

机器学习

Machine Learning

北京航空航天大学计算机学院

School of Computer Science and Engineering, Beihang University

刘庆杰 陈佳鑫

2025年春季学期

Spring 2025

课前回顾

内容提要

- 2.1 数据集的划分方法
- 2.2 模型的性能度量
- 2.3 概率统计基础
- 2.4 贝叶斯决策理论
- 2.5 参数化概率密度估计方法
- 2.6 非参数概率密度估计方法
- 2.7* 矩阵理论基础（拓展）

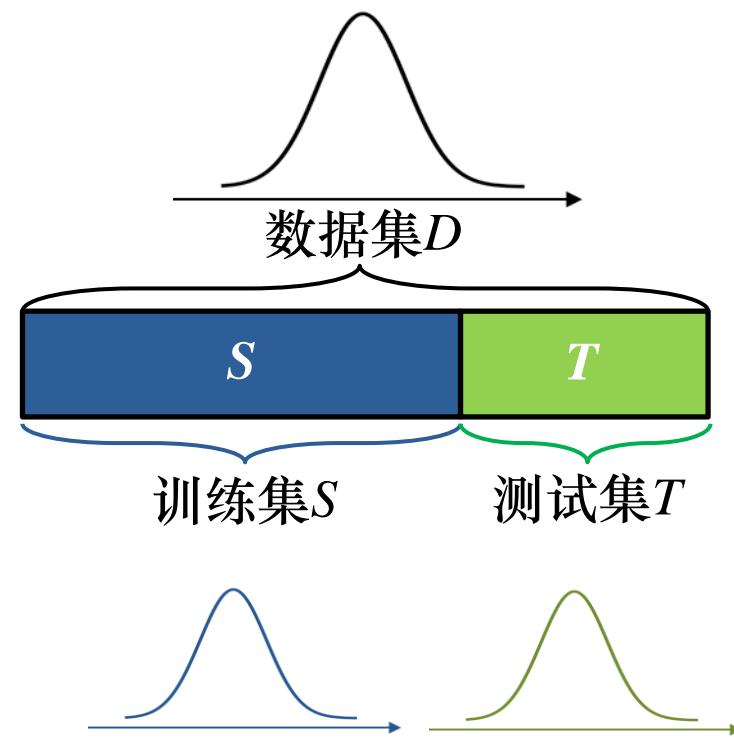
数据集划分的基本要求

● 模型评估

➤ 1. 数据集划分：分为训练集和测试集两部分

- 原则：测试集应尽可能与训练集互斥，测试样本不在训练集中出现
- 目标：将数据集 D 划分为训练集 S 和测试集 T 两部分，在训练集上建立模型，在测试集上评估性能
- 假设：测试样本从原样本真实分布中独立同分布采样得到
- 方法：留出法、自助法、交叉验证法

➤ 2. 性能度量：模型在测试集（新样本）上进行度量，也叫泛化性能



模型的性能度量

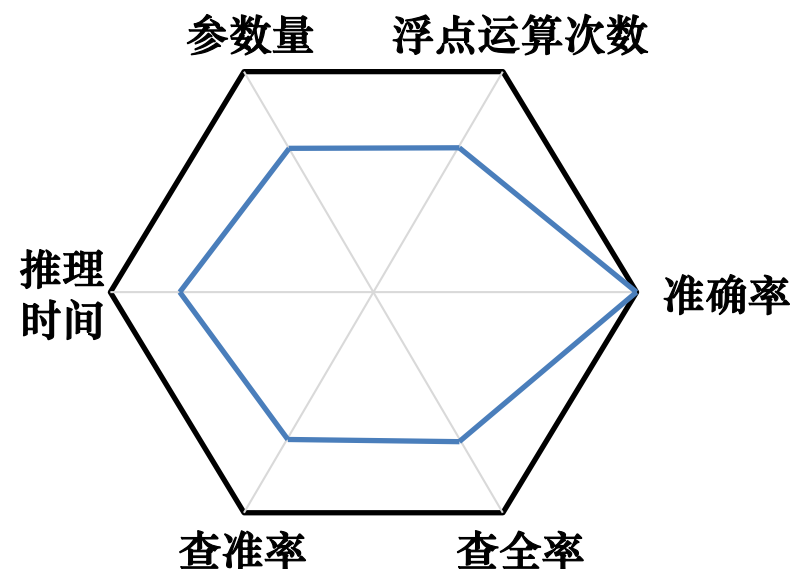
● 模型评估

- 1. 数据集划分：分为训练集和测试集两部分
- 2. **性能度量**：性能度量：模型在测试集（新样本）上进行度量，也叫泛化性能
 - 常用性能度量：错误率和精度（分类任务）
 - 仅能评估是否正确分类，无法提供更全面的评估

示例1：发动机合格检测

精度可以评估“检测出有多少发动机是合格的”

无法评估“检测出的合格发动机有多少是真正合格的”



一些基本概念

- **概率 (Probability)**

- 对随机事件发生可能性大小的度量

- **联合概率 (Joint Probability)**

- A和B共同发生的概率，称事件A和B的联合概率，记作 $P(A, B)$ 或 $P(A \cap B)$

- **条件概率 (Conditional Probability)**

- 事件B已发生的条件下，事件A发生的概率，记作 $P(A|B)$

贝叶斯公式

● 贝叶斯公式 (Bayes' Theorem)

- 贝叶斯公式给出了“结果”事件A已经发生的条件下，“原因”事件B的条件概率，对结果的任何观测都将增加我们对原因事件B的真正分布的知识

后验概率：给定观测数据后，某事件发生的概率

先验概率：在没有观测数据之前，某事件发生的初始概率

似然概率：给定事件发生的情况下，观测数据出现的概率

$$P(B_i|A) = \frac{P(B_i)P(A|B_i)}{\sum_{j=1}^n P(B_j)P(A|B_j)}$$

证据因子：观测数据的边际概率，即所有可能事件下观测数据的总概率

贝叶斯决策

- 贝叶斯决策是统计决策理论中的一个基本方法，用于解决分类问题

- 已知条件：

- 1、属于一定数量类别的数据，类别为标签为： $\omega_i, i = 1, 2, \dots, c$

- 2、各类别 ω_i 的类先验概率 $P(\omega_i)$ 和类条件概率密度 $P(x|\omega_i)$

- 基本思想：根据贝叶斯公式计算后验概率，基于最大后验概率进行判决

- 判决函数：最大化后验概率

$$x \in \omega_k \text{ 当且仅当 } k = \arg \max_i \{P(\omega_i|x)\}, \text{ 其中 } P(\omega_i|x) = \frac{P(x|\omega_i)P(\omega_i)}{\sum_{j=1}^c P(x|\omega_j)P(\omega_j)}$$

2.5 参数化概率密度估计方法

- 概率密度估计的概念
- 参数化概率密度估计方法的概念
- 极大似然估计方法
- 贝叶斯估计方法
- 估计量的性质与评价标准

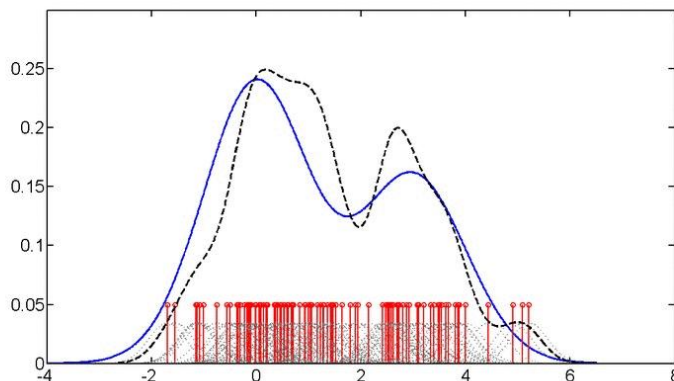
概率密度估计

- 概率密度估计的任务:

- 根据观测样本数据估计类条件概率密度函数 $P(x|\omega_i)$ 和类先验概率 $P(\omega_i)$

- 为什么需要估计概率密度

- 概率密度估计可以建模原始数据分布，辅助精细地了解数据特性。进而帮助识别数据中的异常值、或者生成新数据



- 概率密度估计的方法:

- 参数化方法: 已知概率密度函数的形式, 其中几个参数未知
- 非参数化方法: 概率密度函数的形式未知

参数化概率密度估计

● 参数化概率密度估计的任务

- 已知：概率密度函数的形式，若干参数未知
- 目标：依据样本估计未知参数的值

● 典型方法

- 极大似然估计：把待估计参数看做是确定的量，只是其取值未知。最佳估计就是使产生已观测到样本的概率最大的那个值
- 贝叶斯估计：把待估计参数看做是符合某种先验概率分布的随机变量。对样本进行观测的过程，就是把先验概率密度转化为后验概率密度，从而利用样本信息修正参数的初始估计值的过程

极大似然估计

● 极大似然估计的假设条件

- $P(x|\omega_i)$ 具有某种**确定的解析函数形式**，只有部分参数 θ 未知；
- 参数 θ 通常为向量，如一维正态分布 $N(\mu, \sigma^2)$ 中的 μ 、 σ
- 参数 θ 是确定的未知量，不是随机量
- 各类样本集 $x_i, i=1,2,\dots,c$ 满足**独立同分布**条件(*i.i.d.*)，即 x_i 均为从概率密度函数为 $P(x|\omega_i)$ 的总体分布中独立抽取出来的
- 各类样本只包含本类分布的信息；因此， $P(x|\omega_i)$ 可记为 $P(x|\omega_i; \theta_i)$ 或 $P(x; \theta_i)$

● 基于上述假设，各类条件概率密度可根据各类样本分别估计

极大似然估计

● 似然函数

- 针对一类已知样本 $X = \{x_i, i = 1, 2, \dots, N\}$ ，定义参数 θ 下观测到样本集 X 的联合分布概率密度，称为相对于样本集 X 的 θ 的似然函数

$$l(\theta) = P(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n P(x_i; \theta)$$

● 基本思想

- 在 θ 可能的取值范围内选择使似然函数达到最大的参数值作为参数 θ 的估计值
- 形式化描述为：求 $\hat{\theta}$ ，使得 $l(\hat{\theta}) = \max_{\theta} l(\theta)$
- 如果参数 $\theta = \hat{\theta}$ 时， $l(\theta)$ 最大，则 $\hat{\theta}$ 是最可能的参数估计值。它是关于样本集的函数，记作： $\hat{\theta} = d(x_1, x_2, \dots, x_N) = d(X)$ ，称为极大似然估计量
- 为便于分析求解，实际应用中往往采用对数似然函数： $H(\theta) = \ln l(\theta)$

极大似然估计

● 求解

- 若似然函数连续可微，最大似然函数估计量就是方程 $\frac{dl(\theta)}{d\theta} = 0$ 或 $\frac{dH(\theta)}{d\theta} = 0$ 的解
- 若未知参数不止一个，即 $\theta = [\theta_1, \theta_2, \dots, \theta_s]^T$ ，则需联立以下 s 个方程组求解

$$\frac{dH(\theta)}{d\theta_i} = 0, i = 1, 2, \dots, s$$

● 极大似然函数的求解性质

- 若似然函数连续可导，存在最大值且必要条件方程有唯一解，则该解就是极大似然估计量
- 如果似然函数有多个解，则使似然函数值最大者为极大似然估计量
- 若似然函数单调，可根据极大似然思想，将似然函数最大值点作为参数的极大似然估计值

极大似然估计

● 极大似然估计的求解步骤

- 1、确定需要估计的概率分布 $P(x|\theta)$ ，其中 $\theta = [\theta_1, \theta_2, \dots, \theta_s]^T$
- 2、构造似然函数 $l(\theta) = P(x_1, x_2, \dots, x_N; \theta) = \prod_{i=1}^N P(x_i; \theta)$
- 3、求对数似然函数 $H(\theta) = \ln l(\theta)$
- 4、令 $\frac{dH(\theta)}{d\theta_i} = 0, i = 1, 2, \dots, s$ ，联立求解

极大似然估计

- 示例：单变量正态分布

- 已知

- 参数: $\theta = [\theta_1, \theta_2]$, $\theta_1 = \mu, \theta_2 = \sigma^2$

- 概率密度函数: $P(x|\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$

- 样本集: $X = \{x_1, x_2, \dots, x_N\}$

- 目标: 估计参数 $[\theta_1, \theta_2]$

极大似然估计示例：单变量正态分布

● 解：

➤ 1、求似然函数： $l(\theta) = P(x_1, x_2, \dots, x_n; \theta) = \prod_{i=1}^n P(x_i; \theta)$

➤ 2、求对数似然函数： $H(\theta) = \ln l(\theta) = \sum_{i=1}^N \ln P(x_i; \theta)$

➤ 3、构造方程组：
$$\begin{cases} \frac{\partial H}{\partial \mu} = 0 \\ \frac{\partial H}{\partial \sigma^2} = 0 \end{cases}, \begin{cases} \frac{1}{\sigma^2} [\sum_{i=1}^N x_i - N\mu] = 0 \\ -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^N (x_i - \mu)^2 = 0 \end{cases}$$

➤ 4、联立求解： $\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2, \quad \hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$

贝叶斯估计

● 贝叶斯估计的基本思想

- 把待估计参数 θ 看作是具有先验分布 $P(\theta)$ 的随机变量，其取值与样本集 X 有关，贝叶斯估计利用样本集 X 将 θ 的先验概率分布修正为后验概率分布
- 贝叶斯决策用于分类，计算离散形式的后验概率值，而贝叶斯估计则用于回归，计算连续形式的后验概率密度函数

● 贝叶斯估计损失函数

- 把 θ 估计为 $\hat{\theta}$ 所造成的损失，记为 $\lambda(\hat{\theta}, \theta)$
- 不同于离散形式贝叶斯决策的损失表，由于参数化概率密度估计为连续值估计，因此常采用损失函数，常用平方误差损失函数 $\lambda(\hat{\theta}, \theta) = (\theta - \hat{\theta})^2$

贝叶斯估计

● 贝叶斯估计相关概念

➤ **条件期望损失：** $R(\hat{\theta}|x) = \int_{\Theta} \lambda(\hat{\theta}, \theta) P(\theta|x) d\theta$ ，其中， $x \in E^d$ ， $\theta \in \Theta$ ， E^d 为样本集， Θ 为待估参数集

➤ **期望风险：** $R = \int_{E^d} \int_{\Theta} \lambda(\hat{\theta}, \theta) P(x; \theta) d\theta dx$
 $= \int_{E^d} \int_{\Theta} \lambda(\hat{\theta}, \theta) P(\theta|x) P(x) d\theta dx = \int_{E^d} R(\hat{\theta}|x) P(x) dx$

➤ **贝叶斯估计量：** 使条件期望损失最小的估计量 $\hat{\theta}$ ，即

$$\hat{\theta} = \operatorname{argmin} \left(R(\hat{\theta}|x) \right) = \operatorname{argmin} \left(\int_{\Theta} \lambda(\hat{\theta}, \theta) P(\theta|x) d\theta \right)$$

➤ **定理2.5.1：** 若采用平方误差损失函数，则 θ 的贝叶斯估计量是在给定样本集 X 时 θ 的条件期望，即 $\hat{\theta} = E(\theta|x) = \int_{\Theta} \theta P(\theta|x) d\theta$

贝叶斯估计

● 贝叶斯估计步骤

➤ 1、确定参数 θ 所遵从的先验分布： $P(\theta)$

➤ 2、求样本集的联合分布： $P(X|\theta) = \prod_{i=1}^N P(x_i|\theta)$

➤ 3、求 θ 的后验概率分布：
$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{\int_{\Theta} P(X|\theta)P(\theta)d\theta}$$

➤ 4、求 θ 的贝叶斯估计量（定理2.5.1）： $\hat{\theta} = \int_{\Theta} \theta P(\theta|x)d\theta$

贝叶斯估计

- 示例：单变量正态分布

- 已知

- 参数： $\theta = [\theta_1, \theta_2]$, $\theta_1 = \mu, \theta_2 = \sigma^2$, 其中 θ_2 已知

- 概率密度函数： $P(x|\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$

- 样本集： $X = \{x_1, x_2, \dots, x_N\}$

- 目标：估计均值 θ_1

贝叶斯估计示例：单变量正态分布

● 解：（仅保留关键步骤）

- 1、确定概率密度函数形式： $P(x|\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$
- 2、设估计量 μ 遵从先验分布： $\mu \sim \mathcal{N}(\mu_0, \sigma_0^2)$
- 3、根据观测样本 x 求得 μ 的后验分布： $\mu|x \sim \mathcal{N}\left(\frac{\mu_0\tau_0 + n\bar{x}\tau}{\tau_0 + n\tau}, \frac{1}{\tau_0 + n\tau}\right)$ ，其中 $\tau = \frac{1}{\sigma^2}$ ， $\tau_0 = \frac{1}{\sigma_0^2}$
- 4、 μ 的贝叶斯估计量为

$$E(\mu|x) = \frac{\mu_0\tau_0 + n\bar{x}\tau}{\tau_0 + n\tau} = w\mu_0 + (1-w)\bar{x}, \text{ 其中 } n \text{ 为样本数, } w = \frac{\tau_0}{\tau_0 + n\tau}$$

- 5、当 $n = 0$ 时，估计量 $\hat{\mu} = \mu_0$ ，当 $n \rightarrow \infty$ 时，估计量 $\hat{\mu} = \bar{x}$

估计量的性质与评价标准

● 估计量的性质

无偏性	渐进无偏性	有效性	一致性
$E[\hat{\theta}(x_1, x_2, \dots, x_N)] = \theta$	$E[\hat{\theta}(x_N)] \xrightarrow{N \rightarrow \infty} \theta$	对估计 $\hat{\theta}_1, \hat{\theta}_2$, 若方差 $\sigma^2(\hat{\theta}_1) < \sigma^2(\hat{\theta}_2)$, 则 $\hat{\theta}_1$ 估计更有效	当样本数无穷多时, 每一次估计都在概率意义上任意接近真实值, 即: $\forall \varepsilon > 0, \lim_{N \rightarrow \infty} P(\hat{\theta}_N - \theta > \varepsilon) = 0$

- 结合无偏性和有效性, 要求估计量能够在多次估计中, 以较小的方差平均地表示真实值
- 极大似然估计是无偏的, 在样本充足时有效性更好, 一致性更强;
- 贝叶斯估计可能有偏, 在样本量有限且有合理的先验信息时更有效。

2.6 非参数概率密度估计方法

- 什么是非参数估计?
- Parzen窗算法
- k近邻算法

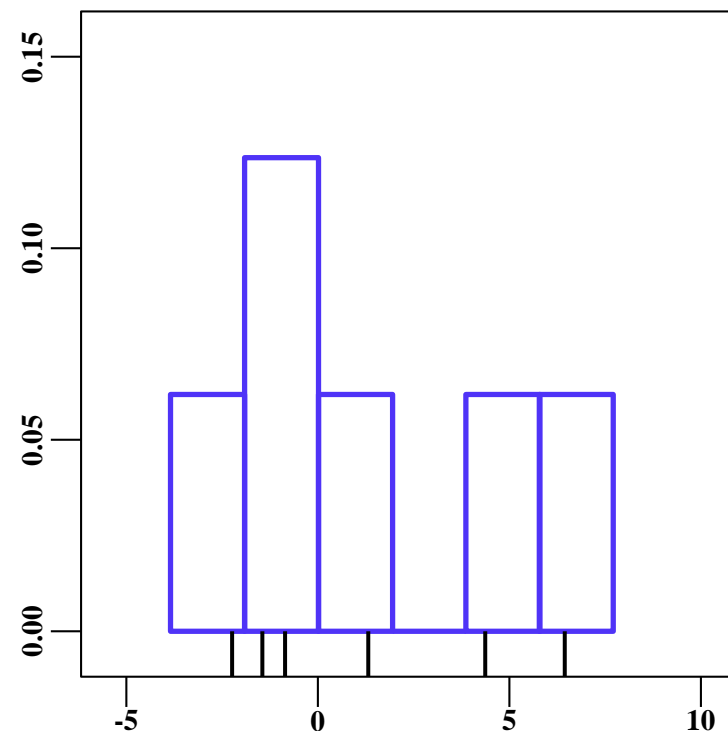
什么是非参数估计?

- 非参数估计

- 概率密度函数的形式未知的模型，直接依赖于数据本身来进行推断和估计

- 常见的非参数估计方法

- Parzen窗法
- k近邻法



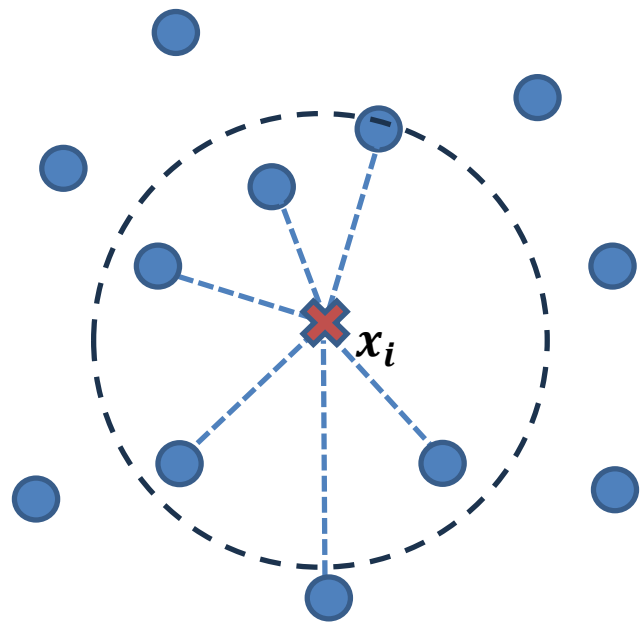
非参数估计方法

● 非参数估计方法

➤ 基本思路：要估计 x_i 点的密度，可把相关样本在该点的“贡献”相加近似作为其概率密度，进而可以以此方法估计每个点的概率密度

➤ 具体流程

- 1.计算 x_i 点处与概率密度相关的贡献点
- 2.确定贡献点对 x_i 点处的贡献
- 3.重复步骤1-2计算所有点处的概率密度



非参数估计方法

● 非参数概率密度估计

- 假设 N 为样本总数，以 x_i 为中心的区域 R （足够小，体积为 V ）内的 k 个点估计 x_i 的概率密度 $p(x)$ 有贡献，则 R 中落入 k 个样本的概率为：

$$P_R = k/N = \int_R p(x) dx = \hat{p}(x)V$$

- 估计得到的概率密度 $\hat{p}(x_i)$ 为： $\hat{p}(x_i) = k/NV$
- 当满足以下条件时，概率估计 $\hat{p}(x_i)$ 收敛于 $p(x_i)$ ：

- 贡献点的区域大小越小越好 $\lim_{N \rightarrow \infty} V_N = 0$
- 贡献点越多越好 $\lim_{N \rightarrow \infty} k_N = \infty$
- 贡献点与总样本数比例越小越好 $\lim_{N \rightarrow \infty} \frac{k_N}{N} = 0$

非参数估计方法

● Parzen法

- 使区域体积序列 V_N 以 N 的某个函数的关系不断缩小
- 同时限制 k_N 和 k_N/N

有限的 N , V 选择很敏感



● k近邻法

- 使落入区域样本数 k_N 为 N 的某个函数
- 选择不同的 V_N 使区域包含 x 的 k_N 个近邻

动态变化 V 的取值

非参数估计方法——Parzen窗法

- Parzen窗法使用窗函数对贡献点进行选择

- 窗函数

- 形式: $k(x, x_i)$, 反映 x_i 对 $p(x)$ 的贡献同时进行区域选择

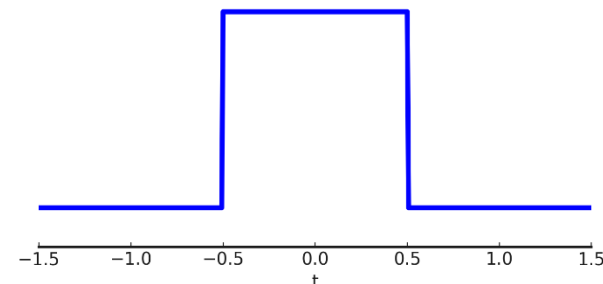
- 条件: $k(x, x_i) \geq 0, \int k(x, x_i) dx = 1$

