

姓名：方盛俊

学号：201300035

一. (20 points) 没有免费的午餐定理

1. 根据教材 1.4 节“没有免费的午餐”定理, 所有学习算法的期望性能都和随机胡猜一样, 是否还有必要继续进行研究机器学习算法?
2. 教材 1.4 节在论述“没有免费的午餐”定理时, 默认使用了“分类错误率”作为性能度量来对分类器进行评估. 若换用其他性能度量 ℓ , 则教材中式 (1.1) 将改为

$$E_{ote}(\mathcal{L}_a|X, f) = \sum_h \cdot \sum_{\mathbf{x} \in \mathcal{X}-X} P(\mathbf{x})\ell(h(\mathbf{x}), f(\mathbf{x}))P(h|\mathcal{X}, \mathcal{L}_a) \quad (1)$$

试证明“没有免费的午餐定理”仍成立.

解:

1. 依然有必要继续研究机器学习算法.

NFL 定理有一个重要的前提, 即所有的”问题”出现的机会相同, 即 f 是均匀分布的, 但是实际上并不是这样. 对于某一类具体的问题来说, f 一般都不是均匀分布的. 因此, 面对一个具体的问题, 一般都会有一个具体的机器学习算法, 能够取得比其他机器学习算法更好的结果.

2. 对于一个有 k 种分类结果的 k 分类器来说, 其真实目标函数可以是任何函数 $\mathcal{X} \mapsto \{1, 2, \dots, k\}$, 函数空间为 $\{1, 2, \dots, k\}^{|\mathcal{X}|}$.

将 f 视为遵循均匀分布, 则我们有预测正确的概率 $\Pr(h(\mathbf{x}) = f(\mathbf{x})) = \frac{1}{k}$, 是一个常数. 预测错误的概率是 $\Pr(h(\mathbf{x}) \neq f(\mathbf{x})) = \frac{k-1}{k}$.

对于分类问题, 任意一种性能度量 $\ell(h(\mathbf{x}), f(\mathbf{x}))$ 的输出只有两种结果, 设预测正确的结果为 c_1 , 预测错误的结果为 c_0 , 则有分布列 $\Pr(\ell(h(\mathbf{x}), f(\mathbf{x})) = c_1) = \frac{1}{k}$, $\Pr(\ell(h(\mathbf{x}), f(\mathbf{x})) = c_0) = \frac{k-1}{k}$.

那么有期望 $\mathbb{E}[\ell(h(\mathbf{x}), f(\mathbf{x}))] = \frac{1}{k} \cdot c_1 + \frac{k-1}{k} \cdot c_0$, 我们可以设 $c = \frac{1}{k} \cdot c_1 + \frac{k-1}{k} \cdot c_0$. 以二分类问题的错误率为例子, 此时 $c = \frac{1}{2} \cdot 0 + \frac{2-1}{2} \cdot 1 = \frac{1}{2}$. 则我们有

$$\begin{aligned} \sum_f E_{ote}(\mathcal{L}_a|X, f) &= \sum_f \sum_h \sum_{\mathbf{x} \in \mathcal{X}-X} P(\mathbf{x}) \ell(h(\mathbf{x}), f(\mathbf{x})) P(h|X, \mathcal{L}_a) \\ &= \sum_{\mathbf{x} \in \mathcal{X}-X} P(\mathbf{x}) \sum_h P(h|X, \mathcal{L}_a) \sum_f \ell(h(\mathbf{x}), f(\mathbf{x})) \\ &= \sum_{\mathbf{x} \in \mathcal{X}-X} P(\mathbf{x}) \sum_h P(h|X, \mathcal{L}_a) c k^{|\mathcal{X}|} \\ &= c k^{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}-X} P(\mathbf{x}) \sum_h P(h|X, \mathcal{L}_a) \\ &= c k^{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}-X} P(\mathbf{x}) \cdot 1 \end{aligned}$$

可以看出, 总误差依然与学习算法无关, 因此对于其他的性能度量来说, NFL 定理依然成立.

