

# **Neural Networks (CSC372)**

## **Unit 2: Rosenblatt's Perceptron**

**(3 Hrs.)**

**Reference: Simon Haykin (3rd Edition)**

**By Kiran Bagale,**

**Department of IT, 2025**

# Learning Objectives

- Understand the structure and working of Rosenblatt's Perceptron
- Explain the Perceptron Convergence Theorem
- Compare Perceptron with Bayes Classifier
- Apply the Batch Perceptron Algorithm for classification tasks

# Before going to Rosenblatt's Perceptron

- In the formative years of neural networks (1943–1958), several researchers stand out for their pioneering contributions:
  - McCulloch and Pitts (1943) for introducing the idea of neural networks as computing machines.
  - Hebb(1949) for postulating the first rule for self-organized learning.
  - Rosenblatt (1958) for proposing the perceptron as the first model for learning with a teacher (i.e., supervised learning).

# Introduction to Perceptron

- Developed by Frank Rosenblatt (1958)
- Inspired by the biological neuron
- Simplest form of neural Network
- Used for binary classification tasks
- Processes input via weighted sum and activation function

## Assumptions

1. Binary classification (*i. e.*  $y_i \in \{-1, +1\}$ )
2. Data is linearly separable

# Perceptron Model

- Rosenblatt's perceptron is built around a nonlinear neuron, namely, the *McCulloch–Pitts model* of a neuron.

Input vector:  $x = [x_1, x_2, \dots, x_m]$

Weight vector:  $w = [w_1, w_2, \dots, w_m]$

Output:  $y = \varphi(w \cdot x + b)$

## Perceptron Goal :

Classify input vector  $(x_1, x_2, \dots, x_m)$  into two classes:

- If output  $y=+1$ , assign to class **c1**
- If output  $y=-1$ , assign to class **c2**

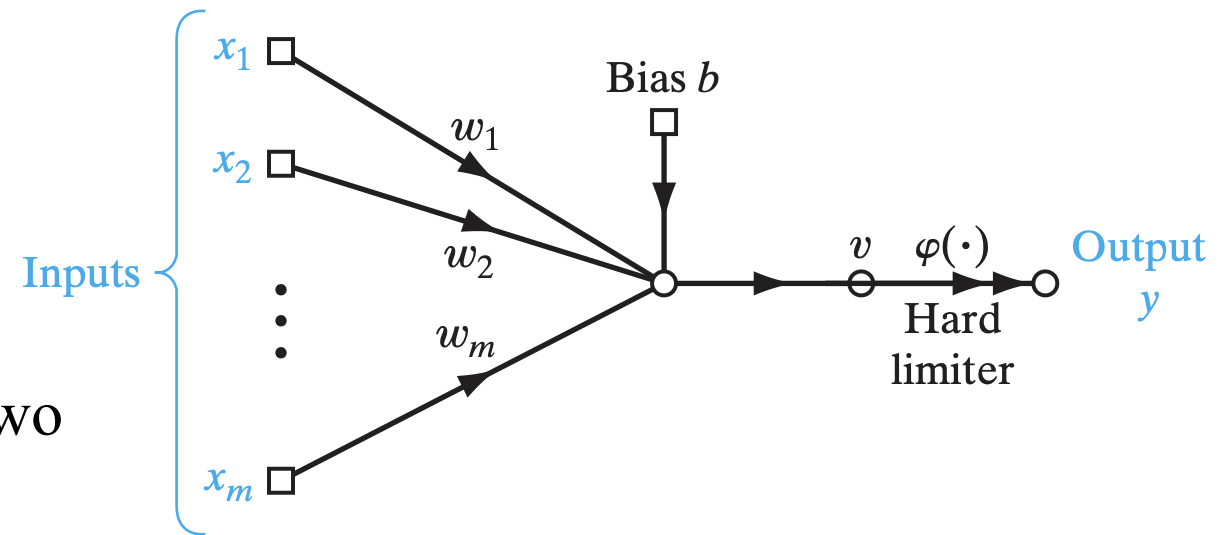


FIG 2.1: Signal-flow graph of the perceptron.

**Learning rule:** *Adjust weights to minimize classification error*

# Classifier

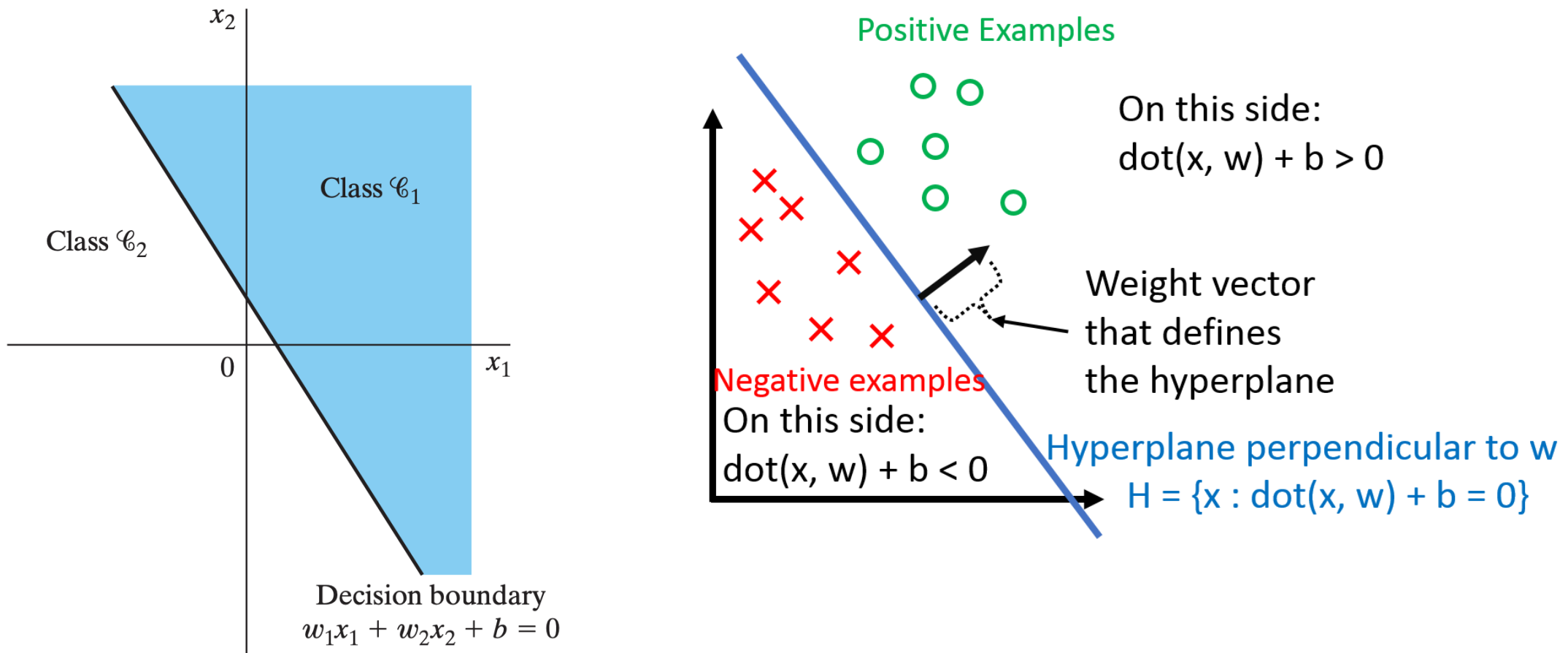


FIG 2.2: Illustration of the hyperplane (in this example, a straight line) as decision boundary for a two-dimensional, two-class pattern-classification problem.