- 窗函数选择

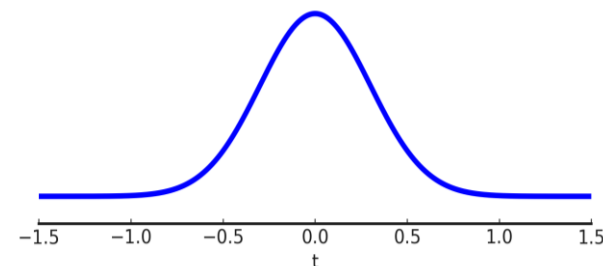
- 方窗函数、正态窗函数、指数窗函数等

- 窗宽选择

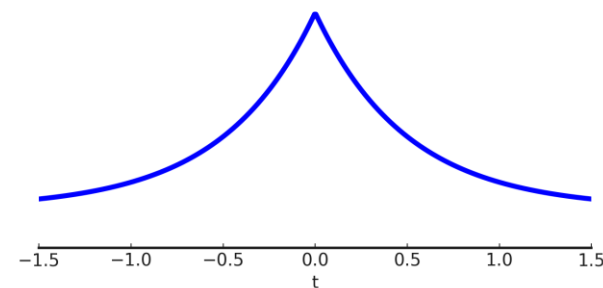
- 原则: 样本数多则选小窗宽; 样本数少则选大窗宽



方窗函数



正态窗函数

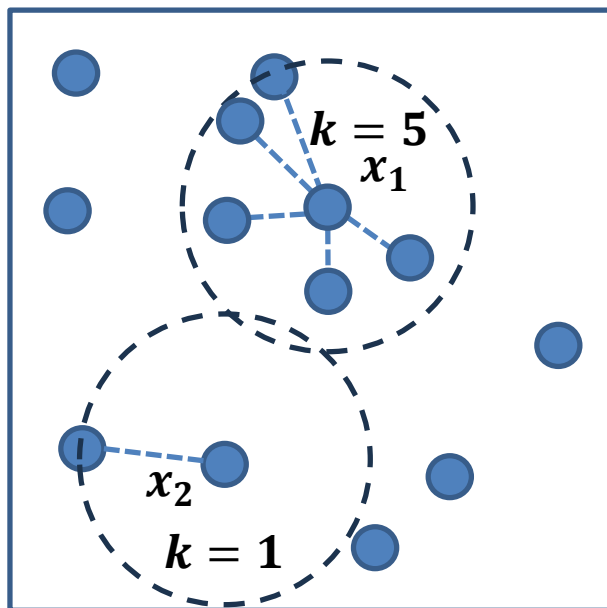


指数窗函数

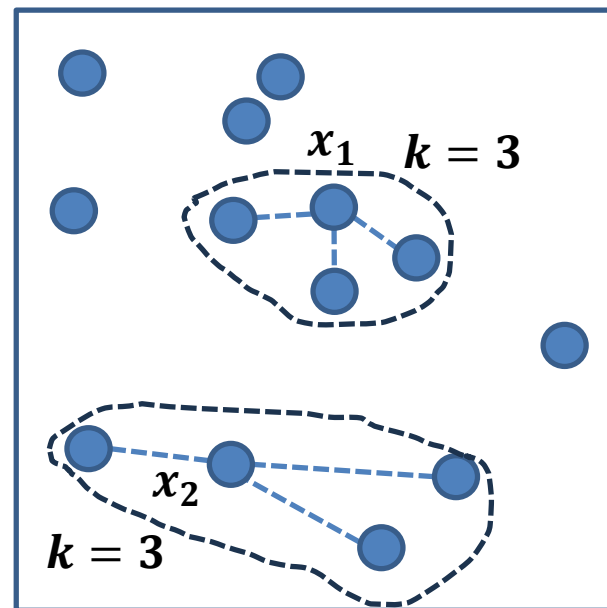
非参数估计方法——k近邻法

- k近邻法使用k近邻算法对贡献点进行选择
 - 选择样本 x_i 一定范围内确定个数的 k 个样本后根据 k/NV 计算概率密度
 - k近邻法更适用于样本分布不均匀的数据

Parzen窗法



k近邻法



k近邻算法

● k近邻算法流程：

- 选择一个正整数 k ，表示需要考虑的邻近样本的数量
- 对于待预测的样本，计算其与训练集中所有样本之间的距离
 - 常用的距离度量包括欧氏距离、曼哈顿距离、余弦相似度等
- 选择 k 个最近邻居：根据计算得到的距离，选择距离待预测样本最近的 k 个邻居

● k 值选择：

- k 值决定了决策的局部性， k 值越大，模型越平滑，越小则越敏感

小结

- Parzen窗法与k近邻法均使用广泛
- 在样本分布不均匀时k近邻法比Parzen窗法表现更好
- 在高维空间中k近邻法比Parzen窗法更易应用，且可通过技术手段缓解维度问题
- k近邻法和Parzen窗法在边界附近都可能会遇到估计偏差

第3章：线性模型

Chapter 3: Linear Models

3.1 什么是线性回归

- “回归”的起源
- 线性回归的概念

“回归”起源

- 研究父母身高与子女身高之间的关系

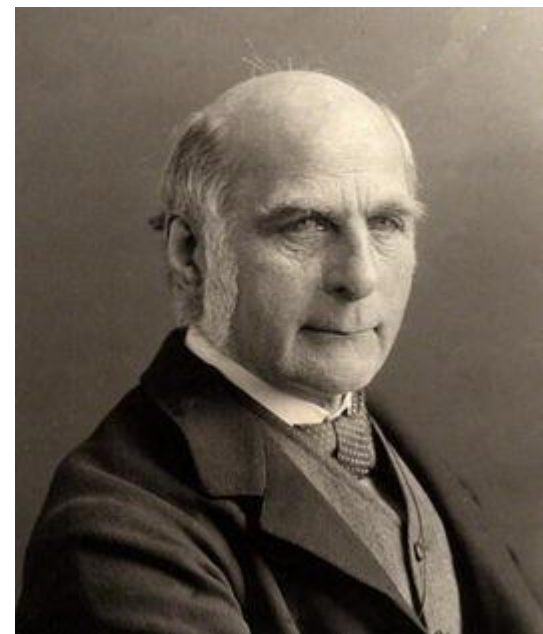
【英国生物学家Francis Galton】

1822-1911

若父母身高高于平均身高，则其子女身高倾向于倒退生长，即会比其父母身高矮一些而更接近于大众平均身高。

若父母身高小于平均身高，则其子女身高倾向于向上生长，即会更接近于大众平均身高。

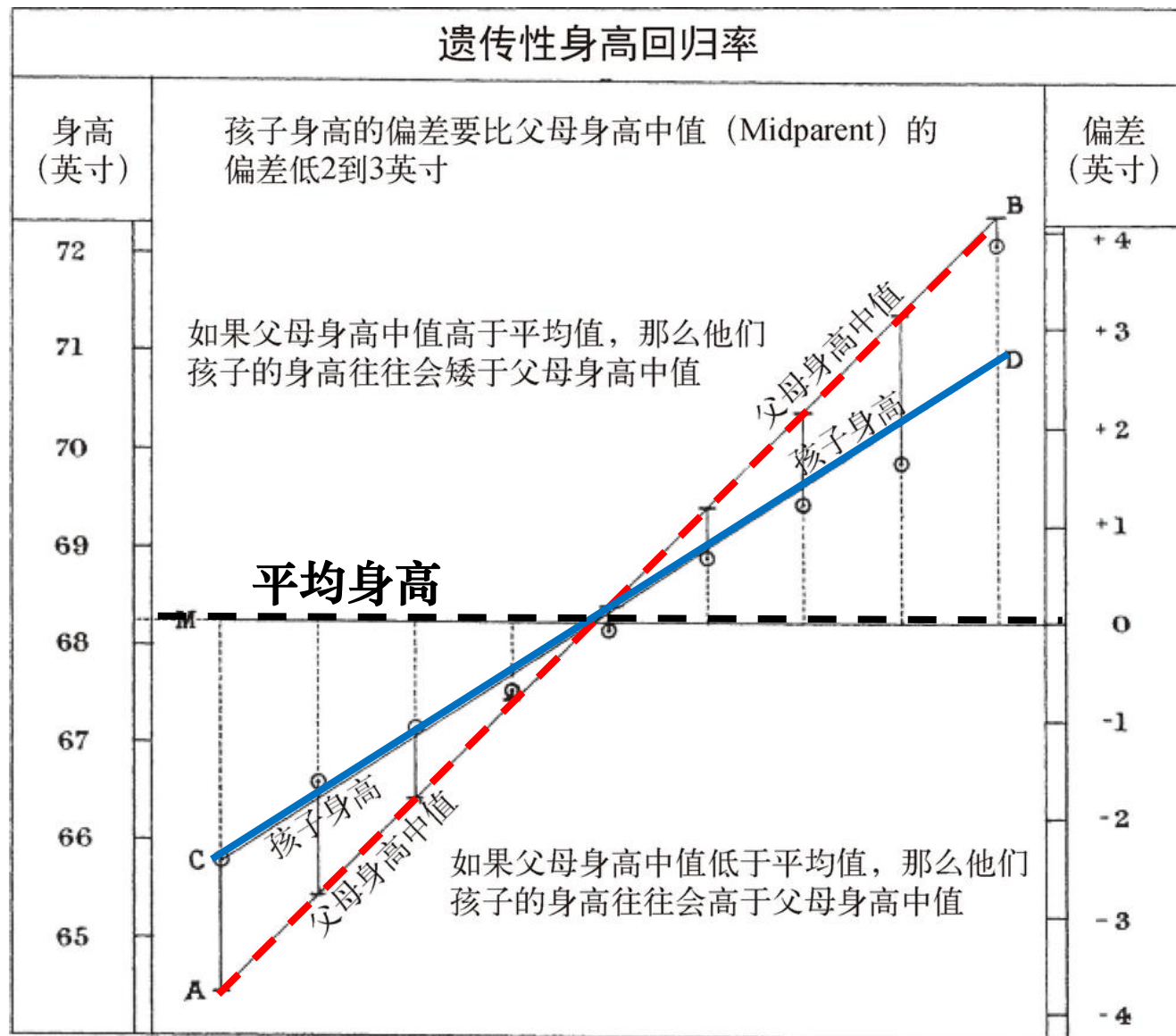
此现象被Galton称之为回归现象，即Regression



“回归”起源

● 父母身高与子女身高之间关系

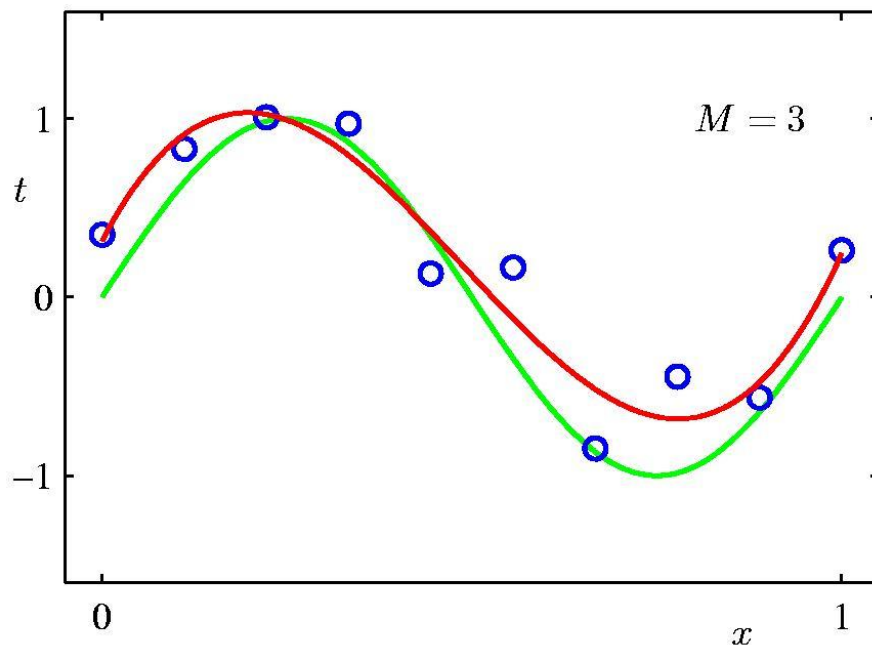
后代的身高倾向于
“回归”到**平均值**



回归的概念

● 回归

- 回归属于**监督学习**，其目标是将输入向量分配至由一个或多个**连续变量**组成的输出*
- 数学表示：给定输入数据 x 和一个连续型输出 y ，目标是找到一个函数 $f: x \rightarrow y$ ，使 $f(x)$ 与 y 间的差距尽可能小



回归举例-房价估计

- 房屋售价与其多种因素有关，现统计了一些房屋的售价和对应的因素数据（包括面积、卧室数量、层数、房龄等）

面积(m ²)	卧室数量(间)	层数	房龄(年)	价格(万元)
210.4	5	1	45	460
141.6	3	2	40	232
153.4	3	2	30	315
85.2	2	1	36	178
...

- 如何利用已知数据构建回归模型，实现对房屋价格的估计？

面积(m ²)	卧室数量(间)	层数	房龄(年)	价格(万元)
150	3	2	30	?

测试样本

回归举例-房价估计

训练样本

面积(m ²) x_1	卧室数量(间) x_2	层数 x_3	房龄(年) x_4	价格(万元) 连续变量 y	} N
210.4	5	1	45	460	
141.6	3	2	40	232	
153.4	3	2	30	315	
85.2	2	1	36	178	
...	

监督学习

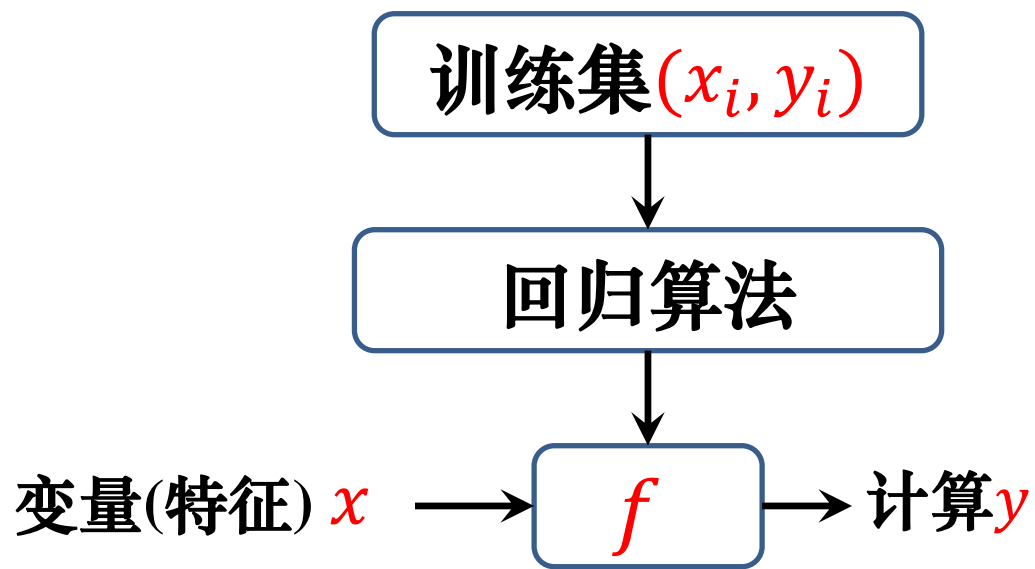
N : 训练样本个数

x : 输入变量/“特征”

y : 输出变量/目标变量

线性回归的概念

● 问题建模



如何表示 f ?

● 线性回归

➤ 假设函数 f 为输入 x 的线性函数

$$f(x) = w_0 + w_1x_1 + \cdots + w_mx_m$$

$$= w_0 + \sum_{i=1}^m w_ix_i$$

$$= \sum_{i=0}^m w_ix_i$$

增加一维 $x_0 = 1$
表示截距项，
转为齐次形式

➤ 写成向量形式

$$f(x) = \mathbf{w}^T \mathbf{x}$$

线性回归的概念

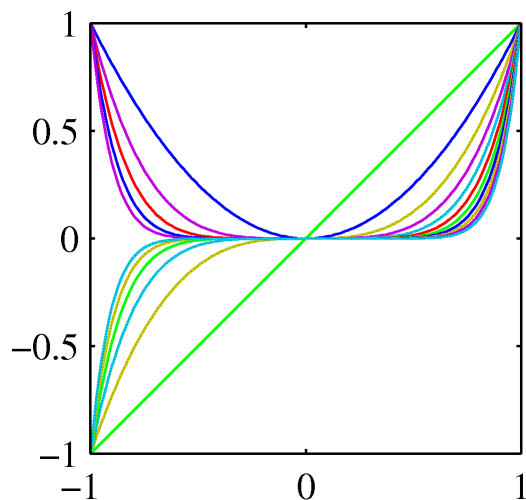
$$y(x, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(x) = \mathbf{w}^T \boldsymbol{\phi}(x)$$

基函数

最简单的情况下: $\phi_j(x) = x_j$

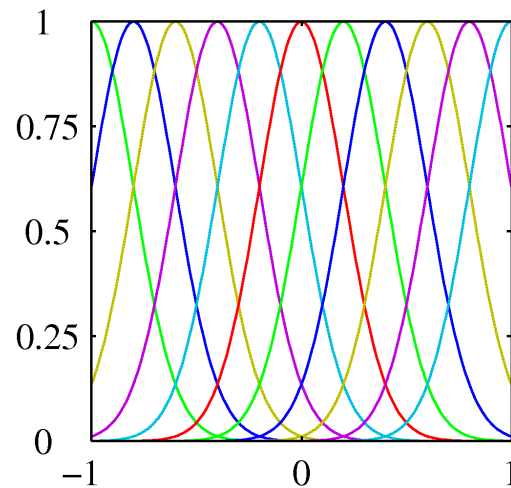
● 多项式基函数

$$\phi_j(x) = x^j$$



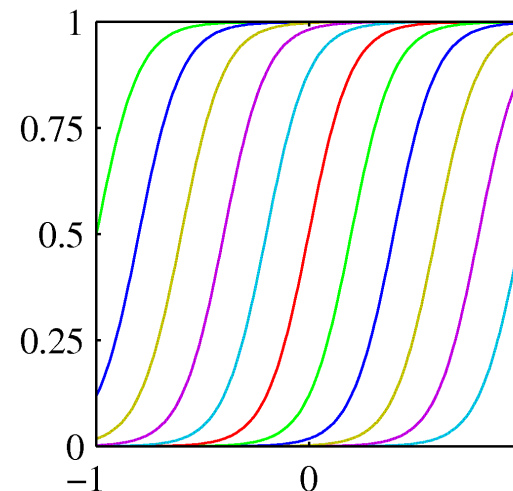
● 高斯基函数

$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$



● Sigmoid基函数

$$\phi_j(x) = \sigma\left\{-\frac{x - \mu_j}{s}\right\} \quad \sigma(a) = \frac{1}{1 + \exp(-a)}$$



3.2 线性回归的求解

- 标准方程组
- 梯度下降法
- 最大似然估计
- 贝叶斯估计

问题回顾

训练样本

面积(m ²) x_1	卧室数量(间) x_2	层数 x_3	房龄(年) x_4	价格(万元) 连续变量 y	} N
210.4	5	1	45	460	
141.6	3	2	40	232	
153.4	3	2	30	315	
85.2	2	1	36	178	
...	

N : 训练样本个数

x : 输入变量/“特征”

y : 输出变量/目标变量

● 线性回归问题建模

➤ 假设函数 f 为输入 x 的线性函数

$$f(x) = w_0 + w_1x_1 + \cdots + w_mx_m$$

$$= w_0 + \sum_{i=1}^m w_ix_i$$

$$= \sum_{i=0}^m w_ix_i$$

➤ 写成向量形式

$$f(x) = \mathbf{w}^T \mathbf{x}$$

问题分析

- 问题本质：确定模型中的参数 w
- 基本思想：基于训练集最小化预测值 $\hat{y} = f(x)$ 与真实输出值 y 的差异
- 定义目标函数 (又叫代价函数)
 - 例如：平均绝对误差、均方误差和均方根误差
- 以基于均方误差的目标函数为例 $\hat{w} = \arg \min_{\mathbf{w}} J(\mathbf{w}) = \arg \min_{\mathbf{w}} \left[\frac{1}{2} \sum_{i=1}^N (f(x_i) - y_i)^2 \right]$
- 求解方法：标准方程组法、梯度下降法

标准方程组法(Normal Equations)

- 目标函数 $J(\mathbf{w}) = \sum_{i=1}^N (f(x_i) - y_i)^2 = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2$

- 写成矩阵形式 $J(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$

- 其中

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix} = \begin{bmatrix} x_{10} & x_{11} & \cdots & x_{1m} \\ x_{20} & x_{21} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N0} & x_{N1} & \cdots & x_{Nm} \end{bmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}^T$$

标准方程组法(Normal Equations)

- 对目标函数直接求导，并令其导数等于0，求得极值

- 求导

$$\begin{aligned}\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}) &= \frac{\partial}{\partial \mathbf{w}} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= 2\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0\end{aligned}$$

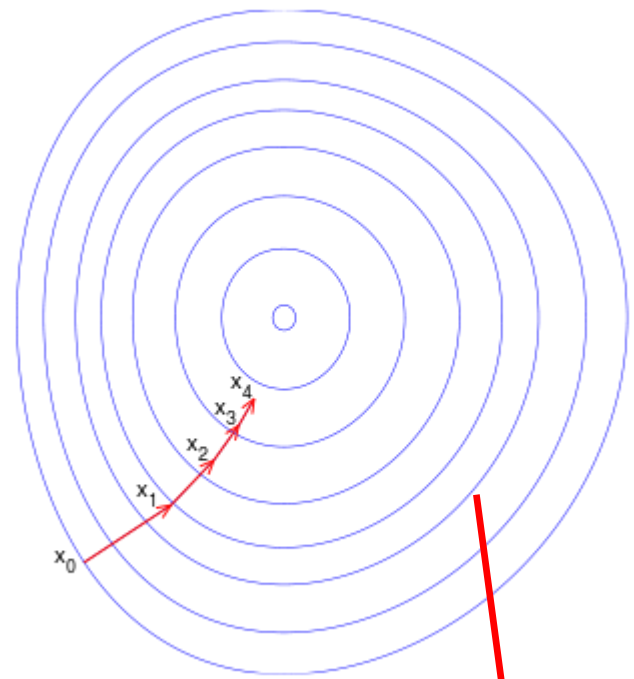
$$\Rightarrow \mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

- 得到模型的参数

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

梯度下降法(Gradient Descent)

- 梯度下降是一个最优化算法，是求解无约束优化问题最基础的方法之一
- 基本思想
 - 更新参数让目标函数沿**负梯度方向**下降
- 一般流程
 - 首先对参数 w 赋值
 - 更新 w 的值，使得 $J(w)$ 按**梯度下降**的方向减少直到收敛
 - 越接近目标值，搜索步长越小，前进越慢



等值线

梯度下降法(Gradient Descent)

- 损失函数: $J(\mathbf{w}) = \sum_{i=1}^N (f(x_i) - y_i)^2$

- 目标: $\min_{\mathbf{w}} J(\mathbf{w})$

- 梯度下降法的步骤

- 给定初始值 w^0 ，这个值可以是随机生成的，也可以是一个全零的向量

- 更新 \mathbf{w} 使得 $J(\mathbf{w})$ 越来越小， α 为学习率或更新步长

$$w_j^t = w_j^{t-1} - \alpha \frac{\partial}{\partial w_j} J(\mathbf{w}) = w_j^{t-1} - \alpha \sum_{i=1}^N (f(x_i) - y_i) \cdot x_{i,j}$$

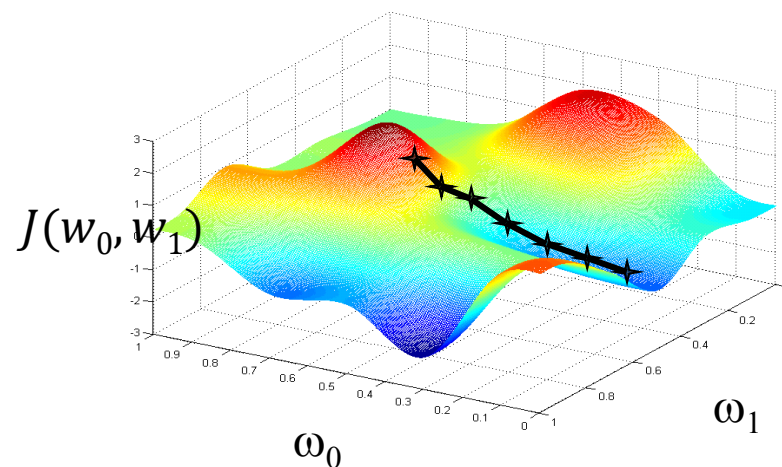
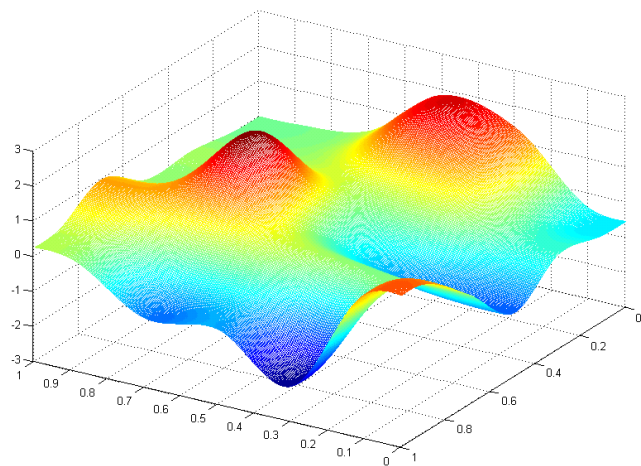
- 同时更新函数值的各维 $f(x_i) = [\mathbf{w}^t]^T \mathbf{x}_i$

批处理梯度下降(Batch Gradient Descent)

- 在梯度下降方法中，每次更新都利用所有数据

$$w_j^t = w_j^{t-1} - \alpha \sum_{i=1}^N (f(x_i) - y_i) \cdot x_j \rightarrow \text{批处理}$$

- 对于凸函数，可以达到全局最优
- 在大样本条件下，批处理梯度下降的迭代速度很慢



随机梯度下降(Stochastic Gradient Descent)

- 基本思想：每次只用一个样本 (x_r, y_r)

