

《机器学习基础》



线性模型

大纲

▶ 线性分类模型

- ▶ Logistic Regression
- ▶ Softmax Regression
- ▶ Perceptron
- ▶ SVM

示例：图像分类

▶ 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



示例：图像分类

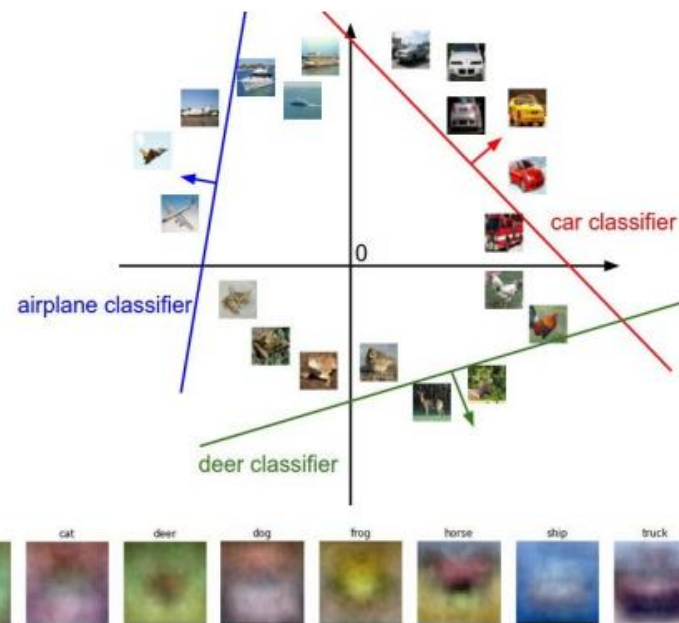
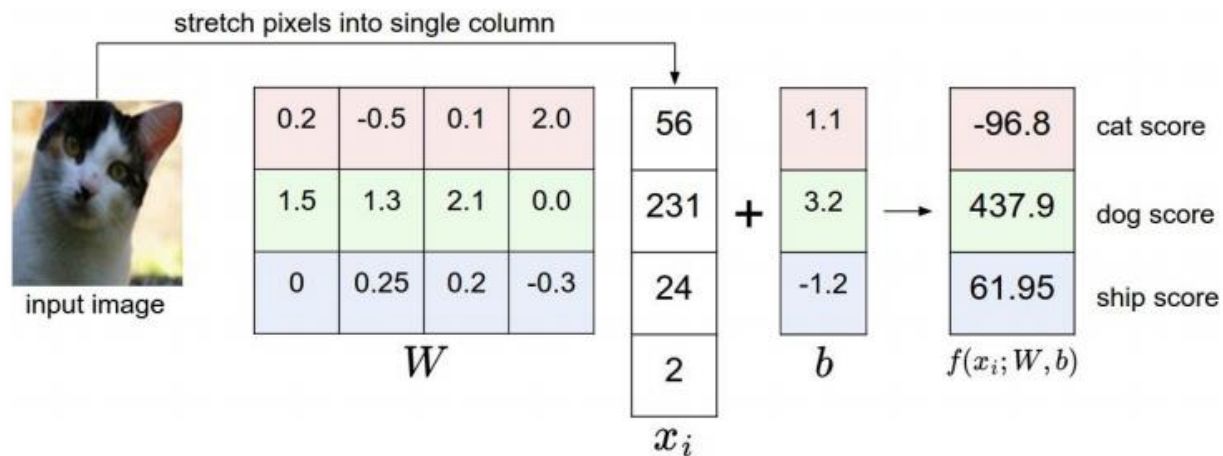


[32x32x3]

array of numbers 0...1
(3072 numbers total)

