# Practical Machine Learning

**Lecture 2**
**Linear models for classification and regression**

Dr. Suyong Eum

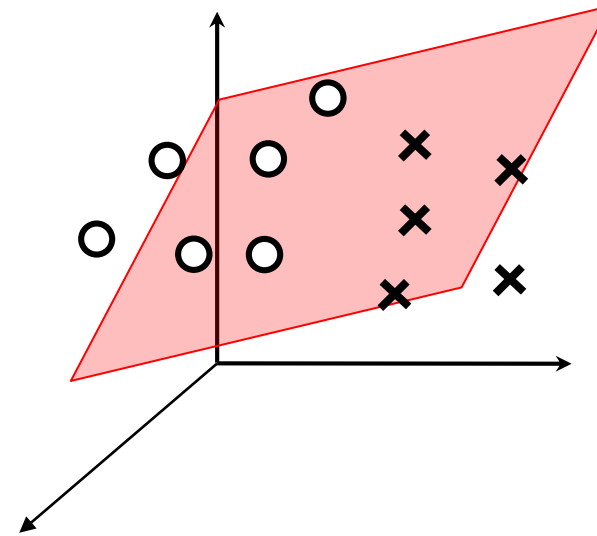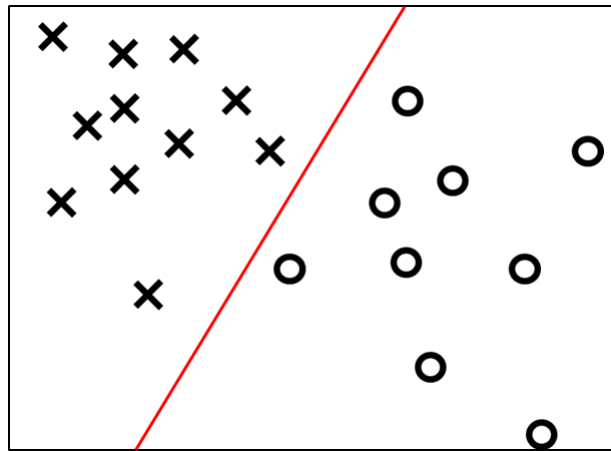OSAKA UNIVERSITY

# Lecture Outline

❑ **Linear classification**

- Perceptron algorithm

❑ **Linear regression**

- Mean Square Error (MSE)

- Normal Equation

- Gradient descent

❑ **Logistic regression**

- Cross Entropy Error (CEE)

# Classification

- ❑ Decision boundary (surfaces)
- ❑ Decision regions
- ❑ (D-1)-dimensional hyperplane within the D-dimensional input space

1) Probabilistic approach
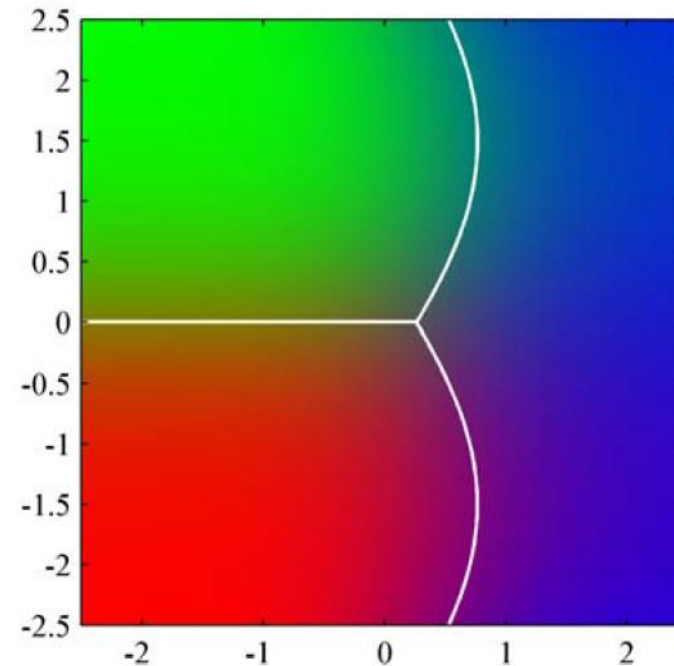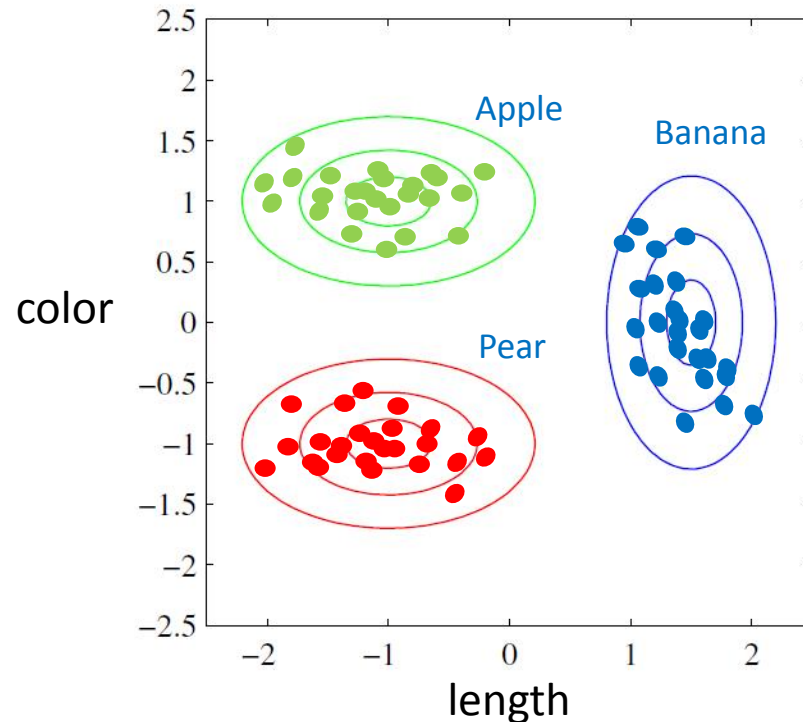   - Focus on the development of a probability model, e.g., cancer/non-cancer

2) Deterministic approach
   - Focus on the determination of a decision boundary

Christopher M. Bishop Pattern Recognition and Machine Learning, 2001

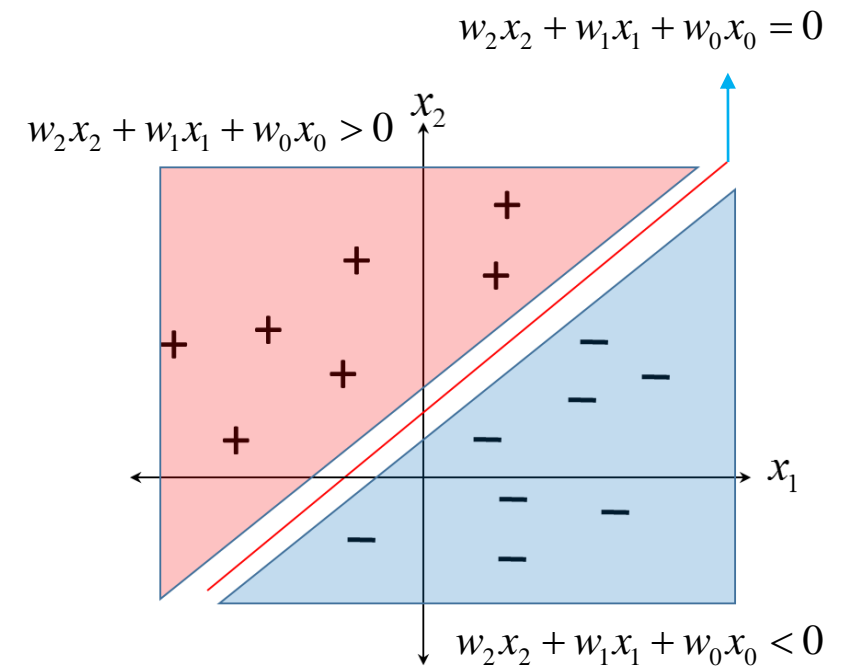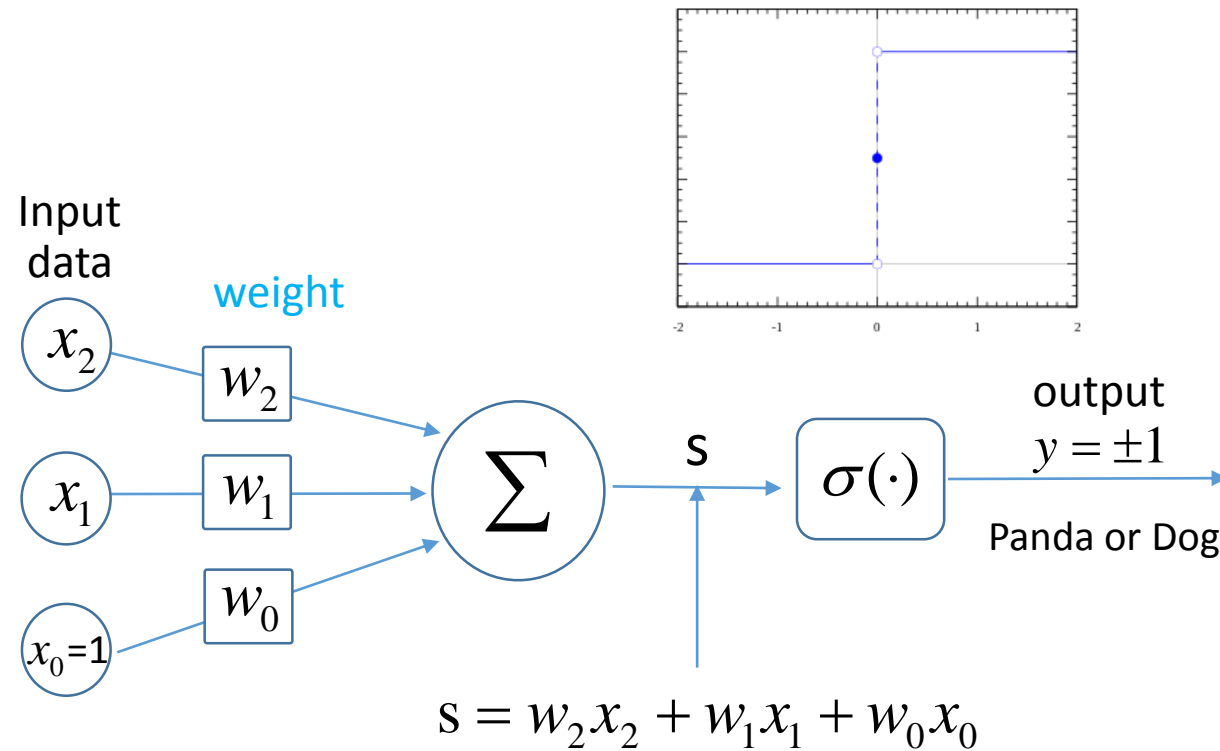❑ Rosenblatt (1962) introduced the perceptron algorithm.

