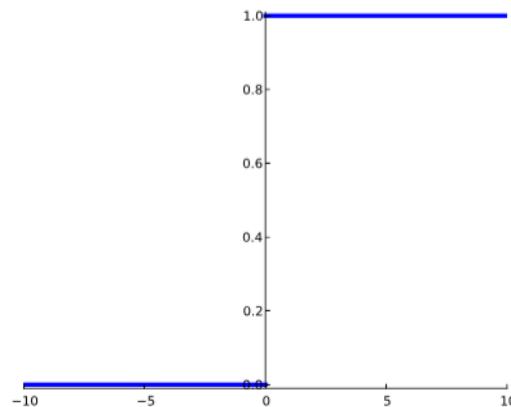
A perceptron outputs either ‘0’ or ‘1’ based on the sum of the inputs and the weights, as well as the threshold function (e.g.  $\text{sgn}(\cdot)$ ).

$$y = \text{sgn}(w_0 + w_1x_1 + w_2x_2)$$



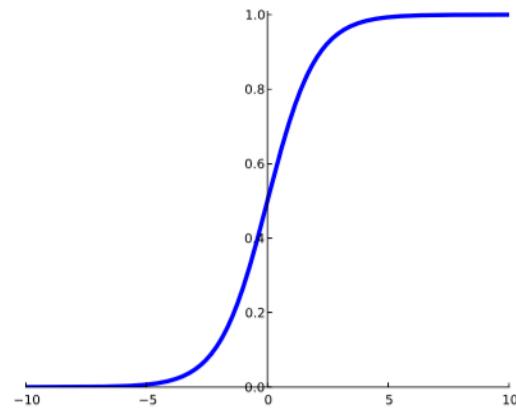
The threshold function  $\text{sgn}(\cdot)$  is a discontinuous function. It cannot be differentiated. Therefore, it cannot be used together with the gradient descent, or other optimization techniques.

$$\text{sgn}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$



The sigmoid function is an S-shaped function. It can be in place of the function  $\text{sgn}(\cdot)$ .

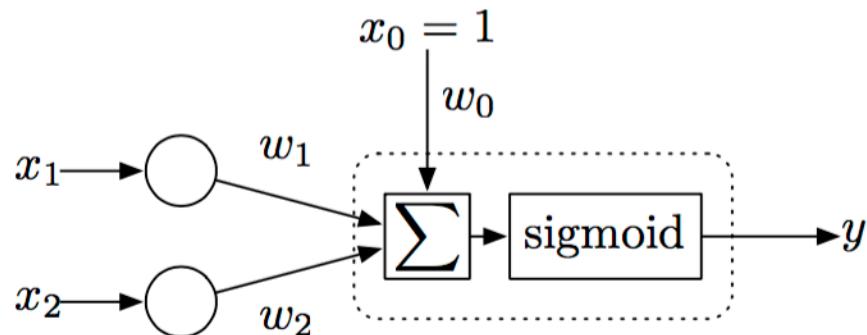
$$\text{sigmoid}(x) = \frac{1}{1 + e^{-\beta x}}$$



### 12.4.3 Sigmoid Unit

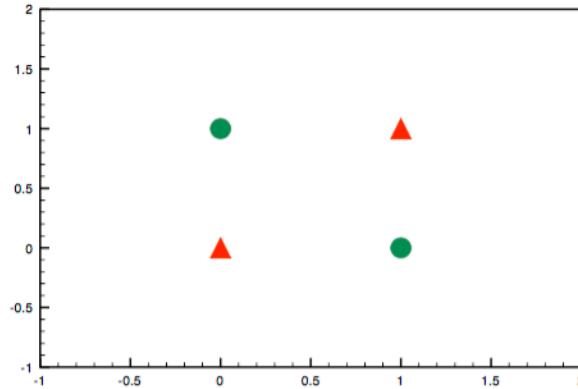
Now, an output of a perceptron can be calculated from:

