

线性拟合——有监督回归

- 如果把样本数据采用矩阵的形式记为

$$\mathbf{X} = \begin{bmatrix} \hat{\mathbf{x}}_1^T \\ \hat{\mathbf{x}}_2^T \\ \vdots \\ \hat{\mathbf{x}}_n^T \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{x}}_1^T & 1 \\ \hat{\mathbf{x}}_2^T & 1 \\ & \vdots \\ \hat{\mathbf{x}}_n^T & 1 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{f} = \begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix}$$

则一组测量数据方程组可整体写为

$$\begin{cases} f(\hat{\mathbf{x}}_1) = \hat{\mathbf{w}}^T \hat{\mathbf{x}}_1 \\ f(\hat{\mathbf{x}}_2) = \hat{\mathbf{w}}^T \hat{\mathbf{x}}_2 \\ \vdots \\ f(\hat{\mathbf{x}}_n) = \hat{\mathbf{w}}^T \hat{\mathbf{x}}_n \end{cases} \quad \longrightarrow \quad \mathbf{f} = \mathbf{X}\hat{\mathbf{w}}$$

- 构建均方误差损失函数

$$\begin{aligned} J(\hat{\mathbf{w}}) &= \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 = (\mathbf{f} - \mathbf{y})^T (\mathbf{f} - \mathbf{y}) \\ &= (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \end{aligned}$$

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

- 岭回归：在线性回归基础上引入正则化项

$$J(\hat{\mathbf{w}}) = (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) + \lambda \|\mathbf{w}\|^2$$

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

- 更一般地，对于任意单调可逆函数 g $y = g(\mathbf{w}^T \mathbf{x} + b)$

- 只要令 $\tilde{y} = g^{-1}(y) = \mathbf{w}^T \mathbf{x} + b$

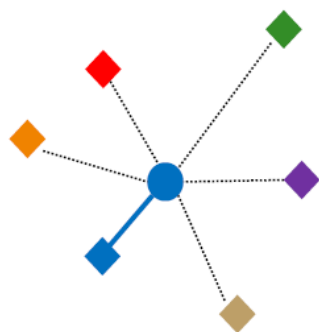
就可以使用标准的线性拟合算法（最小二乘或最大似然）进行求解。

K-means 聚类——无监督聚类

- k均值聚类目标是**最小化类内点到类中心距离

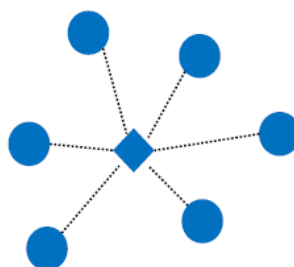
$$J = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = \sum_{\mathbf{x} \in C_1} \|\mathbf{x} - \boldsymbol{\mu}_1\|^2 + \sum_{\mathbf{x} \in C_2} \|\mathbf{x} - \boldsymbol{\mu}_2\|^2 + \dots + \sum_{\mathbf{x} \in C_k} \|\mathbf{x} - \boldsymbol{\mu}_k\|^2$$

- 算法分两步交替迭代，直到收敛
 - 类别划分：把每个样本划分到离它最近的中心点类
 - 计算均值：把每类的样本平均值作为新的类别中心



类别划分

◆ 类中心
● 样本



计算均值

$$j = \arg \min_{i=1,2,\dots,k} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$

$$\boldsymbol{\mu}_i = \sum_{\mathbf{x}_j \in C_i} \mathbf{x}_j$$

输入：样本集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$;
聚类簇数 k .

过程：

1: 从 D 中随机选择 k 个样本作为初始均值向量 $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k\}$

初始类中心

2: **repeat**

3: 令 $C_i = \emptyset$ ($1 \leq i \leq k$)

4: **for** $j = 1, \dots, m$ **do**

5: 计算样本 \mathbf{x}_j 与各均值向量 $\boldsymbol{\mu}_i$ ($1 \leq i \leq k$) 的距离: $d_{ji} = \|\mathbf{x}_j - \boldsymbol{\mu}_i\|_2$;

6: 根据距离最近的均值向量确定 \mathbf{x}_j 的簇标记: $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ji}$;

类别划分

7: 将样本 \mathbf{x}_j 划入相应的簇: $C_{\lambda_j} = C_{\lambda_j} \cup \{\mathbf{x}_j\}$;

8: **end for**

9: **for** $i = 1, \dots, k$ **do**

10: 计算新均值向量: $\boldsymbol{\mu}'_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$;

11: **if** $\boldsymbol{\mu}'_i \neq \boldsymbol{\mu}_i$ **then**

12: 将当前均值向量 $\boldsymbol{\mu}_i$ 更新为 $\boldsymbol{\mu}'_i$

计算均值

13: **else**

14: 保持当前均值向量不变

15: **end if**

16: **end for**

17: **until** 当前均值向量均未更新

18: **return** 簇划分结果

输出：簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$

高斯混合模型——无监督聚类

- 高斯混合模型(Gaussian Mixture Model, GMM): 假设分布中包括了 k 个高斯分布相叠加

$$p_{\mathcal{M}}(\mathbf{x}) = \sum_{i=1}^k \alpha_i p(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \quad 0 \leq \alpha_i \leq 1, \quad \sum_{i=1}^k \alpha_i = 1$$

