# Lecture 12

# Machine Learning

*(handwritten)* k nearest neighbors
neural networks

## 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.*

*(handwritten)* ( learn from training examples)

### 12.1.1 Supervised Machine Learning

Given a set of input-output pairs

*(handwritten)* an input vector

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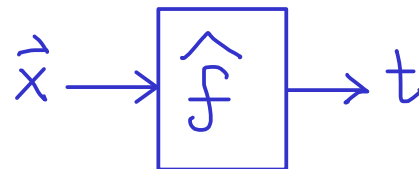
where,

*(handwritten)* a target output

$\mathcal{D}$ 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 $\hat{f} \colon \mathcal{X} \to \mathcal{T}$ that predicts the value of $t$ from an unknown $\mathbf{x}$.

*(handwritten)* $\vec{x} \to \boxed{\hat{f}} \to t$

*(handwritten)* this function is constructed automatically by an ML algorithm.

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

*(handwritten: ← dimension of input vectors)*

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

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 $\mathcal{X}$.
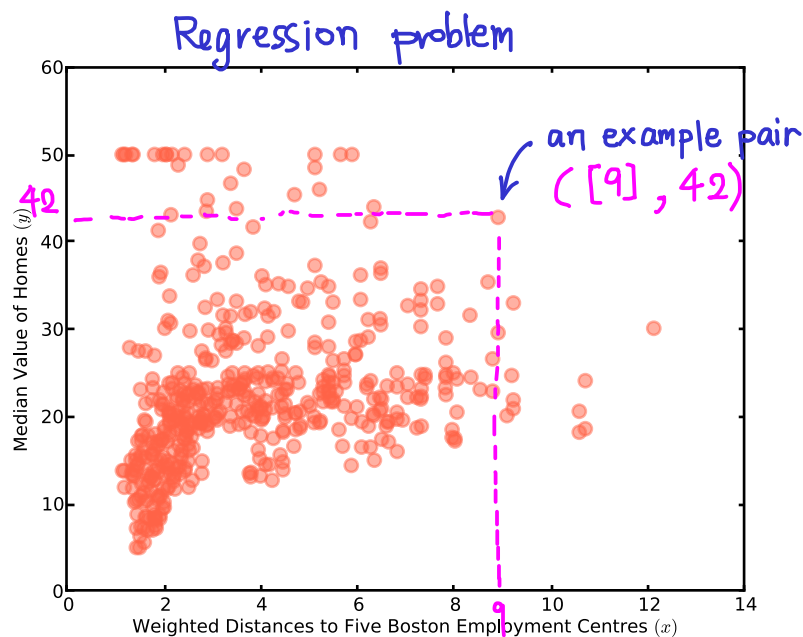
## 12.1.2 Classification and Regression

*(handwritten: → #possible outputs is finite)*

*(handwritten: → #possible outputs is infinite)*

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.[1]



*(handwritten annotations: Regression problem; an example pair ([9], 42); value 42 marked on y-axis)*

**Example 12.2**   Predict the class of iris using its shape.[2]

---

[1]Housing Data Set – `http://archive.ics.uci.edu/ml/datasets/Housing`

[2]Iris Data Set – `http://archive.ics.uci.edu/ml/datasets/Iris`

Classification problem                    $\left(\begin{bmatrix} 7.9 \\ 6.2 \end{bmatrix}, \text{virginica}\right)$



**Example 12.3**   Recognize handwritten digits.[3]

Classification problem



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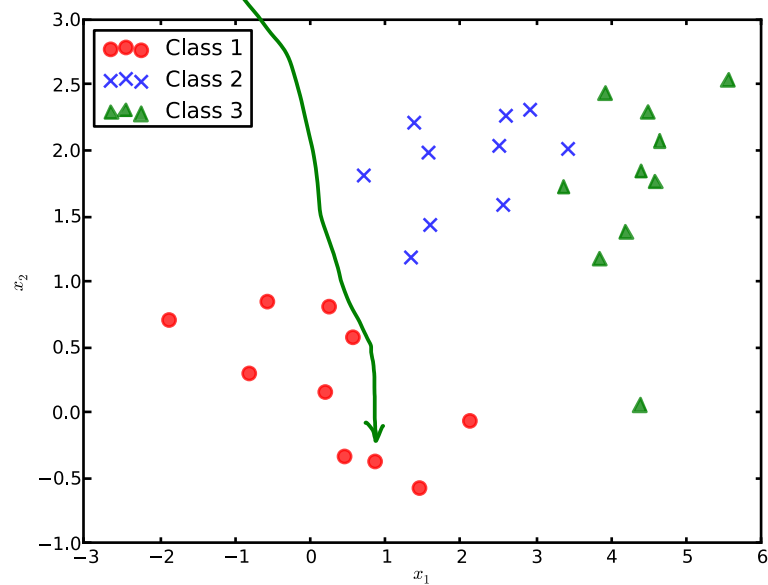
---

[3]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

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

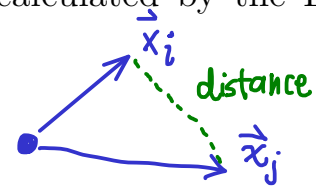*[handwritten annotations: input, a training set, output =0]*

This is known as a **MAP estimate**.[4]

# 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.

*[handwritten annotations: $\vec{x}_i$, distance, $\vec{x}_j$]*

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

*[handwritten annotations:*
$$\vec{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iD} \end{bmatrix} \quad \vec{x}_j = \begin{bmatrix} x_{j1} \\ x_{j2} \\ \vdots \\ x_{jD} \end{bmatrix}$$
*]*

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

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

---

[4]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 |

$x_1$ $x_2$

nearest neighbors

1. $[0.0, 0.0]^\mathsf{T}$ using $K = 3$

$$d\left(\begin{bmatrix}0.0\\0.0\end{bmatrix}, \begin{bmatrix}0.5\\0.5\end{bmatrix}\right) = \sqrt{(0.0-0.5)^2 + (0.0-0.5)^2} = \sqrt{0.5} \quad * \to 0$$

$$d\left(\begin{bmatrix}0.0\\0.0\end{bmatrix}, \begin{bmatrix}0.2\\1.0\end{bmatrix}\right) = \sqrt{(0.0-0.2)^2 + (0.0-1.0)^2} = \sqrt{1.04} \quad * \to 0$$

$$d\left(\begin{bmatrix}0.0\\0.0\end{bmatrix}, \begin{bmatrix}1.0\\0.8\end{bmatrix}\right) = \sqrt{(0.0-1.0)^2 + (0.0-0.8)^2} = \sqrt{1.64} \quad * \to 0$$

$$d\left(\begin{bmatrix}0.0\\0.0\end{bmatrix}, \begin{bmatrix}2.0\\1.5\end{bmatrix}\right) = \sqrt{(0.0-2.0)^2 + (0.0-1.5)^2} = \sqrt{6.25} \quad \text{Predict "0"}$$

$$d\left(\begin{bmatrix}0.0\\0.0\end{bmatrix}, \begin{bmatrix}3.0\\1.0\end{bmatrix}\right) = \sqrt{(0.0-3.0)^2 + (0.0-1.0)^2} = \sqrt{10} \quad \text{for } \begin{bmatrix}0.0\\0.0\end{bmatrix}$$

$$d\left(\begin{bmatrix}0.0\\0.0\end{bmatrix}, \begin{bmatrix}4.0\\3.2\end{bmatrix}\right) = \sqrt{(0.0-4.0)^2 + (0.0-3.2)^2} = \sqrt{26.24}$$