- From the model, we find that the hard limiter input, or induced local field, of the neuron is,

$$y_i = v = \sum_{i=1}^m w_i x_i + b$$

*$b$  is the bias term (without the bias term, the hyperplane that  $\mathbf{w}$  defines would always have to go through the origin). Dealing with  $b$  can be a pain, so we 'absorb' it into the feature vector  $\mathbf{w}$  by adding one additional constant dimension.*

- Under this convention,

$$\begin{aligned} x_i \text{ becomes } \begin{bmatrix} x_i \\ 1 \end{bmatrix} \text{ and} \\ w \text{ becomes } \begin{bmatrix} w \\ b \end{bmatrix} \end{aligned}$$

- We can verify that

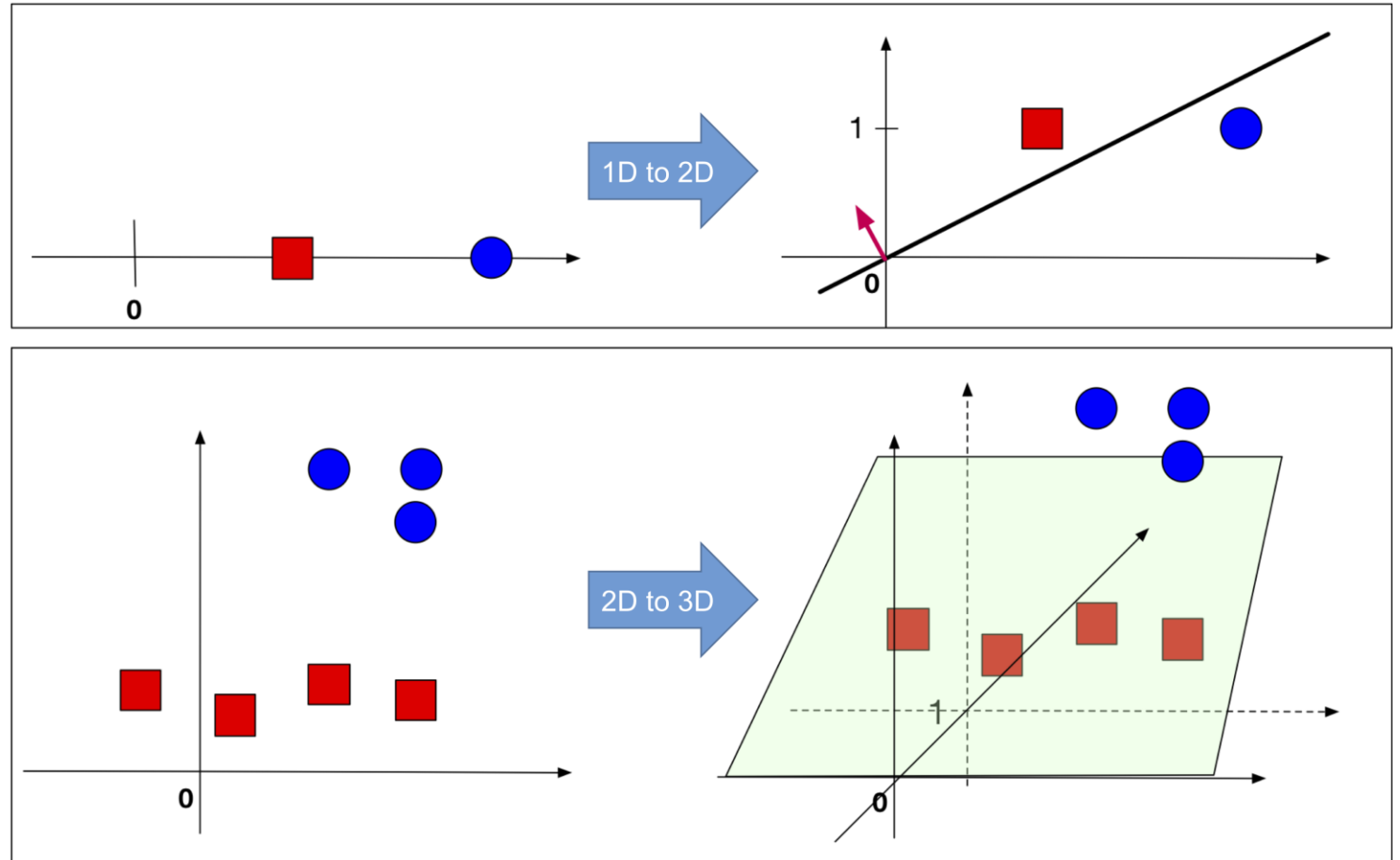
$$\begin{bmatrix} x_i \\ 1 \end{bmatrix}^\top \begin{bmatrix} w \\ b \end{bmatrix} = w^\top x_i + b$$

Using this, we can simplify the above formulation of  $f(x_i)$  to

$$f(x_i) = \varphi(w^\top x)$$

.

**Observation:** Note that  $y_i(\mathbf{w}^\top \mathbf{x}_i) > 0 \Leftrightarrow \mathbf{x}_i$  is classified correctly where 'classified correctly' means that  $\mathbf{x}^i$  is on the correct side of the hyperplane defined by  $\mathbf{w}$ . Also, note that the left side depends on  $y_i \in \{-1, +1\}$  (it wouldn't work if, for example  $y_i \in \{0, +1\}$ ).



