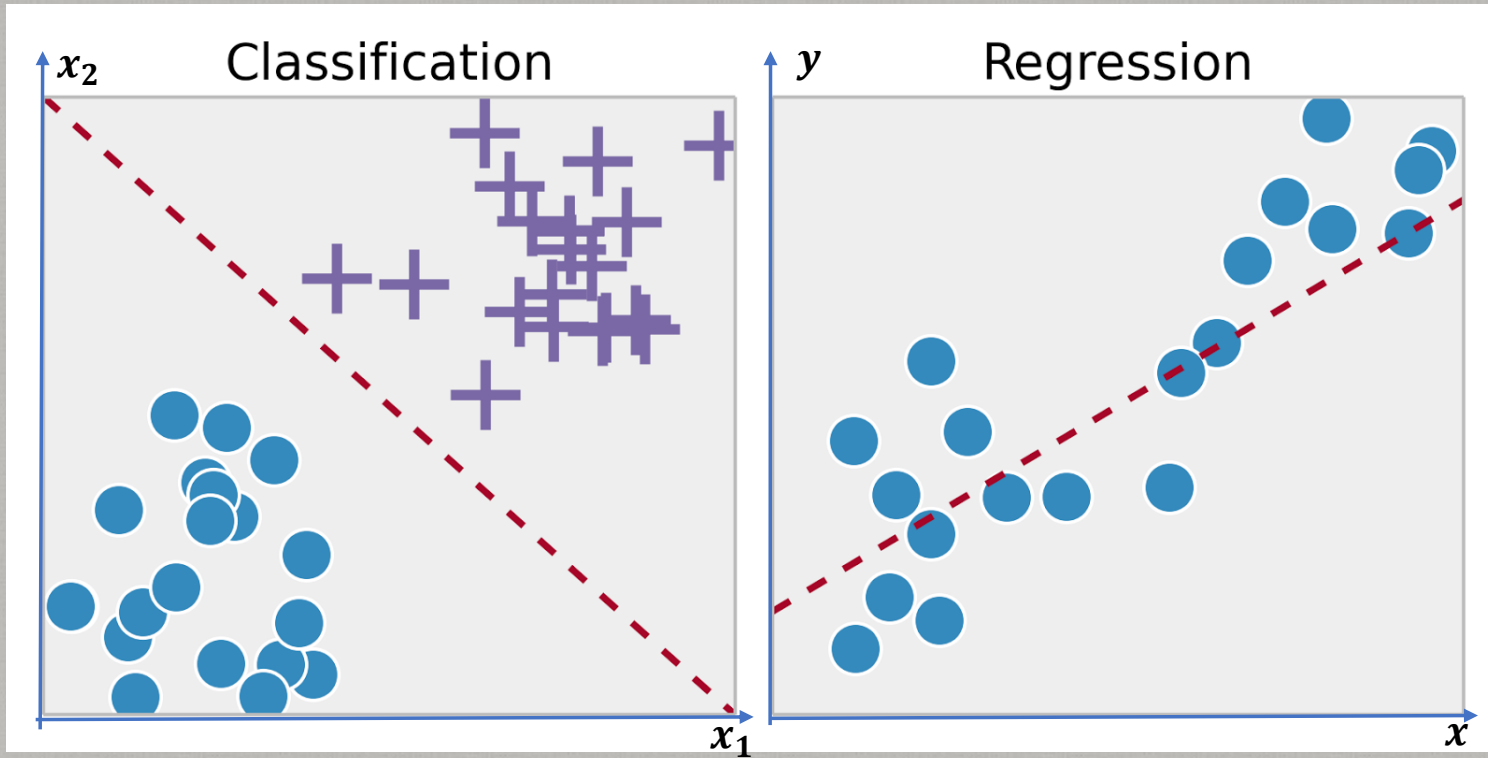


Regression vs Classification



Predict a nominal class label

Predict a continuous quantity

The Classification Problems

- ***The classification problems:*** given a training data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ of n points, where each data points $\mathbf{x}_i \in \mathcal{R}^d$ is paired with a known discrete class label y_i which represent one of K discrete classes \mathcal{C}_k , where $k = 1, 2, \dots, K$. The goal is to train a classifier which, when fed any arbitrary d -dimensional data point, classifies that data point as one of the K discrete classes.

■ **Target values in classification problems:**

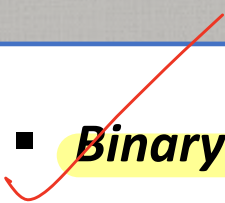
- In a binary classification problem, $y \in \{0,1\}$ or $\{-1,1\}$, such that $y = 1$ represents class \mathcal{C}_1 and $y = 0$, or ($y = -1$) represents class \mathcal{C}_2
- In multi-class classification, it is convenient to use a **1-of-K (one-hot) coding scheme** where \mathbf{y} becomes a vector of length K such that if the class is \mathcal{C}_j then all the elements of \mathbf{y} are zero except the element y_j , which takes the value of 1. For instance, let $K = 4$ classes, the pattern from class 2 would be given the target vector

$$\mathbf{y} = [0 \quad 1 \quad 0 \quad 0]^T$$

- The value of y_j can be interpreted as the probability that the class is \mathcal{C}_j .

❑ Target values in classification problems:

- We prefer one-hot encoding over representing $y = 1, 2, 3, 4$ (**ordinal values**) for C_1, C_2, C_3 and C_4 (**nominal values**) is that the later implies order as $4 > 3 > 2 > 1$.
 - Let the target value of an input x be $y = 4$. Classifier#1 classifies the input as $y_1 = 1$, classifier#2 classifies the same input as $y_2 = 3$.
 - Since $y - y_1 = 3 > y - y_2 = 1$, this implies that classifier#2 is better than classifier#1. This is not true. Both classifiers are equally bad in the classification problem.
- **Mean squared error** is no longer a popular loss function in classification problem.

- 
- **Binary classification (dichotomizer) ($K = 2$)** and **multi-class classification ($K > 2$)**
 - Many applications are binary classification problems
 - Some classification methods can only handle binary classification problems (Perceptron, SVM, etc.)
 - Some classification methods can naturally handle multi-class classification problems (nearest neighbor, ANN, decision trees, etc.)
 - In many cases, a multi-class classification problem can be decomposed as multiple binary classification problems.

■ ***Binary classification for multi-class classification***

- Multi-class classification problem can be turned into multiple binary classification problems.

- One-vs-Rest (OvR)

For example, a multi-class classification problem with three classes: “red”, “blue” and “green”. This could be divided into three binary classification problems as follows:

Binary classifier #1: red vs (blue and green)

Binary classifier #2: blue vs (red and green)

Binary classifier #3: green vs (red and blue)

- One-vs-One (OvO)

Using the same example, the problem is split into multiple binary classification problems as follows:

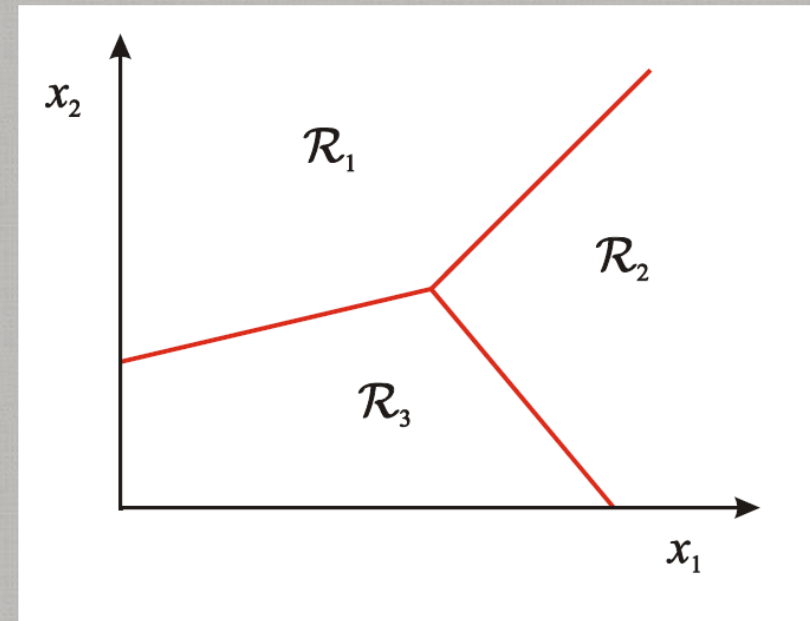
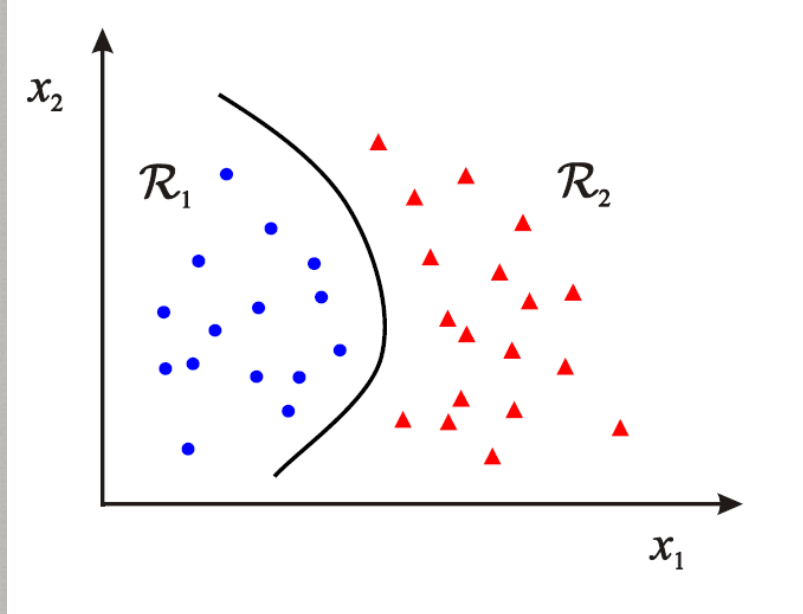
Binary classifier #1: red vs blue

Binary classifier #2: red vs green

Binary classifier#3: blue vs green

Decision regions and Decision boundaries

- In most common scenario, the classes are taken to be disjoint so that **each input is assigned to one and only one class (no class label ambiguity!)** 歧义
- The input space is thereby divided into **decision regions R_k** such that a point falling in R_k is assigned to class C_k . The boundaries between these decision regions are called **decision boundaries**.



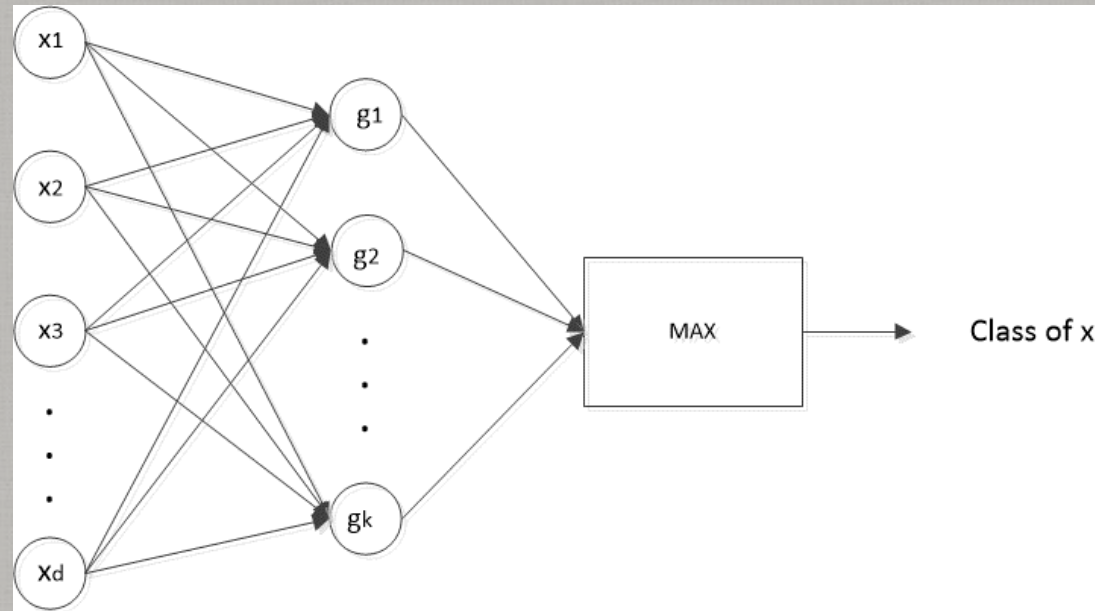
Discriminant Functions 判别函数

- One of the most useful way to represent classifiers is in terms of a set of **discriminant functions** $g_i(\mathbf{x}), i = 1, \dots, K$.
- The classifier is said to assign a feature vector \mathbf{x} to class \mathcal{C}_i if

$$g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \text{for all } j \neq i$$

Discriminant Functions

- The classifier is viewed as a network that computes K discriminant functions and selects the category corresponding to the largest discriminant.