二. (15 points) 线性回归

给定包含 m 个样例的数据集 $\mathbf{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, 其中 $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{id}) \in \mathbb{R}^d$, $y_i \in \mathbb{R}$ 为 \mathbf{x}_i 的实数标记. 针对数据集 \mathbf{D} 中的 m 个示例, 教材 3.2 节所介绍的“线性回归”模型要求该线性模型的预测结果和其对应的标记之间的误差之和最小:

$$\begin{aligned} (\mathbf{w}^*, b^*) &= \frac{1}{2} \arg \min_{(\mathbf{w}, b)} \sum_{i=1}^m (f(\mathbf{x}_i) - y_i)^2 \\ &= \frac{1}{2} \arg \min_{(\mathbf{w}, b)} \sum_{i=1}^m (y_i - (\mathbf{w}^\top \mathbf{x}_i + b))^2. \end{aligned} \quad (2)$$

即寻找一组权重 (\mathbf{w}, b) , 使其对 \mathbf{D} 中示例预测的整体误差最小.¹ 定义

¹公式 ?? 中系数 $\frac{1}{2}$ 是为了化简后续推导. 有时也会乘上 $\frac{1}{m}$ 以计算均方误差 (Mean Square Error), 由于平均误差和误差和在优化

$y = [y_1; \dots, y_m] \in \mathbb{R}^m$, 且 $\mathbf{X} = [\mathbf{x}_1^\top; \mathbf{x}_2^\top; \dots; \mathbf{x}_m^\top] \in \mathbb{R}^{m \times d}$, 请将线性回归的优化过程使用矩阵进行表示.

解:

原式可用矩阵表示为

$$(\mathbf{w}^*, b^*) = \arg \min_{(\mathbf{w}, b)} \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{1}b)^\top (\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{1}b)$$

$$\text{令 } E = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{1}b)^\top (\mathbf{y} - \mathbf{X}\mathbf{w} - \mathbf{1}b)$$

对 \mathbf{w} 求导得

$$\frac{\partial E}{\partial \mathbf{w}} = \mathbf{X}^\top (\mathbf{X}\mathbf{w} + \mathbf{1}b - \mathbf{y})$$

对 b 求导得

$$\frac{\partial E}{\partial b} = \mathbf{1}^\top (\mathbf{X}\mathbf{w} + \mathbf{1}b - \mathbf{y})$$

令上面两式同时等于零即可得到 \mathbf{w} 和 b 的最优解的闭式解.

当 $\mathbf{X}^\top \mathbf{X}$ 是满秩矩阵时, 则有

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

令 $\mathbf{T} = (\mathbf{X}^\top \mathbf{X})^{-1}$, 则

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

代入有

$$\begin{aligned} b^* &= \mathbf{1}^\top \mathbf{y} - \mathbf{1}^\top \mathbf{X} \mathbf{w}^* \\ &= \mathbf{1}^\top \mathbf{y} - \mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top (\mathbf{y} - \mathbf{1}b^*) \\ &= \mathbf{1}^\top \mathbf{y} - \mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{y} + \mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{1} b^* \end{aligned}$$

最后有

过程中只相差一个常数, 不影响优化结果, 因此在后续讨论中省略这一系数.

$$b^* = \frac{\mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{y} - \mathbf{1}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{1} - 1}$$

最后将 b^* 带回即可求出

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

三. (25 points) 正则化

在实际问题中, 我们常常会遇到示例相对较少, 而特征很多的场景. 在这类情况中如果直接求解线性回归模型, 较少的示例无法获得唯一的模型参数, 会具有多个模型能够”完美”拟合训练集中的所有样例, 实现插值 (interpolation). 此外, 模型很容易过拟合. 为缓解这些问题, 常在线性回归的闭式解中引入正则化项 $\Omega(\mathbf{w})$, 通常形式如下:

$$\mathbf{w}_{\text{Ridge}}^*, b_{\text{Ridge}}^* = \arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{X} \mathbf{w} + \mathbf{1} b - \mathbf{y}\|_2^2 + \lambda \Omega(\mathbf{w}). \quad (3)$$

