

# 机器学习导论

## 习题五

学号, 作者姓名, 邮箱

2021 年 6 月 23 日

### 学术诚信

本课程非常重视学术诚信规范，助教老师和助教同学将不遗余力地维护作业中的学术诚信规范的建立。希望所有选课学生能够对此予以重视。<sup>1</sup>

- (1) 允许同学之间的相互讨论，但是**署你名字的工作必须由你完成**，不允许直接照搬任何已有的材料，必须独立完成作业的书写过程；
- (2) 在完成作业过程中，对他人工作（出版物、互联网资料）中文本的直接照搬（包括原文的直接复制粘贴及语句的简单修改等）都将视为剽窃，剽窃者成绩将被取消。**对于完成作业中有关键作用的公开资料，应予以明显引用；**
- (3) 如果发现作业之间高度相似将被判定为互相抄袭行为，**抄袭和被抄袭双方的成绩都将被取消**。因此请主动防止自己的作业被他人抄袭。

### 作业提交注意事项

- (1) 请在L<sup>A</sup>T<sub>E</sub>X模板中**第一页填写个人的姓名、学号、邮箱信息**；
- (2) 本次作业需提交该pdf文件，pdf文件名格式为**学号\_姓名.pdf**，例如190000001\_张三.pdf，**需通过教学立方提交**。
- (3) 未按照要求提交作业，或提交作业格式不正确，将会**被扣除部分作业分数**；
- (4) 本次作业提交截止时间为**6月6日23:55:00**。

---

<sup>1</sup>参考尹一通老师高级算法课程中对学术诚信的说明。

# 1 [30pts] PCA

$\mathbf{x} \in \mathbb{R}^D$  是一个随机向量, 其均值和协方差分别是  $\boldsymbol{\mu}_x = \mathbb{E}(\mathbf{x}) \in \mathbb{R}^D$ ,  $\Sigma_x = \mathbb{E}(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \in \mathbb{R}^{D \times D}$ 。定义随机变量  $y_i = \mathbf{u}_i^\top \mathbf{x} + a_i \in \mathbb{R}, i = 1, \dots, d \leq D$  为  $\mathbf{x}$  的主成分, 其中  $\mathbf{u}_i \in \mathbb{R}^D$  是单位向量( $\mathbf{u}_i^\top \mathbf{u}_i = 1$ ),  $a_i \in \mathbb{R}$ ,  $\{y_i\}_{i=1}^d$  是互不相关的零均值随机变量, 它们的方差满足  $\text{Var}(y_1) \geq \text{Var}(y_2) \geq \dots \geq \text{Var}(y_d)$ 。假设  $\Sigma_x$  没有重复的特征值, 请证明:

1. [5pts]  $a_i = -\mathbf{u}_i^\top \boldsymbol{\mu}_x, i = 1, \dots, d$ 。

2. [10pts]  $\mathbf{u}_1$  是  $\Sigma_x$  最大的特征值对应的特征向量。

提示: 写出要最大化的目标函数, 写出约束条件, 使用拉格朗日乘子法。

3. [15pts]  $\mathbf{u}_2^\top \mathbf{u}_1 = 0$ , 且  $\mathbf{u}_2$  是  $\Sigma_x$  第二大特征值对应的特征向量。

提示: 由  $\{y_i\}_{i=1}^d$  是互不相关的零均值随机变量可推出  $\mathbf{u}_2^\top \mathbf{u}_1 = 0$ 。 $\mathbf{u}_2^\top \mathbf{u}_1 = 0$  可作为第二小问的约束条件之一。

**Solution.** 此处用于写解答(中英文均可)

1. 由  $y_i$  是零均值的随机变量可得

$$\mathbb{E}[y_i] = \boldsymbol{\mu}_i^\top \mathbb{E}[\mathbf{x}] + \mathbb{E}[a_i] = \mathbf{u}_i^\top \boldsymbol{\mu}_x + a_i = 0, \quad i = 1, \dots, d$$

$$\text{则 } a_i = -\mathbf{u}_i^\top \boldsymbol{\mu}_x, i = 1, \dots, d$$

2.

$$\text{Var}(y_1) = \mathbb{E}[(\mathbf{u}_1^\top (\mathbf{x} - \boldsymbol{\mu}_x))^2] = \mathbb{E}[\mathbf{u}_1^\top (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^\top \mathbf{u}_1] = \mathbf{u}_1^\top \Sigma_x \mathbf{u}_1$$

则优化问题为

$$\max_{\mathbf{u}_1 \in \mathbb{R}^D} \mathbf{u}_1^\top \Sigma_x \mathbf{u}_1 \quad \text{s.t.} \quad \mathbf{u}_1^\top \mathbf{u}_1 = 1$$

由拉格朗日函数

$$\mathcal{L} = \mathbf{u}_1^\top \Sigma_x \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^\top \mathbf{u}_1)$$