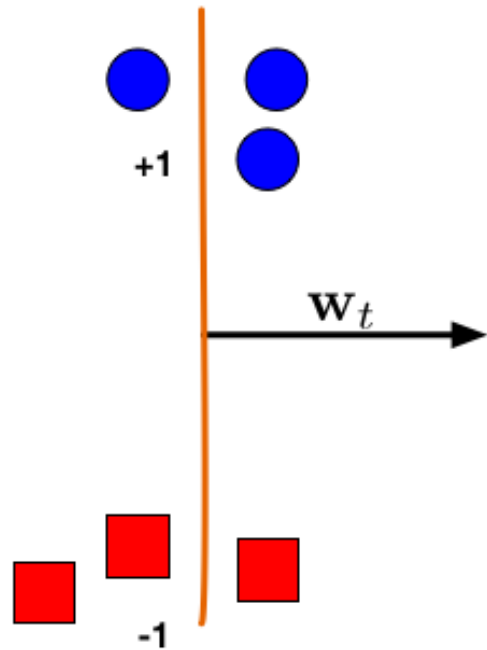
(Left:) The original data is 1-dimensional (top row) or 2-dimensional (bottom row). There is no hyperplane that passes through the origin and separates the red and blue points. (Right:) After a constant dimension was added to all data points such a hyperplane exists.



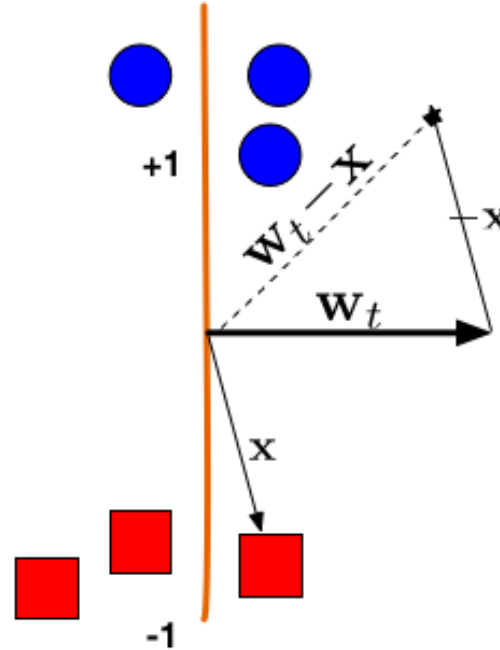
# Perceptron Algorithm

```
i.   Initialize  $\vec{w} = \vec{0}$ 
ii.  while True do
     $m = 0$                                 //count the miss classifications, m
    for  $(x_i, y_i) \in D$  do                //loop over each (data, label) pair in the dataset, D
        if  $y_i(\vec{w}^T \cdot \vec{x}_i) \geq 0$  then    //if the pair  $(x_i, y_i)$  is misclassified
             $\vec{w} \leftarrow \vec{w} + y_i \vec{x}_i$     //update the weight vector
             $m \leftarrow m + 1$ 
        end if
    end for
    if  $m = 0$  then                            //if the recent w gives 0 misclassifications
        break                                //break out the while loop
    end if
end while
```

# Geometric Intuition

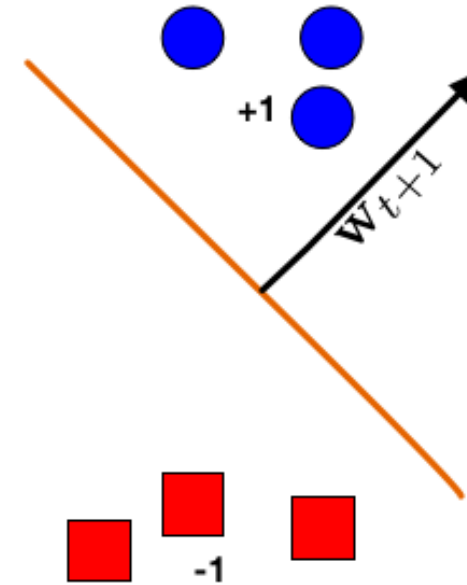


The hyperplane defined by  $w_t$  misclassifies one red (-1) and one blue (+1) point.



The red point  $x$  is chosen and used for an update. Because its label is -1 we need to **subtract**  $x$  from  $w_t$ .

**Fig: lustration of a Perceptron update.**



The updated hyperplane  $w_{t+1} = w_t - x$  separates the two classes and the Perceptron algorithm has converged.

# Perceptron Convergence Theorem

- Proves that perceptron learning algorithm converges in finite steps
- Assumes data is linearly separable
- Guarantees finding a solution (if it exists)
- Let us assume,  $\exists \mathbf{w}^*$  such that  $\mathbf{y}_i(\mathbf{w}^{*T} \mathbf{x}) > 0 \quad \forall (\mathbf{x}_i, \mathbf{y}_i) \in D$
- The bias  $b(n)$  is treated as a synaptic weight driven by a fixed input equal to 1. Thus, Input dimension is  $(m+1)$  by 1 given as,

$$\mathbf{x}(n) = [+1, x_1(n), x_2(n), \dots, x_m(n)]^T$$

where  $n$  denotes the time step in applying the algorithm. Correspondingly, we define the  $(m + 1) -$  by  $-1$  weight vector as,

$$\mathbf{w}(n) = [b, w_1(n), w_2(n), \dots, w_m(n)]^T$$

Accordingly, the linear combiner output is written in the compact form

$$\begin{aligned}v(n) &= \sum_{i=0}^m w_i(n)x_i(n) \\ &= \mathbf{w}^T(n)\mathbf{x}(n)\end{aligned}$$

Adjust the weight vector  $\mathbf{w}$  in such a way that the two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are linearly separable. i.e.

$\mathbf{w}^T\mathbf{X} > \mathbf{0}$  for every input vector  $\mathbf{x}$  belonging to class  $\mathcal{C}_1$   
 $\mathbf{w}^T\mathbf{X} < \mathbf{0}$  for every input vector  $\mathbf{x}$  belonging to class  $\mathcal{C}_2$