其中, $\lambda > 0$ 为正则化参数. 正则化表示了对模型的一种偏好, 例如 $\Omega(\mathbf{w})$ 一般对模型的复杂度进行约束, 因此相当于从多个在训练集上表现同等预测结果的模型中选出模型复杂度最低的一个.

考虑岭回归 (ridge regression) 问题, 即设置公式(??)中正则项 $\Omega(\mathbf{w}) = \|\mathbf{w}\|_2^2$. 本题中将对岭回归的闭式解以及正则化的影响进行探讨.

1. 请给出岭回归的最优解 $\mathbf{w}_{\text{Ridge}}^*$ 和 b_{Ridge}^* 的闭式解表达式, 并使用矩阵形式表示, 分析其最优解和原始线性回归最优解 \mathbf{w}_{LS}^* 和 b_{LS}^* 的区别;
2. 请证明对于任何矩阵 \mathbf{X} , 下式均成立

$$(\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_m)^{-1} \mathbf{X} = \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}. \quad (4)$$

请思考, 上述的结论是否能够帮助岭回归的计算, 在何种情况下能够带来帮助?

3. 针对波士顿房价预测数据 (`boston`), 编程实现原始线性回归模型和岭回归模型, 基于闭式解在训练集上构建模型, 计算测试集上的均方误差 (Mean Square Error, MSE). 请参考 `LinearRegression.py` 进行模型构造.

```

1 from sklearn.datasets import load_boston
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4 import matplotlib.pyplot as plt
5
6 X, y = load_boston(return_X_y=True)
7 trainx, testx, trainy, testy = train_test_split(
8     X, y, test_size=0.33, random_state=42)
9
10
11 def linReg(X_train: np.ndarray, y_train: np.ndarray) -> np.ndarray:
12     # linear regression
13     ones = np.ones(X_train.shape[0]).reshape(-1, 1)
14     X_splash = np.hstack((X_train, ones))
15     T: np.ndarray = np.linalg.inv(X_splash.transpose() @ X_splash)
16     return T @ X_splash.transpose() @ y_train
17
18
19 def linRegMSE(X_train: np.ndarray, y_train: np.ndarray, X_test: np.ndarray, y_test:
20     np.ndarray) -> float:
21     weights = linReg(X_train, y_train).reshape(-1, 1)
22     ones = np.ones(X_test.shape[0]).reshape(-1, 1)
23     X_splash = np.hstack((X_test, ones))
24     E: np.ndarray = y_test.reshape(-1, 1) - X_splash @ weights
25     # return E.transpose() @ E / X_test.shape[0]
26     return np.mean(E ** 2)
27
28 def reportLinRegMSE():
29     return linRegMSE(trainx, trainy, testx, testy)
30
31
32 def ridgeReg(X_train: np.ndarray, y_train: np.ndarray, lmbd: float) -> np.ndarray:
33     # ridge regression
34     ones = np.ones(X_train.shape[0])
35     eye = np.eye(X_train.shape[1])
36     T: np.ndarray = np.linalg.inv(
37         X_train.transpose() @ X_train + 2 * lmbd * eye)
38     S: np.ndarray = ones @ X_train @ T @ X_train.transpose() - ones
39     b = (S @ y_train) / (S @ ones)
40     w = T @ X_train.transpose() @ (y_train - b * ones)
41     return np.hstack((w, b))
42
43
44 def ridgeRegMSE(X_train: np.ndarray, y_train: np.ndarray, X_test: np.ndarray, y_test
45     : np.ndarray, lmbd: float) -> float:
46     weights = ridgeReg(X_train, y_train, lmbd).reshape(-1, 1)
47     ones = np.ones(X_test.shape[0]).reshape(-1, 1)
48     X_splash = np.hstack((X_test, ones))
49     E: np.ndarray = y_test.reshape(-1, 1) - X_splash @ weights
50     # return E.transpose() @ E / X_test.shape[0]
51     return np.mean(E ** 2)
52
53 def reportRidgeRegMSE(lmbd):
54     return ridgeRegMSE(trainx, trainy, testx, testy, lmbd)
55
56 print(linReg(trainx, trainy))
57 print(ridgeReg(trainx, trainy, 0))
58 print(reportLinRegMSE())
59 print(reportRidgeRegMSE(0))
60 lmbds = np.arange(0, 2.1, 0.2)
61 mses = [reportRidgeRegMSE(lmbd) for lmbd in lmbds]
62 print(lmbds)
63 print(mses)
64 plt.plot(lmbds, mses)