2. $[1.0, 1.0]^\mathsf{T}$ using $K = 5$

$$d\left(\begin{bmatrix}1.0\\1.0\end{bmatrix}, \begin{bmatrix}0.5\\0.5\end{bmatrix}\right) = \sqrt{0.5} \quad * \to 0$$

$$d\left(\begin{bmatrix}1.0\\1.0\end{bmatrix}, \begin{bmatrix}0.2\\1.0\end{bmatrix}\right) = \sqrt{0.64} \quad * \to 0$$

$$d\left(\begin{bmatrix}1.0\\1.0\end{bmatrix}, \begin{bmatrix}1.0\\0.8\end{bmatrix}\right) = \sqrt{0.04} \quad * \to 0$$

$$d\left(\begin{bmatrix}1.0\\1.0\end{bmatrix}, \begin{bmatrix}2.0\\1.5\end{bmatrix}\right) = \sqrt{1.25} \quad * \to 1$$

$$d\left(\begin{bmatrix}1.0\\1.0\end{bmatrix}, \begin{bmatrix}3.0\\1.0\end{bmatrix}\right) = \sqrt{4} \quad * \to 1$$

$$d\left(\begin{bmatrix}1.0\\1.0\end{bmatrix}, \begin{bmatrix}4.0\\3.2\end{bmatrix}\right) = \sqrt{13.84}$$

Predict "0" for $\begin{bmatrix}1.0\\1.0\end{bmatrix}$

input $\to$ calculate distances to all training vectors $\longrightarrow$ find vectors with k smallest distances $\downarrow$ voting

141

$$\text{input} = \begin{bmatrix}2.0\\2.5\end{bmatrix}$$

$$d\left(\begin{bmatrix}2.0\\2.5\end{bmatrix}, \begin{bmatrix}0.5\\0.5\end{bmatrix}\right) = \sqrt{8.5}$$

$$d\left(\begin{bmatrix}2.0\\2.5\end{bmatrix}, \begin{bmatrix}0.2\\1.0\end{bmatrix}\right) = \sqrt{5.49}$$

$$0 \leftarrow * \quad d\left(\begin{bmatrix}2.0\\2.5\end{bmatrix}, \begin{bmatrix}1.0\\0.8\end{bmatrix}\right) = \sqrt{3.89}$$

$$1 \leftarrow * \quad d\left(\begin{bmatrix}2.0\\2.5\end{bmatrix}, \begin{bmatrix}2.0\\1.5\end{bmatrix}\right) = \sqrt{1}$$

$$1 \leftarrow * \quad d\left(\begin{bmatrix}2.0\\2.5\end{bmatrix}, \begin{bmatrix}3.0\\1.0\end{bmatrix}\right) = \sqrt{3.25}$$

$$d\left(\begin{bmatrix}2.0\\2.5\end{bmatrix}, \begin{bmatrix}4.0\\3.0\end{bmatrix}\right) = \sqrt{4.25}$$

## 12.3 Model Selection

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

A basic way to evaluate a classifier is to use the *misclassification rate* on the training set. It can be defined as:
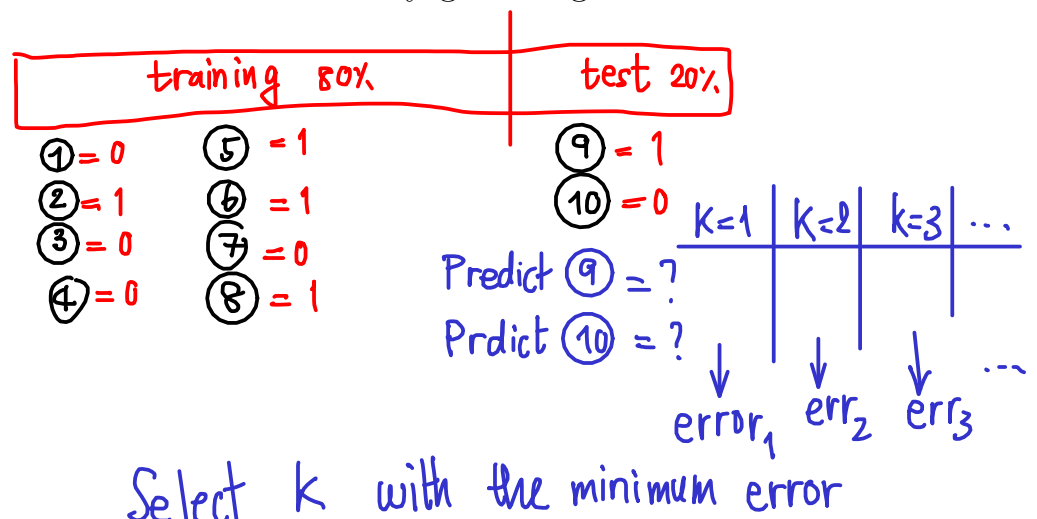
$$\text{err}(f, \mathcal{D}) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{I}(f(\mathbf{x}_i) \neq t_i)$$

However, the error on the training set should not represent the *generalization error*.
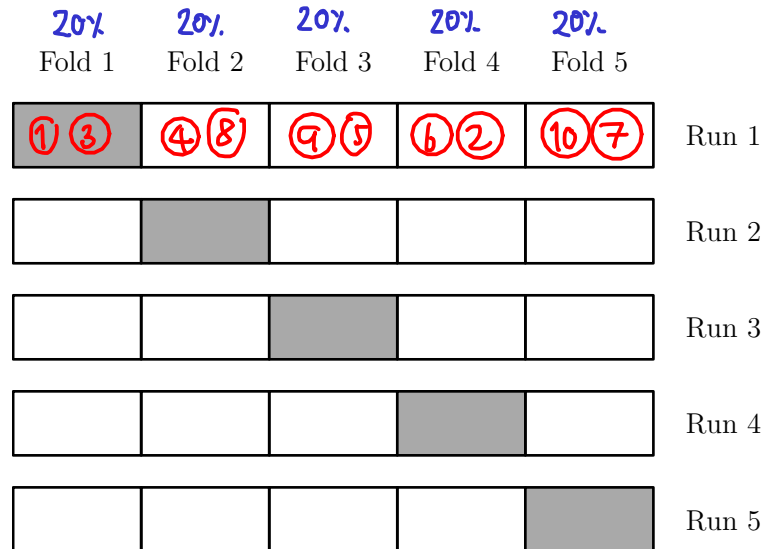
### 12.3.1 Holdout Method