**Frank Rosenblatt**
**1928–1969**

Rosenblatt's perceptron played an important role in the history of machine learning. Initially, Rosenblatt simulated the perceptron on an IBM 704 computer at Cornell in 1957, but by the early 1960s he had built special-purpose hardware that provided a direct, parallel implementation of perceptron learning. Many of his ideas were encapsulated in "Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms" published in 1962. Rosenblatt's work was criticized by Marvin Minksy, whose objections were published in the book "Perceptrons", co-authored with Seymour Papert. This book was widely misinterpreted at the time as showing that neural networks were fatally flawed and could only learn solutions for linearly separable problems. In fact, it only proved such limitations in the case of single-layer networks such as the perceptron and merely conjectured (incorrectly) that they applied to more general network models. Unfortunately, however, this book contributed to the substantial decline in research funding for neural computing, a situation that was not reversed until the mid-1980s. Today, there are many hundreds, if not thousands, of applications of neural networks in widespread use, with examples in areas such as handwriting recognition and information retrieval being used routinely by millions of people.
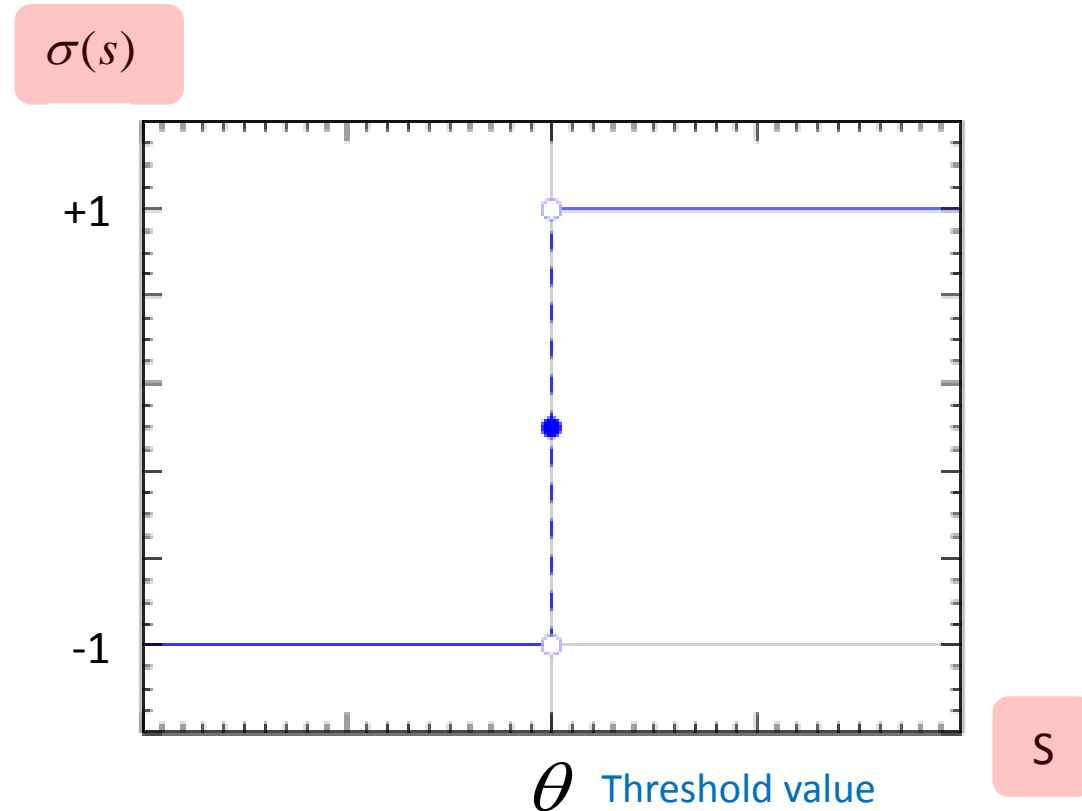
Christopher M. Bishop Pattern Recognition and Machine Learning, 2001

❑ Simplest form of a neural network used for a linear classification.



Input
data

weight

$x_2$

$w_2$

$x_1$

$w_1$

$x_0 = 1$

$w_0$

$\sum$

s

$\sigma(\cdot)$

output
$y = \pm 1$

Panda or Dog

$$\mathrm{s} = w_2 x_2 + w_1 x_1 + w_0 x_0$$

$$w_2 x_2 + w_1 x_1 + w_0 x_0 = 0$$

$$w_2 x_2 + w_1 x_1 + w_0 x_0 > 0$$

$x_2$

$x_1$

$$w_2 x_2 + w_1 x_1 + w_0 x_0 < 0$$

❑ It uses a step function as an activation function: Yes/No question.
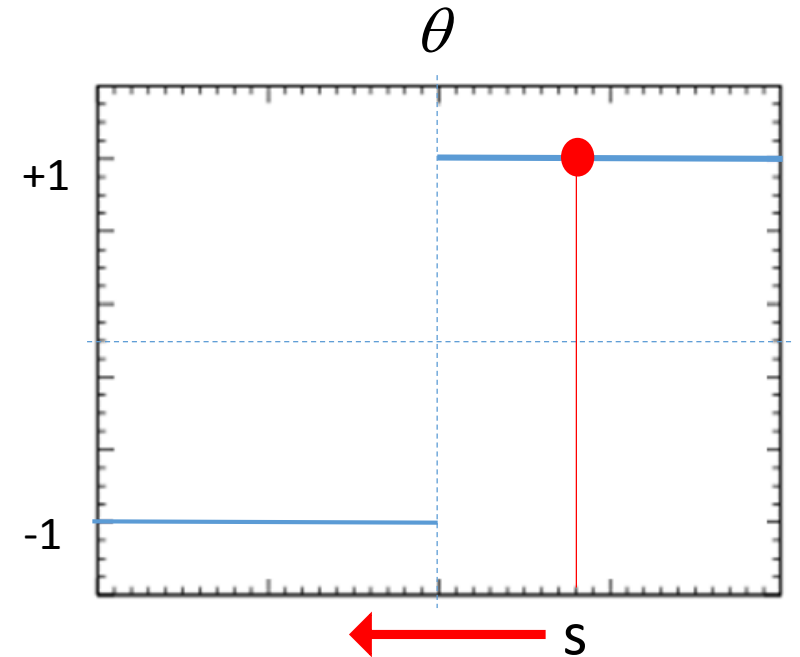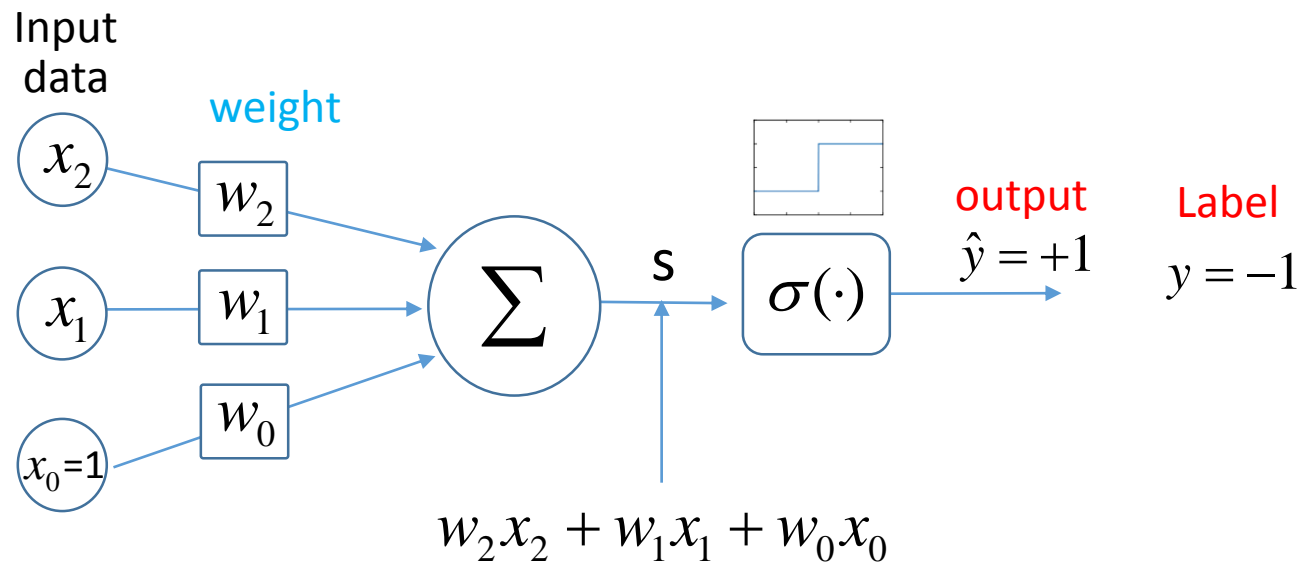
❑ The step function is discontinuous: any problem?