```

```
65 plt.xlabel('lambda')
66 plt.ylabel('MSE')
67 plt.show()
```

- (a) 对于线性回归模型, 请直接计算测试集上的 MSE;
- (b) 对于岭回归问题, 请考察不同正则项权重 λ 的取值范围, 并观察训练集 MSE、测试集 MSE 和 λ 的取值的关系, 总结变化的规律; 除示例代码中使用到的 sklearn 库函数外, 不能使用其他的 sklearn 函数, 需要基于 numpy 实现线性回归模型和 MSE 的计算.

解:

$$1. \text{ 令 } E = \frac{1}{2} \|Xw + 1b - y\|_2^2 + \lambda \|w\|_2^2$$

对 w 求导得

$$\frac{\partial E}{\partial w} = X^\top (Xw + 1b - y) + 2\lambda w$$

对 b 求导得

$$\frac{\partial E}{\partial b} = 1^\top (Xw + 1b - y)$$

当 $(X^\top X + 2\lambda I_d)$ 是满秩矩阵时, 令该式等于零即可得

$$w_{\text{Ridge}}^* = (X^\top X + 2\lambda I_d)^{-1} X^\top (y - 1b_{\text{Ridge}}^*)$$

令 $T = (X^\top X + 2\lambda I_d)^{-1}$, 则

$$w_{\text{Ridge}}^* = TX^\top (y - 1b_{\text{Ridge}}^*)$$

代入有

$$\begin{aligned} b_{\text{Ridge}}^* &= 1^\top y - 1^\top X w_{\text{Ridge}}^* \\ &= 1^\top y - 1^\top X T X^\top (y - 1b_{\text{Ridge}}^*) \\ &= 1^\top y - 1^\top X T X^\top y + 1^\top X T X^\top 1b_{\text{Ridge}}^* \end{aligned}$$

最后有

$$b_{\text{Ridge}}^* = \frac{\mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{y} - \mathbf{1}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{1} - 1}$$

最后将 b_{Ridge}^* 带回即可求出

$$\mathbf{w}_{\text{Ridge}}^* = \mathbf{T} \mathbf{X}^\top (\mathbf{y} - \mathbf{1} \frac{\mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{y} - \mathbf{1}^\top \mathbf{y}}{\mathbf{1}^\top \mathbf{X} \mathbf{T} \mathbf{X}^\top \mathbf{1} - 1})$$

加上第二题的结果, 可以看出, 使用矩阵 \mathbf{T} 抽象之后, 最优解和原始最优解的形式是一样的, 不同的是 \mathbf{T} 的值, 二者的矩阵 \mathbf{T} 括号内相差了一个 $2\lambda \mathbf{I}_d$ 项.

所以我们可以认为 \mathbf{w}_{LS}^* 和 b_{LS}^* 是 $\lambda = 0$ 的 $\mathbf{w}_{\text{Ridge}}^*$ 和 b_{Ridge}^* 特殊情况.