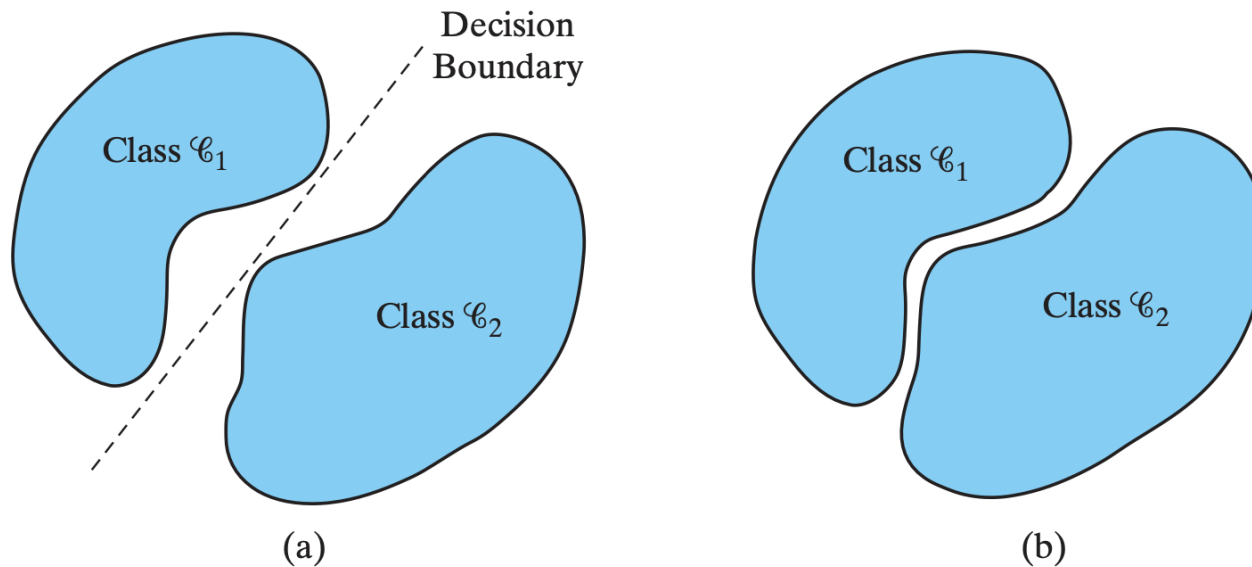


Fig: (a) A pair of linearly separable patterns. (b) A pair of non-linearly separable patterns.

Now, suppose that we rescale each data point and the  $\mathbf{w}^*$  such that  $\|\mathbf{w}^*\| = 1$  and  $\|x_i\| \leq 1 \quad \forall x_i \in D$

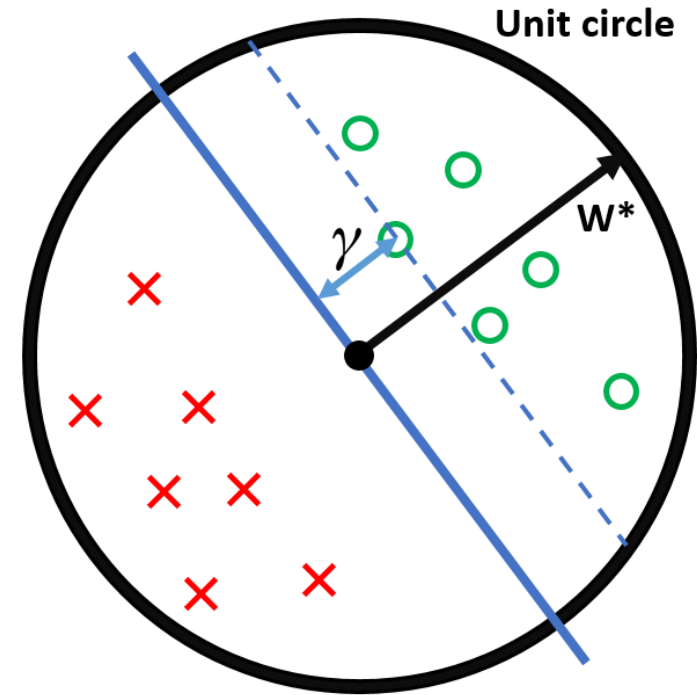
Let us define the margin  $\gamma$  of the hyperplane  $\mathbf{w}^*$  as

$$\gamma = \min_{(x_i, y_i) \in D} |\mathbf{w}^{*T} x_i|$$

To summarize our setup:

- All inputs  $x_i$  live within the unit sphere.
- There exists a separating hyperplane defined by  $\mathbf{w}^*$ , with  $\|\mathbf{w}^*\|=1$  (i.e.  $\mathbf{w}^*$  lies exactly on the unit sphere).
- $\gamma$  is the distance from this hyperplane (blue) to the closest data point.
- **Theorem:** If all of the above holds, then the Perceptron algorithm makes at most  $1/\gamma^2$  mistakes.

Or “For any finite set of linearly separable labeled examples, the Perceptron Learning Algorithm (PLA) will terminate after a finite number of iterations.”



## Proof:

Keeping what we defined above, consider the effect of an update ( $\mathbf{w}$  becomes  $\mathbf{w} + y\mathbf{x}$ ) on the two terms  $\mathbf{w}^T \mathbf{w}^*$  and  $\mathbf{w}^T \mathbf{w}$ . We will use two facts:

- $y(\mathbf{x}^T \mathbf{w}) \leq 0$ : This holds because  $\mathbf{x}$  is misclassified by  $\mathbf{w}$  - otherwise we wouldn't make the update.
- $y(\mathbf{x}^T \mathbf{w}^*) > 0$ : This holds because  $\mathbf{w}^*$  is a separating hyperplane and classifies all points correctly.

1. Consider the effect of an update on  $\mathbf{w}^T \mathbf{w}^*$ :

$$(\mathbf{w} + y\mathbf{x})^T \mathbf{w}^* = \mathbf{w}^T \mathbf{w}^* + y(\mathbf{x}^T \mathbf{w}^*) \geq \mathbf{w}^T \mathbf{w}^* + \gamma$$

The inequality follows from the fact that, for  $\mathbf{w}^*$ , the distance from the hyperplane defined by  $\mathbf{w}^*$  to  $\mathbf{x}$  must be at least  $\gamma$  (i. e.  $y(\mathbf{x}^T \mathbf{w}^*) = |\mathbf{x}^T \mathbf{w}^*| \geq \gamma$ ).

This means that for each update,  $\mathbf{w}^T \mathbf{w}^*$  grows by at least  $\gamma$ .

2. Consider the effect of an update on  $\mathbf{w}^\top \mathbf{w}$ :

$$(\mathbf{w} + \mathbf{y}\mathbf{x})^\top (\mathbf{w} + \mathbf{y}\mathbf{x}) = \mathbf{w}^\top \mathbf{w} + 2\mathbf{y}(\mathbf{w}^\top \mathbf{x}) + \mathbf{y}^2(\mathbf{x}^\top \mathbf{x}) \leq \mathbf{w}^\top \mathbf{w} + 1$$

The inequality follows from the fact that

- $2\mathbf{y}(\mathbf{w}^\top \mathbf{x}) < 0$  as we had to make an update, meaning  $\mathbf{x}$  was misclassified
- $0 \leq \mathbf{y}^2(\mathbf{x}^\top \mathbf{x}) \leq 1$  as  $\mathbf{y}^2=1$  and all  $\mathbf{x}^\top \mathbf{x} \leq 1$  (because  $\|\mathbf{x}\| \leq 1$ ).

This means that for each update,  $\mathbf{w}^\top \mathbf{w}$  grows by at most 1.