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.
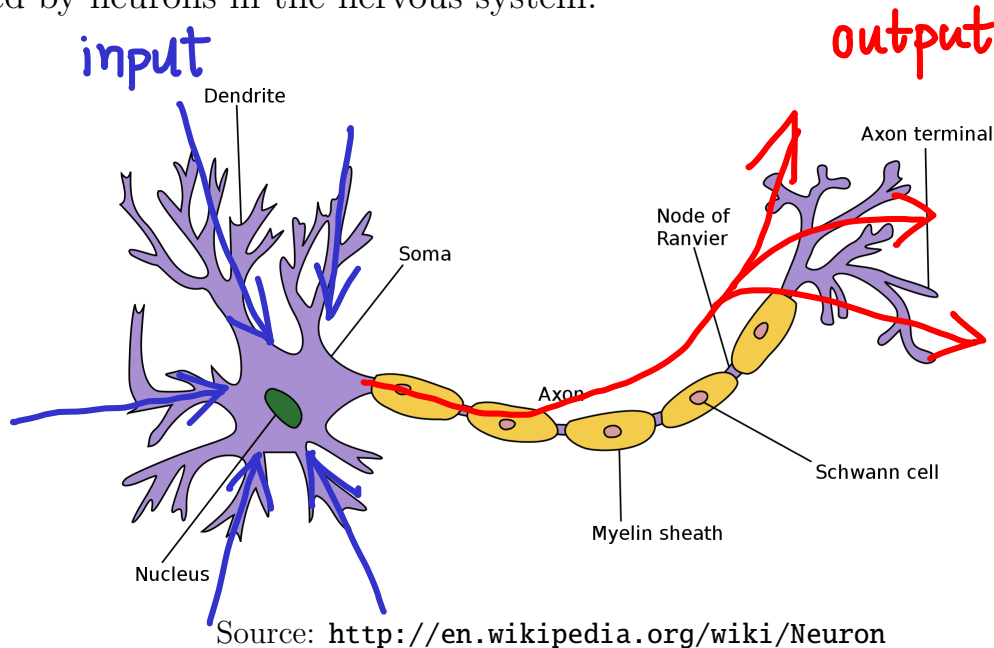
## 12.3.2 Cross Validation Method

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.
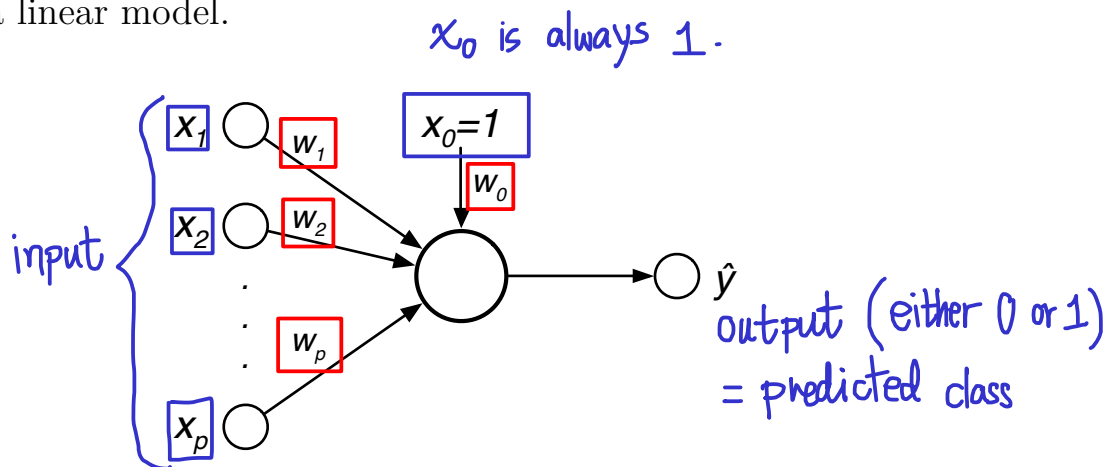
# 12.4 Artificial Neural Networks

**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 $(x_0, x_1, x_2, \ldots, x_p)^\mathsf{T}$. A weight is attached to each input as a vector $(w_0, w_1, w_2, \ldots, w_p)^\mathsf{T}$. 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

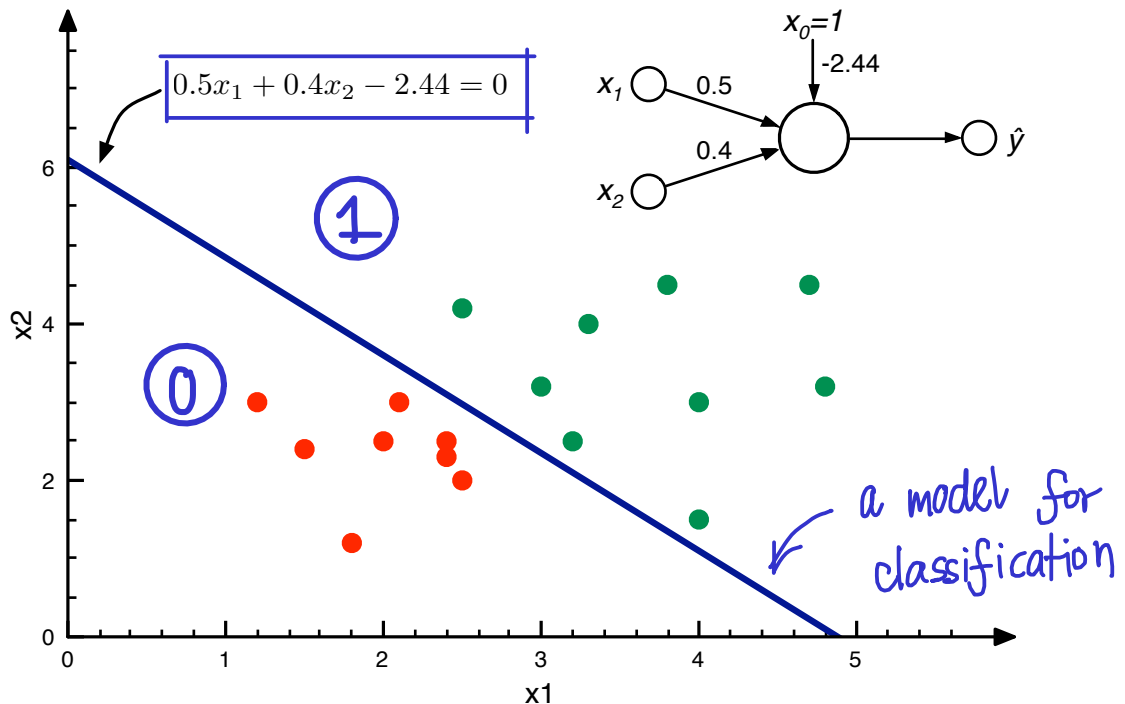Given an input vector $(x_1, x_2, \ldots, x_p)^T$, we have

predicted ← output   $\boxed{y} = \text{sgn}(\sum_{i=0}^{p} w_i x_i)$

$$= \text{sgn}(w_0 + \boxed{w_1 x_1} + \cdots + \boxed{w_p x_p})$$

weighted sum of the input vector

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

↓ the sign function

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.
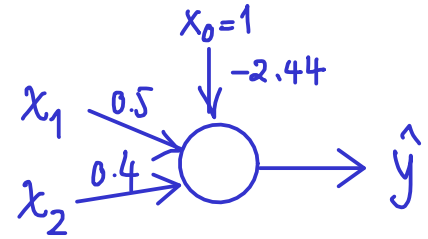
$\hat{y} = \text{sgn}(-2.44 + 0.5x_1 + 0.4x_2)$



$\boxed{0.5x_1 + 0.4x_2 - 2.44 = 0}$

$x_0 = 1$
$-2.44$
$x_1$   0.5
0.4
$x_2$
$\hat{y}$

①    ⓪

a model for classification

a perceptron = a line that classified the space
(with appropriate    into two areas $(1, 0)$
weights)