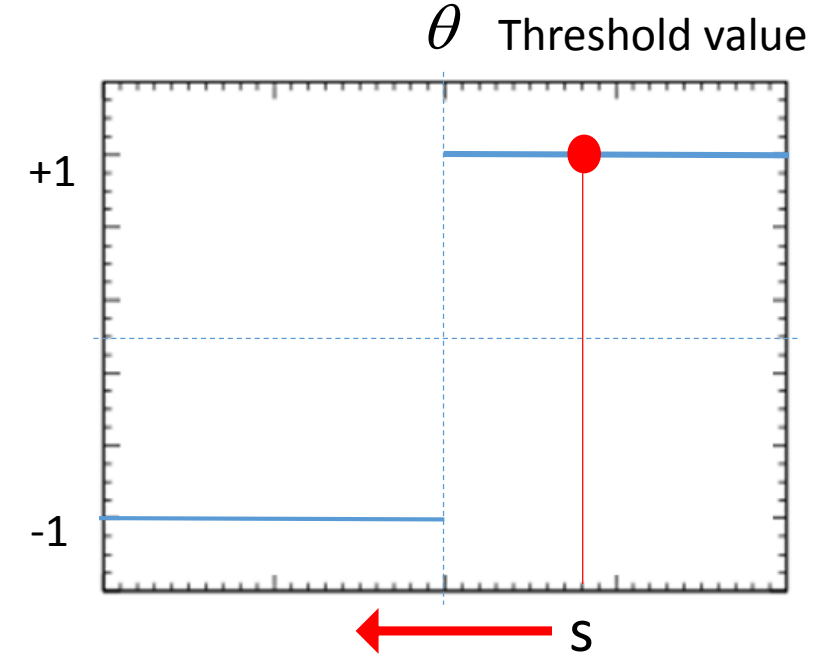$$\sigma(s) = \begin{cases} +1 & if\ s > \theta \\ -1 & otherwise \end{cases}$$

$\sigma(s)$
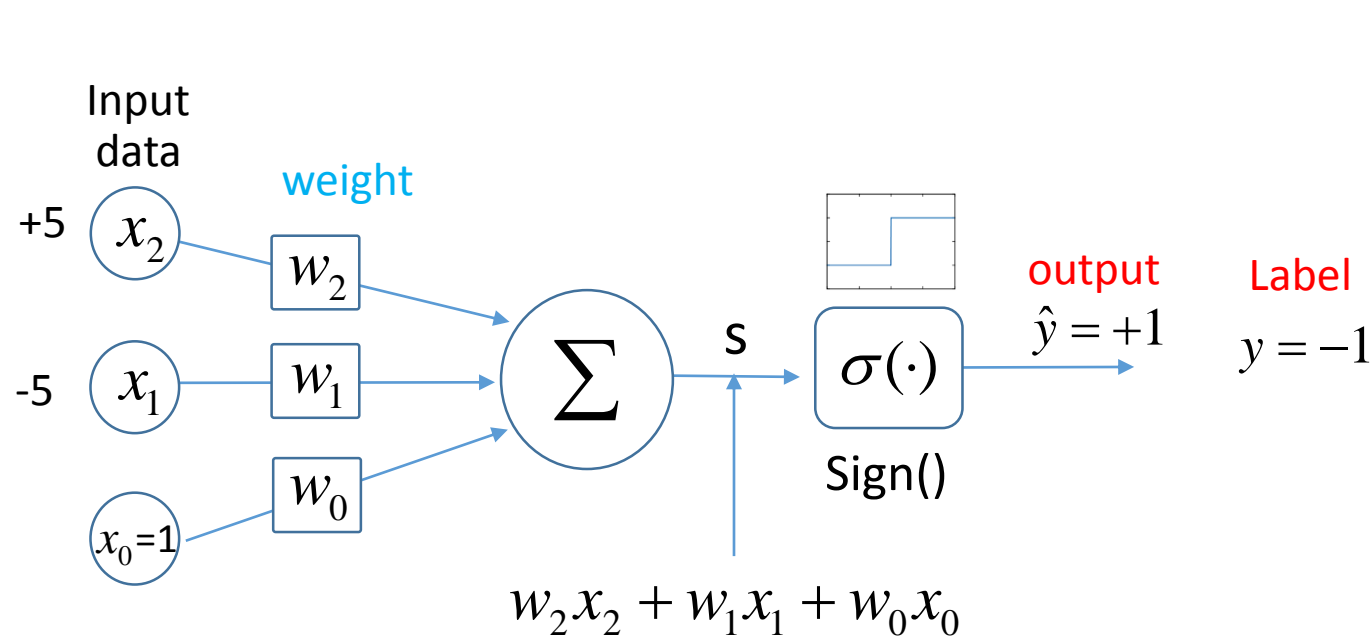
+1

-1

$\theta$  Threshold value

s

❑ Assuming that a data point X (*x1, x2*) has a label (-1)

❑ Assuming that it is misclassified to be (+1)

❑ Weight should be updated.

Input data

weight

$x_2$

$w_2$

$x_1$

$w_1$

$x_0 = 1$

$w_0$

$\sum$

s

$\sigma(\cdot)$

output

$\hat{y} = +1$

Label

$y = -1$

$w_2 x_2 + w_1 x_1 + w_0 x_0$

$\theta$

+1

-1

s

$$\mathrm{w}^{(\text{new})} = \mathrm{w}^{(\text{old})} + \gamma \cdot y_n \cdot \mathrm{X_n}$$

$\theta$  Threshold value

Input data

weight

+5  $x_2$

-5  $x_1$

$x_0 = 1$

$w_2$

$w_1$

$w_0$

$\sum$

s

$\sigma(\cdot)$

Sign()

output

$\hat{y} = +1$

Label

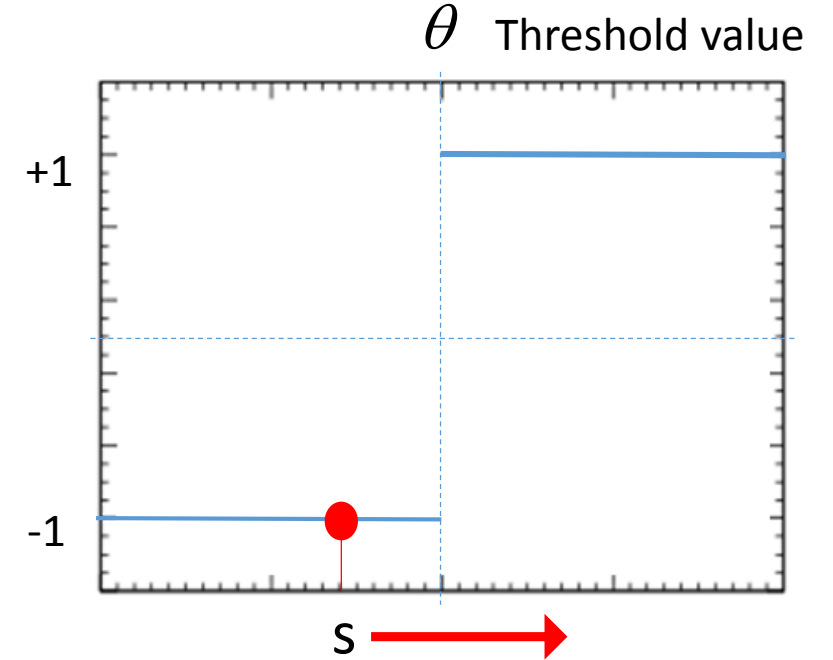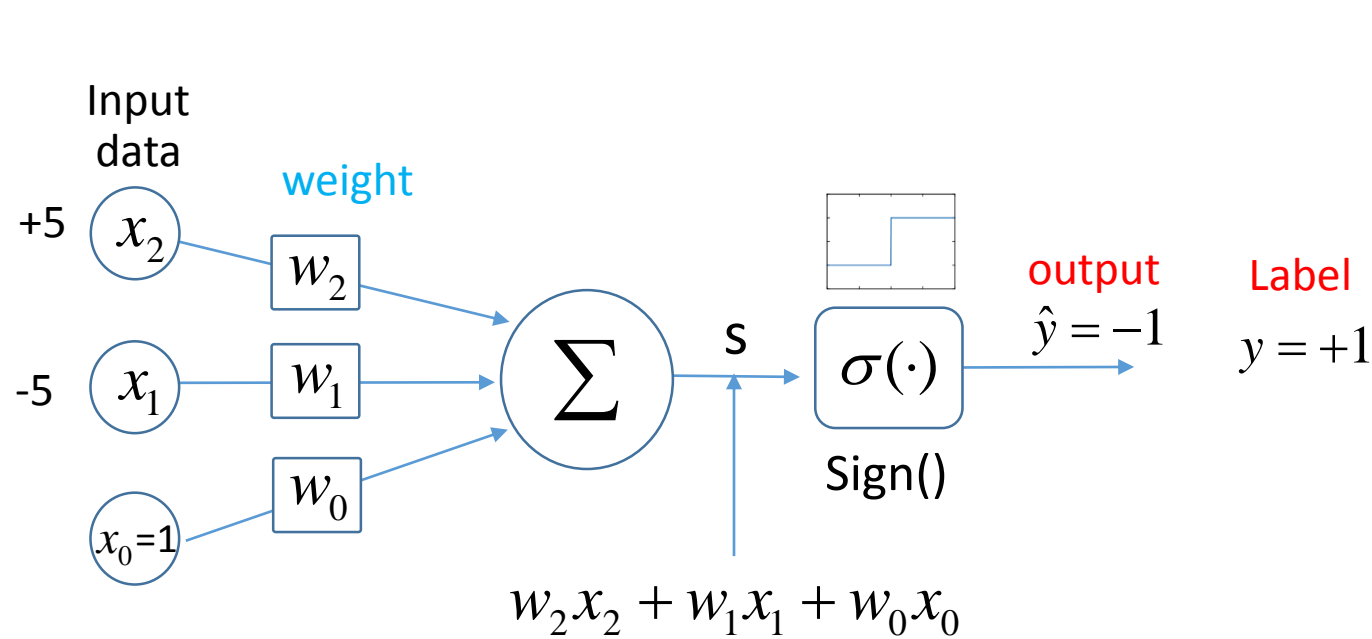$y = -1$