- A1: 极大似然估计

$$\frac{\partial LL(\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^n \frac{\alpha_i p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{p_{\mathcal{M}}(\mathbf{x}_j)} = 0 \quad \Rightarrow \quad \boldsymbol{\mu}_i = \frac{\sum_{j=1}^n \gamma_{ji} \mathbf{x}_j}{\sum_{j=1}^n \gamma_{ji}} \quad \text{加权均值}$$

$$\boldsymbol{\Sigma}_i = \frac{\sum_{j=1}^n \gamma_{ji} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^T}{\sum_{j=1}^n \gamma_{ji}} \quad \text{加权方差} \quad \alpha_i = \sum_{j=1}^n \gamma_{ji} \quad \text{该类平均后验概率}$$

- 从而依次可求出每个高斯成分的参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)\}$, $i = 1, 2, \dots, k$
- A2: 贝叶斯定理: 给定样本 $\mathbf{x}_j \in D$, 则它属于第 i 个高斯成分的概率

$$\gamma_{ji} = \frac{\alpha_i \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{l=1}^k \alpha_l p(\mathbf{x}_j | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}$$

输入: 样本集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$;
高斯混合成分个数 k .

过程:

1: 初始化高斯混合分布的模型参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \mid 1 \leq i \leq k\}$

随机初始化

2: repeat

3: for $j = 1, \dots, m$ do

4: 根据(9.30)计算 \mathbf{x}_j 由各混合成分生成的后验概率, 即
 $\gamma_{ji} = p_{\mathcal{M}}(z_j = i | \mathbf{x}_j)$ ($1 \leq i \leq k$)

E步: 计算后验概率

5: end for

6: for $i = 1, \dots, k$ do

7: 计算新均值向量: $\boldsymbol{\mu}'_i = \frac{\sum_{j=1}^m \gamma_{ji} \mathbf{x}_j}{\sum_{j=1}^m \gamma_{ji}}$;

8: 计算新协方差矩阵: $\boldsymbol{\Sigma}'_i = \frac{\sum_{j=1}^m \gamma_{ji} (\mathbf{x}_j - \boldsymbol{\mu}'_i)(\mathbf{x}_j - \boldsymbol{\mu}'_i)^T}{\sum_{j=1}^m \gamma_{ji}}$;

M步: 求解模型参数

9: 计算新混合系数: $\alpha'_i = \frac{\sum_{j=1}^m \gamma_{ji}}{m}$;

10: end for

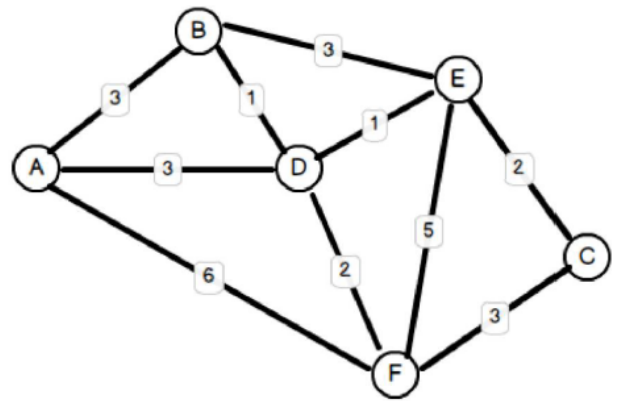
11: 将模型参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \mid 1 \leq i \leq k\}$ 更新为 $\{(\alpha'_i, \boldsymbol{\mu}'_i, \boldsymbol{\Sigma}'_i) \mid 1 \leq i \leq k\}$

12: until 满足停止条件

谱聚类——无监督聚类

- 假设有数据集 $D = \{\mathbf{x}_i\}_{i=1}^n$
- 图包括节点集(顺序无关)和对应的边集
- 图的邻接矩阵 \mathbf{W} , 表示节点的相似性

$$\mathbf{W} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 3 & 0 & 3 & 0 & 6 \\ 3 & 0 & 0 & 1 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 & 3 \\ 3 & 1 & 0 & 0 & 1 & 2 \\ 0 & 3 & 2 & 1 & 0 & 5 \\ 6 & 0 & 3 & 2 & 5 & 0 \end{bmatrix} \end{matrix}$$



- 常用全连接图(任意两定点点都相连)或 k NN连接图(只连接最近 k 顶点)
- 边的权重采用 $\mathbf{W}_{ij} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$
- 有了图的邻接矩阵 \mathbf{W}