- 随机梯度下降的特点

- 更新速度较快 $w_j^t = w_j^{t-1} - \alpha(f(x_r) - y_r) \cdot x_r$

- 更适用于大样本数据情况

- 对于更复杂的优化目标，可以跳出局部最优解

- 变化形式

- 在线学习(Online Learning)，每次“看”一个样本，对所有样本循环使用多次
(一次循环称为一个Epoch)

- 每次可以看一些样本(Mini-Batch)

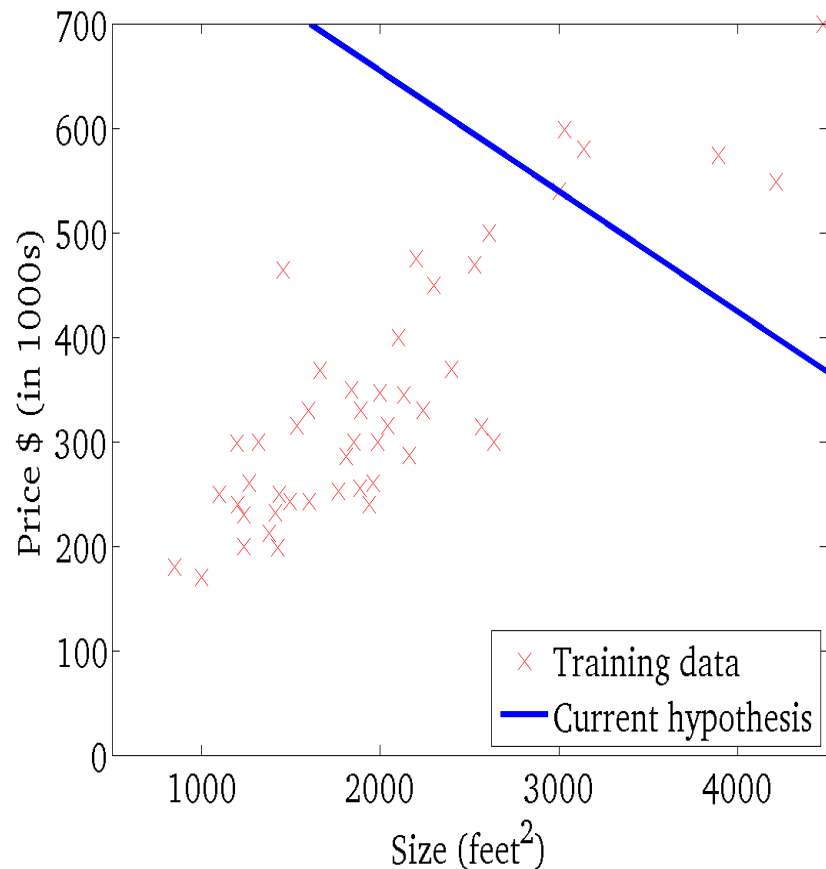
随机梯度下降

- 问题描述：如何利用已知数据构建回归模型，实现对房屋价格的估计（简化房屋价格的影响因素，只保留房屋面积）
- 问题建模：构建 $f_{\mathbf{w}}(x) = w_1x + w_0$ 来根据房屋面积 x 估计房屋价格 y
- 代价函数： $J(\mathbf{w}) = \sum_{i=1}^N (f(x_i) - y_i)^2$
- 问题求解：使用随机梯度下降法求解

随机梯度下降

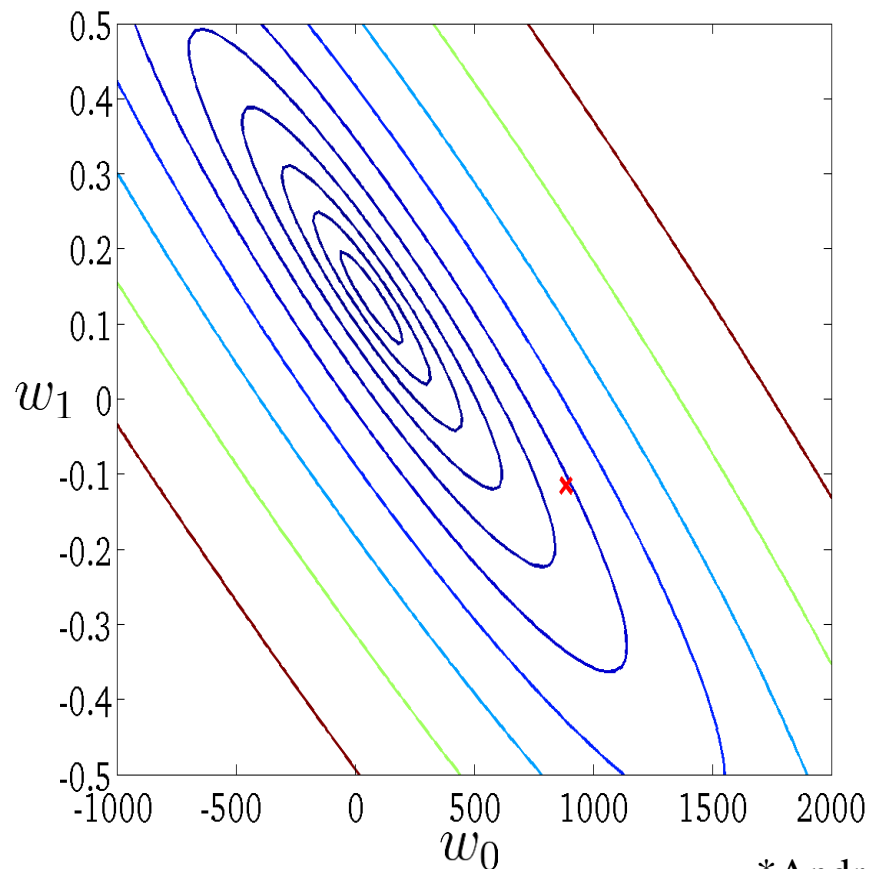
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

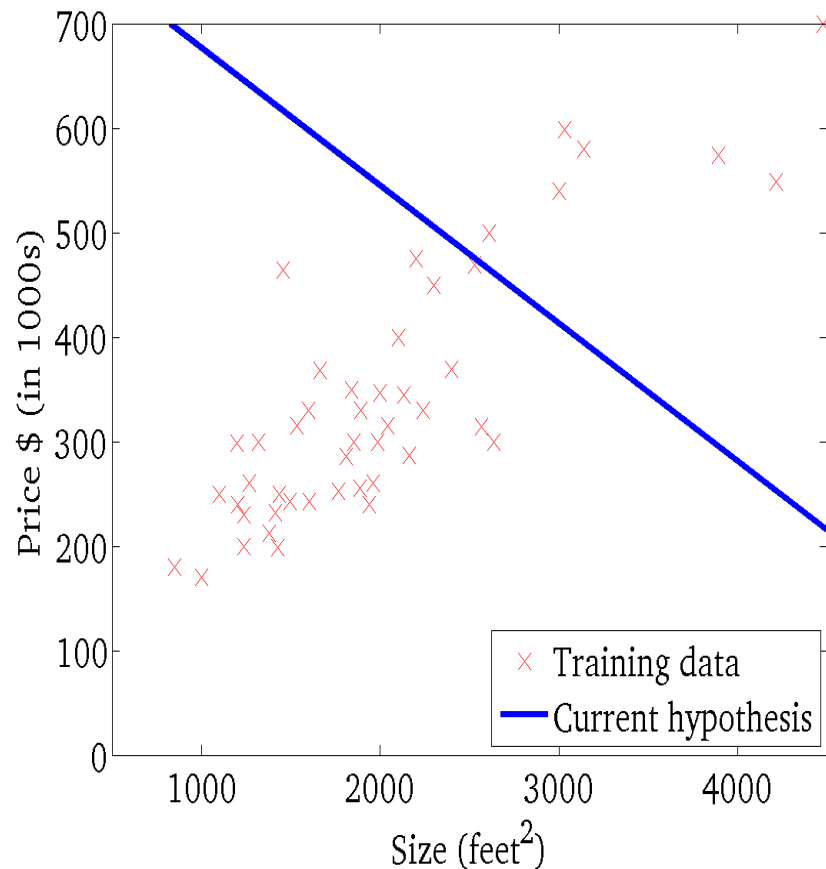
(关于参数 w_1, w_0 的函数)



随机梯度下降

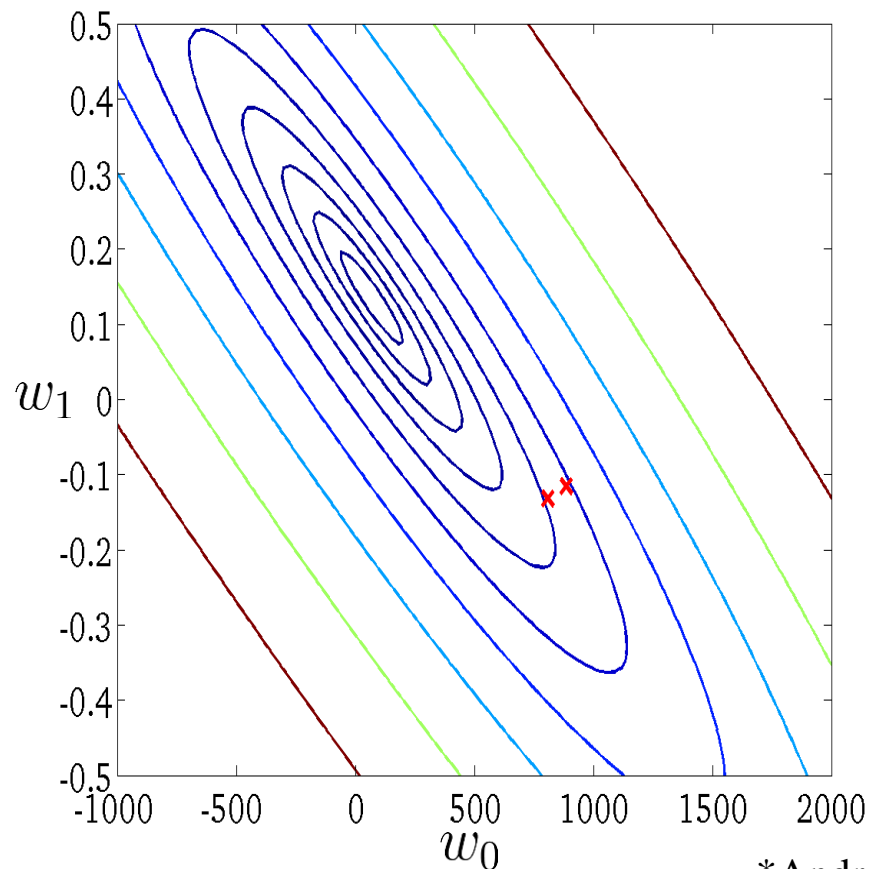
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

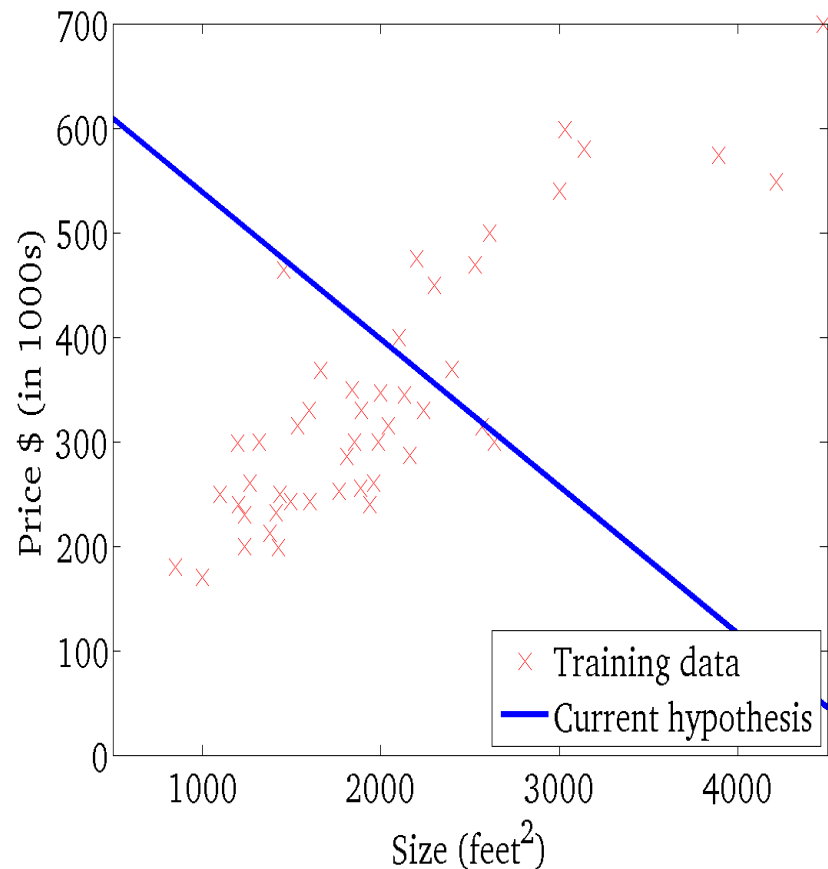
(关于参数 w_1, w_0 的函数)



随机梯度下降

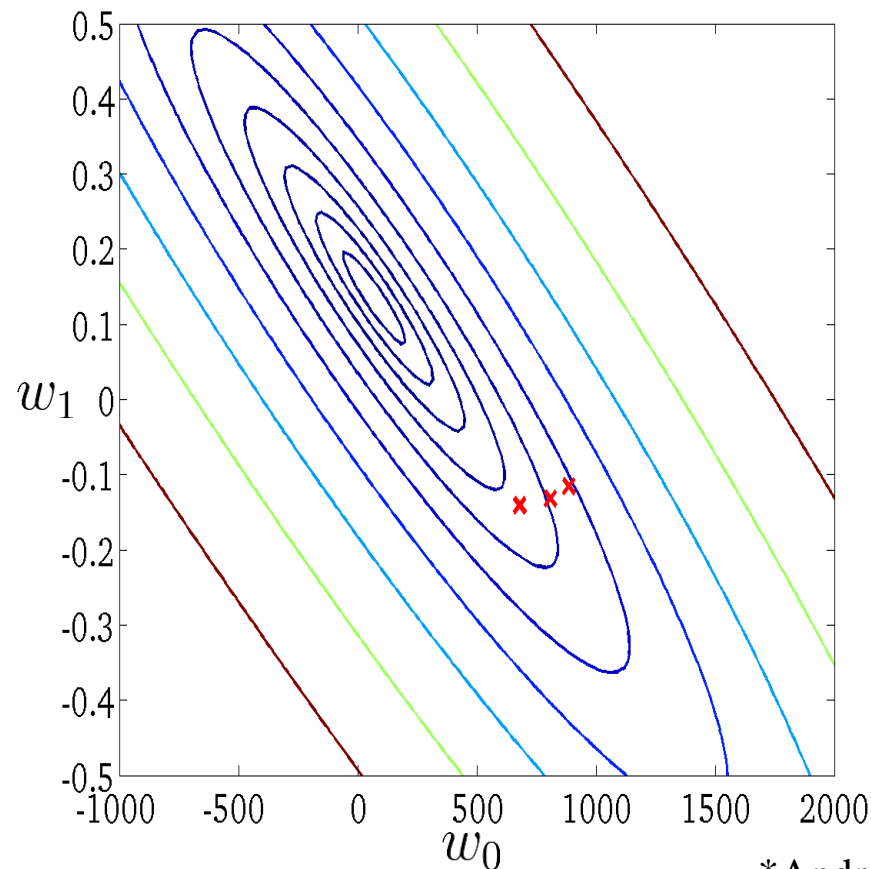
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

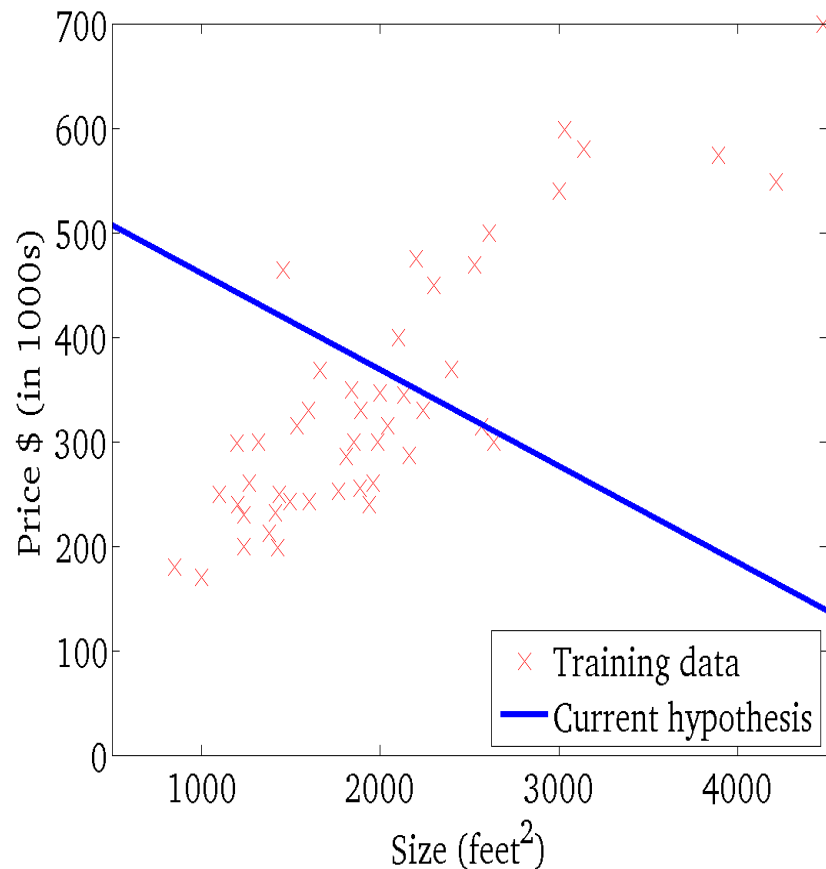
(关于参数 w_1, w_0 的函数)



随机梯度下降

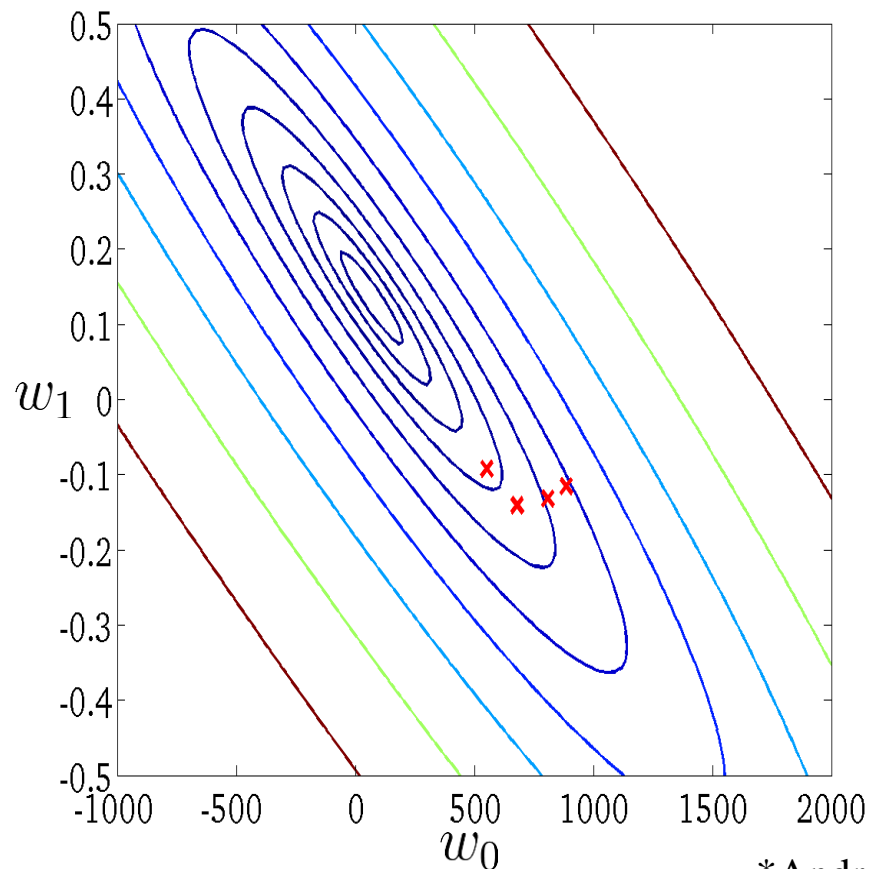
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

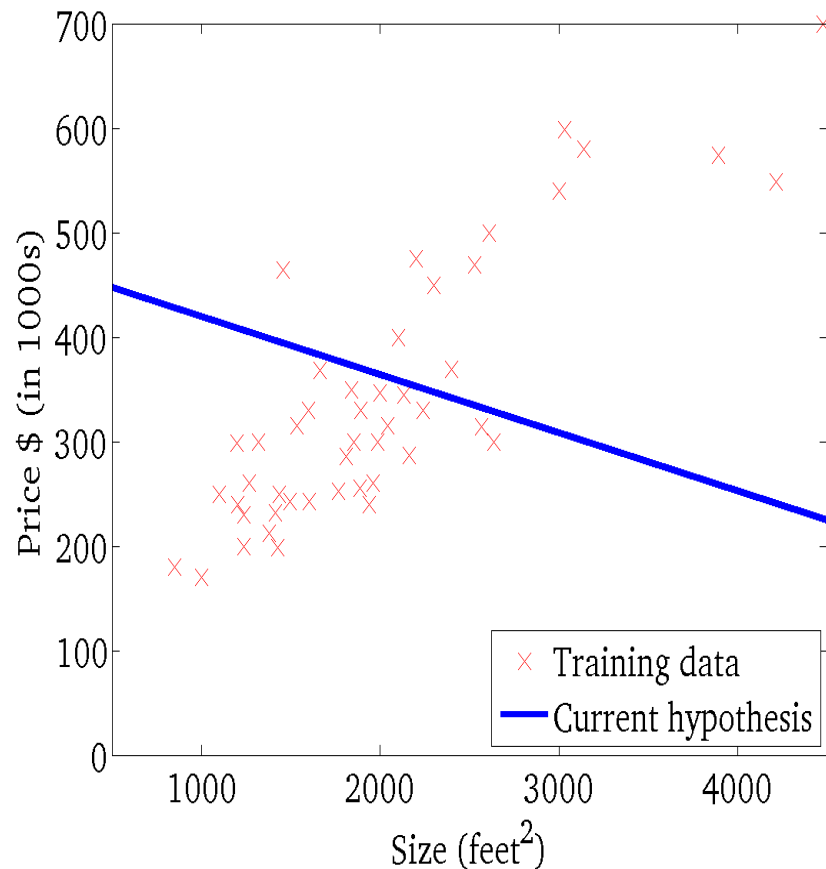
(关于参数 w_1, w_0 的函数)



随机梯度下降

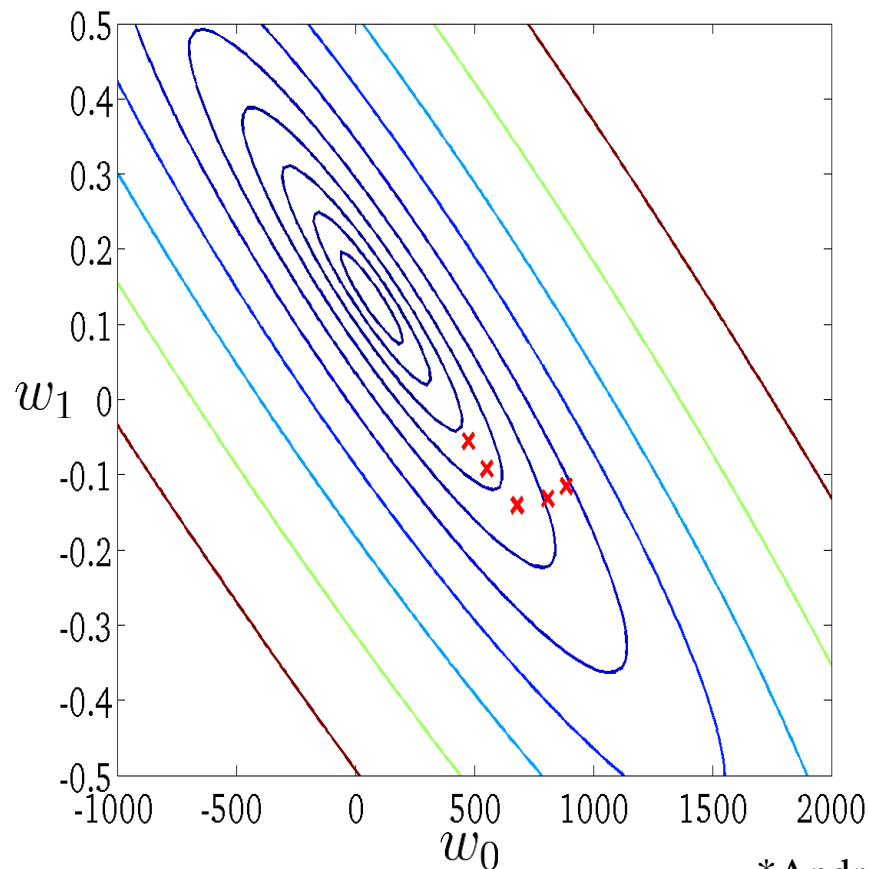
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