image parameters
 $f(\mathbf{x}, \mathbf{W})$

10 numbers, indicating
class scores



示例：图像分类、目标检测、实例分割

Classification



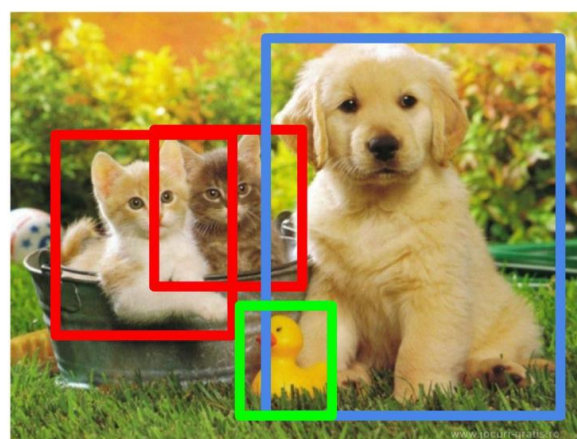
CAT

**Classification
+ Localization**



CAT

Object Detection



CAT, DOG, DUCK

**Instance
Segmentation**



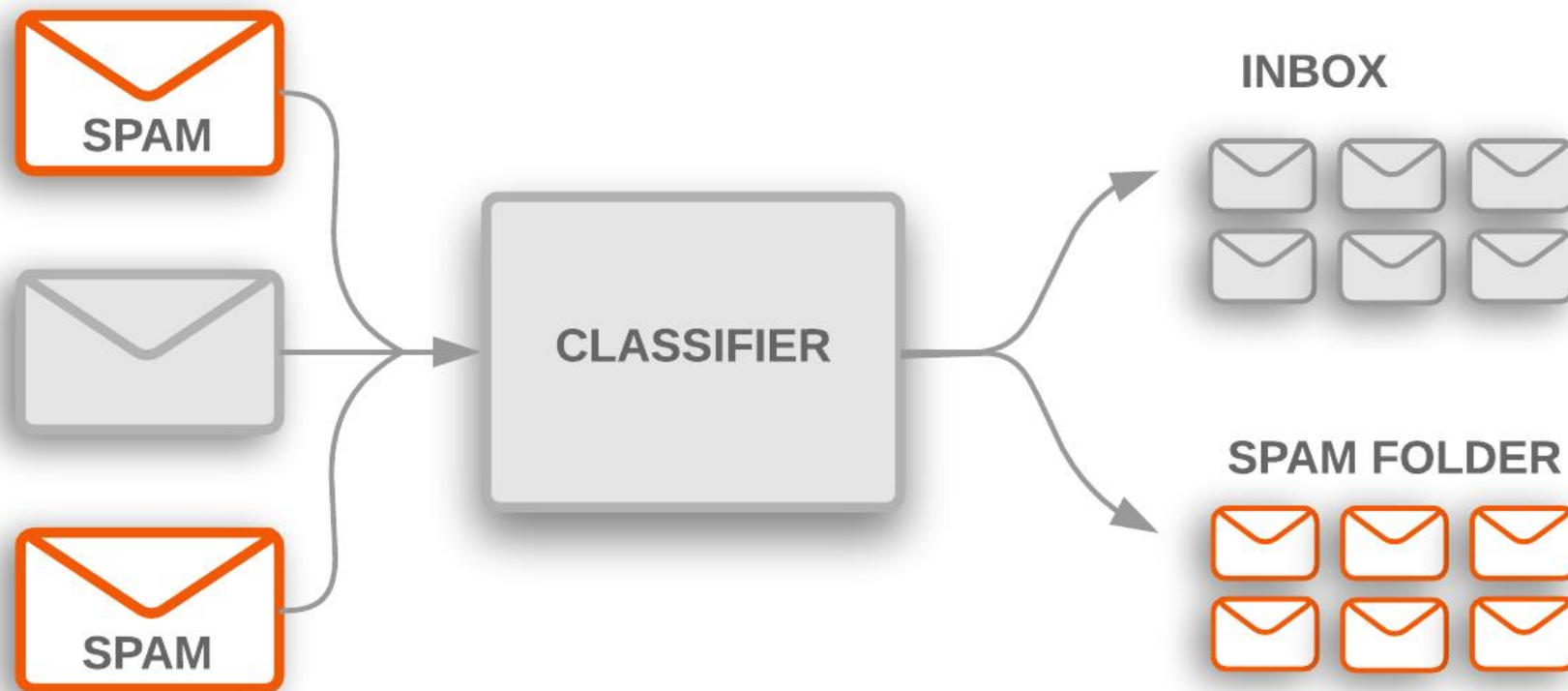
CAT, DOG, DUCK

Single object

Multiple objects

<https://medium.com/zylapp/review-of-deep-learning-algorithms-for-object-detection-c1f3d437b852>

示例：垃圾邮件过滤

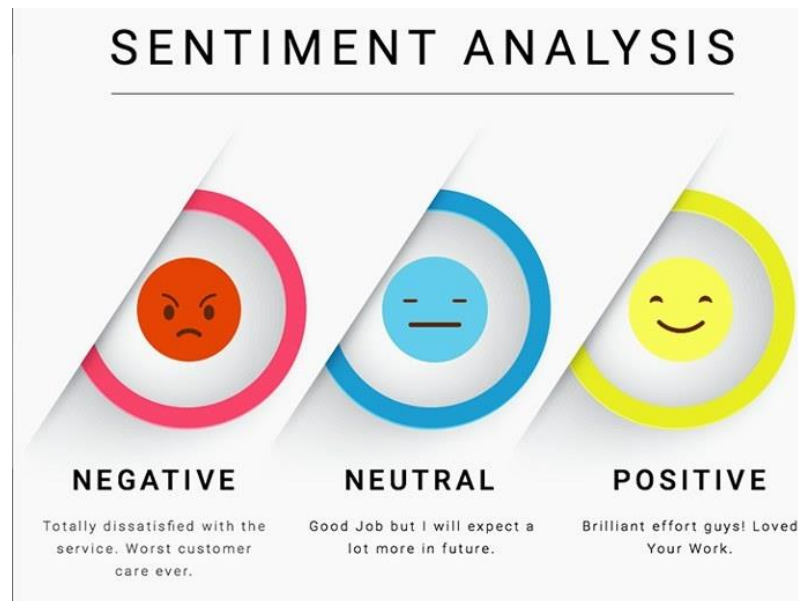


示例：文档归类



<https://towardsdatascience.com/automated-text-classification-using-machine-learning-3df4f4f9570b>

示例：情感分类



Review (X)

"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?"

Rating (Y)



示例：文本分类

- ▶ 将样本 x 从文本形式转为向量形式
- ▶ 词袋模型 (Bag-of-Words, BoW) 模型

the dog is on the table

| | | | | | | | |
|-----|-----|-----|----|-----|----|-------|-----|
| 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 |
| are | cat | dog | is | now | on | table | the |

比如两个文本“我 喜欢 读书”和“我 讨厌 读书”中共有“我”、“喜欢”、“讨厌”、“读书”四个词，它们的BoW表示分别为

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

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

示例：文本情感分类

根据文本内容来判断文本的相应类别

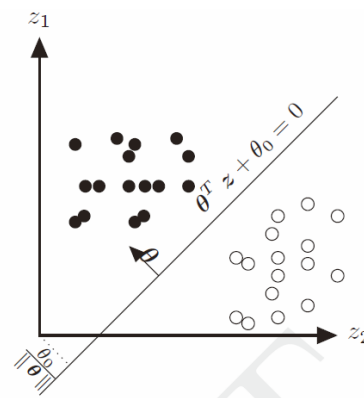
D_1 : “我喜欢读书”

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 |

+

-





线性分类模型

线性模型

▶ 线性模型 (Linear Model)

