

模式识别第二次作业提交版

201300086 史浩男 人工智能学院

一、 (3.2) K-means

(a)公式抽象

不妨假设 μ_i 为聚类中心，我们的目标就是把数据点根据到中心距离分类，形式化目标函数如下：

$$D = \sum_{j=1}^M \sum_{i=1}^K \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \quad (1)$$

因此只需要最小化这个目标函数，就能求出对应的 γ_{ij} 和 μ_i ，从而完成聚类

$$\arg \min_{\gamma_{ij}, \mu_i} \sum_{i=1}^K \sum_{j=1}^M \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \quad (2)$$

(b)迭代规则

固定 μ_i ：

此时只需找到每个数据点距离最近的中心是哪个，所有中心都是不变的，因此表达式为：

$$\gamma_{ij} = \begin{cases} 1, & \text{if } i = \arg \min_i \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

固定 γ_{ij} ：

此时所有类别包含哪些点已经确定，只需在每个类中找到类中点最近的中心位置，可以解出：

$$\boldsymbol{\mu}_i = \frac{\sum_{j=1}^M \gamma_{ij} \mathbf{x}_j}{\sum_{j=1}^M \gamma_{ij}} \quad (4)$$

(c)证明收敛

只需证明，（b）中的两个迭代步骤，都会使目标函数D不增（单调递减有下界的函数必然收敛。而如果两个步骤都不增不减，说明已经收敛。如果至少有一个是递减的，那么满足条件单调递减有下界，一定会最终收敛）

固定 μ_i ：

$$\begin{aligned} D' - D &= \sum_{i=1}^K \sum_{j=1}^M \gamma'_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 - \sum_{i=1}^K \sum_{j=1}^M \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \\ &= \sum_{i=1}^K \sum_{j=1}^M (\gamma'_{ij} - \gamma_{ij}) \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \\ &= \sum_{i=1}^K \left(\|\mathbf{x}_j - \boldsymbol{\mu}'_i\|^2 - \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \right) \\ &\leq 0 \end{aligned} \quad (5)$$

固定 γ_{ij} :

$$\begin{aligned}
D'_j - D_j &= \sum_{i=1}^K \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}'_i\|^2 - \sum_{i=1}^K \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \\
&= \sum_{i=1}^K \gamma_{ij} \left(\|\mathbf{x}_j - \boldsymbol{\mu}'_i\|^2 - \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \right) \\
&= \sum_{i=1}^K \gamma_{ij} (\mathbf{x}_j - \boldsymbol{\mu}'_i + \mathbf{x}_j - \boldsymbol{\mu}_i)^T (\mathbf{x}_j - \boldsymbol{\mu}'_i - \mathbf{x}_j + \boldsymbol{\mu}_i) \\
&= \sum_{i=1}^K \gamma_{ij} (\boldsymbol{\mu}_i - \boldsymbol{\mu}'_i)^T (2\mathbf{x}_j - \boldsymbol{\mu}'_i - \boldsymbol{\mu}_i) \\
&= (\boldsymbol{\mu}_i - \boldsymbol{\mu}'_i)^T \left(2 \left(\sum_{i=1}^K \gamma_{ij} \mathbf{x}_j \right) - \left(\sum_{i=1}^K \gamma_{ij} \right) (\boldsymbol{\mu}'_i + \boldsymbol{\mu}_i) \right) \\
&= \left(\sum_{i=1}^K \gamma_{ij} \right) (\boldsymbol{\mu}_i - \boldsymbol{\mu}'_i)^T \left(2 \frac{\sum_{i=1}^K \gamma_{ij} \mathbf{x}_j}{\sum_{i=1}^K \gamma_{ij}} - \boldsymbol{\mu}'_i - \boldsymbol{\mu}_i \right) \\
&= \left(\sum_{i=1}^K \gamma_{ij} \right) (\boldsymbol{\mu}_i - \boldsymbol{\mu}'_i)^T (2\boldsymbol{\mu}'_i - \boldsymbol{\mu}'_i - \boldsymbol{\mu}_i) \\
&\leq 0
\end{aligned} \tag{6}$$

因此我们简洁地证明了，两个迭代步骤都使目标函数D不增。

综上，一定会最终收敛

二、 (4.2) LR

(a, b)优化问题与重写

$$\arg \min_{\boldsymbol{\beta}} (y_i - x_i^T \boldsymbol{\beta})^2 \tag{7}$$

矩阵表示重写：

$$\arg \min_{\boldsymbol{\beta}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \tag{8}$$

(c)求解

$$\frac{\partial E}{\partial \boldsymbol{\beta}} = 2\mathbf{X}^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) = \mathbf{0} \tag{9}$$

