

# Homework 2

Instructor: 吴建鑫

Name: 方盛俊, StudentId: 201300035

## 1. 习题一

(a)

题目中引入了  $\gamma_{ij}$ , 其标识着样本  $\mathbf{x}_j$  是否被分到了第  $i$  组, 同时引入了  $\boldsymbol{\mu}_i$  作为第  $i$  组的代表, 因此我们使用  $\gamma_{ij}$  和  $\boldsymbol{\mu}_i$  对 K-means 的目标进行形式化.

K-means 的目标是每组的样本彼此相似, 即属于相同组的一对样本之间的距离很小. 为了形式化地定义 K-means 的目标, 我们认为  $\boldsymbol{\mu}_i$  表示了聚类的中心, 我们的目标是找到数据点分别属于的聚类, 以及一组向量  $\{\boldsymbol{\mu}_i\}$ , 使得每个数据点和它最近的向量  $\boldsymbol{\mu}_i$  之间的距离的平方和最小. 因此, 我们就可以定义一个目标函数

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

其表示每个数据点和它被分配到的向量  $\boldsymbol{\mu}_i$  之间距离的平方和, 只需要最小化目标函数  $J$  即可求解出最优  $\gamma_{ij}$  和  $\boldsymbol{\mu}_i$ , 进而完成样本的聚类, 即有

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

(b)

第  $i$  步:

在固定  $\boldsymbol{\mu}_i$  的情况下, 每个样本  $\mathbf{x}_j$  对应的  $J$  的每一个贡献成分

$$J_j = \sum_{i=1}^K \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

互相独立, 因此我们只要最优化每一个  $J_j$  即可最优化最终的  $J$ .

由于  $\sum_{i=1}^K \gamma_{ij} = 1$  且  $\gamma_{ij} \in \{0, 1\}$ , 因此  $\{\gamma_{ij}, i = 1, 2, \dots, K\}$  中只有一个元素为 1, 其他元素均为零.

则我们易知最优化问题

$$\arg \min_i \sum_{j=1}^M \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

与最优化问题

$$\arg \min_i \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

等价.

因此我们可以简单地将数据点对应的聚类设置为最近的聚类中心, 形式化地表达即为

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

第 ii 步:

在固定  $\gamma_{ij}$  的情况下,  $J$  是一个关于  $\boldsymbol{\mu}_i$  的二次函数, 为了最小化  $J$ , 我们只需令  $J$  关于  $\boldsymbol{\mu}_i$  的导数等于零, 即可有最小值, 即

$$\frac{\partial J}{\partial \boldsymbol{\mu}_i} = 2 \sum_{j=1}^M \gamma_{ij} (\mathbf{x}_j - \boldsymbol{\mu}_i) = 0$$

解出  $\boldsymbol{\mu}_i$  即可得这一步的更新规则, 即

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

(c)

对于第 i 步:

由于我们有更新后的  $\gamma'_{ij}$ :

$$\gamma'_{ij} = \begin{cases} 1, & \text{if } i = \arg \min_i \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2 \\ 0, & \text{otherwise} \end{cases}$$

因此新目标函数值减去原目标函数值为

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

其中  $\boldsymbol{\mu}'_i$  是离  $\mathbf{x}_j$  最近的簇, 而  $\boldsymbol{\mu}_i$  可能是任意一个簇, 因此我们有  $J' \leq J$ . 即第 i 步会使目标函数  $J$  的值减小或持平.

对于第 ii 步:

对于任意一个簇  $i$  来说, 它的簇中心原来是  $\boldsymbol{\mu}_i$ , 可能是任意一个向量, 之后被优化为

$$\boldsymbol{\mu}'_i = \frac{\sum_{j=1}^K \gamma_{ij} \mathbf{x}_j}{\sum_{j=1}^K \gamma_{ij}}$$

原来的目标函数可以改写为

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

即交换求和符号, 这样我们只需要证明  $J_j = \sum_{i=1}^K \gamma_{ij} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$  降低或不增加即可.

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

因此我们有  $J'_j - J_j$ , 也即有第  $j$  步会使目标函数  $J$  的值降低或不增加.