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

$x_1 \quad x_2$

1. $[0.0, 0.0]^\mathsf{T}$

$$\hat{y} = \text{sgn}(-2.44 + (0.5)(0) + (0.4)(0))$$
$$= \text{sgn}(-2.44) = 0$$

2. $[4.0, 4.0]^\mathsf{T}$

$$\hat{y} = \text{sgn}(-2.44 + (0.5)(4.0) + (0.4)(4.0))$$
$$= \text{sgn}(1.16) = 1$$

3. $[3.0, 2.0]^\mathsf{T}$

$$\hat{y} = \text{sgn}(-2.44 + (0.5)(3.0) + (0.4)(2.0))$$
$$= \text{sgn}(-0.14) = 0$$

$x_0 = 1$

$x_1 \xrightarrow{0.5}$   $-2.44$

$x_2 \xrightarrow{0.4}$   $\longrightarrow \hat{y}$

Examples $\longrightarrow$ Training $\longrightarrow$

to find the most appropriate weights from a set of labeled examples

$x_1, x_2$
$\downarrow$
Model
$\downarrow$
$\hat{y}$
predicted output

Examples $\searrow$

$x_1, x_2 \nearrow$   k-NN $\longrightarrow \hat{y}$ predicted output

inputs — target output

**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 |

↳ Training examples

1. Draw the structure of a perceptron for learning from this dataset.



$X_0 = 1$

$x_1$   $w_1 = ?$   $w_0 = ?$

$\hat{y}$

$x_2$   $w_2 = ?$

We need to find an appropriate vector of weights from the given dataset.

2. Find the weights that makes the perceptron correctly classify the train-
   ing examples. *manually*



$x_2$

$(\begin{bmatrix} 0 \\ 1 \end{bmatrix}, 0)$   $(\begin{bmatrix} 1 \\ 1 \end{bmatrix}, 1)$

$(1.5, 0)$ We place a line here to make it
correctly classify the training
examples.

$(0, 1.5)$

$x_1$

$(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, 0)$

$(\begin{bmatrix} 1 \\ 0 \end{bmatrix}, 0)$

$\rightarrow w_0 + w_1 x_1 + w_2 x_2 = 0$

Let's find the weights from this line.

From $(0, 1.5)$, we have   $w_0 + w_1(0) + w_2(1.5) = 0$

$$w_0 + 1.5 w_2 = 0 \quad \text{——— (1)}$$

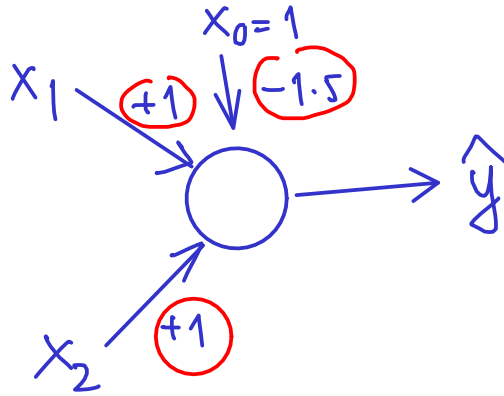From $(1.5, 0)$, we have   $w_0 + w_1(1.5) + w_2(0) = 0$

$$w_0 + 1.5 w_1 = 0 \quad \text{——— (2)}$$

147

$(2) - (1)$   $1.5 w_1 - 1.5 w_2 = 0 \quad \longrightarrow \quad w_1 = w_2$

when $w_2 = +1 \longrightarrow w_1 = +1$

substitute $w_2$ to (1) $\rightsquigarrow$ $w_0 = -1.5$

$$-1.5 + x_1 + x_2 = 0$$

$x_0 = 1$

$x_1$ $(+1)$ $\downarrow$ $(-1.5)$

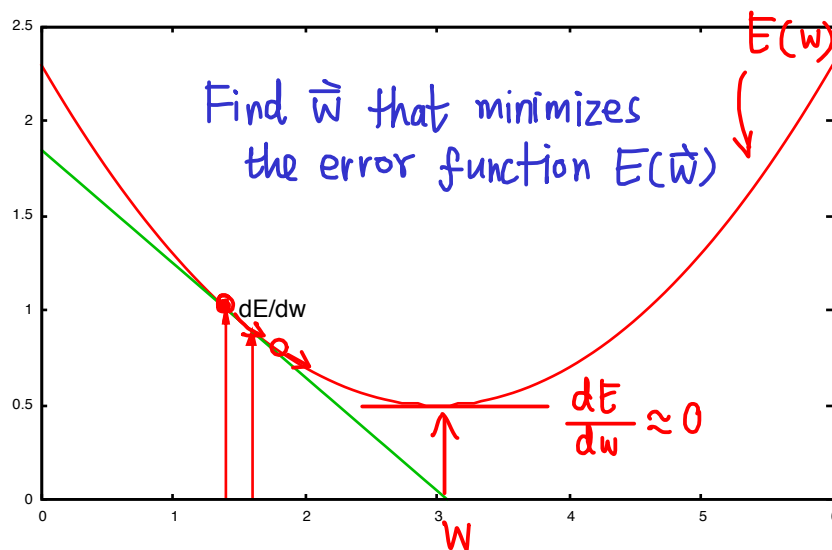$\rightarrow \hat{y}$

$x_2$ $(+1)$

## 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 $\mathbf{w}$ *[weight vector]* i.e. $E(\mathbf{w})$. *[→ error function]*

2. Randomly select initial values for $\mathbf{w}$.

3. Update weights using $\boxed{\mathbf{w} - \eta \frac{dE}{d\mathbf{w}}}$ where $\eta$ is a positive small number.

4. Iteratively update until $\frac{dE}{d\mathbf{w}} \approx 0$



*Find $\vec{w}$ that minimizes the error function $E(\vec{w})$*

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

Error function

target output

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

predicted output

$$= \frac{1}{2} \sum_{d=1}^{N} \left( t_d - \sum_{i=0}^{p} w_i x_{di} \right)$$

$\longrightarrow$ a perceptron

We compute gradient of the objective function which is denoted as

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

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

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \left( \frac{1}{2} \sum_{d=1}^{N} (t_d - y_d)^2 \right)$$

$$= \frac{1}{2} \sum_{d=1}^{N} \frac{\partial}{\partial w_i} \left( (t_d - y_d)^2 \right)$$

$$= \frac{1}{2} \sum_{d=1}^{N} 2(t_d - y_d) \frac{\partial}{\partial w_i} \left( t_d - y_d \right)$$

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

$$= \sum_{d=1}^{N} (t_d - y_d)(-x_{di}) = \nabla E(w_i)$$

a new weight $w_i$

the current $w_i$

From the gradient of $E$, we have

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

We then update $w_i$ by

input

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

$\rightarrow$ Weight updating rule

target output    predicted output

feed an input vector $\longrightarrow$ make a prediction

if there is a error $\longrightarrow$ update weights

$(t \neq y)$

*→ find the most appropriate weights for a perceptron.*

A <u>training algorithm</u> for perceptron is as below:

1. Start with random weights.

*↗ $t \neq y$*

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

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

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 $([2.5, 2.0]^{\mathsf{T}}, 0)$ to the perceptron, then

$$y = \mathrm{sgn}(0.18 \times 2.5 + 0.2 \times 2.0 - 0.4)$$
$$= \mathrm{sgn}(0.45)$$
$$= 1 \quad \leftarrow \text{predicted output}$$