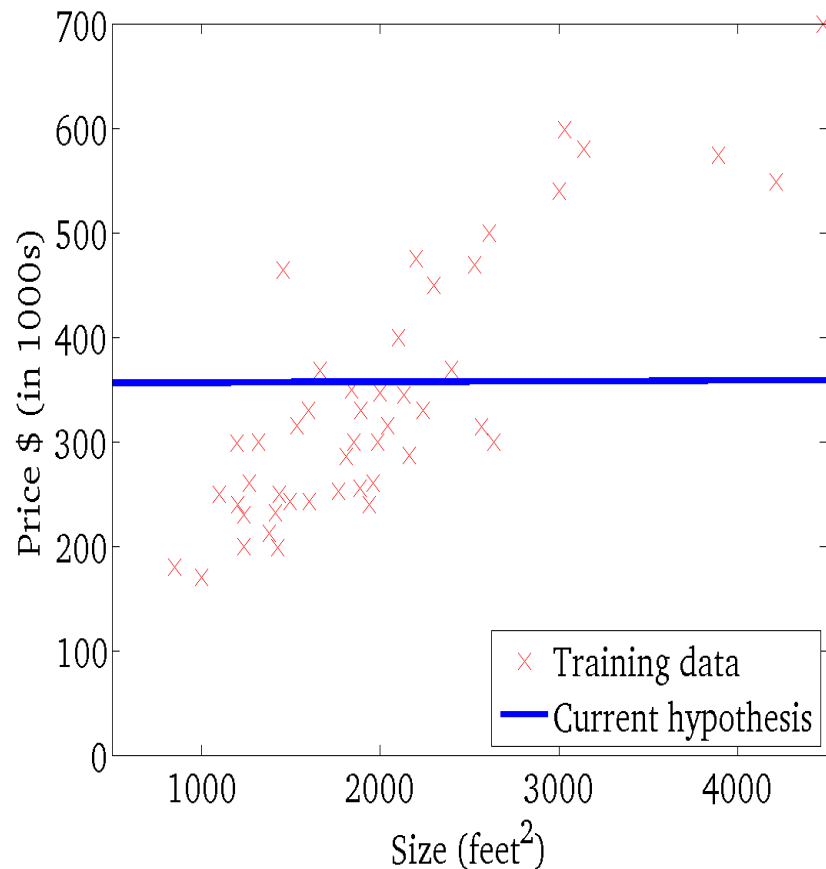
(关于参数 w_1, w_0 的函数)



随机梯度下降

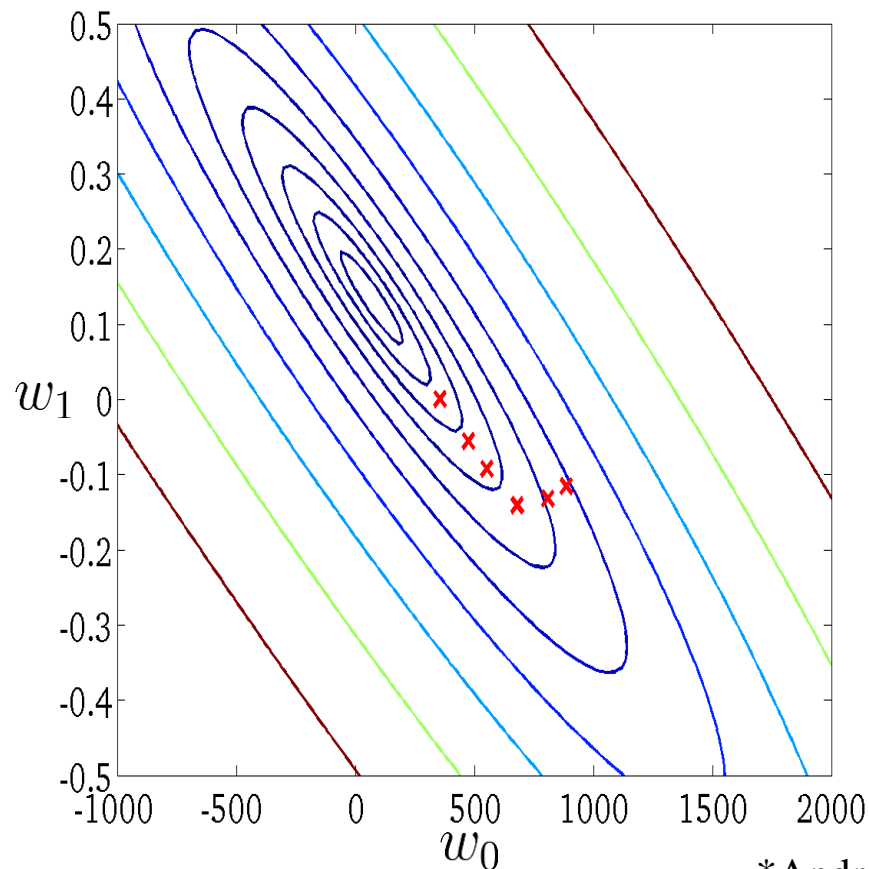
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

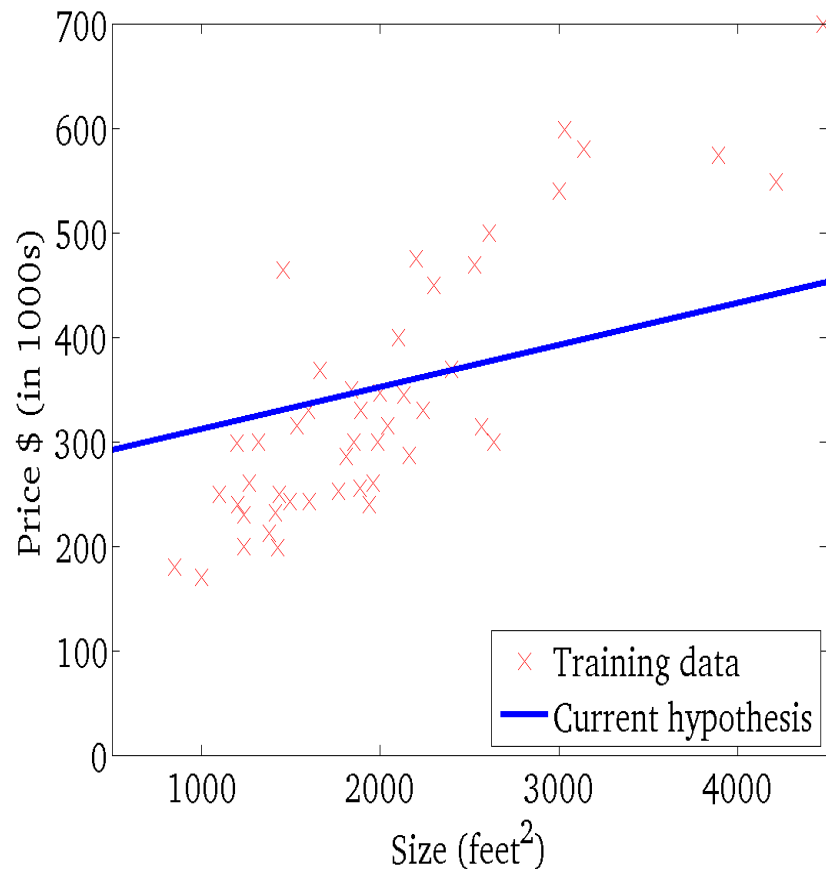
(关于参数 w_1, w_0 的函数)



随机梯度下降

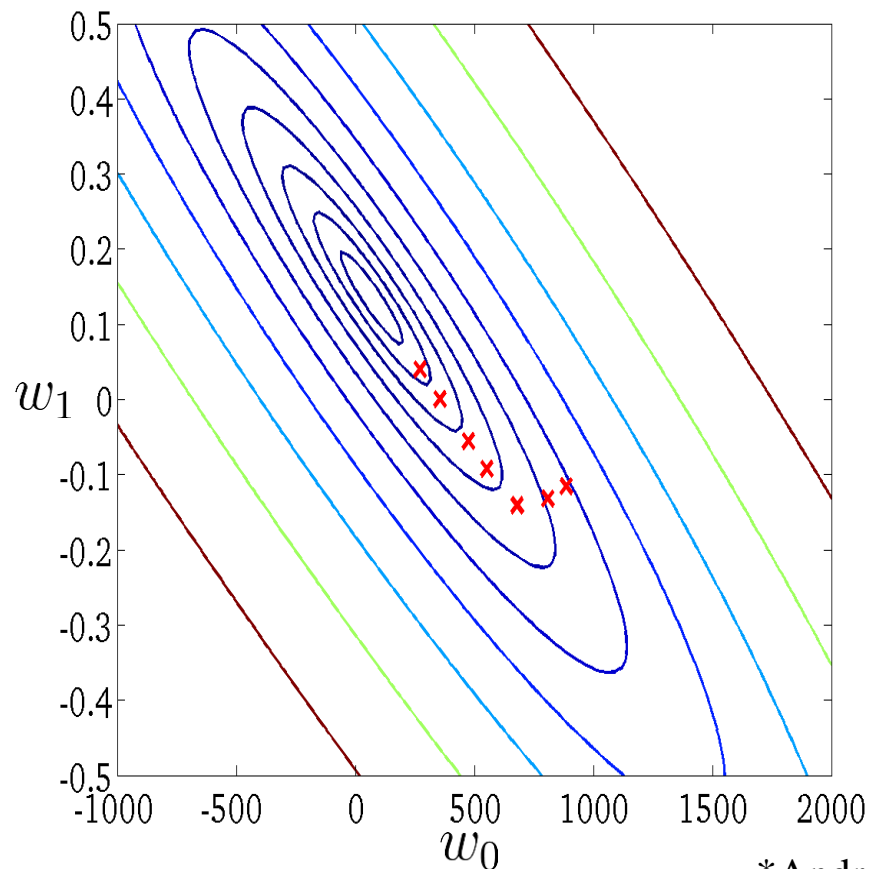
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

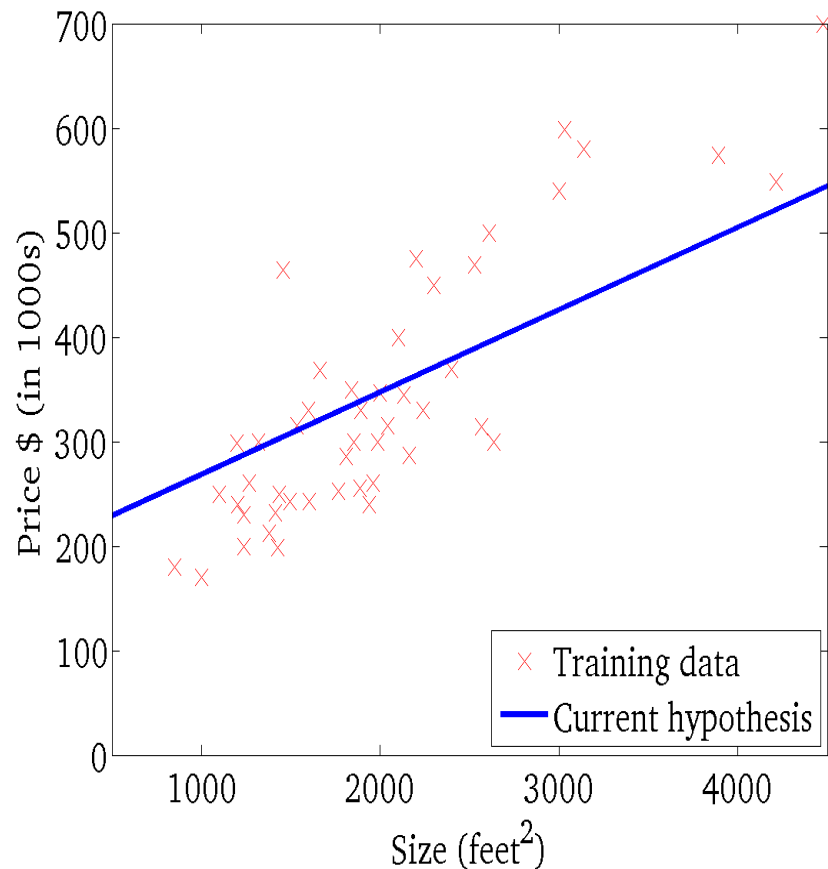
(关于参数 w_1, w_0 的函数)



随机梯度下降

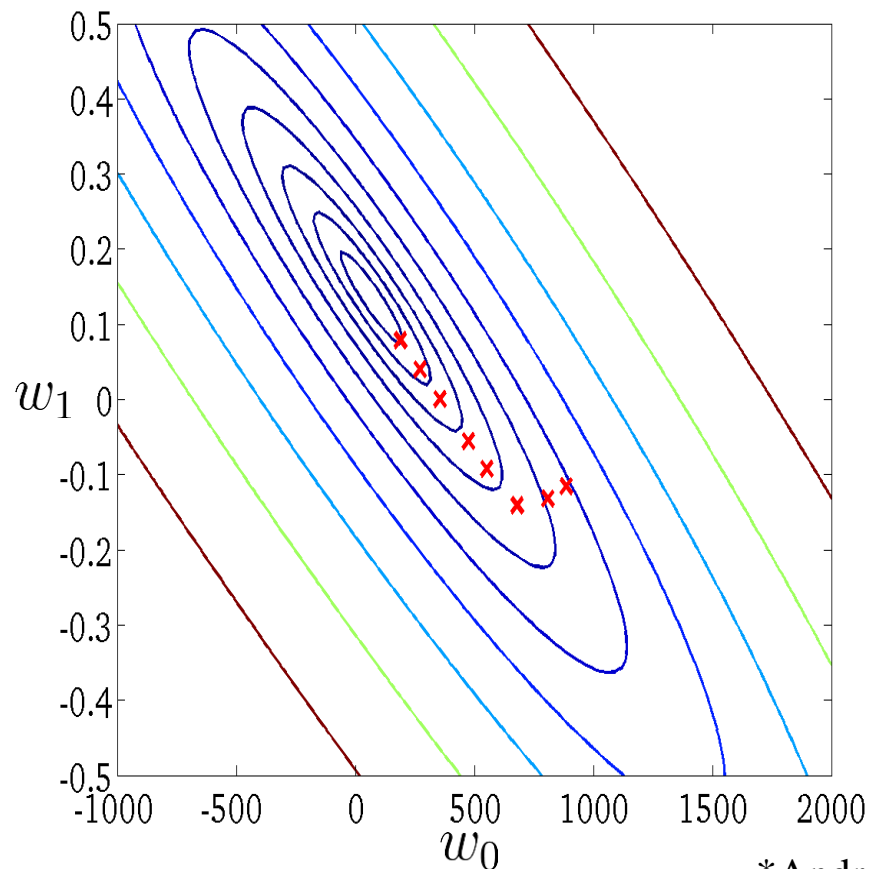
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

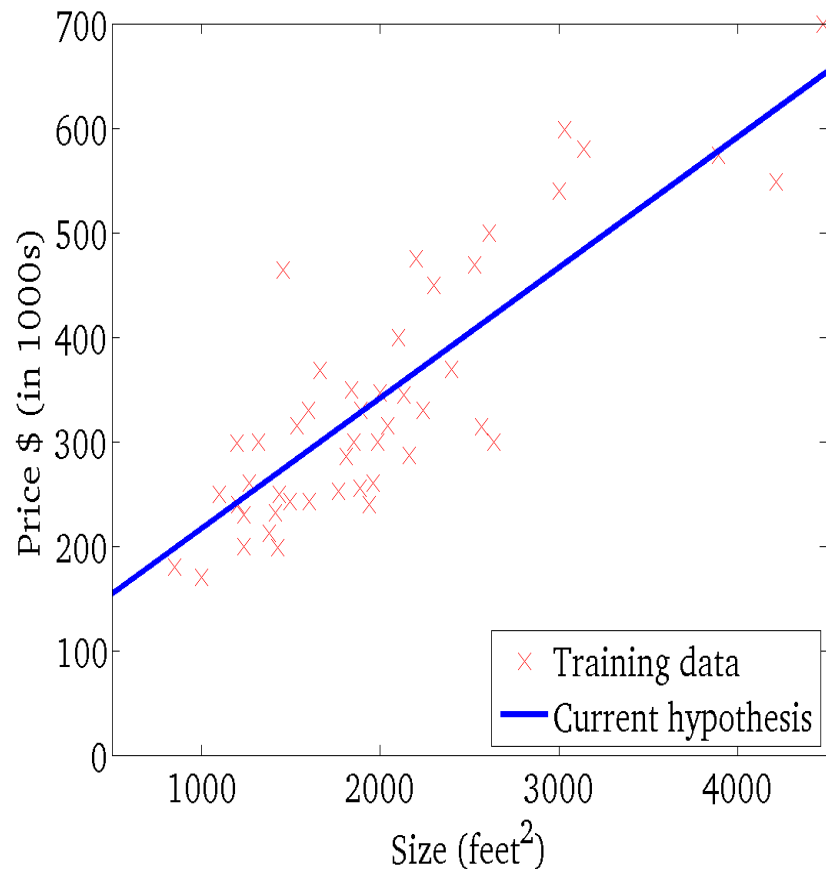
(关于参数 w_1, w_0 的函数)



随机梯度下降

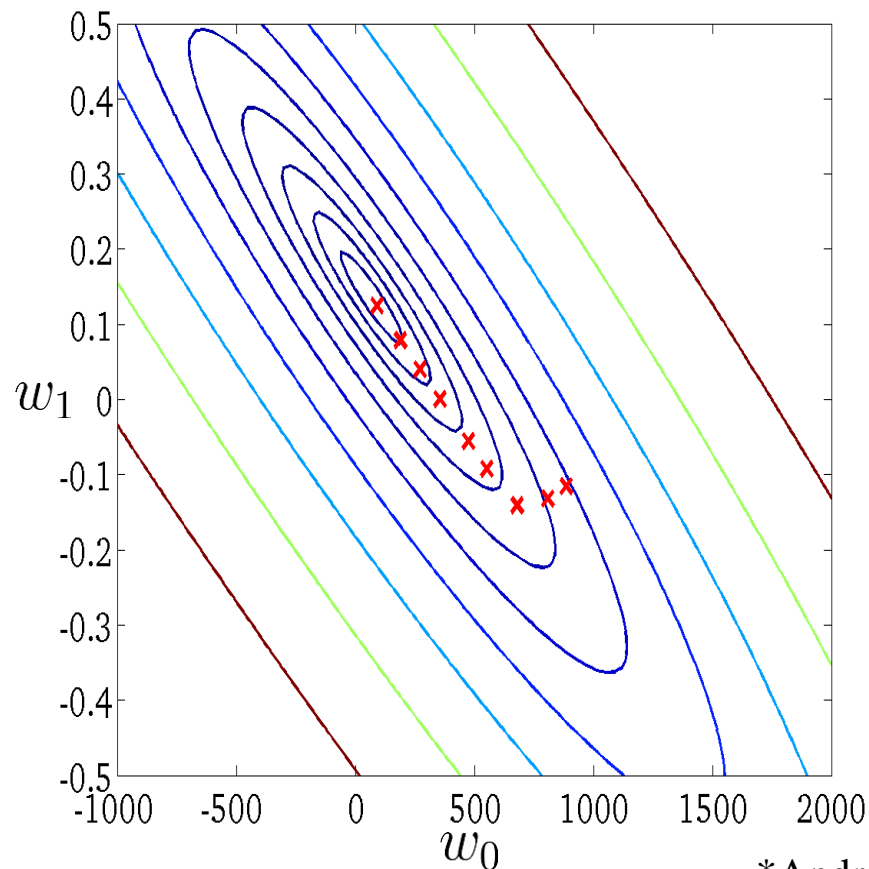
$$f_{\mathbf{w}}(x)$$

(对于固定的 w_1, w_0 , 这是关于 x 的函数)



$$J(\mathbf{w})$$

(关于参数 w_1, w_0 的函数)



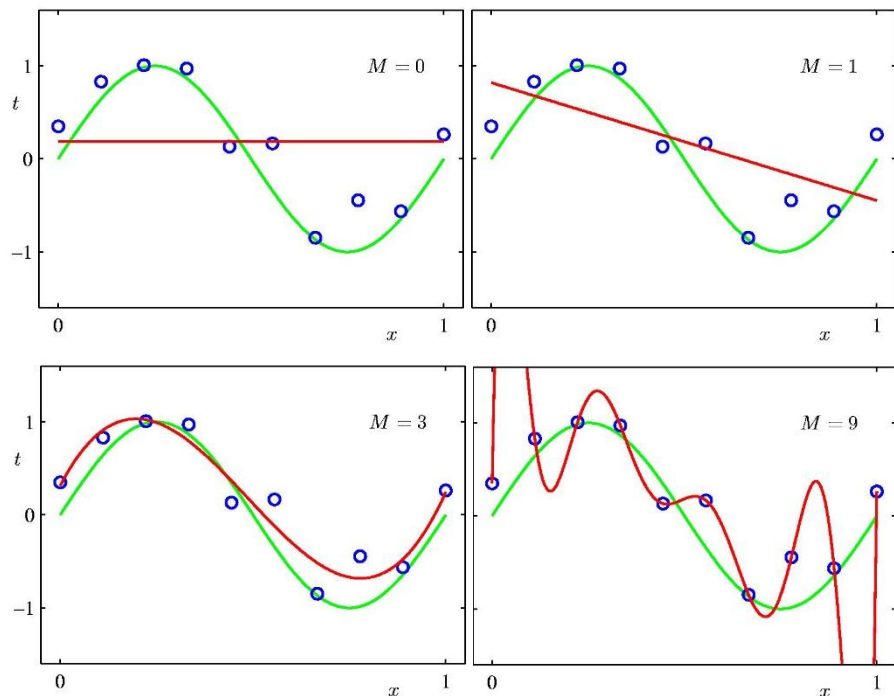
标准方程组 v.s. 梯度下降

	标准方程组	梯度下降
超参数	不需要超参 α	需要选择 α
迭代次数	不需要迭代多次	需要迭代多次
归一化	无需数据归一化	需要数据归一化
适用情况	样本量大时不适用 需要计算 $(\mathbf{X}^T \mathbf{X})^{-1}$	样本量非常大时也适用

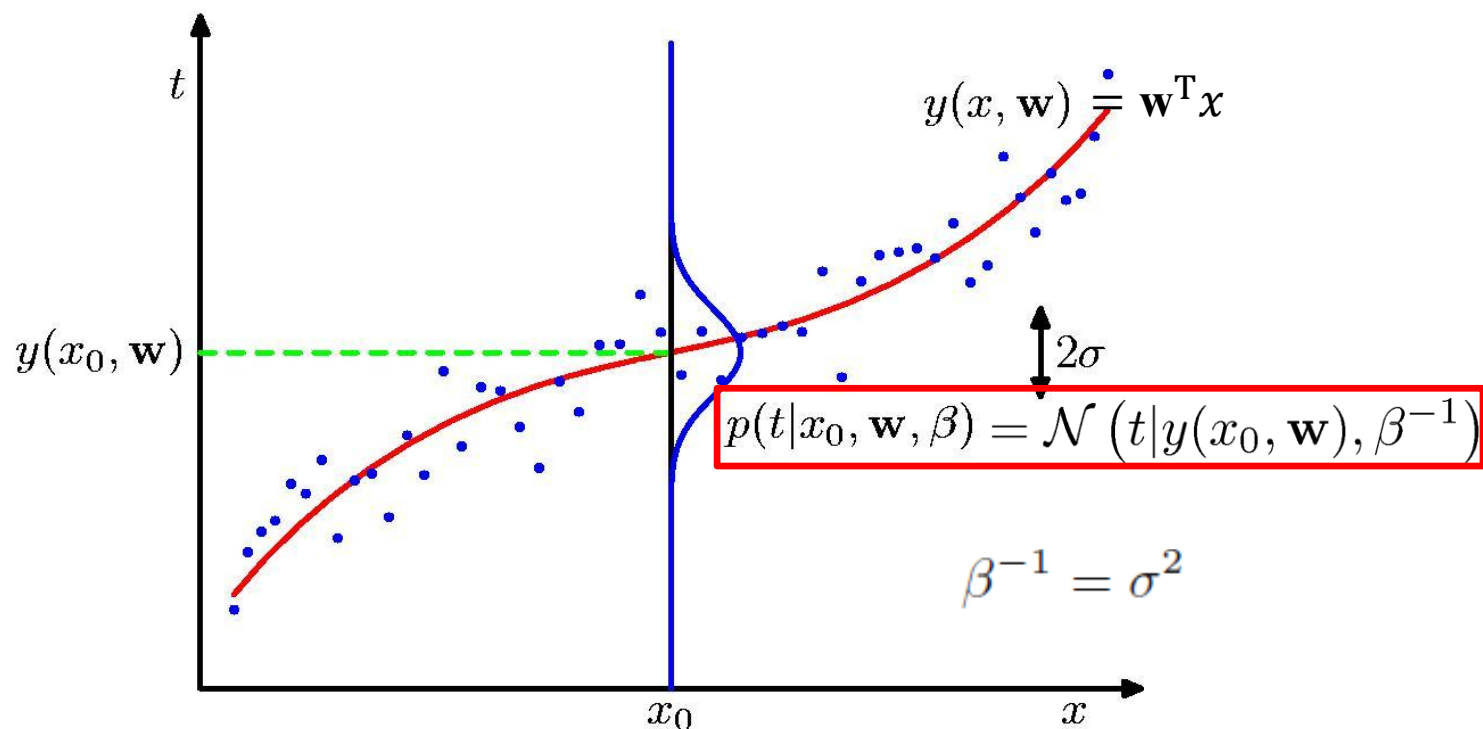
样本量较小时选用标准方程组求解
样本量较大时选用梯度下降法求解

从概率视角看回归

- 第一章的曲线拟合中，目标变量值出现不确定性，可用概率分布表示
- 假定目标值 t 服从以 $y(x, \mathbf{w})$ 为均值， σ^2 为方差的高斯分布



