



# **SUPERVISED LEARNING AND PERCEPTRON**

**INT301 Bio-computation, Week 2, 2025**



# Outline

- McCulloch-Pitts neuron
- Hebb's learning rule 只能说明相关性 a input with  $x^t$ , 但我们希望 model
- **Supervised learning model:**  
**Perceptron** 能找到因果性如 a 现象 / 特征  $\rightarrow$  这是什么类别

# Recall: Machine learning and ANN



- Like human learning from past experiences.
- A computer does not have “experiences”.
- A computer system learns from data, which represent some “past experiences” of an application domain.
- **Our focus:** learn a target function that can be used to predict the values of a discrete class attribute, e.g., yes or no, and high or low.
- The task is commonly called: supervised learning.

# The data and the goal

- **Data:** A set of data records (also called examples, instances or cases) described by
  - $k$  attributes:  $A_1, A_2, \dots, A_k$
  - a class: Each example is labelled with a pre-defined class.
- **Goal:** To learn a classification model from the data that can be used to predict the classes of new (future, or test) cases/instances.



# An example: data (loan application)

attributes

class  
Approve or not

ID	Age	Has Job	Own_House	Credit_Rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

# An example: the learning task

- **Learn a classification model** from the data
- Use the model to classify future loan applications into
  - **Yes (approve)** and
  - **No (disapprove)**
- What is the class for following case/instance?

Age	Has_Job	Own_house	Credit-Rating	Class
young	false	false	good	?

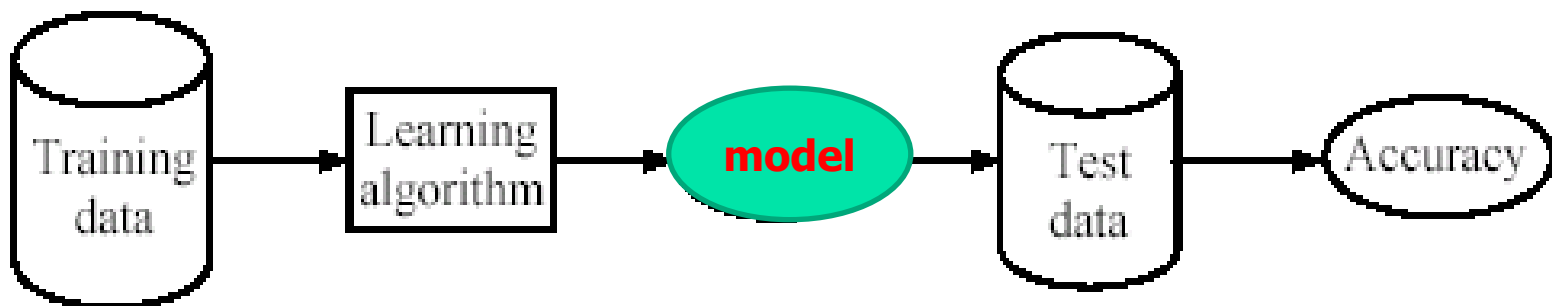
# Supervised vs. unsupervised Learning

- **Supervised learning:** classification is seen as supervised learning from examples.
  - **Supervision:** The data (observations, measurements, etc.) are labeled with pre-defined classes. It is like that a “teacher” gives the classes (supervision).
  - Test data are classified into these classes too.
- **Unsupervised learning (e.g. clustering)**
  - Class labels of the data are unknown
  - Given a set of data, the task is to establish the existence of classes or clusters in the data

# Supervised learning process: two steps

- **Learning (training)**: Learn a model using the training data
- **Testing**: Test the model using **unseen** test data to assess the model accuracy

$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}};$$



Step 1: Training      Step 2: Testing



# What do we mean by learning?



## ■ Given

- a data set  $D$
- a task  $T$
- a performance measure  $M$

a computer system is said to **learn** from  $D$  to perform the task  $T$  if after learning the system's performance on  $T$  improves as measured by  $M$ .

- In other words, the learned model helps the system to perform  $T$  better as compared to no learning.

# Fundamental assumption of learning

学到/验证的规律一致

↑  
测试和训练的样本分布相同

**Assumption:** The distribution of training examples is **identical** to the distribution of test examples (including future unseen examples).

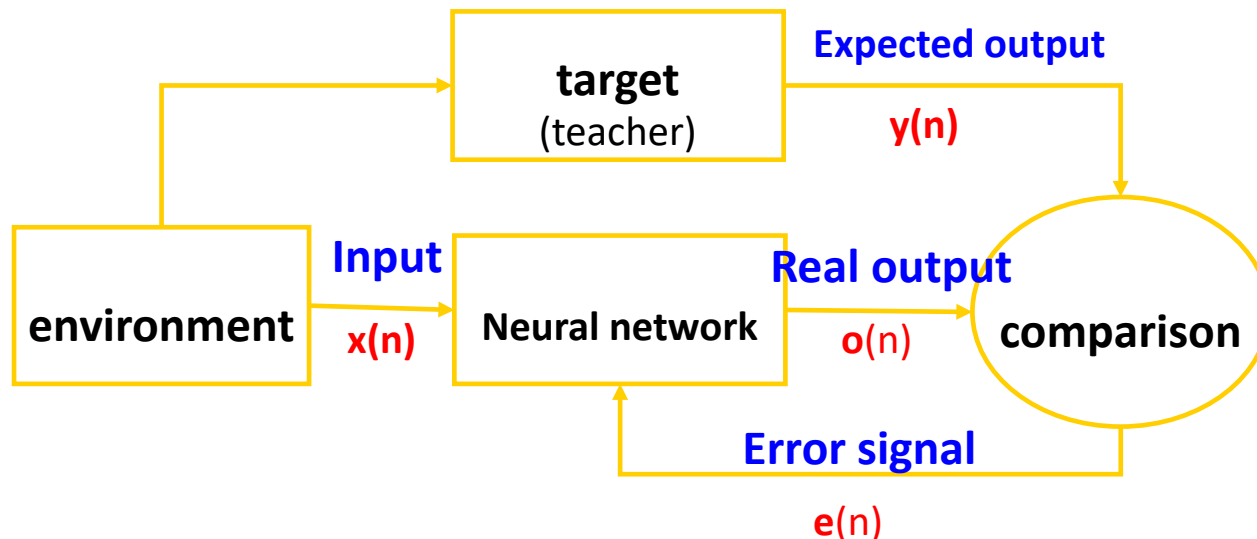
- In practice, this assumption is often violated to certain degree.
- Strong violations will clearly result in poor classification accuracy.
- To achieve good accuracy on the test data, training examples must be sufficiently representative of the test data.

# Perceptron (1958) 感知器

- Rosenblatt (1958) explicitly considered the problem of **pattern recognition**, where a “teacher” is essential.
- Perceptrons are neural networks that change with “experience” using *error-correcting rule*.  
→ 纠错规则
- According to the rule, weight of a response unit changes when it makes erroneous response to stimuli presented to the network.

# ANN for Pattern Recognition

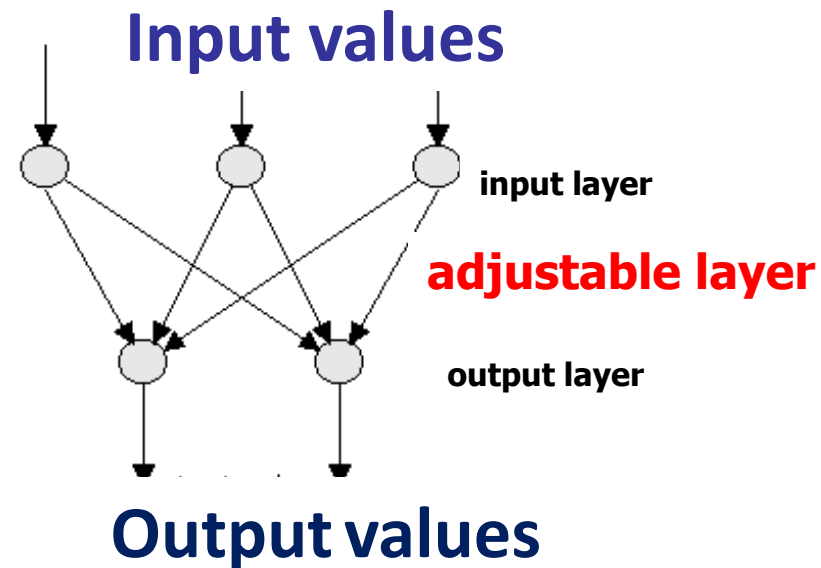
- ◆ **Training data**: set of sample pairs  $(x, y)$ .
- ◆ Network (model, classifier) **adjusts its connection weights according to the errors** between target and network output



# Perceptron (1958)

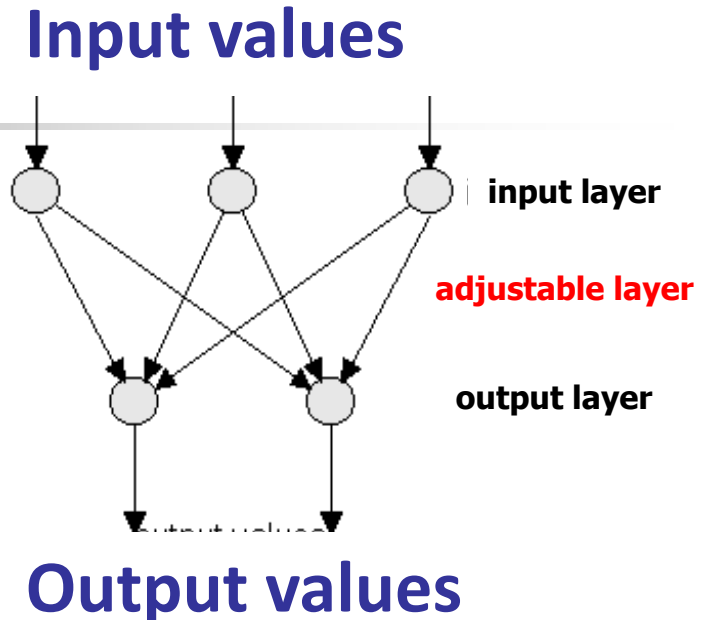
网络单元：感知器最简单结构由两层理想化神经元组成

- The simplest architecture of perceptron comprises two layers of idealised "neurons", which we shall call *"units" of the network*.
- There are
  - one layer of input units, and
  - one layer of output units.in the perceptron



# Perceptron (1958)

- The two layers are fully interconnected, *i.e.*, every input unit is connected to every output unit
- Thus, *processing elements* of the perceptron are the *abstract neurons*
- Each processing element has the same input comprising total input layer, but individual outputs with individual connections and therefore different weights of connections.