- One popular example:

$$g_i(\mathbf{x}) = P(\mathcal{C}_i|\mathbf{x})$$

$$P(\mathcal{C}_i|\mathbf{x}) = \frac{P(\mathbf{x}|\mathcal{C}_i)P(\mathcal{C}_i)}{P(\mathbf{x})}$$

$$P(\mathbf{x}) = \sum_{j=1}^n P(\mathbf{x}|\mathcal{C}_j)P(\mathcal{C}_j)$$

The maximum discriminant function becomes the maximum posterior probability

后验概率

- Discriminant function is not unique!**

单调

$g_i(\mathbf{x})$ can be replaced by $f(g_i(\mathbf{x}))$ if $f(\cdot)$ is **monotonically increasing function!**

$$g_i(\mathbf{x}) = P(\mathcal{C}_i|\mathbf{x}) = \frac{P(\mathbf{x}|\mathcal{C}_i)P(\mathcal{C}_i)}{\sum_{i=1}^K P(\mathbf{x}|\mathcal{C}_i)P(\mathcal{C}_i)}$$

$$g_i(\mathbf{x}) = P(\mathbf{x}|\mathcal{C}_i)P(\mathcal{C}_i)$$

$$g_i(\mathbf{x}) = \ln P(\mathbf{x}|\mathcal{C}_i) + \ln P(\mathcal{C}_i)$$

- ***Discriminant function for two-class case:***

Instead of using two discriminant functions g_1 and g_2 and assigning \mathbf{x} to \mathcal{C}_1 if $g_1(\mathbf{x}) > g_2(\mathbf{x})$, it is more common to define a single discriminant function

$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

And to use the following decision rule:

Assign \mathbf{x} to \mathcal{C}_1 if $g(\mathbf{x}) > 0$, otherwise assign \mathbf{x} to \mathcal{C}_2

Examples:

$$g(\mathbf{x}) = P(\mathcal{C}_1|\mathbf{x}) - P(\mathcal{C}_2|\mathbf{x})$$

$$g(\mathbf{x}) = \ln \left(\frac{P(\mathbf{x}|\mathcal{C}_1)}{P(\mathbf{x}|\mathcal{C}_2)} \right) + \ln \left(\frac{P(\mathcal{C}_1)}{P(\mathcal{C}_2)} \right)$$

Performance Metrics for Classifiers

混淆矩阵

- **Confusion Matrix:** a natural way to represent classification results.

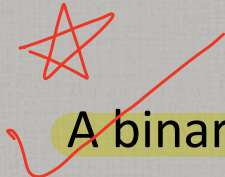
A binary classifier case:

| classifier | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | positive | 10 | 2 |
| | negative | 3 | 5 |

- 13 positive instances and 7 negative instances
- 10 positive instances are classified as positive
- 3 positive instances are classified as negative
- 5 negative instances are classified as negative
- 2 negative instances are classified as positive
- Diagonal numbers represent number of instance classified correctly
- Off-diagonal numbers represent classification errors

- A multi-class classifier case: (iris data set)

| classifier | | Actual class | | |
|-----------------|------------|--------------|------------|-----------|
| | | setosa | versicolor | virginica |
| Predicted class | setosa | 13 | 0 | 0 |
| | versicolor | 0 | 10 | 6 |
| | virginica | 0 | 0 | 9 |



A binary classifier case:

| Classifier | | Actual class | |
|-----------------|----------|--------------|---------------------------|
| | | positive | negative |
| Predicted class | Positive | TP (hit) | FP (false alarm) |
| | negative | FN (miss) | TN (correct rejection) |

- **Condition positive (P):** the number of real positive instances in the data.
- **Condition negative (N):** the number of real negative instances in the data.
- **True positive (TP) (hit):** the number of real positive instances that are classified as positive.
- **True negative (TN) (correct rejection):** the number of negative instances that are classified as negative.
- **False positive (FP) (false alarm) (type I error):** number of real negative instances that are classified as positive.
- **False negative (FN) (miss) (type II error):** number of positive instances that are classified as negative.

Accuracy and Error Rate:

| Classifier | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | TP | FP |
| | negative | FN | TN |

- **Accuracy** = $\frac{\text{\# of instances classified correctly}}{\text{\# of instances tested}}$

$$= \frac{TP + TN}{TP + FP + FN + TN}$$

- **Error rate** = $1 - \text{Accuracy}$

$$= \frac{FP + FN}{TP + FP + FN + TN}$$

Example: Find the accuracy of the following two classifiers

| Classifier#1 | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | 90 | 10 |
| | negative | 10 | 90 |

$$\text{Accuracy} = \frac{90+90}{200} = 90\%$$

| Classifier#2 | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | 10 | 10 |
| | negative | 10 | 170 |

$$\text{Accuracy} = \frac{10+170}{200} = 90\%$$

Problems with Accuracy and Error rate

- ***The problem with class imbalance:***

Accuracy simply treat all instances the same. Accuracy maybe fine if you deal with balanced (or approximated balanced) data sets. It can be misleading when deal with imbalanced data sets.

- In real applications domain, imbalanced data sets are the norm!
For example, cancer diagnose, intrusion detection, defects detection, etc.
The positive examples to negative examples ratio can be 1:10 or less!
- Need more metrics for classifier performance evaluation!!!

Precision

| Classifier | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | TP | FP |
| | negative | FN | TN |

- **Precision** = $\frac{TP}{TP+FP}$

- Precision measures what portion of instances that are classified as positive are actually positive.
- Precision gives us information about its performance with respect to False Positive (false alarm).
- If the focus of the application is on **minimizing False Positive (false alarm)**, we should make precision close to 100%.
- For example, spam email detector.

| Classifier#1 | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | 90 | 10 |
| | negative | 10 | 90 |

$$\text{Accuracy} = \frac{90+90}{200} = 90\%$$

$$\text{Precision} = \frac{90}{90+10} = 90\%$$

| Classifier#2 | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | 10 | 10 |
| | negative | 10 | 170 |

$$\text{Accuracy} = \frac{10+170}{200} = 90\%$$

$$\text{Precision} = \frac{10}{10+10} = 50\%$$

Recall or Sensitivity

| Classifier | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | TP | FP |
| | negative | FN | TN |

- **Recall (sensitivity)** = $\frac{TP}{TP+FN}$
- **Recall** measures what portion of positive instances that are classified as positive
- **Recall** gives us information about a classifier performance with respect to **False Negative (miss)**
- *If the application focus more on minimizing False Negative (miss), we should make recall close to 100% and without precision being too bad.*
- *For example, cancer diagnose, information retrieval*

Performance comparison of two classifiers

| Classifier#2 | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | 10 | 10 |
| | negative | 10 | 170 |

$$\text{Recall} = \frac{10}{10+10} = 50\%$$

$$\text{precision} = \frac{10}{10+10} = 50\%$$

$$\text{Accuracy} = \frac{10+170}{200} = 90\%$$

| Classifier#2 | | Actual class | |
|-----------------|----------|--------------|----------|
| | | positive | negative |
| Predicted class | Positive | 20 | 20 |
| | negative | 0 | 160 |

$$\text{Recall} = \frac{20}{15+5} = 100\%$$

$$\text{precision} = \frac{20}{20+20} = 50\%$$

$$\text{Accuracy} = \frac{20+160}{200} = 90\%$$

F-score: a combination of precision and recall

- F_1 score (balanced F score) is the *harmonic mean* of precision and recall

$$F_1 = \frac{2}{recall^{-1} + precision^{-1}} = \frac{2 * precision * recall}{precision + recall}$$

- F_β score: $\beta > 0$ is chosen such that recall is considered β times as important as precision

$$F_\beta = (1 + \beta^2) \frac{precision * recall}{(\beta^2) * precision + recall}$$

Other metrics

- **True positive rate (TPR)(hit rate)** $= \frac{TP}{TP+FN} = \frac{TP}{\text{total positives}} = \text{recall} = \text{sensitivity}$
 $= 1 - \text{false negative rate (FNR)(miss rate)}$

- **False positive rate (FPR)(false alarm rate)** $= \frac{FP}{FP+TN} = \frac{FP}{\text{total negatives}}$
 $= 1 - \text{true negative rate (TNR)}$

- **Specificity = true negative rate (TNR)** $= \frac{TN}{FP+TN}$

Generative vs. Discriminative Classifier

- There are two main types of classification models: **generative models and discriminative models**
- Posterior probability plays critical roles in classification problems
- **Generative models** have strong roots in probabilistic modeling. As in Bayes Theorem, generative models **explicitly** use prior probability and likelihood to calculate the posterior probabilities

$$\hat{y} = \underset{i}{\operatorname{argmax}} P(\mathcal{C}_i|\mathbf{x}) = \underset{i}{\operatorname{argmax}} P(\mathbf{x}|\mathcal{C}_i)P(\mathcal{C}_i)$$

e.g., Naïve Bayes, Hidden Markov Models (HMMs), Bayesian Networks, etc.

Generative vs. Discriminative Classifier

- On the other hand, in ***discriminative models***, we bypass learning a generative model altogether and directly learn a decision boundary.
- **Discriminative models** are parameterized by weights that either form a posterior distribution without considering the prior or conditional distribution or directly form a hard decision boundary without considering any probabilities in the first place.

e.g., Perceptron, logistic regression, SVM, artificial neural networks (ANN), nearest neighbor, etc.

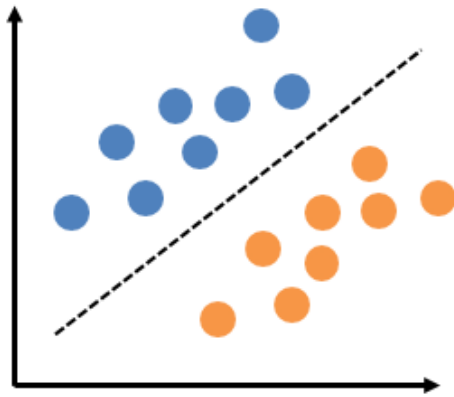
Linearly separable data

Binary linearly separable case:

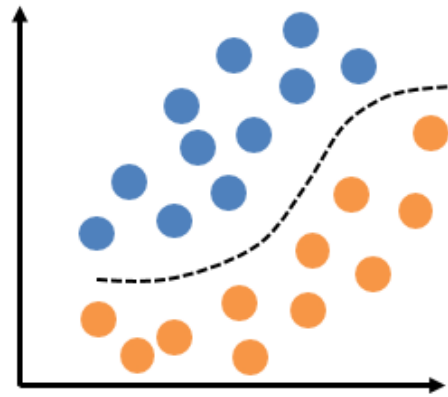
- Let's assume that we have two sets of points on **2D Euclidian plane**, one set of points are colored blue and other colored red.
- These two sets are **linearly separable** if there exists at least one **line** in the plane with all the blue points on one side and all the red points on the other side.
- This definition generalizes to **higher-dimensional Euclidian spaces** if **line** is replaced by **hyperplane** 超平面