第一章中示例：曲线拟合



最大似然估计与贝叶斯估计

● 最大似然估计

- 把待估计的参数 \mathbf{w} 看做是**确定的量**，只是其取值未知。最佳估计就是使得观测到的样本的概率最大的那个值。

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$$

● 贝叶斯估计

- 把待估计的参数 \mathbf{w} 看做是符合某种先验概率分布的**随机变量**。对样本进行观测的过程，就是把先验概率密度转化为后验概率密度，从而利用样本信息修正了对参数的初始估计值。

$$p(t|\mathbf{x}, \mathbf{t}, \alpha, \beta) = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}(t | m_N^T \mathbf{x}, \sigma_N^2(\mathbf{x}))$$

最大似然估计(Maximum Likelihood Estimation)

已知训练数据 \mathbf{x} 和 \mathbf{t} ，现有新输入变量 x ，需要预测目标变量 t

似然函数表示为 $p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$

对数似然函数为 $\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\boxed{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$ 均方误差

通过最大化对数似然函数解得 $\boxed{\mathbf{w}_{ML}}$ 等价于以均方误差为损失函数的线性回归的解

此时高斯分布的参数 $\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{ML}) - t_n\}^2$

得到概率分布 $p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t | y(x, \mathbf{w}_{ML}), \beta_{ML}^{-1})$

MAP: 离贝叶斯更近一步

假定 \mathbf{w} 的先验分布为

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

其中 α 是分布精度， M 是 \mathbf{w} 的阶数

根据贝叶斯定理 $p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$$

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w} | 0, \alpha^{-1}\mathbf{I})$$

取对数可知，最大化 $\ln p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta)$ 即最大化

$$-\ln \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

MAP: 离贝叶斯更近一步

最大化 $p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta)$ 即最小化

$$\beta \tilde{E}(\mathbf{w}) = \underbrace{\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\text{平方误差}} + \underbrace{\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}}_{\text{正则化项}}$$

正则平方误差函数

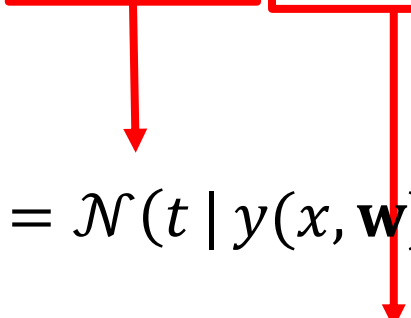
令 $\lambda = \frac{\alpha}{\beta}$ 为正则化系数，令 $\nabla_{\mathbf{w}} \tilde{E}(\mathbf{w}) = 0$ 得解

$$\mathbf{w}_{MAP} = (\mathbf{x}^T \mathbf{x} + \lambda \mathbf{I})^{-1} \mathbf{x}^T \mathbf{t}$$

等价于以正则平方误差为损失函数的线性回归的解

贝叶斯估计

预测分布简单表示为

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}$$


其中

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\alpha)$$

预测分布涉及两个高斯分布的卷积，可采用以下形式

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), s^2(x))$$

贝叶斯估计

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}(t | m(x), s^2(x))$$

代入

$$y(x, \mathbf{w}) = \sum_{j=0}^{M-1} w_j x_j = \mathbf{w}^T \mathbf{x}$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w}|\alpha)$$

使用边际高斯和条件高斯的转换可知

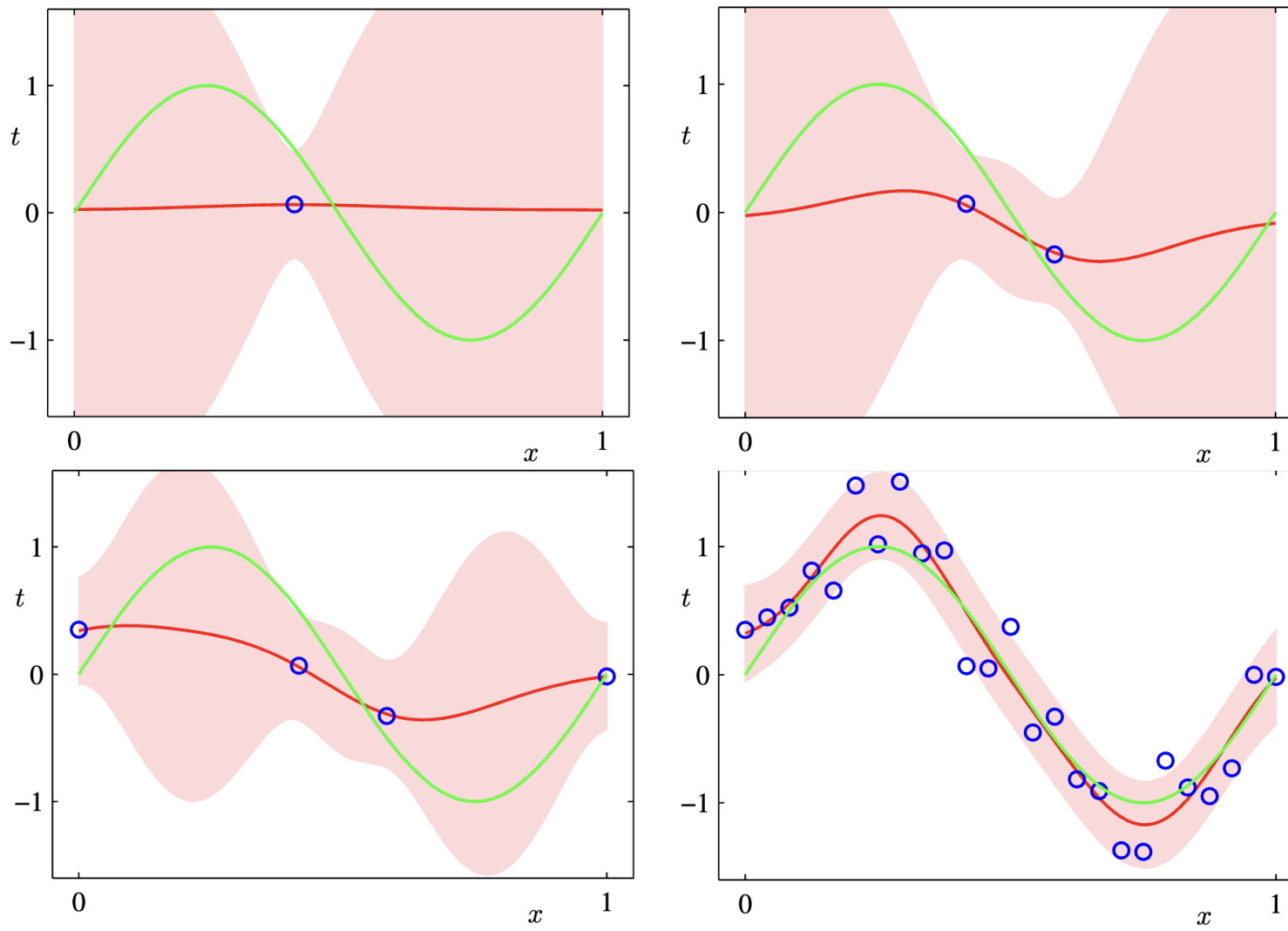
$$m(x) = \beta x^T S \sum_{n=1}^N x_n t_n \quad s^2(x) = \beta^{-1} + x^T S x$$

其中

$$S^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N x_n x_n^T$$

贝叶斯估计

$$p(t|x_0, \mathbf{w}, \beta) = \mathcal{N}(t | \mathbf{w}^T \phi(x_0), \beta^{-1})$$



3.3 什么是线性分类器？

- 线性分类器的概念
- 线性判别函数的一般形式
- 线性判别函数的几何理解
- 广义线性判别函数
- 线性判别函数的齐次化

从贝叶斯分类说起

● 贝叶斯公式：

$$P(w_i|\mathbf{x}) = \frac{P(w_i)P(\mathbf{x}|w_i)}{\sum_{j=1}^n P(w_j)P(\mathbf{x}|w_j)}$$

后验概率

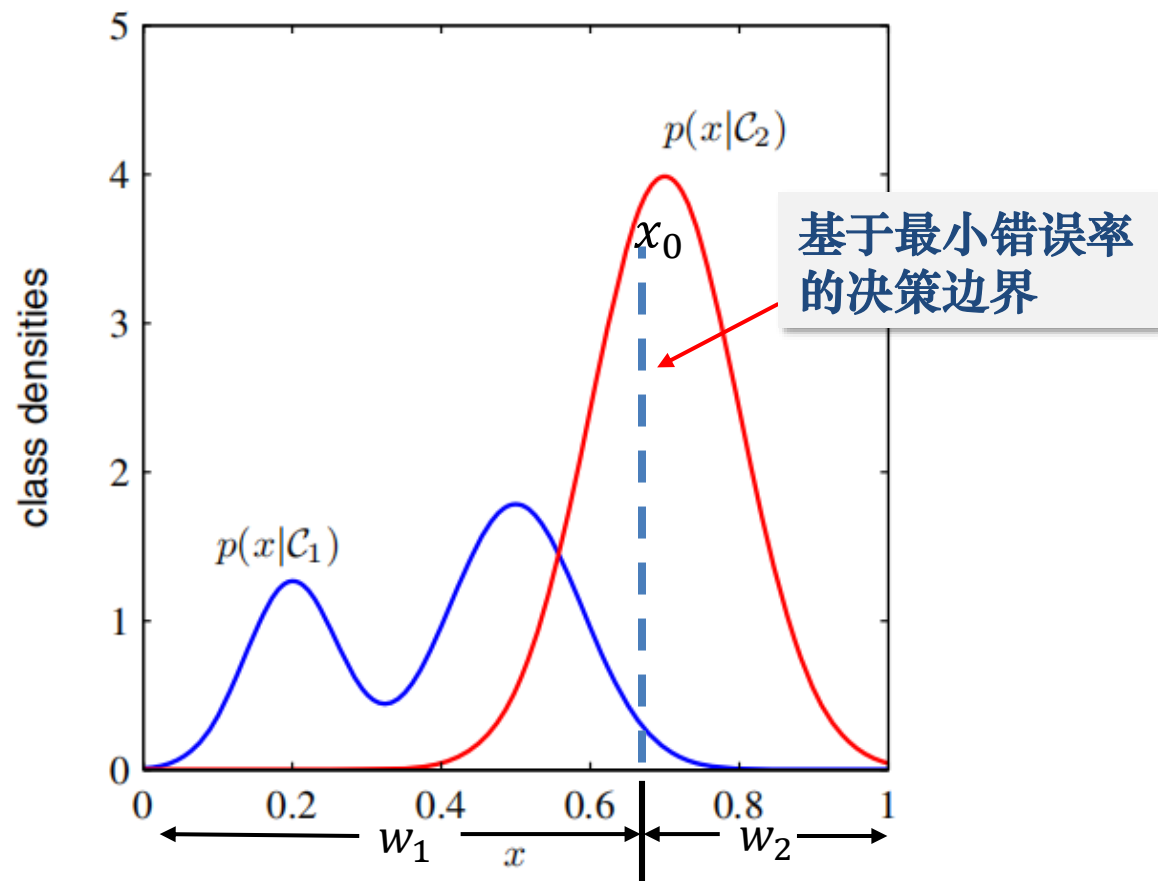
先验概率

似然概率(类条件概率)

证据因子

贝叶斯决策根据贝叶斯公式计算后验概率，基于最大后验概率进行判决

● 决策规则



几个问题

- 已知类条件概率密度 $p(\mathbf{x}|w_i)$ 和先验概率 $P(w_i)$ ，计算后验概率 $p(w_i|\mathbf{x})$ 进行决策

若类条件概率密度参数未知？

- 已知类条件概率密度 $p(\mathbf{x}|w_i)$ 的参数表达式，利用样本估计 $p(\mathbf{x}|w_i)$ 的未知参数，再利用贝叶斯定理将其转化成后验概率 $p(w_i|\mathbf{x})$ 进行决策

- 非参数方法估计

若类条件概率密度形式难以确定？

需要大量样本

线性分类器

- 利用样本集直接设计分类器

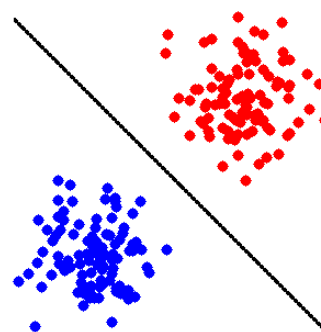
线性回归: $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ 利用样本估计 \mathbf{w} , 对于给定 \mathbf{x} , 计算 y

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$

$$f(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$

$$\mathbf{w}^T \mathbf{x} + w_0 \geq 0 \quad C_1$$

$$\mathbf{w}^T \mathbf{x} + w_0 \leq 0 \quad C_2$$



- 将分类器设计问题转化为求准则函数极值的问题；准则函数：分类器设计的某些要求的函数形式

线性判别函数的一般形式

- 二分类情况下线性判别函数一般形式为 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

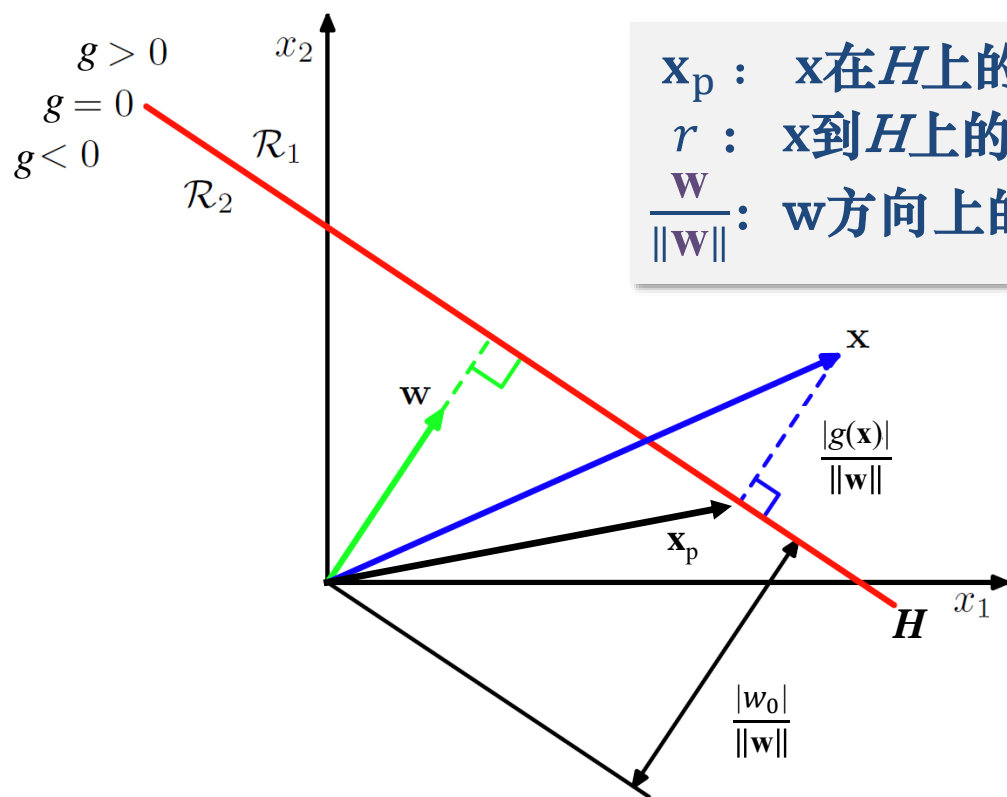
其中 $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$ 是特征向量/样本向量； $\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$ 是权向量； w_0 是阈值权 (常数)

➤ 分类时，如果 $\begin{cases} g(\mathbf{x}) > 0, \mathbf{x} \in w_1 \\ g(\mathbf{x}) < 0, \mathbf{x} \in w_2 \\ g(\mathbf{x}) = 0, \text{可将}\mathbf{x}\text{分到任意一类或拒绝} \end{cases}$

➤ $g(\mathbf{x}) = 0$ 定义了一个决策面，当 $g(\mathbf{x})$ 为线性函数时，决策面就是超平面

线性判别函数的几何理解

- 如果 \mathbf{x}_1 和 \mathbf{x}_2 都在决策面 H 上, 则有 $\mathbf{w}^T \mathbf{x}_1 + w_0 = \mathbf{w}^T \mathbf{x}_2 + w_0$, 即 $\mathbf{w}^T (\mathbf{x}_1 - \mathbf{x}_2) = 0$, 说明 \mathbf{w} 和超平面 H 上任一向量正交, 即 \mathbf{w} 是 H 的法向量
- 判别函数 $g(\mathbf{x})$ 可看成是特征空间中某点 \mathbf{x} 到超平面 H 距离的一种代数度量



\mathbf{x}_p : \mathbf{x} 在 H 上的投影向量
 r : \mathbf{x} 到 H 上的垂直距离
 $\frac{\mathbf{w}}{\|\mathbf{w}\|}$: \mathbf{w} 方向上的单位向量

$$\begin{aligned}\mathbf{x} &= \mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \\ g(\mathbf{x}) &= \mathbf{w}^T \left(\mathbf{x}_p + r \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0 \\ &= \mathbf{w}^T \mathbf{x}_p + r \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} + w_0 \\ &= r \|\mathbf{w}\| \\ r &= \frac{g(\mathbf{x})}{\|\mathbf{w}\|}\end{aligned}$$

线性判别函数的几何理解

- 若 \mathbf{x} 为原点，则 $g(\mathbf{x}) = w_0$ ；从原点到超平面 H 的距离 $r_0 = \frac{w_0}{\|\mathbf{w}\|}$

如果 $w_0 > 0$ ，则原点在 H 的正侧

如果 $w_0 < 0$ ，则原点在 H 的负侧

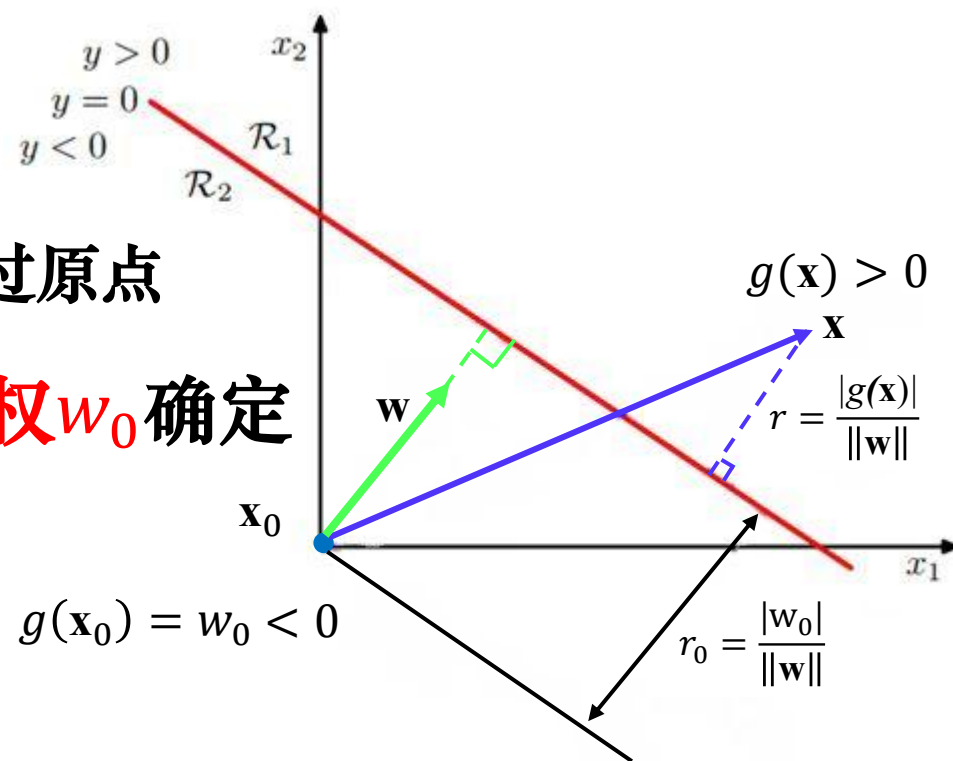
如果 $w_0 = 0$ ，则 $g(\mathbf{x})$ 具有齐次形式，超平面 H 通过原点

- 超平面方向由**权向量** \mathbf{w} 决定；位置由**阈值权** w_0 确定

判别函数 $g(\mathbf{x})$ 正比于 \mathbf{x} 点到超平面的代数距离

当 \mathbf{x} 在 H 的正侧时， $g(\mathbf{x}) > 0$

当 \mathbf{x} 在 H 的负侧时， $g(\mathbf{x}) < 0$



广义线性判别函数

- 两类别问题, X 是一维样本空间

不适用于非凸和多连通区域划分

若 $x < a$ 或 $x > b$, $x \in w_1$; 若 $a < x < b$, $x \in w_2$

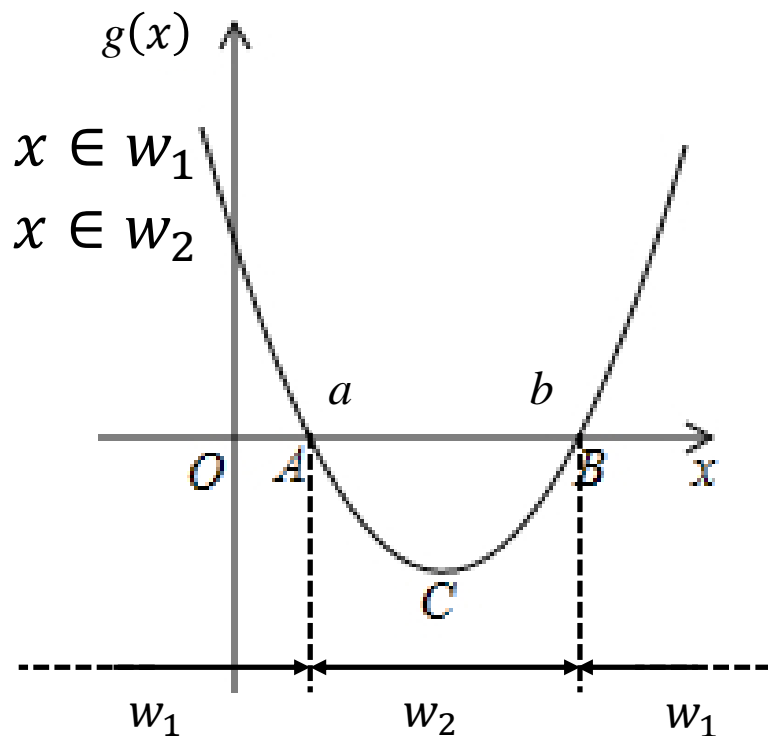
如果建立 $g(x) = (x - a)(x - b)$ 决策规则 $\begin{cases} g(x) > 0, x \in w_1 \\ g(x) < 0, x \in w_2 \end{cases}$

一般形式 $g(x) = c_0 + c_1x + c_2x^2$, 可转化为

$$g(x) = \mathbf{a}^T \mathbf{y} = \sum_{i=1}^3 a_i y_i \quad \text{其中 } \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} \quad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} c_0 \\ c_1 \\ c_2 \end{bmatrix}$$

称 $g(x) = \mathbf{a}^T \mathbf{y}$ 为 **广义线性判别函数**, \mathbf{a} 为 **广义权向量**

当 $g(x)$ 取二次形式时, 也称为 **二次判别分析(QDA)***



利用线性函数的简单性解决复杂问题
维数大大增加 \rightarrow “维数灾难”

*更多内容详见Christopher M. Bishop, Pattern Recognition and Machine Learning, P199-200

线性判别函数的齐次化

- 线性判别函数 $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$

改写成 $g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i = \sum_{i=1}^d a_i y_i = \mathbf{a}^T \mathbf{y}$ 称为**线性判别函数的齐次简化**

其中 $\mathbf{y} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$ 称为**增广样本向量**； $\mathbf{a} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$ 称为**增广权向量**

$\hat{d} = d + 1$ ；**y**比**x**增加一维，保持样本空间欧式距离不变，变换后的样本仍全部位于**d**维子空间(原**X**空间)中

方程 $\mathbf{a}^T \mathbf{y} = 0$ 在Y空间确定了一个通过原点的超平面 \hat{H} ，它对d维子空间的划分与原决策面 $\mathbf{w}^T \mathbf{x} + w_0 = 0$ 对原X空间的划分完全相同。Y空间中任意一点y到 \hat{H} 的距离：

$$r = \frac{g(\mathbf{x})}{\|\mathbf{a}\|} = \frac{\mathbf{a}^T \mathbf{y}}{\|\mathbf{a}\|}$$

线性分类器设计的一般过程

- 利用训练样本建立线性判别函数

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

$$g(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i = \sum_{i=1}^d a_i y_i = \mathbf{a}^T \mathbf{y}$$