# Perceptron (1958)

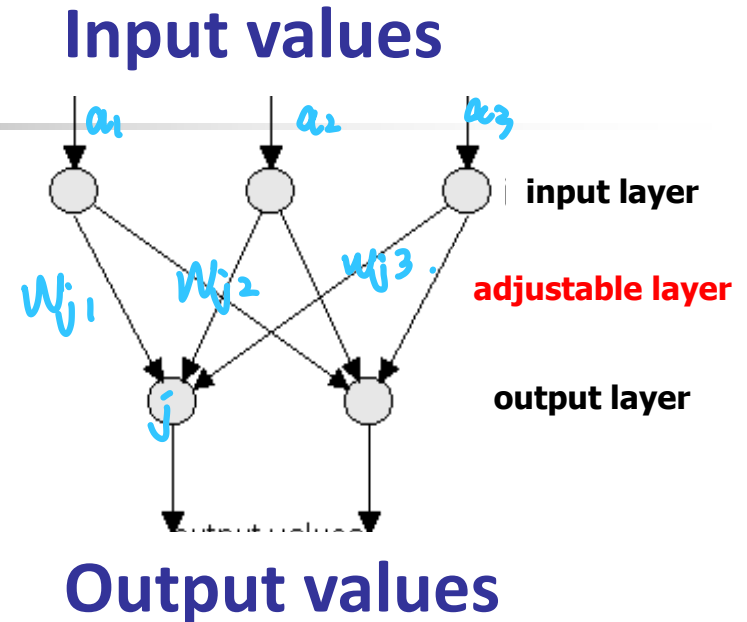
The total input to the output unit  $j$  is

$$S_j = \sum_{i=0}^n w_{ji} a_i$$

$a_i$  : input value from the  $i$ th input unit

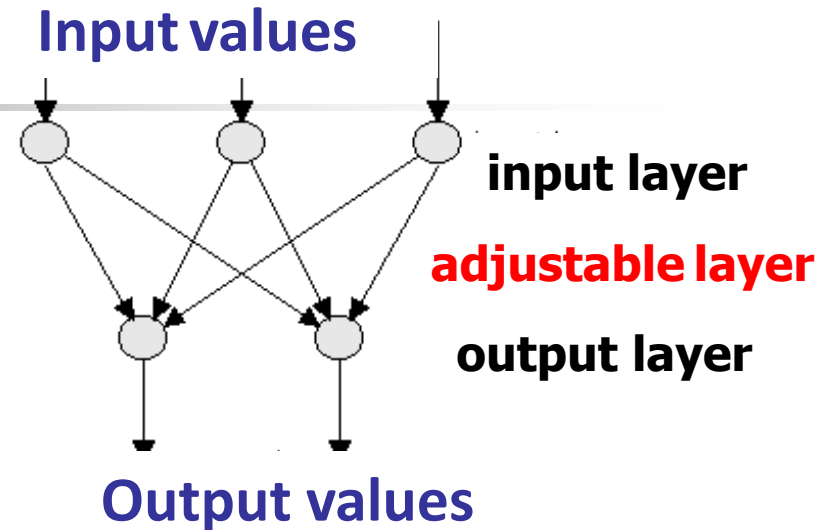
$w_{ji}$  : the weight of connection btw  $i$ -th input and  $j$ -th output units

- The sum is taken **over all  $n+1$  inputs** units connected to the output unit  $j$ .
- There is special **bias input** unit **number 0** in the input layer.



# Perceptron (1958)

$$S_j = \sum_{i=0}^n w_{ji} a_i$$



- There is a special **bias input** unit **number 0** in the input layer.
- The bias unit always produces inputs  $a_0$  of the fixed values of +1.
- The input  $a_0$  of **bias unit** functions as a constant value in the sum.
- The **bias unit connection** to output unit  $j$  has a weight  $w_{j0}$  adjusted in the same way as all the other weights

# Perceptron (1958)

- The output value  $X_j$  of the output unit  $j$  depends on whether the weighted sum is above or below the unit's threshold value.
- $X_j$  is defined by the unit's threshold activation function.

$$X_j = f(S_j) = \begin{cases} 1, S_j \geq \theta_j \\ 0, S_j < \theta_j \end{cases}$$

## Definition:

the ordered set of instant outputs of all units in the output

layer  $X = \{X_0, X_1, \dots, X_n\}$  ↪ 输出的是一个向量

constitutes an **output vector** of the network



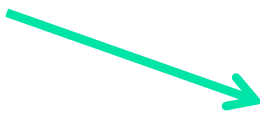

# Perceptron (1958)

---

- The instant output  $X_j$  of the j-th unit in the output layer constitutes the j-th component of the output vector.
- Weight  $w_{ji}$  of connections between the two layers are changed according to **perceptron learning rule**, so the network is more likely to produce the desired output in response to certain inputs.
- The process of weights adjustment is called **perceptron learning (or training)**.

# Perceptron Training

Every processing element computes an output according its state and threshold:

$$S_j = \sum_{i=0}^n w_{ji} a_i$$

$$X_j = f(S_j) = \begin{cases} 1, S_j \geq \theta_j \\ 0, S_j < \theta_j \end{cases}$$


The network instant outputs  $X_j$  are then compared to the desired outputs specified in the training set.

$$e_j = (t_j - X_j)$$

The error of an output unit is the difference between the target output and the instant one

# Perceptron Training

**The error are computed and used to re-adjust the values of the weights of connections.**

$$S_j = \sum_{i=0}^n w_{ji} a_i$$

$$X_j = f(S_j) = \begin{cases} 1, S_j \geq \theta_j \\ 0, S_j < \theta_j \end{cases}$$

$$e_j = (t_j - X_j)$$

The weights re-adjustment is done in such a way that the network is – on the whole – more likely to give the desired response next time.



# Perceptron Updating of the Weights

The goal of the training session is to arrive at a single set of weights that allow each of the mappings in the training set to be done successfully by the network.

## 1. Compute *error* of every output unit

$$e_j = (t_j - X_j)$$

where

$t_j$  is the target value for output unit  $j$

$X_j$  is the instant output produced by output unit  $j$

# Perceptron Updating of the Weights

Having the errors computed,

## 2. Update the weights

$$\underset{\text{new}}{w_{ji}} = \underset{\text{old}}{w_{ji}} + \Delta w_{ji}$$

权重按数值而非向量计算

where

$$\Delta w_{ji} = C e_j a_i = C(t_j - X_j) a_i$$

**Perceptron  
learning rule**

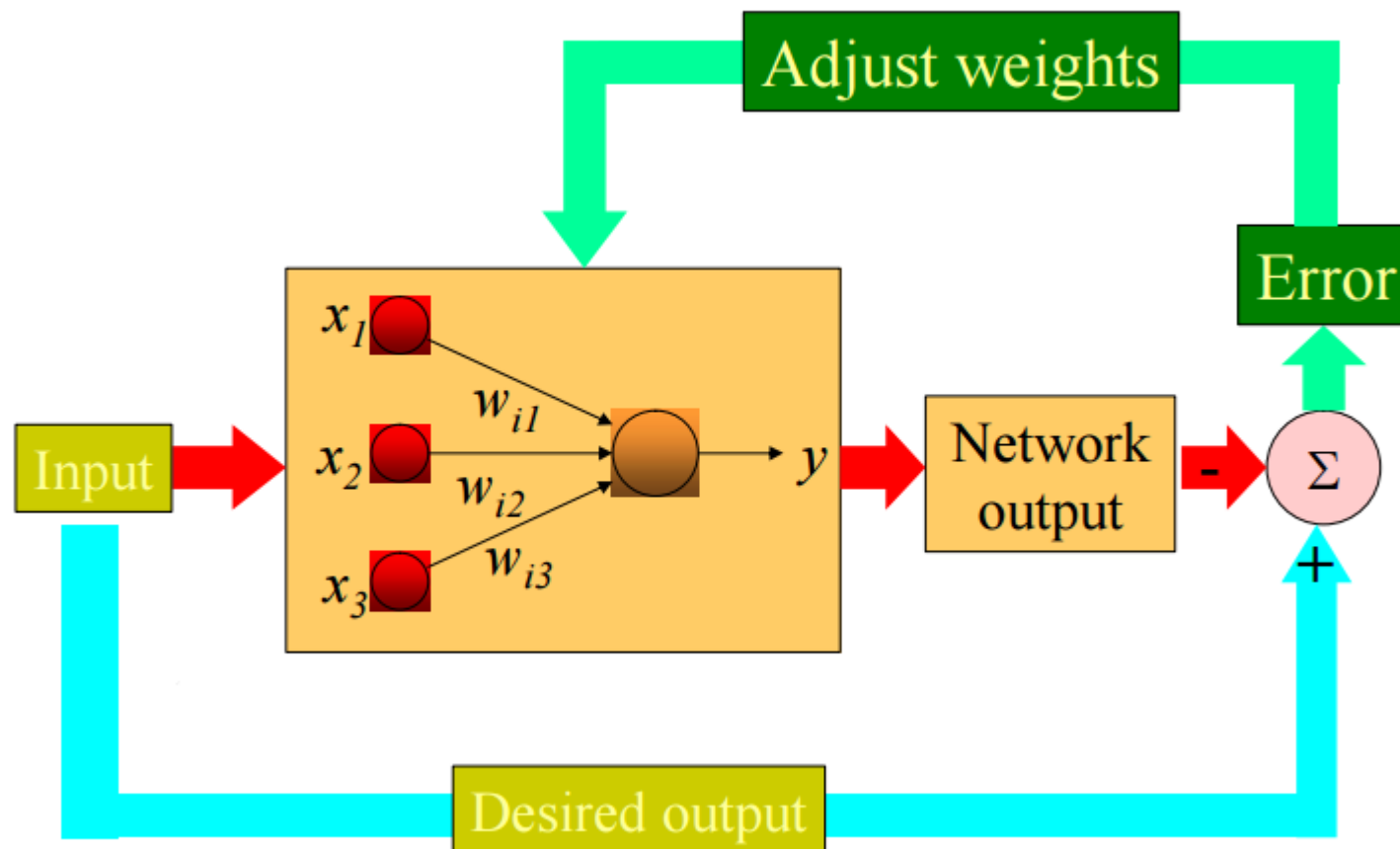
# Perceptron training

- A sequential learning procedure for updating the weights.
- Perceptron training algorithm (**delta rule**)


$$\Delta w = \text{learning rate} \times (\text{teacher} - \text{output}) \times \text{input}$$

**error**

# Perceptron



# Example



- Define our “features”:

Taste	Sweet = 1, Not_Sweet = 0
Seeds	Edible = 1, Not_Edible = 0
Skin	Edible = 1, Not_Edible = 0

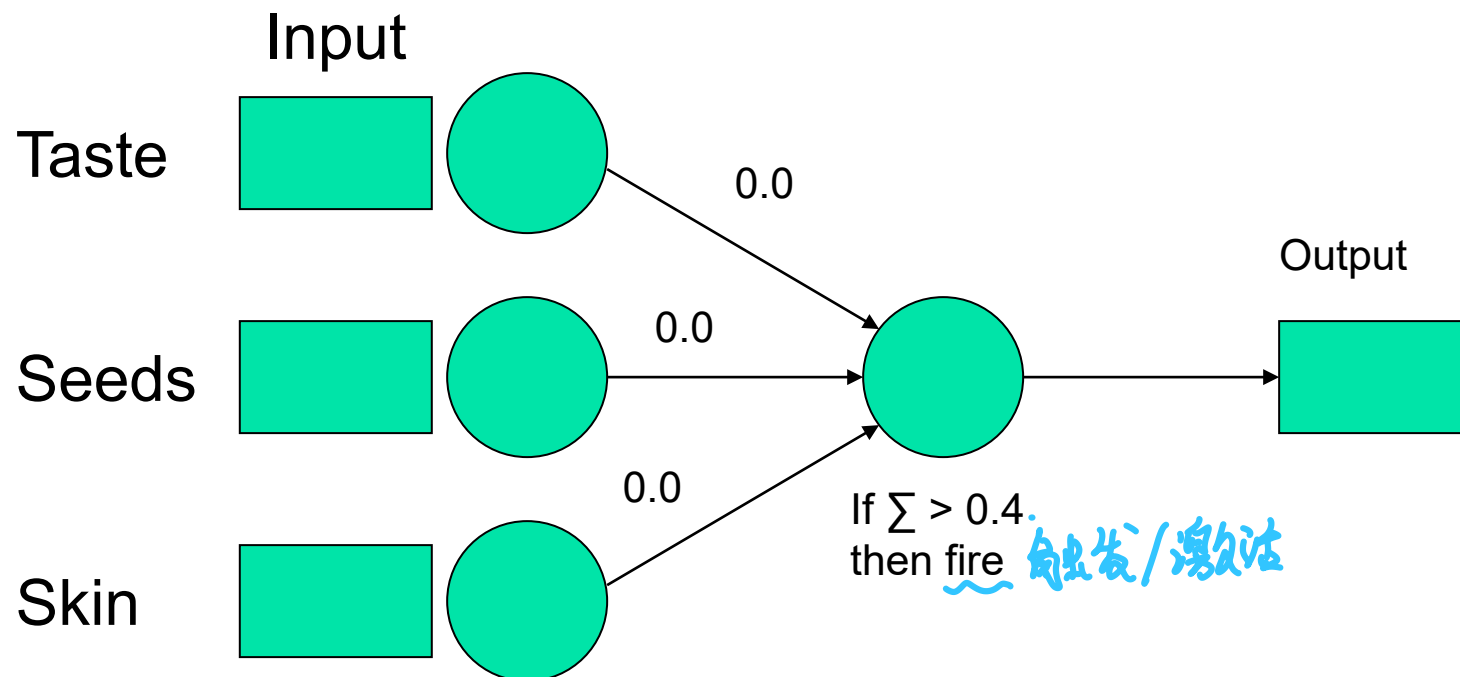
For output:

`Good_Fruit = 1`

`Not_Good_Fruit = 0`

# Example

Let's start with no knowledge:





# Example



- To train the perceptron, we will show it each example and have it categorize each one.
- Since it's starting with no knowledge, it is going to make mistakes. When it makes a mistake, we are going to adjust the weights to make that mistake less likely in the future.