**下面证明 Lloyd 算法会停止:**

假设算法不会在有限步内停止, 则目标函数的值  $J$  一直在变化.

由 (1) 可知, 目标函数  $J$  的值降低或不增加, 又因为  $J$  一直在变化, 可以将一系列  $J$  的值视作严格单调递减数列, 由于  $J \geq 0$ , 有下界, 因此一定收敛.

并且我们可知,  $J$  的值由簇的分类  $i$  和簇中心  $\boldsymbol{\mu}_i$  唯一确定, 而  $\boldsymbol{\mu}_i = \frac{\sum_{j=1}^K \gamma_{ij} \mathbf{x}_j}{\sum_{j=1}^K \gamma_{ij}}$ , 因此  $\boldsymbol{\mu}_i$  也由  $i$  唯一确定, 也即  $J$  由  $i$  唯一确定, 其中  $j = 1, \dots, k$ .

由于样本  $\mathbf{x}_j$  是有限个的, 因此  $i$  的划分方式是有限个的, 也就是  $J$  的取值是离散的有限个的值, 再由  $J$  是严格单调递减数列且有下界可知, 一定会有一个最小值  $J_{\min} \leq J$ , 对于任何一个  $J$  值来说. 因此  $J$  一定会在  $J_{\min}$  的时候停止, 与假设矛盾.

因此算法会在有限步内停止, 即能够收敛.

## 2. 习题二

(a)

由线性回归的模型假设

$$y = \mathbf{x}^T \boldsymbol{\beta} + \epsilon$$

可得平方误差

$$\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$$

则线性回归任务可以表示为优化问题

$$\arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$$

(b)

使用  $\mathbf{X}$  和  $\mathbf{y}$  重写优化问题即有

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

(c)

$$\text{令平方误差 } E = \sum_{i=1}^n \epsilon_i^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

为了最小化  $E$ , 我们对  $\boldsymbol{\beta}$  求导且令其等于零向量有

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

由于  $\mathbf{X}^T \mathbf{X}$  是可逆的, 则有最优

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

(d)

由于  $d > n$ , 且我们知道  $\mathbf{X}$  是一个  $n \times d$  矩阵, 则其秩  $\text{rank}(\mathbf{X}) \leq n < d$ .

由矩阵的性质可知  $\text{rank}(\mathbf{X}^T \mathbf{X}) = \text{rank}(\mathbf{X}) \leq n < d$ .

而我们又知道  $\mathbf{X}^T \mathbf{X}$  是一个  $d \times d$  的矩阵, 因此  $\mathbf{X}^T \mathbf{X}$  必然不满秩, 因此  $\mathbf{X}^T \mathbf{X}$  不可逆.

(e)

该正则项带来的影响为减少模型的复杂度. 在实际问题中, 我们常常会遇到示例相对较少, 而特征较多的情况, 这种情况下  $\mathbf{X}^T \mathbf{X}$  不一定可逆, 因此无法获得唯一的模型参数, 会有多个模型能够 "完美" 拟合训练集中的所有样例. 在加入正则化项之后, 由于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, 因此相对于在多个训练集上表现痛的预测结果的模型中选出模型复杂度最低的一个.

(f)

加入正则化项后得到的岭回归的优化问题为

$$\arg \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta$$

为了最小化  $E$ , 我们对  $\beta$  求导并令其等于零向量可得

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

解得最优

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

(g)

当  $\mathbf{X}^T \mathbf{X}$  不可逆时, 我们无法求解普通线性回归, 也即无法获得唯一的模型参数, 会有多个模型能够 "完美" 拟合训练集中的所有样例. 加入岭回归正则项后,  $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  几乎总是可逆的, 进而可以求解.

在加入正则化项之后, 由于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, 因此相对于在多个训练集上表现痛的预测结果的模型中选出模型复杂度最低的一个.

(h)

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  $\arg \min_{\beta} \beta^T \beta$ , 进而只能解出  $\beta = \mathbf{0}$ .

(i)

不能, 因为联合优化的话岭回归会退化为普通线性回归, 即一定有  $\lambda = 0$ .