Binary Linearly separable data sets

Linear

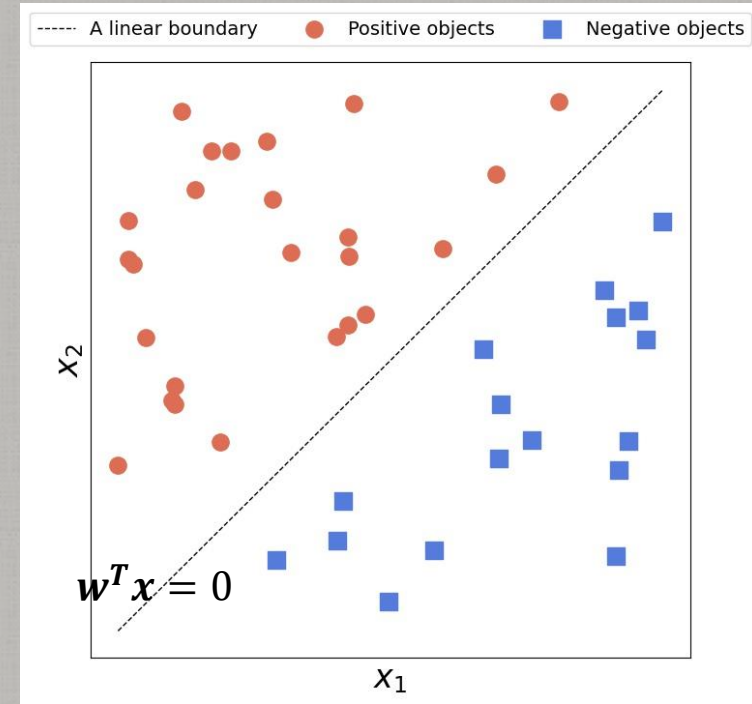


Nonlinear



Equation of the linear boundary

- The equation of the hyperplane: $w^T x = 0$
 - Any point x on the hyperplane: $w^T x = 0$
 - Any point x above the hyperplane: $w^T x > 0$
 - Any point x below the hyperplane: $w^T x < 0$



The Perceptron Model (a discriminative model)

感知机

Binary linearly separable case:

- The **perceptron classifier** classifies an input x into one of two classes: C_1 and C_2 . The target values used in this problem is $y = 1$ for C_1 and $y = -1$ for C_2 .
- The **discriminant function** is in a **linear form**: $g(x) = w^T x$, ✓

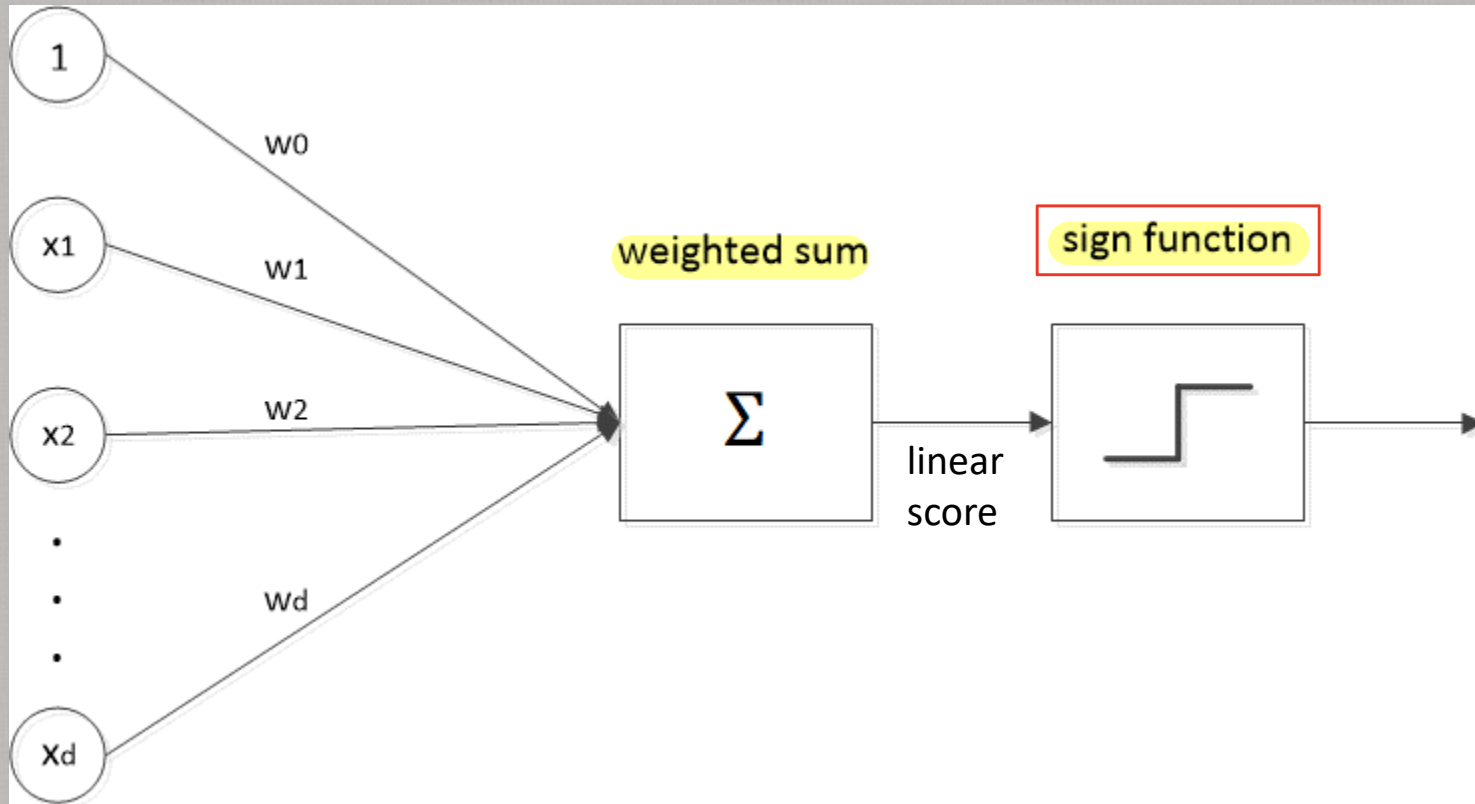
$$w^T x_i > 0, \quad x_i \text{ in } C_1$$

$$w^T x_i \leq 0, \quad x_i \text{ in } C_2$$

- w is the weight vector to be determined.
- $w^T x$ is called the **linear score** of the input x

The Perceptron Model - a graphical explanation

$$y(x) = f(w^T x)$$



Binary linearly separable case: ✓

- Notice that in this algorithm, we are seeking the weight vector \mathbf{w} such that for any input \mathbf{x}_i in the given data set,

$$\mathbf{w}^T \mathbf{x}_i > 0, \quad \mathbf{x}_i \text{ in } C_1 \ (y_i = 1)$$

$$\mathbf{w}^T \mathbf{x}_i \leq 0, \quad \mathbf{x}_i \text{ in } C_2 \ (y_i = -1)$$

This is equivalent to

$$\mathbf{w}^T \mathbf{x}_i y_i \geq 0, \quad \text{any } \mathbf{x}_i \in D$$

Where $y_i \in \{1, -1\}$ is the target value of the input \mathbf{x}_i . $\mathbf{w}^T \mathbf{x}_i y_i$ is the **normalized linear score** of input \mathbf{x}_i .



The training algorithm of Perceptron Model

How to find the weight vector that can completely classifies all the data points in the training data set.

准则

The Perceptron criterion : the error (loss) function:

Hence, we have the following perceptron criterion error function:

$$E_p(\mathbf{w}) = - \sum_{i \in M} \mathbf{w}^T \mathbf{x}_i y_i$$

Where, M denote the **set of all misclassified instances**.

The perceptron criterion associates zero error with any instance that is correctly classified; for a misclassified instance \mathbf{x}_i , it tries to minimize the quantity $-\mathbf{w}^T \mathbf{x}_i y_i$

Show that for any **misclassified** instance, $-\mathbf{w}^T \mathbf{x}_i y_i \geq 0$

- An instance \mathbf{x}_i is misclassified using the decision rule:

$$\mathbf{w}^T \mathbf{x}_i > 0, \quad \mathbf{x}_i \text{ in } C_1$$

$$\mathbf{w}^T \mathbf{x}_i \leq 0, \quad \mathbf{x}_i \text{ in } C_2$$

- Case 1: \mathbf{x}_i is a positive instance ($y_i = +1$) but classified as a negative instance. $\mathbf{w}^T \mathbf{x}_i \leq 0$, hence, $\mathbf{w}^T \mathbf{x}_i y_i \leq 0$
- Case 2: \mathbf{x}_i is a negative instance ($y_i = -1$) but classified as a positive instance. $\mathbf{w}^T \mathbf{x}_i > 0$, hence, $\mathbf{w}^T \mathbf{x}_i y_i < 0$



The Gradient Descent Procedure

To find the weight vector \mathbf{w} that minimize the *perceptron criterion error function* $E_p(\mathbf{w})$, we use the following **gradient descent procedure**:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta * \nabla E_p(\mathbf{w}^k)$$

Where,

- k is an integer that indexes the steps of the algorithm;
- η is a positive scale factor or learning rate that set the step size;
- $\nabla E_p(\mathbf{w})$ is the **gradient** vector of the *perceptron criterion error function*.

- To find the weight vector \mathbf{w} that minimize the *perceptron criterion error function* $E_p(\mathbf{w})$, we use the following **gradient descent procedure**:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \nabla E_p(\mathbf{w}^k)$$

- Gradient descent is based on the observation that if a multi-variable function $E_p(\mathbf{w})$ is defined and differentiable in a neighborhood of a point \mathbf{w} , then $E_p(\mathbf{w})$ decreases fastest if one goes from \mathbf{w} in the direction of the negative gradient of E_p at \mathbf{w} , $-\nabla E_p(\mathbf{w})$.
- It follows that if $\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \nabla E_p(\mathbf{w}^k)$, and for η **small enough**, one starts with a guess \mathbf{w}^0 for the local minimum of E_p , and consider the sequence $\mathbf{w}^0, \mathbf{w}^1, \mathbf{w}^2, \dots$, we will have a monotonic sequence

$$E_p(\mathbf{w}^0) \geq E_p(\mathbf{w}^1) \geq E_p(\mathbf{w}^2) \geq \dots, \quad \downarrow$$

so the sequence $\mathbf{w}^k, k \geq 0$ converges to the desired local minimum.

- To find the weight vector \mathbf{w} that minimize the *perceptron criterion error function* $E_p(\mathbf{w})$, we use the following **gradient descent procedure**:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \eta \nabla E_p(\mathbf{w}^k)$$

Where, k is an integer that indexes the steps of the algorithm; η^k is a positive scale factor or learning rate that set the step size; $\nabla E_p(\mathbf{w})$ is the gradient vector of the *perceptron criterion error function*.

$$\nabla E_p(\mathbf{w}) = - \sum_{i \in M} \mathbf{x}_i y_i$$

- Substitute this into weight update expression, we obtain:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + \eta \sum_{i \in M} \mathbf{x}_i y_i$$

The Perceptron training Algorithms

When should the weight vector be updated?