- ▶ 是机器学习中应用最广泛的模型，指通过样本特征的线性组合来进行预测的模型。
- ▶ 线性判别函数：

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

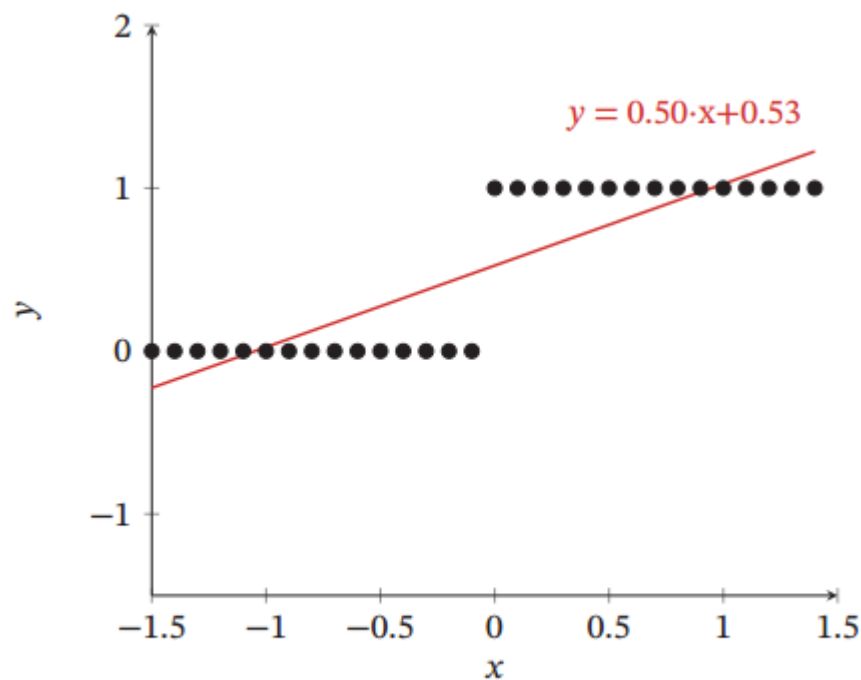
- ▶ 线性的决策函数：

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

线性决策函数

- ▶ 线性的决策函数预测分类问题的输出效果差，因为输出目标 y 是一些离散的标签

$$y = f(x; w)$$



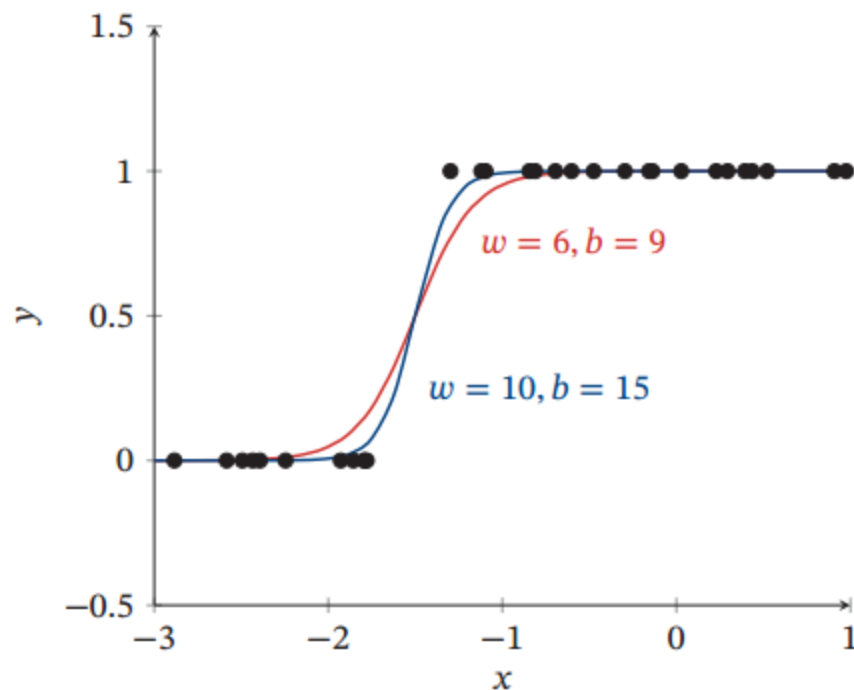
(a) 线性回归

非线性决策函数

- ▶ 引入一个非线性的决策函数来更好的进行拟合

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

- ▶ 把线性函数的值域从实数区间“挤压”到了(0,1)之间，可以用来表示概率。



线性模型

- ▶ Logistic 回归
- ▶ Softmax 回归
- ▶ 感知器
- ▶ 支持向量机



二分类与Logistic 回归

二分类

- ▶ 二分类问题的类别标签只有两种取值
- ▶ 通常可以设为 $\{+1, -1\}$ 或 $\{0, 1\}$.
- ▶ 决策函数

$$g(f(\mathbf{x}; \mathbf{w})) = \begin{cases} 1 & \text{if } f(\mathbf{x}; \mathbf{w}) > 0 \\ 0 & \text{if } f(\mathbf{x}; \mathbf{w}) < 0 \end{cases}$$

▶ 或

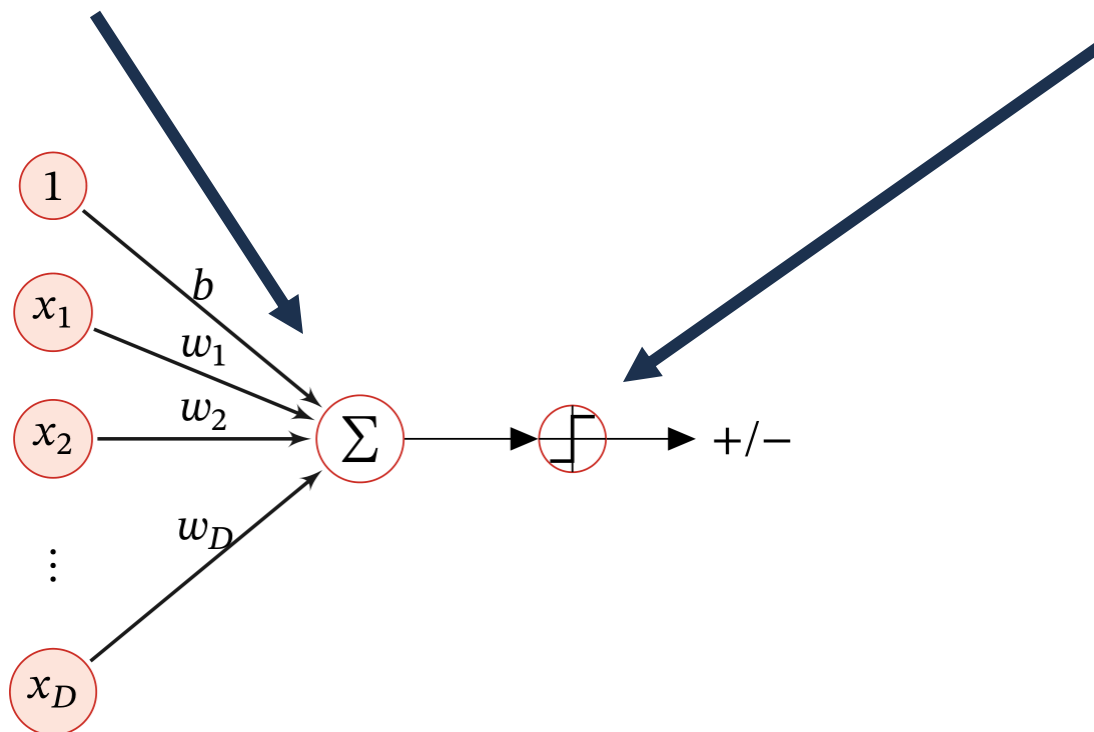
$$\begin{aligned} g(f(\mathbf{x}; \mathbf{w})) &= \text{sgn}(f(\mathbf{x}; \mathbf{w})) \\ &\triangleq \begin{cases} +1 & \text{if } f(\mathbf{x}; \mathbf{w}) > 0, \\ -1 & \text{if } f(\mathbf{x}; \mathbf{w}) < 0. \end{cases} \end{aligned}$$

二分类模型的基本框架

► 线性判别函数 + 二分类决策函数

$$f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^\top \mathbf{x} + b$$

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



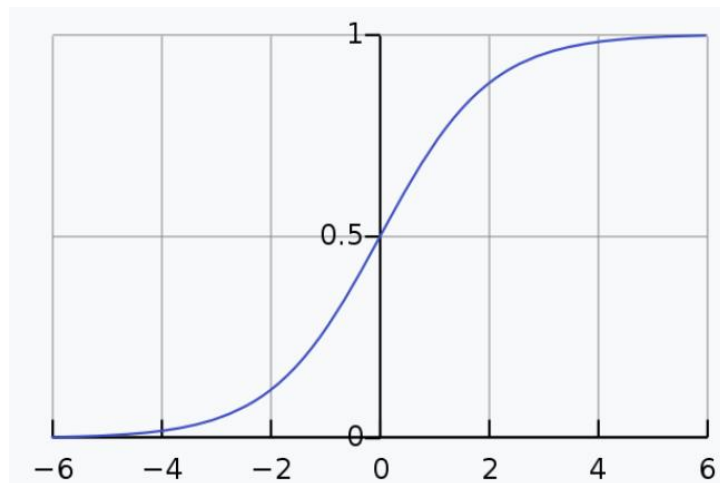
二分类模型的基本框架

- ▶ 逻辑回归的总体思路就是，先用逻辑函数把线性回归的结果 $(-\infty, \infty)$ 映射到 $(0, 1)$ ，再通过决策边界建立与分类的概率联系。
- ▶ 存在问题
 - ▶ 函数 g 非连续，无法求导 $g(f(\mathbf{x}; \mathbf{w})) = \text{sgn}(f(\mathbf{x}; \mathbf{w}))$
 - ▶ 学习速率较慢

Logistic函数与回归

► Logistic函数(Sigmoid函数)