我们可以使用反证法来说明:

假设我们联合优化后得到最优  $\beta^*$  和非零的最优  $\lambda^*$ , 对应的岭回归损失函数值为

$$E^* = (\mathbf{y} - \mathbf{X}\beta^*)^T (\mathbf{y} - \mathbf{X}\beta^*) + \lambda^* \beta^{*T} \beta^*$$

但是我们令  $\lambda = 0$  有

$$(\mathbf{y} - \mathbf{X}\beta^*)^T (\mathbf{y} - \mathbf{X}\beta^*) < (\mathbf{y} - \mathbf{X}\beta^*)^T (\mathbf{y} - \mathbf{X}\beta^*) + \lambda^* \beta^{*T} \beta^* = E^*$$

则与  $E^*$  是最小岭回归损失矛盾, 因此一定有  $\lambda = 0$ .

### 3. 习题三

(a)

下标	类别标记	得分	查准率 P	查全率 R	AUC-PR	AP
0	-	-	1.0000	0.0000	-	-
1	1	1.0	1.0000	0.2000	0.2000	0.2000

下标	类别标记	得分	查准率 P	查全率 R	AUC-PR	AP
2	2	0.9	0.5000	0.2000	0.0000	0.0000
3	1	0.8	0.6667	0.4000	0.1167	0.1333
4	1	0.7	0.7500	0.6000	0.1417	0.1500
5	2	0.6	0.6000	0.6000	0.0000	0.0000
6	1	0.5	0.6667	0.8000	0.1267	0.1333
7	2	0.4	0.5714	0.8000	0.0000	0.0000
8	2	0.3	0.5000	0.8000	0.0000	0.0000
9	1	0.2	0.5556	1.0000	0.1056	0.1111
10	2	0.1	0.5000	1.0000	0.0000	0.0000
-	-	-	-	-	0.6906	0.7278

## (b)

如表格所示. 由于 AUC-PR 和 AP 都是对 PR 曲线的总结, 因此它们的值确实彼此相似. 但是由于所采取的计算方法不同, 精度估计不同, 最后得到的结果也会有所差异, 例如此处的 AP 就比 AUC-PR 要大 0.0372.

同理我们也可以通过数学推导的方式证明相似:

$$AUC\_PR = \sum_{i=1}^n (r_i - r_{i-1}) \frac{p_i + p_{i-1}}{2}$$

$$AP = \sum_{i=1}^n (r_i - r_{i-1}) p_i$$

两式相减可得

$$AP - AUC\_PR = \sum_{i=1}^n (r_i - r_{i-1}) p_i - \sum_{i=1}^n (r_i - r_{i-1}) \frac{p_i + p_{i-1}}{2}$$

$$= \sum_{i=1}^n \frac{1}{2} (r_i - r_{i-1}) (p_i - p_{i-1})$$

可以看出 AP 总是比 AUC-PR 大一点, 但是不会过分地大.

## (c)

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

## (d)

代码如下:

```

values = [1.0, 0.9, 0.8, 0.7, 0.6, 0.5, 0.4, 0.3, 0.2, 0.1]
labels = [1, 2, 1, 1, 2, 1, 2, 2, 1, 2]
# swap 9 and 10
# labels = [1, 2, 1, 1, 2, 1, 2, 2, 2, 1]

P = [1.0]
R = [0.0]
TPR = [0.0]
FPR = [0.0]

for i in range(1, len(values) + 1):
    P_counter = Counter(labels[:i])
    N_counter = Counter(labels[i:])
    TP = P_counter.get(1, 0)
    FP = P_counter.get(2, 0)
    FN = N_counter.get(1, 0)
    TN = N_counter.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))]
AUC_PR_SUM = sum(AUC_PR)
AP = [(R[i] - R[i - 1]) * P[i] for i in range(1, len(R))]
AP_SUM = sum(AP)

print('P:', [%.4f % f for f in P])
print('R:', [%.4f % f for f in R])
print('AUC_PR:', [%.4f % f for f in AUC_PR])
print('AUC_PR_SUM:', '%.4f' % AUC_PR_SUM)
print('AP:', [%.4f % f for f in AP])
print('AP_SUM:', '%.4f' % AP_SUM)

```

## 4. 习题四

(a)

以下推导均为选定一个训练集上的样本  $\mathbf{x}$  的情况下进行的.