由于 $\mathbf{X}^T \mathbf{X}$ 可逆，解出：

$$\boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \tag{10}$$

(d)维度大于样本导致不可解

由矩阵的性质可知： $\text{rank}(\mathbf{X}^T \mathbf{X}) = \text{rank}(\mathbf{X}) \leq n < d$ ，而 $\mathbf{X}^T \mathbf{X}$ 是一个 $d \times d$ 的矩阵，不满秩的矩阵必然不可逆

(e)正则化项作用

正则化项度量了模型复杂度，是用于对抗过拟合的关键手段。正则化表示了对模型的一种偏好, 可以对模型的复杂度进行约束, 因此可以在性能相同的模型中，选择出模型复杂度最低的一个。

(f)求解岭回归优化问题

$$\arg \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta$$
$$\frac{\partial E}{\partial \beta} = 2\mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}) + 2\lambda \beta = \mathbf{0} \quad (11)$$

解得最优

$$\beta^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (12)$$

(g)正则化在可逆方面的作用

加入岭回归正则项后, $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ 几乎总是可逆, 总能求解, 解决了无法获得唯一的模型参数的问题。

同时正则化是用于对抗过拟合的关键手段

(h)极端 λ 的影响

如果 $\lambda = 0$, 岭回归退化为普通线性回归

如果 $\lambda = \infty$, 则优化问题变为 $\arg \min_{\beta} \beta^T \beta$, 解出 $\beta = \mathbf{0}$

(i) λ 为什么必须是超参数

因为正则化项 $\lambda \beta^T \beta$ 恒正, 目标函数中只有正则化项中出现了 λ , 最优化目标函数时一定会将 λ 优化为0, 失去了正则化的意义

三、 (4.5) AUC

(a)

下标	标记	得分	P	R	AUC-PR	AP
0			1	0		
1	1	1	1	0.2	0.2	0.2
2	2	0.9	0.5	0.2	0	0
3	1	0.8	0.67	0.4	0.1167	0.1333
4	1	0.7	0.75	0.6	0.1417	0.15
5	2	0.6	0.6	0.6	0	0
6	1	0.5	0.67	0.8	0.1267	0.1333

下标	标记	得分	P	R	AUC-PR	AP
7	2	0.4	0.57	0.8	0	0
8	2	0.3	0.5	0.8	0	0
9	1	0.2	0.56	1	0.1056	0.111
10	2	0.1	0.5	1	0	0

(b)AP&PR

相似是正常的，而且AP比PR总是稍微大一点点

原因是他们的计算方式只有细微区别：

$$\begin{aligned}
 AP - PR &= \sum_{i=1}^n (r_i - r_{i-1}) p_i - \sum_{i=1}^n (r_i - r_{i-1}) \frac{p_i + p_{i-1}}{2} \\
 &= \sum_{i=1}^n \frac{1}{2} (r_i - r_{i-1}) (p_i - p_{i-1})
 \end{aligned} \tag{13}$$

(c)

交换了第 9 行和第 10 行的类别标记之后, AUC-PR=0.6794, AP=0.7167.

(d)

代码:

```

from collections import Counter

v = [1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
label = [1, 2, 1, 1, 2, 1, 2, 2, 1, 2]#[1, 2, 1, 1, 2, 1, 2, 2, 1, 2]

P = [1.0]
R = [0.0]
TPR = [0.0]
FPR = [0.0]
for i in range(1, len(v) + 1):
    pos_count = Counter(label[:i])
    neg_count = Counter(label[i:])
    TP = pos_count.get(1, 0)
    FP = pos_count.get(2, 0)
    FN = neg_count.get(1, 0)
    TN = neg_count.get(2, 0)
    P.append(TP / (TP + FP))
    R.append(TP / (TP + FN))
AUC_PR = [0.5 * (R[i] - R[i - 1]) * (P[i] + P[i - 1]) for i in range(1, len(R))]
AP = [(R[i] - R[i - 1]) * P[i] for i in range(1, len(R))]

print('P:', [*P])

```

```
print('R:', [*R])
print('AUC_PR:', [*AUC_PR])
print('AP:', [*AP])
```

四、 (4.6) KNN

(a)偏置-方差分解

首先给出误差表达式

$$\mathbb{E}_D[(y - f(\mathbf{x}; D))] = \mathbb{E}_D[(F(\mathbf{x}) - f(\mathbf{x}; D) + \epsilon)^2] = \mathbb{E}_D[(F(\mathbf{x}) - f(\mathbf{x}; D))^2] + \sigma^2 \quad (14)$$