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



► 优点

- 将任意实数映射到 $[0,1]$ 范围内
- 任意阶可导
- 在0附近几乎是线性的，但在两端趋于平缓
- 将异常值压向0或1

其他常用的激活函数

Hyperbolic tangent

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

Arctangent function

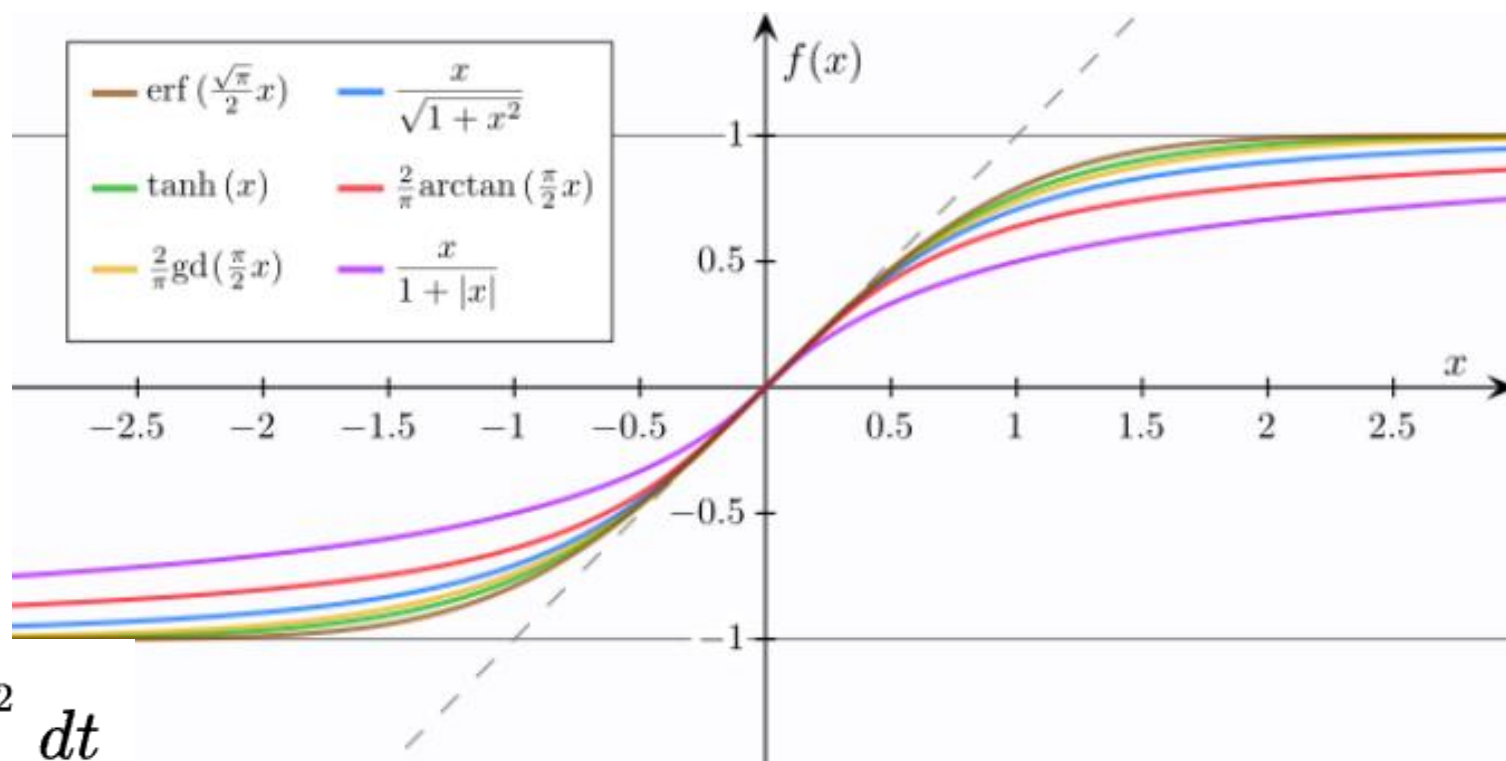
$$f(x) = \arctan x$$

Gudermannian function

$$f(x) = \operatorname{gd}(x) = \int_0^x \frac{dt}{\cosh t}$$

Error function

$$f(x) = \operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$$



Logistic函数与回归

► 对于一个二元逻辑回归，我们只需要确保

$$P(y = 1) + P(y = 0) = 1$$

$$P(y = 1) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$P(y = 0) = 1 - \sigma(z) = 1 - \frac{1}{1 + \exp(-z)} = \frac{\exp(-z)}{1 + \exp(-z)}$$

► 根据 Sigmoid 函数的性质

$$1 - \sigma(x) = \sigma(-x)$$

推导过程

► 将分类问题看作条件概率估计问题

- 基于已有数据以及学习权重类别标签的条件概率 $p(y = c|x)$ 最大化。
- 以二分类为例，

$$p(y = 1|\mathbf{x}) = g(f(\mathbf{x}; \mathbf{w}))$$

- 标签 $y = 1$ 的后验概率为

$$\begin{aligned} p(y = 1|\mathbf{x}) &= \sigma(\mathbf{w}^\top \mathbf{x}) \\ &\triangleq \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} \end{aligned}$$

推导过程

▶ 标签 $y = 0$ 的后验概率为

$$\begin{aligned} p(y = 0|\mathbf{x}) &= 1 - p(y = 1|\mathbf{x}) \\ &= \frac{\exp(-\mathbf{w}^\top \mathbf{x})}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}. \end{aligned}$$

▶ 在二分类模型中，事件发生与不发生的概率之比 $p/(1-p)$ 称为事件的几率。我们有

$$\mathbf{w}^\top \mathbf{x} = \log \frac{p(y = 1|\mathbf{x})}{1 - p(y = 1|\mathbf{x})}$$

风险函数

- ▶ 我们使用极大似然估计法来求解. 我们把单个样本看做一个事件, 那么这个事件发生的概率就是:

$$P(y|\mathbf{x}) = \begin{cases} p, y = 1 \\ 1 - p, y = 0 \end{cases}$$

- ▶ 它等价于:

$$P(y_i|\mathbf{x}_i) = p^{y_i} (1 - p)^{1-y_i}$$

- ▶ 则组样本的概率:

$$\begin{aligned} P &= P(y_1|\mathbf{x}_1)P(y_2|\mathbf{x}_2)P(y_3|\mathbf{x}_3)\dots P(y_m|\mathbf{x}_m) \\ &= \prod_{i=1}^m p^{y_i} (1 - p)^{1-y_i} \end{aligned}$$

风险函数

▶ 则组样本的概率：

$$\begin{aligned} P &= P(y^1 | \mathbf{x}^1) P(y^2 | \mathbf{x}^2) P(y^3 | \mathbf{x}^3) \dots P(y^m | \mathbf{x}^m) \\ &= \prod_{i=1}^m p^{y^{(i)}} (1 - p)^{1 - y^{(i)}} \end{aligned}$$

▶ 则似然函数为

$$L(w) = \prod_{i=1}^m P \left(y^{(i)} \mid \mathbf{w}^\top \mathbf{x}^{(i)} \right) = \prod_{i=1}^m \left(\sigma \left(\mathbf{w}^\top \mathbf{x}^{(i)} \right) \right)^{y^{(i)}} \left(1 - \sigma \left(\mathbf{w}^\top \mathbf{x}^{(i)} \right) \right)^{1 - y^{(i)}}$$