其中  $\lambda_1$  为拉格朗日乘子。

令上式对  $\mathbf{u}_1$  和  $\lambda_1$  求导分别等于0, 可得

$$\Sigma_x \mathbf{u}_1 = \lambda_1 \mathbf{u}_1 \quad \text{and} \quad \mathbf{u}_1^\top \mathbf{u}_1 = 1$$

因此  $\mathbf{u}_1$  是  $\lambda_1$  所对应的  $\Sigma_x$  的特征向量。

要最大化  $\mathbf{u}_1^\top \Sigma_x \mathbf{u}_1 = \lambda_1 \mathbf{u}_1^\top \mathbf{u}_1 = \lambda_1$ , 则  $\lambda_1$  应为  $\Sigma_x$  最大的特征值,  $\mathbf{u}_1$  是  $\Sigma_x$  最大的特征向量。

3. 由  $y_1, y_2$  不相关, 得  $\text{Cov}(y_1, y_2) = \mathbb{E}(y_1 y_2) - \mathbb{E}(y_1)\mathbb{E}(y_2) = 0$ 。

又  $y_1, y_2$  均值为0, 所以  $\mathbb{E}(y_1 y_2) = \mathbb{E}(y_1)\mathbb{E}(y_2) = 0$ 。即

$$\mathbb{E}[(\mathbf{u}_1^\top (\mathbf{x} - \boldsymbol{\mu}_x))(\mathbf{u}_2^\top (\mathbf{x} - \boldsymbol{\mu}_x))] = \mathbb{E}[\mathbf{u}_1^\top (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^\top \mathbf{u}_2] = \mathbf{u}_1^\top \Sigma_x \mathbf{u}_2 = \lambda_1 \mathbf{u}_1^\top \mathbf{u}_2 = 0$$

由于  $\lambda_1 \neq 0$ , 因此  $\mathbf{u}_1^\top \mathbf{u}_2 = 0$ 。

类似于  $\mathbf{u}_1$  的求解过程, 求解  $\mathbf{u}_2$  对应的最大化目标是  $\mathbf{u}_2^\top \Sigma_x \mathbf{u}_2$ , 多一个约束条件  $\mathbf{u}_1^\top \mathbf{u}_2 = 0$ 。即优化问题为

$$\max_{\mathbf{u}_2 \in \mathbb{R}^D} \mathbf{u}_2^\top \Sigma_x \mathbf{u}_2 \quad \text{s.t.} \quad \mathbf{u}_2^\top \mathbf{u}_2 = 1 \quad \text{and} \quad \mathbf{u}_1^\top \mathbf{u}_2 = 0$$

定义拉格朗日函数

$$\mathcal{L} = \mathbf{u}_2^\top \Sigma_x \mathbf{u}_2 + \lambda_2 (1 - \mathbf{u}_2^\top \mathbf{u}_2) + \gamma \mathbf{u}_1^\top \mathbf{u}_2$$

对  $\mathbf{u}_2, \lambda_2, \gamma$  分别求导令导数等于0, 得

$$\Sigma_x \mathbf{u}_2 + \frac{\gamma}{2} \mathbf{u}_1 = \lambda_2 \mathbf{u}_2, \quad \mathbf{u}_2^\top \mathbf{u}_2 = 1 \quad \text{and} \quad \mathbf{u}_1^\top \mathbf{u}_2 = 0$$

第一个式子左乘以  $\mathbf{u}_1^\top$  可得,  $\mathbf{u}_1^\top \Sigma_x \mathbf{u}_2 + \frac{\gamma}{2} \mathbf{u}_1^\top \mathbf{u}_1 = \lambda_2 \mathbf{u}_1^\top \mathbf{u}_2 + \frac{\gamma}{2} = \lambda_2 \mathbf{u}_1^\top \mathbf{u}_2$ , 因此  $\gamma = 2(\lambda_2 - \lambda_1) \mathbf{u}_1^\top \mathbf{u}_2 = 0$ , 故  $\Sigma_x \mathbf{u}_2 = \lambda_2 \mathbf{u}_2$ 。则最大化的目标为  $\mathbf{u}_2^\top \Sigma_x \mathbf{u}_2 = \lambda_2 = \text{Var}(y_2)$ 。

因为  $\Sigma_x$  没有重复的特征值, 所以当  $\lambda_2$  为  $\Sigma_x$  第二大的特征值,  $\mathbf{u}_2$  为  $\Sigma_x$  第二大的特征值对应的特征向量时,  $\mathbf{u}_2^\top \Sigma_x \mathbf{u}_2$  取得最大值。

## 2 [30pts] Clustering

考虑  $p$  维特征空间里的混合模型

$$g(x) = \sum_{k=1}^K \pi_k g_k(x)$$

其中  $g_k = N(\mu_k, \mathbf{I} \cdot \sigma^2)$ ,  $\mathbf{I}$  是单位矩阵,  $\pi_k > 0$ ,  $\sum_k \pi_k = 1$ 。  $\{\mu_k, \pi_k\}, k = 1, \dots, K$  和  $\sigma^2$  是未知参数。