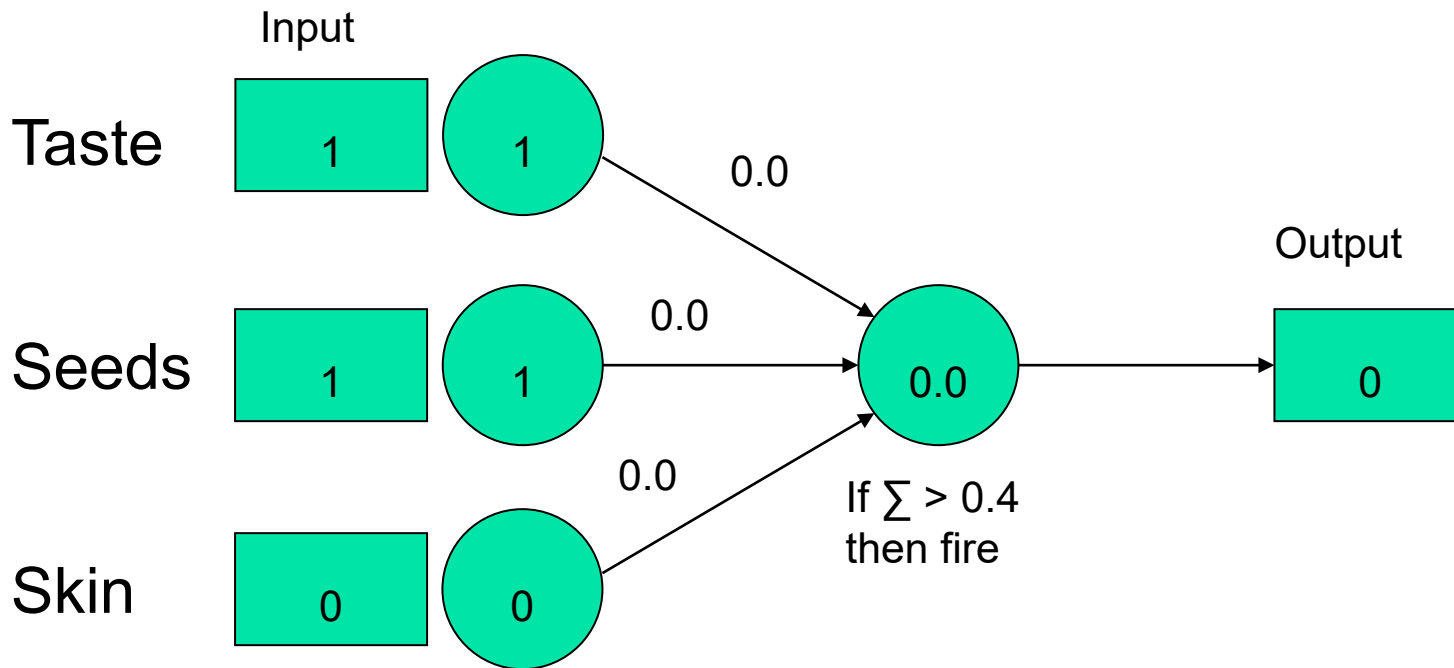
# Example



- When we adjust the weights, we're going to take relatively small steps to be sure we don't over-correct and create new problems.
- We're going to learn the category "good fruit" defined as anything that is sweet.
  - `Good fruit = 1`
  - `Not good fruit = 0`

# Example

Show it a banana:





# Example

---

- In this case we have:  
 $(1 \times 0) = 0$   
 $+ (1 \times 0) = 0$   
 $+ (0 \times 0) = 0$
- It adds up to 0.0.
- Since that is less than the threshold (0.40), we responded “no.”
- Is that correct? No.



# Example

---

- Since we got it wrong, we need to change the weights. We'll do that using the **delta rule** (delta for change).

$$\Delta w = \text{learning rate} \times (\text{teacher} - \text{output}) \times \text{input}$$

# Example

The three parts of that are:

- **Learning rate:** We set that ourselves. Set large enough that learning happens in a reasonable amount of time; and also small enough to avoid too fast. Here pick 0.25.
- **(teacher - output):** The teacher knows the correct answer (e.g., that a banana should be a good fruit). In this case, the teacher says 1, the output is 0, so  $(1 - 0) = 1$ .
- **Input:** That's what came out of the node whose weight we're adjusting. For the first node. 1.





# Example

---

- To pull it together:
  - Learning rate: 0.25.
  - (teacher - output): 1.
  - input: 1.

$$\Delta w = 0.25 \times 1 \times 1 = 0.25.$$

- Since it's a  $\Delta w$ , it's telling us how much to change the first weight. In this case, we're adding 0.25 to it.



# Example

---

Let's think about the delta rule:

(teacher - output)

- If we get the categorization right, (teacher - output) will be zero (the right answer minus itself).
- In other words, if we get it right, we won't change any of the weights. As far as we know we have a good solution, why would we change it?



# Example

---

Let's think about the delta rule:

(teacher - output)

- If we get the categorization wrong, (teacher - output) will either be -1 or +1.
  - If we said “yes” when the answer was “no”, we’re too high on the weights and we will get a (teacher - output) of -1 which will result in reducing the weights.
  - If we said “no” when the answer was “yes”, we’re too low on the weights and this will cause them to be increased.



# Example

---

Let's think about the delta rule:

- Input:
  - If the node whose weight we're adjusting sent in a 0, then it didn't participate in making the decision. In that case, it shouldn't be adjusted. Multiplying by zero will make that happen.
  - If the node whose weight we're adjusting sent in a 1, then it did participate and we should change the weight (up or down as needed) if the corresponding output wrong.

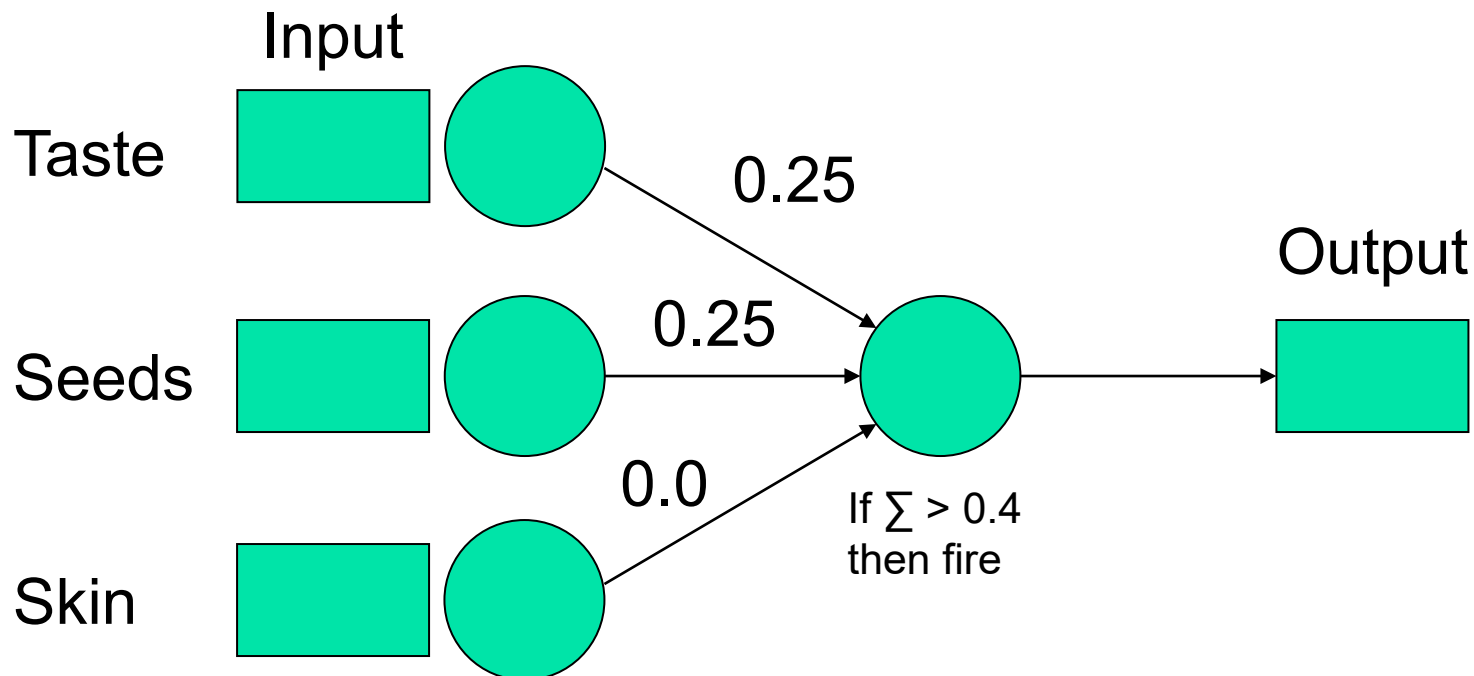
# Example

- How do we change the weights for banana?

Feature:	Learning rate:	(teacher - output):	Input:	$\Delta w$
taste	0.25 $\times$	1 $\times$	1 $=$	+0.25
seeds	0.25	1	1	+0.25
skin	0.25	1	0	0

# Example

Here it is with the adjusted weights:





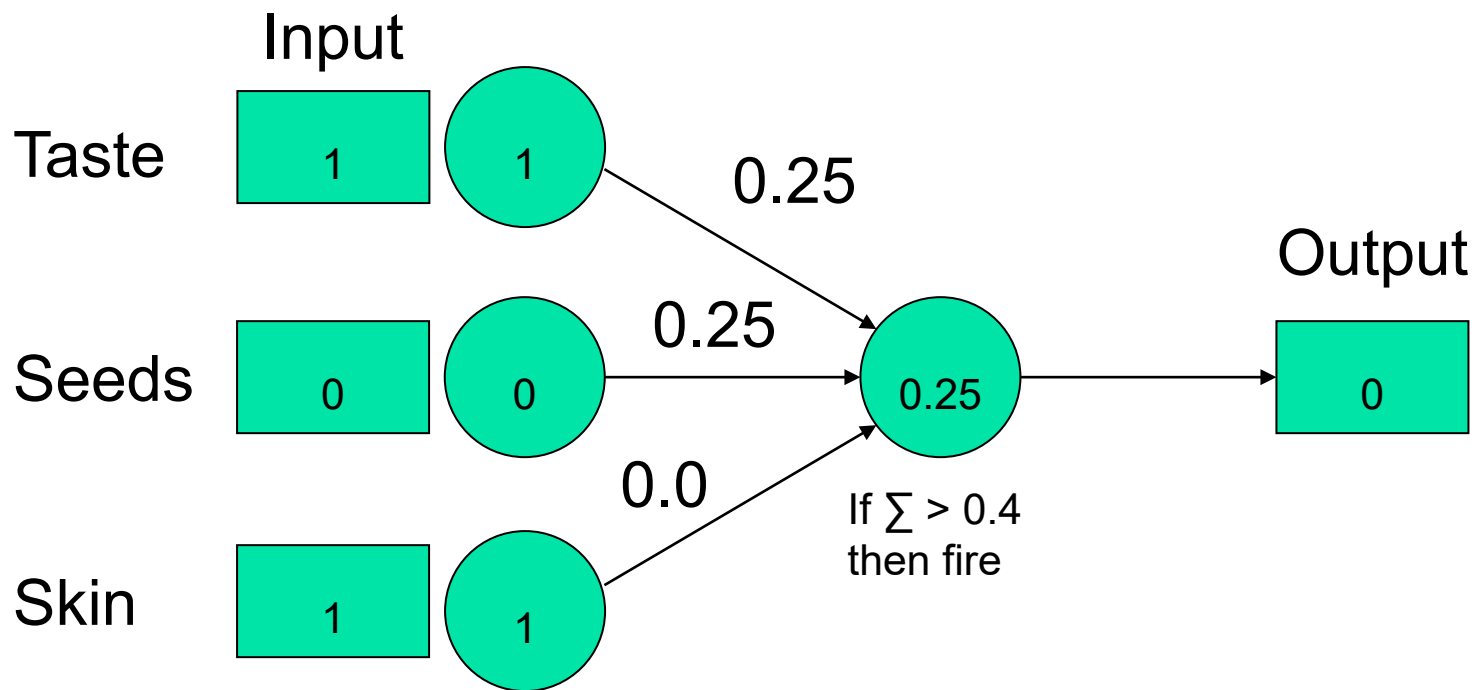
# Example

---

- To continue training, we show it the next example, adjust the weights...
- We will keep cycling through the examples until we go all the way through one time without making any changes to the weights. At that point, the concept is learned.