梯度下降

► 等式两边同时取对数：

$$\mathcal{R}(\mathbf{w}) = -\frac{1}{N} \sum_{i=1}^N \left(y^{(i)} \log \left(\sigma(\mathbf{w}^\top \mathbf{x}^{(i)}) \right) + (1 - y^{(i)}) \log \left(1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(i)}) \right) \right).$$

► 交叉熵损失函数，模型在训练集的风险函数为

$$\begin{aligned} l(w) = \log L(w) &= \sum_{i=1}^m \left(y^{(i)} \log \left(\sigma \left(\mathbf{w}^\top x^{(i)} \right) \right) + \left(1 - y^{(i)} \right) \log \left(1 - \sigma \left(\mathbf{w}^\top x^{(i)} \right) \right) \right) \\ &= -\frac{1}{N} \sum_{n=1}^N \left(y^{(n)} \log \hat{y}^{(n)} + (1 - y^{(n)}) \log(1 - \hat{y}^{(n)}) \right) \quad \hat{y}^{(n)} = \sigma(\mathbf{w}^\top \mathbf{x}^{(n)}) \end{aligned}$$

推导过程

► 梯度为：

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

$$\frac{\partial \mathcal{R}(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \left(\mathbf{x}^{(n)} \cdot \left(\sigma(\mathbf{w}^T \mathbf{x}^{(n)}) - y^{(n)} \right) \right)$$

$$\begin{aligned}\sigma'(z) &= \left(\frac{1}{1 + e^{-z}} \right)' = \frac{e^{-z}}{(1 + e^{-z})^2} \\ &= \frac{1 + e^{-z} - 1}{(1 + e^{-z})^2} \\ &= \frac{1}{(1 + e^{-z})} \left(1 - \frac{1}{(1 + e^{-z})} \right) \\ &= \sigma(z)(1 - \sigma(z))\end{aligned}$$

推导过程

► 采用梯度下降法，Logistic 回归的训练过程为：

► 初始化 $w_0 \leftarrow 0$ ，然后通过下式来迭代更新参数：

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



多分类与 Softmax 回归

多分类 (Multi-class Classification)

► 分类的类别数 C 大于 2，一般需要多个线性判别函数

► “argmax” 方式：共需要 C 个判别函数

$$f_c(\mathbf{x}; \mathbf{w}_c) = \mathbf{w}_c^\top \mathbf{x} + b_c, \quad c \in \{1, \dots, C\}$$

► “argmax” 方式的预测函数定义为

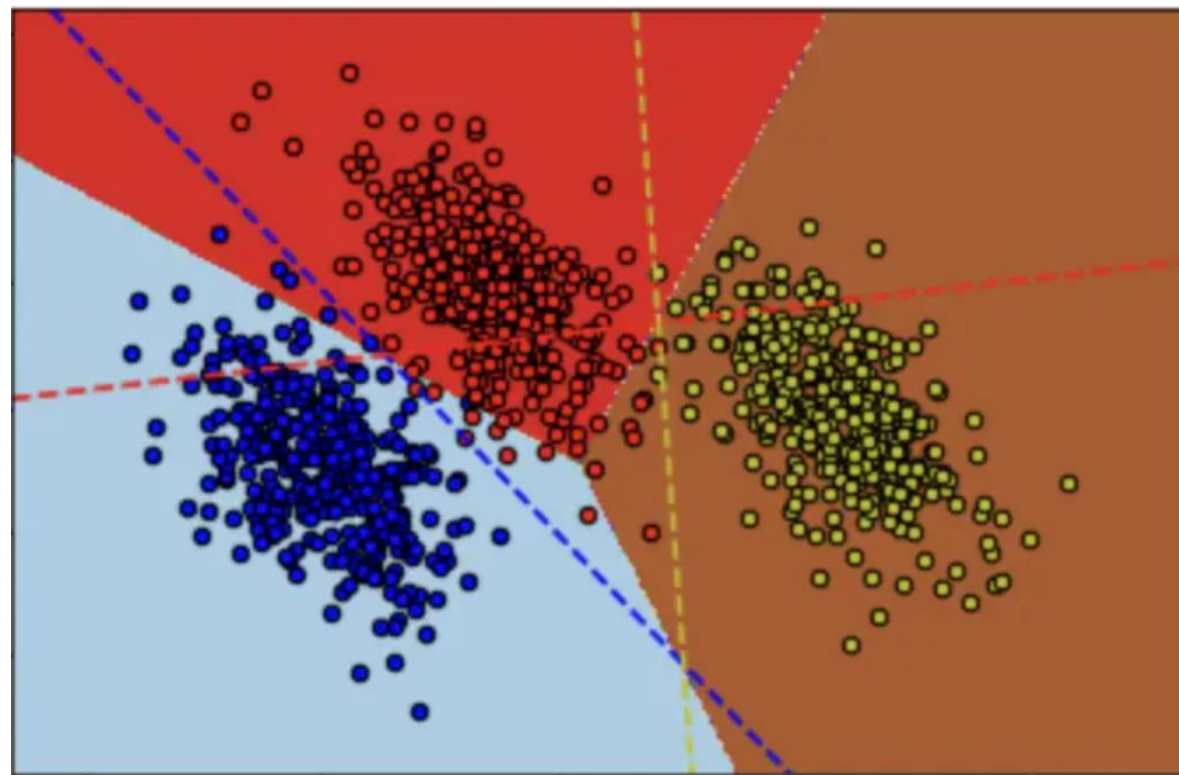
$$y = \arg \max_{c=1}^C f_c(\mathbf{x}; \mathbf{w}_c)$$

多分类 (Multi-class Classification)

▶ 相邻两类 i 和 j 的决策边界实际上是由

$$f_i(x; w_i) - f_j(x; w_j) = 0$$

▶ 决定，其法向量为 $w_i - w_j$.



Softmax回归

► Softmax函数

$$\text{softmax}(x_k) = \frac{\exp(x_k)}{\sum_{i=1}^K \exp(x_i)}$$

► 给定一个样本 x ，Softmax回归预测的属于类别 c 的条件概率为

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

Softmax回归

► 多分类

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

► Softmax回归的决策函数可以表示为

$$\begin{aligned} \hat{y} &= \arg \max_{c=1}^C p(y = c | \mathbf{x}) \\ &= \arg \max_{c=1}^C \frac{\exp(\mathbf{w}_c^\top \mathbf{x})}{\sum_{c'=1}^C \exp(\mathbf{w}_{c'}^\top \mathbf{x})} \end{aligned}$$

Softmax回归

▶ 向量形式可以写为

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

- ▶ 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,
- ▶ $\hat{\mathbf{y}} \in \mathbb{R}^C$ 为所有类别的预测条件概率组成的向量

参数学习

- ▶ 采用交叉熵损失函数，Softmax回归模型的风险函数为

$$\begin{aligned}\mathcal{R}(\mathbf{W}) &= -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \mathbf{y}_c^{(n)} \log \hat{\mathbf{y}}_c^{(n)} \\ &= -\frac{1}{N} \sum_{n=1}^N (\mathbf{y}^{(n)})^\top \log \hat{\mathbf{y}}^{(n)}, \quad \hat{\mathbf{y}}^{(n)} = \text{softmax}(\mathbf{W}^\top \mathbf{x}^{(n)})\end{aligned}$$