- 图的度矩阵, 对角矩阵, 每个对角元

$$\mathbf{D}_{ii} = \sum_{j=1}^n \mathbf{W}_{ij} \quad \text{节点 } i \text{ 的度}$$

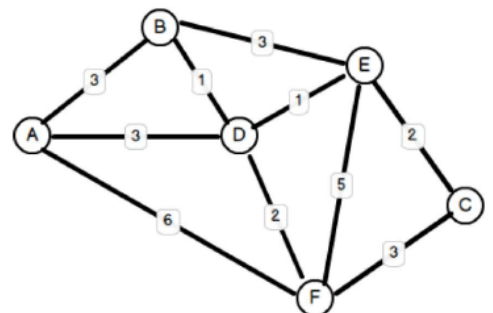
- 图的拉普拉斯矩阵, 正定矩阵

$$\mathbf{L} = \mathbf{D} - \mathbf{W}$$

- 归一化拉普拉斯矩阵, 正定矩阵

$$\mathcal{L} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$$

$$\mathbf{W} = \begin{matrix} & \begin{matrix} A & B & C & D & E & F \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} & \begin{bmatrix} 0 & 3 & 0 & 3 & 0 & 6 \\ 3 & 0 & 0 & 1 & 3 & 0 \\ 0 & 0 & 0 & 0 & 2 & 3 \\ 3 & 1 & 0 & 0 & 1 & 2 \\ 0 & 3 & 2 & 1 & 0 & 5 \\ 6 & 0 & 3 & 2 & 5 & 0 \end{bmatrix} \end{matrix}$$



- k 类情况, 求解 $\mathbf{Lz} = \lambda \mathbf{Dz}$ 最小的 k 个特征值对应的特征向量 ($\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$), 排成一个 $n \times k$ 矩阵
- 每行代表一个样本 (映射样本), 对样本使用聚类算法如 k -mean聚类

层次聚类（了解）——无监督聚类

- 层次聚类试图在不同层次对数据集进行划分，从而形成树形的聚类结构。
- 数据集划分采用“自底向上”的聚合策略，或采用“自顶向下”的分拆策略。
- 自底向上的层次聚类算法
 - a) 首先，将数据集中的**每一个样本**看做一个初始聚类簇，
 - b) 然后在每一步中找出距离最近的两个聚类簇**进行合并**
 - c) 该过程不断重复，直到达到预设的聚类簇的个数
- 自顶向下的层次聚类算法
 - a) 首先，将数据集中的**所有样本**看做一个初始聚类簇，
 - b) 然后在每一步中找出最大的聚类**进行拆分**（例如，可先找到该类中最远的两个样本，然后依照距离对该类中剩余样本进行划分）
 - c) 该过程不断重复，直到达到预设的聚类簇的个数

朴素贝叶斯——有监督分类

- 朴素贝叶斯分类器

$$c^* = \arg \max_{c \in Y} P(c) \prod_{i=1}^d p(x_i | c)$$

- 对于给定的一组样本 $D_c = \{(\mathbf{x}_i, y_i = c)\}_{i=1}^m$ ，如果是离散属性

$$P(c) = \frac{|D_c|}{|D|}, \quad P(x_i | c) = \frac{|D_{c, x_i}|}{|D_c|}$$

D_c 中第*i*个属性取值为 x_i 的样本数

如果第*i*个属性是连续属性，则利用最大似然估计每个属性的类概率密度函数 $p(x_i | c)$ ，然后可利用朴素贝叶斯分类器

线性分类器——有监督分类

- 如何把线性拟合 $z = \mathbf{w}^T \mathbf{x} + b$ 的输出转换为0-1分类？
 - 阶跃函数

$$y = \begin{cases} 0, & z < 0; \\ 0.5, & z = 0; \\ 1, & z > 0, \end{cases}$$

对数回归分类器——有监督分类

$$p(y=1|\mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{x} + b}}{1 + e^{\mathbf{w}^T \mathbf{x} + b}} \quad p(y=0|\mathbf{x}) = \frac{1}{1 + e^{\mathbf{w}^T \mathbf{x} + b}}$$

$$l(\mathbf{w}, b) = \sum_{i=1}^n \ln p(y_i = j | \mathbf{x}_i, \mathbf{w}, b) = \sum_{i=1}^n \left[y_i (\mathbf{w}^T \mathbf{x}_i + b) - \ln(1 + e^{\mathbf{w}^T \mathbf{x}_i + b}) \right]$$

$$\begin{cases} \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \lambda \Delta \mathbf{w} = \mathbf{w}^{(t)} - \lambda \frac{\partial l(\mathbf{w}, b)}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^{(t)}, b=b^{(t)}} \\ b^{(t+1)} = b^{(t)} - \lambda \Delta b = b^{(t)} - \lambda \frac{\partial l(\mathbf{w}, b)}{\partial b} \Big|_{\mathbf{w}=\mathbf{w}^{(t)}, b=b^{(t)}} \end{cases}$$

式中

$$\begin{cases} \frac{\partial l(\mathbf{w}, b)}{\partial \mathbf{w}} = - \sum_{i=1}^n [\mathbf{x}_i y_i - \mathbf{x}_i p(y_i = 1 | \mathbf{x}_i, \mathbf{w}, b)] \\ \frac{\partial l(\mathbf{w}, b)}{\partial b} = - \sum_{i=1}^n [y_i - p(y_i = 1 | \mathbf{x}_i, \mathbf{w}, b)] \end{cases}$$

支持向量机——有监督分类

- 有约束优化问题（原优化问题）

$$\min_{\mathbf{x} \in \mathbb{R}^d} f_0(\mathbf{x})$$

目标函数

$$\text{s. t. } f_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, p$$

不等式约束

$$g_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, q$$

等式约束

- 拉格朗日对偶函数

$$g(\lambda, \tau) = \inf_{\mathbf{x} \in D} L(\mathbf{x}, \lambda, \tau) = \left(f_0(\mathbf{x}) + \lambda^T \mathbf{F}(\mathbf{x}) + \tau^T \mathbf{G}(\mathbf{x}) \right) \Big|_{\mathbf{x}^* \leftarrow \nabla L_{\mathbf{x}}(\mathbf{x}^*, \lambda, \tau) = 0}$$