$w_2 x_2 + w_1 x_1 + w_0 x_0$

+1

-1

s

$$\mathrm{w}^{(\text{new})} = \mathrm{w}^{(\text{old})} + \gamma \cdot y_n \cdot \mathrm{x_n}$$

$$w_2^{(\text{new})} = w_2^{(\text{old})} + \gamma \cdot (-1) \cdot (x_2 = 5)$$

$$w_1^{(\text{new})} = w_1^{(\text{old})} + \gamma \cdot (-1) \cdot (x_1 = -5)$$

$$w_0^{(\text{new})} = w_0^{(\text{old})} + \gamma \cdot (-1) \cdot (x_0)$$

Input data

weight

$+5$   $x_2$   $w_2$

$-5$   $x_1$   $w_1$

$x_0=1$   $w_0$

$\sum$   s

$\sigma(\cdot)$   Sign()

output   $\hat{y} = -1$

Label   $y = +1$

$w_2 x_2 + w_1 x_1 + w_0 x_0$

$\theta$   Threshold value

$+1$

$-1$

s

$$\mathrm{W}^{(\text{new})} = \mathrm{W}^{(\text{old})} + \gamma \cdot y_n \cdot \mathrm{X}_n$$

$$w_2^{(\text{new})} = w_2^{(\text{old})} + \gamma \cdot (+1) \cdot (x_2 = 5)$$

$$w_1^{(\text{new})} = w_1^{(\text{old})} + \gamma \cdot (+1) \cdot (x_1 = -5)$$

$$w_0^{(\text{new})} = w_0^{(\text{old})} + \gamma \cdot (+1) \cdot (x_0)$$

11

❑ You will see another type of update rule in a literature, which is same.

$$\mathrm{w}^{(\text{new})} = \mathrm{w}^{(\text{old})} + \gamma \cdot (\hat{y}_n - y_n) \cdot \mathrm{x}_n$$

$$\mathrm{w}^{(\text{new})} = \mathrm{w}^{(\text{old})} + \gamma \cdot y_n \cdot \mathrm{x}_n$$

$$n \in A$$

$$n \in M$$

Where "A" denotes the set of all data points

Where "M" denotes the set of all misclassified patterns

1) Random initialization of parameters

$$\text{random} \sim w = \{w_0, w_1, w_2, \ldots, w_d\}$$

2) Searching misclassified data points using the decision boundary ($y_n$)

$$y_n = w_0 x_0 + w_1 x_1 + w_2 x_2 + \ldots + w_d x_d$$

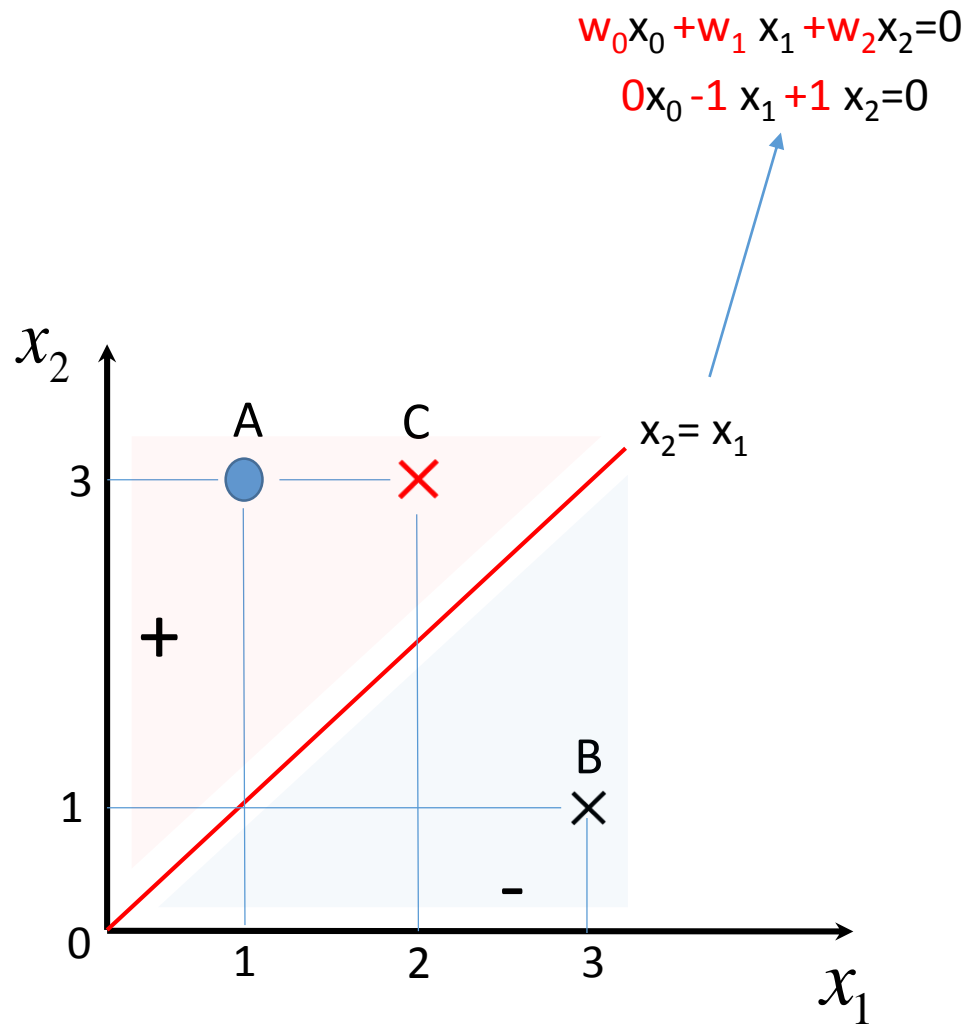2-1) Updating the parameters using the misclassified data points

$$w^{(new)} = w^{(old)} + \gamma \cdot y_n \cdot x_n$$

2-2) Go to step 2)



$x_2$

$y_n = 0$      $y_n > 0$

$y_n < 0$

$x_1$

13

$w_0 x_0 + w_1 x_1 + w_2 x_2 = 0$

$0 x_0 - 1 x_1 + 1 x_2 = 0$
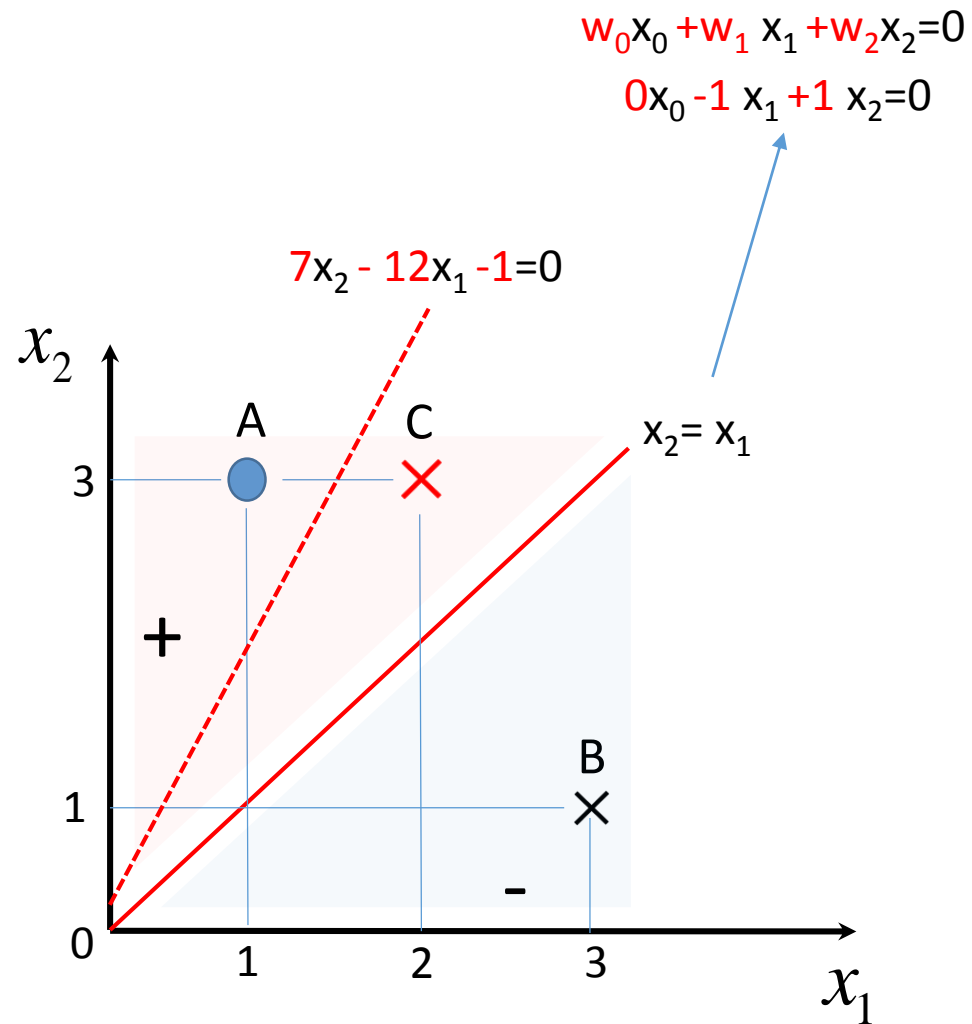
$x_2 = x_1$

A   C

+

B

-

$x_2$

$x_1$

- ❑  Current decision boundary is determined with w=(0, -1, +1)
- ❑  Data point $x_C$ (1, 2, 3) is misclassified
- ❑  Learning rate: $\gamma = 0.1$
- ❑  Let's update the parameter