$$y = \text{sigmoid}(w_0 + w_1x_1 + w_2x_2)$$



## 12.5 Linear Separability and Multilayer Perceptron

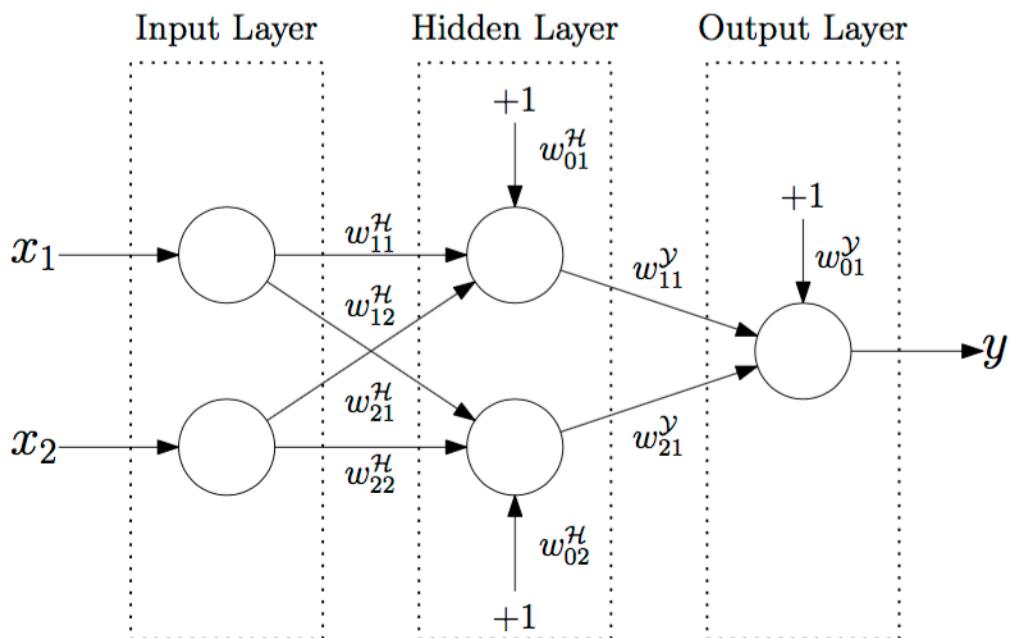
Since a perceptron represents a linear model for classification, it works very well on a training set that is **linearly separable**. However, a perceptron cannot deal with the training set that cannot be separated by just a line.



How can the perceptron be improved to handle this situation?

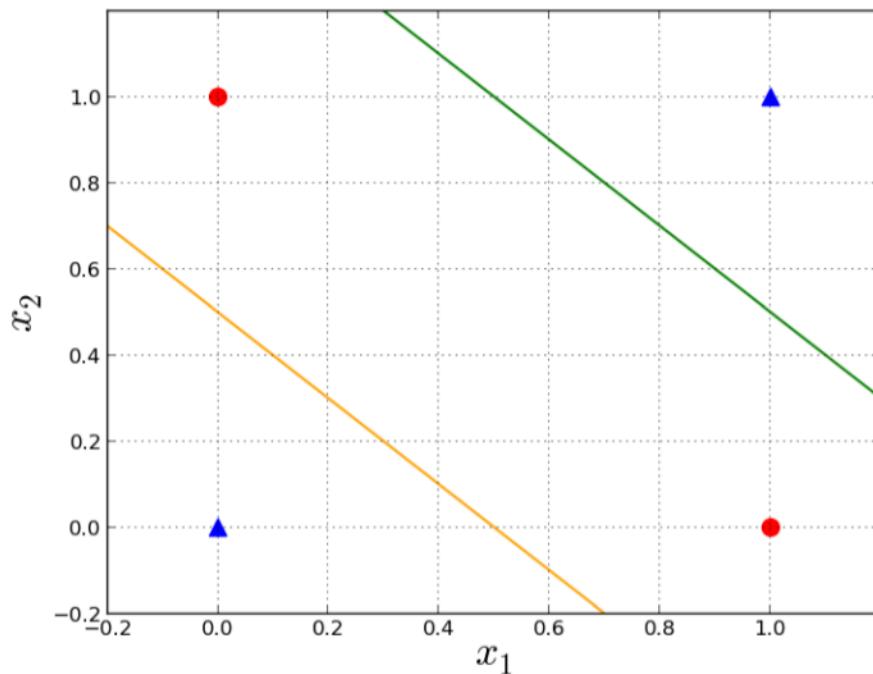
To solve the non-linearly separable problem, we connect multiple sigmoid units to become a neural network.