- ▶ 风险函数 $\mathcal{R}(\mathbf{W})$ 关于 \mathbf{W} 的梯度为

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

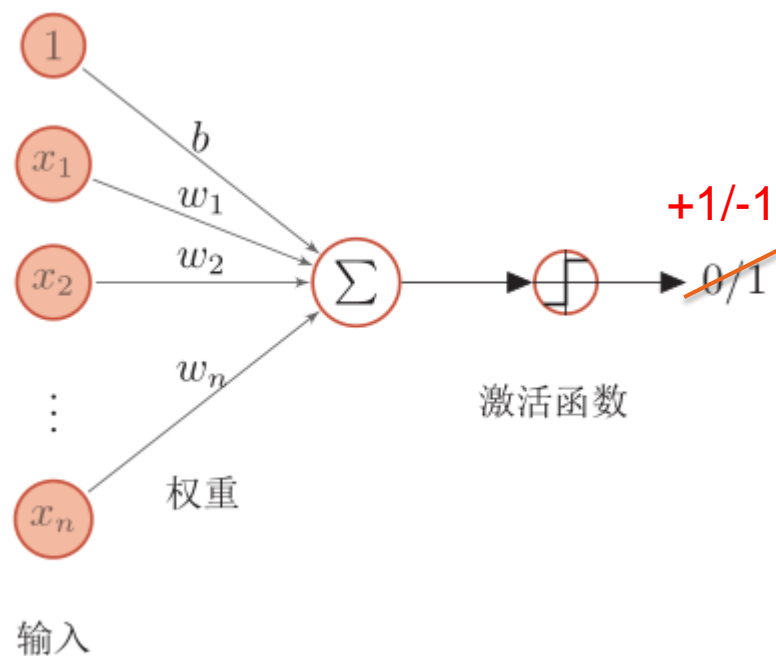


感知器

感知器

- ▶ 由 Frank Roseblatt 于 1957 年提出，是一种广泛使用的线性分类器。感知器可谓是最简单的人工神经网络，只有一个神经元
- ▶ 模拟生物神经元行为的机器，有与生物神经元相对应的部件，如权重（突触）、偏置（阈值）及激活函数（细胞体），输出为+1或-1。

$$\hat{y} = \text{sgn}(\mathbf{w}^T \mathbf{x}) = \begin{cases} +1 & \text{当 } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{当 } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases},$$



感知器

▶ 学习算法

- ▶ 感知器的学习算法是一种错误驱动的在线学习算法：
- ▶ 先初始化一个权重向量 $\mathbf{w} \leftarrow \mathbf{0}$ （通常是全零向量）；
- ▶ 每次分错一个样本 (\mathbf{x}, y) 时，即

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

- ▶ 用这个样本来更新权重

$$\mathbf{w} \leftarrow \mathbf{w} + y\mathbf{x}$$

- ▶ 根据感知器的学习策略，可以反推出感知器的损失函数为

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

感知器

► 采用随机梯度下降，其每次更新的梯度为

$$\frac{\partial \mathcal{L}(\mathbf{w}; \mathbf{x}, y)}{\partial \mathbf{w}} = \begin{cases} 0 & \text{if } y\mathbf{w}^\top \mathbf{x} > 0, \\ -y\mathbf{x} & \text{if } y\mathbf{w}^\top \mathbf{x} < 0. \end{cases}$$

