



# 第3讲：线性模型

---

师 佳

化学化工学院

化学工程与生物工程系

# 纲要

## ➤ 分类问题示例

## ➤ 线性分类模型

### ➤ 概述

### ➤ Logistic 回归

### ➤ Softmax 回归

### ➤ 感知器

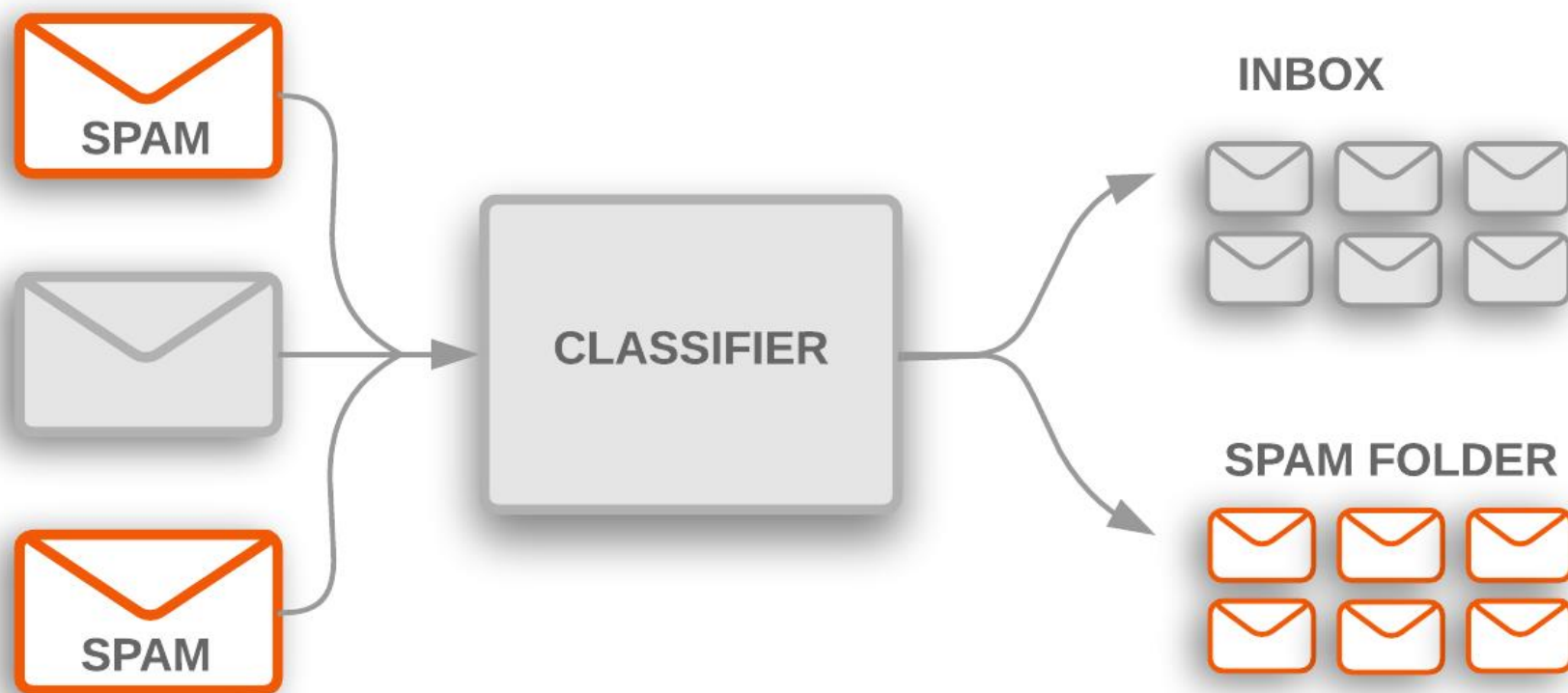
### ➤ 支持向量机 (SVM)

## ➤ 总结

# 分类问题示例

## ➤ 两分类问题

### ➤ 垃圾邮件过滤



# 分类问题示例

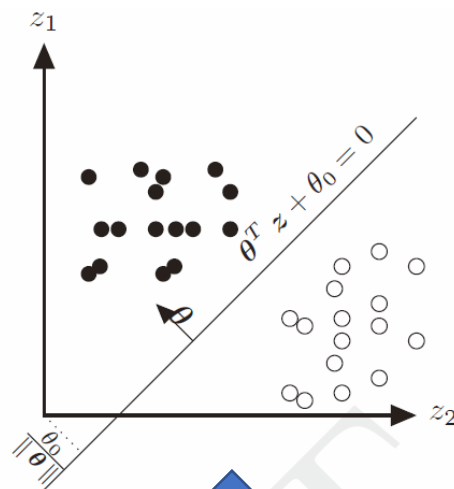
## ➤ 两分类问题

### ➤ 文本语意分类

$D_1$ : “我喜欢读书”

$D_2$ : “我讨厌读书”