The units are arranged into more than two layers. The technique is called *multilayer perceptron*.

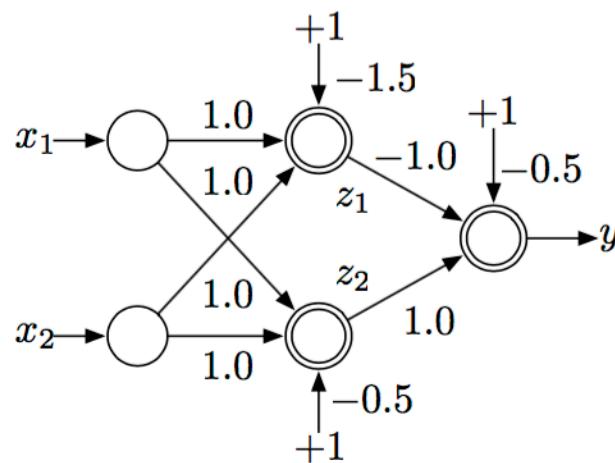


Each perceptron works as a linear discriminant. The inputs are classified by the linear discriminants in the hidden layer. Either 0 or 1 are yielded from each unit in the hidden layer. The outputs from the hidden layer are fed as inputs to the output layer. Here, the perceptron in the output is a linear discriminate that classifies the pre-classified inputs.

**Example 12.6** We uses two lines to correctly discriminate the examples:



These two lines can be represented as the following network:

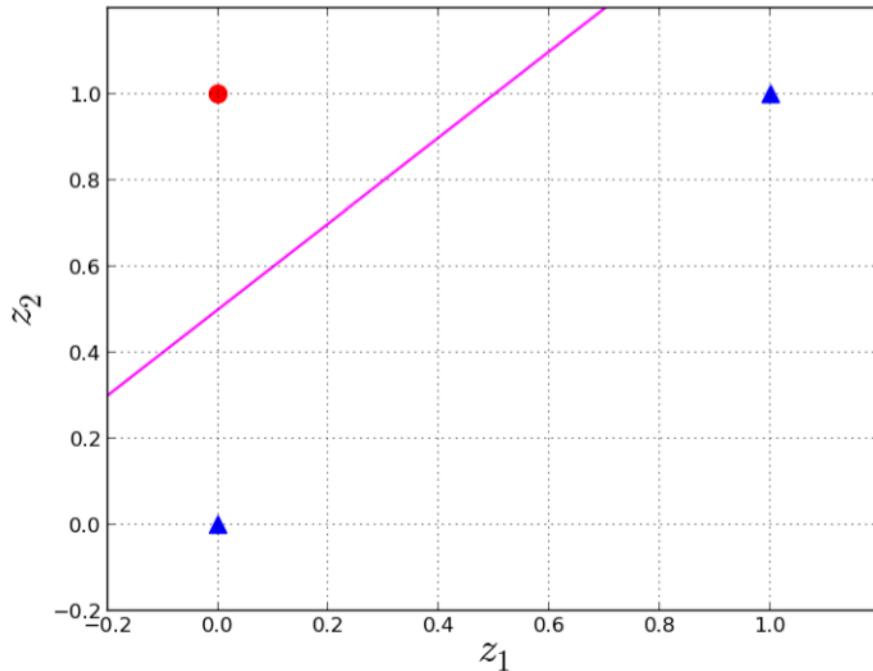


$$z_1 = \text{sigmoid}(x_1 + x_2 - 0.5)$$
$$z_2 = \text{sigmoid}(x_1 + x_2 - 1.5)$$