感知器的学习过程

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:  $\mathbf{w}_0 \leftarrow 0, k \leftarrow 0, t \leftarrow 0$ ;  
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  $\mathbf{w}_k^\top (y^{(n)} \mathbf{x}^{(n)}) \leq 0$  then  
7        $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + y^{(n)} \mathbf{x}^{(n)}$ ;  
8        $k \leftarrow k + 1$ ;  
9     end  
10     $t \leftarrow t + 1$ ;  
11    if  $t = T$  then break;  
12  end  
13 until  $t = T$ ;  
    输出:  $\mathbf{w}_k$ 
```

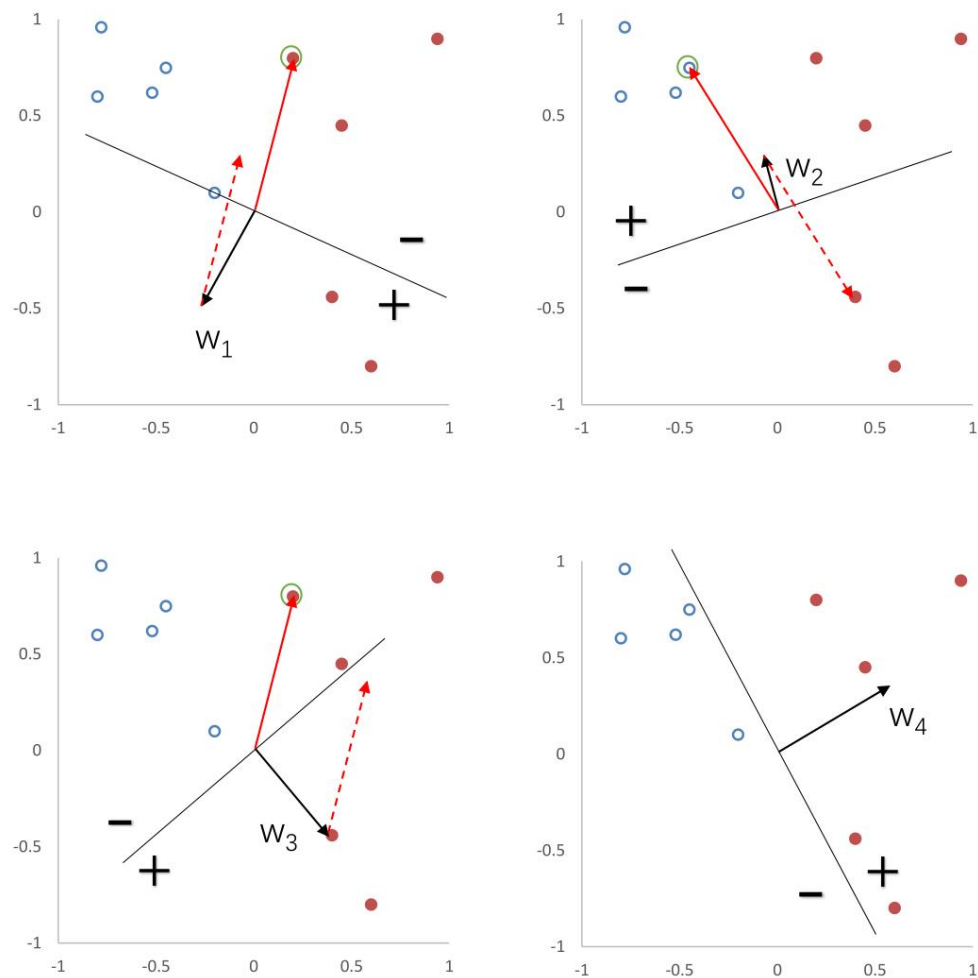
表示分错

// 达到最大迭代次数

对比Logistic回归的更新方式:

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

感知器参数学习的更新过程





支持向量机

支持向量机

- ▶ 是在分类与回归分析中分析数据的监督式学习模型与相关的学习算法。
- ▶ 支持向量机的分类方法，是在一组分布中找出一个超平面作为决策边界，使模型在数据上的分类误差尽量接近于零
- ▶ 支持向量机在线性和非线性分类中，效果都非常好。

支持向量机

- ▶ 给定一个二分类器数据集

$$\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N \quad y_n \in \{+1, -1\}$$

- ▶ 如果两类样本是线性可分的，即存在一个超平面