- **Batch Training Algorithm** 批训练
- **Sequential Training Algorithm** 顺序训练

Some terminologies about gradient descent algorithm

批量大小 { Batch Gradient Descent 全批量梯度下降, 批量大小等于整个训练集大小
Mini-batch Gradient Descent 小批量梯度下降, 批量大小介于1和训练集大小之间 (如32, 64)
Stochastic Gradient descent, SGD 随机梯度下降, 批量大小为1, 速度快, 方向不准, 收敛路径曲折

- **Batch size** is the number of training instances to work through before updating the weight vector 每次梯度更新时使用多个样本来计算梯度
- Training dataset can be divided into one or more batches
- An **epoch** is one cycle through the entire training dataset 训练周期
 - It usually takes more than a few epochs for a gradient descent algorithm to converge
- **Batch gradient descent:** batch size == N
- **Sequential (stochastic) gradient descent:** batch size == 1 + shuffling the training data set every epoch

Batch Training Algorithm:

Batch size = N

Begin initialize \mathbf{w} with random values, and $k(\# \text{ of epoch}) \leftarrow 0$

do $k \leftarrow k + 1$

$M_k = \{\}, j = 0$ (index of data points)

do $j \leftarrow j + 1$

Classify \mathbf{x}_j using current \mathbf{w}

if \mathbf{x}_j is misclassified then append \mathbf{x}_j, y_j to M_k

until $j = N$

$\mathbf{w} \leftarrow \mathbf{w} + \eta \sum_{i \in M_k} \mathbf{x}_i y_i$

until $M_k == \{\}$

return \mathbf{w}

End

在每次迭代中，批训练使用整个训练集的数据来计算误差并更新权重

Sequential Training Algorithm:

Begin initialize \mathbf{w} with random values, and $k \leftarrow 0$ (index of points in the training dataset)

do $k \leftarrow k \bmod N$

Classify \mathbf{x}_k using current \mathbf{w}

if \mathbf{x}_k is misclassified, $\mathbf{w} = \mathbf{w} + \eta \mathbf{x}_k y_k$

$k \leftarrow k + 1$

until all instances in the training data set are properly classified

return \mathbf{w}

End

顺序训练在每次迭代中使用
单个样本来更新权重

N is the size of the training dataset

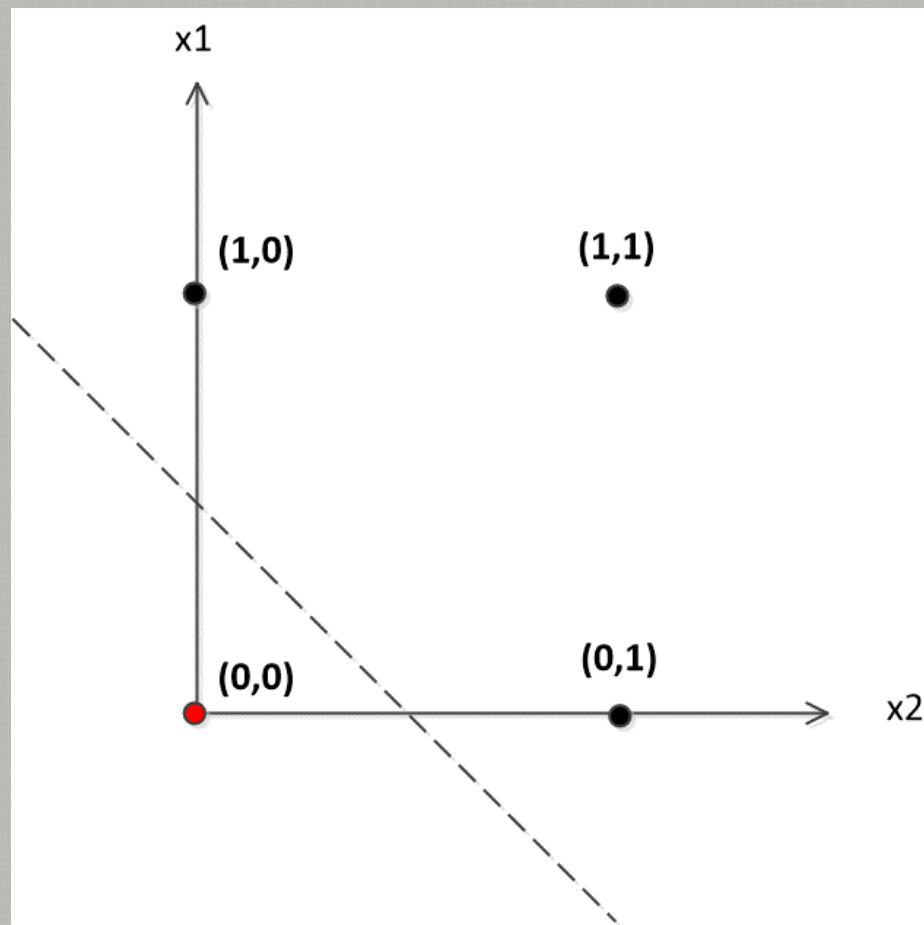
只要两个输入中至少有一个是1, 输出就为1

The logic **OR** problem: 或

$$y = x_1 \text{ or } x_2$$

$$x_1, x_2, y \in \{0,1\}$$

| x_1 | x_2 | y |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |



当两个输入不同时,输出为1, 当两个输入相同时,输出为0

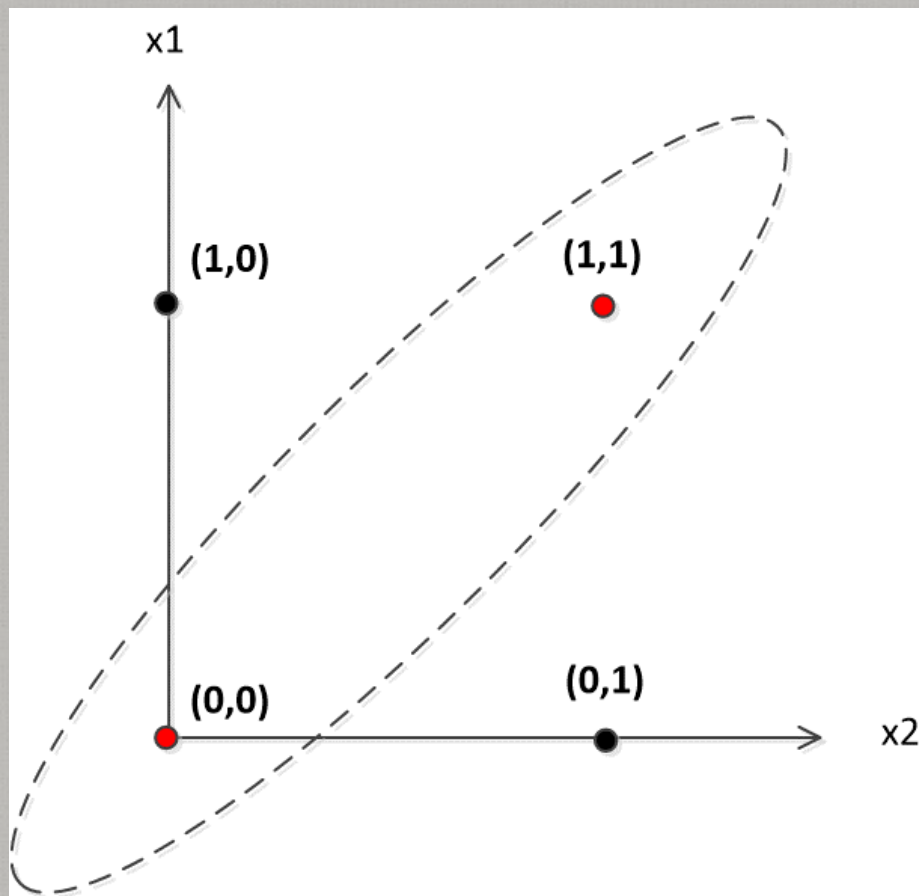
The logic **XOR** (exclusive **OR**) problem:

$$y = x_1 \oplus x_2$$

异或

$$x_1, x_2, y \in \{0,1\}$$

| x_1 | x_2 | y |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |





Training a Perceptron for the logic OR problem (Sequential training)

The weight vector is $\mathbf{w} = [w_0 \quad w_1 \quad w_2]^T$

the extended feature vector is $\mathbf{x} = [1 \quad x_1 \quad x_2]^T$

Let's set the initial value of \mathbf{w} as $\mathbf{w}^0 = [0 \quad 0 \quad 0]^T$

The learning rate is chosen as $\eta = 1$

Sequential weight update rule:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + \mathbf{x}_i y_i \quad \text{if } \mathbf{x}_i \text{ is misclassified}$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

符号函数

$$\text{sign}(x) = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases}$$

$$\vec{x}_i = (1, x_1, x_2)^T \quad \vec{w}^k = (w_0, w_1, w_2)^T$$

Step 1: feed the first training instance $\mathbf{x}_1 = [1 \ 0 \ 0]^T$, $y_1 = -1$

$$(\mathbf{w}^0)^T \mathbf{x}_1 = [0 \ 0 \ 0] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 0, \quad \hat{y}_1 = \text{sign}((\mathbf{w}^0)^T \mathbf{x}_1) = -1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^1 = \mathbf{w}^0$$

$$w = w + \eta x_k y_k$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step2: feed the second training instance $\mathbf{x}_2 = [1 \ 0 \ 1]^T$, $y_2 = 1$

$$(\mathbf{w}^1)^T \mathbf{x}_2 = [0 \ 0 \ 0] \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = 0, \quad \hat{y}_2 = \text{sign}((\mathbf{w}^1)^T \mathbf{x}_2) = -1, \quad \text{missclassified}$$

Then,

$$w^{k+1} = w^k + x_i y_i \quad \mathbf{w}^2 = \mathbf{w}^1 + \mathbf{x}_2 y_2 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step3: feed the third training instance $\mathbf{x}_3 = [1 \ 1 \ 0]^T$, $y_3 = 1$

$$(\mathbf{w}^2)^T \mathbf{x}_3 = [1 \ 0 \ 1] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 1, \quad \hat{y}_3 = \text{sign}((\mathbf{w}^2)^T \mathbf{x}_3) = 1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^3 = \mathbf{w}^2$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step4: feed the fourth training instance $\mathbf{x}_4 = [1 \ 1 \ 1]^T$, $y_4 = 1$

$$\hat{y}_4(\mathbf{w}^3)^T \mathbf{x}_4 = [1 \ 0 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 2, \quad \hat{y}_4 = \text{sign}((\mathbf{w}^3)^T \mathbf{x}_4) = 1, \quad \text{correct}$$

Then,

$$\mathbf{w}^4 = \mathbf{w}^3$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step5: feed the first training instance $\mathbf{x}_1 = [1 \ 0 \ 0]^T$, $y_1 = -1$

$$(\mathbf{w}^4)^T \mathbf{x}_1 = [1 \ 0 \ 1] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 1, \quad \hat{y}_1 = \text{sign}((\mathbf{w}^4)^T \mathbf{x}_1) = 1, \quad \text{missclassified}$$

Then,

$$\mathbf{w}^5 = \mathbf{w}^4 + \mathbf{x}_1 y_1 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step6: feed the second training instance $\mathbf{x}_2 = [1 \ 0 \ 1]^T$, $y_2 = 1$

$$(\mathbf{w}^5)^T \mathbf{x}_2 = [0 \ 0 \ 1] \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = 1, \quad \hat{y}_2 = \text{sign}\left((\mathbf{w}^5)^T \mathbf{x}_2\right) = 1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^6 = \mathbf{w}^5$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step7: feed the third training instance $\mathbf{x}_3 = [1 \ 1 \ 0]^T$, $y_3 = 1$

$$(\mathbf{w}^6)^T \mathbf{x}_3 = [0 \ 0 \ 1] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 0, \quad \hat{y}_3 = \text{sign}((\mathbf{w}^6)^T \mathbf{x}_3) = -1, \quad \text{missclassified}$$

Then,

$$\mathbf{w}^7 = \mathbf{w}^6 + \mathbf{x}_3 y_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step8: feed the fourth training instance $\mathbf{x}_4 = [1 \ 1 \ 1]^T$, $y_4 = 1$

$$(\mathbf{w}^7)^T \mathbf{x}_4 = [1 \ 1 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 3, \quad \hat{y}_4 = \text{sign}((\mathbf{w}^7)^T \mathbf{x}_4) = 1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^8 = \mathbf{w}^7$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step9: feed the first training instance $\mathbf{x}_1 = [1 \ 0 \ 0]^T$, $y_1 = -1$

$$(\mathbf{w}^8)^T \mathbf{x}_1 = [1 \ 1 \ 1] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 1, \quad \hat{y}_1 = \text{sign}((\mathbf{w}^8)^T \mathbf{x}_1) = 1, \quad \text{missclassified}$$

Then,

$$\mathbf{w}^9 = \mathbf{w}^8 + \mathbf{x}_1 y_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step10: feed the second training instance $\mathbf{x}_2 = [1 \ 0 \ 1]^T, y_2 = 1$

$$(\mathbf{w}^9)^T \mathbf{x}_2 = [0 \ 1 \ 1] \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = 1, \quad \hat{y}_2 = \text{sign}((\mathbf{w}^9)^T \mathbf{x}_2) = 1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^{10} = \mathbf{w}^9$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step11: feed the third training instance $\mathbf{x}_3 = [1 \ 1 \ 0]^T$, $y_3 = 1$

$$(\mathbf{w}^{10})^T \mathbf{x}_3 = [0 \ 1 \ 1] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 1, \quad \hat{y}_3 = \text{sign}((\mathbf{w}^{10})^T \mathbf{x}_3) = 1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^{11} = \mathbf{w}^{10}$$

| x_1 | x_2 | y |
|-------|-------|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step12: feed the fourth training instance $\mathbf{x}_4 = [1 \ 1 \ 1]^T, y_4 = 1$

$$(\mathbf{w}^{11})^T \mathbf{x}_4 = [0 \ 1 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = 2, \quad \hat{y}_4 = \text{sign}((\mathbf{w}^{11})^T \mathbf{x}_4) = 1, \quad \text{correct}$$