- 对偶优化问题

$$\max g(\lambda, \tau)$$

$$\text{s.t. } \lambda_i \geq 0, \quad i = 1, \dots, p$$

- KKT条件：最优解的必要条件(凸优化时也是充分必要条件)

- 原始问题约束 $f_i(\mathbf{x}) \leq 0, \quad i = 1, 2, \dots, p; \quad g_j(\mathbf{x}) = 0, \quad j = 1, 2, \dots, q$
- 对偶问题约束 $\lambda_i \geq 0, \quad i = 1, \dots, p$
- 互补松弛 $\lambda_i f_i(\mathbf{x}) = 0, \quad i = 1, \dots, p$
- 梯度消失 $\nabla L_{\mathbf{x}}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\tau}) = 0$

- 最大间隔: 寻找参数 \mathbf{w} 和 b , 使得 γ 最大, 同时满足正确分类的约束条件.

$$\begin{aligned} \arg \max_{\mathbf{w}, b} \quad & \frac{2}{\|\mathbf{w}\|} \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, i = 1, 2, \dots, n \end{aligned}$$

- 第一步：引入拉格朗日乘子 $\alpha_i \geq 0$ 得到拉格朗日函数

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i (y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1)$$

- 第二步：令 $L(\mathbf{w}, b, \boldsymbol{\alpha})$ 对 \mathbf{w} 和 b 的偏导为零可得

$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i, \quad \sum_{i=1}^n \alpha_i y_i = 0$$

- 第三步：回代到 $L(\mathbf{w}, b, \boldsymbol{\alpha})$ 中得到拉格朗日对偶函数, 从而可得对偶优化问题

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^n \alpha_i \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, n \end{aligned}$$

对偶问题

- 基本思路：不断执行如下两个步骤直至收敛。
 - 第一步：选取一对需更新的变量 α_i 和 α_j 。
 - 第二步：固定 α_i 和 α_j 以外的参数，求解对偶问题更新 α_i 和 α_j 。
- 上面第二步：仅考虑 α_i 和 α_j 时，对偶问题的约束变为

$$\alpha_i y_i + \alpha_j y_j = - \sum_{k \neq i, j} \alpha_k y_k, \quad \alpha_i \geq 0, \quad \alpha_j \geq 0$$

用一个变量表示另一个变量，再代入对偶问题中可得一个单变量的二次规划，该问题具有闭式解。（舍弃不符合约束条件的负数解）

- 偏移项 b ：通过支持向量满足的方程（蓝线）来确定 $y_i f(\mathbf{x}_i) = 1$
- 最终判决 $y = \text{sign}[f(\mathbf{x}_i)]$ ， $\text{sign}[x]$ 表示取 x 符号

- 引入松弛变量

$$\begin{aligned} \min_{\mathbf{w}, b, \xi} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i & \min_{\alpha} \quad & \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^n \alpha_i \\ \text{s.t.} \quad & y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i & \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad i = 1, 2, \dots, n \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, n \end{aligned}$$

- 设样本 \mathbf{x} 映射后的向量为 $\phi(\mathbf{x})$ ，划分超平面为 $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$

$$\text{原始问题} \begin{cases} \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, 2, \dots, n \end{cases}$$

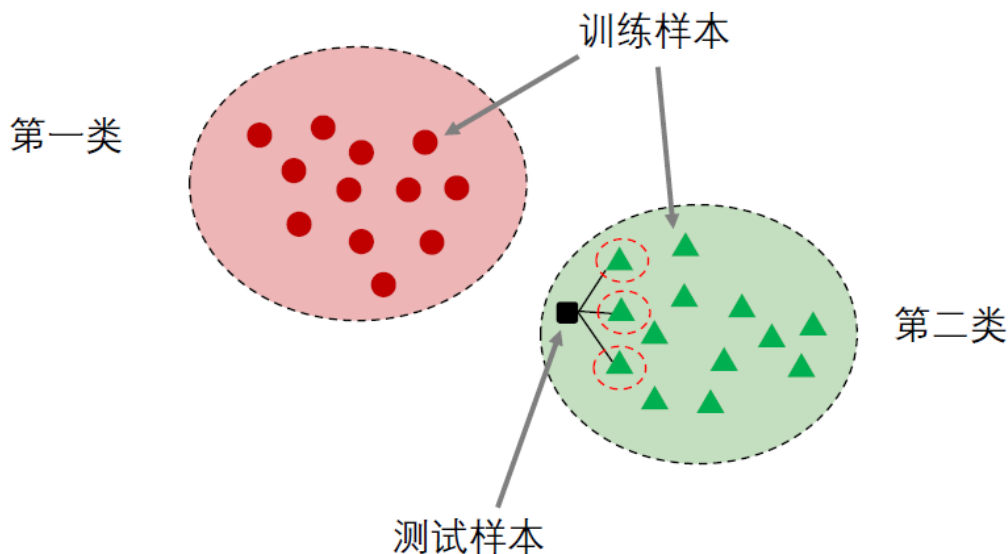
$$\text{对偶问题} \begin{cases} \min_{\alpha} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) - \sum_{i=1}^n \alpha_i \\ \text{s.t.} \quad \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad i = 1, 2, \dots, n \end{cases}$$

$$\text{预测} \quad f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b = \sum_{j=1}^n \alpha_j y_j \phi(\mathbf{x}_j)^T \phi(\mathbf{x}) + b$$