我们知道误差是

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

因为噪声  $\epsilon$  独立于其他随机变量, 则有

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

其中, 由于独立性有

$$\mathbb{E}_D[\epsilon^2] = (\mathbb{E}_D[\epsilon])^2 + \text{Var}(\epsilon) = \sigma^2$$

以及有

$$\mathbb{E}_D[(F(\mathbf{x}) - f(\mathbf{x}; D))\epsilon] = \mathbb{E}_D[(F(\mathbf{x}) - f(\mathbf{x}; D))]\mathbb{E}_D[\epsilon] = 0$$

我们进一步展开  $\mathbb{E}_D[(F(\mathbf{x}) - f(\mathbf{x}; D))^2]$  可得

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

对于右边表达式的第一项, 由于  $F(\mathbf{x})$  是确定的, 与训练集  $D$  无关, 即  $\mathbb{E}_D[F(\mathbf{x})] = F(\mathbf{x})$ , 则有

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

对于右边表达式的第二项, 由于  $F(\mathbf{x})$  是确定的, 在这里可以视为一个常数, 不影响方差, 因此我们有

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

综合上面的结果, 我们可以得到偏置-方差分解

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

其中第一项为真实标记  $F(\mathbf{x})$  与所有训练集  $D$  下期望输出标记的偏差, 第二项为  $f(\mathbf{x}; D)$  关于训练集  $D$  的方差, 第三项为噪声  $\epsilon$  关于训练集  $D$  的方差.

**(b)**

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

**(c)**

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

**(d)**

方差项为

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

当  $k$  增加时, 在训练集上寻找的最近邻数就越多, 就越有可能使得  $y_{nn(i)}$  与  $\mathbb{E}_D[F(\mathbf{x}_{nn(i)})]$  项接近, 同时方差项的系数  $\frac{1}{k^2}$  会不断减小, 因此方差项总体会变得平滑, 并且总体上会不断减小.

**(e)**

偏置的平方项为



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

当  $k$  增加时, 由于寻找的最近邻数变得更多, 因此  $\frac{1}{k} \sum_{i=1}^k \mathbb{E}_D[F(\mathbf{x}_{nn(i)})]$  会逐渐远离  $F(\mathbf{x})$ , 导致偏置的平方项会不断增大.

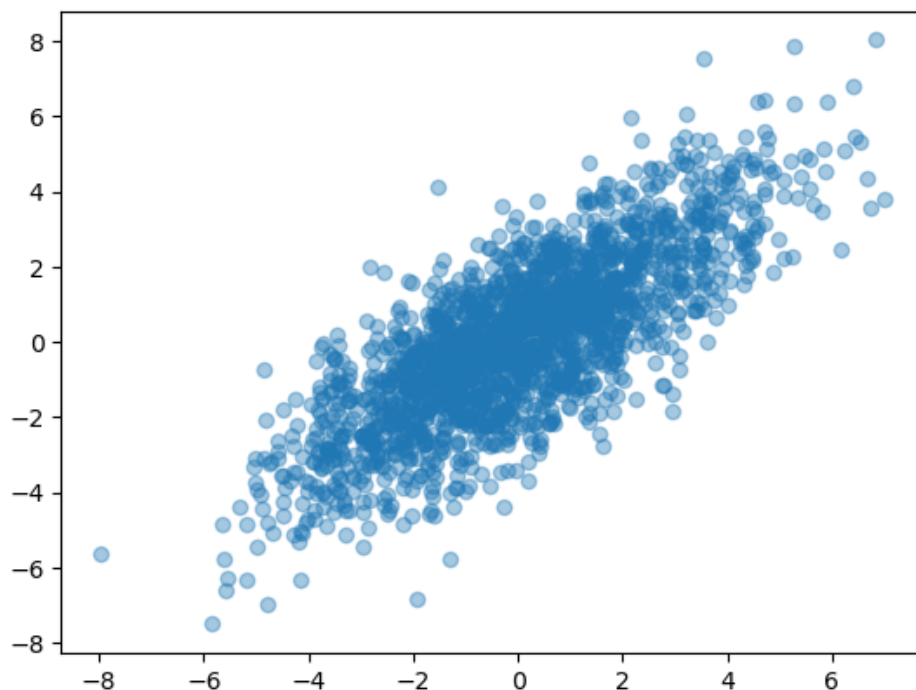
当  $k = n$  时,  $\frac{1}{k} \sum_{i=1}^k \mathbb{E}_D[F(\mathbf{x}_{nn(i)})]$  会变为样本均值, 此时偏置的平方项会最大.