Then, no weight update

$$\mathbf{w}^{12} = \mathbf{w}^{11}$$

| x1 | x2 | y |
|----|----|--------|
| 0 | 0 | 0 (-1) |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

Step13: feed the first training instance $\mathbf{x}_1 = [1 \ 0 \ 0]^T$, $y_1 = -1$

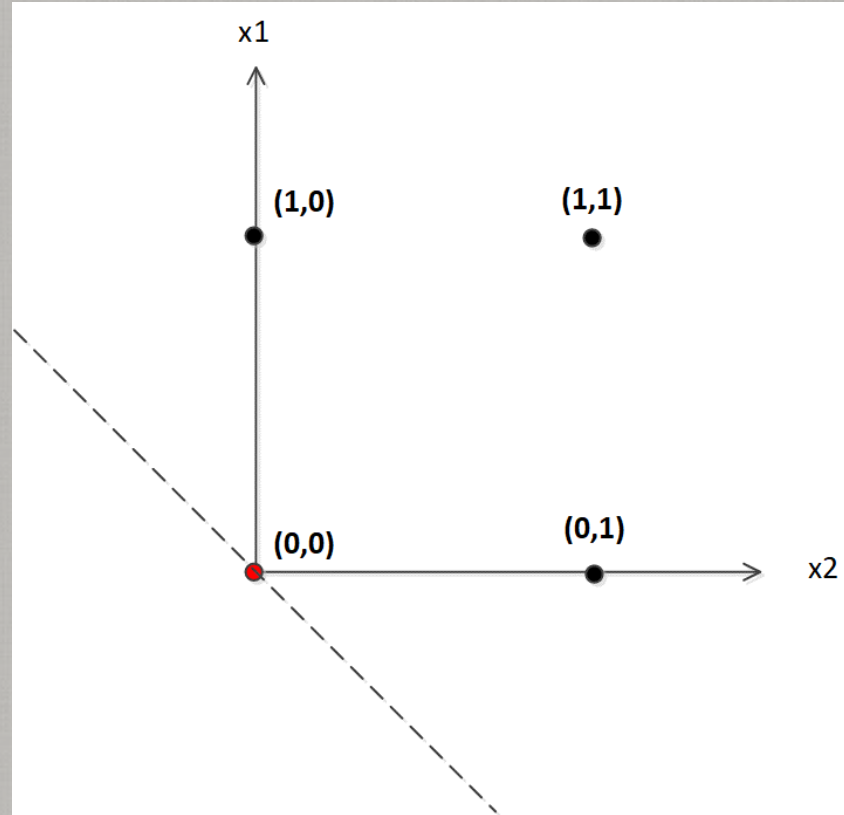
$$(\mathbf{w}^{12})^T \mathbf{x}_1 = [0 \ 1 \ 1] \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = 0, \quad \hat{y}_1 = \text{sign}((\mathbf{w}^{12})^T \mathbf{x}_1) = -1, \quad \text{correct}$$

Now, We have all the instances classified correctly using the weight:

$$\mathbf{w} = [0 \ 1 \ 1]^T$$

i.e.,

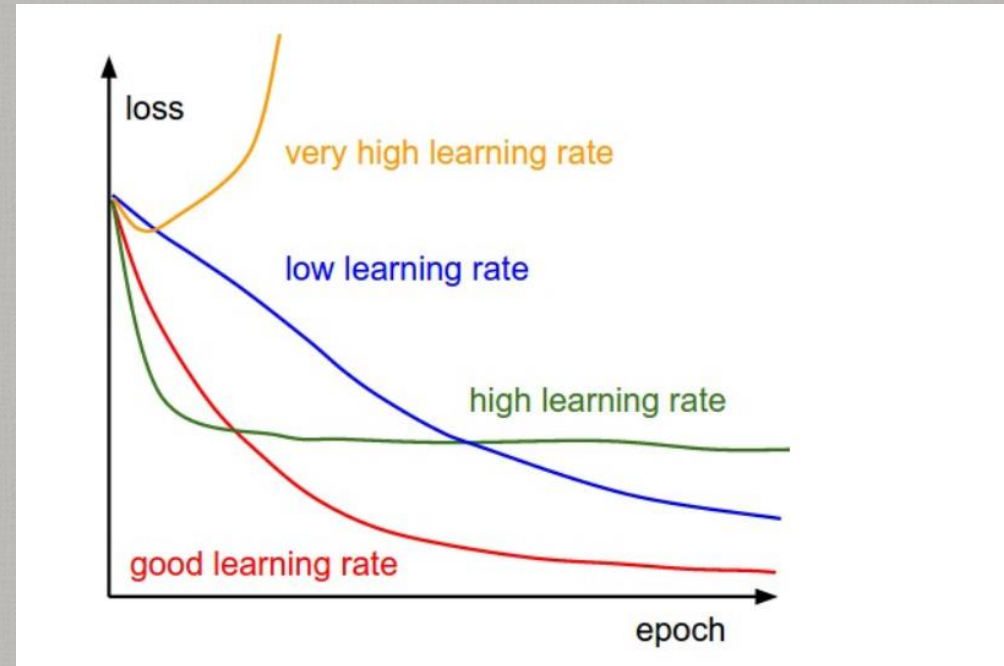
$$\mathbf{w}^T \mathbf{x} = x_1 + x_2$$





Learning rate in the gradient descent algorithm

- if the learning rate η is too small, it would take long time for the weights to converge, and the algorithm becomes computationally expensive
- If the learning rate η is too large, the weights may fail to converge and overshoot the minimum
- The most commonly used learning rates are: 0.001, 0.01, 0.1
- Fixed learning rate vs changing learning rate



Binary linearly separable case:

- Notice that in this algorithm, we are seeking the weight vector \mathbf{w} such that for an instance \mathbf{x}_i

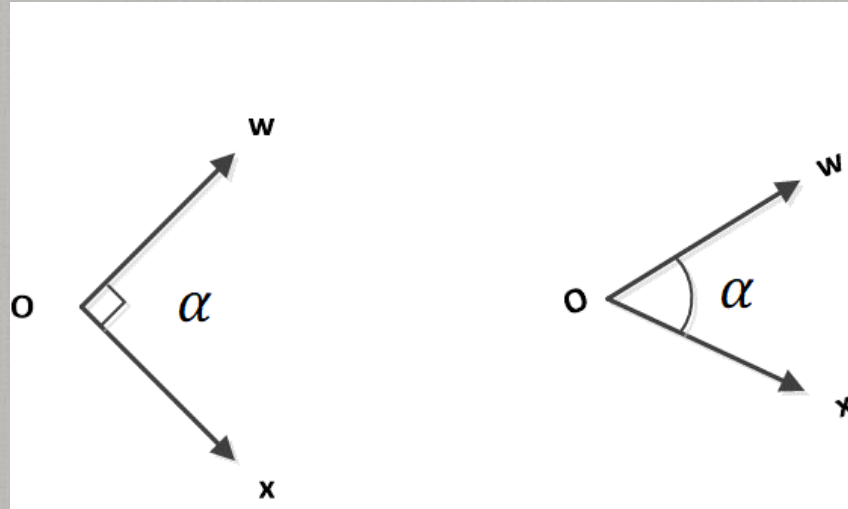
$$\mathbf{w}^T \mathbf{x}_i > 0, \quad \mathbf{x}_i \text{ in } C_1 \quad (y_i = 1)$$

$$\mathbf{w}^T \mathbf{x}_i \leq 0, \quad \mathbf{x}_i \text{ in } C_2 \quad (y_i = -1)$$

This is equivalent to

$$\mathbf{w}^T \mathbf{x}_i y_i \geq 0, \quad \text{any } \mathbf{x}_i$$

Where $y_i \in \{1, -1\}$ is the target value of the input \mathbf{x}_i . $\mathbf{w}^T \mathbf{x}_i y_i$ is the **normalized linear score** of input \mathbf{x}_i .



- Notice that

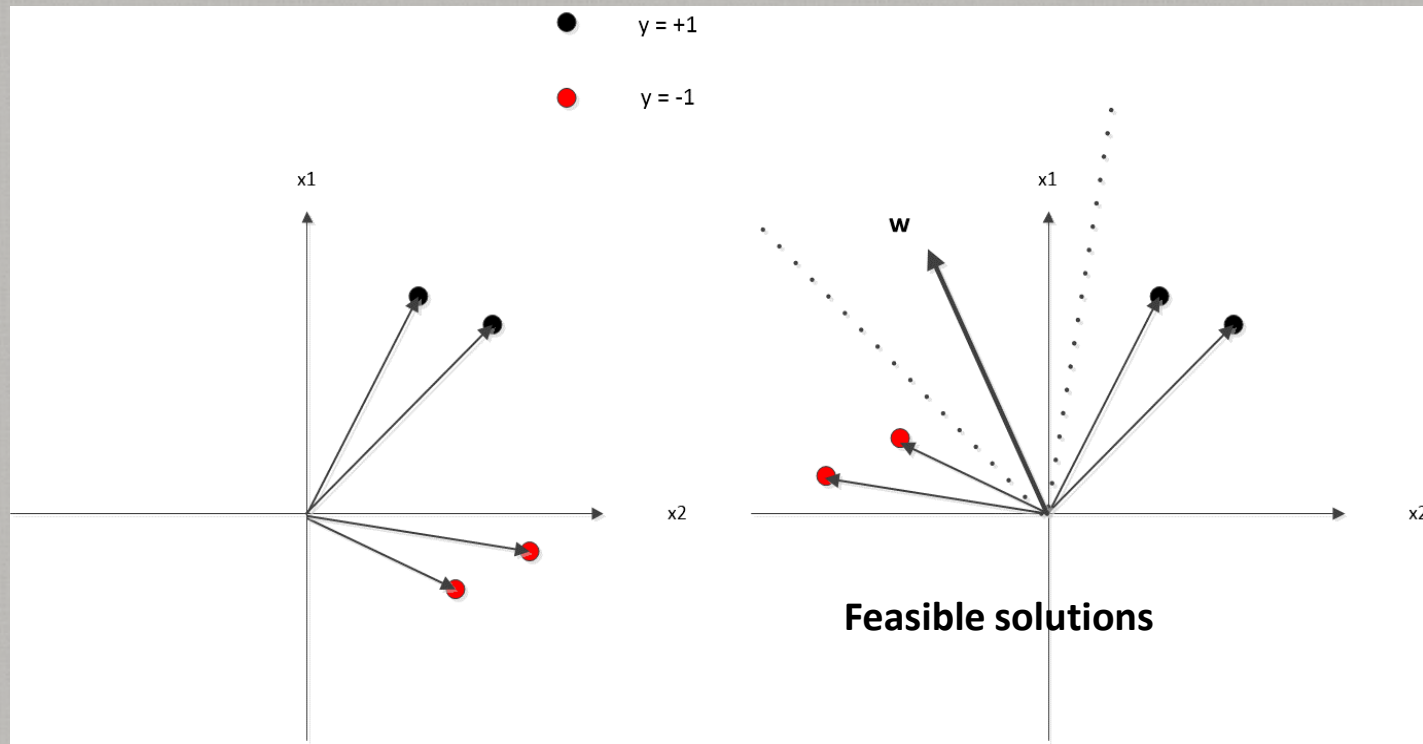
$$\mathbf{w}^T \mathbf{x} = \|\mathbf{w}\| \|\mathbf{x}\| \cos \alpha$$