展开可得

$$\mathbb{E}_D[(F(\mathbf{x}) - f(\mathbf{x}; D))^2] = (\mathbb{E}_D[F(\mathbf{x}) - f(\mathbf{x}; D)])^2 + \text{Var}(F(\mathbf{x}) - f(\mathbf{x}; D)) \quad (15)$$

由于 $F(\mathbf{x})$ 是确定的, 与训练集 D 无关, 即 $\mathbb{E}_D[F(\mathbf{x})] = F(\mathbf{x})$, 则上式进一步简化为:

$$\begin{aligned} & (\mathbb{E}_D[F(\mathbf{x}) - f(\mathbf{x}; D)])^2 + \text{Var}(F(\mathbf{x}) - f(\mathbf{x}; D)) \\ &= (F(\mathbf{x}) - \mathbb{E}_D[f(\mathbf{x}; D)])^2 + \text{Var}(f(\mathbf{x}; D)) \\ &= (F(\mathbf{x}) - \mathbb{E}_D[f(\mathbf{x}; D)])^2 + \mathbb{E}_D[(f(\mathbf{x}; D) - \mathbb{E}_D[f(\mathbf{x}; D)])^2] \end{aligned} \quad (16)$$

综上,得到偏置-方差分解

$$\begin{aligned} & \mathbb{E}_D[(y - f(\mathbf{x}; D))] \\ &= (F(\mathbf{x}) - \mathbb{E}_D[f(\mathbf{x}; D)])^2 + \mathbb{E}_D[(f(\mathbf{x}; D) - \mathbb{E}_D[f(\mathbf{x}; D)])^2] + \sigma^2 \end{aligned} \quad (17)$$

(b)带入, 缩写

$$\mathbb{E}[f] = \mathbb{E}\left[\frac{1}{k} \sum_{i=1}^k y_{nn(i)}\right] = \frac{1}{k} \sum_{i=1}^k \mathbb{E}[F(\mathbf{x}_{nn(i)}) + \epsilon] = \frac{1}{k} \sum_{i=1}^k \mathbb{E}[F(\mathbf{x}_{nn(i)})] \quad (18)$$

(c)x,y带入f

$$\begin{aligned} & (F(\mathbf{x}) - \mathbb{E}_D[f(\mathbf{x}; D)])^2 + \mathbb{E}_D[(f(\mathbf{x}; D) - \mathbb{E}_D[f(\mathbf{x}; D)])^2] + \sigma^2 \\ &= \left(F(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \mathbb{E}_D[F(\mathbf{x}_{nn(i)})]\right)^2 + \mathbb{E}_D\left[\frac{1}{k} \left(\sum_{i=1}^k y_{nn(i)} - \frac{1}{k} \sum_{i=1}^k \mathbb{E}_D[F(\mathbf{x}_{nn(i)})]\right)^2\right] + \sigma^2 \\ &= \left(F(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \mathbb{E}_D[F(\mathbf{x}_{nn(i)})]\right)^2 + \frac{1}{k^2} \mathbb{E}_D\left[\left(\sum_{i=1}^k (y_{nn(i)} - \mathbb{E}_D[F(\mathbf{x}_{nn(i)})])\right)^2\right] + \sigma^2 \end{aligned} \quad (19)$$

(d)方差项与k

$$\frac{1}{k^2} \mathbb{E}_D \left[\left(\sum_{i=1}^k (y_{nn(i)} - \mathbb{E}_D [F(\mathbf{x}_{nn(i)})]) \right)^2 \right] \quad (20)$$

k增大时，方差项系数变小，找到的最近邻更多，方差整体减小

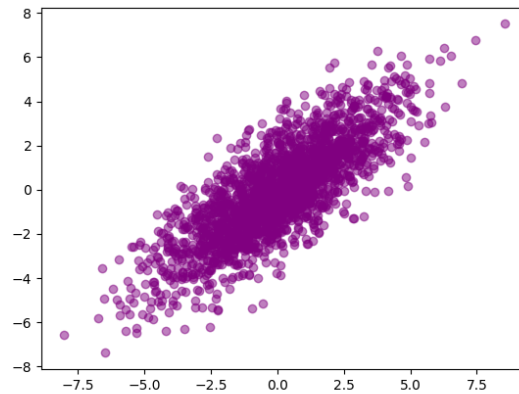
(e)偏差平方项与k

$$\left(F(\mathbf{x}) - \frac{1}{k} \sum_{i=1}^k \mathbb{E}_D [F(\mathbf{x}_{nn(i)})] \right)^2 \quad (21)$$