设有数据  $x_1, x_2, \dots, x_N \sim g(x)$ ,

1. [10pts] 请写出数据的对数似然。
2. [15pts] 请写出求解极大似然估计的 EM 算法。
3. [5pts] 请简要说明如果  $\sigma$  的值已知, 并且  $\sigma \rightarrow 0$ , 那么该 EM 算法就相当于 K-means 聚类。

**Solution.** 此处用于写解答(中英文均可)

1.

$$\begin{aligned} LL(D) &= \log \left( \prod_{i=1}^N g(x_i) \right) \\ &= \sum_{i=1}^N \log \left( \sum_{k=1}^K \pi_k g_k(x_i) \right) \end{aligned}$$

2. 算法仿照书210页图9.6, 其中关键的公式:

E步, 样本  $x_i$  由第  $k$  个成分生成的概率

$$\gamma_{ik} = \frac{\hat{\pi}_k \hat{g}_k(x_i)}{\sum_{l=1}^K \hat{\pi}_l \hat{g}_l(x_i)}$$

M步, 计算新均值向量

$$\hat{\mu}_k = \frac{\sum_{i=1}^N \gamma_{ik} x_i}{\sum_{i=1}^N \gamma_{ik}}$$

一般计算新协方差矩阵,

$$\Sigma_k = \frac{\sum_{i=1}^N \gamma_{ik} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T}{\sum_{i=1}^N \gamma_{ik}}$$

不过这里方差是标量, 因此应为

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^N \sum_{k=1}^K \gamma_{ik} (x_i - \hat{\mu}_k)^T (x_i - \hat{\mu}_k)}{pN}$$

计算新混合系数

$$\hat{\pi}_k = \frac{\sum_{i=1}^N \gamma_{ik}}{N}$$

3.

$$\gamma_{ik} = \frac{\exp\left(-\frac{1}{2\sigma^2} \|x_i - \mu_k\|^2\right)}{\sum_{l=1}^K \exp\left(-\frac{1}{2\sigma^2} \|x_i - \mu_l\|^2\right)}$$

当  $\sigma \rightarrow 0$ , 对于离  $x_i$  最近的成分  $k$ ,  $\gamma_{ik} \rightarrow 1$ , 对于其他成分  $\gamma_{il} \rightarrow 0$ 。

则E步等价于k-means中划入均值向量距离最近的簇。M步等价于重新计算新的簇的均值。

### 3 [40pts] Ensemble Methods

(1) [10pts] GradientBoosting[Friedman, 2001] 是一种常用的 Boosting 算法, 请简要分析其与 AdaBoost 的异同。

(2) [10pts] 请简要说明随机森林为何比决策树 Bagging 集成的训练速度更快。

(3) [20pts] Bagging 产生的每棵树是同分布的, 那么  $B$  棵树均值的期望和其中任一棵树的期望是相同的。因此, Bagging 产生的偏差和其中任一棵树的偏差相同, Bagging 带来的性能提升来自于方差的降低。

我们知道, 方差为  $\sigma^2$  的  $B$  个独立同分布的随机变量, 其均值的方差为  $\frac{1}{B}\sigma^2$ 。如果这些随机变量是同分布的, 但不是独立的, 设两两之间的相关系数  $\rho > 0$ , 请推导均值的方差为  $\rho\sigma^2 + \frac{1-\rho}{B}\sigma^2$ 。

**Solution.** 此处用于写解答(中英文均可)

(1) 共同点：都是用弱学习器逐步拟合没学好的部分，组合在一起成为一个较好的集成模型。

不同点：Adaboost调整没学好的样本的权重，GradientBoosting调整target拟合残差。

言之有理亦可。

(2) 随机森林相对于Bagging决策树的关键区别在于，在选择划分属性时，首先随机选择一个属性集的子集，再在这个子集中寻找最优属性。

由于一般随机选择的属性子集规模比所有属性集小(如推荐值 $k = \log_2 d$ )，训练时只需考察这个较小的子集，从而训练速度更快。

(3)

$$\begin{aligned}\text{Var}\left(\frac{\sum_{i=1}^B X_i}{B}\right) &= \frac{1}{B^2} \text{Var}\left(\sum x_i\right) \\ &= \frac{1}{B^2} \sum_{i=1}^B \text{Var}(X_i) + \frac{1}{B^2} \sum_{i \neq j}^B \text{Cov}(X_i, X_j) \\ &= \frac{\sigma^2}{B} + \frac{B-1}{B} \sigma^2 \rho \\ &= \sigma^2 \rho + \frac{1-\rho}{B} \sigma^2.\end{aligned}$$

(备注：随着 $B$ 的增加，第二项逐渐消失，但第一项消不掉，所以树与树之间的相关性限制了平均带来的好处。)

## 参考文献

[Friedman, 2001] Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232.