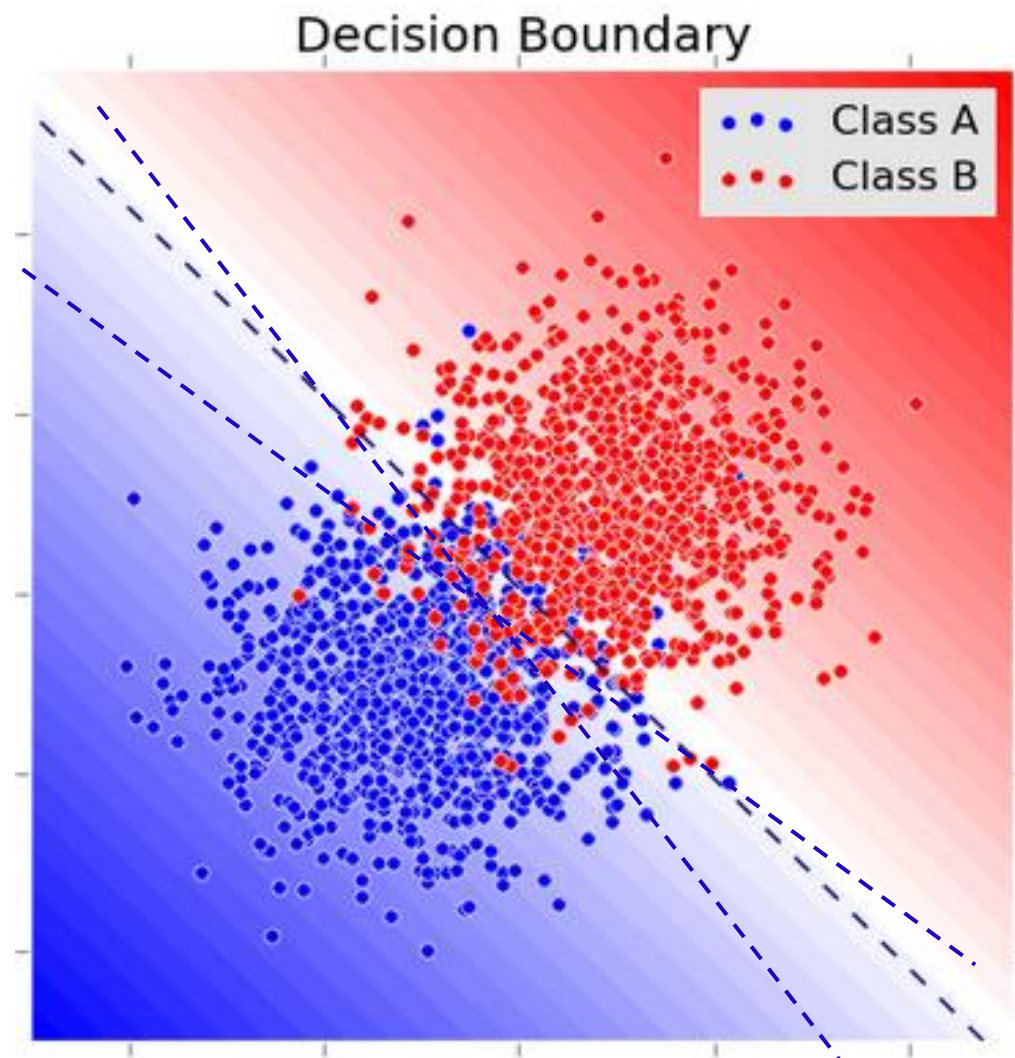
最好的结果一般出现在准则函数的极值点上，所以将分类器设计问题转化为求准则函数极值 \mathbf{w}^* , w_0^* 或 \mathbf{a}^* 的问题。

- 步骤1：具有类别标志的样本集 $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ 或其增广样本集 \mathcal{Y}
- 步骤2：确定准则函数 \mathcal{J} ，满足① \mathcal{J} 是样本集和 \mathbf{w} , w_0 或 \mathbf{a} 的函数；② \mathcal{J} 的值反映分类器的性能，其极值对应“最好”的决策
- 步骤3：优化求解准则函数极值 \mathbf{w}^* , w_0^* 或 \mathbf{a}^*

最终得到线性判别函数： $g(\mathbf{x}) = \mathbf{w}^{*T} \mathbf{x} + w_0$ 或 $g(\mathbf{x}) = \mathbf{a}^{*T} \mathbf{y}$ ，对于未知类别样本 \mathbf{x}_k ，计算 $g(\mathbf{x}_k)$ 并通过决策规则判断其类别

准则函数的设计

- Fisher准则
- 感知机准则
- 最小二乘准则



3.4 Fisher准则

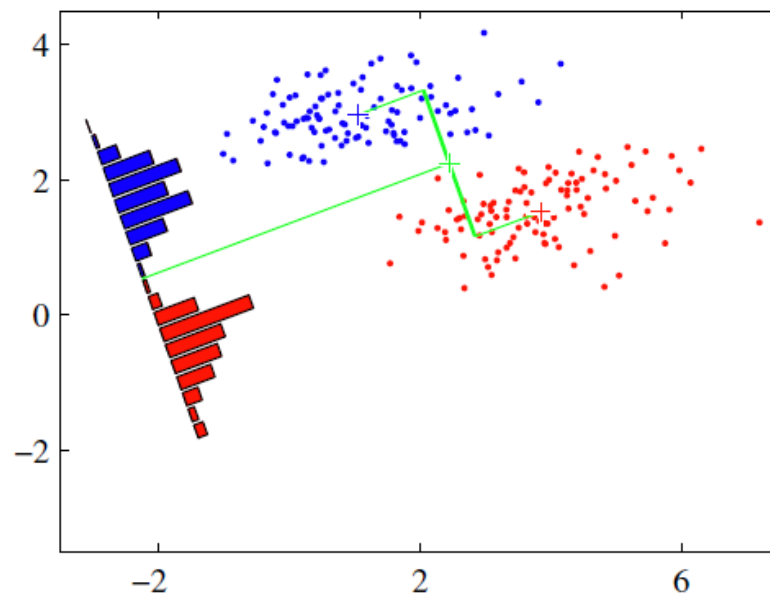
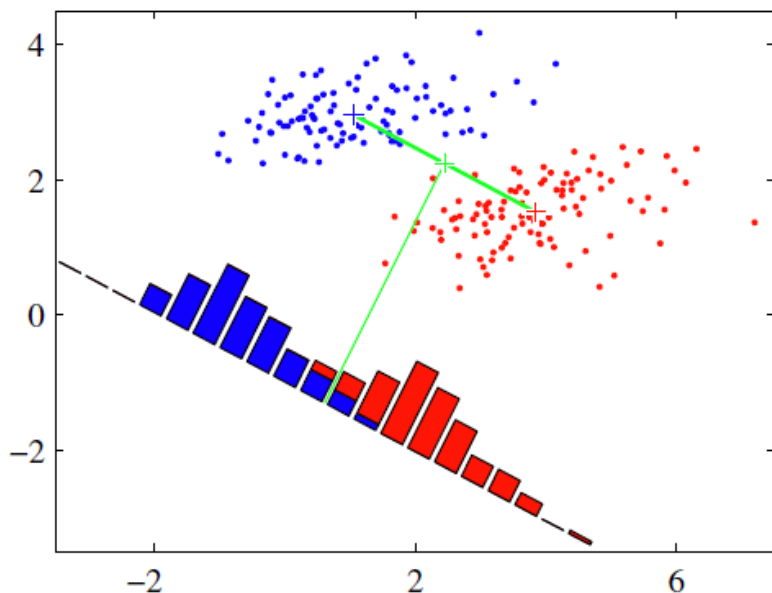
- Fisher准则的概念
- Fisher准则的求解

Fisher准则的概念

- 考虑把 d 维空间的样本投影到一条直线上形成一维空间。在一般情况下总可以找到某个方向，使样本在这个方向的直线上的投影分开得最好

【R. A. Fisher, 1936】

**Fisher准则就是要解决
如何根据实际情况找到这条最好的、最易于分类的投影方向的问题**



Fisher准则的推导

● 寻找最好投影方向 w^*

➤ 以二分类问题为例， d 维样本 x_1, x_2, \dots, x_N ，其中 N_1 个属于 w_1 类记为子集 X_1 ， N_2 个属于 w_2 类记为子集 X_2

➤ 在 d 维 X 空间

各类样本的均值向量 m_i
$$m_i = \frac{1}{N_i} \sum_{x \in X_i} x, i = 1, 2$$

样本类内离散度矩阵 S_i 和总类内离散度矩阵 S_W

$$S_i = \sum_{x \in X_i} (x - m_i)(x - m_i)^T, i = 1, 2 \quad S_W = S_1 + S_2$$

样本类间离散度矩阵 S_b

$$S_W = P(w_1)S_1 + P(w_2)S_2$$

$$S_b = (m_1 - m_2)(m_1 - m_2)^T \quad S_b = P(w_1)P(w_2)(m_1 - m_2)(m_1 - m_2)^T$$

Fisher准则

- 寻找最好投影方向 w^*

- 以二分类问题为例, d 维样本 x_1, x_2, \dots, x_N , 其中 N_1 个属于 w_1 类记为子集 X_1 , N_2 个属于 w_2 类记为子集 X_2
- 在一维 Y 空间 $y_n = w^T x_n$

各类样本均值 \widetilde{m}_i

$$\widetilde{m}_i = \frac{1}{N_i} \sum_{y \in \eta_i} y, i = 1, 2$$

样本类内离散度 \tilde{s}_i^2 和总类内离散度 \tilde{s}_w

$$\tilde{s}_i^2 = \sum_{y \in \eta_i} (y - \widetilde{m}_i)^2, i = 1, 2$$

$$\tilde{s}_w = \tilde{s}_1^2 + \tilde{s}_2^2$$

Fisher准则

- 希望投影后在一维 Y 空间中各类样本尽可能分开，即两类均值之差越大越好
- 同时希望各类样本内部尽量密集，即类内离散度越小越好



➤ Fisher准则定义为类间方差与类内方差之比：

$$J_F(\mathbf{w}) = \frac{(\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)^2}{\tilde{\mathbf{S}}_1^2 + \tilde{\mathbf{S}}_2^2}$$

两类均值之差越大越好

类内离散度越小越好

Fisher准则

● 求 $J_F(\mathbf{w})$ 取得最大值的 \mathbf{w} $J_F(\mathbf{w}) = \frac{(\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)^2}{\tilde{\mathbf{S}}_1^2 + \tilde{\mathbf{S}}_2^2}$

$$\tilde{\mathbf{m}}_i = \frac{1}{N_i} \sum_{y \in \eta_i} y = \frac{1}{N_i} \sum_{\mathbf{x} \in \mathcal{X}_i} \mathbf{w}^T \mathbf{x} = \mathbf{w}^T \left(\frac{1}{N_i} \sum_{\mathbf{x} \in \mathcal{X}_i} \mathbf{x} \right) = \mathbf{w}^T \mathbf{m}_i$$

$$(\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)^2 = (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 = \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} = \mathbf{w}^T \mathbf{S}_b \mathbf{w}$$

$$\tilde{\mathbf{S}}_i^2 = \sum_{y \in \eta_i} (y - \tilde{\mathbf{m}}_i)^2 = \sum_{\mathbf{x} \in \mathcal{X}_i} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mathbf{m}_i)^2 = \mathbf{w}^T \sum_{\mathbf{x} \in \mathcal{X}_i} (\mathbf{x} - \mathbf{m}_i) (\mathbf{x} - \mathbf{m}_i)^T \mathbf{w} = \mathbf{w}^T \mathbf{S}_i \mathbf{w}$$

$$\tilde{\mathbf{S}}_1^2 + \tilde{\mathbf{S}}_2^2 = \mathbf{w}^T (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{w} = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \quad \rightarrow \quad J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

➤ Lagrange乘子法求解:

$$\mathbf{w}^* = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

➤ 分类时确定分界阈值 y_0 , 与 $y = \mathbf{w}^{*T} \mathbf{x}$ 比较进行决策

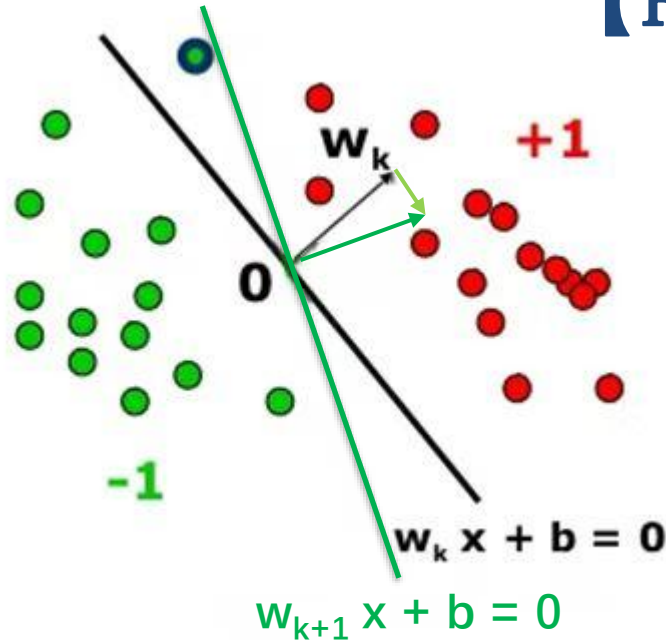
3.5 感知机准则

- 感知机准则的概念
- 感知机准则的求解

感知机准则

- 感知准则是一种自学习判别函数生成方法，由于 Rosenblatt 试图将其用于脑模型**感知机**，因此得名
- 该方法对随意给定的判别函数初始值，通过样本分类训练过程逐步对其修正直至最终确定

【F. Rosenblatt, 1957】



几个基本概念

● 线性可分性

- 一组容量为 N 的样本集 $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$, 其中 \mathbf{y}_n 为 \hat{d} 维增广样本向量, 分别来自 w_1 类和 w_2 类, 如果存在权向量 \mathbf{a} , 使得对于任何 $\mathbf{y} \in w_1$, 都有 $\mathbf{a}^T \mathbf{y} > 0$, 而对于任何 $\mathbf{y} \in w_2$, 都有 $\mathbf{a}^T \mathbf{y} < 0$, 则称这组样本为线性可分的, 反之亦然

● 样本的规范化

$$\begin{cases} \mathbf{a}^T \mathbf{y}_i > 0, \mathbf{y}_i \in w_1 \\ \mathbf{a}^T \mathbf{y}_j < 0, \mathbf{y}_j \in w_2 \end{cases}$$



$$\mathbf{y}'_n = \begin{cases} \mathbf{y}_i & , \mathbf{y}_i \in w_1 \\ -\mathbf{y}_j & , \mathbf{y}_j \in w_2 \end{cases}$$

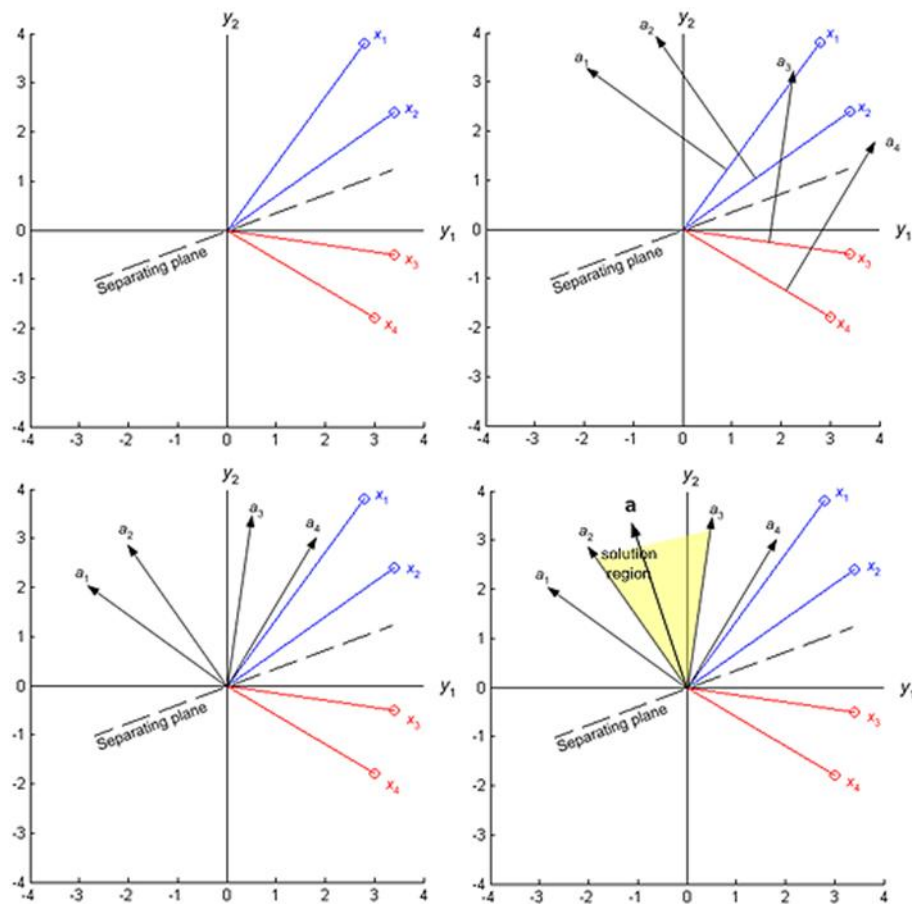
规范化增广样本向量



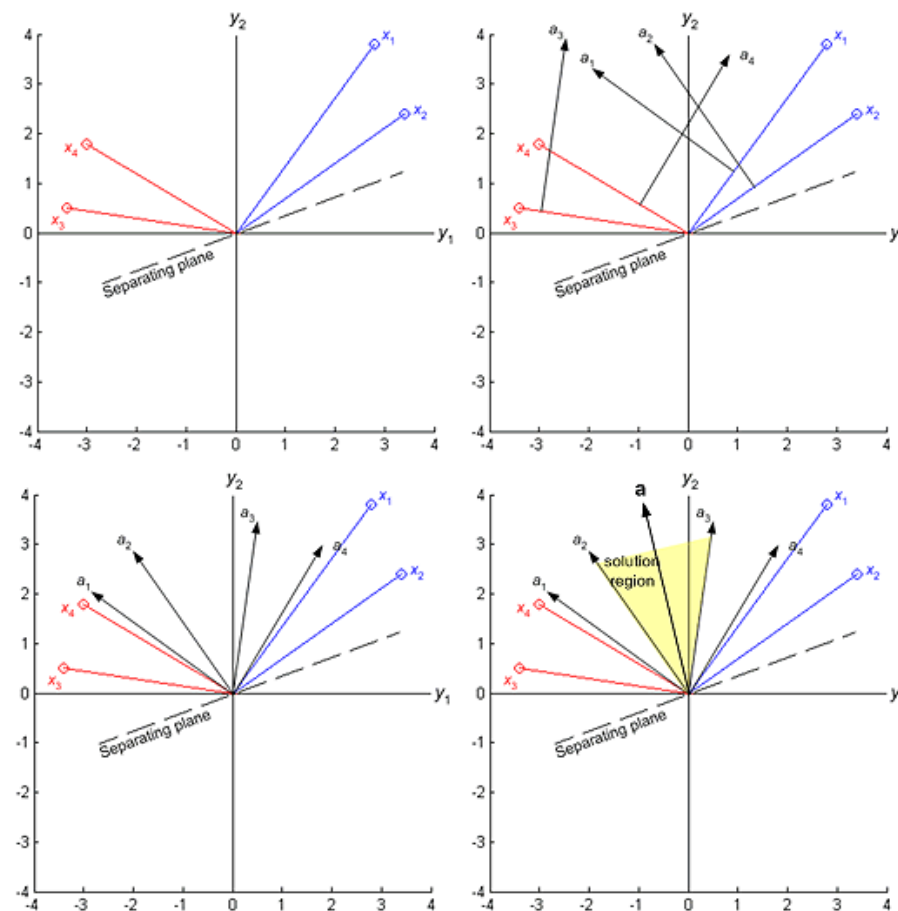
$$\mathbf{a}^T \mathbf{y}'_i > 0, n = 1, 2, \dots, N$$

几个基本概念

● 解向量和解区



未规范化



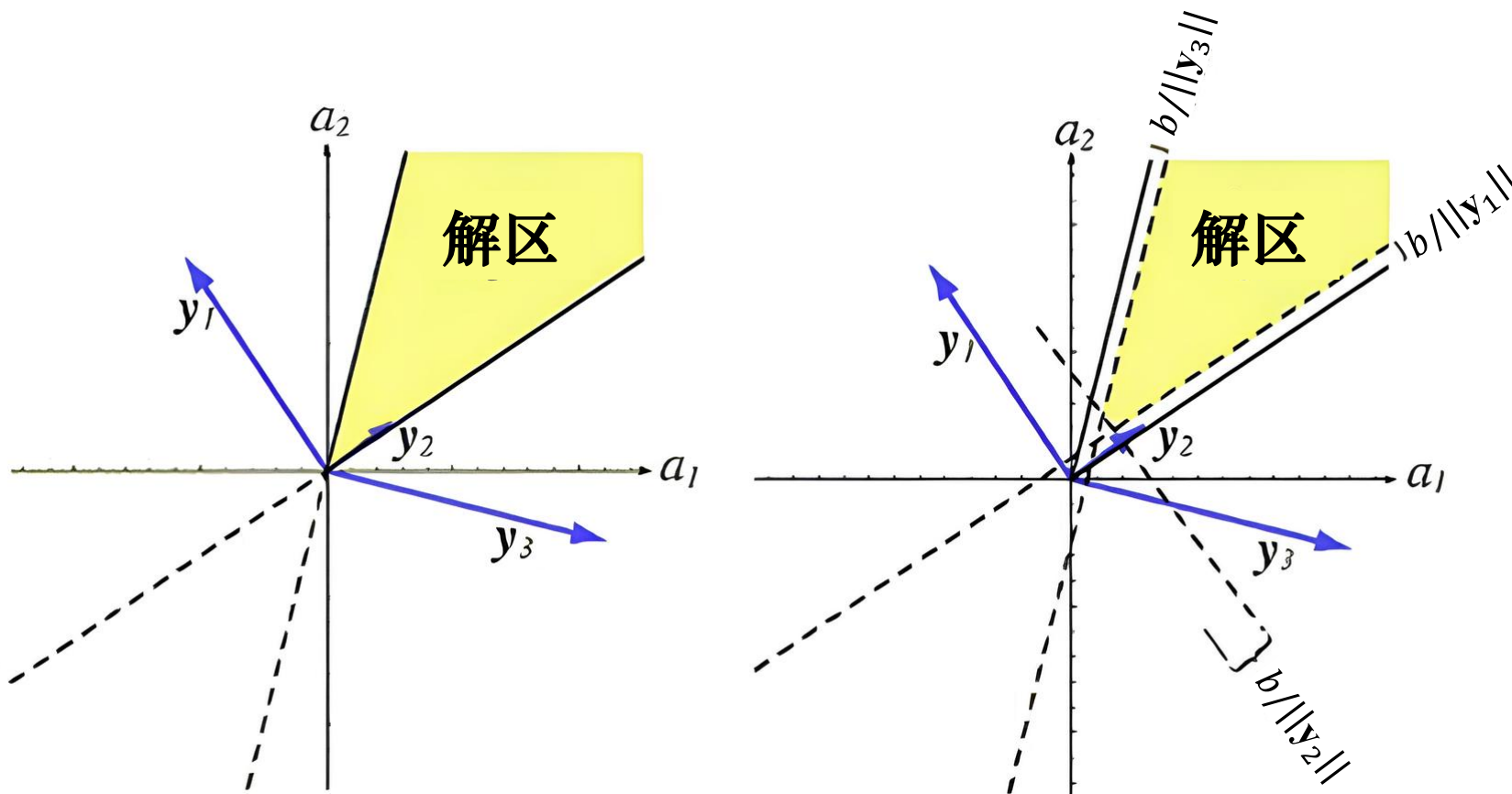
规范化

规范化以后，解向量均满足 $a^T y_n > 0$ ，解区并没有变化

几个基本概念

● 对解区的限制

- 使解向量更可靠 $\mathbf{a}^T \mathbf{y}_n \geq b > 0$ ，避免解收敛到解区边界的某点上



感知机准则

● 寻找解向量 \mathbf{a}^*

➤ 对于一组样本 $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$, 其中 \mathbf{y}_n 是规范化增广样本向量, 使得:

$$\mathbf{a}^T \mathbf{y}_n > 0, n = 1, 2, \dots, N$$

➤ 对于线性可分问题, 构造准则函数 $J_P(\mathbf{a}) = \sum_{\mathbf{y} \in \eta^k} (-\mathbf{a}^T \mathbf{y})$

η^k 是被权向量 \mathbf{a} 错分的样本集合, 即当 \mathbf{y} 被错分时, 就有 $\mathbf{a}^T \mathbf{y}_n \leq 0$