	我	喜欢	讨厌	读书
$D_1$	1	1	0	1
$D_2$	1	0	1	1



$$\mathbf{x}_1 = [1 \ 1 \ 0 \ 1]^T,$$

$$\mathbf{x}_2 = [1 \ 0 \ 1 \ 1]^T.$$

爱读书

不爱读书

# 分类问题示例

## ➤ 多分类问题

- 数据集: CIFAR-10
- 60000张32x32色彩图像, 共10类
- 每类6000张图像

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck



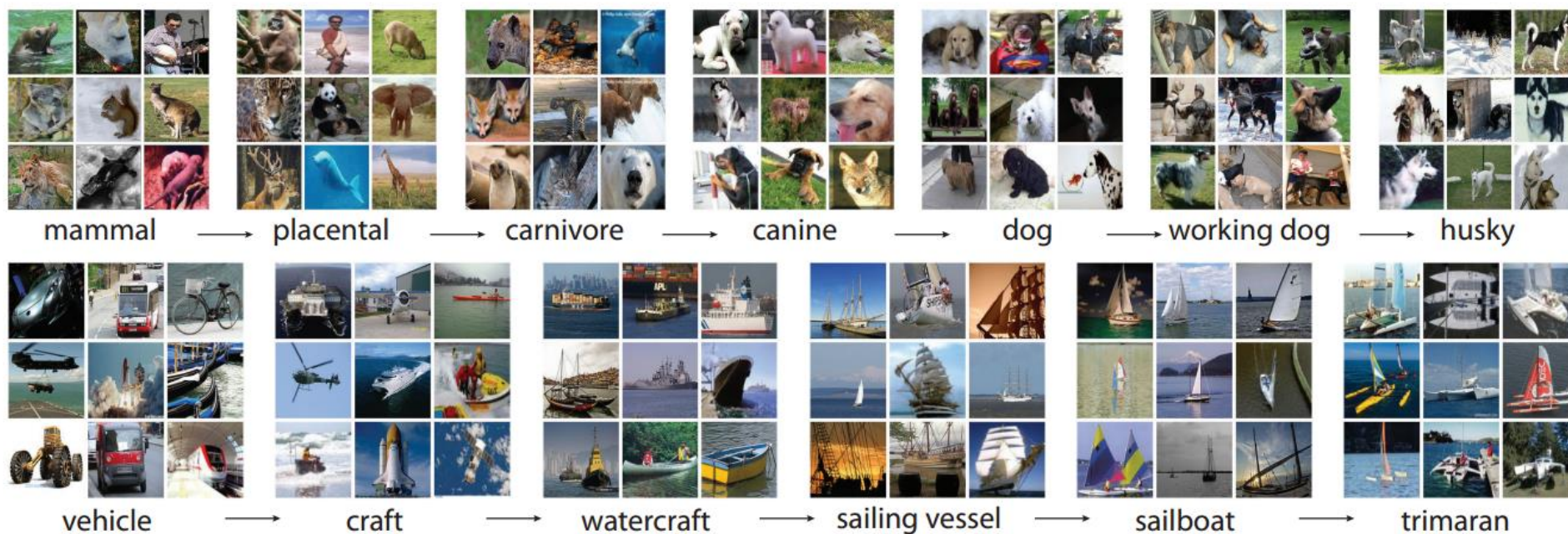


# 分类问题示例

## ➤ 多分类问题

➤ 数据集: ImageNet

➤ 14,197,122 images, 21841 synsets



# 分类问题示例

## ➤ 多分类问题

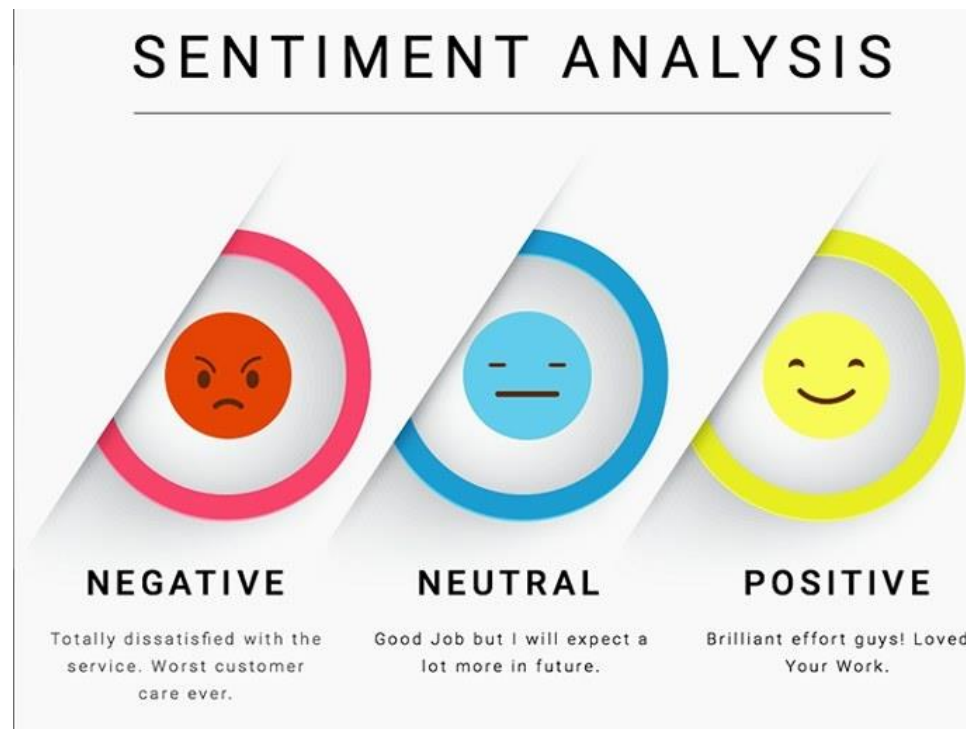
### ➤ 文档归类



# 分类问题示例

## ➤ 多分类问题

### ➤ 文本情感分类



Review (X)

Rating (Y)

"This movie is fantastic! I really like it because it is so good!"



"Not to my taste, will skip and watch another movie"



"This movie really sucks! Can I get my money back please?"





# 分类问题示例

## ➤ 分类问题的扩展

### ➤ 图像分类、目标检测、实例分割

**Classification**



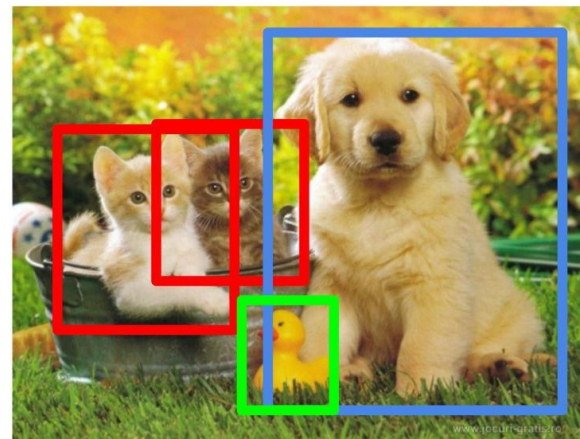
CAT

**Classification  
+ Localization**



CAT

**Object Detection**



CAT, DOG, DUCK

**Instance  
Segmentation**



CAT, DOG, DUCK

Single object

Multiple objects

# 线性分类模型

## ➤概述

➤ **线性分类模型**：给定一个  $D$  维样本  $\mathbf{x} = [x^1, \dots, x^D]^\top$ ，其线性组合函数为

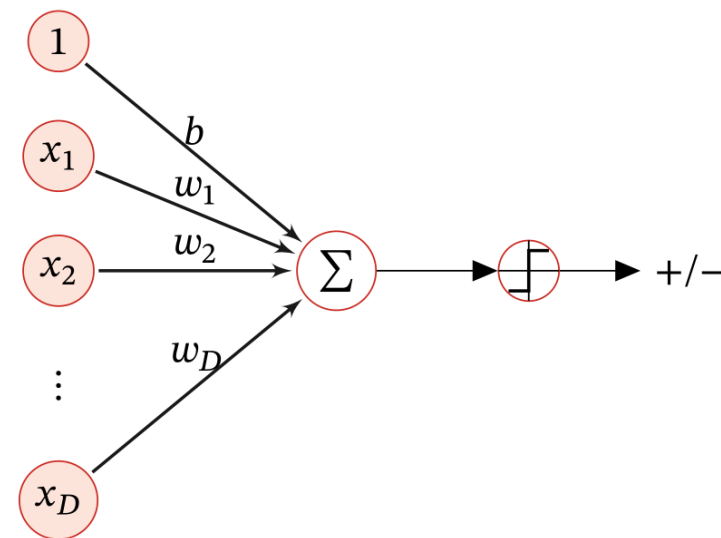
$$\begin{aligned} f(\mathbf{x}; \mathbf{w}) &= w_1 x_1 + w_2 x_2 + \dots + w_D x_D + b \\ &= \mathbf{w}^\top \mathbf{x} + b, \end{aligned}$$