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

$$w_0 \leftarrow -0.40 + \big((0.01)(0-1)(1.0)\big) = -0.410$$
$$w_1 \leftarrow +0.18 + \big((0.01)(0-1)(2.5)\big) = +0.155$$
$$w_2 \leftarrow +0.20 + \big((0.01)(0-1)(2.0)\big) = +0.180$$

current $w_i$ $\quad \eta \quad t \quad y \quad x_i \quad$ new $w_i$
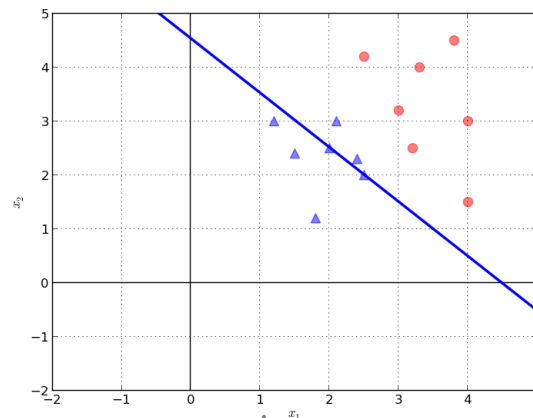
*#2 ~ #14*

After applying all the remaining 13 examples, we will get weights as:
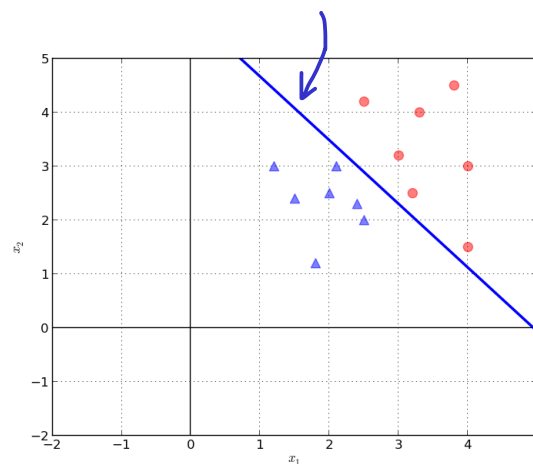
$$w_0 = -0.440; \quad w_1 = 0.098; \quad w_2 = 0.097$$



*#1 ~ #14 (5 times)*

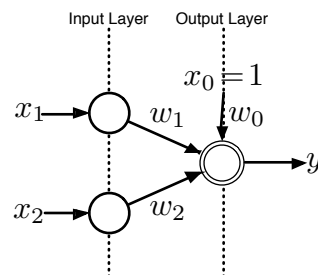After updating the weights for five rounds, the weights become:

$$w_0 = -0.450; \quad w_1 = 0.091; \quad w_2 = 0.077$$
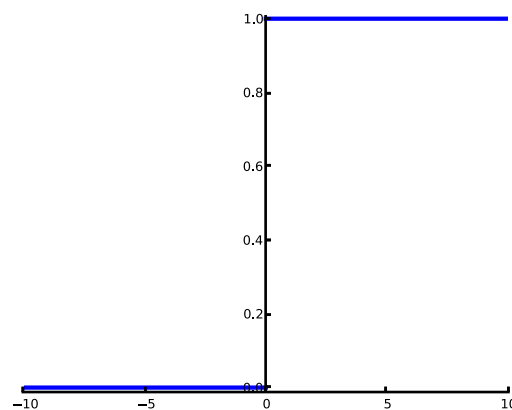
## 12.4.2 Threshold Function

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. $\mathrm{sgn}(\cdot)$).

$$y = \mathrm{sgn}(w_0 + w_1 x_1 + w_2 x_2)$$



The theshold function $\mathrm{sgn}(\cdot)$ is a discontinuous function. It cannot be differentiated. Therefore, it cannot be used together with the gradient descent, or other optimization techniques.
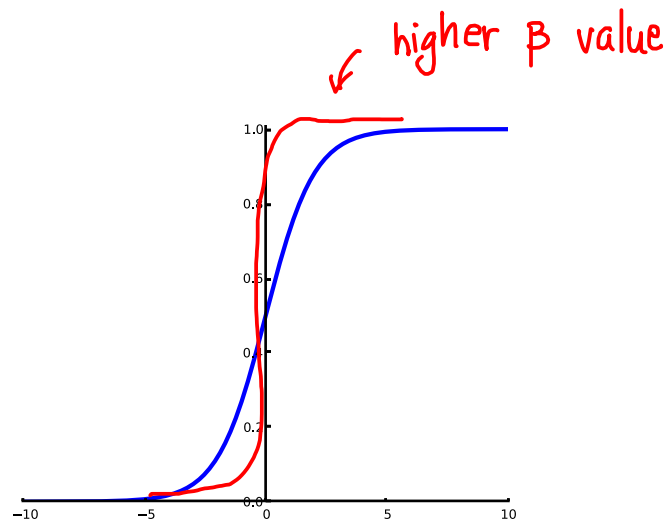
$$\mathrm{sgn}(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$



The sigmoid function is an S-shaped function. It can be in place of the function $\mathrm{sgn}(\cdot)$.
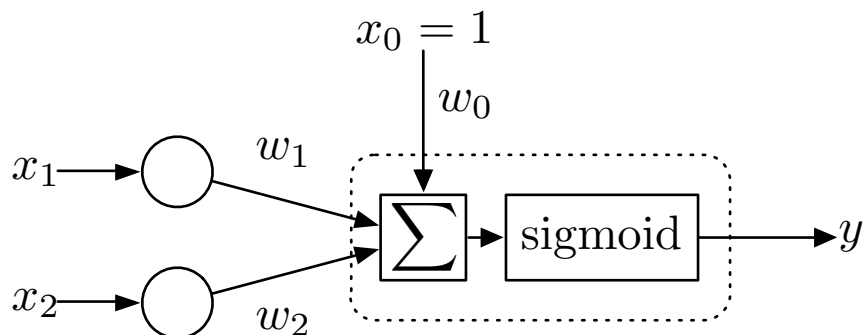
$$\mathrm{sigmoid}(x) = \frac{1}{1 + e^{-\beta x}}$$

$\beta = $ a positive constant

### 12.4.3 Sigmoid Unit

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

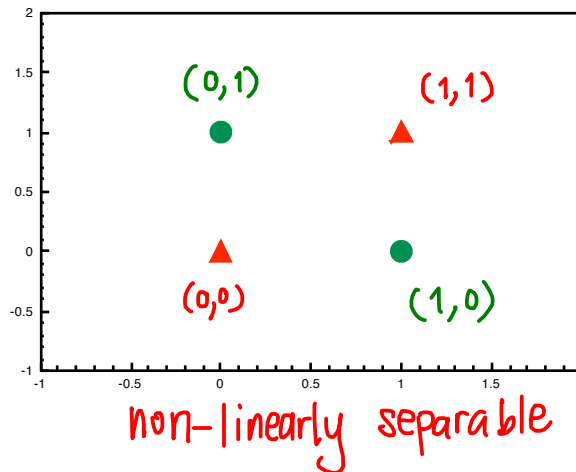$$y = \text{sigmoid}(w_0 + w_1x_1 + w_2x_2)$$