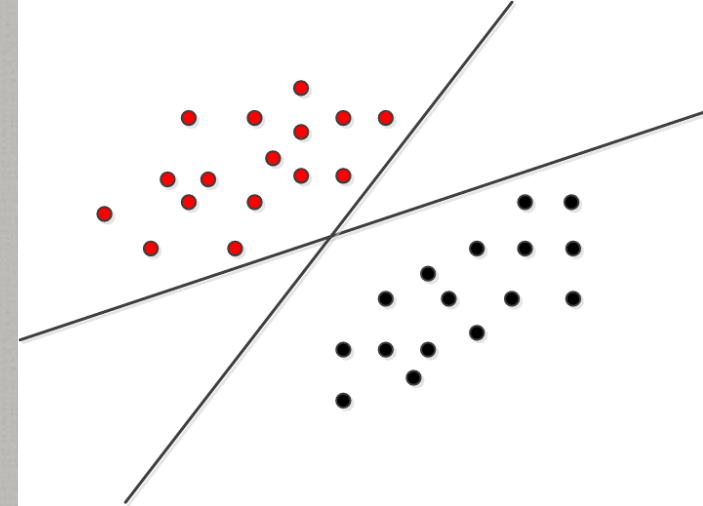
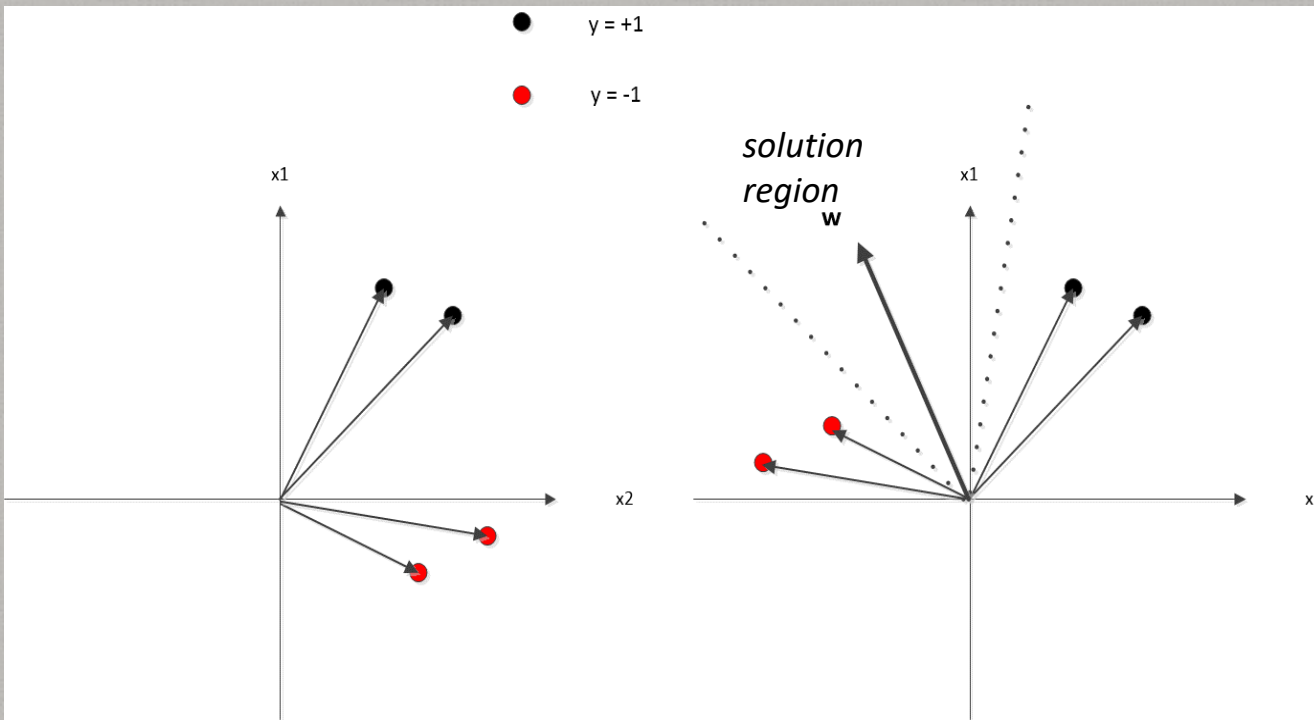
- When $\alpha = \frac{\pi}{2}$, we have $\mathbf{w}^T \mathbf{x} = 0$, \mathbf{w} and \mathbf{x} are orthogonal to each other.
- When $|\alpha| < \frac{\pi}{2}$, we have $\mathbf{w}^T \mathbf{x} > 0$

- we are seeking the weight vector \mathbf{w} such that for an instance \mathbf{x}_i

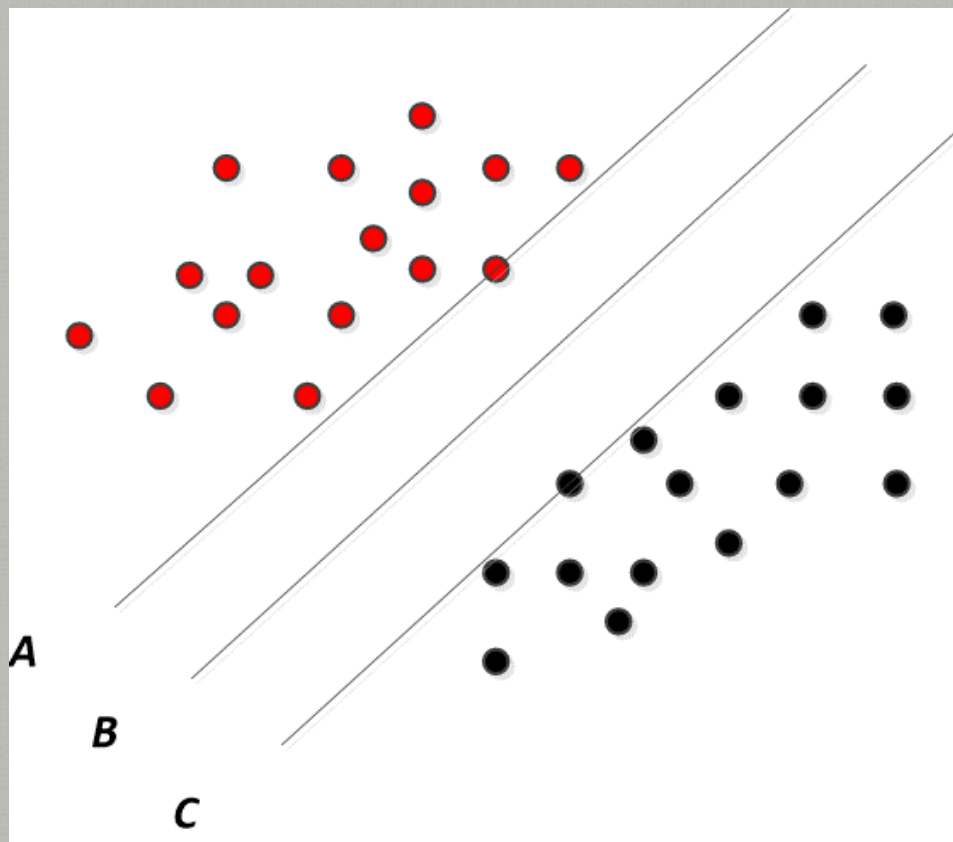
$$\mathbf{w}^T \mathbf{x}_i y_i \geq 0, \quad \text{any } \mathbf{x}_i, i = 1, \dots, N$$

- Any weight vector \mathbf{w} within the **solution region** can classify all the instances correctly. \mathbf{w} is called a **solution vector**.





- Any weight vector w within the **solution region** can classify all the instances correctly. w is called a **solution vector**.
- Solution vectors are not unique!



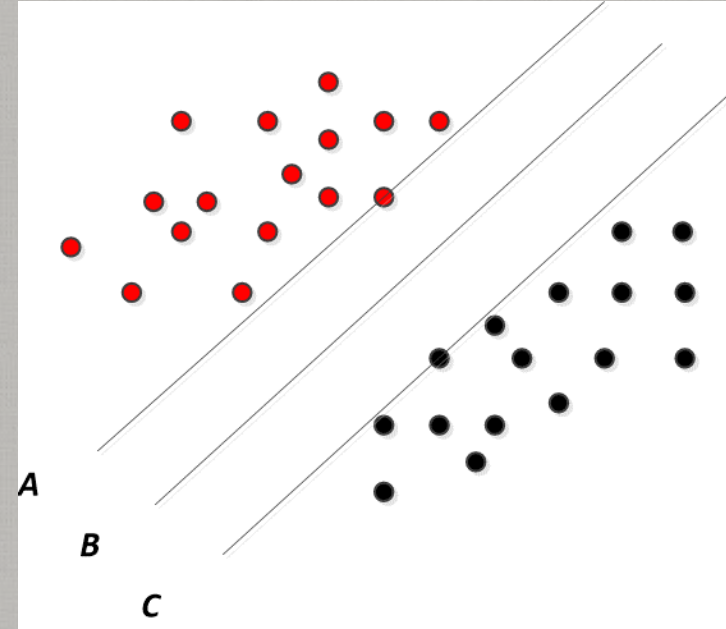
强加

- There are several ways to impose additional requirements to constrain the solution vector.
- One possibility is to seek minimum length weight vector satisfying

$$\mathbf{w}^T \mathbf{x}_i y_i \geq b, \quad \text{for all } i$$

Where $b > 0$ is a constant called **margin**.

- ***The result solution is more likely to classify new test instances correctly***

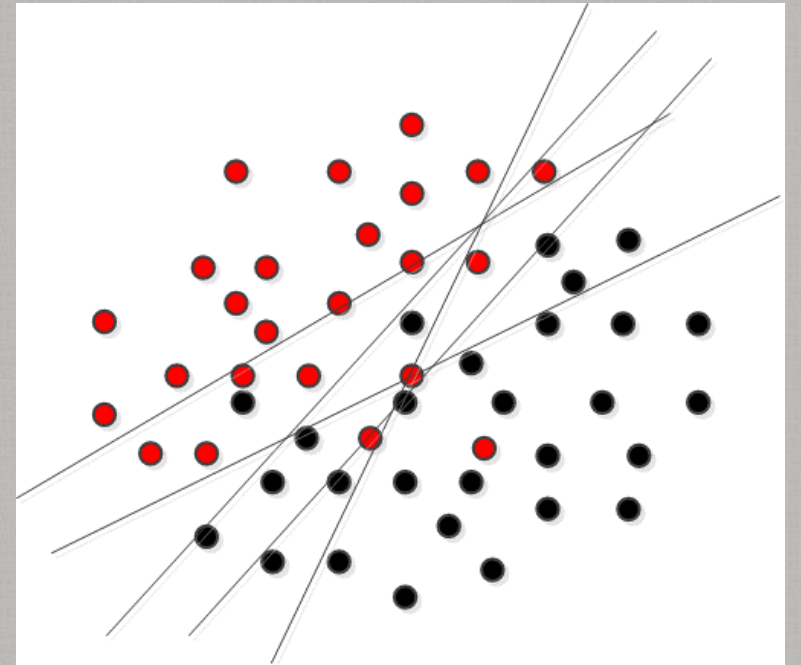


Non-linearly separable Behavior

When apply the Perceptron model to a not linearly separable data set, what is going to happen?

- No guarantee of the convergence
- No guarantee of the performance of the trained model

What should we do?



The Voted and Averaged Perceptron

Training:

Input: a labeled training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$

number of epochs T

Output: a list of weighted Perceptrons $\{(\mathbf{w}_1, c_1), \dots, (\mathbf{w}_k, c_k)\}$

Initialize: $k \leftarrow 0, \mathbf{w}_1 = 0, c_1 = 0$

Repeat T times: { %training stops after T epochs

for $i = 1, \dots, m$

 compute prediction $\hat{y} = \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$

if $y_i \neq \hat{y}$

$c_k = c_k + 1$ (times of “survive”)

else

$\mathbf{w}_{k+1} = \mathbf{w}_k + \mathbf{x}_i y_i$

$c_{k+1} = 1$

$k = k + 1$

}

Return $\{(\mathbf{w}_1, c_1), \dots, (\mathbf{w}_k, c_k)\}$

Test:

Given: a list of weighted Perceptrons $\{(\mathbf{w}_1, c_1), \dots, (\mathbf{w}_k, c_k)\}$
an unseen test instance \mathbf{x}

The **voted perceptron** computes a predicted label \hat{y} as follows:

$$s = \sum_{i=1}^k c_i \text{sign}(\mathbf{w}_i^T \mathbf{x}); \quad \hat{y} = \text{sign}(s)$$

The **averaged perceptron** computes a predicted label \hat{y} as:

$$\hat{y} = \text{sign} \left(\sum_{i=1}^k c_i (\mathbf{w}_i^T \mathbf{x}) \right)$$

How to enhance the Perceptron model to handle non-linearly separable data sets?