引入一个**非线性的决策函数** (Decision Function)  $g(\cdot)$ 来预测输出目标

$$y = g(f(\mathbf{x}; \mathbf{w}))$$

如符号函数

$$\begin{aligned} g(f(\mathbf{x}; \mathbf{w})) &= \text{sgn}(f(\mathbf{x}; \mathbf{w})) \\ &\triangleq \begin{cases} +1 & \text{if } f(\mathbf{x}; \mathbf{w}) > 0, \\ -1 & \text{if } f(\mathbf{x}; \mathbf{w}) < 0. \end{cases} \end{aligned}$$



# 线性分类模型

## ➤概述

➤ 二分类问题所有满足 $f(x;w)=0$ 的点组成一个**分割超平面**将特征空间一分为二，划分成两个区域，每个区域对应一个类别。

➤ 每个样本点到决策平面的**有向距离** (Signed Distance) 为

$$\gamma = \frac{f(x;w)}{\|w\|}.$$

➤ 给定 $N$  个样本的训练集 $\mathcal{D}=\{(x^{(n)},y^{(n)})\}_{n=1,\dots,N}$ ，其中 $y^{(n)}\in\{+1,-1\}$ ，线性模型试图学习到参数 $w^*$ ，使得对于每个样本 $(x^{(n)},y^{(n)})$ 尽量满足

$$y^{(n)} f(x^{(n)}; w^*) > 0, \quad \forall n \in [1, N]$$

**定义 3.1 – 两类线性可分：** 对于训练集  $\mathcal{D} = \{(x^{(n)}, y^{(n)})\}_{n=1}^N$ ，如果存在权重向量  $w^*$ ，对所有样本都满足  $yf(x; w^*) > 0$ ，那么训练集  $\mathcal{D}$  是线性可分的。

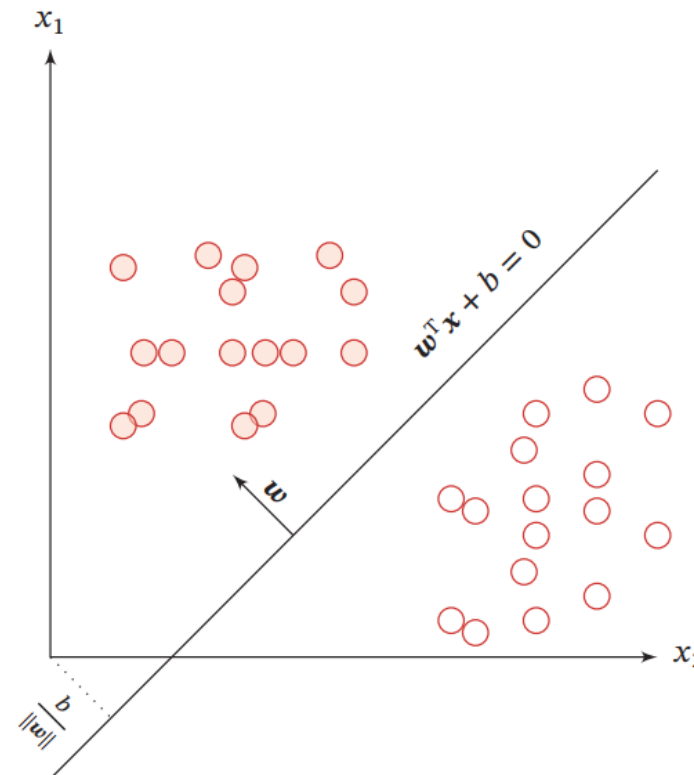
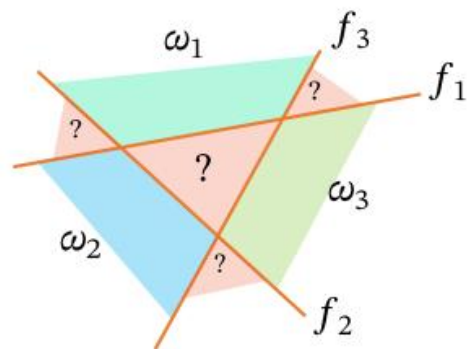


图 3.2 二分类的决策边界示例

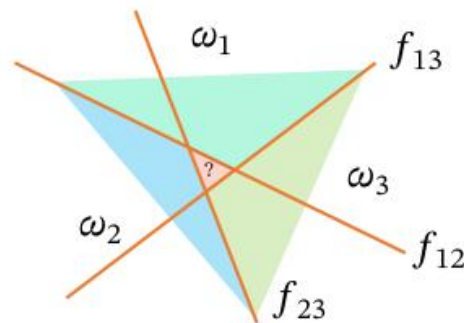
# 线性分类模型

## ➤ 概述

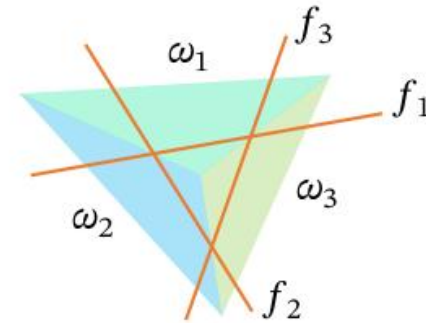
### ➤ 多分类问题



(a) “一对其余”方式



(b) “一对一”方式



(c) “argmax”方式

- “argmax” 方式：共需要 $C$ 个判别函数  $f_c(\mathbf{x}; \mathbf{w}_c) = \mathbf{w}_c^\top \mathbf{x} + b_c$ ,  $c \in \{1, \dots, C\}$   
“argmax” 方式的预测函数定义为：

$$y = \arg \max_{c=1}^C f_c(\mathbf{x}; \mathbf{w}_c)$$

**定义 3.2** – 多类线性可分：对于训练集  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ , 如果存在  $C$  个权重向量  $\mathbf{w}_1^*, \dots, \mathbf{w}_C^*$ , 使得第  $c$  ( $1 \leq c \leq C$ ) 类的所有样本都满足  $f_c(\mathbf{x}; \mathbf{w}_c^*) > f_{\tilde{c}}(\mathbf{x}, \mathbf{w}_{\tilde{c}}^*), \forall \tilde{c} \neq c$ , 那么训练集  $\mathcal{D}$  是线性可分的。