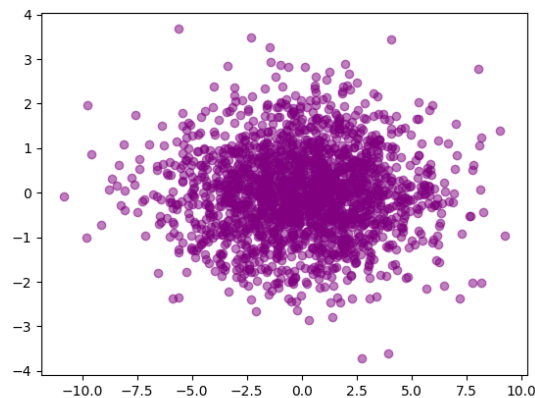
k增大时，这两项的差会越来越大，导致偏差增大。尤其是当k=n时，偏差达到最大，方差达到最小 (0)

五、 (5.3) 编程：PCA&白化

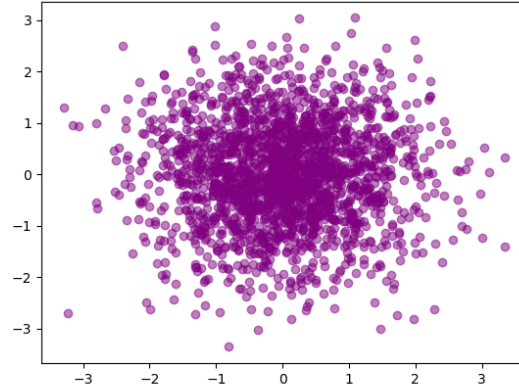
(a)



(b)



(c)



(d)PCA本质是旋转

因为PCA 本质是将数据视作了一个多维空间中的类球形, 并把这个球的各个轴按照各方差最大的方向, 旋转对齐到坐标轴上。用数学方式解释, 就是把数据乘上一个旋转矩阵,

PCA 旋转这一操作有效的原因: PCA数据降维的本质, 就是在对齐到坐标轴上后, 把短轴对应纬度去掉, 保留几个长轴对应的维度, 进而得到新的降维后数据。由于已经进行旋转对齐, 所以去除短轴这一过程很简单, 只需比较轴长短即可。

六、 (6.3) 条件数

(a)矩阵 2-范数 = σ_{max}

矩阵 2-范数等于其最大奇异值, 可知 $\|X\|_2 = \sigma_1$, 且由矩阵的逆的性质可知 $\|X^{-1}\|_2 = \frac{1}{\sigma_n}$

$$\kappa_2(X) = \|X\|_2 \|X^{-1}\|_2 = \frac{\sigma_1}{\sigma_n} \quad (22)$$

(b)病态线性系统

我们想要解释的是, 在 $\kappa_2(A)$ 很大的情况下, 稍微改变 A 或 b 就会使 x 有很大的改变.

已知 $\Delta x = A^{-1} \Delta b$, $\|b\| \leq \|A\| \|x\|$, $\|\Delta x\| \leq \|A^{-1}\| \|\Delta b\|$

相乘再除以 $\|b\| \|x\|$ 可得

$$\frac{\|\Delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\Delta b\|}{\|b\|} = \kappa_2(A) \frac{\|\Delta b\|}{\|b\|} \quad (23)$$

再进行扰动 ΔA 可得

$$\begin{aligned} (A + \Delta A)(x + \Delta x) &= b = Ax \\ A\Delta x &= -\Delta A(x + \Delta x) \\ \Delta x &= -A^{-1}\Delta A(x + \Delta x) \end{aligned} \quad (24)$$

因此我们使用范数不等式并两边除以 $\|x + \Delta x\|$ 有

$$\frac{\|\Delta \mathbf{x}\|}{\|\mathbf{x} + \Delta \mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\Delta \mathbf{A}\|}{\|\mathbf{A}\|} = \kappa_2(\mathbf{A}) \frac{\|\Delta \mathbf{A}\|}{\|\mathbf{A}\|} \quad (25)$$

由此表达式发现，即使当有较小的扰动 $\Delta \mathbf{A}$ 或者 $\Delta \mathbf{b}$ 的时候, 也会带来较大的 $\Delta \mathbf{x}$, 小的输入变换就会导致较大的输出变化

这一定程度上说明了病态系统的原因

(c)良态正交矩阵

正交矩阵的逆等于其转置，有相同特征值

$$\kappa_2(\mathbf{X}) = \|\mathbf{W}\|_2 \|\mathbf{W}^{-1}\|_2 = \|\mathbf{W}\|_2 \|\mathbf{W}^T\|_2 = (\|\mathbf{W}\|_2)^2 = 1 \quad (26)$$