w=(0, -1, +1)          $x_c$=(1, 2, 3)
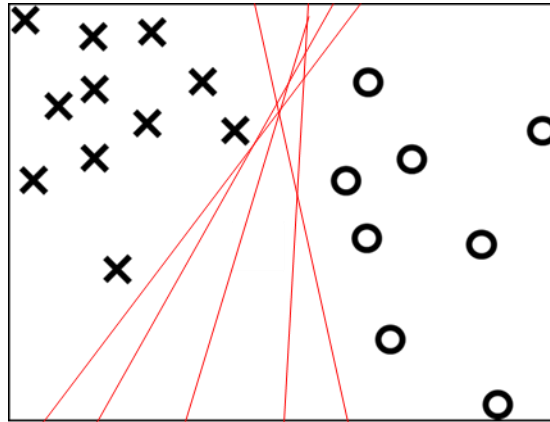
$$\mathrm{w}^{(new)} = \mathrm{w}^{(old)} + \gamma \cdot y_n \cdot \mathrm{x}_n$$

$$= (0, -1, +1) + (0.1)(-1)(1, 2, 3)$$

$$= (-0.1, \quad -1.2, \quad 0.7)$$

$w_0 x_0 + w_1 x_1 + w_2 x_2 = 0$

$0 x_0 - 1 x_1 + 1 x_2 = 0$

$7x_2 - 12x_1 - 1 = 0$

$x_2$

A    C    $x_2 = x_1$

3

+

B

1

−

0    1    2    3

$x_1$

☐ Current decision boundary is determined with w=(0, -1, +1)

☐ Data point $x_C$ (1, 2, 3) is misclassified

☐ Learning rate: $\gamma = 0.1$

☐ Let's update the parameter

w=(0, -1, +1)        $x_c$=(1, 2, 3)

$$\mathrm{w}^{(new)} = \mathrm{w}^{(old)} + \gamma \cdot y_n \cdot \mathrm{x}_n$$

$$= (0, -1, +1) + (0.1)(-1)(1, 2, 3)$$

$$= (-0.1, \quad -1.2, \quad 0.7)$$

☐ New decision boundary is

$7x_2 - 12x_1 - 1 = 0$

❑ The perceptron algorithm is guaranteed to find an exact solution within a finite number of iteration if given data set is linearly separable.

- Slow convergence: cannot tell its feasibility until it's convergence

- Does not converge if there is not any solution

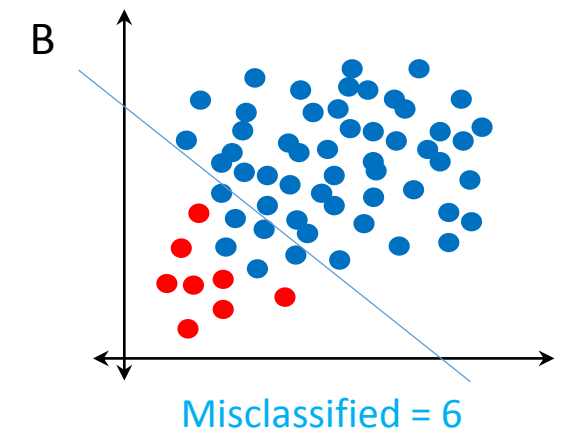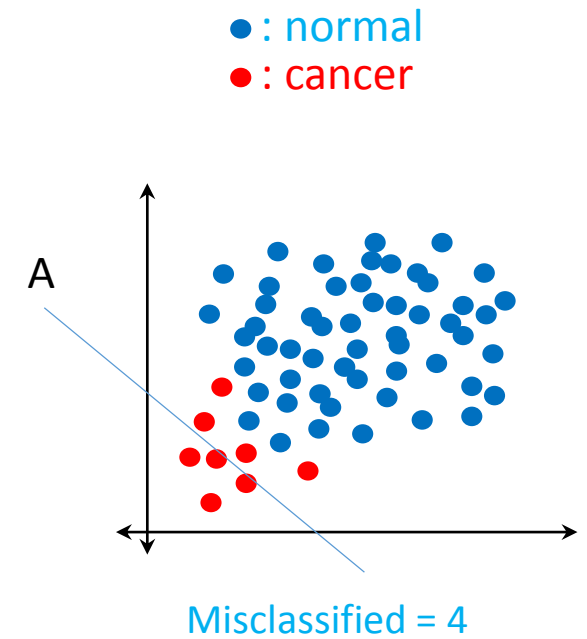- Many solutions exist: converge to one depending on an initial and the order of data feeding

A   B   C   D

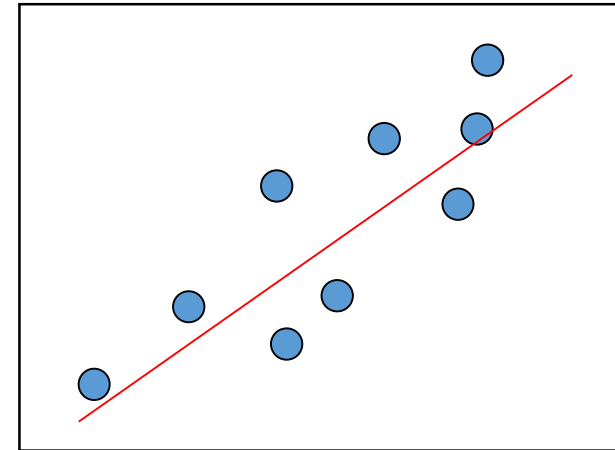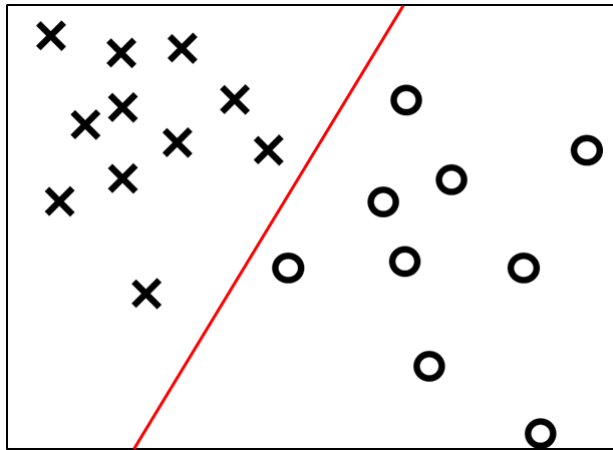| # of binary classifications to be carried out (when there is K classes) | One to One | One to rest |
|---|---|---|
| # of binary classifications to be carried out (when there is K classes) | ❏ K(K-1)/2<br>- A to B<br>- A to C<br>- A to D<br>- B to C<br>- B to D<br>- C to D | ❏ K<br>- A to B/C/D<br>- B to A/C/D<br>- C to A/B/D<br>- D to A/B/C |
| Feature | ❏ Complexity high | ❏ Class imbalance problem |