## 5. 习题五

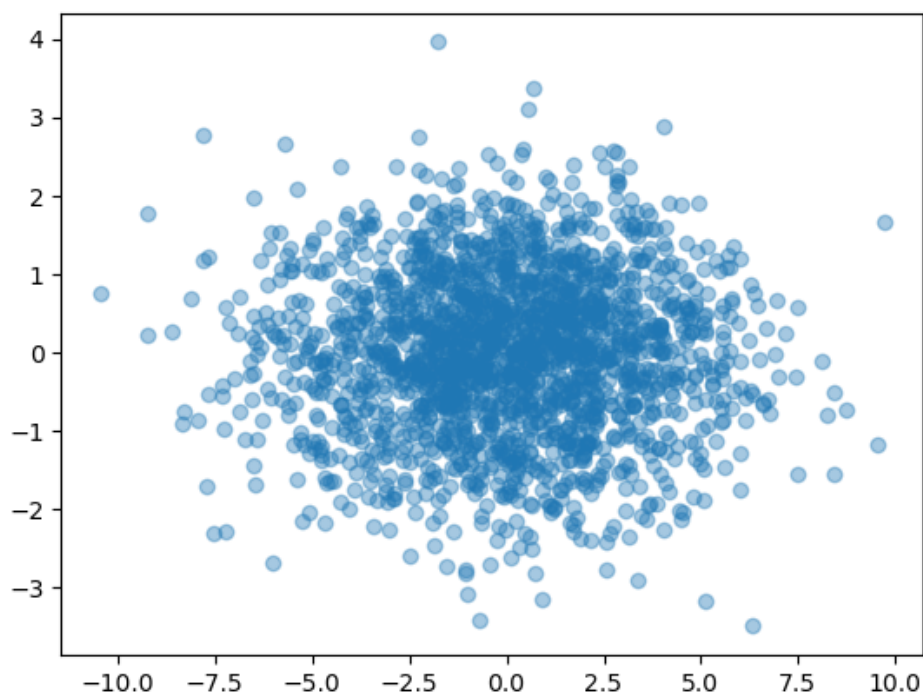
### (a) (b) (c) 代码

见代码文件 `problem_5.ipynb`.

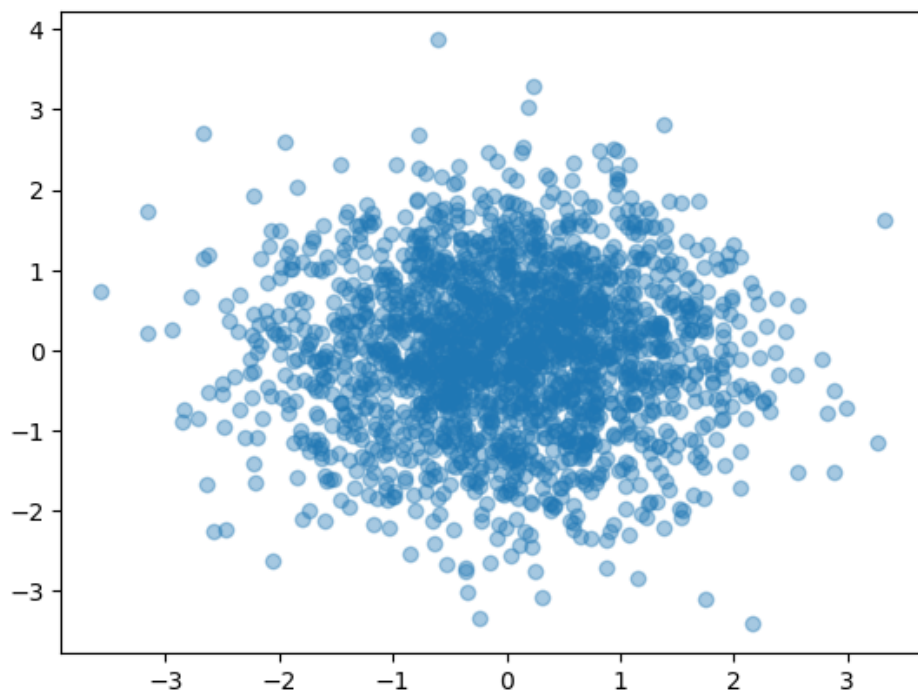
#### (a)