只以内积的形式出现，不需计算出映射，只要设计核函数 $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$

kNN——有监督分类

- $k > 1$ 时，投票决定， k 个样本中得票最多的类作为最终类



- Given:

- training examples $\{x_i, y_i\}$

- x_i ... attribute-value representation of examples
- y_i ... class label: {ham,spam}, digit {0,1,...9} etc.

- testing point x that we want to classify

- Algorithm:

- compute distance $D(x, x_i)$ to every training example x_i
- select k closest instances $x_{i1} \dots x_{ik}$ and their labels $y_{i1} \dots y_{ik}$
- output the class y^* which is most frequent in $y_{i1} \dots y_{ik}$

- 用 $\mathcal{S} = \{(\mathbf{x}_i, \mathbf{x}_j) \mid y_i = y_j\}$ 表示所有同类样本对， $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{x}_j) \mid y_i \neq y_j\}$ 表示所有异类样本对，则度量学习目标如下

$$\max \quad g(\mathbf{A}) = \sum_{(\mathbf{x}_l, \mathbf{x}_m) \in \mathcal{D}} \|\mathbf{x}_l - \mathbf{x}_m\|_{\mathbf{A}}$$

$$\text{s. t.} \quad f(\mathbf{A}) = \sum_{(\mathbf{x}_i, \mathbf{x}_j) \in \mathcal{S}} \|\mathbf{x}_i - \mathbf{x}_j\|_{\mathbf{A}}^2 \leq 1 \quad \leftarrow C_1$$

Iterate

Iterate

$$A := \arg \min_{A'} \{\|A' - A\|_F : A' \in C_1\}$$

$$A := \arg \min_{A'} \{\|A' - A\|_F : A' \in C_2\}$$

until A converges

$$A := A + \alpha(\nabla_A g(A)) \perp \nabla_A f$$

until convergence

决策树——有监督分类

Function *TreeGenerate* (输入数据集 $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, 属性集 A)

- Step 1: 从 A 中选择一个**最优属性** a^* , 创建节点 a^*
- Step 2: 按 a^* 相同取值 (离散型)或区间(连续型)把 D 划分为若干互斥子集 D_1, D_2, \dots, D_k .
- Step 3: 从 A 中移除 $a^* : A_a = A - \{a^*\}$, 对每个子集 D_k 递归调用自身 *TreeGenerate* (输入数据集 D_k , 属性集 A_a)

输出: 以为 a^* 根节点的一棵决策树

Algorithm 1 决策树学习基本算法

输入:

- 训练集 $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$;
- 属性集 $A = \{a_1, \dots, a_d\}$.

过程: 函数 *TreeGenerate*(D, A)

```
1: 生成结点 node;
2: if  $D$  中样本全属于同一类别  $C$  then
3:   将 node 标记为  $C$  类叶结点; return (1)
4: end if
5: if  $A = \emptyset$  OR  $D$  中样本在  $A$  上取值相同 then
6:   将 node 标记叶结点, 其类别标记为  $D$  中样本数最多的类; return (2)
7: end if
8: 从  $A$  中选择最优划分属性  $a_*$ ;
9: for  $a_*$  的每一个值  $a_*^v$  do
10:  为 node 生成每一个分枝; 令  $D_v$  表示  $D$  中在  $a_*$  上取值为  $a_*^v$  的样本子集;
11:  如果  $D_v$  为空 then
12:    将分枝结点标记为叶结点, 其类别标记为  $D$  中样本最多的类; return (3)
13:  else
14:    以 TreeGenerate( $D_v, A - \{a_*\}$ ) 为分枝结点
15:  end if
16: end for
```

输出: 以 node 为根结点的一棵决策树

三种情况下递归回终止

(1) 当前结点包含的样本全部属于同一类别

(2) 当前属性集为空, 或所有样本在所有属性上取值相同

(3) 当前结点包含的样本集合为空

- 数据集 $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ 包含 k 个类样本, 每类样本所占比例 p_k , 则数据集 D 的信息熵定义为

$$\text{Ent}(D) = - \sum_{k=1}^{|Y|} p_k \log_2 p_k$$

- 数据集越纯, $\text{Ent}(D)$ 越小

$$\text{Gain}(D, a) = \text{Ent}(D) - \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