3. Now we know that after  $M$  updates the following two inequalities must hold:

$$1) \mathbf{w}^\top \mathbf{w}^* \geq M\gamma$$

$$2) \mathbf{w}^\top \mathbf{w} \leq M.$$

We can then complete the proof:

$$M\gamma \leq \mathbf{w}^\top \mathbf{w}^*$$

By (1)

$$= \|\mathbf{w}\| \cos(\theta)$$

by definition of inner-product, where  $\theta$  is the angle between  $\mathbf{w}$  and  $\mathbf{w}^*$ .

$$\leq \|\mathbf{w}\|$$

by definition of  $\cos$ , we must have  $\cos(\theta) \leq 1$ .

$$= \sqrt{\mathbf{w}^\top \mathbf{w}}$$

by definition of  $\|\mathbf{w}\|$

$$\leq \sqrt{M}$$

By (2)

$$\Rightarrow M\gamma \leq \sqrt{M}$$

$$\Rightarrow M^2\gamma^2 \leq M$$

$$\Rightarrow M \leq 1/\gamma^2$$

*And hence, the number of updates  $M$  is bounded from above by a constant.*



# Limitations of Perceptron

- Fails to solve non-linearly separable problems (e.g., XOR)
- Can only classify linearly separable datasets
- Motivated development of multilayer perceptron

# RELATION BETWEEN THE PERCEPTRON AND BAYES CLASSIFIER FOR A GAUSSIAN ENVIRONMENT

- The perceptron bears a certain relationship to a classical pattern classifier known as the Bayes classifier.
- When the environment is Gaussian, the Bayes classifier reduces to a linear classifier.

$$\mathcal{X} = \mathcal{H}$$

$$\mathcal{C}_j = \mathcal{C}_i = \mathcal{C}_j$$

## Bayes Classifier

- *Bayes classifier*, or *Bayes hypothesis testing procedure*
  - minimize the *average risk*, denoted by  $\mathfrak{R}$ .

- For a two-class problem, represented by classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , the average risk is defined by Van Trees (1968) as,

$$\begin{aligned} \mathcal{R} = & c_{11}p_1 \int_{\mathcal{H}_1} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} + c_{22}p_2 \int_{\mathcal{H}_2} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} \quad \left. \vphantom{\int_{\mathcal{H}_1}} \right\} \text{Correctly classified} \\ & + c_{21}p_1 \int_{\mathcal{H}_2} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} + c_{12}p_2 \int_{\mathcal{H}_1} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} \quad \left. \vphantom{\int_{\mathcal{H}_1}} \right\} \text{missclassified} \end{aligned} \quad (1.23)$$

Where,

$p_i$  - *prior probability* that the observation vector corresponds to an object in class  $\mathcal{C}_i$ , with  $i=1,2$ , and

$$p_1 + p_2 = 1$$

$\mathcal{C}_{ij}$  cost of deciding in favor of class  $\mathcal{C}_i$  represented by subspace  $H_i$  when class  $\mathcal{C}_i$  is true (i.e., observation vector  $\mathbf{x}$  corresponds to an object in class  $\mathcal{C}_1$ ), with  $i, j = 1, 2$

$p_{\mathbf{x}}(\mathbf{x}|\mathcal{C}_i)$  conditional probability density function (pdf) of the random vector  $\mathbf{X}$ , given that the observation vector  $\mathbf{x}$  corresponds to an object in class  $\mathcal{C}_i$ , with  $i = 1, 2$ .

- To minimize the risk,
  - each observation vector  $\mathbf{x}$  must be assigned in the overall observation space  $\mathcal{H}$  to either  $\mathcal{H}1$  or  $\mathcal{H}2$ . Thus,
 
$$\mathcal{H} = \mathcal{H}1 + \mathcal{H}2 \dots \dots (1.24)$$

we may rewrite Eq. (1.23) in the equivalent form

$$\begin{aligned} \mathcal{R} = & c_{11}p_1 \int_{\mathcal{X}_1} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} + c_{22}p_2 \int_{\mathcal{X}-\mathcal{X}_1} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} \\ & + c_{21}p_1 \int_{\mathcal{X}-\mathcal{X}_1} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} + c_{12}p_2 \int_{\mathcal{X}_1} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} \end{aligned} \quad (1.25)$$

where  $c_{11} < c_{21}$  and  $c_{22} < c_{12}$ . We now observe the fact that

$$\int_{\mathcal{X}} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)d\mathbf{x} = \int_{\mathcal{X}} p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)d\mathbf{x} = 1 \quad (1.26)$$

Hence, Eq. (1.25) reduces to

$$\begin{aligned}\mathcal{R} = & c_{21}p_1 + c_{22}p_2 \\ & + \int_{\mathcal{X}_1} [p_2(c_{12} - c_{22}) p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2) - p_1(c_{21} - c_{11}) p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)] d\mathbf{x}\end{aligned}\tag{1.27}$$

The first two terms on the right-hand side of Eq. (1.27) represent a fixed cost. Since the requirement is to minimize the average risk  $\mathcal{R}$ , we may therefore deduce the following strategy from Eq.(1.27) for optimum classification:

1. Assign any value of  $\mathbf{x}$  to class **C1** if the expression inside the square brackets is **-ve**, because this will reduce the overall risk  $\mathcal{R}$ .
2. Assign any value of  $\mathbf{x}$  to class **C2** if the expression is **+ve**, because this will increase the overall risk if included in class **C1**.
3. If the expression is **0**, it doesn't affect the risk, so those values of  $\mathbf{x}$  can be assigned to either class. Here, we choose to assign them to **class C2**.

On this basis, we may now formulate the Bayes classifier as follows:

*If the condition*

$$p_1(c_{21} - c_{11}) p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1) > p_2(c_{12} - c_{22}) p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)$$

*holds, assign the observation vector  $\mathbf{x}$  to subspace  $\mathcal{X}_1$  (i.e., class  $\mathcal{C}_1$ ). Otherwise assign  $\mathbf{x}$  to  $\mathcal{X}_2$  (i.e., class  $\mathcal{C}_2$ ).*

To simplify matters, define

$$\Lambda(\mathbf{x}) = \frac{p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_1)}{p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_2)} \quad (1.28)$$

and

$$\xi = \frac{p_2(c_{12} - c_{22})}{p_1(c_{21} - c_{11})} \quad (1.29)$$

The quantity  $\Lambda(\mathbf{x})$ , the ratio of two conditional probability density functions, is called the *likelihood ratio*.