2. 将恒等式

$$\mathbf{X} = \mathbf{X}$$

右端乘上 $(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}$ 即 \mathbf{I}_d 则有

$$\mathbf{X} = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}$$

将 \mathbf{X} 乘入括号内则有

$$\mathbf{X} = (\mathbf{X} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{X})(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}$$

再将 \mathbf{X} 提出至右侧则有

$$\mathbf{X} = (\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_m) \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}$$

最后则有

$$(\mathbf{X} \mathbf{X}^\top + \lambda \mathbf{I}_m)^{-1} \mathbf{X} = \mathbf{X} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1}$$

式子成立.

这个结论能够帮助岭回归的计算.

当样例的维度 d 大于样例数目 m 的时候, 将 $\mathbf{X}(\mathbf{X}^\top \mathbf{X} + 2\lambda I_d)^{-1}$ 转为 $(\mathbf{X}\mathbf{X}^\top + 2\lambda I_m)^{-1}\mathbf{X}$ 能够将矩阵求逆求 \mathbf{T} 的矩阵维度减少, 进而加快矩阵求逆的速度.

3. (a)

算出线性回归模型测试集上的 MSE 结果为 20.724023437336253.

(b)

取不同的权重所得结果为

λ	0.0	0.2	0.4	0.6	0.8	1.0	1.2	1.4	1.6	1.8	2.0
MSE	20.72	20.92	21.09	21.21	21.31	21.39	21.4	21.50	21.54	21.57	21.60

Table 1: lambda

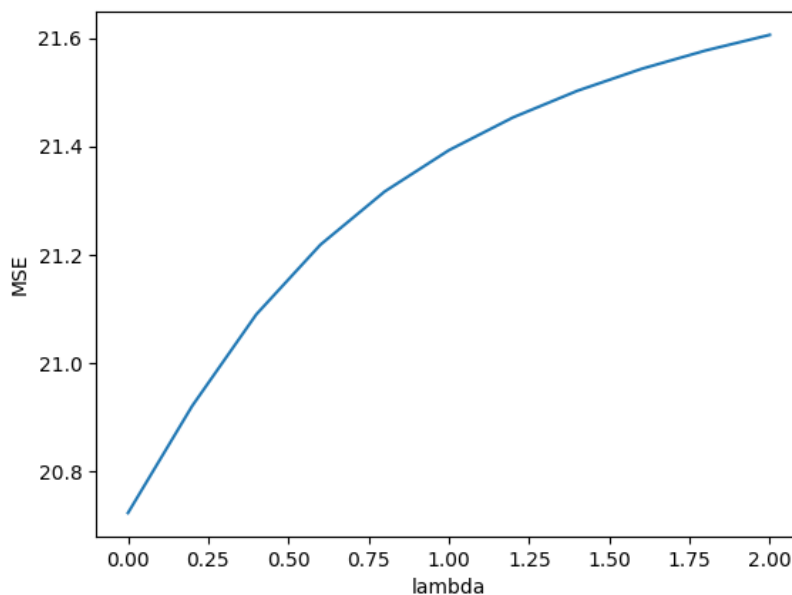


Figure 1: P-R 曲线

可以看出, MSE 随着 λ 的增大而缓慢增大, 所以我们可以选择一个相对较小但仍然为正数的 λ 值.

四. (20 points) 线性判别分析

教材 3.4 节介绍了“线性判别分析”模型 LDA (Linear Discriminative Analysis), 本题首先针对 LDA 从分布假设的角度进行推导和分析. 考虑 N 分类问题, 训练集 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, 其中, 第 n 类样例从高斯分布 $\mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ 中独立同分布采样得到 (其中, $n = 1, 2, \dots, N$). 记该类样例数量为 m_n . 类别先验为 $p(y = n) = \pi_n$, 反映了各类别出现的概率. 若 $\mathbf{x} \in \mathbb{R}^d \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 则其概率密度函数为

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{d}{2}} \det(\boldsymbol{\Sigma})^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (5)$$

假设不同类别的条件概率为高斯分布, 当不同类别的协方差矩阵 $\boldsymbol{\Sigma}_n$ 相同时, 对于类别的预测转化为类别中心之间的线性问题, 下面对这一模型进行进一步分析. 假设 $\boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}$, 分析 LDA 的分类方式以及参数估计步骤.