# Example

Show it a pear:





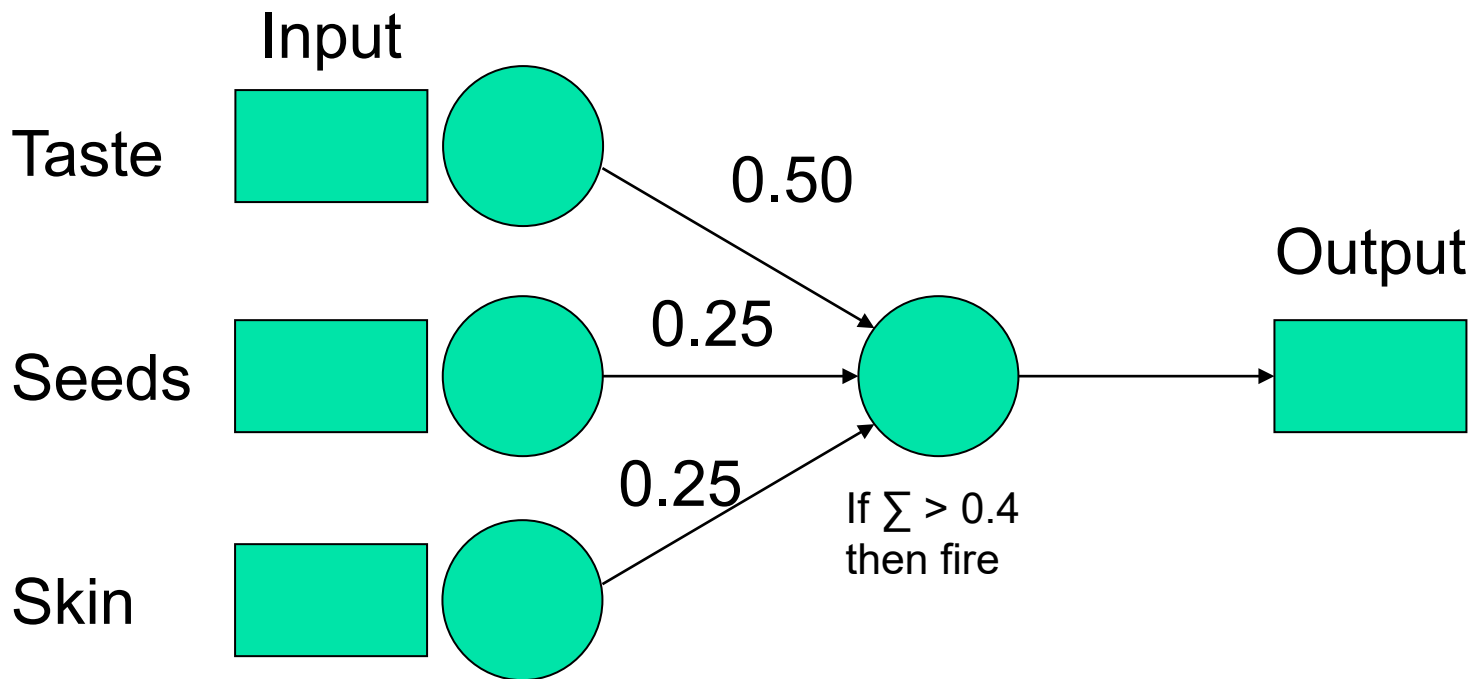
# Example

- How do we change the weights for pear?

Feature:	Learning rate:	(teacher - output):	Input:	$\Delta w$
taste	0.25	1	1	+0.25
seeds	0.25	1	0	0
skin	0.25	1	1	+0.25

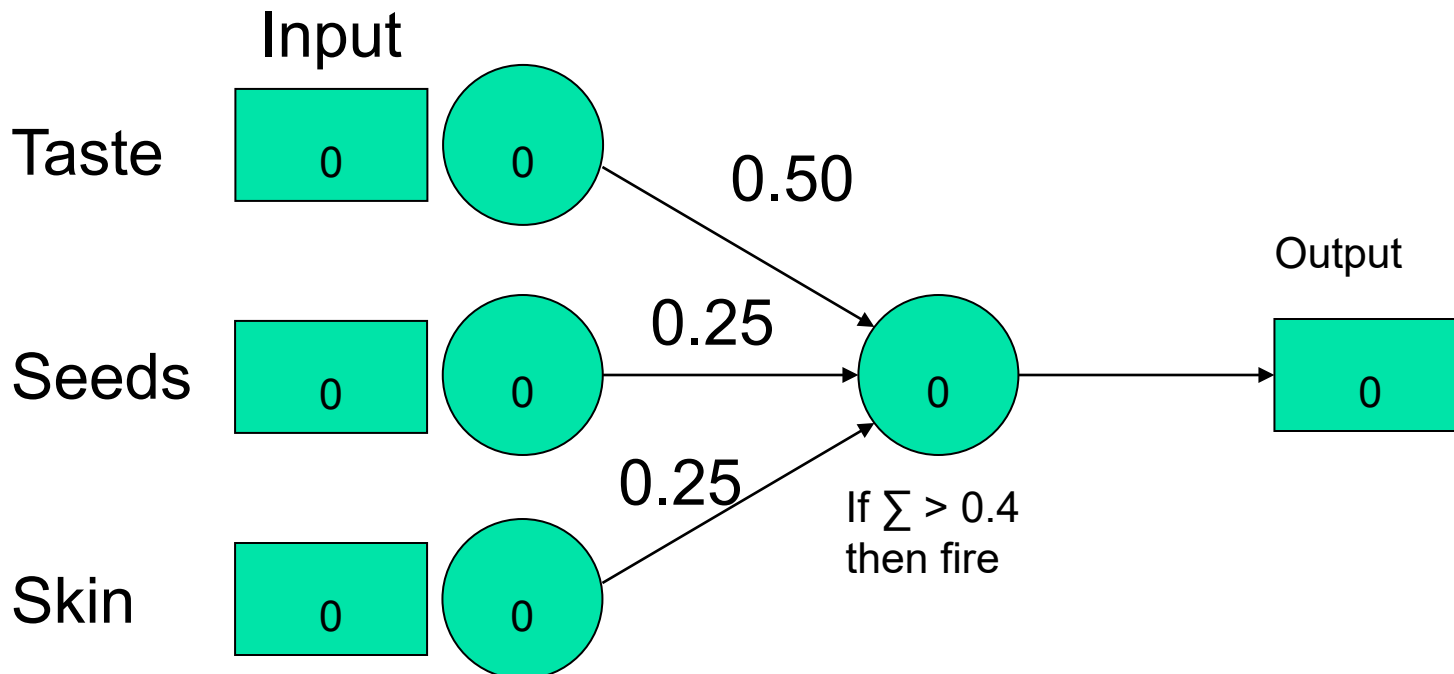
# Example

Here it is with the adjusted weights:



# Example

Show it a lemon:



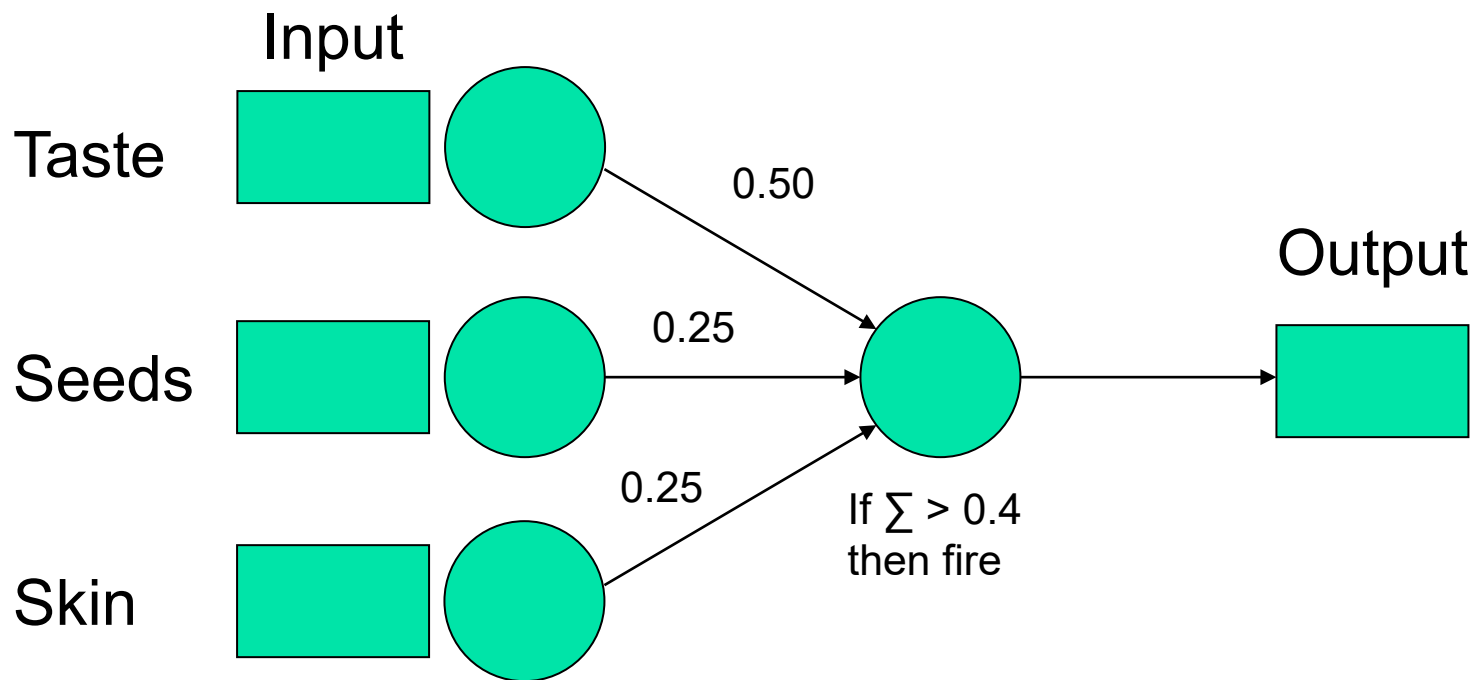
# Example

- How do we change the weights for lemon?