## 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.
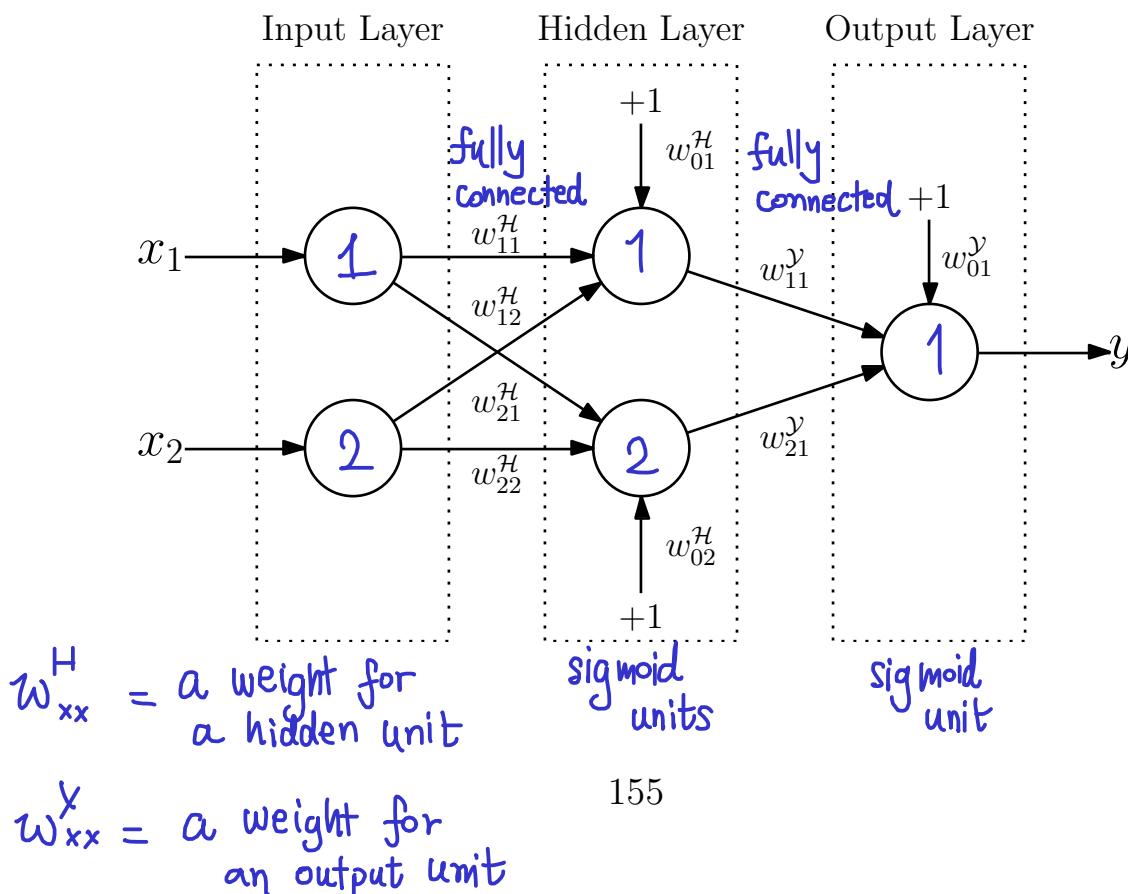
XOR Problem

| $x_1$ | $x_2$ | $t$ |
|:---:|:---:|:---:|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

a perceptron
$\equiv$ a line
　(a hyperplane)

non-linearly separable

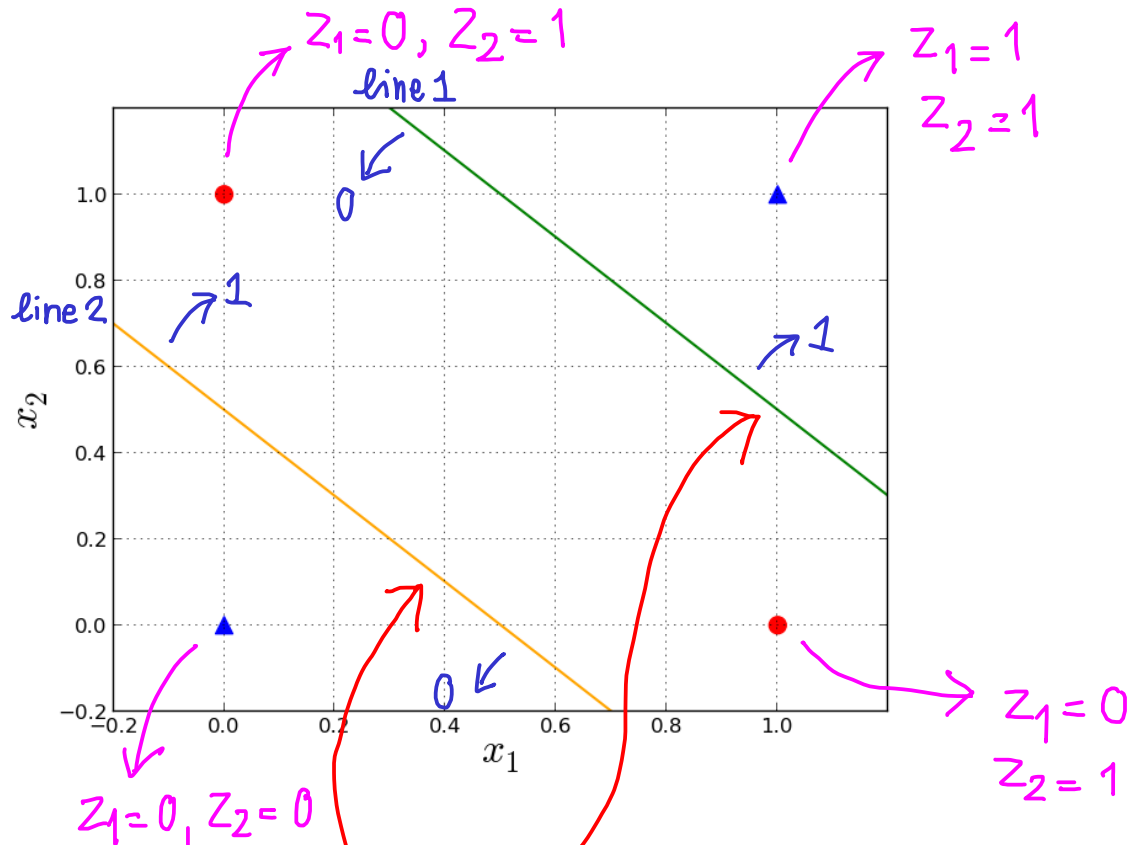How can the perceptron be improved to handle this situation?