因此 $J_P(\mathbf{a}) \geq 0$, 仅当 \mathbf{a} 为解向量或在解区边界时 $J_P(\mathbf{a}) = 0$

也就是说, 当且仅当 η_k 为空集时 $J_P^*(\mathbf{a}) = \min J_P(\mathbf{a}) = 0$

此时无错分样本, 这时的 \mathbf{a} 就是解向量 \mathbf{a}^*

感知机准则函数的求解

- 求使 $J_P(\mathbf{a})$ 达到最小值的 \mathbf{a}^*

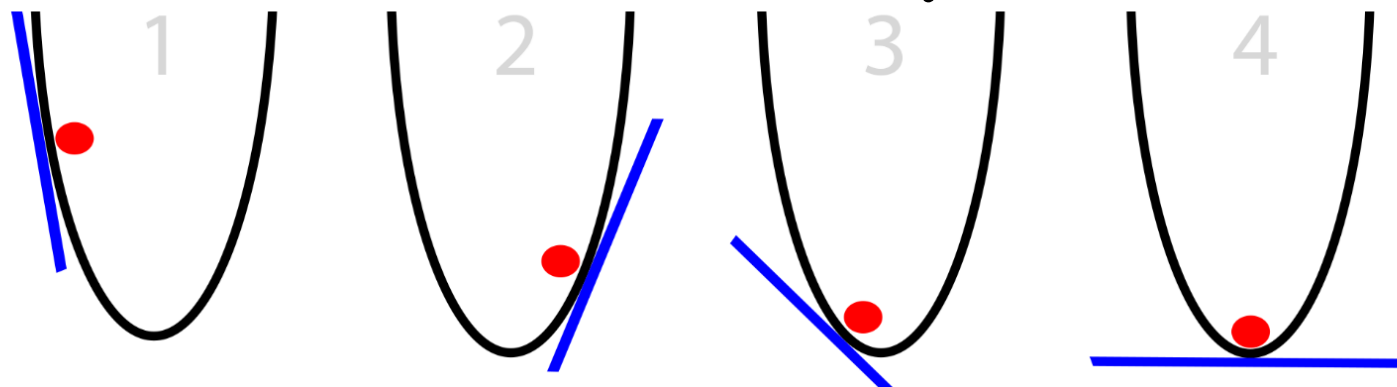
➤ 采用梯度下降法求解 $J_P(\mathbf{a}) = \sum_{\mathbf{y} \in \eta_k} (-\mathbf{a}^T \mathbf{y})$

$$\nabla J_P(\mathbf{a}) = \frac{\partial J_P(\mathbf{a})}{\partial \mathbf{a}} = \sum_{\mathbf{y} \in \eta_k} (-\mathbf{y})$$

$$\mathbf{a}(k+1) = \mathbf{a}(k) - \rho_k \nabla J$$

$$\mathbf{a}(k+1) = \mathbf{a}(k) + \rho_k \sum_{\mathbf{y} \in \eta_k} \mathbf{y}$$

梯度下降法迭代公式



感知机准则-示例

● 二分情况下的样本集

➤ 类别1:

$$\mathbf{x}_1 = \begin{pmatrix} -2 \\ 2 \end{pmatrix}, \mathbf{x}_2 = \begin{pmatrix} -2 \\ -2 \end{pmatrix}$$

$$\mathbf{y}_1 = \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix}$$

➤ 类别2:

$$\mathbf{x}_3 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \mathbf{x}_4 = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$$

$$\mathbf{y}_3 = \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix}, \mathbf{y}_4 = \begin{pmatrix} -1 \\ -2 \\ 1 \end{pmatrix}$$

➤ 初始化权重向量:

$$\mathbf{a}(1) = \begin{pmatrix} 0 \\ 2 \\ 1 \end{pmatrix}$$

$$\mathbf{a}(1)^T = (0 \quad 2 \quad 1)$$

感知机准则-示例

● 迭代过程

➤ (1) $\mathbf{y}^{(1)T} = (1 \quad -2 \quad 2)$

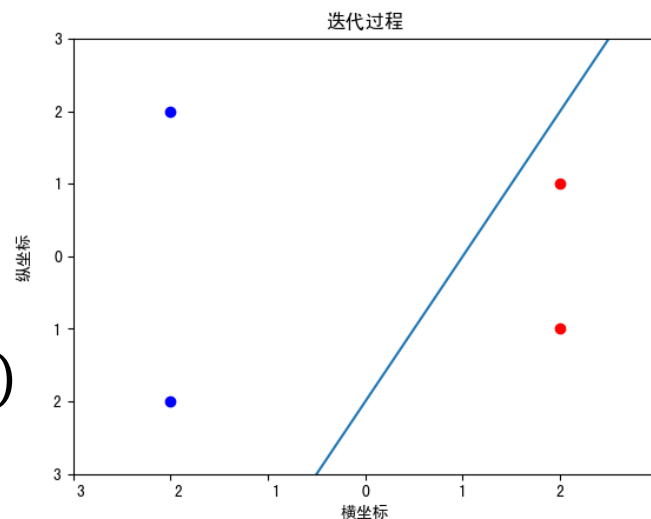
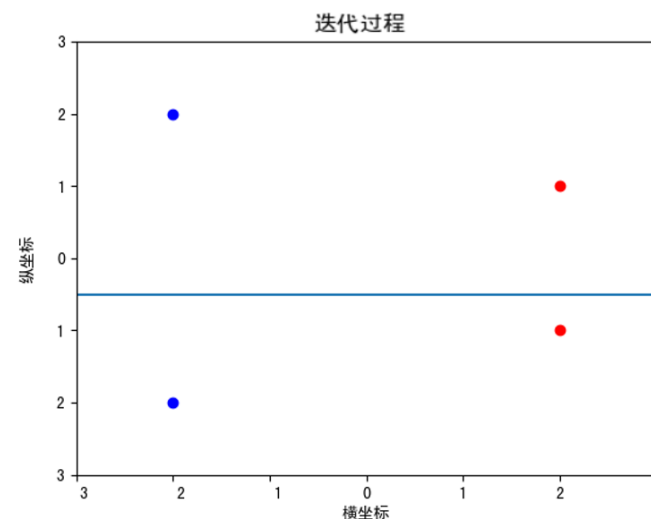
$$\mathbf{a}(1)^T \mathbf{y}^{(1)} = (0 \quad 2 \quad 1) \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix} = -2 < 0$$

$$\mathbf{a}(2)^T = (0 \quad 2 \quad 1) + (1 \quad -2 \quad 2) = (1 \quad 0 \quad 3)$$

➤ (2) $\mathbf{y}^{(2)T} = (1 \quad -2 \quad -2)$

$$\mathbf{a}(2)^T \mathbf{y}^{(2)} = (1 \quad 0 \quad 3) \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix} = -5 < 0$$

$$\mathbf{a}(3)^T = (1 \quad 0 \quad 3) + (1 \quad -2 \quad -2) = (2 \quad -2 \quad 1)$$



感知机准则-示例

● 迭代过程

➤ (3) $\mathbf{y}^{(3)T} = (-1 \quad -2 \quad -1)$

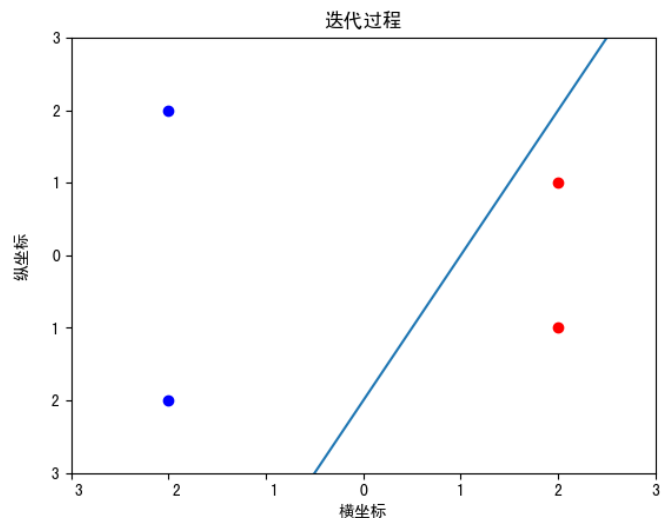
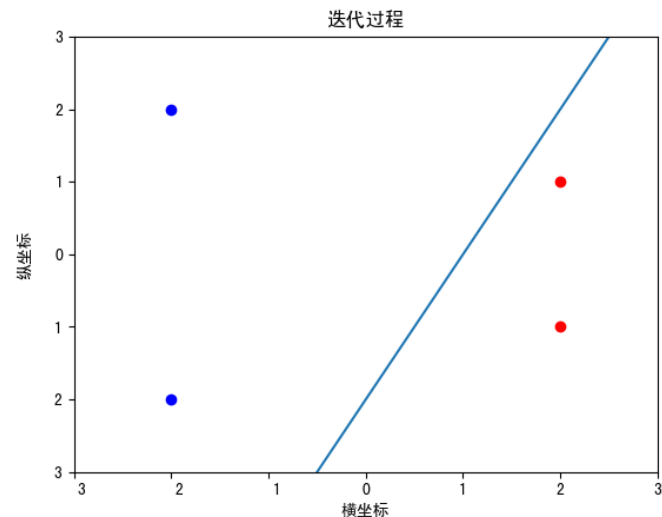
$$\mathbf{a}(3)^T \mathbf{y}^{(3)} = (2 \quad -2 \quad 1) \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix} = 1 > 0$$

$$\mathbf{a}(4)^T = (2 \quad -2 \quad 1) \quad (\text{没有变化})$$

➤ (4) $\mathbf{y}^{(4)T} = (2 \quad -2 \quad 1)$

$$\mathbf{a}(4)^T \mathbf{y}^{(4)} = (2 \quad -2 \quad 1) \begin{pmatrix} -1 \\ -2 \\ 1 \end{pmatrix} = 3 > 0$$

$$\mathbf{a}(5)^T = (2 \quad -2 \quad 1) \quad (\text{没有变化})$$



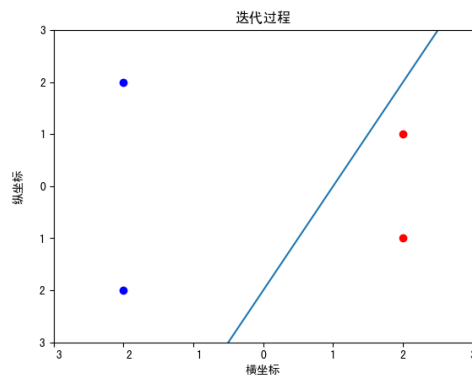
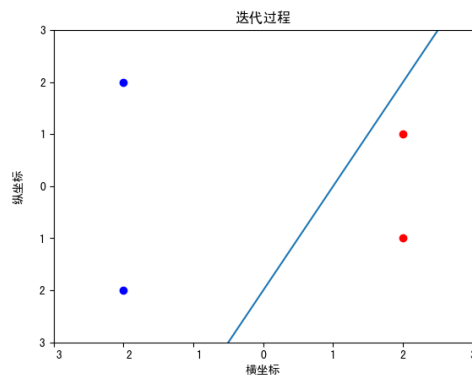
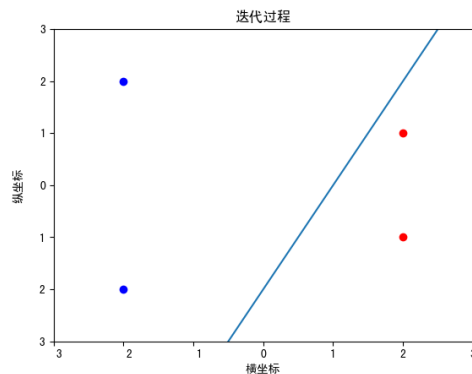
感知机准则-示例

● 迭代过程

➤ (5) $\mathbf{y}^{(5)T} = (2 \quad -2 \quad 1)$
 $\mathbf{a}^{(5)T} \mathbf{y}^{(5)} = (2 \quad -2 \quad 1) \begin{pmatrix} 1 \\ -2 \\ 2 \end{pmatrix} = 8 > 0$
 $\mathbf{a}^{(6)T} = (2 \quad -2 \quad 1)$ (没有变化)

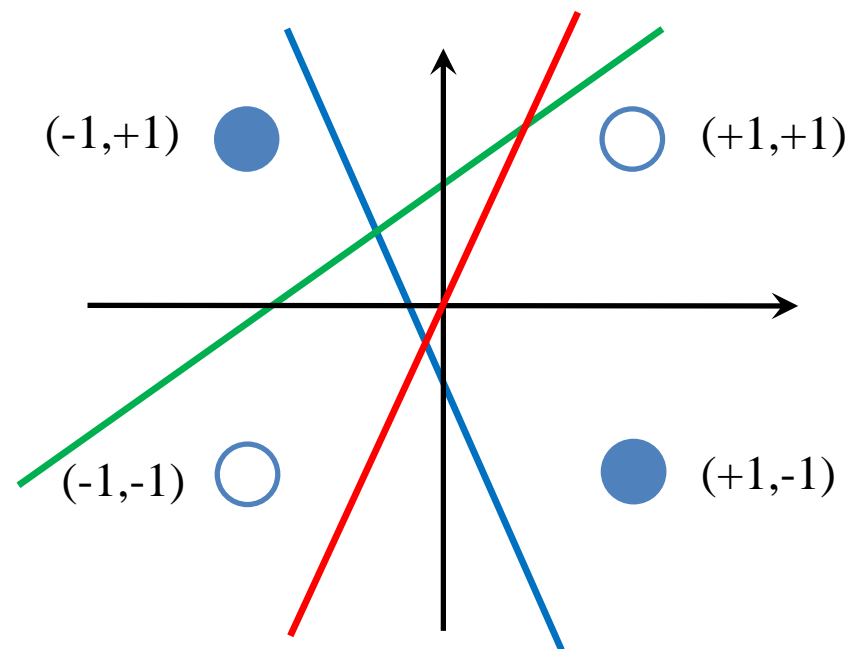
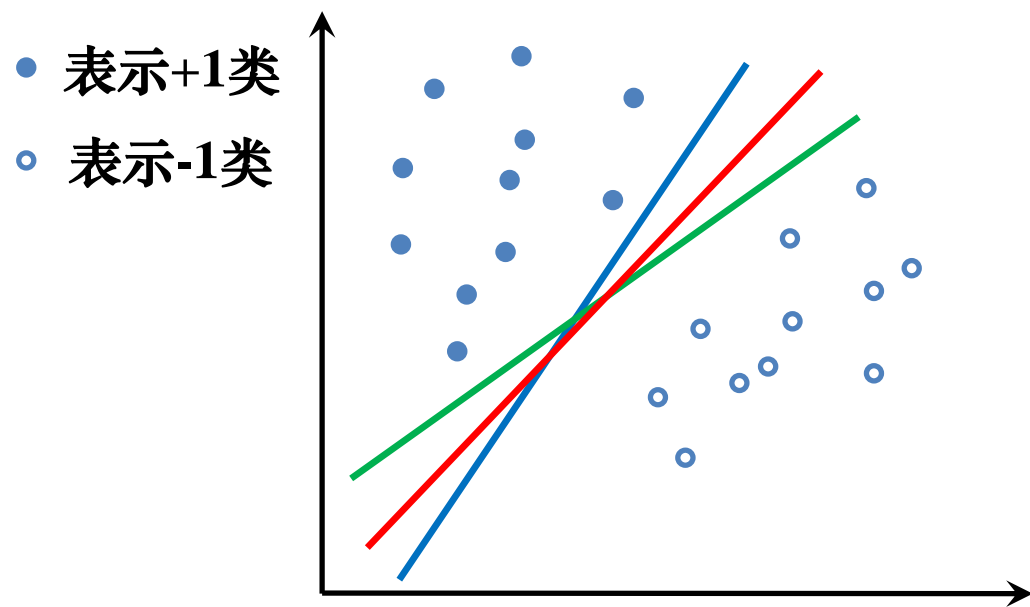
➤ (6) $\mathbf{y}^{(6)T} = (1 \quad -2 \quad 2)$
 $\mathbf{a}^{(6)T} \mathbf{y}^{(6)} = (2 \quad -2 \quad 1) \begin{pmatrix} 1 \\ -2 \\ -2 \end{pmatrix} = 4 > 0$
 $\mathbf{a}^{(7)T} = (2 \quad -2 \quad 1)$ (没有变化)

➤ (7) $\mathbf{y}^{(7)T} = (-1 \quad -2 \quad -1)$
 $\mathbf{a}^{(7)T} \mathbf{y}^{(7)} = (2 \quad -2 \quad 1) \begin{pmatrix} -1 \\ -2 \\ -1 \end{pmatrix} = 1 > 0$
 $\mathbf{a}^{(8)T} = (2 \quad -2 \quad 1)$ (没有变化)



感知机准则

- 在线性可分情形下，感知机准则一定收敛，但收敛结果不唯一
- 在线性不可分情况下不收敛



3.6 最小二乘准则

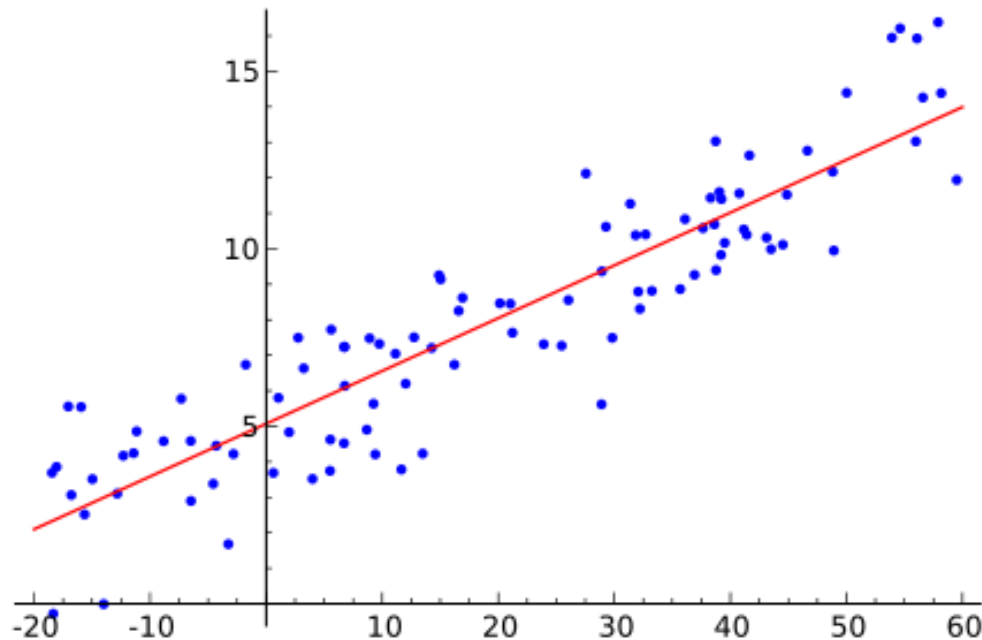
- 最小二乘准则的概念
- 最小二乘准则的求解

最小二乘准则

- 最小二乘法(最小平方误差法)通过最小化误差的平方和寻找数据的最佳函数匹配，即使求得的数据与实际数据之间误差的平方和最小

【A.-M. Legendre, 1806提出】

【C. Gauss, 1809提出, 1829证明】



最小二乘准则

- 目标：寻找最优投影方向 \mathbf{a}^*

$$\mathbf{a}^T \mathbf{y}_n > 0 \quad \longrightarrow \quad \mathbf{a}^T \mathbf{y}_n = b_n > 0 \quad b_n \text{ 是任意给定的正常数}$$

- 方程组形式： $\mathbf{Y}\mathbf{a} = \mathbf{b}$

$$\mathbf{Y} = \begin{bmatrix} \mathbf{y}_1^T \\ \mathbf{y}_2^T \\ \vdots \\ \mathbf{y}_N^T \end{bmatrix} = \begin{bmatrix} y_{11} & y_{12} & \cdots & y_{1\hat{d}} \\ y_{21} & y_{22} & \cdots & y_{2\hat{d}} \\ \cdots & \cdots & \cdots & \cdots \\ y_{N1} & y_{N2} & \cdots & y_{N\hat{d}} \end{bmatrix} \quad \mathbf{b} = [b_1 \quad b_2 \quad \cdots \quad b_N]$$

\mathbf{y}^n 是规范化增广向量样本； \mathbf{Y} 是 $N \times \hat{d}$ 维矩阵，通常 $N > \hat{d}$ ，一般为列满秩阵
 \mathbf{b} 是 N 维向量， $b_n > 0$ ， $n=1, 2, \dots, N$

- 方程数多于未知数的矛盾方程组通常没有精确解
- 定义误差向量： $\mathbf{e} = \mathbf{Y}\mathbf{a} - \mathbf{b}$ 以及平方误差准则函数

$$J_s(\mathbf{a}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{n=1}^N (\mathbf{a}^T \mathbf{y}_n - b_n)^2$$

最小二乘准则

- 求使 $J_S(\mathbf{a})$ 最小的 \mathbf{a}^* (最小二乘近似解/伪逆解/均方误差(MSE)解)

- 采用解析法求伪逆解

$$J_S(\mathbf{a}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{n=1}^N (\mathbf{a}^T \mathbf{y}_n - b_n)^2$$

$$\nabla J_S(\mathbf{a}) = \sum_{n=1}^N 2(\mathbf{a}^T \mathbf{y}_n - b_n) \mathbf{y}_n = 2\mathbf{Y}^T (\mathbf{Y}\mathbf{a} - \mathbf{b})$$

$$\text{令 } \nabla J_S(\mathbf{a}) = 0 \text{ 得 } \mathbf{Y}^T \mathbf{Y} \mathbf{a}^* = \mathbf{Y}^T \mathbf{b}$$

矩阵 $\mathbf{Y}^T \mathbf{Y}$ 是 $\hat{d} \times \hat{d}$ 方阵，一般非奇异

$$\text{唯一解 } \mathbf{a}^* = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} = \mathbf{Y}^+ \mathbf{b}$$

其中 $\hat{d} \times N$ 矩阵 $\mathbf{Y}^+ = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$ 是 \mathbf{Y} 的左逆矩阵

- 如何选 \mathbf{b} ?

$$\mathbf{b} = \begin{bmatrix} N/N_1 \\ \cdots \\ N/N_1 \\ N/N_2 \\ \cdots \\ N/N_2 \end{bmatrix} \quad \begin{matrix} N_1 \uparrow \\ \\ N_2 \uparrow \end{matrix}$$

\mathbf{a}^* 等价于 Fisher 解

$$g_0(x) = P(w_1|\mathbf{x}) - P(w_2|\mathbf{x})$$

$$N \rightarrow \infty, \mathbf{b} = [1, 1, \cdots, 1]^T$$

以最小均方误差逼近贝叶斯判别函数

最小二乘准则

- 求使 $J_S(\mathbf{a})$ 最小的 \mathbf{a}^* (最小二乘近似解/伪逆解/均方误差(MSE)解)

$$\mathbf{a}^* = \mathbf{Y}^+ \mathbf{b}$$

$$\mathbf{Y}^+ = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$$

➤ 问题：①要求 $\mathbf{Y}^T \mathbf{Y}$ 非奇异；②求 \mathbf{Y}^+ 计算量大同时可能引入较大误差

➤ 采用梯度下降法求解

$$J_S(\mathbf{a}) = \|\mathbf{e}\|^2 = \|\mathbf{Y}\mathbf{a} - \mathbf{b}\|^2 = \sum_{n=1}^N (\mathbf{a}^T \mathbf{y}_n - \mathbf{b}_n)^2$$

$$\nabla J_S(\mathbf{a}) = 2\mathbf{Y}^T(\mathbf{Y}\mathbf{a} - \mathbf{b})$$

$$\begin{cases} \mathbf{a}(1), \text{随机值} \\ \mathbf{a}(k+1) = \mathbf{a}(k) - \rho_k \mathbf{Y}^T(\mathbf{Y}\mathbf{a} - \mathbf{b}) \end{cases}$$

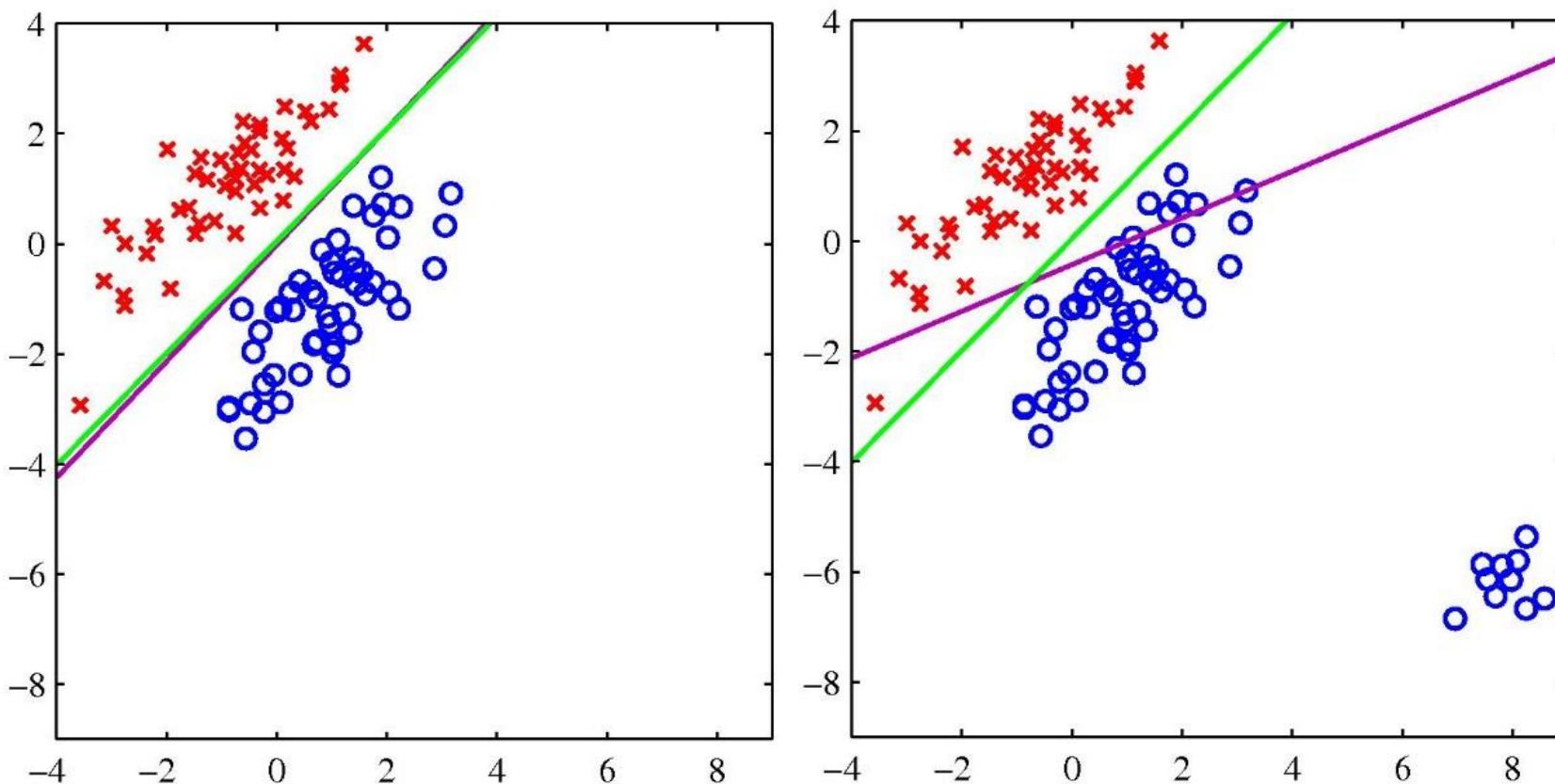
可以证明，选择 $\rho_k = \frac{\rho_1}{k}$ ， ρ_1 是任意常数，该算法权向量收敛于使 $\nabla J_S(\mathbf{a}) = 2\mathbf{Y}^T(\mathbf{Y}\mathbf{a} - \mathbf{b}) = 0$ 的权向量 \mathbf{a}^*