(b)



(c)



(d)

为了推导保留所有维度的 PCA, 也即多维 PCA 的特殊情况, 我们先从 PCA 降维到 0 维开始, 依次推广到多维.

**PCA 降维到 0 维:**

即找到固定的向量  $\mathbf{m}$  使得  $\min_{\mathbf{m}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{m}\|^2$ . 通过对  $\mathbf{m}$  求导等于零我们可得  $\sum_{i=1}^n (\mathbf{m}^* - \mathbf{x}_i) = 0$ , 解得  $\mathbf{m}^* = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$ .

**PCA 降维到 1 维:**

在 PCA 降维到 0 维的基础上, 我们进一步准备 PCA 降维到 1 维, 即有找到  $\omega$  使得  $x_i \approx \bar{x} + a_i \omega$  最近似, 即最小化残差  $x_i - (\bar{x} + a_i \omega)$ , 其中  $\bar{x}$  即是我们 PCA 降维到 0 维的结果.

令  $\mathbf{a} = (a_1, a_2, \dots, a_n)^T$ , 我们定义优化目标  $J(\omega, \mathbf{a}) = \frac{1}{n} \sum_i^n \|x_i - (\bar{x} + a_i \omega)\|^2 = \frac{1}{n} \sum_{i=1}^n (a_i^2 \|\omega\|^2 + \|x_i - \bar{x}\|^2 - 2a_i \omega^T (x_i - \bar{x}))$ .

令对  $a_i$  的偏导等于零即可得

$$\frac{\partial J}{\partial a_i} = \frac{2}{n} (a_i \|\omega\|^2 - \omega^T (x_i - \bar{x})) = 0, \forall i$$

可以解得  $a_i = \frac{(x_i - \bar{x})^T \omega}{\|\omega\|^2} = \frac{(x_i - \bar{x})^T \omega}{\|\omega\|} \cdot \frac{1}{\|\omega\|}$ .

因此我们可以得到  $a_i$  即为向量  $x_i - \bar{x}$  在向量  $\omega$  上的投影乘上  $1/\|\omega\|$ . 并且容易发现,  $J(\omega, \mathbf{a})$  等价于  $J(c\omega, \frac{1}{c}\mathbf{a})$ , 因此为了简化计算, 我们可以令  $\|\omega\| = 1$ , 这也是后续我们推导出 PCA 是数据通过减去均值平移后得一个旋转这个结论的基础.

这时  $a_i$  的取值就简化为  $x_i - \bar{x}$  到单位向量  $\omega$  的投影  $a_i = (x_i - \bar{x})^T \omega$ .

重新带入优化目标则有  $J(\omega, \mathbf{a}) = \frac{1}{n} \sum_{i=1}^n (\|x_i - \bar{x}\|^2 - a_i^2)$ . 因此只需要最大化  $\frac{1}{n} \sum_{i=1}^n a_i^2$  即可得到最优参数.