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.* = a neural network



Input Layer　　Hidden Layer　　Output Layer

fully connected　fully connected

$w_{01}^{\mathcal{H}}$　$w_{01}^{\mathcal{Y}}$

$w_{11}^{\mathcal{H}}$　$w_{12}^{\mathcal{H}}$　$w_{21}^{\mathcal{H}}$　$w_{22}^{\mathcal{H}}$

$w_{11}^{\mathcal{Y}}$　$w_{21}^{\mathcal{Y}}$

$w_{02}^{\mathcal{H}}$

sigmoid units　sigmoid unit

$w_{xx}^{H}$ = a weight for a hidden unit

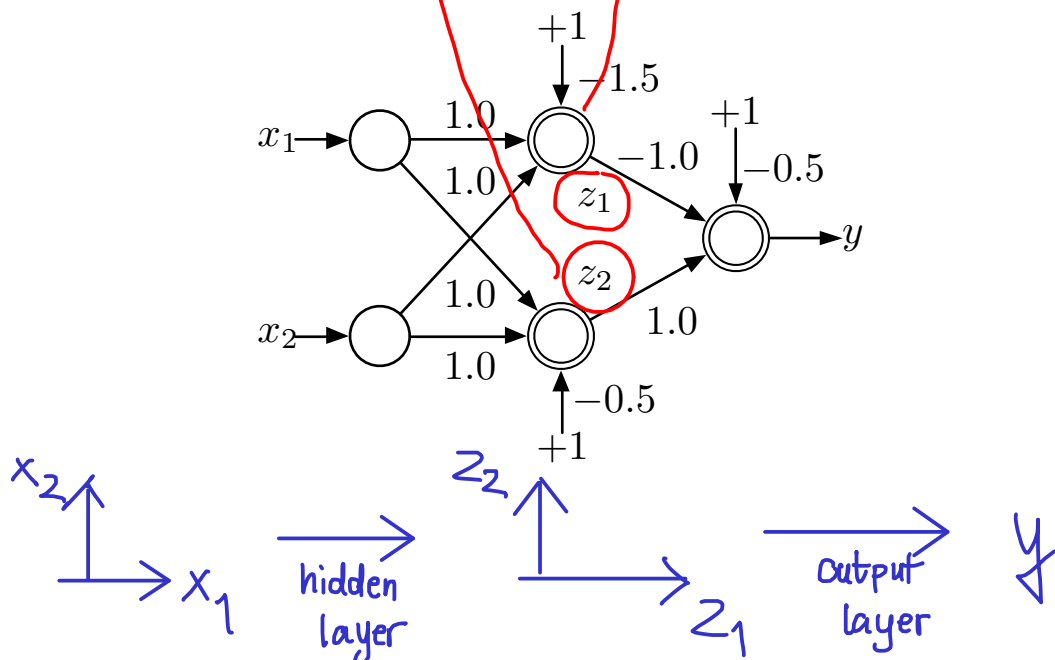$w_{xx}^{Y}$ = a weight for an output unit

155

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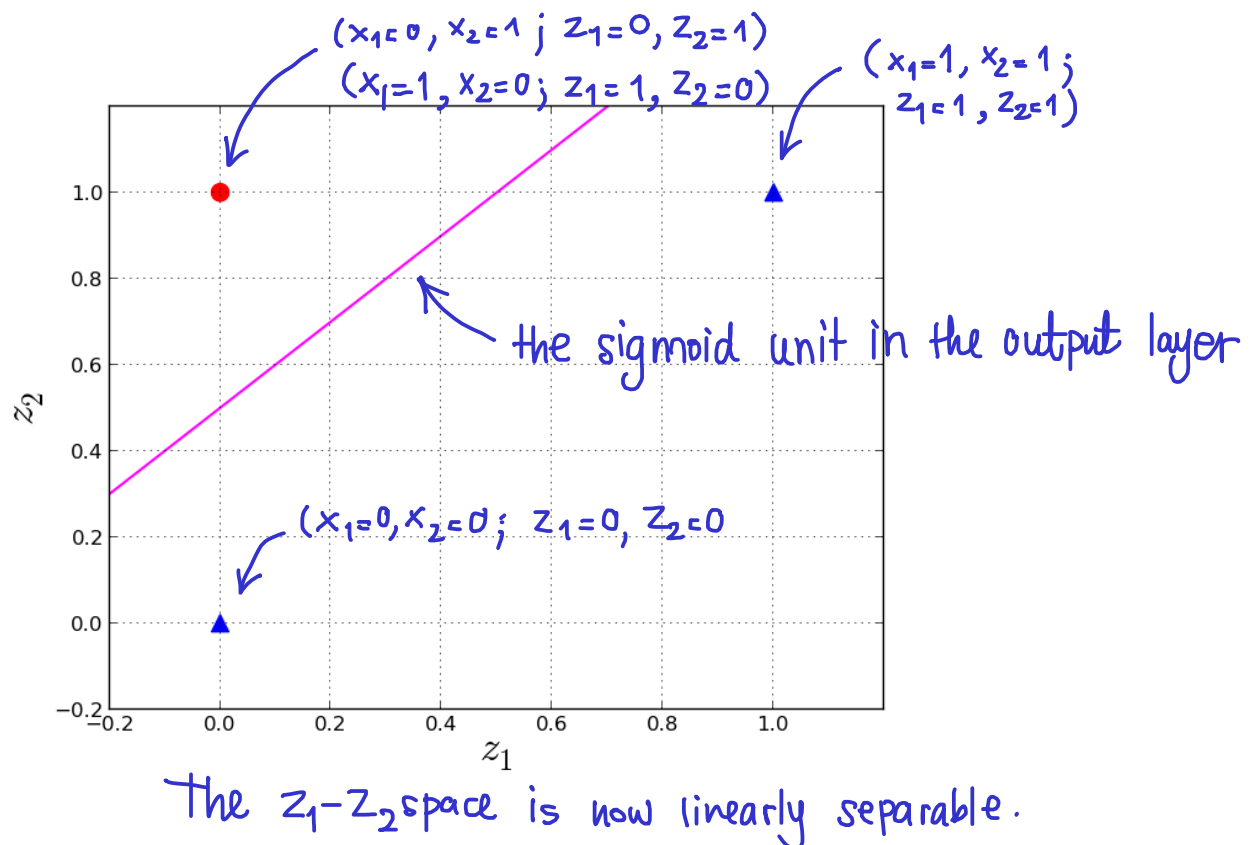
$$z_1 = \text{sigmoid}(x_1 + x_2 - 0.5)$$
$$z_2 = \text{sigmoid}(x_1 + x_2 - 1.5)$$

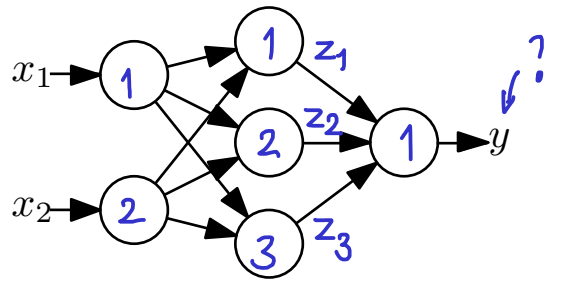From the network, we have two intermediate outputs (i.e. $z_1$, and $z_2$) from the perceptrons in the hidden layer: *hidden layer*

| $x_1$ | $x_2$ | $z_1$ | $z_2$ | $t$ |
|---|---|---|---|---|
| 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:



*($x_1 = 0$, $x_2 = 1$ ; $z_1 = 0$, $z_2 = 1$)*
*($x_1 = 1$, $x_2 = 0$ ; $z_1 = 1$, $z_2 = 0$)*
*($x_1 = 1$, $x_2 = 1$ ; $z_1 = 1$, $z_2 = 1$)*
*the sigmoid unit in the output layer*
*($x_1 = 0$, $x_2 = 0$ ; $z_1 = 0$, $z_2 = 0$*
*The $z_1 - z_2$ space is now linearly separable.*