1. 样例 \mathbf{x} 的后验概率 $p(y = n | \mathbf{x})$ 表示了样例属于第 n 类的可能性, 当计算样例针对 N 个类别的后验概率后, 找出后验概率最大的类别对样例的标记进行预测, 即 $\arg \max_n p(y = n | \mathbf{x})$. 等价于考察 $\ln p(y = n | \mathbf{x})$ 的大小, 请证明在此假设下,

$$\arg \max_y p(y | \mathbf{x}) = \arg \max_n \underbrace{\mathbf{x}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_n - \frac{1}{2} \boldsymbol{\mu}_n^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_n + \ln \pi_n}_{\delta_n(\mathbf{x})}. \quad (6)$$

其中 $\delta_n(\mathbf{x})$ 为 LDA 在分类时的判别函数.

2. 在 LDA 模型中, 需要估计各类别的先验概率, 以及条件概率中高斯分布的参数. 针对二分类问题 ($N = 2$), 使用如下方式估计类别先验、均值与协方差矩阵:

$$\hat{\pi}_n = \frac{m_n}{m}; \quad \hat{\boldsymbol{\mu}}_n = \frac{1}{m_n} \sum_{y_i=n} \mathbf{x}_i, \quad (7)$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{m - N} \sum_{n=1}^N \sum_{y_i=n} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_n) (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_n)^\top. \quad (8)$$

LDA 使用这些经验量替代真实参数, 计算判别式 $\delta_n(\mathbf{x})$ 并按照第??问中的准则做出预测. 请证明:

$$\mathbf{x}^\top \hat{\boldsymbol{\Sigma}}^{-1} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) > \frac{1}{2} (\hat{\boldsymbol{\mu}}_2 + \hat{\boldsymbol{\mu}}_1)^\top \hat{\boldsymbol{\Sigma}}^{-1} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) - \ln(m_2/m_1) \quad (9)$$

时 LDA 将样例预测为第 2 类. 请分析这一判别方式的几何意义.

3. 在 LDA 中, 对样例 \mathbf{x} 的判别可视为在投影的空间中和某个阈值进行比较. 上述推导通过最大后验概率的方法得到对投影后样例分布的需求, 而 Fisher 判别分析 (Fisher Discriminant Analysis, FDA) 也是一种常见的线性判别分析方法, 直接对样例投影后数据的分布情况进行约束. FDA 一般通过广义瑞利商进行求解, 请基于教材 3.4 节对“线性判别分析”的介绍, 对广义瑞利商的性质进行分析, 探讨 FDA 多分类推广的性质. 下面请说明对于 N 类分类问题, FDA 投影的维度最多为 $N - 1$, 即投影矩阵 $\mathbf{W} \in \mathbb{R}^{d \times (N-1)}$.

提示: 矩阵的秩具有如下性质: 对于矩阵 $\mathbf{A} \in \mathbb{R}^{m \times n}$, 矩阵 $\mathbf{B} \in \mathbb{R}^{n \times r}$, 则

$$\text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B}) - n \leq \text{rank}(\mathbf{AB}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}. \quad (10)$$

对于任意矩阵 \mathbf{A} , 以下公式成立

$$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^\top) = \text{rank}(\mathbf{AA}^\top) = \text{rank}(\mathbf{A}^\top\mathbf{A}). \quad (11)$$

解:

1. 由贝叶斯公式可知

$$\begin{aligned} & p(y = n|\mathbf{x}) \\ = & \frac{p(y = n)p(\mathbf{x}|y = n)}{p(\mathbf{x})} \\ = & \frac{\pi_n((2\pi)^{\frac{d}{2}} \det(\mathbf{\Sigma})^{\frac{1}{2}})^{-1} \exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_n)^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_n))}{p(\mathbf{x})} \\ = & \frac{\pi_n((2\pi)^{\frac{d}{2}} \det(\mathbf{\Sigma})^{\frac{1}{2}})^{-1} \exp(\mathbf{x}^\top \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_n - \frac{1}{2} \boldsymbol{\mu}_n^\top \mathbf{\Sigma}^{-1} \boldsymbol{\mu}_n - \frac{1}{2} \mathbf{x}^\top \mathbf{\Sigma}^{-1} \mathbf{x})}{p(\mathbf{x})} \end{aligned}$$

并且 $\arg \max_n p(y = m|\mathbf{x})$ 可以转化为

$$\arg \max_n \ln p(y = m|\mathbf{x})$$

即有

$$\begin{aligned} \arg \max_n \ln \pi_n + \ln((2\pi)^{\frac{d}{2}} \det(\Sigma)^{\frac{1}{2}})^{-1} + \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_n \\ - \frac{1}{2} \boldsymbol{\mu}_n^\top \Sigma^{-1} \boldsymbol{\mu}_n - \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} - \ln p(\mathbf{x}) \end{aligned}$$

去除与 n 无关的项即可得

$$\arg \max_n \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_n - \frac{1}{2} \boldsymbol{\mu}_n^\top \Sigma^{-1} \boldsymbol{\mu}_n + \ln \pi_n$$

2. 将式子

$$\mathbf{x}^\top \hat{\Sigma}^{-1} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) > \frac{1}{2} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1)^\top \hat{\Sigma}^{-1} (\hat{\boldsymbol{\mu}}_2 - \hat{\boldsymbol{\mu}}_1) - \ln \frac{m_2}{m_1}$$

拆开可得

$$\mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_2 - \mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 > \frac{1}{2} \hat{\boldsymbol{\mu}}_2^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_2 - \frac{1}{2} \hat{\boldsymbol{\mu}}_1^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 - \ln \frac{m_2}{m_1}$$

进行移项

$$\mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_2 - \frac{1}{2} \hat{\boldsymbol{\mu}}_2^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_2 + \ln \frac{m_2}{m} > \mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 - \frac{1}{2} \hat{\boldsymbol{\mu}}_1^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 + \ln \frac{m_1}{m}$$

将 $\ln \frac{m_n}{m}$ 替换为 $\ln \hat{\pi}_n$ 可得

$$\mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_2 - \frac{1}{2} \hat{\boldsymbol{\mu}}_2^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_2 + \ln \hat{\pi}_2 > \mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 - \frac{1}{2} \hat{\boldsymbol{\mu}}_1^\top \hat{\Sigma}^{-1} \hat{\boldsymbol{\mu}}_1 + \ln \hat{\pi}_1$$

即有 $n^* = \arg \max_n \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_n - \frac{1}{2} \boldsymbol{\mu}_n^\top \Sigma^{-1} \boldsymbol{\mu}_n + \ln \pi_n = 2$

所以此时 LDA 将样例预测为第 2 类.

几何意义:

$\frac{1}{2}(\hat{\boldsymbol{\mu}}_2 + \hat{\boldsymbol{\mu}}_1)$ 是 $\hat{\boldsymbol{\mu}}_1$ 和 $\hat{\boldsymbol{\mu}}_2$ 这两个类别中心的中点.

而 $x^\top \hat{\Sigma}^{-1}(\hat{\mu}_2 - \hat{\mu}_1)$ 可以化为 $(\hat{\Sigma}^{-\frac{1}{2}}x)^\top (\hat{\Sigma}^{-\frac{1}{2}}(\hat{\mu}_2 - \hat{\mu}_1))$

即将特殊的多维正态分布缩放旋转成了标准多维正态分布, 然后衡量离第 1 类的类别中心 $\hat{\mu}_1$ 的远近.