● : normal
● : cancer
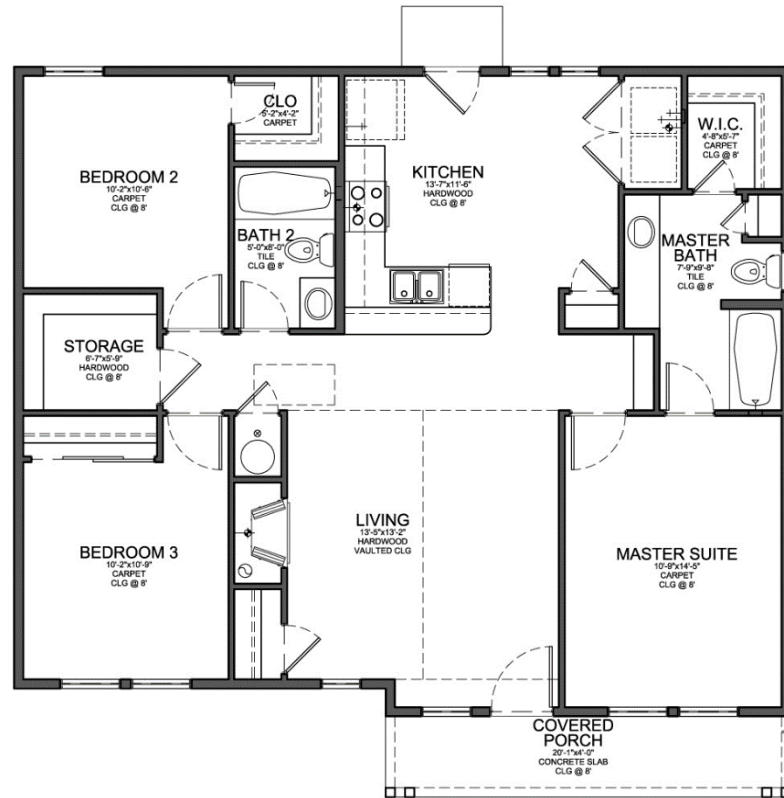
A

Misclassified = 4

B

Misclassified = 6

# Regression

# Regression

❑ Given data set, regression fits them to a certain function

❑ Classification vs Regression

- Classification: given a data, its output is a discretized label

- Regression: given a data, its output is a continuous value
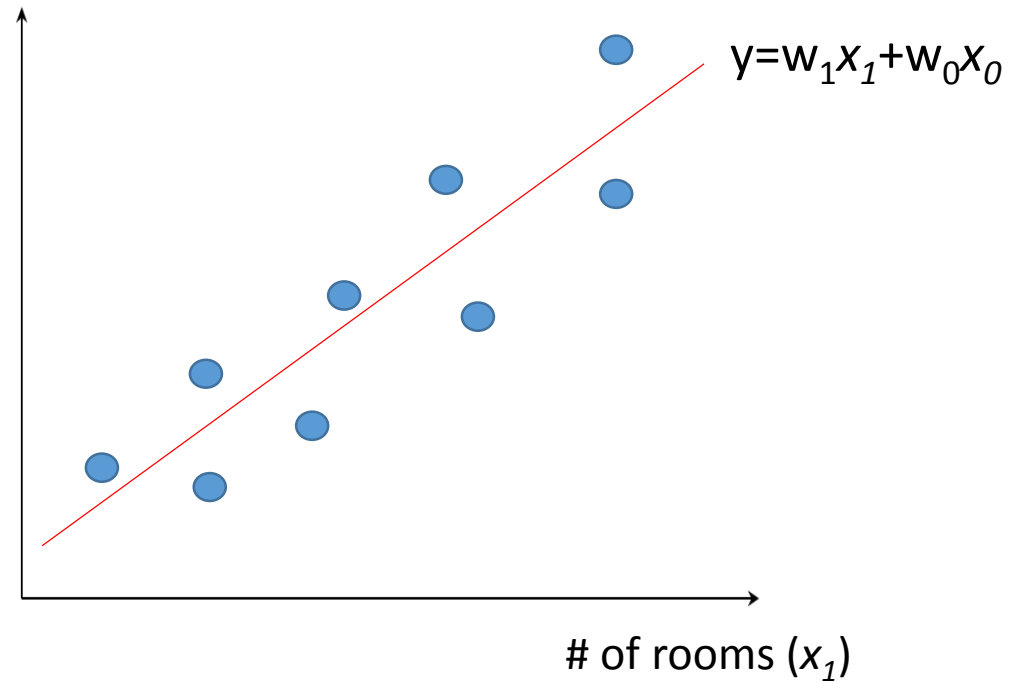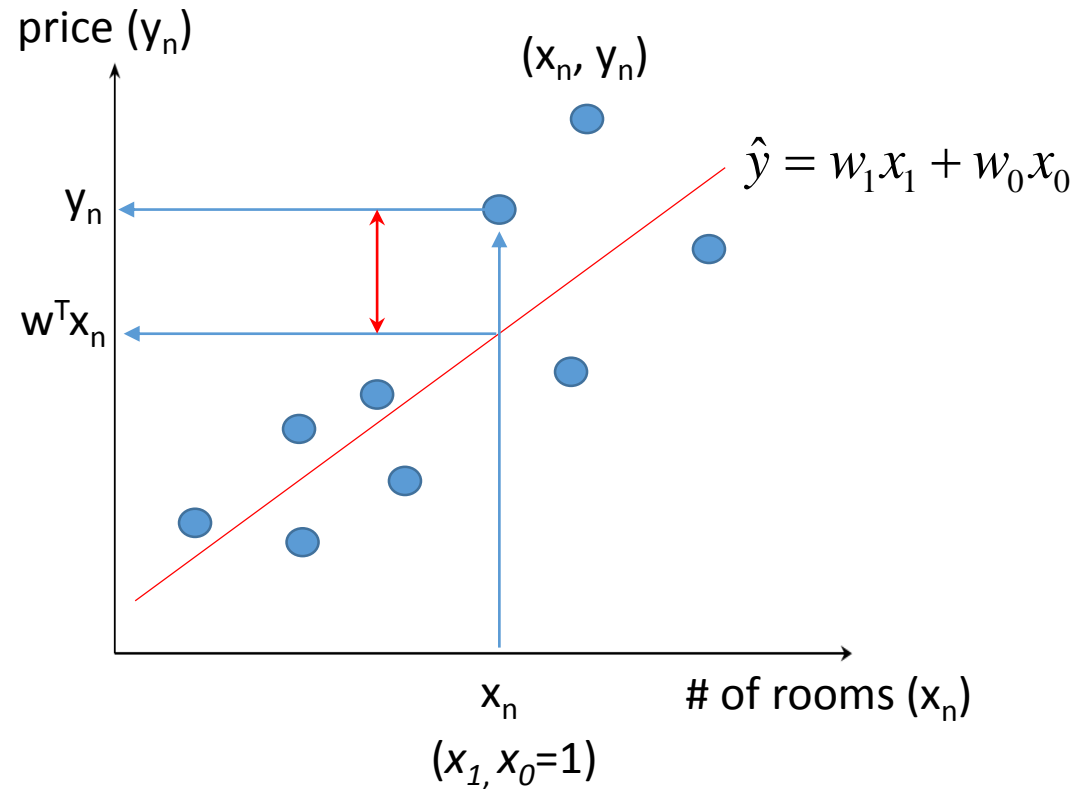
- Both belongs to supervised learning: input + label

❑ How much is the house below?



Price (y)

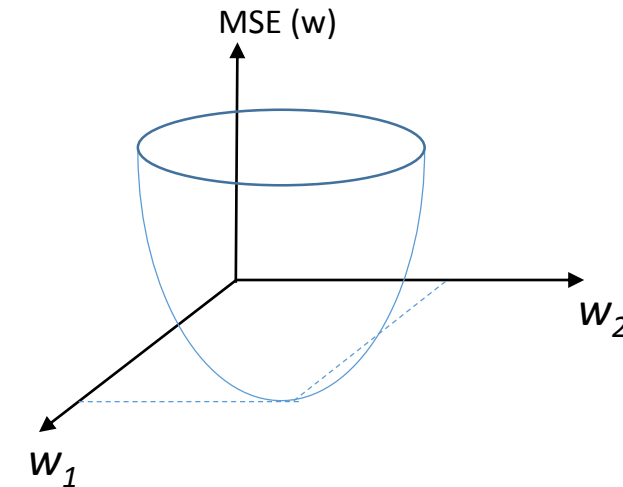$y=w_1x_1+w_0x_0$

# of rooms ($x_1$)

$$\mathrm{MSE}(\mathrm{w}) = \frac{1}{2} \sum_{n=1}^{N} (\mathrm{w}^{\mathrm{T}} \mathrm{x}_n - \mathrm{y}_n)^2$$

price ($y_n$)

($x_n$, $y_n$)

$\hat{y} = w_1 x_1 + w_0 x_0$

$y_n$

$\mathrm{w}^{\mathrm{T}}\mathrm{x}_n$

$x_n$

# of rooms ($x_n$)

($x_1$, $x_0 = 1$)

MSE (w)

$w_2$

$w_1$

1) Analytical approach: normal equations
2) Systematical approach: gradient descent

21

$$MSE(w) = \frac{1}{2}(XW - Y)^T (XW - Y)$$

(1x1) = [(nxm)(mx1)-(nx1)]$^{-T}$ [(nxm)(mx1)-(nx1)]

$$= \frac{1}{2}((XW)^T - Y^T)(XW - Y)$$

$$= \frac{1}{2}(W^T X^T - Y^T)(XW - Y)$$

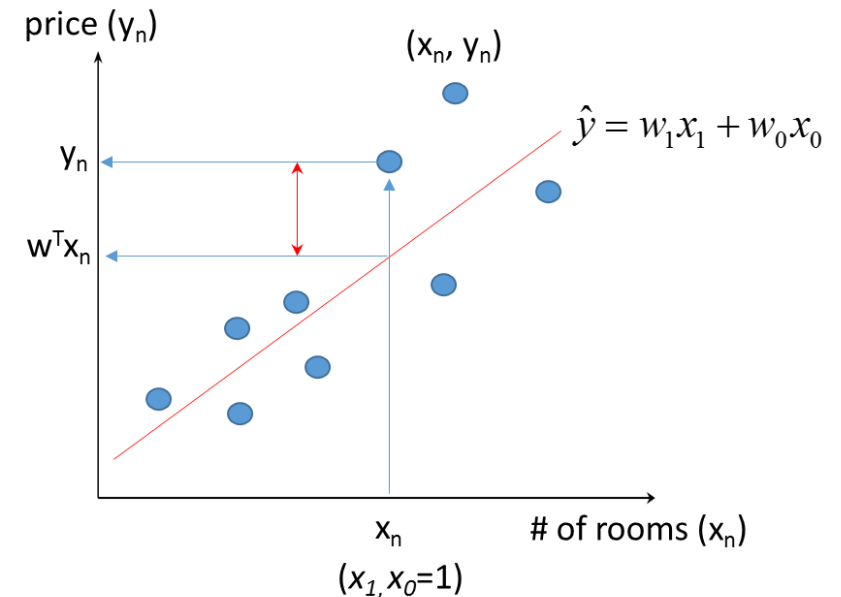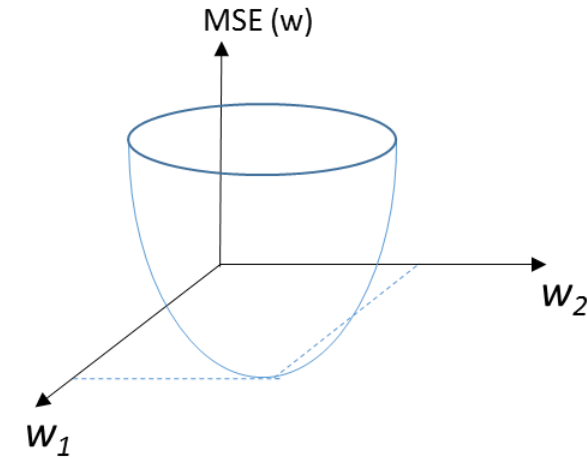$$= \frac{1}{2}(W^T X^T XW - Y^T XW - W^T X^T Y + Y^T Y)$$

$$Y^T XW = (W^T X^T Y)^T$$

$$= \frac{1}{2}(W^T X^T XW - 2W^T X^T Y + Y^T Y)$$

$$\frac{dE}{dW^T} = \frac{1}{2}(2X^T XW - 2X^T Y) = 0$$

$$X^T XW = X^T Y \qquad W = (X^T X)^{-1} X^T Y$$

(mx1)   [(mxn)(nxm)]$^{-1}$(mxn)(nx1)