- 连续属性时, 可选取时 $\text{Gain}(D, a)$ 最大化属性值作为划分点

$$\text{Gain_ratio}(D, a) = \frac{\text{Gain}(D, a)}{\text{IV}(a)}$$

- 其中

$$\text{IV}(a) = - \sum_{v=1}^V \frac{|D^v|}{|D|} \log_2 \frac{|D^v|}{|D|}$$

- 数据集 D 的纯度除了用信息熵外，还可以用基尼值衡量

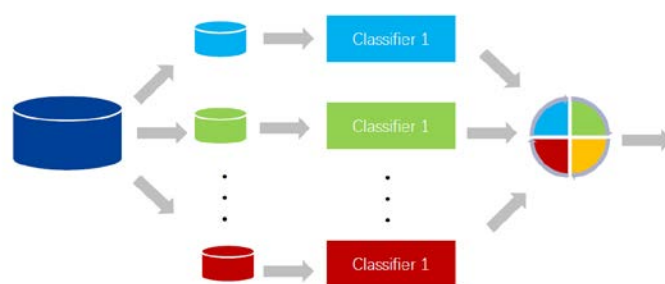
$$\text{Gini}(D) = \sum_{k=1}^{|Y|} \sum_{k' \neq k} p_k p_{k'} = 1 - \sum_{k=1}^{|Y|} p_k^2$$

- 基尼值越小，数据集纯度越高。
- 属性 a 的基尼指数定义为：

$$\text{Gini_index}(D, a) = \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Gini}(D^v)$$

- **最优属性**是应选择那个使划分后**基尼指数最小**的属性作为最优划分属性
- 预剪枝：构建树的过程中，计算划分前和划分后验证精度，如有提升，则保留该分支；否则不进行划分
- 后剪枝：构建树后，从底往上逐个考察中间节点，如果删除后验证精度有下降，则保留；否则就删除该中间节点

随机森林——有监督分类



- 给定数据 $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ，随机森林包括三步：
 - **Bagging**：随机可重复地从 D 中选取 n 个样板组成一个训练集，此过程重复 T 次获得 T 个训练集 D_1, D_2, \dots, D_T
 - **训练**：每个训练集上学习一个“随机决策树”，即每个节点从随机选的 k 个属性中选取一个最优的属性用于划分
 - **预测**：对新样本，用 T 个决策树预测的结果投票(分类)或平均(回归)决定最终结果

神经网络——有监督分类

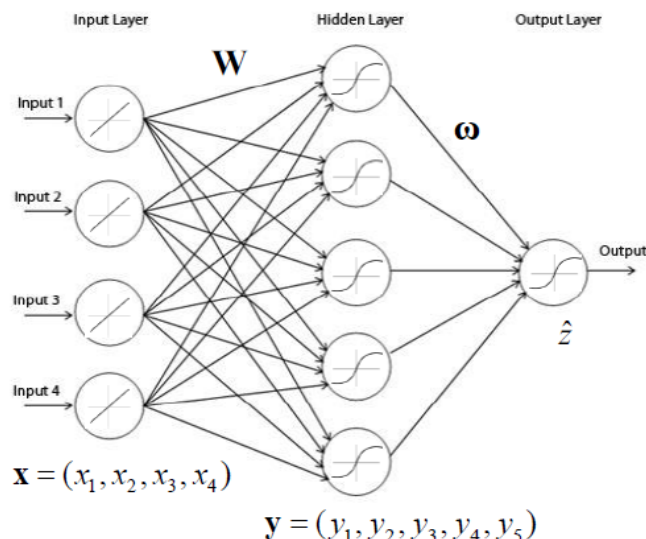
$$\begin{cases} \hat{z} = f\left(\sum_{i=1}^5 \omega_i y_i + b_o\right) \\ y_i = f\left(\sum_{j=1}^4 w_{ji} x_j + b_i\right), i=1..5 \end{cases}$$

$$\boldsymbol{\omega}^T = (\omega_1, \omega_2, \dots, \omega_4)$$

$$\mathbf{y}^T = (y_1, y_2, \dots, y_5)$$

$$\mathbf{b}^T = (b_1, b_2, \dots, b_5)$$

$$\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_5)$$



$$\begin{cases} \hat{z} = f(\boldsymbol{\omega}^T \mathbf{y} + b_o) & \text{网络输出} \\ \mathbf{y} = f(\mathbf{W}^T \mathbf{x} + \mathbf{b}) & \text{隐含层输出 (函数 } f \text{ 元素运算函数)} \end{cases} \begin{cases} \mathbf{y} = f(\mathbf{x}) \\ y_i = f(x_i) \end{cases}$$

	输入：训练数据集 $(\mathbf{x}_1, z_1), (\mathbf{x}_2, z_2), \dots, (\mathbf{x}_n, z_n)$ ，学习步长 λ
1	随机初始化 $\boldsymbol{\omega}, b_o, \mathbf{W}, \mathbf{b}$ ，设置计数器 $t = 0$
2	随机选择训练数据 (\mathbf{x}_i, z_i)
3	前向计算 $\mathbf{y} = f(\mathbf{W}^T \mathbf{x}_i + \mathbf{b})$ ， $\hat{z}_i = f(\boldsymbol{\omega}^T \mathbf{y} + b_o)$ ， $J(\mathbf{x}_i; z_i) = \frac{1}{2}(\hat{z}_i - z_i)^2$
4	反向计算 $\Delta \boldsymbol{\omega}, \Delta b_o, \Delta \mathbf{W}, \Delta \mathbf{b}$
5	变量更新 $\begin{cases} \boldsymbol{\omega}^{(t+1)} = \boldsymbol{\omega}^{(t)} - \lambda \Delta \boldsymbol{\omega} \\ b_o^{(t+1)} = b_o - \lambda \Delta b_o \end{cases}, \begin{cases} \mathbf{W}^{(t+1)} = \mathbf{W}^{(t)} - \lambda \Delta \mathbf{W} \\ \mathbf{b}^{(t+1)} = \mathbf{b} - \lambda \Delta \mathbf{b} \end{cases}$
6	如未满足终止条件， $t = t+1$ ，重复第2步