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



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

$$Z_1 = \text{sigmoid}(w_{01}^H + w_{11}^H x_1 + w_{21}^H x_2) = \text{sigmoid}(0.5 + (-0.5)(1.0) + (0.1)(1.0))$$
$$= \text{sigmoid}(0.1) = \frac{1}{1 + e^{-(1)(0.1)}} = 0.52$$

$$Z_2 = \text{sigmoid}(w_{02}^H + w_{12}^H x_1 + w_{22}^H x_2) = \text{sigmoid}(1 + (0.2)(1) + (-0.2)(1))$$
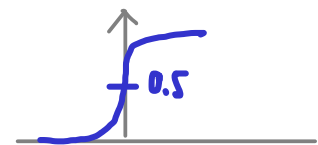$$= \text{sigmoid}(1) = 0.73$$

$$Z_3 = \text{sigmoid}(w_{03}^H + w_{13}^H x_1 + w_{23}^H x_2) = \text{sigmoid}(-1 + (0.1)(1) + (0.5)(1))$$
$$= \text{sigmoid}(-0.4) = 0.40$$

$$y = \text{sigmoid}(w_{01}^Y + w_{11}^Y Z_1 + w_{21}^Y Z_2 + w_{31}^Y Z_3)$$
$$= \text{sigmoid}(-1 + (1.5)(0.52) + (-1)(0.73) + (1)(0.40))$$
$$= \text{sigmoid}(-0.58) = 0.36$$

using the threshold of 0.5, the predicted output is 0

$(0.36 < 0.5)$

## 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:**   *Given an input vector, compute the predicted output*

    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:**   *Update $\vec{w}^H$ and $\vec{w}^Y$ according to the error*  $\delta^Y$  $(t - y)$

    1. compute the <u>error at the output</u> using:

$$\delta_k^{\mathcal{Y}} = (t_k - y_k)y_k(1 - y_k)$$
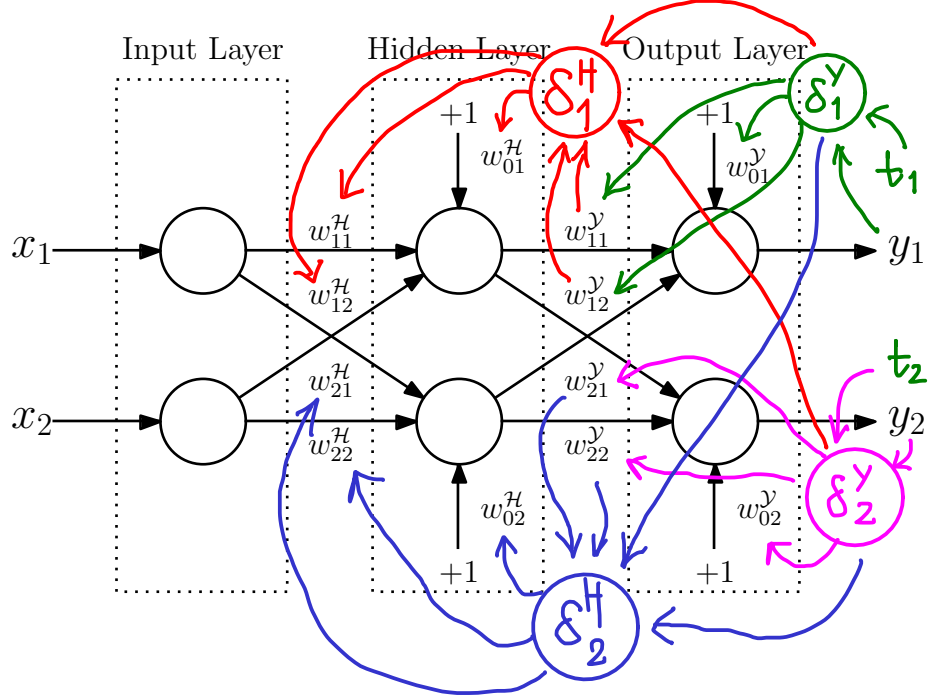
    2. update the output layer weights using:

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

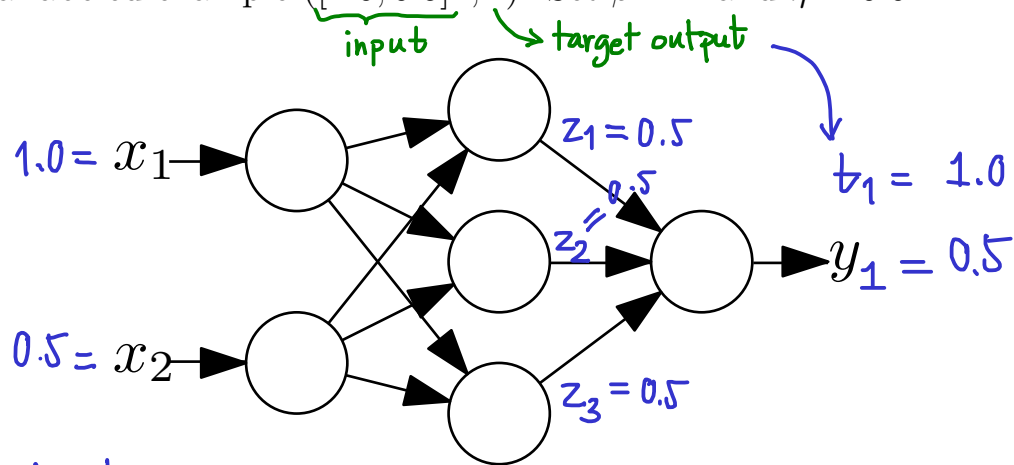    3. compute the <u>error in the hidden layer</u> using:

$\delta^H$

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

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 $([1.0, 0.5]^T, 1)$. Set $\beta = 1$ and $\eta = 0.01$.



Backward phase

$$\delta_1^Y = (t_1 - y_1)(y_1)(1 - y_1) = 0.125$$

$$W_{11}^Y = 0 + \eta\,\delta_1^Y Z_1 = 0 + (0.01)(0.125)(0.5) = 0.000625$$

$$W_{21}^Y = 0 + \eta\,\delta_1^Y Z_2 = 0.000625$$

$$W_{31}^Y = 0 + \eta\,\delta_1^Y Z_3 = 0.000625$$

$$W_{01}^Y = 0 + \eta\,\delta_1^Y (1) = 0.00125$$

$$\delta_1^H = Z_1(1 - Z_1) \sum_{\textcircled{k}} W_{1k}^Y \delta_K^Y$$

$$\textcircled{k} \rightarrow \text{\#output units}$$

$$= (0.5)(1 - 0.5)(0.000625)(0.125) = 1.95 \times 10^{-5}$$

$$W_{01}^H = 0 + \eta\,\delta_1^H (1) = 0 + (0.01)(1.95 \times 10^{-5})(1) = 1.95 \times 10^{-7}$$

$$W_{11}^H = 0 + \eta\,\delta_1^H x_1 = 0 + (0.01)(1.95 \times 10^{-5})(1) = 1.95 \times 10^{-7}$$

$$W_{21}^H = 0 + \eta\,\delta_1^H x_2 = 0 + (0.01)(1.95 \times 10^{-5})(0.5) = 9.75 \times 10^{-8}$$

$$\delta_2^H \rightarrow W_{02}^H, \quad W_{12}^H, \quad W_{22}^H$$

$$\delta_3^H \rightarrow W_{03}^H, \quad W_{13}^H, \quad W_{23}^H$$

# 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.