MSE (w)

$w_2$

$w_1$

price ($y_n$)

($x_n$, $y_n$)

$\hat{y} = w_1 x_1 + w_0 x_0$

$y_n$

$w^T x_n$

$x_n$                    # of rooms ($x_n$)

($x_1$, $x_0$=1)

22

$$W = (X^T X)^{-1} X^T Y$$

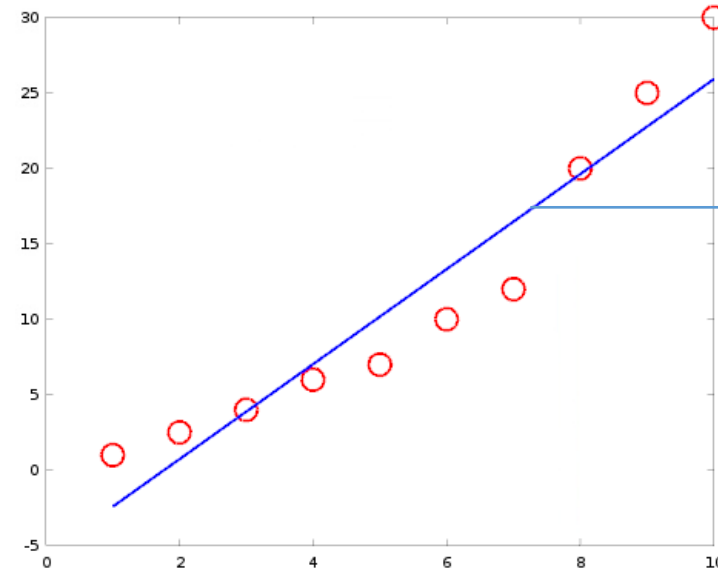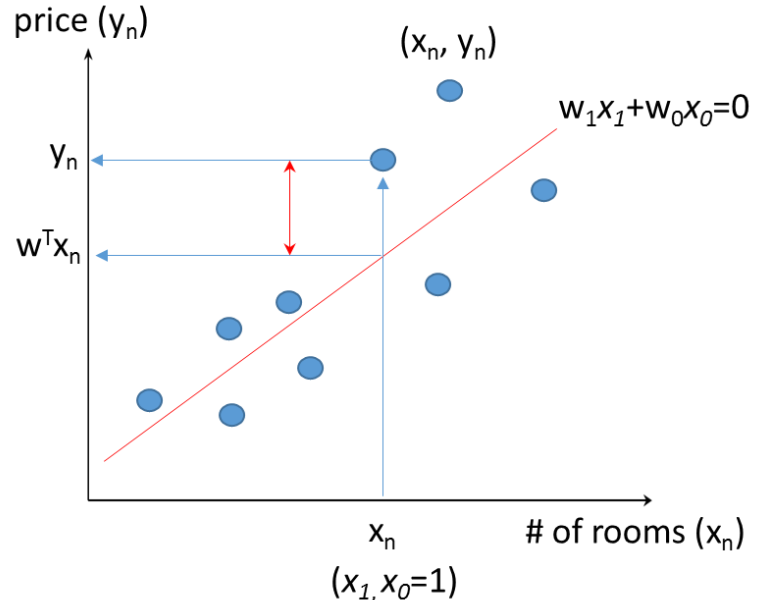(mx1)    [(mxn)(nxm)]$^{-1}$(mxn)(nx1)

```
>> xt
xt =

    1   2   3   4   5   6   7   8   9   10
    1   1   1   1   1   1   1   1   1   1

>> yt
yt =

    1.0  2.5  4.0   6.0   7.0  10.0  12.0   20.0  25.0  30.0

>> w = inverse(xt*x)*xt*y
w =

    3.1485
   -5.5667
```
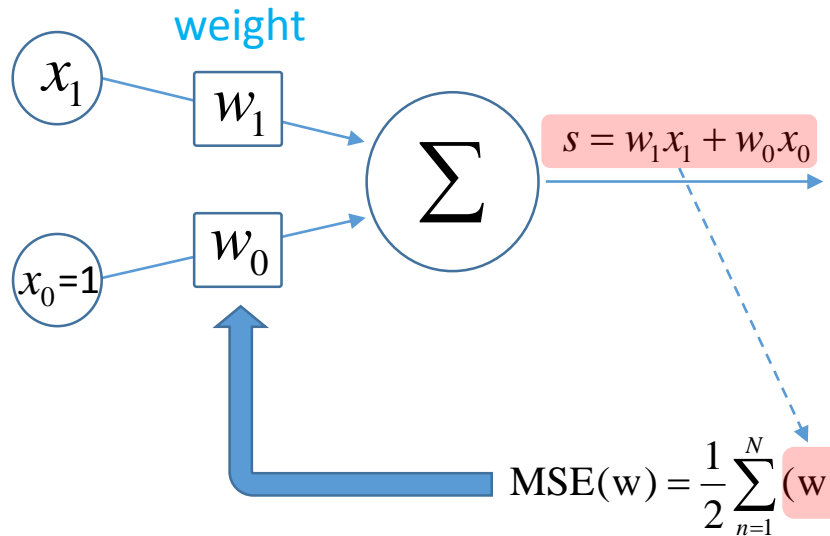
❑ transpose()
❑ Inverse()

price ($y_n$)

($x_n$, $y_n$)

$w_1 x_1 + w_0 x_0 = 0$

$y_n$

$w^T x_n$

$x_n$

# of rooms ($x_n$)

($x_{1,} x_0 = 1$)

3.1485 $x_1$ -5.5667

23

Input data

weight

Label: $y_n$

```
>> x
x =

    1      1
    2      1
    3      1
    4      1
    5      1
    6      1
    7      1
    8      1
    9      1
   10      1
```

$x_1$

$w_1$

$\Sigma$

$s = w_1 x_1 + w_0 x_0$

$x_0 = 1$

$w_0$

```
>> transpose(y)
ans =

    1.0000
    2.5000
    4.0000
    6.0000
    7.0000
   10.0000
   12.0000
   20.0000
   25.0000
   30.0000
```

MSE (w)

$$\text{MSE}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} (\mathbf{w}^{\text{T}} \mathbf{x}_n - \mathbf{y}_n)^2$$
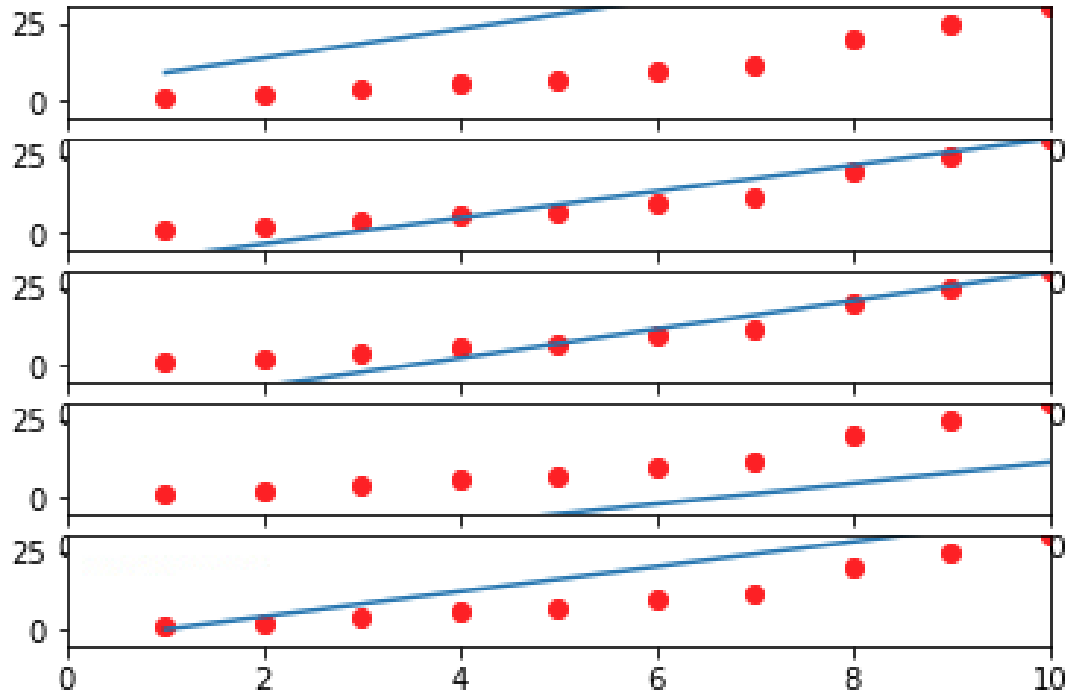
$w_2$