- 交叉熵损失 (Cross Entropy Loss): $J(\mathcal{Z}, \hat{\mathcal{Z}}) = \sum_{i=1}^n H(\mathbf{z}_i, \hat{\mathbf{z}}_i)$

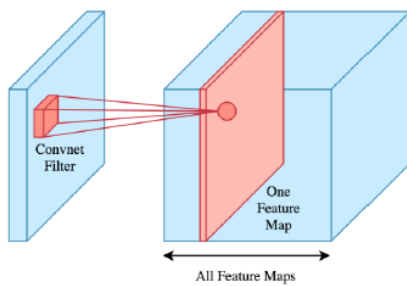
其中交叉熵 $H(q, p)$ 衡量概率分布 q, p 的差异性

$$H(\mathbf{q}, \mathbf{p}) = -\sum_{j=1}^c q_j \log p_j \quad H(q, p) = \int q(x) \log p(x) dx$$

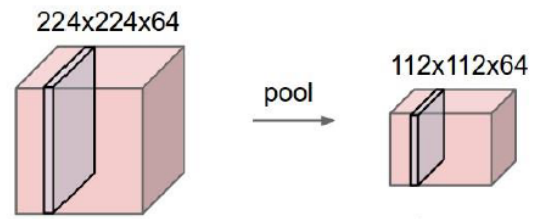
- 如何把网络输出输出 \hat{z}_i 转化为类别概率分布？

$$P(j | \hat{\mathbf{z}}) = \frac{e^{\hat{z}_j}}{\sum_{k=1}^c e^{\hat{z}_k}} \quad \sum_{j=1}^c P(j | \hat{\mathbf{z}}) = \sum_{j=1}^c \frac{e^{\hat{z}_j}}{\sum_{k=1}^c e^{\hat{z}_k}} = \frac{1}{\sum_{k=1}^c e^{\hat{z}_k}} \sum_{j=1}^c e^{\hat{z}_j} = 1$$

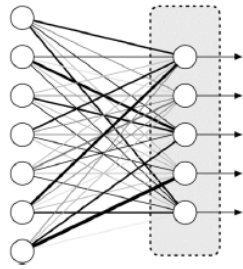
- 只要给定输入输出和上游导数，就可以求下游导数，即梯度反向传播



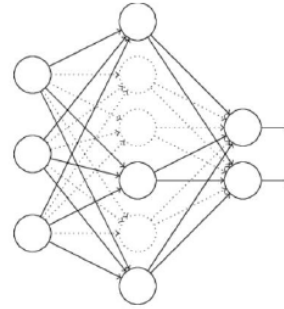
卷积层(convolutional layer)



聚集层(Pooling layer)



全连接层(fully connected layer)



丢弃层(dropout layer)

线性判别分析 LDA——无监督降维

- 数据集 $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n, y \in \{0, 1\}$ ，二分类问题
- 投影前，每类的均值和协方差矩阵

$$\begin{cases} \mathbf{u}_0 = \frac{1}{n_0} \sum_{y_i=0} \mathbf{x}_i, & \boldsymbol{\Sigma}_0 = \frac{1}{n_0 - 1} \sum_{y_i=0} \mathbf{x}_i \mathbf{x}_i^T \\ \mathbf{u}_1 = \frac{1}{n_1} \sum_{y_i=1} \mathbf{x}_i, & \boldsymbol{\Sigma}_1 = \frac{1}{n_1 - 1} \sum_{y_i=1} \mathbf{x}_i \mathbf{x}_i^T \end{cases}$$

- 投影后，每类的均值和协方差（投影到**直线**，因此均值协方差都是实数）

$$\begin{cases} \hat{u}_0 = \mathbf{w}^T \mathbf{u}_0, & \hat{\Sigma}_0 = \mathbf{w}^T \boldsymbol{\Sigma}_0 \mathbf{w} \\ \hat{u}_1 = \mathbf{w}^T \mathbf{u}_1, & \hat{\Sigma}_1 = \mathbf{w}^T \boldsymbol{\Sigma}_1 \mathbf{w} \end{cases}$$

- 定义类内离散度矩阵和类间离散度矩阵

$$\mathbf{S}_w = \boldsymbol{\Sigma}_0 + \boldsymbol{\Sigma}_1 \quad \mathbf{S}_b = (\mathbf{u}_0 - \mathbf{u}_1)(\mathbf{u}_0 - \mathbf{u}_1)^T$$