Feature:	Learning rate:	(teacher - output):	Input:	$\Delta w$
taste	0.25	0	0	0
seeds	0.25	0	0	0
skin	0.25	0	0	0

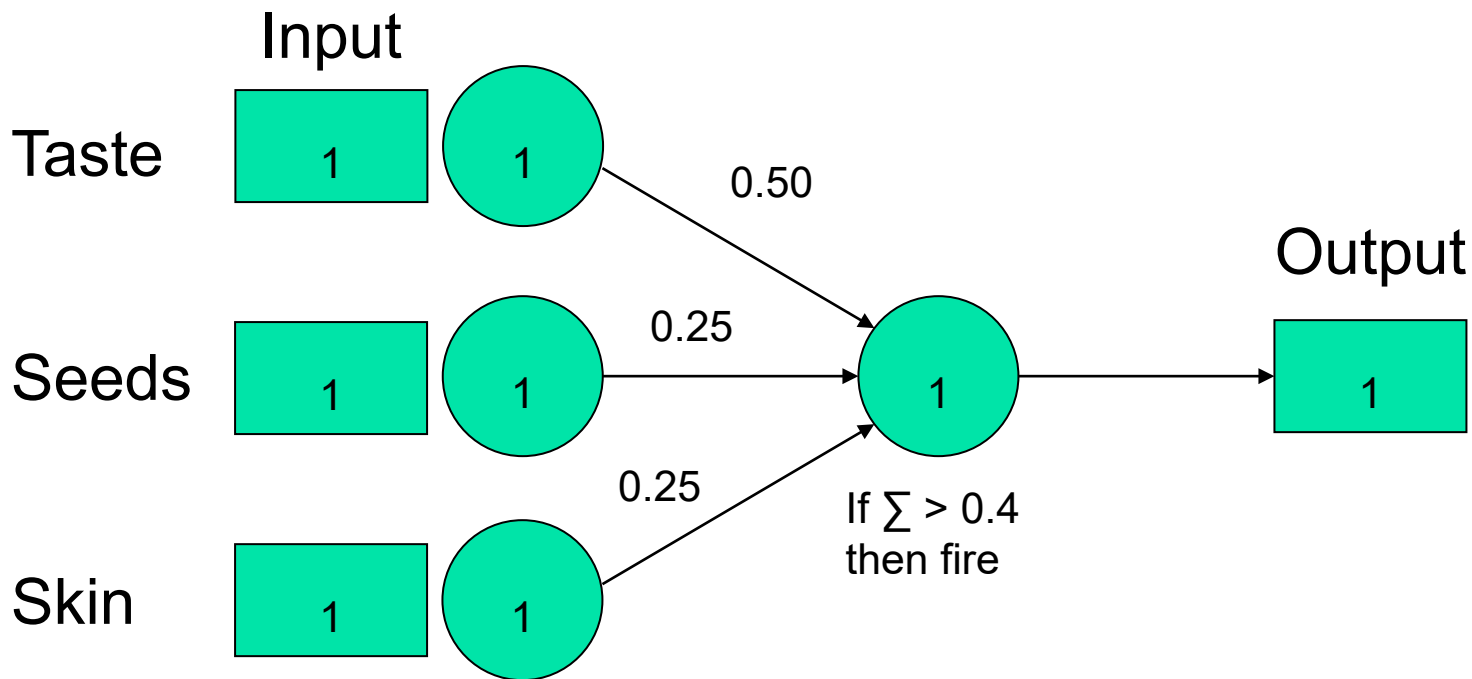
# Example

Here it is with the adjusted weights:



# Example

Show it a strawberry:





# Example

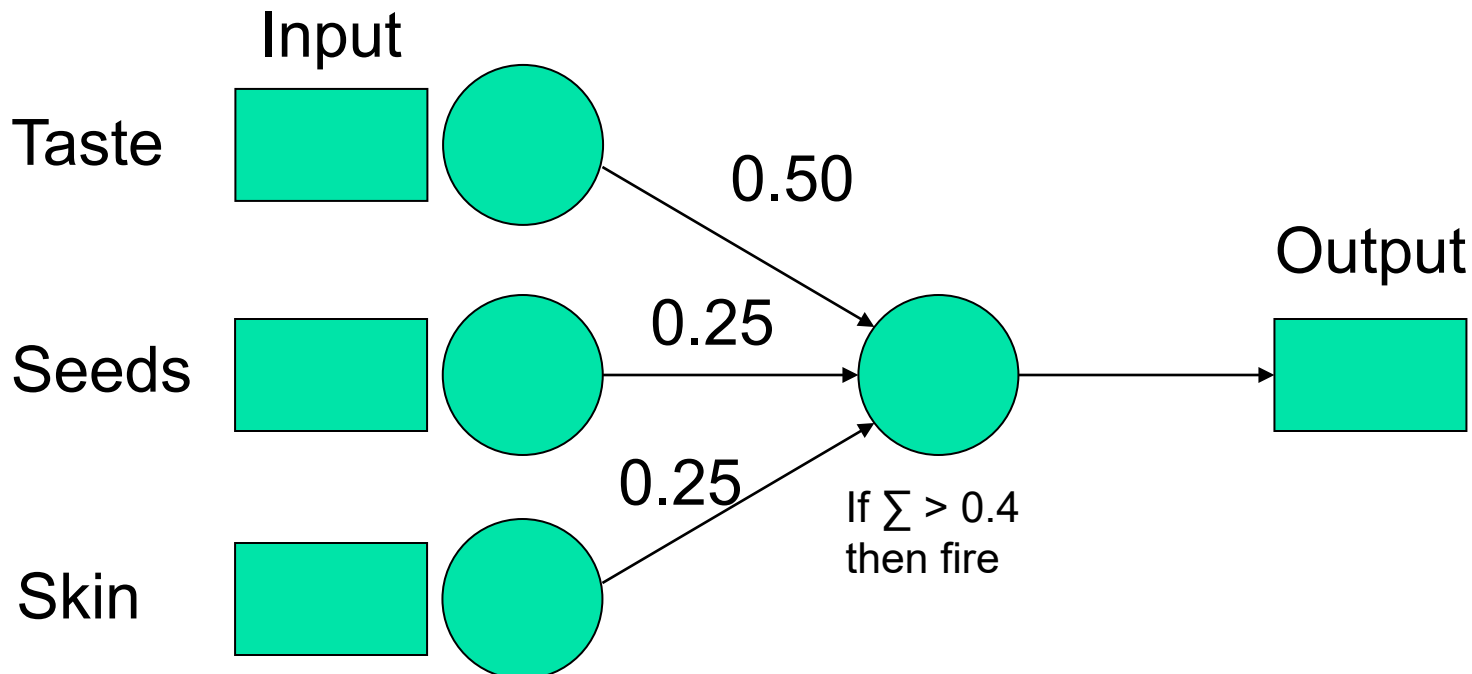
---

- How do we change the weights for strawberry?

Feature:	Learning rate:	(teacher - output):	Input:	$\Delta w$
taste	0.25	0	1	0
seeds	0.25	0	1	0
skin	0.25	0	1	0

# Example

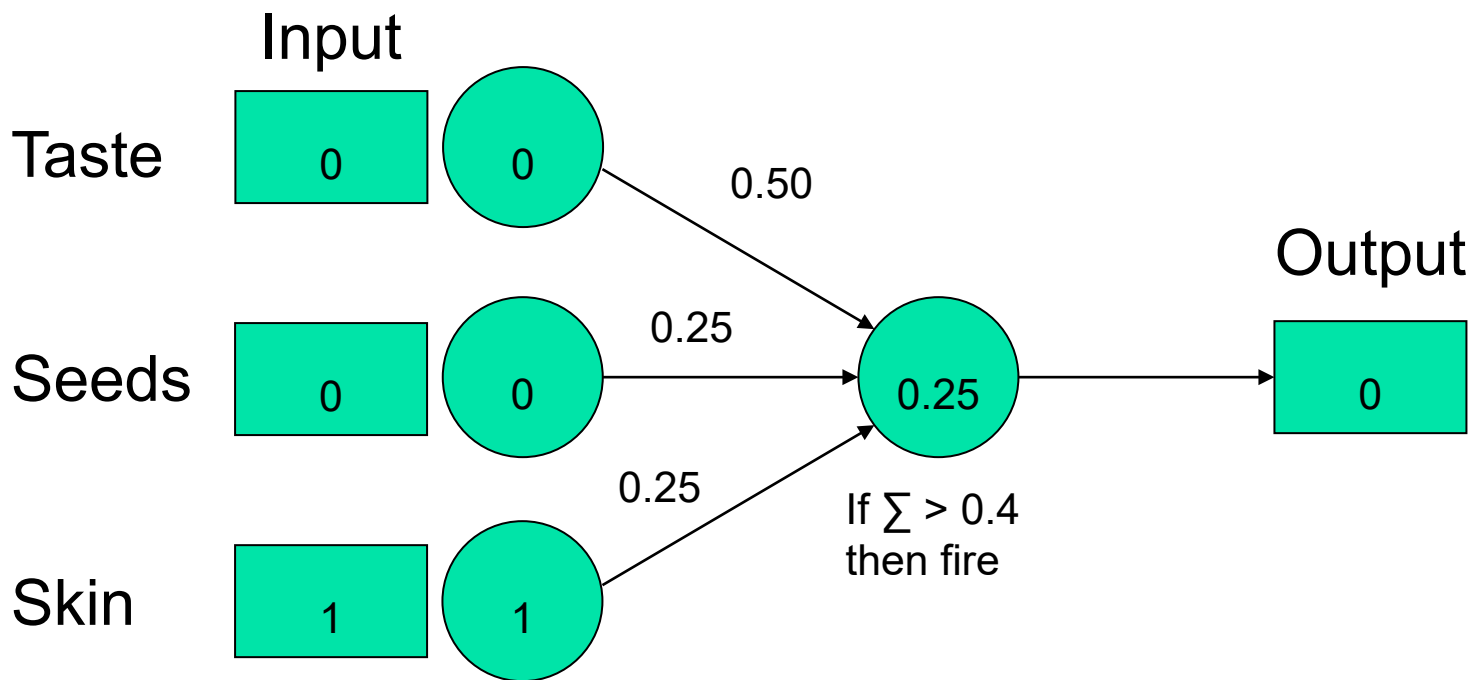
Here it is with the adjusted weights:





# Example

Show it a green apple:





# Example

---

- If you keep going, you will see that this perceptron can correctly classify the examples that we have.

# Perceptron Rule:

## Further Discussion

- A weight of connection changes *only if* both the input value and the error of the output unit are not equal to 0.
  - If the output is correct ( $y_e = o_e$ ) the weights  $w_i$  are not changed
  - If the output is incorrect ( $y_e \neq o_e$ ) the weights  $w_i$  are changed such that the output of the perceptron for the new weights is *closer* to  $y_e$ .
- The algorithm **converges** to the correct classification, if
  - the training data is **linearly separable**;
  - and the learning rate is sufficiently small, usually set below 1, which **determines the amount of correction made in a single iteration.**

# Perceptron convergence Theorem



For any data set *that's linearly separable*, the learning rule is guaranteed to find a solution in a finite number of steps.

## Assumptions:

- At least one such set of weights,  $w^*$ , exists
- There are a finite number of training patterns.
- The threshold function is uni-polar (output is 0 or 1).