# 线性分类模型

## ➤ Logistic回归

➤ 用途: **两分类**问题

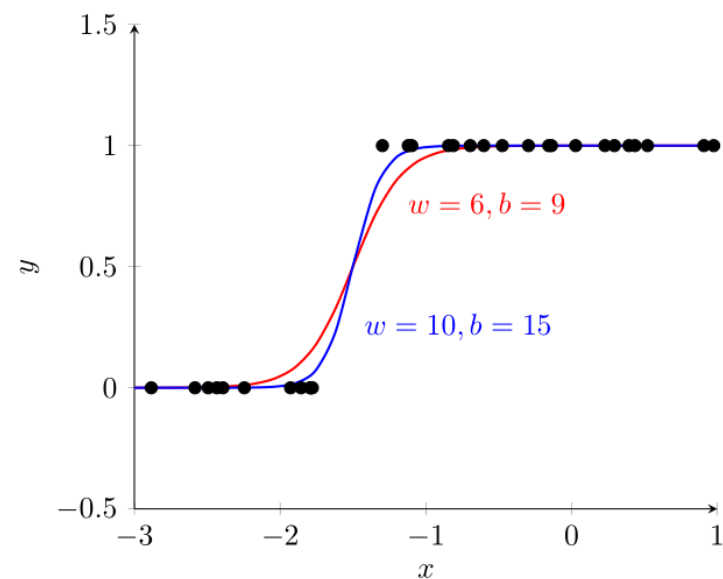
➤ 模型: 引入非线性函数  $g: \mathbb{R}^D \rightarrow (0, 1)$  来预测类别标签的**后验概率**  $p(y=1|x)$

$$p(y=1|x) = g(f(x; w)) = \sigma(w^\top x) \\ \triangleq \frac{1}{1 + \exp(-w^\top x)}$$

➤ 学习准则: **交叉熵损失函数**

➤ 给定  $N$  个训练样本  $\{(x^{(n)}, y^{(n)})\}_{n=1, \dots, N}$ , 用Logistic回归模型对每个样本  $x^{(n)}$  进行预测, 输出其标签为1的后验概率, 记为  $\hat{y}^{(n)}$ , **经验风险函数**定义为:

$$\begin{aligned} \mathcal{R}(w) &= -\frac{1}{N} \sum_{n=1}^N \left( p_r(y^{(n)} = 1|x^{(n)}) \log \hat{y}^{(n)} + p_r(y^{(n)} = 0|x^{(n)}) \log(1 - \hat{y}^{(n)}) \right) \\ &= -\frac{1}{N} \sum_{n=1}^N \left( y^{(n)} \log \hat{y}^{(n)} + (1 - y^{(n)}) \log(1 - \hat{y}^{(n)}) \right) \end{aligned}$$



# 线性分类模型

## ➤ Logistic回归

### ➤ 优化算法：梯度下降法（牛顿法）

➤ 经验风险函数 $\mathcal{R}(\mathbf{w})$ 关于参数 $\mathbf{w}$ 的偏导数（梯度）为：

$$\begin{aligned}\frac{\partial \mathcal{R}(\mathbf{w})}{\partial \mathbf{w}} &= -\frac{1}{N} \sum_{n=1}^N \left( y^{(n)} \frac{\hat{y}^{(n)}(1 - \hat{y}^{(n)})}{\hat{y}^{(n)}} \mathbf{x}^{(n)} - (1 - y^{(n)}) \frac{\hat{y}^{(n)}(1 - \hat{y}^{(n)})}{1 - \hat{y}^{(n)}} \mathbf{x}^{(n)} \right) \\ &= -\frac{1}{N} \sum_{n=1}^N \left( y^{(n)}(1 - \hat{y}^{(n)}) \mathbf{x}^{(n)} - (1 - y^{(n)}) \hat{y}^{(n)} \mathbf{x}^{(n)} \right) \\ &= -\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (y^{(n)} - \hat{y}^{(n)}).\end{aligned}$$

➤ Logistic 回归的训练过程为：初始化  $\mathbf{w}_0 \leftarrow \mathbf{0}$ ，然后通过下式来迭代更新参数

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (y^{(n)} - \hat{y}_{\mathbf{w}_t}^{(n)})$$

$$\begin{aligned}\frac{\partial}{\partial \mathbf{w}} \log \hat{y}^{(n)} &= \frac{\partial}{\partial \mathbf{w}} \log \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{(n)})} \\ &= \frac{1}{\hat{y}^{(n)}} \frac{\partial}{\partial \mathbf{w}} \left( \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x}^{(n)})} \right) \\ &= \frac{1}{\hat{y}^{(n)}} \frac{\exp(-\mathbf{w}^T \mathbf{x}^{(n)})}{\left(1 + \exp(-\mathbf{w}^T \mathbf{x}^{(n)})\right)^2} \mathbf{x}^{(n)} \\ &= \frac{\hat{y}^{(n)}(1 - \hat{y}^{(n)})}{\hat{y}^{(n)}} \mathbf{x}^{(n)}\end{aligned}$$



# 线性分类模型

## ➤ Softmax回归 (多分类Logistic回归)

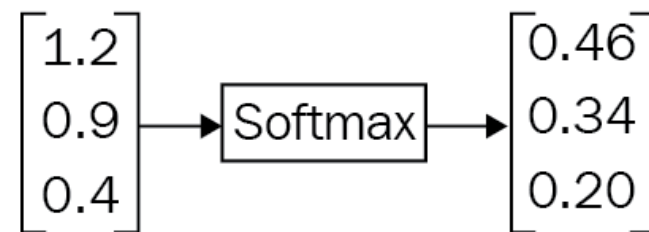
➤ 用途: **多分类问题**

➤ 模型: 对于多类问题, 类别标签  $y \in \{1, 2, \dots, C\}$  可以有  $C$  个取值. 给定一个样本  $x$ , Softmax 回归预测的属于类别  $c$  的 **条件概率** 为

$$p(y = c | x) = \text{softmax}(\mathbf{w}_c^\top \mathbf{x}) = \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x})}$$

向量表示

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{W}^\top \mathbf{x}) = \frac{\exp(\mathbf{W}^\top \mathbf{x})}{\mathbf{1}_C^\top \exp(\mathbf{W}^\top \mathbf{x})}$$



其中  $\mathbf{W} = [\mathbf{w}^1, \dots, \mathbf{w}^C]$  是由  $C$  个类的权重向量组成的矩阵,  $\mathbf{1}^C$  为  $C$  维的全1向量,  $\mathbf{y} \in \mathbb{R}^C$  为所有类别的预测条件概率组成的向量, 取预测函数为:

$$\hat{y} = \arg \max_{c=1}^C p(y = c | x) = \arg \max_{c=1}^C \mathbf{w}_c^\top \mathbf{x}.$$