再令对  $\omega$  的偏导等于零即可得

$$\frac{\partial J}{\partial \omega} = \frac{2}{n} \sum_{i=1}^n (a_i^2 \omega - a_i (x_i - \bar{x})) = \mathbf{0}$$

也就有

$$\frac{1}{n} \left( \sum_{i=1}^n a_i^2 \right) \omega = \frac{1}{n} \sum_{i=1}^n a_i (x_i - \bar{x}) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \omega = \text{Cov}(\mathbf{x}) \omega$$

也即有

$$\text{Cov}(\mathbf{x}) \omega = \frac{\sum_{i=1}^n a_i^2}{n} \omega$$

即  $\omega$  是  $\text{Cov}(\mathbf{x})$  的特征向量, 而  $\frac{1}{n} \sum_{i=1}^n a_i^2$  是该特征向量对应的特征值.

我们不妨令  $\text{Cov}(\mathbf{x})$  对应的特征值为  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ , 以及每个特征值对应的特征向量  $\omega_1, \omega_2, \dots, \omega_d$ , 其中  $d$  为样本  $x_i$  的维度.

为了最大化  $\frac{1}{n} \sum_{i=1}^n a_i^2$  以取得最优参数, 我们应该选择  $\text{Cov}(\mathbf{x})$  中最大的特征值  $\lambda_1$  对应的特征向量  $\omega_1$ , 且有特征向量是单位向量, 即  $\|\omega_1\| = 1$ .

## PCA 降维到多维:

为了将 PCA 从 1 维推广到多维, 我们由谱分解即可得

$$\text{Cov}(\mathbf{x}) = \sum_{i=1}^d \lambda_i \omega_i \omega_i^T$$

我们构造一个矩阵  $\mathbf{W} = [\boldsymbol{\omega}_1 \ \boldsymbol{\omega}_2 \ \cdots \ \boldsymbol{\omega}_d]$ . 由于实对称矩阵  $\text{Cov}(\mathbf{x})$  的特征向量满足  $\boldsymbol{\omega}_i^T \boldsymbol{\omega}_j = 0, i \neq j$ , 因此我们有  $\mathbf{W}$  是正交矩阵, 即有  $\mathbf{W}\mathbf{W}^T = \mathbf{W}^T\mathbf{W} = \mathbf{I}$ .

因此我们有

$$\mathbf{x}_i - \bar{\mathbf{x}} = \mathbf{W}\mathbf{W}^T(\mathbf{x}_i - \bar{\mathbf{x}}) = \boldsymbol{\omega}_1^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_1 + \cdots + \boldsymbol{\omega}_d^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_d$$

为了将 PCA 降维推广到多维, 我们用原样本减去已经进行了 1 维 PCA 得到的结果, 即

$$\mathbf{x}'_i = \mathbf{x}_i - (\bar{\mathbf{x}} + \boldsymbol{\omega}_1^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_1) = \sum_{j=2}^d \boldsymbol{\omega}_j^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_j$$

可以看出, 如果我们在新数据集  $\{\mathbf{x}_i = \sum_{j=2}^d \boldsymbol{\omega}_j^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_j, i = 1, \cdots, n\}$  上做 1 维 PCA 变换, 则对于每个新样本会得到 PCA 结果  $\boldsymbol{\omega}_2^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_2$ , 以此类推, 直到最后一维.

我们逐一进行 PCA 的过程, 其本质和我们直接选取前  $d'$  个最大特征值和特征向量是等价的, 其中  $d'$  为要降到的维度.

如果我们直接保留所有的维度, 其实质就是保留整个  $\mathbf{W}$ , 也就是有

$$\mathbf{x}_i - \bar{\mathbf{x}} = \sum_{j=1}^d \boldsymbol{\omega}_j^T(\mathbf{x}_i - \bar{\mathbf{x}})\boldsymbol{\omega}_j = \mathbf{W}^T(\mathbf{x}_i - \bar{\mathbf{x}})\mathbf{W}$$

即我们最终得到的新表示为

$$\mathbf{y}_i = \mathbf{W}^T(\mathbf{x}_i - \bar{\mathbf{x}})$$

**以上即为多维 PCA 的推导, 基于这个推导我们有下面结论:**

由于  $\mathbf{W}$  是一个  $d \times d$  的正交矩阵, 因此我们可知保留所有维度的 PCA 是数据在进行平移之后的一个旋转.