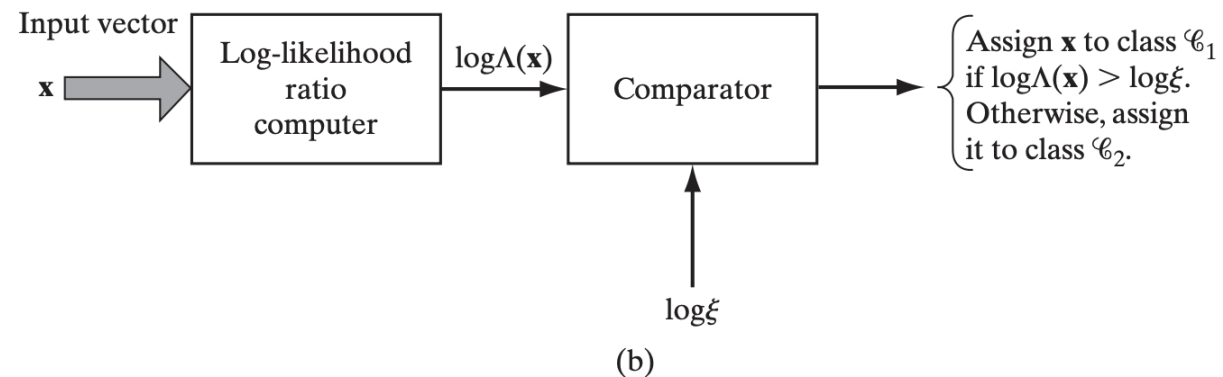
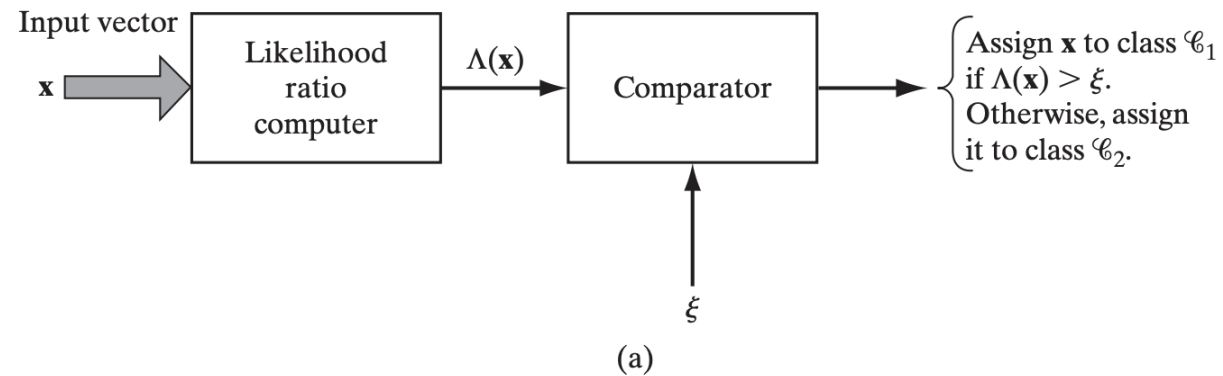
The quantity  $\xi$  is called the *threshold* of the test.

*Note that both  $\xi(\mathbf{x})$  and  $\Lambda(\mathbf{x})$  are always positive.*

Bayes classifier by stating the following

*If, for an observation vector  $\mathbf{x}$ , the likelihood ratio  $\Lambda(\mathbf{x})$  is greater than the threshold  $\xi$ , assign  $\mathbf{x}$  to class  $\mathcal{C}_1$ . Otherwise, assign it to class  $\mathcal{C}_2$ .*

The data processing involved in designing the Bayes classifier is confined entirely to the computation of the likelihood ratio  $\Lambda(\mathbf{x})$ .



*Fig: Two equivalent implementations of the Bayes classifier: (a) Likelihood ratio test, (b) Log-likelihood ratio test.*

# Bayes Classifier for a Gaussian Distribution

- Consider a case of a two-class problem, for which the underlying distribution is Gaussian.
- The random vector  $\mathbf{X}$  has a mean value that depends on whether it belongs to class  $c_1$  or class  $c_2$ , but the covariance matrix of  $\mathbf{X}$  is the same for both classes.

$$\begin{aligned}\text{Class } \mathcal{C}_1: \quad \mathbb{E}[\mathbf{X}] &= \boldsymbol{\mu}_1 \\ &\mathbb{E}[(\mathbf{X} - \boldsymbol{\mu}_1)(\mathbf{X} - \boldsymbol{\mu}_1)^T] = \mathbf{C} \\ \text{Class } \mathcal{C}_2: \quad \mathbb{E}[\mathbf{X}] &= \boldsymbol{\mu}_2 \\ &\mathbb{E}[(\mathbf{X} - \boldsymbol{\mu}_2)(\mathbf{X} - \boldsymbol{\mu}_2)^T] = \mathbf{C}\end{aligned}$$

The covariance matrix  $\mathbf{C}$  is nondiagonal, which means that the samples drawn from classes  $c_1$  and  $c_2$  are *correlated*. It is assumed that  $\mathbf{C}$  is nonsingular, so that its inverse matrix  $\mathbf{C}^{-1}$  exists.



- The conditional probability density function of  $\mathbf{X}$  as the multivariate Gaussian distribution

$$p_{\mathbf{X}}(\mathbf{x}|\mathcal{C}_i) = \frac{1}{(2\pi)^{m/2}(\det(\mathbf{C}))^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)\right), \quad i = 1, 2 \quad (1.30)$$

where  $m$  is the dimensionality of the observation vector  $\mathbf{x}$ .

The two classes  $C_1$  and  $C_2$  are equiprobable:

$$p_1 = p_2 = \frac{1}{2} \quad (1.31)$$

Misclassifications carry the same cost, and no cost is incurred on correct classifications:

$$c_{21} = c_{12} \quad \text{and} \quad c_{11} = c_{22} = 0 \quad (1.32)$$

- By substituting Eq. (1.30) into (1.28) and taking the natural logarithm, we get (after simplifications)

$$\begin{aligned}\log \Lambda(\mathbf{x}) &= -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{C}^{-1} \mathbf{x} + \frac{1}{2} (\boldsymbol{\mu}_2^T \mathbf{C}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^T \mathbf{C}^{-1} \boldsymbol{\mu}_1)\end{aligned}\quad (1.33)$$

- By substituting Eqs. (1.31) and (1.32) into Eq. (1.29) and taking the natural logarithm, we get

$$\log \xi = 0 \quad (1.34)$$

Equations (1.33) and (1.34) state that the Bayes classifier for the problem at hand is a *linear classifier*, as described by the relation

$$y = \mathbf{w}^T \mathbf{x} + b \quad (1.35)$$

Where,

$$\mathbf{y} = \log \Lambda(\mathbf{x}) \quad (1.36)$$

$$\mathbf{w} = \mathbf{C}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (1.37)$$

$$b = \frac{1}{2}(\boldsymbol{\mu}_2^T \mathbf{C}^{-1} \boldsymbol{\mu}_2 - \boldsymbol{\mu}_1^T \mathbf{C}^{-1} \boldsymbol{\mu}_1) \quad (1.38)$$

On the basis of Eq. (1.35), we may now describe the log-likelihood ratio test for our two-class problem as follows:

*If the output  $y$  of the linear combiner (including the bias  $b$ ) is positive, assign the observation vector  $\mathbf{x}$  to class  $C1$ . Otherwise, assign it to class  $C2$ .*

$\mu_1$  and  $\mu_2$  the covariance matrix  $\mathbf{C}$ .