# 线性分类模型

## ➤ Softmax回归（多分类Logistic回归）

### ➤ 学习准则：交叉熵损失函数

- 用 $C$ 维的one-hot向量 $\mathbf{y} \in \{0, 1\}^C$ 来表示类别标签

$$\mathbf{y} = [I(1 = c), I(2 = c), \dots, I(C = c)]^\top$$

- Softmax回归模型的经验风险损失函数为

$$\mathcal{R}(\mathbf{W}) = -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \mathbf{y}_c^{(n)} \log \hat{\mathbf{y}}_c^{(n)} = -\frac{1}{N} \sum_{n=1}^N (\mathbf{y}^{(n)})^\top \log \hat{\mathbf{y}}^{(n)}$$

### ➤ 优化算法：梯度下降法

- 经验风险函数 $\mathcal{R}(\mathbf{w})$ 关于参数 $\mathbf{w}$ 的偏导数（梯度）为：

$$\frac{\partial \mathcal{R}(\mathbf{W})}{\partial \mathbf{W}} = -\frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (\mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)})^\top$$

- Softmax回归的训练过程为：初始化 $\mathbf{W}_0 \leftarrow \mathbf{0}$ ，然后通过下式进行迭代更新

$$\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t + \alpha \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)} (\mathbf{y}^{(n)} - \hat{\mathbf{y}}_{\mathbf{W}_t}^{(n)})^\top \right)$$

# 线性分类模型

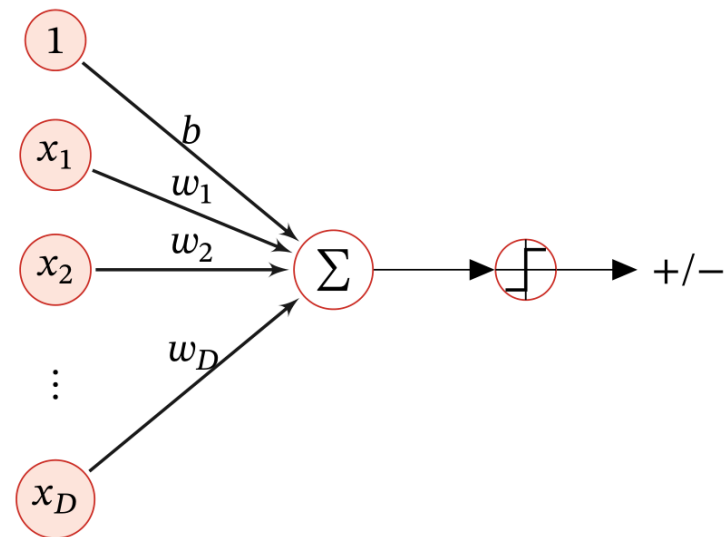
## ➤ 感知器

➤ 用途：二分类问题

➤ 模型：符号函数

$$g(f(\mathbf{x}; \mathbf{w})) = \text{sgn}(f(\mathbf{x}; \mathbf{w}))$$

$$\triangleq \begin{cases} +1 & \text{if } f(\mathbf{x}; \mathbf{w}) > 0, \\ -1 & \text{if } f(\mathbf{x}; \mathbf{w}) < 0. \end{cases}$$



➤ 给定 $N$ 个样本的训练集： $\{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1, \dots, N}$ ，其中 $y^{(n)} \in \{+1, -1\}$ ，感知器学习算法试图找到一组参数 $\mathbf{w}^*$ ，使得对于每个样本 $(\mathbf{x}^{(n)}, y^{(n)})$ 有

$$y^{(n)} \mathbf{w}^{*\top} \mathbf{x}^{(n)} > 0, \quad \forall n \in \{1, \dots, N\}$$

➤ 学习准则：感知器的损失函数为

$$\mathcal{L}(\mathbf{w}; \mathbf{x}, y) = \max(0, -y \mathbf{w}^\top \mathbf{x})$$

# 线性分类模型

## 感知器

### 优化算法：梯度下降法

每一个样本的梯度  $\frac{\partial \mathcal{L}(w; x, y)}{\partial w} = \begin{cases} 0 & \text{if } yw^\top x > 0, \\ -yx & \text{if } yw^\top x < 0. \end{cases}$

算法 3.1 两类感知器的参数学习算法

```
输入: 训练集  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ , 最大迭代次数  $T$ 

1 初始化:  $w_0 \leftarrow 0, k \leftarrow 0, t \leftarrow 0$ ;
2 repeat
3   对训练集  $\mathcal{D}$  中的样本随机排序;
4   for  $n = 1 \cdots N$  do
5     选取一个样本  $(\mathbf{x}^{(n)}, y^{(n)})$ ;
6     if  $w_k^\top (y^{(n)} \mathbf{x}^{(n)}) \leq 0$  then
7        $w_{k+1} \leftarrow w_k + y^{(n)} \mathbf{x}^{(n)}$ ;
8        $k \leftarrow k + 1$ ;
9     end
10     $t \leftarrow t + 1$ ;
11    if  $t = T$  then break;           // 达到最大迭代次数
12  end
13 until  $t = T$ ;
输出:  $w_k$ 
```

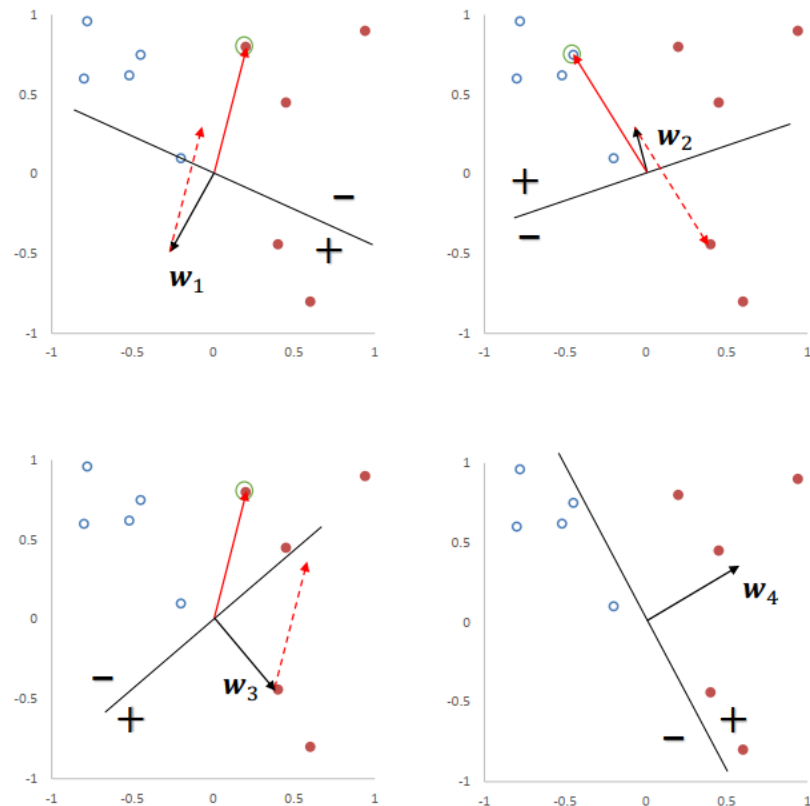


图 3.5 感知器参数学习的更新过程

# 线性分类模型

## ➤ 感知器

### ➤ 感知器收敛性

- 当数据集是两类线性可分时，对于训练集  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1, \dots, N}$ ，其中  $\mathbf{x}^{(n)}$  为样本的增广特征向量， $y^{(n)} \in \{-1, 1\}$ ，那么存在一个正的常数  $\gamma (\gamma > 0)$  和权重向量  $\mathbf{w}^*$ ，并且  $\|\mathbf{w}^*\| = 1$ ，对所有  $n$  都满足  $(\mathbf{w}^*)^\top (y^{(n)} \mathbf{x}^{(n)}) \geq \gamma$ 。可以证明如下定理。