$$\mathbf{w}^\top \mathbf{x} + b = 0$$

- ▶ 将两类样本分开，那么对于每个样本都有

$$y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b) > 0.$$

- ▶ 每个样本 $\mathbf{x}^{(n)}$ 到分割超平面的距离为

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

支持向量机

- ▶ 支持向量机的目标是寻找一个超平面 (\mathbf{w}^*, b^*) 使得 γ 最大, 即

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

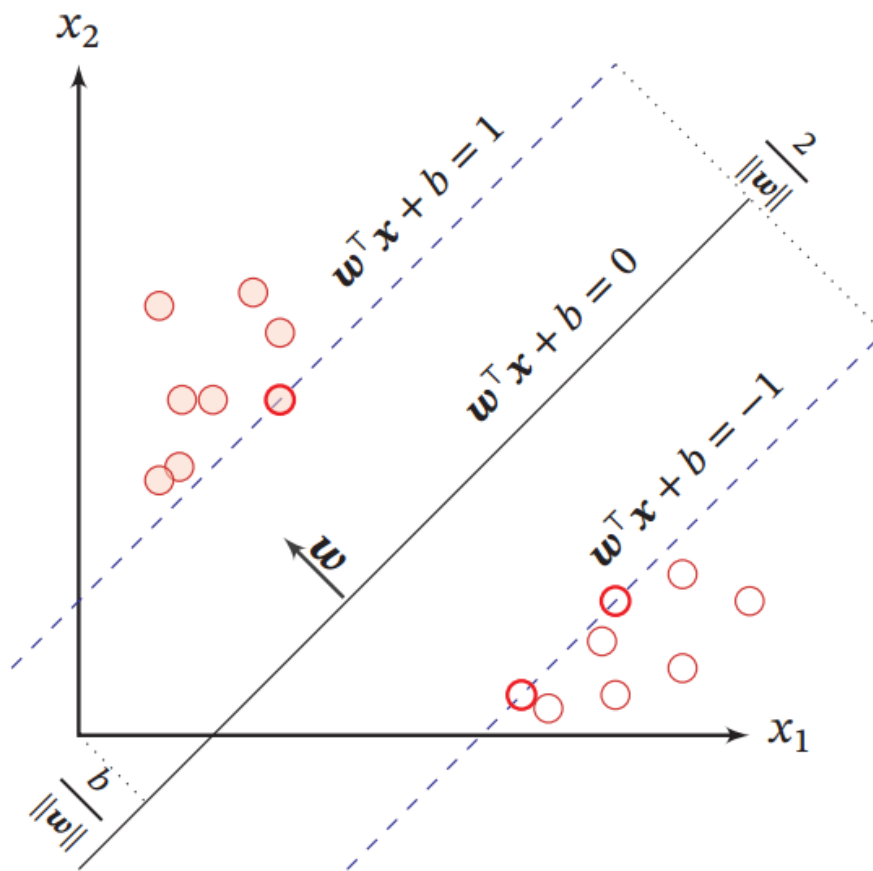
- ▶ 等价于

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

- ▶ 数据集中所有满足 $y^{(n)}(\mathbf{w}^\top \mathbf{x}^{(n)} + b) = 1$ 的样本点, 都称为支持向量

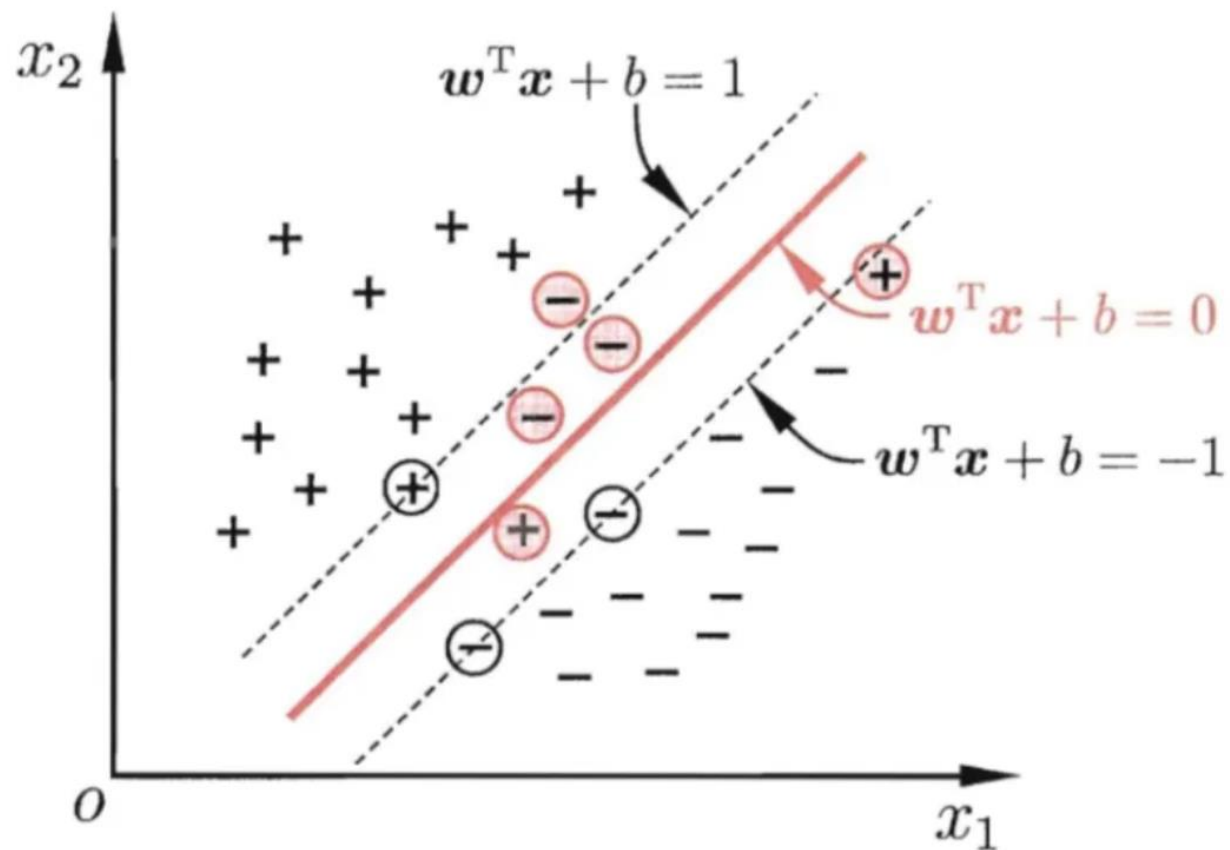
支持向量机

- ▶ 其分割超平面有很多个，但是间隔最大的超平面是唯一的

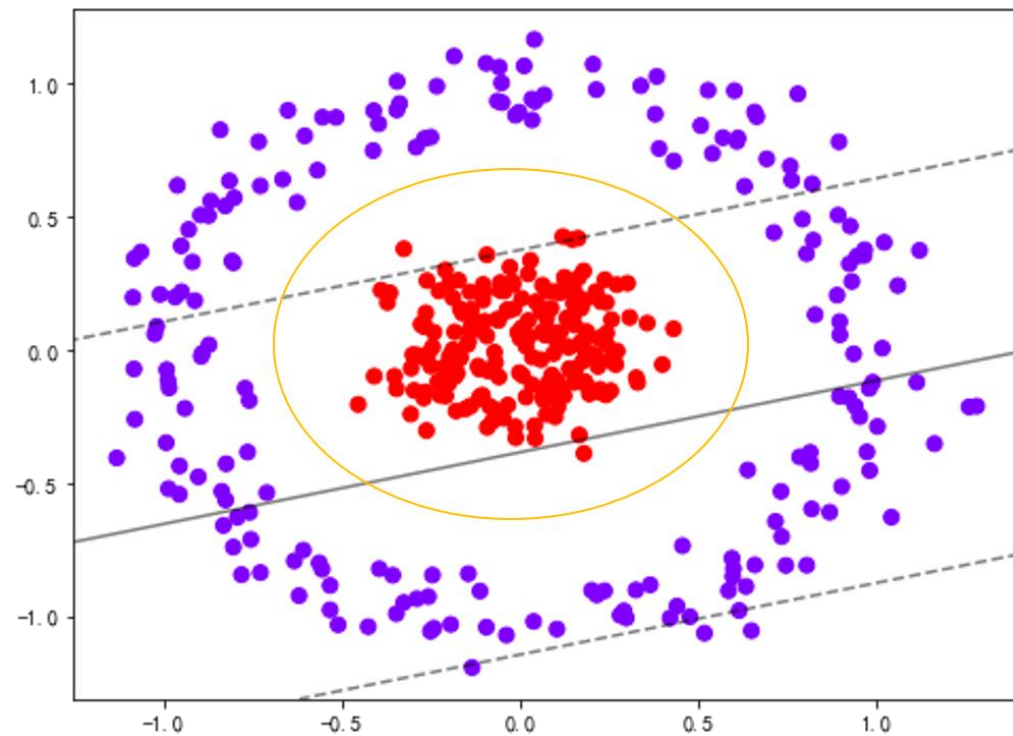


支持向量机

- ▶在最大化间隔的同时，不满足约束的样本应尽可能地少

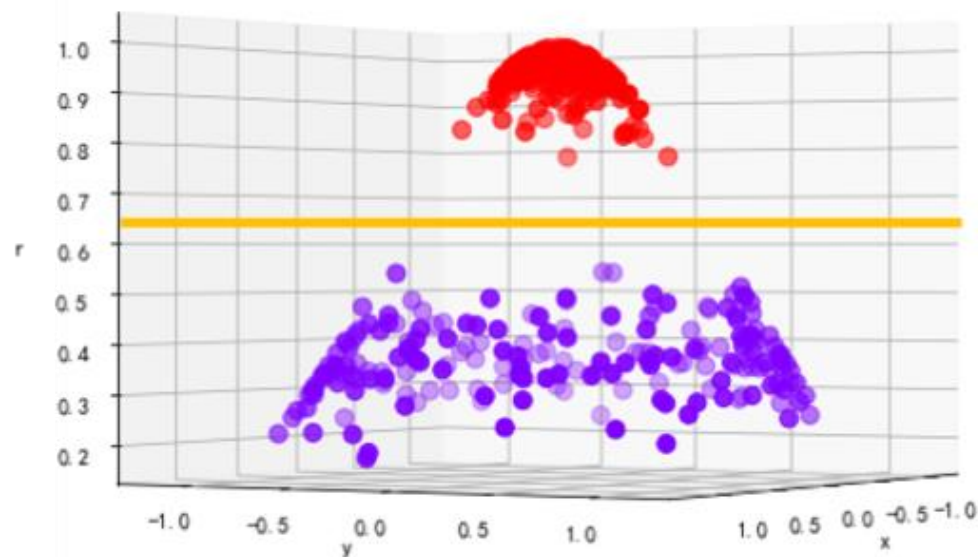
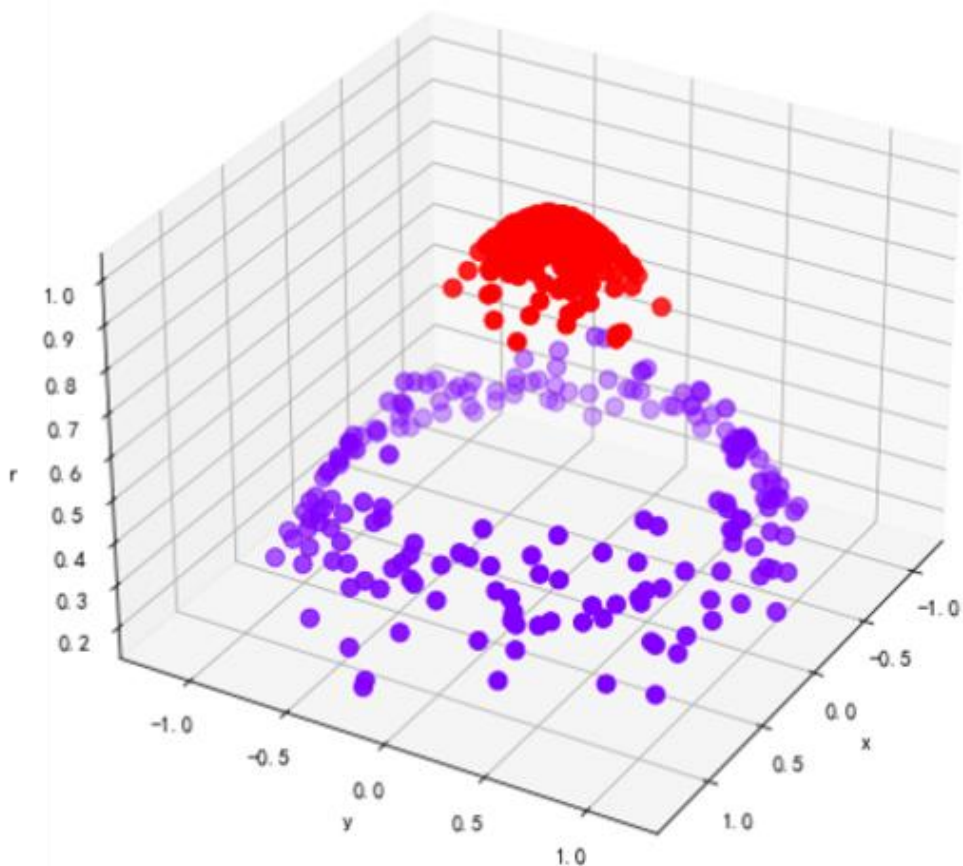


非线性支持向量机



非线性支持向量机

- ▶ 原理利用核技巧，将数据从原始的空间投射到新空间中
- ▶ 这种非线性变换，将非线性问题变换成线性问题





总结

线性分类模型小结

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