# Network Performance for Perceptron

The network performance during training session can be measured by a *root-mean-square (RMS) error value*.

$$\sqrt{\frac{\sum_{i=1}^N (x_i - \hat{x}_i)^2}{N}}$$

均方根误差值  
↓  
评估网络性能

where

N is the number of data points

$x_i$  is the target output

$\hat{x}_i$  is the real/instant output

- **The RMS error is a function of the instant output values only**

# The Network Performance

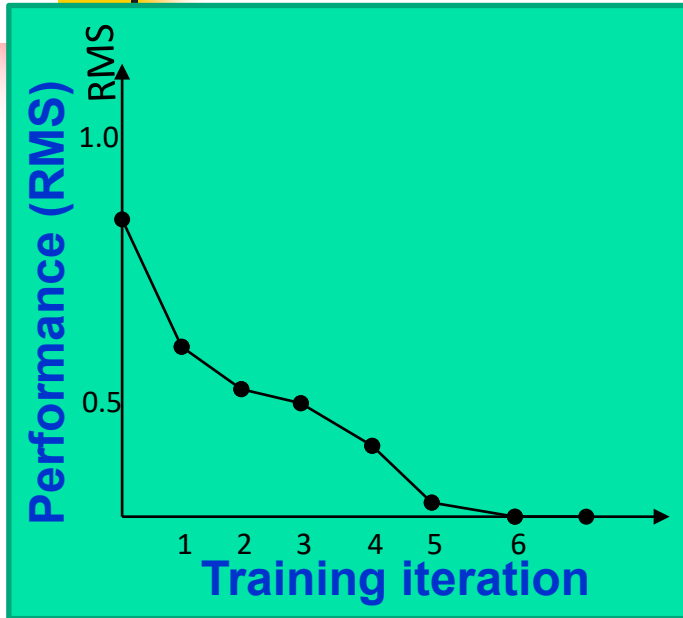
- In turn, the instant outputs  $\hat{x}_i$  are functions of the input values, which are also constants, and of the weights of connections  $w_{ji}$

→ 网络性能好坏只与神经网络权重有关。

So the **performance of the network** measured by the RMS error also **is function of the weights of connections only**

The *best performance of the network* corresponds to the minimum of the RMS error, and we adjust the weights of connections in order to get that minimum.

# RMS on Training Set



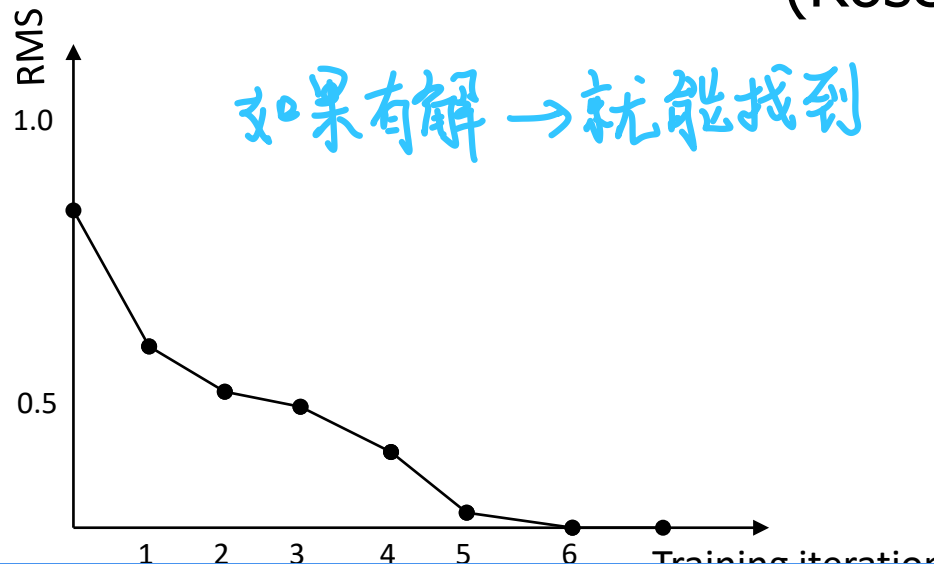
The figure shows a **learning curve**, i.e., dependence of the RMS error on the number of iterations for the training set.

- Initially, the adaptable weights are all set to small random values, and the network does not perform very well.
- As weights are adjusted during training, performance improves; when the error rate is low enough, training stops and the network is said to have **converged**.

# Recall: Perceptron Convergence Theorem

*If a set of weights that allow the perceptron to respond correctly to all of the training patterns exists, **then** the perceptron's learning method will find the set of weights, and it will do it in a finite number of iterations.*

(Rosenblatt, 1962)



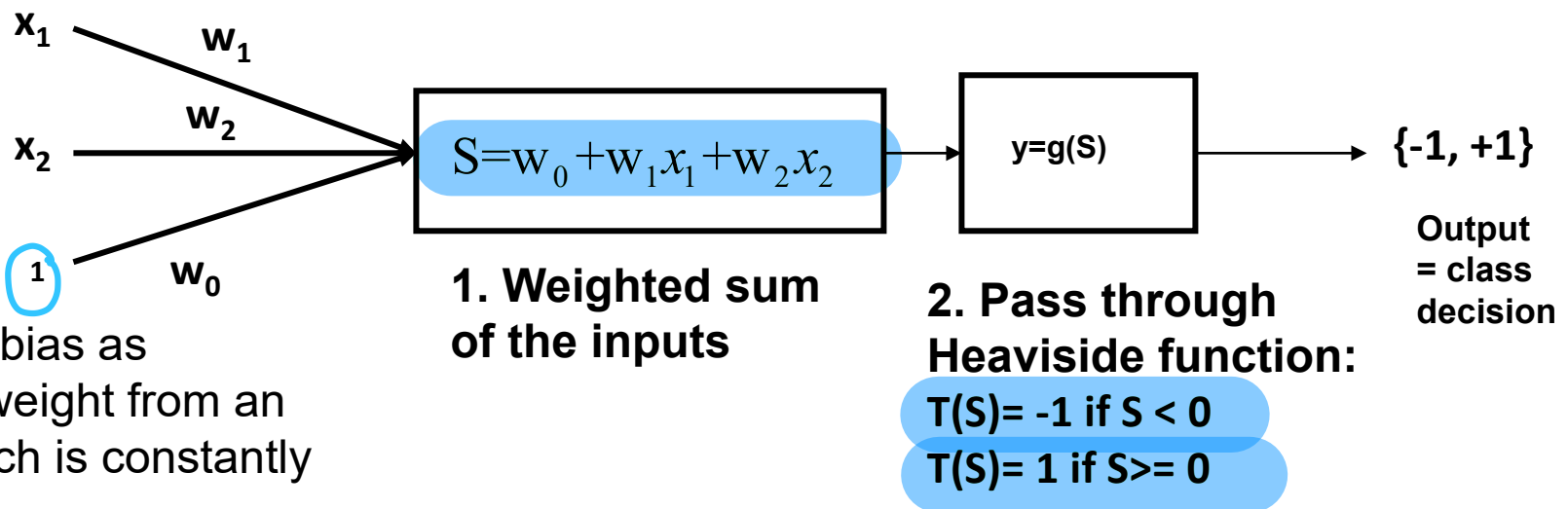


# More on Perceptron Convergence

- There might be another possibility during a training session:
  - eventually performance stops improving, and the RMS error does not get smaller regardless of number of iterations.
- That means the network has **failed** to learn all of the answers correctly.
- *If the training is successful*, the perceptron is said
  - to have *gone through* the supervised learning, and
  - is able to classify patterns *similar to those of the training set*.

# Perceptron As a Classifier

For  $d$ -dimensional data, perceptron consists of  $d$ -weights, a bias, and a thresholding activation function. For 2D data example, we have:



View the bias as another weight from an input which is constantly on

If we group the weights as a vector  $w$ , the net output  $y$  can be expressed as:

$$y = g(w \cdot x + w_0)$$

## Further Discussion

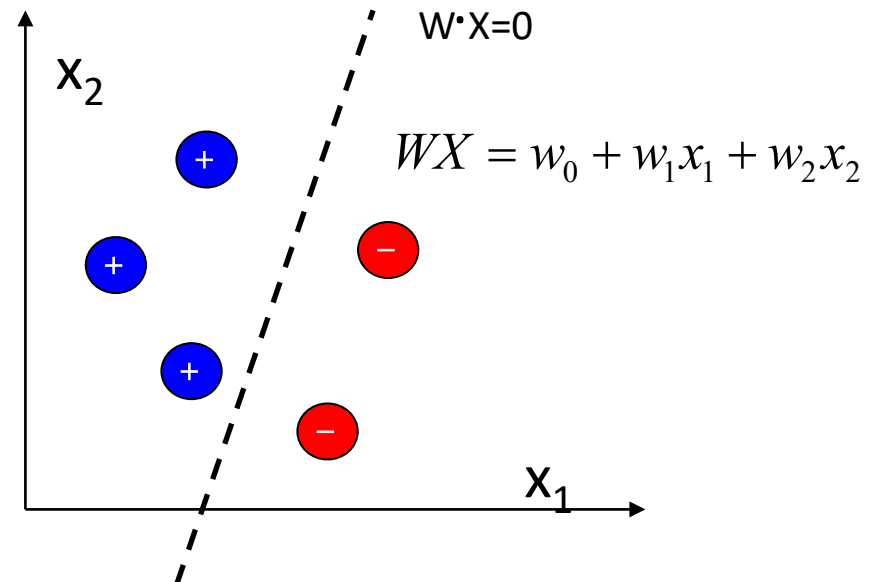
# Perceptron As a Classifier

- A perceptron training is to compute weight vector:

$$W = [w_0, w_1, w_2, \dots, w_p]$$

to correctly classify all the training examples.

E.g.,  
consider when  $p=2$



$WX$  is a **hyperplane** which in 2d is a straight line

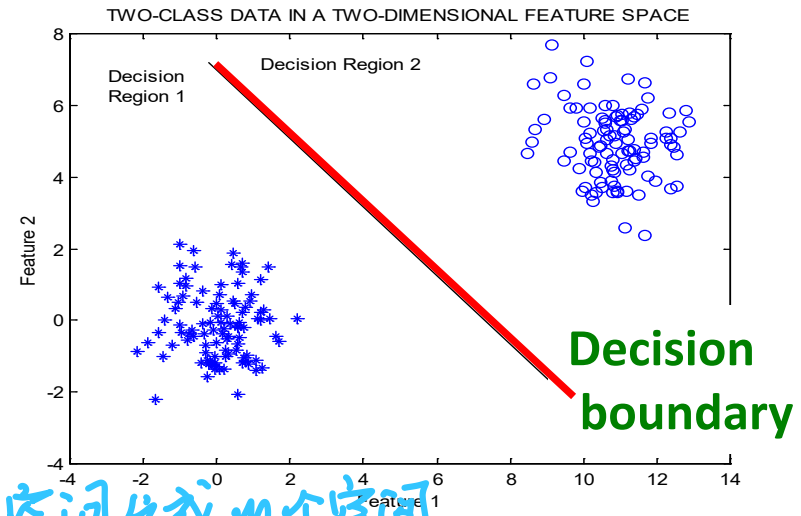
# Neural Network as Classifier

- For 2 classes, view net output as a **discriminant function**  $y(x, w)$ , where:

$y(x, w) = 1$  , if  $x$  in class 1 (C1)

$y(x, w) = -1$ , if  $x$  in class 2 (C2)

## Example



$m$  个类别,  $nn$  作为分类器需要把特征空间分成  $m$  个空间

- For  $m$  classes, a classifier should partition the feature space into  $m$  **decision regions**
  - The **line or curve** separating the classes is the **decision boundary**.
  - In more than 2 dimensions, this is a hyperplane.



# Further on Perceptron Decision Boundary

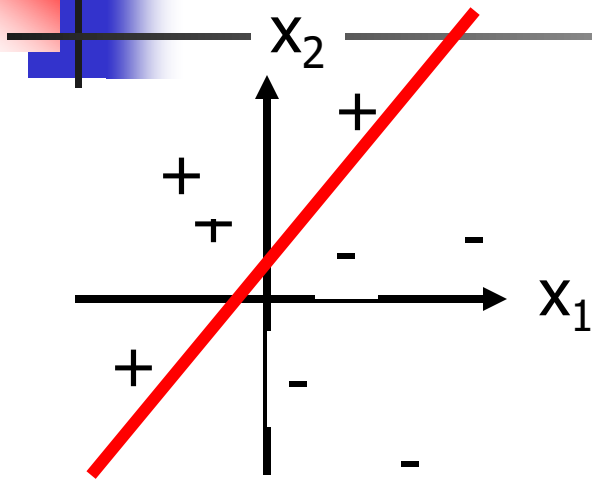
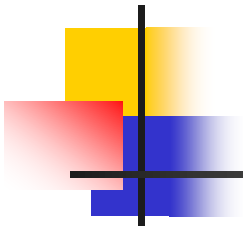
---

A perceptron represents a *hyperplane decision surface* in d-dimensional space, for example, a line in 2D, a plane in 3D, etc.