➤ 不要求 $\mathbf{Y}^T \mathbf{Y}$ 奇异与否，只计算 $\hat{d} \times \hat{d}$ 方阵 $\mathbf{Y}^T \mathbf{Y}$ ，比 $\hat{d} \times N$ 阵 \mathbf{Y}^+ 计算量小

最小二乘准则的局限性

- 对于异常值(Outlier)非常敏感

- 平方误差准则函数会惩罚那些“太正确”的预测，因为它们在决策边界的正确一侧距离太远

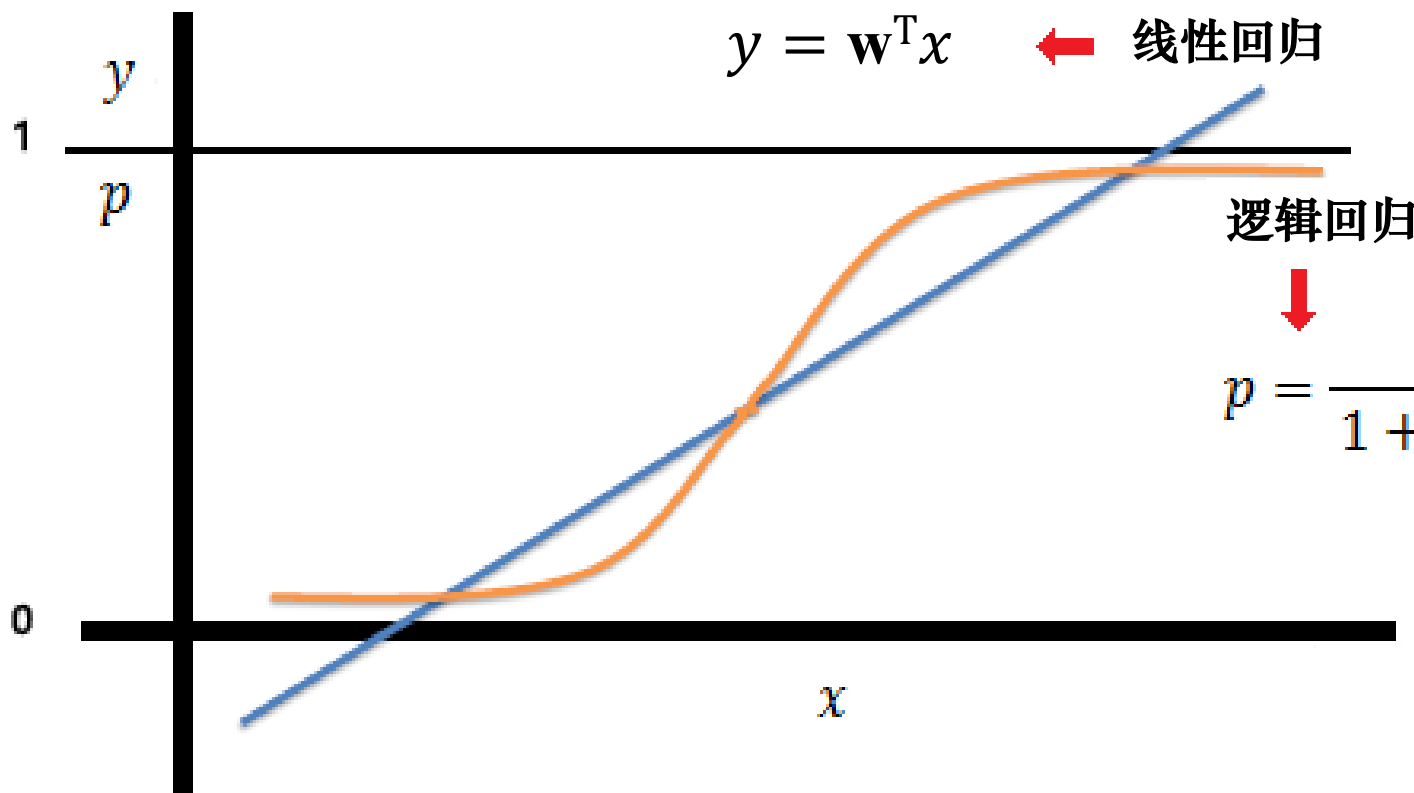


3.7 扩展与讨论

- 逻辑回归
- 二分类和多分类
- 生成式模型和判别式模型

回归和分类

- 线性回归和逻辑回归



Sigmoid函数

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

↓

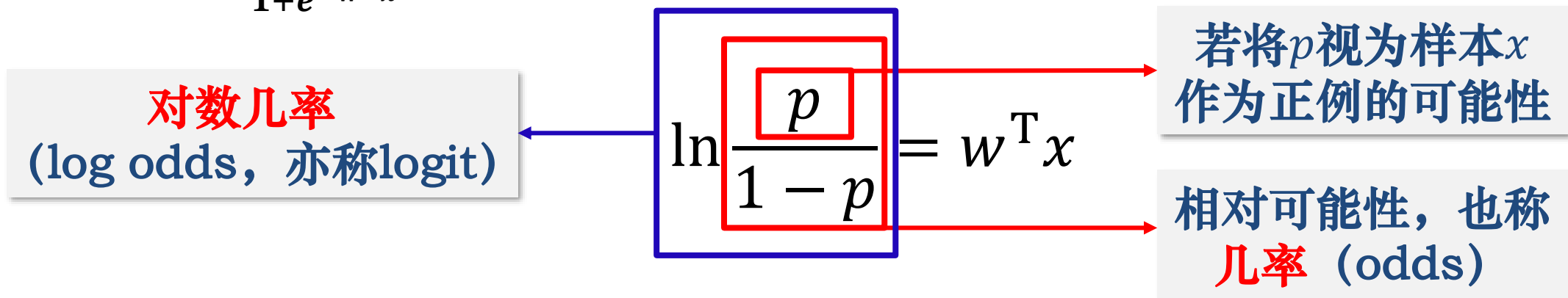
$$p = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$

逻辑回归 (Logistic Regression)

● 将回归问题转为分类问题

➤ 逻辑回归是**概率型非线性回归**，但其本质是线性回归，只是在特征到结果的映射中加入了一层函数映射，即先把特征线性求和，再使用sigmoid函数 σ 预测

➤ 将 $p = \frac{1}{1+e^{-w^T x}}$ 转换形式为

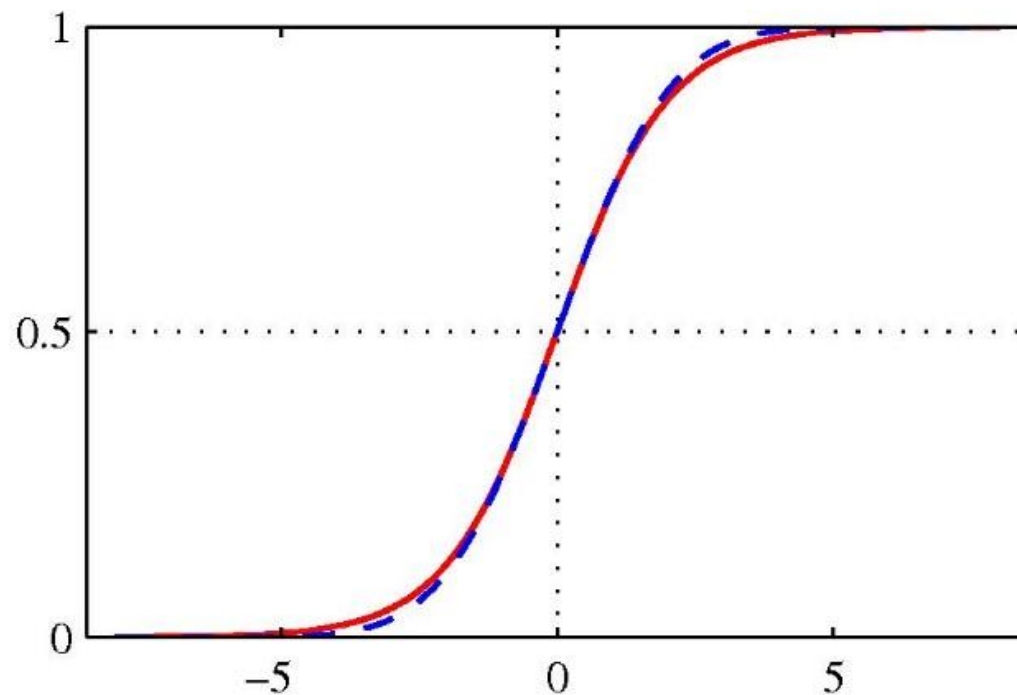


➤ 实际上逻辑回归的表达式 $p = \frac{1}{1+e^{-w^T x}}$ 是在**用线性回归模型的预测结果去逼近真实标记的对数几率**，故逻辑回归又称对数几率回归

逻辑回归 (Logistic Regression)

● Sigmoid函数

- 形似 “S”
- 可以将连续值映射到0和1区间上
- 对称性 $\sigma(-a) = 1 - \sigma(a)$
- 反转性 $a = \ln\left(\frac{\sigma}{1 - \sigma}\right)$
- 可导性 $\frac{d\sigma}{da} = \sigma(1 - \sigma)$



Logit函数

逻辑回归 (Logistic Regression)

● 逻辑回归的求解

- 将后验概率表示为作用于变量 x 的线性函数的Logistic Sigmoid

$$p(C_1|x) = y(x) = \sigma(\mathbf{w}^T x)$$

$$p(C_2|x) = 1 - p(C_1|x)$$

- 数据集 $\{x_n, t_n\}, t_n \in \{0, 1\}, n = 1, 2, \dots, N$

- 似然函数 $p(t|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}, t = (t_1, \dots, t_N)^T, y_n = p(C_1|x_n)$

- 误差函数

$$E(\mathbf{w}) = -\ln p(t|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$y_n = \sigma(a_n), a_n = \mathbf{w}^T x_n$$

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

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) x_n$$

交叉熵(Cross-Entropy)
误差函数

- 令导数等于零求解 $\nabla_{\mathbf{w}} \ln p(t|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T x_n\} x_n^T = 0$

逻辑回归 (Logistic Regression)

- 将回归问题转为分类问题

- 多类问题

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{\sum_j p(x|C_j)p(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

$$a_k = \ln p(x|C_k)p(C_k)$$

归一化指数
Normalized Exponential



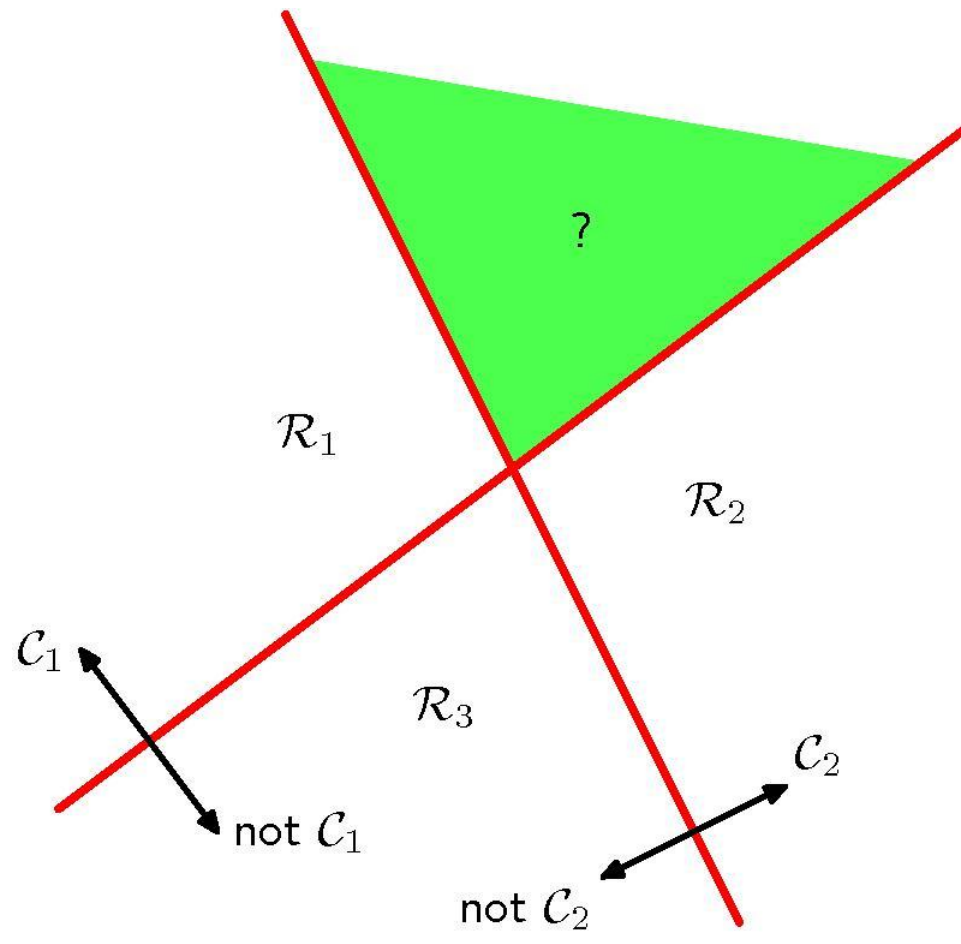
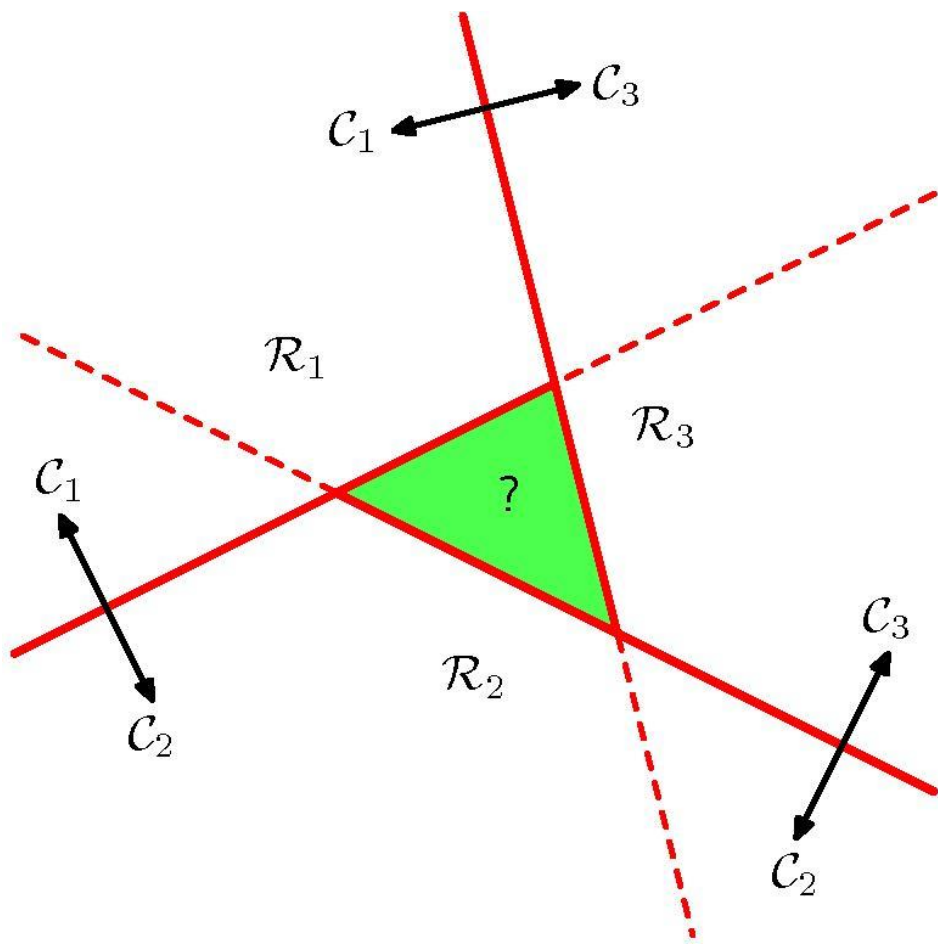
Softmax函数

是一个平滑的max函数，对于所有的 $k \neq j$ ，如果 $a_k \gg a_j$ ，有 $p(C_k|x) \simeq 1$, $p(C_j|x) \simeq 0$

Sigmoid函数可以看成是Softmax函数在 $n = 2$ 的一种特殊情形，前者经常用于二元回归/分类问题，而后者则可以应用于多元回归/分类

二分类与多分类

- 多分类问题可以被分解为多个二分类问题进行求解，常用的分解策略有一对一(1 vs. 1)和一对其余(1 vs. (N-1))



生成式模型和判别式模型

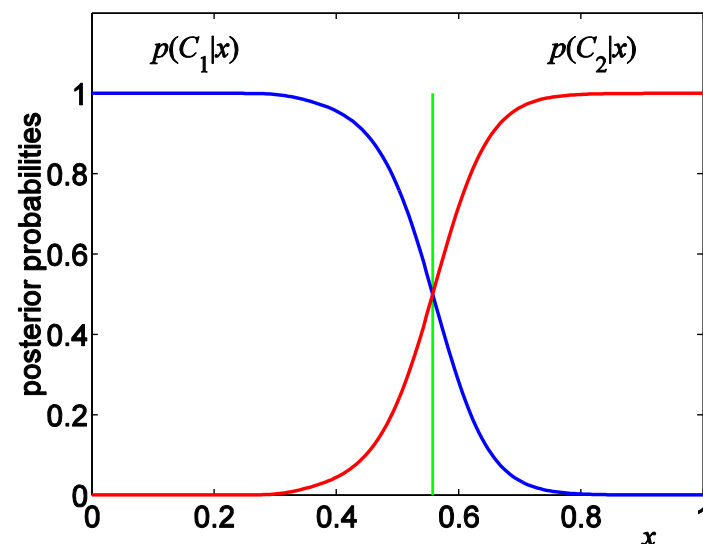
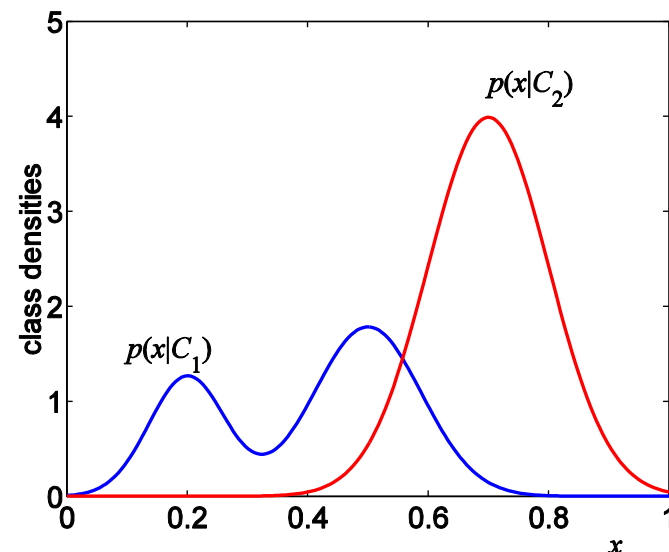
● 生成式模型 (Generative Model)

- 分别对各类的**类条件密度** $p(x|C_k)$ 和**先验概率** $p(C_k)$ 进行建模，之后利用贝叶斯定理计算**后验概率**
- 或者直接对**联合分布** $p(x, C_k)$ 建模得到后验概率

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{\sum_j p(x|C_j)p(C_j)}$$

● 判别式模型 (Discriminative Model)

- 直接对**后验概率** $p(C_k|x)$ 建模



生成式模型和判别式模型

$$p(C_1|x) = \frac{p(x|C_1)p(C_1)}{p(x|C_1)p(C_1) + p(x|C_2)p(C_2)} = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

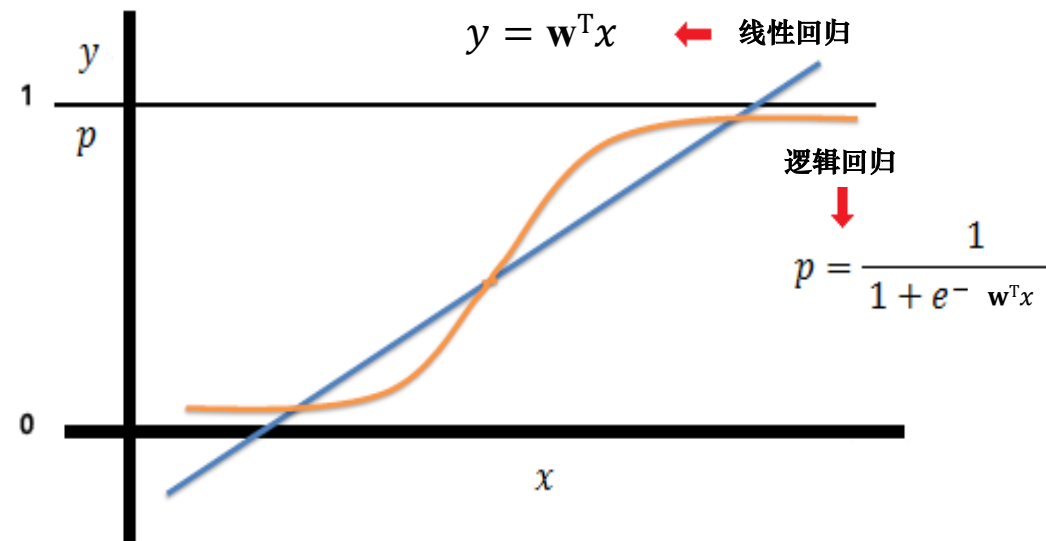
$$a = \ln \frac{p(C_1|x)}{p(C_2|x)} = \ln \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)}$$

- 生成式模型

- 估计类条件密度并计算 a

- 判别式模型

- 把 a 视为线性函数 $a = \mathbf{w}^T \mathbf{x} + w_0$ 直接估计

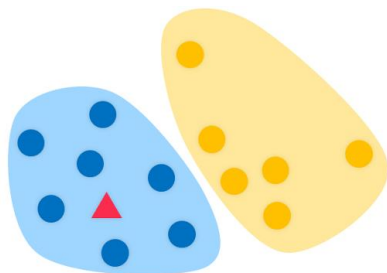


生成式模型和判别式模型

● 生成式模型

优点:

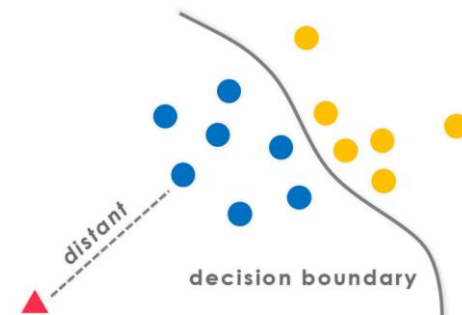
- 信息丰富
- 单类问题灵活性强
- 增量学习
- 合成缺失数据



● 判别式模型

优点:

- 类间差异清晰
- 分类边界灵活
- 学习简单
- 性能较好



缺点:

- 学习过程复杂
- 为分布牺牲分类性能

缺点:

- 不能反映数据特性
- 需要全部数据进行学习

由生成模型可以得到判别模型
但由判别模型得不到生成模型

生成式模型和判别式模型

● 生成式模型代表算法：

- Naive Bayes
- Mixtures of Gaussians
- Hidden Markov Models
- Bayesian Networks
- Deep Belief Network

● 判别式模型代表算法：

- Linear & Logistic Regression
- Support Vector Machine
- Nearest Neighbor
- Conditional Random Fields
- Boosting