**定理 3.1 – 感知器收敛性：** 给定训练集  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ ，令  $R$  是训练集中最大的特征向量的模，即

$$R = \max_n \|\mathbf{x}^{(n)}\|.$$

如果训练集  $\mathcal{D}$  线性可分，两类感知器的参数学习算法 3.1 的权重更新次数不超过  $\frac{R^2}{\gamma^2}$ 。

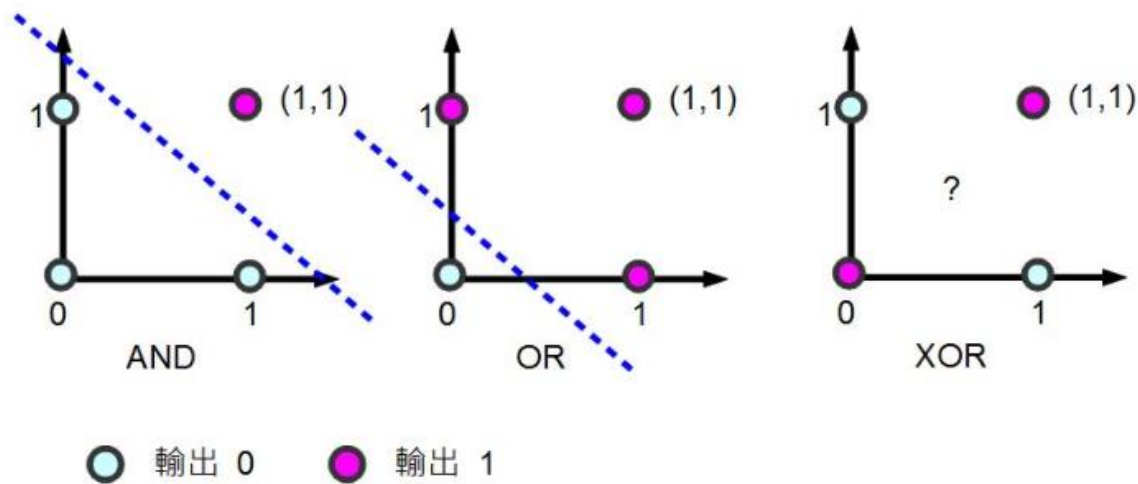
# 线性分类模型

## ➤ 感知器

### ➤ 感知器的缺陷：

- 在数据集线性可分时，感知器虽然可以找到一个超平面把两类数据分开，但**并不能保证其泛化能力**。
- 感知器**对样本顺序比较敏感**。每次迭代的顺序不一致时，找到的分割超平面也往往不一致。
- 如果训练集不是线性可分的，就**永远不会收敛**。

### ➤ 感知器无法解决XOR问题





# 线性分类模型

## ➤ 支持向量机 (Support Vector Machine: SVM)

➤ 用途: **二分类问题**

➤ 模型: 符号函数  $g(f(x; w)) = \text{sgn}(f(x; w))$

$$\triangleq \begin{cases} +1 & \text{if } f(x; w) > 0, \\ -1 & \text{if } f(x; w) < 0. \end{cases}$$

定义**间隔** (Margin)  $\gamma$ 为整个数据集 $D$ 中所有样本到分割超平面的最短距离

$$\gamma = \min_n \gamma^{(n)} = \min_n \frac{|w^\top x^{(n)} + b|}{\|w\|} = \min_n \frac{y^{(n)}(w^\top x^{(n)} + b)}{\|w\|}$$

➤ 学习准则: 支持向量机的目标是**寻找一个超平面**( $w^*$ ,  $b^*$ )**使得 $\gamma$ 最大**, 即

$$\begin{aligned} & \max_{w, b} \quad \gamma \\ & \text{s.t.} \quad \frac{y^{(n)}(w^\top x^{(n)} + b)}{\|w\|} \geq \gamma, \forall n \in \{1, \dots, N\} \end{aligned}$$

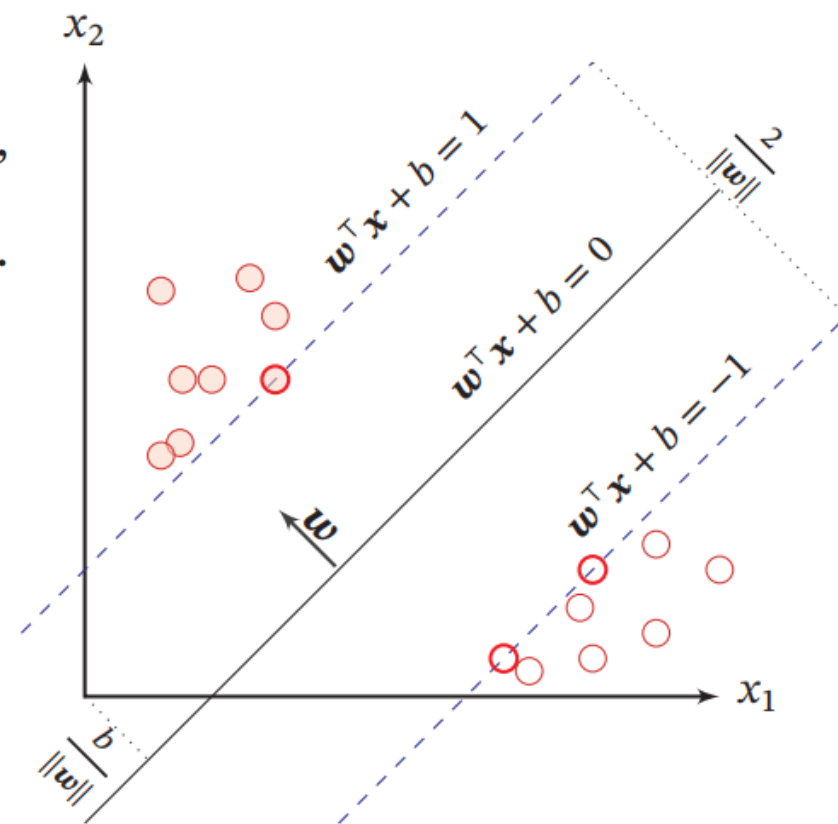


图 3.6 支持向量机示例

# 线性分类模型

## ➤ 补充：约束优化问题

- **约束优化(Constrained Optimization)问题**中变量 $x$ 需要满足一些等式或不等式的约束.
- **线性规划(Linear Programming)问题**:目标函数和所有的约束函数都为线性函数.
- **非线性规划(Nonlinear Programming)问题**:目标函数或任何一个约束函数为非线性函数.
- 约束优化问题可以表示为:

$$\min_x f(x)$$

$$\text{s.t.} \quad \begin{aligned} h_m(x) &= 0, \quad m = 1, \dots, M \\ g_n(x) &\leq 0, \quad n = 1, \dots, N \end{aligned}$$

可行域

$$\mathcal{D} = \text{dom}(f) \cap \bigcap_{m=1}^M \text{dom}(h_m) \cap \bigcap_{n=1}^N \text{dom}(g_n) \subseteq \mathbb{R}^D$$

- 在非线性优化问题中，有一类比较特殊的问题是凸优化(Convex Optimization)问题
  - 变量  $x$  的可行域为凸集(Convex Set)，即对于集合中任意两点，它们的连线全部位于集合内部
  - 目标函数  $f$  也必须为凸函数，即满足

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y), \quad \forall \alpha \in [0, 1]$$

# 线性分类模型

## ➤ 补充：约束优化问题

➤ 约束优化问题通常使用**拉格朗日乘数法**来进行求解问题。

➤ 等式约束优化问题

$$\begin{array}{ll} \min_x & f(x) \\ \text{s.t.} & h_m(x) = 0, \quad m = 1, \dots, M \end{array}$$

拉格朗日函数

$$\Lambda(x, \lambda) = f(x) + \sum_{m=1}^M \lambda_m h_m(x)$$

➤ 如果  $f(x^*)$  是原始约束优化问题的局部最优值，那么存在一个  $\lambda^*$  使得  $(x^*, \lambda^*)$  为拉格朗日函数  $\Lambda(x, \lambda)$  的驻点，满足

$$\begin{array}{l} \frac{\partial \Lambda(x, \lambda)}{\partial x} = 0 \\ \frac{\partial \Lambda(x, \lambda)}{\partial \lambda} = 0 \end{array} \quad \Rightarrow \quad \begin{array}{l} \nabla f(x) + \sum_{m=1}^M \lambda_m \nabla h_m(x) = 0 \\ h_m(x) = 0, \quad \forall m = 1, \dots, M \end{array}$$

➤ 拉格朗日乘数法是将一个有  $D$  个变量和  $M$  个等式约束条件的最优化问题转换为一个有  $D + M$  个变量的函数求驻点的问题

➤ 因为驻点不一定是最小解，所以在实际应用中需根据具体问题来验证是否为最小解

# 线性分类模型

## ➤ 补充：约束优化问题

➤ 约束优化问题通常使用**拉格朗日乘数法**来进行求解问题。

➤ 不等式约束优化问题

$$\begin{array}{ll} \min_{\mathbf{x}} & f(\mathbf{x}) \\ \text{s.t.} & h_m(\mathbf{x}) = 0, \quad m = 1, \dots, M \\ & g_n(\mathbf{x}) \leq 0, \quad n = 1, \dots, N \end{array}$$

拉格朗日函数

$$\Lambda(\mathbf{x}, \mathbf{a}, \mathbf{b}) = f(\mathbf{x}) + \sum_{m=1}^M a_m h_m(\mathbf{x}) + \sum_{n=1}^N b_n g_n(\mathbf{x})$$

➤ 等价于如下优化问题（主问题）

$$\begin{array}{ll} \min_{\mathbf{x}} \max_{\mathbf{a}, \mathbf{b}} & \Lambda(\mathbf{x}, \mathbf{a}, \mathbf{b}), \\ \text{s.t.} & \mathbf{b} \geq 0, \end{array}$$

对偶函数

$$\Gamma(\mathbf{a}, \mathbf{b}) = \inf_{\mathbf{x} \in \mathcal{D}} \Lambda(\mathbf{x}, \mathbf{a}, \mathbf{b}) \leq \Lambda(\tilde{\mathbf{x}}, \mathbf{a}, \mathbf{b}) \leq f(\tilde{\mathbf{x}})$$

令  $p^*$  是原问题的最优值

$$\Gamma(\mathbf{a}, \mathbf{b}) \leq p^*$$

$$\begin{array}{ll} \max_{\mathbf{a}, \mathbf{b}} & \Gamma(\mathbf{a}, \mathbf{b}), \\ \text{s.t.} & \mathbf{b} \geq 0. \end{array}$$

拉格朗日对偶问题

令  $d^*$  表示拉格朗日对偶问题的最优值

弱对偶性  
强对偶性

$$d^* \leq p^*$$

$$d^* = p^*$$

KKT条件

$$\begin{array}{ll} \nabla f(\mathbf{x}^*) + \sum_{m=1}^M a_m^* \nabla h_m(\mathbf{x}^*) + \sum_{n=1}^N b_n^* \nabla g_n(\mathbf{x}^*) = 0 \\ h_m(\mathbf{x}^*) = 0, & m = 1, \dots, M \\ g_n(\mathbf{x}^*) \leq 0, & n = 1, \dots, N \\ b_n^* g_n(\mathbf{x}^*) = 0, & n = 1, \dots, N \\ b_n^* \geq 0, & n = 1, \dots, N \end{array}$$

# 线性分类模型

## ➤ 支持向量机 (Support Vector Machine: SVM)

### ➤ 优化算法:

对偶问题

$$\begin{aligned} \max_{w,b} \quad & \gamma \\ \text{s.t.} \quad & \frac{y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b)}{\|\mathbf{w}\|} \geq \gamma, \forall n \in \{1, \dots, N\} \end{aligned}$$