$w_1$

$$w_1 = w_1 + \gamma \frac{\partial E(\mathbf{w})}{\partial w_1}$$

How to obtain?

backpropagation

$$w_1 = w_1 + \gamma \frac{\partial E(\mathbf{w})}{\partial w_1} = w_1 + \gamma \left[ \frac{\partial E(\mathbf{w})}{\partial s} \frac{\partial s}{\partial w_1} \right]$$
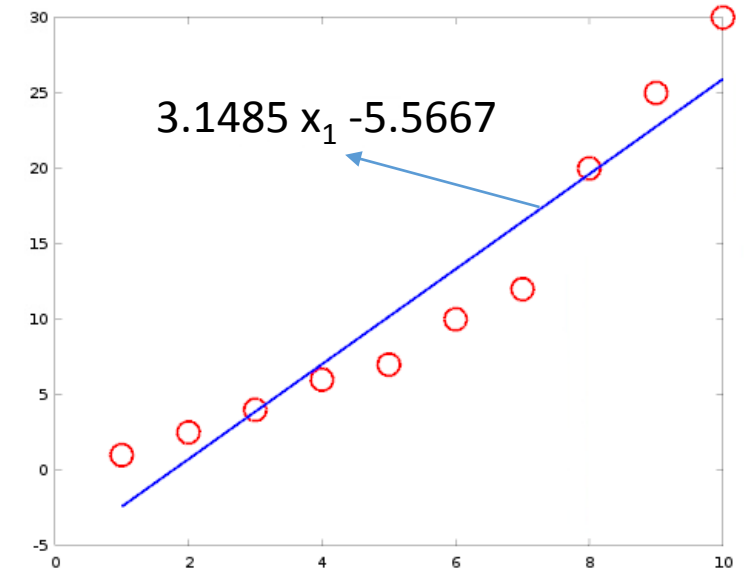
$$w_2 = w_2 + \gamma \frac{\partial E(\mathbf{w})}{\partial w_2}$$

24

iteration
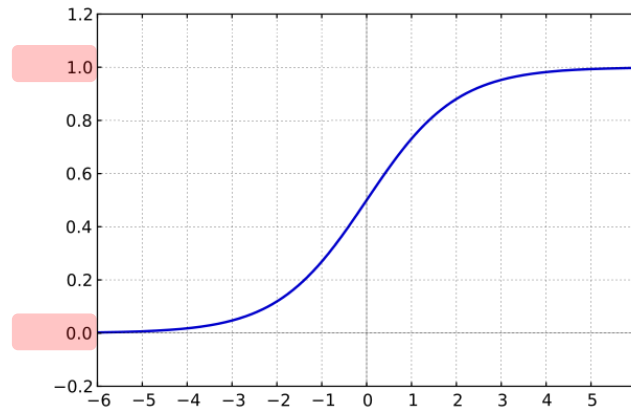
Solution from
gradient descent

$3.1485\ x_1 - 5.5667$

Solution from
normal equation

# Logistic regression

❑ A logistic function is used to fit the given binary data set.
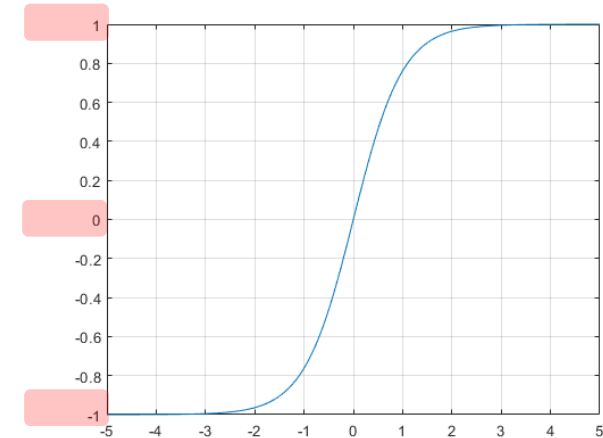
❑ A logistic function is a common "S" shape function.

$$\sigma(x) = \frac{L}{1 + e^{-a(x+b)}}$$

- ▪ L: maximum value of the function
- ▪ a: steepness of the curve
- ▪ b: location of the midpoint



sigmoid function



Hyper tangent function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

❑ Is it a bicycle or motorbike?
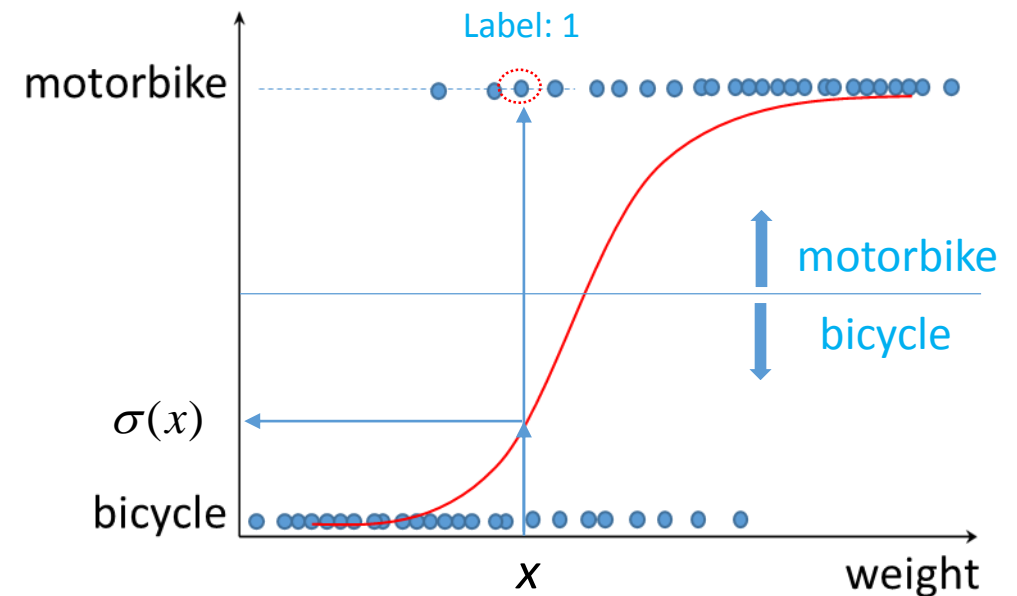
❑ **How to calculate the error given below?**

- Label is one or zero
- Outcome σ(x) is a floating number between [0, 1]

$$\sigma(x) = \frac{1}{1 + e^{-a(x+b)}}$$

Label : $y = (1, 0)$

prediction : $\hat{y} = (\sigma(x), 1 - \sigma(x))$

$$CEE(y, \hat{y}) = -\sum y \log \hat{y}$$



28

$$\text{CEE}(y, \hat{y}) = -\sum y \log \hat{y}$$

prediction : $\hat{y} = (\sigma(x), 1 - \sigma(x))$

Label : $y = (1, 0)$

When the prediction is correct, how correct is it?

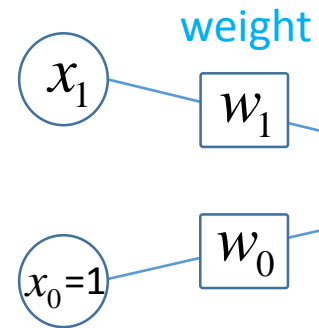| Prediction | Label | Cross Entropy (error) | MSE |
|---|---|---|---|
| 0.1, 0.2, 0.7 | 0, 0, 1 | -ln(0.1)*0-ln(0.2)*0-ln(0.7)*1 = **0.357** | $(0.1-0)^2 +(0.2-0)^2 +(0.7-1)^2$ = **0.14** |
| 0.3, 0.3, 0.4 | 0, 0, 1 | -ln(0.3)*0-ln(0.3)*0-ln(0.4)*1 = 0.916 | $(0.3-0)^2 +(0.3-0)^2 +(0.4-1)^2$ = 0.54 |

When the prediction is wrong, how wrong is it?

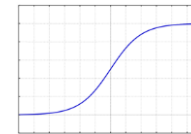| Prediction | Label | Cross Entropy (error) | MSE |
|---|---|---|---|
| 0.1, 0.2, 0.7 | 1, 0, 0 | -ln(0.1)*1-ln(0.2)*0-ln(0.7)*0 = **2.303** | $(0.1-1)^2 +(0.2-0)^2 +(0.7-0)^2$ = **1.34** |
| 0.3, 0.3, 0.4 | 1, 0, 0 | -ln(0.3)*1-ln(0.3)*0-ln(0.4)*0 = 1.204 | $(0.3-1)^2 +(0.3-0)^2 +(0.4-0)^2$ = 0.74 |

$$\sigma(x) = \frac{1}{1+e^{-s}} = \frac{1}{1+e^{-w_1 x_1 - w_0 x_0}}$$

Input data

```
>> x
x =

      1    1
      2    1
      3    1
      4    1
      5    1
      6    1
      7    1
      8    1
      9    1
     10    1
```

weight

$x_1$

$w_1$

$x_0 = 1$

$w_0$

$\sum$

$s = w_1 x_1 + w_0 x_0$

$\sigma(\cdot)$

Sigmoid function

$\hat{y}$

Label: $y_n$

```
>> transpose(y)
ans =

      0
      0
      0
      0
      1
      1
      1
      1
      1
      1
```

$$\text{CEE}(y, \hat{y}) = -\sum y \log \hat{y}$$

Same story
backpropagation

$$w_1 = w_1 + \gamma \frac{\partial E(\text{w})}{\partial w_1} = w_1 + \gamma \left[ \frac{\partial E(\text{w})}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s} \frac{\partial s}{\partial w_1} \right]$$

30