Can nonlinear feature mapping be used in Perceptron model?

Logistic Regression

A discriminative classification model

Let's consider the **binary case**: no need to be linearly separable

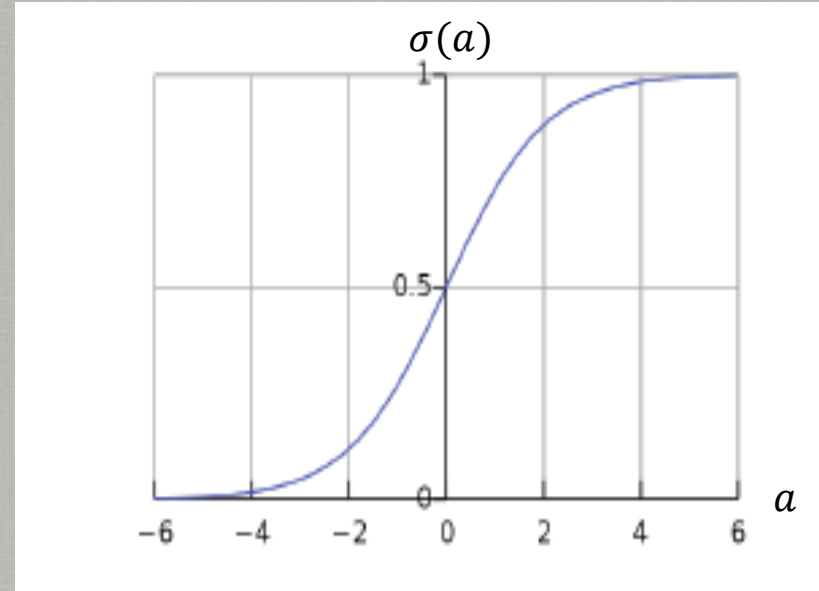
- Given a training data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where $y_i \in \{0,1\}, i = 1, \dots, N$.
- *Logistic regression method directly construct the **posterior probability** that a data point is in class C_1 or C_2 .*

- Logistic regression method directly construct the **posterior probability** that a data point is in class C_1 or C_2 .
- It converts the raw linear “score” $\mathbf{w}^T \mathbf{x}$ to a probability between 0 and 1 by applying a **sigmoid (logistic)** function $\sigma(\mathbf{w}^T \mathbf{x})$, where,

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

- A nice property of the **sigmoid** function

$$\frac{d\sigma(a)}{da} = \sigma(a)(1 - \sigma(a))$$



- To classify an arbitrary point \mathbf{x} , we use the **sigmoid (logistic)** function to output a posterior probability over the class C_1 and C_2 :

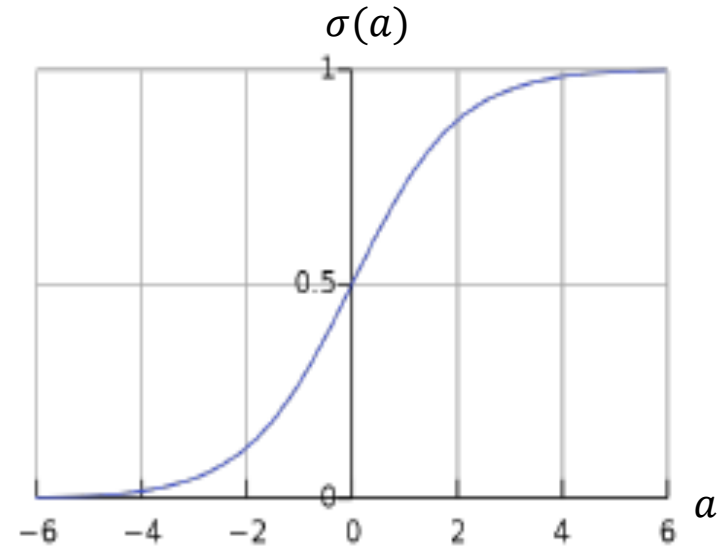
$$P(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}), \quad P(C_2|\mathbf{x}) = 1 - \sigma(\mathbf{w}^T \mathbf{x})$$

- We classify \mathbf{x} as the class with the maximum **posterior probability**:

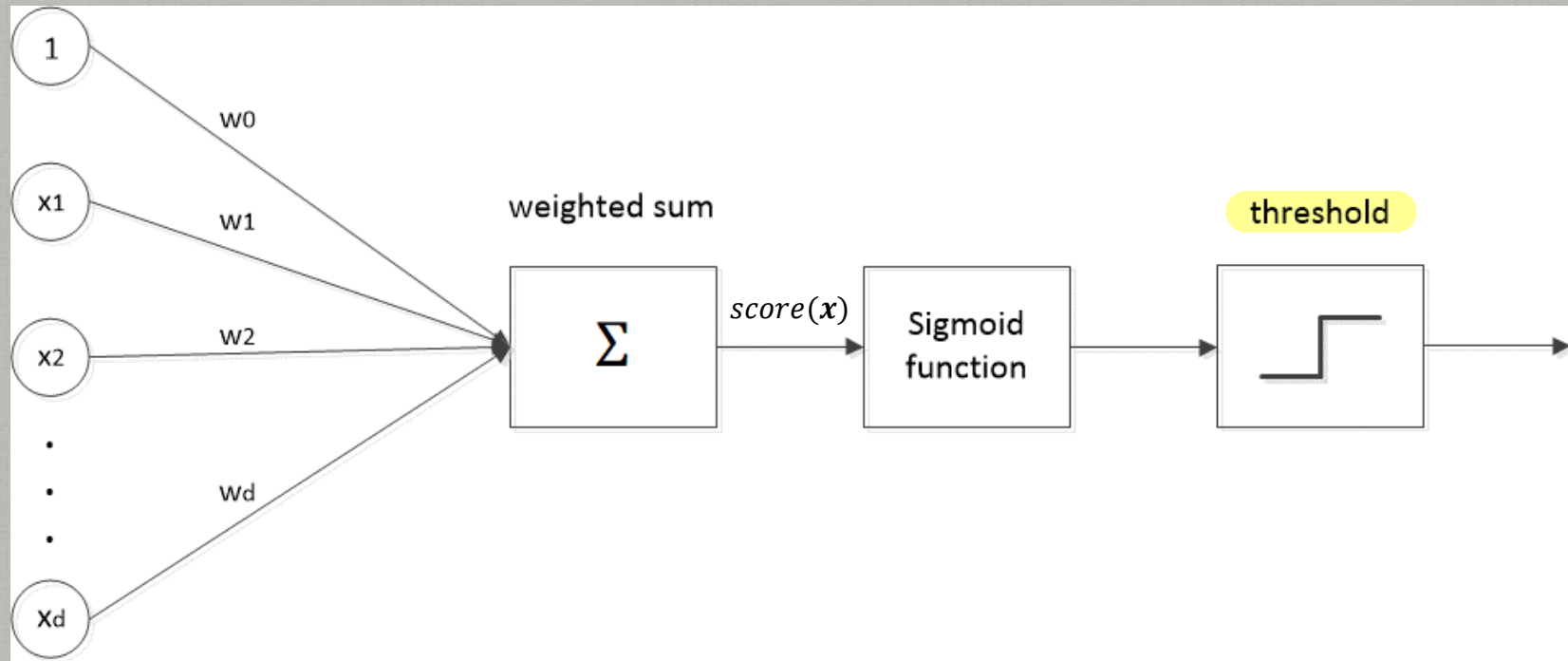
$$\hat{y} = \max_{k=1,2} P(C_k|\mathbf{x}) = \begin{cases} 1, & \text{if } \sigma(\mathbf{w}^T \mathbf{x}) \geq 0.5 \\ 0, & \text{otherwise} \end{cases}$$

- Equivalently, we classify \mathbf{x} as (the decision rule):

$$\hat{y} = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ 0, & \text{otherwise} \end{cases}$$



Graphical model of Binary Logistic Regression



Logistic Regression - A discriminative method

- Perceptron and Logistic regression are both discriminative method
- Perceptron directly work on the discriminant function to find the decision boundary
- Logistic Regression model use posterior probability $P(C_1|\mathbf{x})$ to make classification decision
- Logistic Regression does not calculate the posterior probability using likelihood and prior like in Bayes Theorem. Instead, it constructs the posterior probability by learning from the training data.
- Logistic Regression is a ***discriminative method***

$$P(C_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x}), \quad P(C_2|\mathbf{x}) = 1 - \sigma(\mathbf{w}^T \mathbf{x}), \quad \sigma(a) = \frac{1}{1+e^{-a}}$$

The error (loss) function: (with a maximum likelihood explanation)

Given a training data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, where $y_i \in \{0,1\}, i = 1, \dots, N$. We view each target value y_i as an independent sample from a ***Bernoulli distribution*** $Y_i \sim \text{Bern}(p_i)$, where p_i is a function of \mathbf{x}_i . Thus, we have the following probability:

$$P(Y_i = y_i) = \begin{cases} p_i, & \text{if } y_i = 1 \\ 1 - p_i, & \text{if } y_i = 0 \end{cases}$$

Which can also be written as follows:

$$P(Y_i = y_i) = p_i^{y_i} (1 - p_i)^{1-y_i}$$

Then the ***likelihood function*** can be written as: the probability that the specific sample (the training data set) happens given the hypothesis model (represented by the weight vector \mathbf{w}).

$$p(\mathbf{y}|\mathbf{w}) = P[(Y_1 = y_1) \text{ and } (Y_2 = y_2) \text{ and } \dots \text{ and } (Y_N = y_N)] = \prod_{i=1}^N p_i^{y_i} (1 - p_i)^{1-y_i}$$

Where, $\mathbf{y} = (y_1, \dots, y_N)^T$ and $p_i = P(C_1|\mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i)$.

- ***The error function:***

We can define an error function by taking the **negative** logarithm of the likelihood, which gives the ***cross-entropy*** error function in the form:

$$E(\mathbf{w}) = -\ln p(\mathbf{y}|\mathbf{w}) = -\sum_{i=1}^N \{y_i \ln p_i + (1 - y_i) \ln(1 - p_i)\}$$

- ***This is the maximum likelihood explanation of the error function!***
- ***Minimizing $E(\mathbf{w})$ is equivalent to Maximizing the likelihood function $p(\mathbf{y}|\mathbf{w})$***

Logistic Regression – the training method

The error function:

$$E(\mathbf{w}) = -\ln p(\mathbf{y}|\mathbf{w}) = -\sum_{i=1}^N \{y_i \ln p_i + (1 - y_i) \ln(1 - p_i)\}$$

Where,

$$p_i = P(C_1|\mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$$

The gradient is:

$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w}) &= \nabla_{\mathbf{w}} \left(- \sum_{i=1}^N \{y_i \ln p_i + (1 - y_i) \ln(1 - p_i)\} \right) \\ &= - \sum_{i=1}^N \{y_i \nabla_{\mathbf{w}} \ln p_i + (1 - y_i) \nabla_{\mathbf{w}} \ln(1 - p_i)\} \\ &= - \sum_{i=1}^N \left\{ \frac{y_i}{p_i} \nabla_{\mathbf{w}} p_i - \frac{(1 - y_i)}{1 - p_i} \nabla_{\mathbf{w}} p_i \right\} = - \sum_{i=1}^N \left\{ \frac{y_i}{p_i} - \frac{(1 - y_i)}{1 - p_i} \right\} \nabla_{\mathbf{w}} p_i\end{aligned}$$

Let's find out $\nabla_{\mathbf{w}} p_i$, with $p_i = P(C_1 | \mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}$.