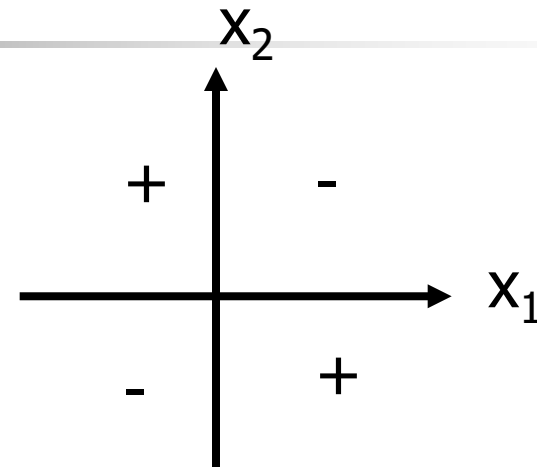
The equation of the hyperplane is  $\mathbf{w} \cdot \mathbf{x}^T = 0$

This is the equation for points in x-space that are **on** the boundary

# Decision boundary of Perceptron



Linearly separable



Non-Linearly separable

- Perceptron is able to represent some useful functions
- But functions that are not linearly separable (e.g. XOR) are not representable

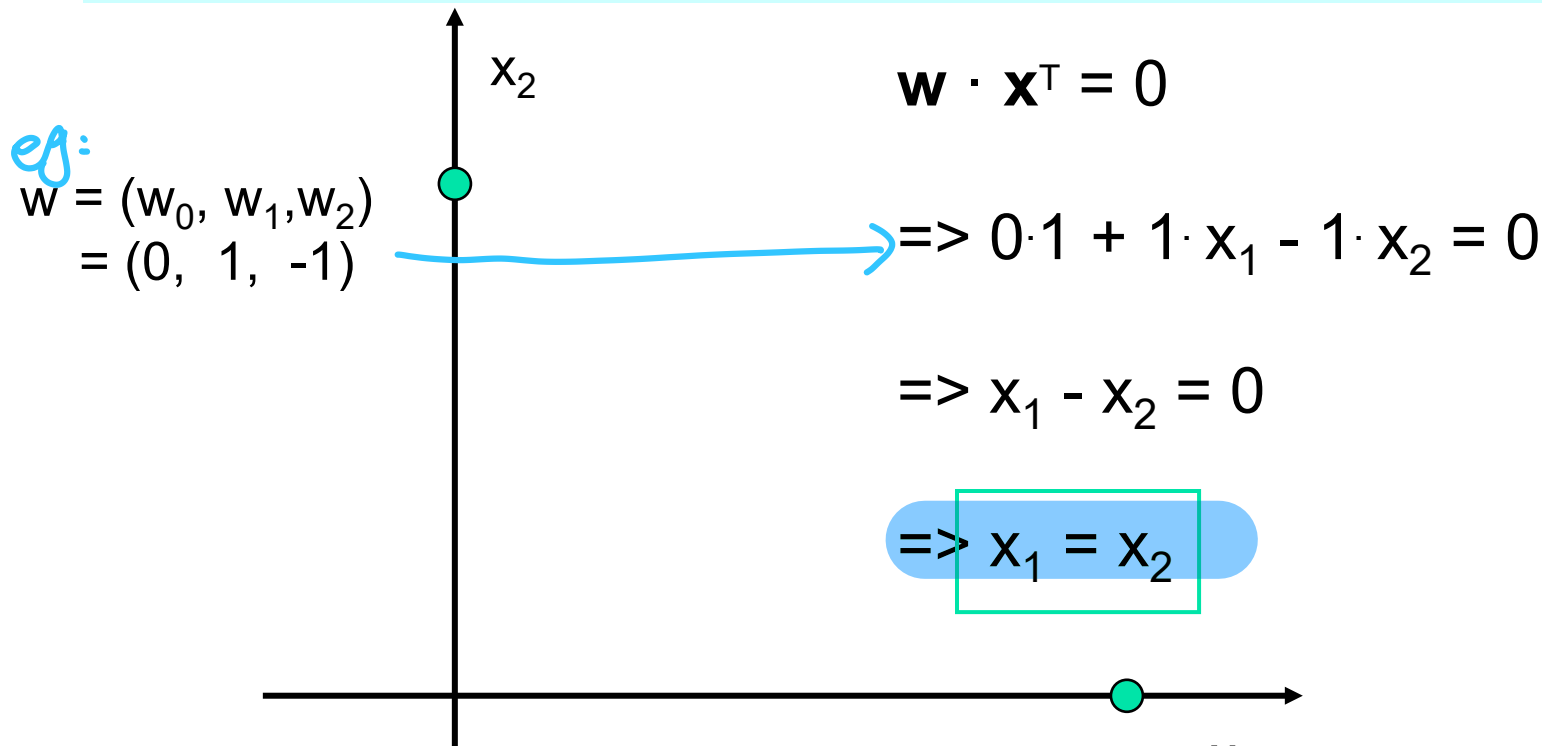
# Example of Perceptron Decision Boundary

## Boundary

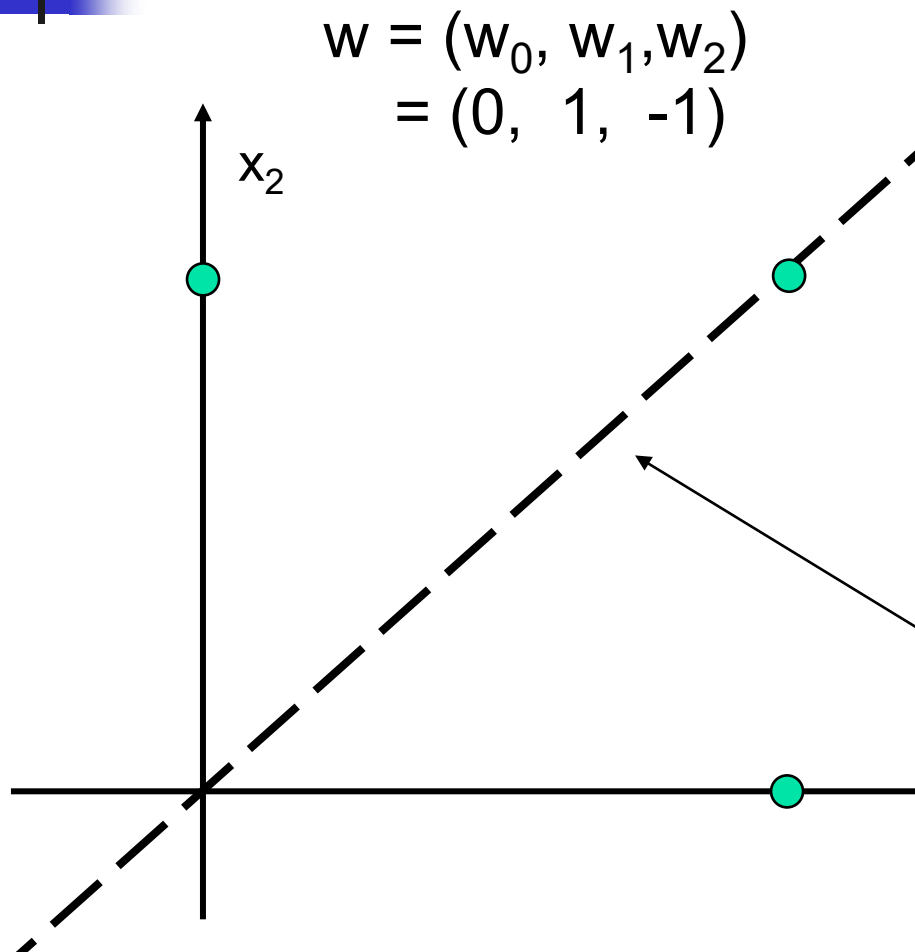
就是 decision boundary

- Decision surface is the surface at which the output of the unit is precisely equal to the threshold, i.e.

$$\sum w_i x_i = \theta$$



# Example of Perceptron Decision Boundary



$$w \cdot x^T = 0$$

$$\Rightarrow 0 \cdot 1 + 1 \cdot x_1 - 1 \cdot x_2 = 0$$

$$\Rightarrow x_1 - x_2 = 0$$

$$\Rightarrow x_1 = x_2$$

**This is the equation  
for the decision boundary**



# Example of Perceptron Decision Boundary

$$\mathbf{w} \cdot \mathbf{x}^T < 0$$

$$\Rightarrow x_1 - x_2 < 0$$

$$\Rightarrow x_1 < x_2$$

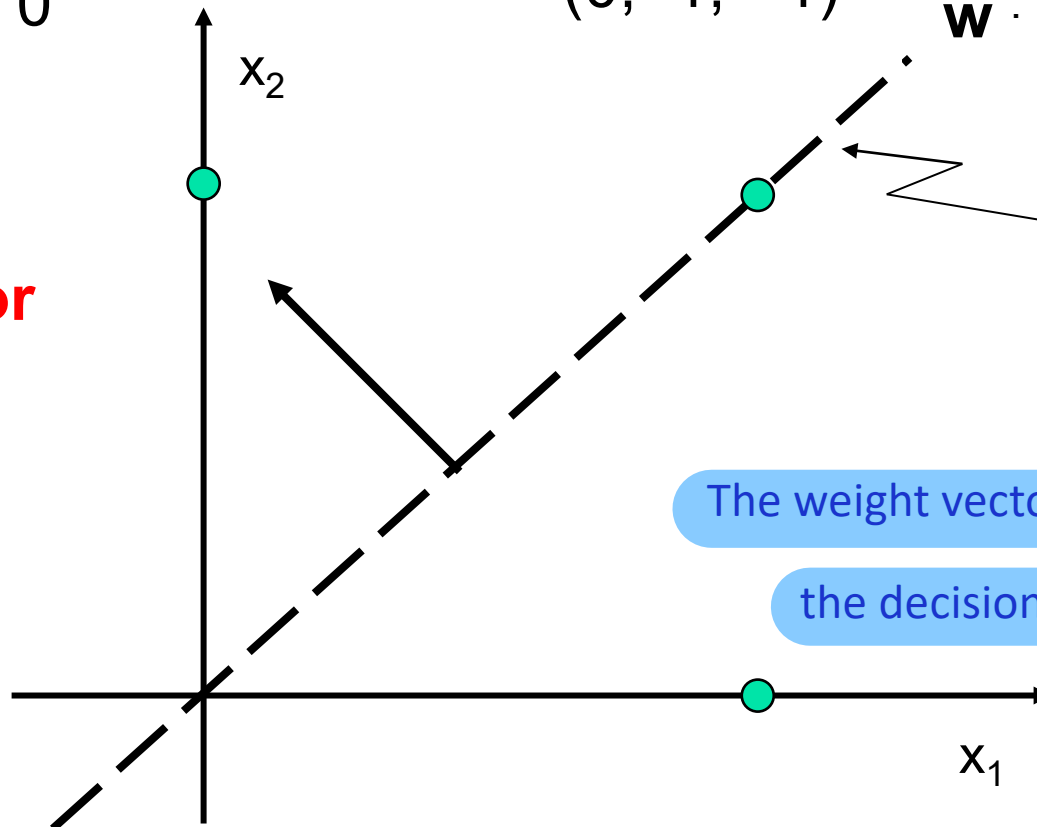
(this is the equation for decision region -1)