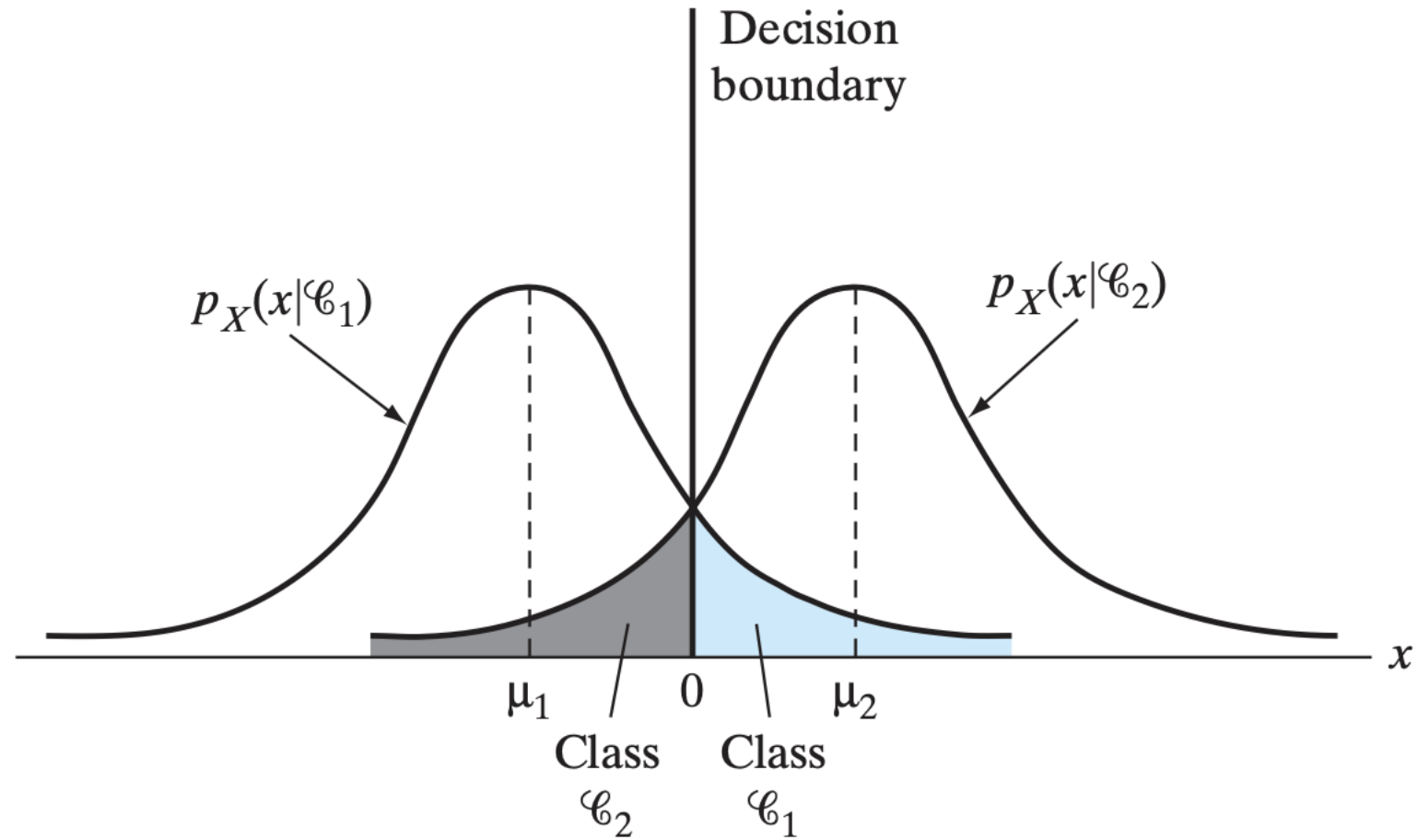


Fig: Two overlapping, one-dimensional Gaussian distributions.

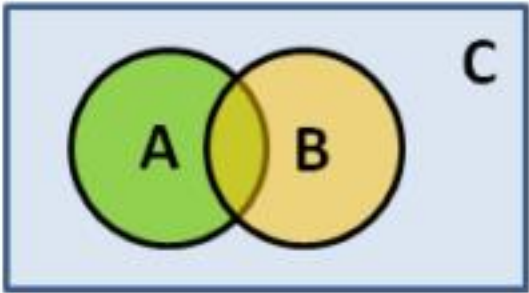
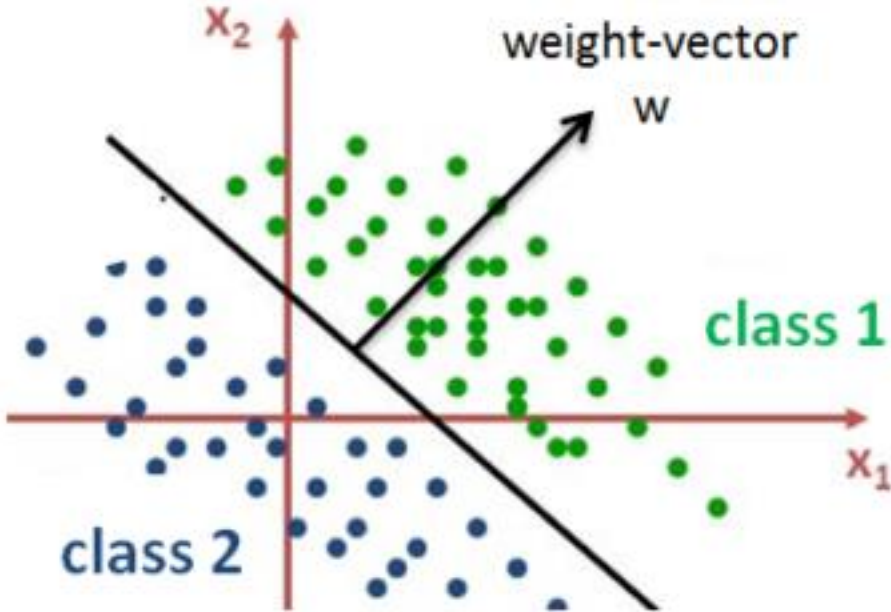
# Perceptron vs. Bayes Classifier