限制  $\|\mathbf{w}\| \cdot \gamma = 1$

$$\begin{aligned} \max_{w,b} \quad & \frac{1}{\|\mathbf{w}\|^2} \\ \text{s.t.} \quad & y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b) \geq 1, \forall n \in \{1, \dots, N\}. \end{aligned}$$

$$\begin{aligned} \max_{\lambda \geq 0} \quad & \Gamma(\lambda) \\ \Gamma(\lambda) = \quad & -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_m \lambda_n y^{(m)} y^{(n)} (\mathbf{x}^{(m)})^\top \mathbf{x}^{(n)} + \sum_{n=1}^N \lambda_n. \end{aligned}$$

$$\begin{aligned} \Lambda(\mathbf{w}, b, \lambda) = \quad & \frac{1}{2} \|\mathbf{w}\|^2 + \sum_{n=1}^N \lambda_n (1 - y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b)), \\ \mathbf{w} = \quad & \sum_{n=1}^N \lambda_n y^{(n)} \mathbf{x}^{(n)}, \quad 0 = \sum_{n=1}^N \lambda_n y^{(n)}. \end{aligned}$$

$$\begin{aligned} \min_{w,b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & 1 - y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b) \leq 0, \quad \forall n \in \{1, \dots, N\}. \end{aligned}$$

$$\begin{aligned} f(\mathbf{x}) &= \text{sgn}(\mathbf{w}^{*\top} \mathbf{x} + b^*) \\ &= \text{sgn} \left( \sum_{n=1}^N \lambda_n^* y^{(n)} (\mathbf{x}^{(n)})^\top \mathbf{x} + b^* \right) \end{aligned}$$

最优预测函数



# 线性分类模型

## ➤ 支持向量机 (Support Vector Machine: SVM)

### ➤ 扩展: 非线性可分样本

➤ 容忍部分不满足约束的样本, 可以引入**松弛变量** (Slack Variable)  $\xi$ , 将优化问题变为

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & 1 - y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b) - \xi_n \leq 0, \quad \forall n \in \{1, \dots, N\} \\ & \xi_n \geq 0, \quad \forall n \in \{1, \dots, N\} \end{aligned}$$



$$\min_{\mathbf{w}, b} \quad \sum_{n=1}^N \max\left(0, 1 - y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b)\right) + \frac{1}{2C} \|\mathbf{w}\|^2,$$

Hinge 损失函数

正则化项



# 线性分类模型

## ➤ 支持向量机 (Support Vector Machine: SVM)

### ➤ 扩展: 非线性可分样本

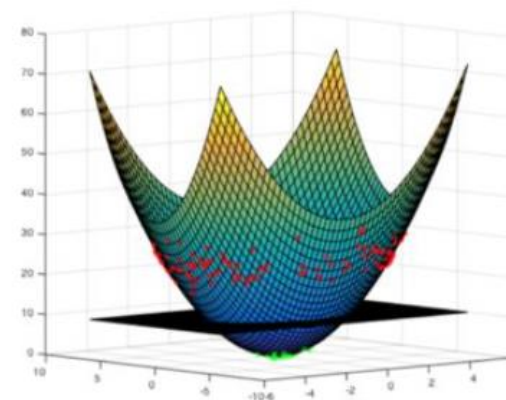
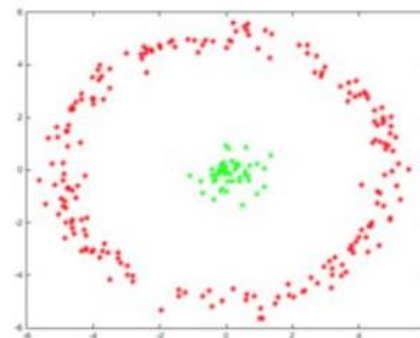
- 可以使用**核函数** (Kernel Function) 隐式地将样本从原始特征空间映射到更高维的空间, 并解决原始特征空间中的线性不可分问题

$$\begin{aligned} f(\mathbf{x}) &= \text{sgn}(\mathbf{w}^{*\top} \phi(\mathbf{x}) + b^*) \\ &= \text{sgn} \left( \sum_{n=1}^N \lambda_n^* y^{(n)} k(\mathbf{x}^{(n)}, \mathbf{x}) + b^* \right) \end{aligned}$$

其中  $k(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^\top \phi(\mathbf{z})$  称为核函数  
例:

$$k(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x}^\top \mathbf{z})^2 = \phi(\mathbf{x})^\top \phi(\mathbf{z}),$$

$$\phi(\mathbf{x}) = [1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1x_2, x_1^2, x_2^2]^\top$$



$$\Phi: \mathbf{R}^2 \rightarrow \mathbf{R}^3$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \Phi(\mathbf{x}) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{bmatrix} \in \mathbf{R}^3$$

# 总结

表 3.1 几种常见的线性模型对比

线性模型	激活函数	损失函数	优化方法
线性回归	-	$(y - \mathbf{w}^T \mathbf{x})^2$	最小二乘、梯度下降
Logistic 回归	$\sigma(\mathbf{w}^T \mathbf{x})$	$y \log \sigma(\mathbf{w}^T \mathbf{x})$	梯度下降
Softmax 回归	$\text{softmax}(\mathbf{W}^T \mathbf{x})$	$y \log \text{softmax}(\mathbf{W}^T \mathbf{x})$	梯度下降
感知器	$\text{sgn}(\mathbf{w}^T \mathbf{x})$	$\max(0, -y\mathbf{w}^T \mathbf{x})$	随机梯度下降
支持向量机	$\text{sgn}(\mathbf{w}^T \mathbf{x})$	$\max(0, 1 - y\mathbf{w}^T \mathbf{x})$	二次规划、SMO 等

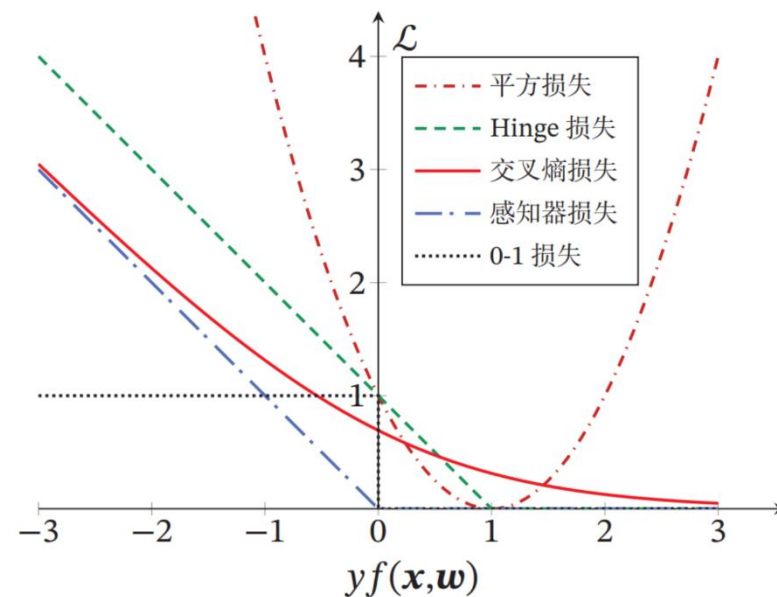


图 3.7 不同损失函数的对比

# 课后作业

## ➤ 完成

- 习题3-1
- 习题3-2
- 习题3-6



谢谢