- 则目标函数简化为 $J = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$

$$\mathbf{w}^* = \mathbf{S}_w^{-1} (\mathbf{u}_0 - \mathbf{u}_1)$$

多维缩放 MDS——无监督降维

- 给定 d 维空间中样本距离矩阵 $\mathbf{D} \in \mathbb{R}^{n \times n}$, \mathbf{D}_{ij} 表示样本 $\mathbf{x}_i, \mathbf{x}_j$ 之间的距离, 通常假设是欧氏距离 $\mathbf{D}_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2$
- 假设降维后对应的样本为 $\{\mathbf{z}_i\}_{i=1}^n \in \mathbb{R}^{d'}$ ($d' \leq d$)
- 求解方法二**: 不直接求解低维嵌入坐标 $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$, 转而求它们的内积矩阵 $\mathbf{B} = \mathbf{Z}^T \mathbf{Z}$, 其中 $\mathbf{Z} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n)$, 并且 $\mathbf{B}_{ij} = \mathbf{z}_i^T \mathbf{z}_j$
- 设低维嵌入坐标 $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$ 均值为0: $\frac{1}{n} \sum_{i=1}^n \mathbf{z}_i = 0$

$$\mathbf{B}_{ij} = \frac{1}{2n} \left(\mathbf{D}_{i\cdot} + \mathbf{D}_{\cdot j} - n \mathbf{D}_{ij}^2 - \frac{1}{n} \mathbf{D}_{\cdot\cdot} \right)$$

- 因此给定 $\mathbf{D} \in \mathbb{R}^{n \times n}$ 后, 矩阵 \mathbf{B} 可求出。在利用特征分解可求出 \mathbf{Z}

输入: 距离矩阵 $\mathbf{D} \in \mathbb{R}^{m \times m}$, 其元素 $dist_{ij}$ 为样本 \mathbf{x}_i 到 \mathbf{x}_j 的距离;
低维空间维数 d' .

过程:

- 根据式(10.7)–(10.9)计算 $dist_{i\cdot}^2, dist_{\cdot j}^2, dist_{\cdot\cdot}^2$;
- 根据式(10.10)计算矩阵 \mathbf{B} ;
- 对矩阵 \mathbf{B} 做特征值分解;
- 取 $\tilde{\mathbf{\Lambda}}$ 为 d' 个最大特征值所构成的对角矩阵, $\tilde{\mathbf{V}}$ 为相应的特征向量矩阵.

输出: 矩阵 $\tilde{\mathbf{V}} \tilde{\mathbf{\Lambda}}^{1/2} \in \mathbb{R}^{m \times d'}$, 每行是一个样本的低维坐标

主成分分析 PCA——无监督降维

Function	PCA(数据集D, 主成分数 k)
1	去中心化 $\mathbf{m} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$; $\hat{\mathbf{x}}_i = \mathbf{x}_i - \mathbf{m}, i = 1, \dots, n$
2	求协方差矩阵 $\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T$
3	特征分解 $\mathbf{C}\mathbf{w} = \lambda\mathbf{w}$ 求出前 k 个特征对 $(\lambda_i, \mathbf{w}_i), i = 1 \dots k$
4	计算前k个主成分 $\mathbf{Y} = \mathbf{W}^T \mathbf{X}, \mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k)$
5	输出 \mathbf{Y}, \mathbf{W}

核化主成分分析 KPCA（了解）——无监督降维

- 映射变换 $\mathbf{z} = \phi(\mathbf{x})$ 后，在特征空间中寻找主成分

$$K_{ij} := (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)) \quad \mathbf{K}\mathbf{a} = \lambda\mathbf{a}$$

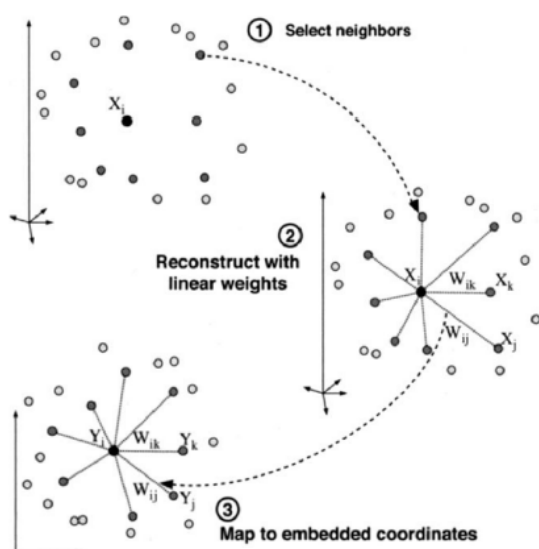
等度量映射 Isomap（了解）——无监督降维

- 计算图上任意两点间的最小图距离 $g_{ij}(\mathbf{x}_i, \mathbf{x}_j)$
- 对图距离应用经典的多维度缩放MDS算法

$$\min E(Y) = \min \sum_{i,j} (\|y_i - y_j\| - d_g(x_i, x_j))^2$$

局部线性嵌入 LLE（了解）——无监督降维

Step 2:



$$\phi(W) = \sum_i \left| \vec{X}_i - \sum_j W_{ij} \vec{X}_j \right|^2$$

$$s.t. \quad \sum_j W_{ij} = 1$$

Step 3:

$$\Phi(Y) = \sum_i \left| \vec{Y}_i - \sum_j W_{ij} \vec{Y}_j \right|^2$$

$$s.t. \quad \sum_i \vec{Y}_i = 0; \quad \frac{1}{N} \sum_i \vec{Y}_i \otimes \vec{Y}_i = I$$