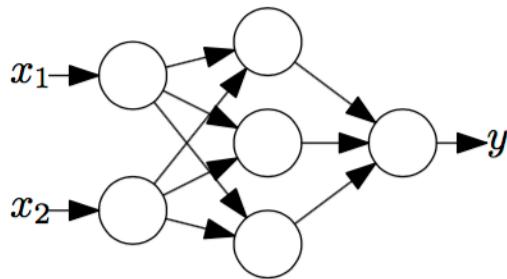
From the network, we have two intermediate outputs (i.e.  $z_1$ , and  $z_2$ ) from the perceptrons in the hidden layer:

$x_1$	$x_2$	$z_1$	$z_2$	$t$
0	0	0	0	0
0	1	0	1	1
1	0	0	1	1
1	1	1	1	0

We can use two lines to correctly discriminate the examples:



**Exercise 12.4** Use the multilayer perceptron with the following weights to predict the output of  $[1.0, 1.0]^\top$ . Set  $\beta = 1$ .



$$\begin{aligned} w_{01}^{\mathcal{H}} &= +0.50, & w_{11}^{\mathcal{H}} &= -0.50, & w_{21}^{\mathcal{H}} &= +0.10, \\ w_{02}^{\mathcal{H}} &= +1.00, & w_{12}^{\mathcal{H}} &= +0.20, & w_{22}^{\mathcal{H}} &= -0.20, \\ w_{03}^{\mathcal{H}} &= -1.00, & w_{13}^{\mathcal{H}} &= +0.10, & w_{23}^{\mathcal{H}} &= +0.50, \\ w_{01}^{\mathcal{Y}} &= -1.00, & w_{11}^{\mathcal{Y}} &= +1.50, & w_{21}^{\mathcal{Y}} &= -1.00, & w_{31}^{\mathcal{Y}} &= 1.00 \end{aligned}$$

### 12.5.1 Backpropagation Algorithm

The gradient descent technique is used to train a multilayer perceptron in the same manner as training a perceptron.

We first define the error function:

$$E(\mathbf{w}^{\mathcal{H}}, \mathbf{w}^{\mathcal{Y}}) = \frac{1}{2} \sum_{d=1}^N (t_d - y_d)^2$$

We can use the chain rule to calculate the gradient

$$\frac{\partial E}{\partial w_{jl}^{\mathcal{H}}} = \frac{\partial E}{\partial y} \frac{\partial y}{\partial z_l} \frac{\partial z_l}{\partial w_{jl}^{\mathcal{H}}}$$

Here, the error propagates from the output back to the inputs. Thus, the algorithm is named *backpropagation*.

For each example in  $D$ :

#### **Forward phase:**

1. compute the output of each perceptron in the hidden layer:

$$z_l = \text{sigmoid}\left(\sum_i w_{il}^{\mathcal{H}} x_i\right)$$

2. compute the output of each perceptron in the output layer:

$$y_k = \text{sigmoid}\left(\sum_l w_{lk}^{\mathcal{Y}} z_l\right)$$

#### **Backward phase:**

1. compute the error at the output using:

$$\delta_k^{\mathcal{Y}} = (t_k - y_k)y_k(1 - y_k)$$

2. update the output layer weights using:

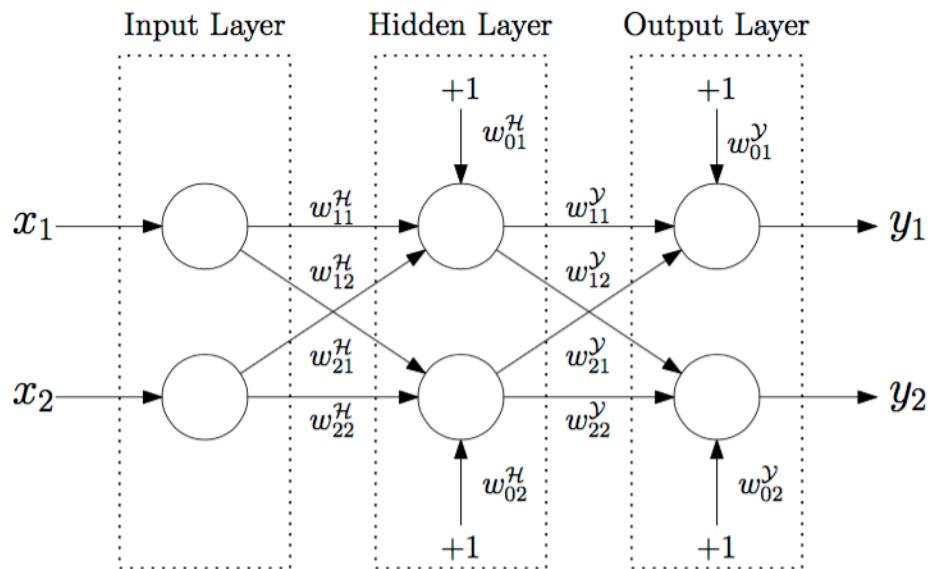
$$w_{lk}^{\mathcal{Y}} \leftarrow w_{lk}^{\mathcal{Y}} + \eta \delta_k^{\mathcal{Y}} z_l$$

3. compute the error in the hidden layer using:

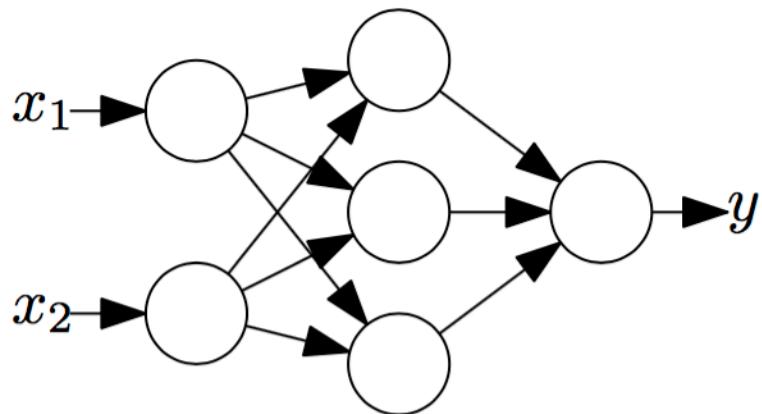
$$\delta_l^{\mathcal{H}} = z_l(1 - z_l) \sum_k w_{lk}^{\mathcal{Y}} \delta_k^{\mathcal{Y}}$$

4. update the hidden layer weights using:

$$w_{il}^{\mathcal{H}} \leftarrow w_{il}^{\mathcal{H}} + \eta \delta_l^{\mathcal{H}} x_i$$



**Exercise 12.5** Given the following multilayer perceptron with all initial weights of 0's. Find the all updated values of the weights of all layers when we feed a labeled example  $([1.0, 0.5]^\top, 1)$ . Set  $\beta = 1$  and  $\eta = 0.01$ .



## References

- Satish Kumar, “Neural networks: a classroom approach”, McGraw-Hill, 2005.
- Stephen Marsland, “Machine Learning: an algorithm perspective”, CRC Press, 2009.
- Ethem Alpaydin, “Introduction to Machine Learning”, The MIT Press, 2004.
- Kevin P. Murphy, “Machine Learning: A Probabilistic Perspective”, The MIS Press, 2012.