- **Bayes Classifier**

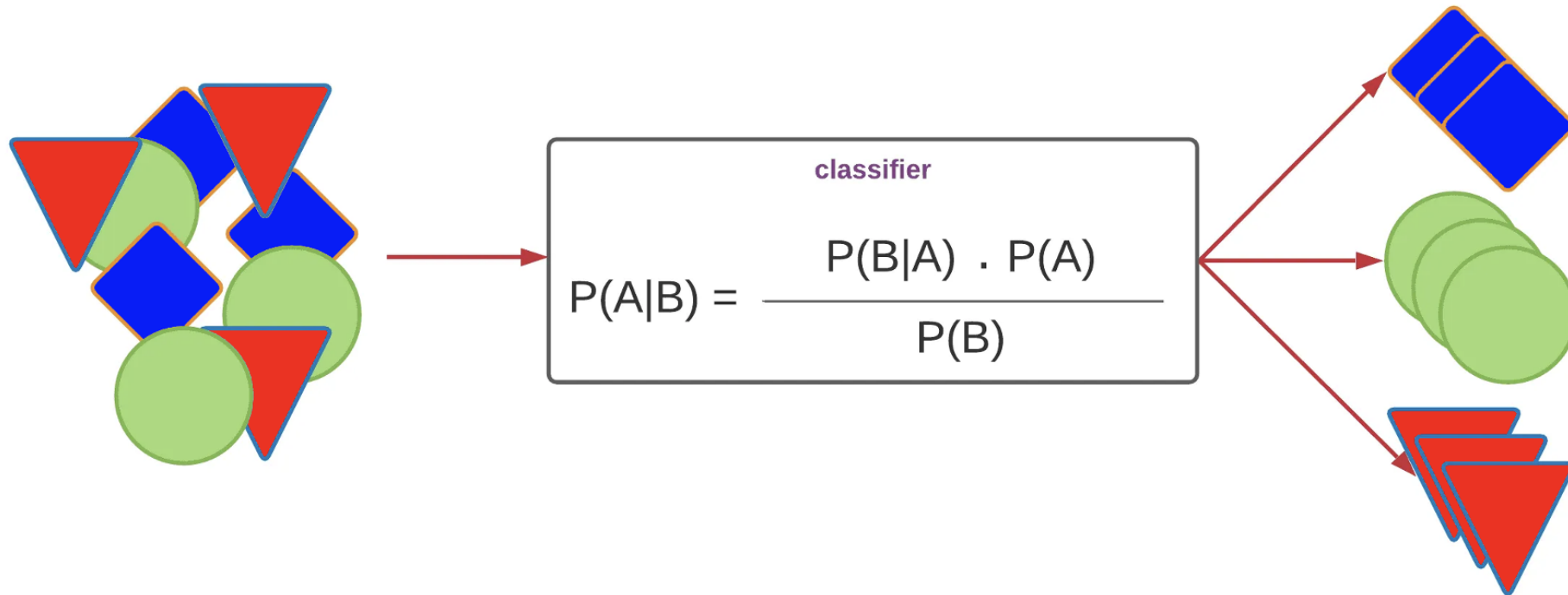
- Minimizes classification error probability, regardless of overlap between class distributions.
- Decision boundary lies where class distributions intersect (e.g., Gaussians).
- **Parametric**: Assumes specific distribution form (e.g., Gaussian), limiting applicability.
- Fixed design; can be made adaptive but requires more storage and computation.

- **Perceptron Convergence Algorithm:**

- **Nonparametric**: Makes no assumptions about input distribution.
- Focuses on classification errors, especially in overlapping regions.
- Performs well with nonlinear or non-Gaussian input distributions.
- **Adaptive and simple**: Requires only storage for weights and bias.

Statistical Approach (Naive Bayes)	Geometrical Approach (Perceptron, SVM)
$P(A B) = \frac{P(B A)P(A)}{P(B)}$  <p>A Venn diagram illustrating set theory. It features a light blue rectangular universal set labeled 'C'. Inside 'C', there are two overlapping circles: a green circle labeled 'A' and a yellow circle labeled 'B'. The intersection of 'A' and 'B' is shaded in a darker green/yellow color.</p>	 <p>A scatter plot illustrating a geometrical approach to classification. The plot has two axes, <math>x_1</math> and <math>x_2</math>, shown as red arrows. Two classes of data points are plotted: 'class 1' (green dots) and 'class 2' (blue dots). A black line represents the decision boundary separating the two classes. A black arrow labeled 'weight-vector w' originates from the decision boundary and points towards the region containing 'class 1'.</p>

# Naive Bayes Classifier



# THE BATCH PERCEPTRON ALGORITHM

1. *Introduces the generalized form of a perceptron cost function.*
2. *Uses the cost function to formulate a batch version of the perceptron convergence algorithm.*

we define the *perceptron cost function* as

$$J(\mathbf{w}) = \sum_{\mathbf{x}(n) \in \mathcal{H}} (-\mathbf{w}^T \mathbf{x}(n) d(n)) \quad (1.39)$$

where  $\mathcal{H}$  is the set of samples  $\mathbf{x}$  *misclassified* by a perceptron using  $\mathbf{w}$  as its weight vector.

- If all the samples are classified correctly, then the set  $\mathcal{H}$  is empty, in which case the cost function  $J(\mathbf{w})$  is zero.
- In any event, the nice feature of the cost function  $J(\mathbf{w})$  is that it is *differentiable* with respect to the weight vector  $\mathbf{w}$ .



Differentiating  $J(\mathbf{w})$  with respect to  $\mathbf{w}$  yields the *gradient vector*,

$$\nabla J(\mathbf{w}) = \sum_{\mathbf{x}(n) \in \mathcal{X}} (-\mathbf{x}(n)d(n)) \quad (1.40)$$

where the *gradient operator*

$$\nabla = \left[ \frac{\partial}{\partial w_1}, \frac{\partial}{\partial w_2}, \dots, \frac{\partial}{\partial w_m} \right]^T \quad (1.41)$$

In the *method of steepest descent*, the adjustment to the weight vector  $\mathbf{w}$  at each time-step of the algorithm is applied in a direction *opposite* to the gradient vector  $\nabla J(\mathbf{w})$ . Accordingly, the algorithm takes the form

$$\begin{aligned} \mathbf{w}(n+1) &= \mathbf{w}(n) - \eta(n) \nabla J(\mathbf{w}) \\ &= \mathbf{w}(n) + \eta(n) \sum_{\mathbf{x}(n) \in \mathcal{X}} \mathbf{x}(n)d(n) \end{aligned} \quad (1.42)$$

which includes the single-sample correction version of the perceptron convergence algorithm as a special case.

# Applications and Experiment

- Simple pattern classification
- Handwritten digit recognition (basic)
- Basis for further neural network models

# Summary

- Perceptron is a foundational model for neural networks
- Perceptron convergence theorem ensures learning under linear separability
- Limitations led to the evolution of deep networks

**Thank You!!!**