Note that $\nabla_a \sigma(a) = \sigma(a)(1 - \sigma(a))$, and from the chain rule of derivative, we have,

$$\begin{aligned}\nabla_{\mathbf{w}} p_i &= \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \mathbf{x}_i) \\ &= \sigma(\mathbf{w}^T \mathbf{x}_i)(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) \mathbf{x}_i \\ &= p_i(1 - p_i) \mathbf{x}_i\end{aligned}$$

Hence, we have,

$$\begin{aligned}\nabla_{\mathbf{w}} E(\mathbf{w}) &= - \sum_{i=1}^N \left\{ \frac{y_i}{p_i} - \frac{(1 - y_i)}{1 - p_i} \right\} p_i (1 - p_i) \mathbf{x}_i \\ &= - \sum_{i=1}^N \{ y_i (1 - p_i) - (1 - y_i) p_i \} \mathbf{x}_i \\ &= - \sum_{i=1}^N (y_i - p_i) \mathbf{x}_i\end{aligned}$$

The gradient descent update is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla_{\mathbf{w}} E(\mathbf{w}_k) = \mathbf{w}_k + \eta \sum_{i=1}^N (y_i - p_i) \mathbf{x}_i$$

Where, $p_i = \sigma(\mathbf{w}_k^T \mathbf{x}_i) = \frac{1}{1 + e^{-\mathbf{w}_k^T \mathbf{x}_i}}$

sigmoid function

Q1: What are the differences between Perceptron training and Logistic regression training algorithms?

differences between batch training and sequential training?

Q2: Do we have batch training and stochastic training for Logistic regression?

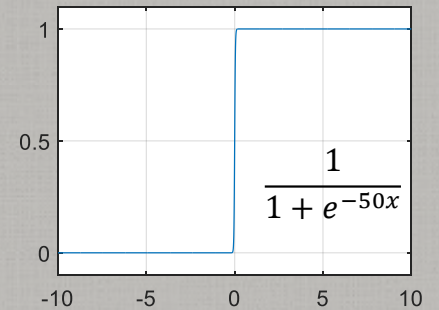
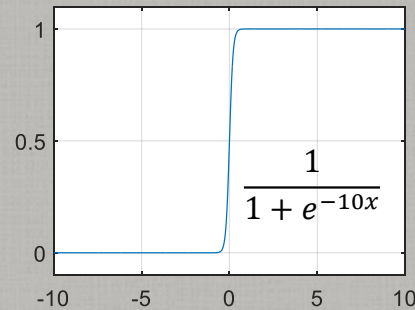
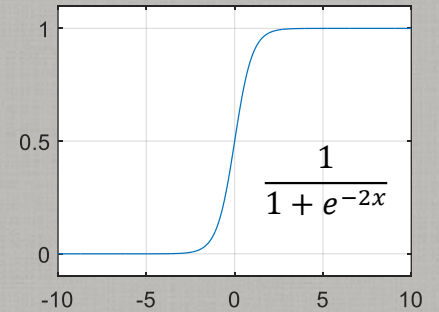
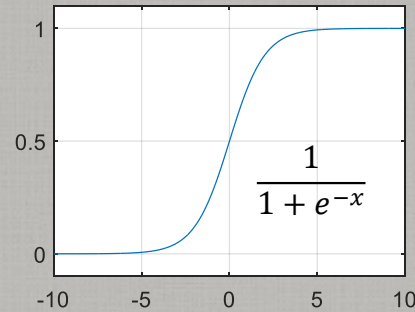
Q3: Can we use feature mapping $\phi(\mathbf{x})$ in Logistic Regression?

Regularized Logistic Regression

- ***When data set is linearly separable, logistic regression weights go to infinite.***

when data set is linearly separable, the conditional likelihood function is perfectly 0 and 1, the error function $E(\mathbf{w})$ goes to zero.

This is corresponding to the steepest sigmoid function $\frac{1}{1+e^{\mathbf{w}^T \mathbf{x}}}$ where the weights go to infinite.



Regularized Logistic Regression

- *Also, using feature mapping $\phi(x)$ may cause overfitting*
- Large number of features in the model trained using a small data set may cause overfitting
- To avoid overfitting, a regularization penalty term can be included into the error function:

$$E(\mathbf{w}) = -\ln p(\mathbf{y}|\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 = -\sum_{i=1}^N \{y_i \ln p_i + (1 - y_i) \ln(1 - p_i)\} + \lambda \|\mathbf{w}\|_2^2$$

Regularized Logistic Regression

This new error function is equivalent to:

$$E(\mathbf{w}) = 1/\lambda \left\{ - \sum_{i=1}^N \{y_i \ln p_i + (1 - y_i) \ln(1 - p_i)\} \right\} + \|\mathbf{w}\|_2^2$$

$$E(\mathbf{w}) = C \left\{ - \sum_{i=1}^N \{y_i \ln p_i + (1 - y_i) \ln(1 - p_i)\} \right\} + \|\mathbf{w}\|_2^2$$

Where, $C = 1/\lambda$ is the inverse of the regularization strength (*sk-learn use regularized LR by default*).

Multinomial Logistic Regression

Multiclass Logistic Regression:

- Let's generalize logistic regression to the case where there are K classes
 - The training data set $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$
 - The one-hot vector encoding:
- If the i^{th} observation belongs to class C_k , instead of $y_i = k$, we use the representation $\mathbf{y}_i = \mathbf{e}_k$, the k^{th} canonical basis vector. For example, if an observation belong to the class C_2 , its label representation would be: ($K = 3, k = 2$)

$$\mathbf{y}_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Now, the target value becomes a K dimensional vector.

Multiclass Logistic Regression:

- Instead of a weight vector, now we have a weight matrix to be determined:

$$\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{w}_K^T \end{bmatrix}$$

Where $\mathbf{w}_k \in R^{d+1}$ is the weight vector for the k^{th} class.

- For each input $\mathbf{x}_i \in R^{d+1}$, each class k is given a “score” $z_k = \mathbf{w}_k^T \mathbf{x}_i$
- In total there are K raw linear scores for an arbitrary input \mathbf{x} :

$$\mathbf{W}\mathbf{x} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \vdots \\ \mathbf{w}_K^T \mathbf{x} \end{bmatrix} = \begin{bmatrix} z_1 \\ \vdots \\ z_K \end{bmatrix}$$

Multinomial Logistic Regression

Multiclass Logistic Regression:

$$Wx = \begin{bmatrix} \mathbf{w}_1^T x \\ \vdots \\ \mathbf{w}_K^T x \end{bmatrix} = \begin{bmatrix} z_1 \\ \vdots \\ z_K \end{bmatrix}$$

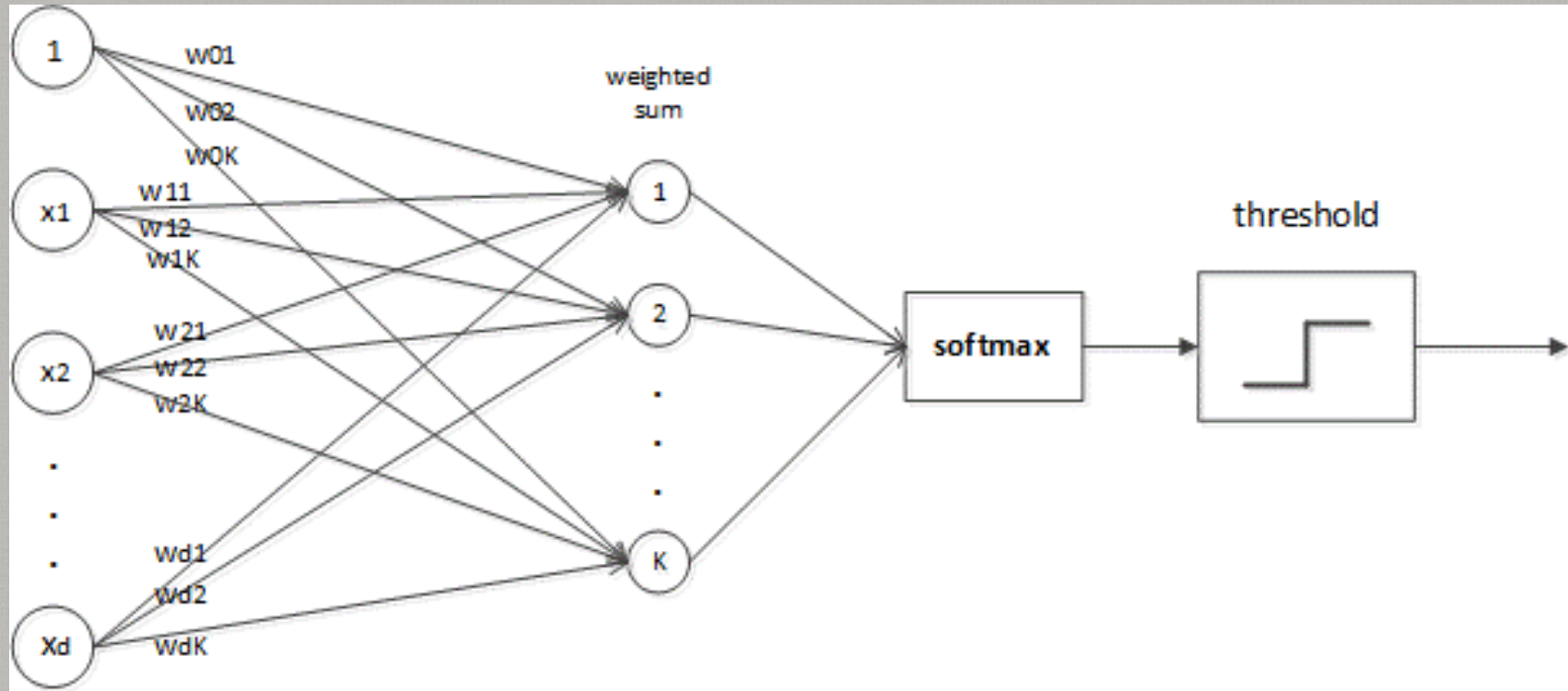
- We now must transform all of these scores into a posterior probability distribution $P(C_k|x)$. For the multiclass case, We use ***softmax*** function:

$$P(C_k|x) = \sigma(z_k) = \frac{e^{z_k}}{\sum_{j=1}^K e^{z_j}} = \frac{e^{\mathbf{w}_k^T x}}{\sum_{j=1}^K e^{\mathbf{w}_j^T x}}, \quad k = 1, 2, \dots, K$$

- We then predict the class with the maximum probability:

$$\hat{y} = \max_k P(C_k|x)$$

Graphical model of multinomial Logistic Regression



The error function for multi-class logistic regression:

- The *likelihood function* is given by:

$$p(\mathbf{y}|\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = \prod_{i=1}^N \prod_{k=1}^K [P(C_k|\mathbf{x}_i)]^{y_{ik}} = \prod_{i=1}^N \prod_{k=1}^K [p_{ik}]^{y_{ik}}$$

Where \mathbf{y} is an $K \times N$ matrix of target values with elements y_{ik} and

$$p_{ik} = P(C_k|\mathbf{x}_i) = \frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}}$$

- Taking the negative logarithm then gives:

$$E(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = -\ln p(\mathbf{y}|\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K) = -\sum_{i=1}^N \sum_{k=1}^K (y_{ik} \ln p_{ik}) = -\sum_{i=1}^N \sum_{k=1}^K \left(y_{ik} \ln \left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}} \right) \right)$$

Which is know as the ***cross-entropy*** error function for the multiclass logistic regression problem.

- Instead of finding the gradient with respect to all of the parameters of the matrix \mathbf{W} , let's find the gradient with respect to one row of \mathbf{W} at a time:

$$\begin{aligned}
 \nabla_{\mathbf{w}_l} E(\mathbf{w}_1, \dots, \mathbf{w}_K) &= \nabla_{\mathbf{w}_l} \left\{ - \sum_{i=1}^N \sum_{k=1}^K \left(y_{ik} \ln \left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}} \right) \right) \right\} \\
 &= - \sum_{i=1}^N \sum_{k=1}^K \left(y_{ik} \nabla_{\mathbf{w}_l} \left\{ \ln \left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}} \right) \right\} \right) \\
 &= - \sum_{i=1}^N \sum_{k=1}^K \left(y_{ik} \nabla_{\mathbf{w}_l} \left\{ \ln \left(\frac{e^{\mathbf{w}_k^T \mathbf{x}_i}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}_i}} \right) \right\} \right) \\
 &= - \sum_{i=1}^N (y_{il} - p_{il}) \mathbf{x}_i
 \end{aligned}$$

Multinomial Logistic Regression

- The gradient descent update for \mathbf{w}_l is then, ($l = 1, \dots, K$)

$$\mathbf{w}_l^{t+1} = \mathbf{w}_l^t + \eta \sum_{i=1}^N (y_{il} - p_{il}) \mathbf{x}_i$$