$$\begin{aligned}\mathbf{w} &= (w_0, w_1, w_2) \\ &= (0, 1, -1)\end{aligned}$$

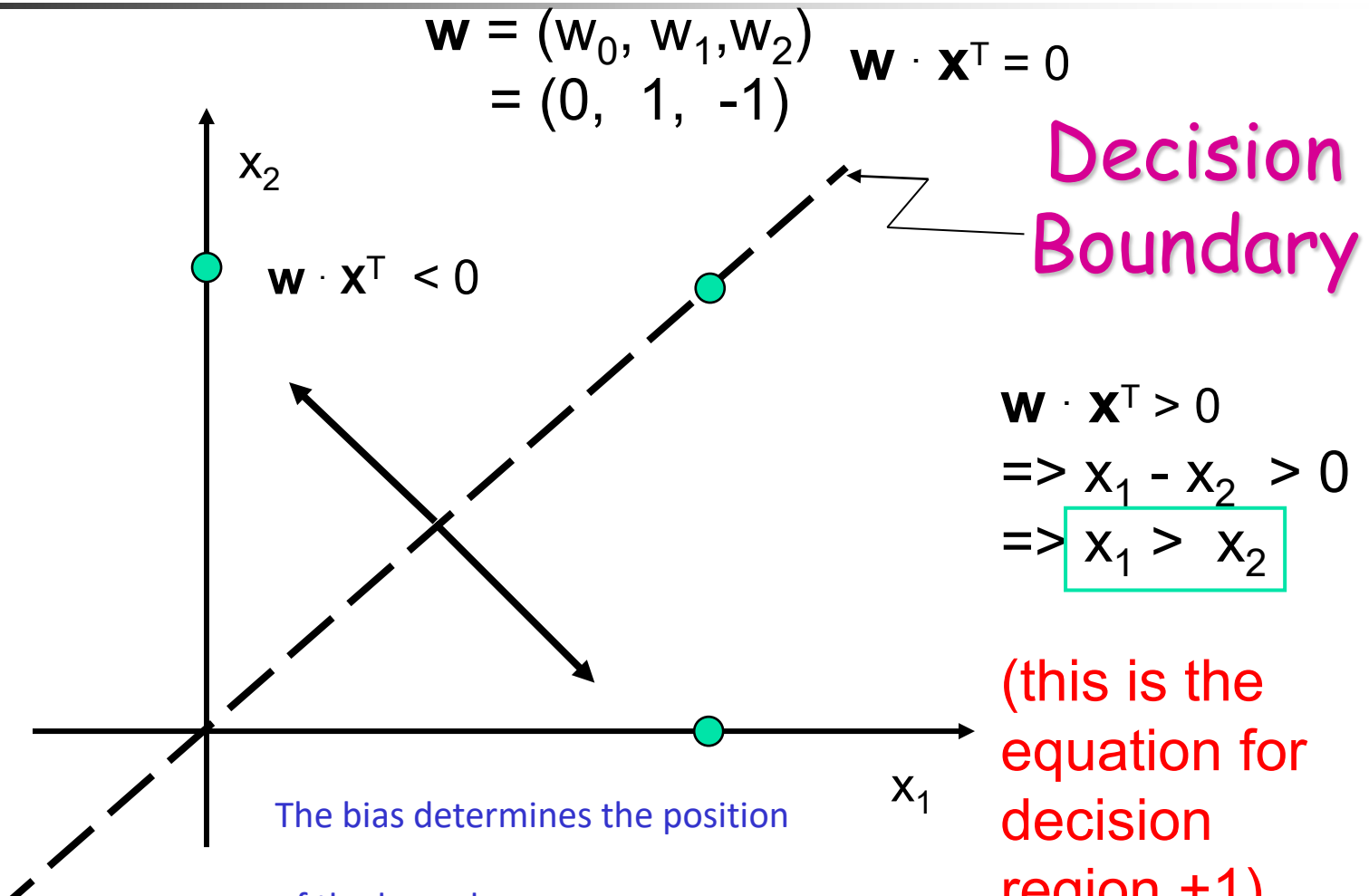
$$\mathbf{w} \cdot \mathbf{x}^T = 0$$

Decision Boundary

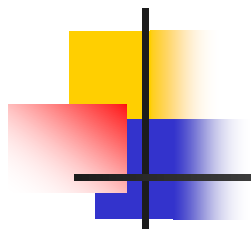
The weight vector is orthogonal to the decision boundary



# Example of Perceptron Decision Boundary



# Linear Separability Problem



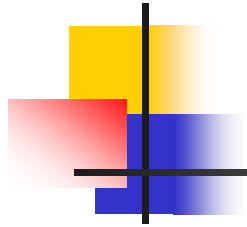
- If two classes of patterns can be separated by a decision boundary, represented by the linear equation

$$b + \sum_{i=1}^n x_i w_i = 0 \quad (wx+b=0)$$

then they are said to be *linearly separable* and the perceptron can correctly classify any patterns

NOTE: without the bias term, the hyperplane will be forced to intersect origin.

# Linear Separability Problem



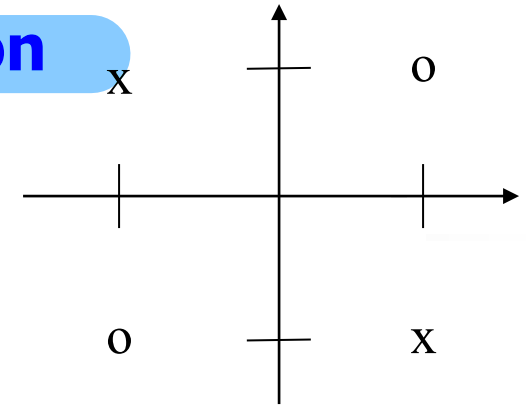
- Decision boundary (i.e.,  $W, b$ ) of linearly separable classes can be determined either by some learning procedures, or by solving linear equation systems based on representative patterns of each classes
- If such a decision boundary does not exist, then the two classes are said to be linearly inseparable.
- Linearly inseparable problems cannot be solved by the simple perceptron network, more sophisticated architecture is needed.

## examples of linearly inseparable classes

### - Logical XOR (exclusive OR) function

patterns (bipolar) decision boundary

$x_1$	$x_2$	$y$
-1	-1	-1
-1	1	1
1	-1	1
1	1	-1



x: class I ( $y = 1$ )  
o: class II ( $y = -1$ )

No line can separate these two classes, as can be seen from the fact that the following linear inequality system has no solution

$$\begin{cases} b - w_1 - w_2 < 0 & (1) \\ b - w_1 + w_2 \geq 0 & (2) \\ b + w_1 - w_2 \geq 0 & (3) \\ b + w_1 + w_2 < 0 & (4) \end{cases}$$

because we have  $b < 0$  from (1) + (4),

and  $b \geq 0$  from (2) + (3) which is a contradiction

# Examples of linearly separable classes

## - Logical **AND** function

$w_0 = 1$  固定!!

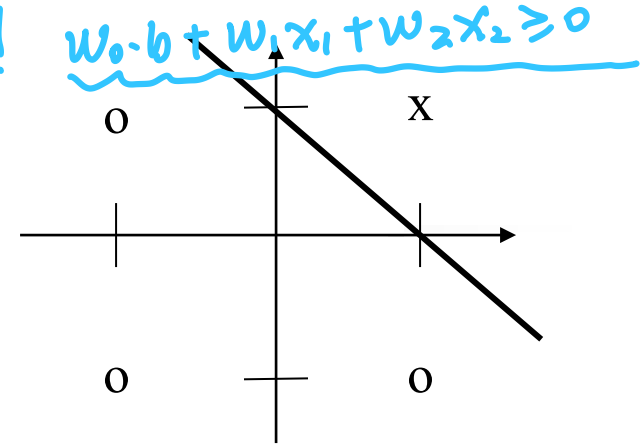
$$w_0 \cdot b + w_1 x_1 + w_2 x_2 \geq 0$$

patterns (bipolar) decision boundary

$x_1$	$x_2$	$y$
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

$$\begin{aligned} w_1 &= 1 \\ w_2 &= 1 \\ b &= -1 \\ \theta &= 0 \end{aligned}$$

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



x: class I ( $y = 1$ )  
o: class II ( $y = -1$ )

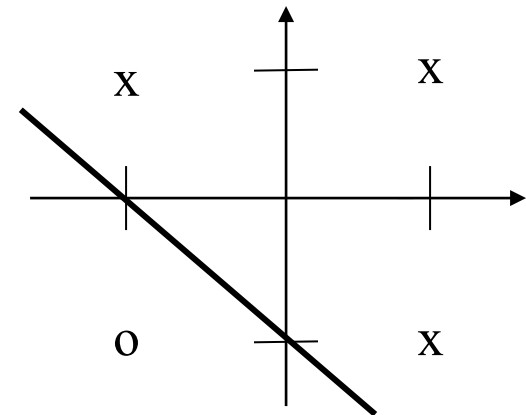
## - Logical **OR** function

patterns (bipolar) decision boundary

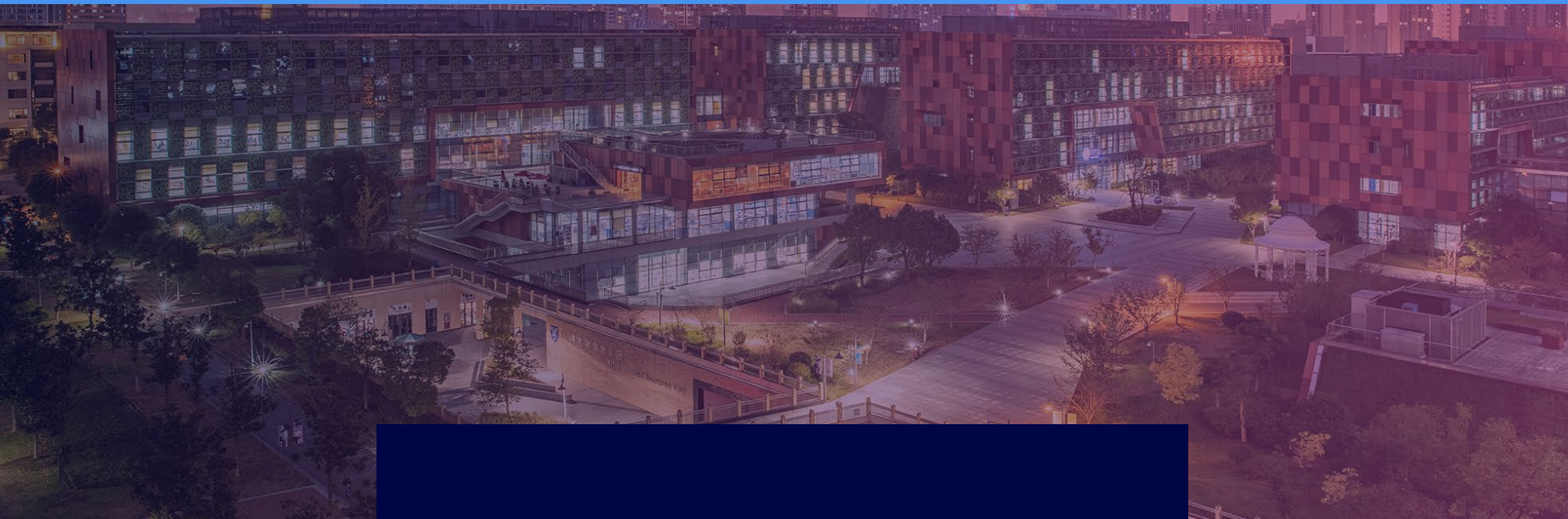
$x_1$	$x_2$	$y$
-1	-1	-1
-1	1	1
1	-1	1
1	1	1

$$\begin{aligned} w_1 &= 1 \\ w_2 &= 1 \\ b &= 1 \\ \theta &= 0 \end{aligned}$$

$$1 + x_1 + x_2 = 0$$



x: class I ( $y = 1$ )  
o: class II ( $y = -1$ )



# THANK YOU



**VISIT US**

[WWW.XJTLU.EDU.CN](http://WWW.XJTLU.EDU.CN)



**FOLLOW US**

@XJTLU



Xi'an Jiaotong-Liverpool University

西交利物浦大學