我们保留所有维度的 PCA, 也就是进行了一个旋转, 之所以有用, 是因为 PCA 本质是将数据视作了一个多维空间中的椭球, PCA 要做的事情就是将这个椭球的各个轴旋转对齐到坐标轴上 (即将各方差最大的方向对齐到坐标轴上). 如果我们需要进行降维, 那就将这个椭球的几个短轴对应的维度去掉, 保留几个长轴对应的维度, 由于数据已经与坐标轴对齐, 我们可以很简单地将短轴对应的那几个维度对应的数字去掉, 进而得到新的降维后数据. 这就是 PCA 的本质, 这也是 PCA 旋转这一操作有效的原因.

## 6. 习题六

(a)

由矩阵 2-范数的定义可知  $\|\mathbf{X}\|_2 = \sigma_1$ , 且由矩阵的逆的性质可知  $\|\mathbf{X}^{-1}\|_2 = \frac{1}{\sigma_n}$ , 因此有

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

(b)

若我们要求解  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , 则有  $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ .

我们要说明的就是, 在  $\kappa_2(\mathbf{A})$  很大的情况下, 稍微改变  $\mathbf{A}$  或  $\mathbf{b}$  就会导致  $\mathbf{x}$  有很大的改变.

推导可知 (本小问中的  $\|\cdot\|$  均指 2-范数  $\|\cdot\|_2$ )

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

由上面的式子可得

$$\begin{aligned} \|\mathbf{b}\| &\leq \|\mathbf{A}\| \|\mathbf{x}\| \\ \|\Delta\mathbf{x}\| &\leq \|\mathbf{A}^{-1}\| \|\Delta\mathbf{b}\| \end{aligned}$$

两式相乘再除以  $\|\mathbf{b}\| \|\mathbf{x}\|$  可得

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

同理我们对  $\mathbf{A}$  进行扰动变为  $\mathbf{A} + \Delta\mathbf{A}$  可得

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

因此我们使用范数不等式并两边除以  $\|\mathbf{x} + \Delta\mathbf{x}\|$  有

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

可以看出, 当有较小的扰动  $\Delta\mathbf{A}$  或者  $\Delta\mathbf{b}$  的时候, 尤其是能够取得等号的时候, 均会带来较大的  $\Delta\mathbf{x}$ , 即较小的输入变换就会导致较大的输出变化. 这种变化对我们使用计算机进行一定精度的矩阵计算尤为不利.

例如用 Matlab 求解下列方程时:

$$\begin{bmatrix} 1 & 1 \\ 1 & 0.999 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 0.999 \end{bmatrix} \Rightarrow \mathbf{x} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

条件数  $\kappa_2(\mathbf{A}) = 3998$ , 可以看出是一个较大的数.

扰动  $\mathbf{A}$  后求解  $\mathbf{x}$ :

$$\begin{bmatrix} 1.001 & 1.001 \\ 1.001 & 0.999 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 0.999 \end{bmatrix} \Rightarrow \mathbf{x} = \begin{bmatrix} 0.499 \\ 0.500 \end{bmatrix}$$

扰动  $\mathbf{b}$  后求解  $\mathbf{x}$ :

$$\begin{bmatrix} 1 & 1 \\ 1 & 0.999 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1.001 \\ 0.999 \end{bmatrix} \Rightarrow \mathbf{x} = \begin{bmatrix} -0.999 \\ +2.000 \end{bmatrix}$$

可以看出, 无论是对  $\mathbf{A}$  还是  $\mathbf{b}$  做一个轻微的扰动, 均会导致解出来的  $\mathbf{x}$  产生较大的变化, 因此这个线性系统是病态的.

### (c)

对于正交矩阵  $\mathbf{W}$ , 我们有  $\mathbf{W}^{-1} = \mathbf{W}^T$ , 且  $\mathbf{W}^T$  与  $\mathbf{W}$  有相同的特征值, 以及  $\mathbf{W}$  的奇异值 (特征值) 绝对值为 1. 因此有

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

因此正交矩阵是良态的, 有较小的条件数.