若大于号成立, 说明 x 比中点 $\frac{1}{2}(\hat{\mu}_1 + \hat{\mu}_2)$ 离 $\hat{\mu}_1$ 更远, 也就是离 $\hat{\mu}_2$ 更近.

但是还要考虑一个因素, 如果第 2 类的样例比第 1 类的样例数目多, 就要加入一定的偏好, 比如这里的对数几率 $\ln(m_2/m_1)$, 当 $m_2 > m_1$ 时是正数, 就能更容易判断为第 2 类.

3. 求解

$$\max_{\mathbf{W}} \frac{\text{tr}(\mathbf{W}^\top \mathbf{S}_b \mathbf{W})}{\text{tr}(\mathbf{W}^\top \mathbf{S}_w \mathbf{W})}$$

最后可以转化为广义特征值问题求解

$$\mathbf{S}_b \mathbf{W} = \lambda \mathbf{S}_w \mathbf{W}$$

即求解

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{W} = \lambda \mathbf{W}$$

也就是求 $\mathbf{S}_w^{-1} \mathbf{S}_b$ 特征值最大的前几个的特征向量.

而我们知道

$$\mathbf{S}_b = \sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top$$

其中我们知道

$$\sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu}) = \sum_{i=1}^N m_i \boldsymbol{\mu}_i - \boldsymbol{\mu} \sum_{i=1}^N m_i = \mathbf{0}$$

也就是

$$\boldsymbol{\mu}_N - \boldsymbol{\mu} = -\frac{1}{m_N} \sum_{i=1}^{N-1} m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})$$

即 $(\boldsymbol{\mu}_N - \boldsymbol{\mu})$ 可以由 $(\boldsymbol{\mu}_i - \boldsymbol{\mu})$ 线性表示出来.
因此我们有

$$\begin{aligned} \text{rank}(\mathbf{S}_w^{-1} \mathbf{S}_b) &\leq \min\{\text{rank}(\mathbf{S}_w^{-1}), \text{rank}(\mathbf{S}_b)\} \\ &\leq \text{rank}(\mathbf{S}_b) \\ &= \text{rank}\left(\sum_{i=1}^N m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top\right) \\ &= \text{rank}\left(\sum_{i=1}^{N-1} m_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top\right) \\ &\leq \sum_{i=1}^{N-1} \text{rank}((\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top) \\ &= \sum_{i=1}^{N-1} \text{rank}(\boldsymbol{\mu}_i - \boldsymbol{\mu}) \\ &= N - 1 \end{aligned}$$

因此 $\mathbf{S}_w^{-1} \mathbf{S}_b$ 最多有 $N - 1$ 个非零特征值, 对应 $N - 1$ 个特征向量. 而 \mathbf{W} 的每个列向量都是 $\mathbf{S}_w^{-1} \mathbf{S}_b$ 的线性无关的特征向量. 所以 FDA 投影的维度最多为 $N - 1$.

五. (20 points) 多分类学习

教材 3.5 节介绍了“多分类学习”的多种方式, 本题针对 OvO 和 OvR 两种多分类学习方法进行分析:

1. 分析两种多分类方法的优劣. 思考这两种多分类推广方式是否存在难以处理的情况?
2. 在 OvR 的每一个二分类子任务中, 目标类别作为正类, 而其余所有类别作为负类. 此时, 是否需要显式考虑正负类别的不平衡带来的影响?

解:

1. OvO 的优势: OvO 的每个分类器仅用到两个类的样例, 在类别很多的时候, OvO 的训练时间开销通常比 OvR 小.

OvO 的劣势: OvO 需要训练 $N(N-1)/2$ 个分类器, 所以存储开销和训练时间开销通常比较大.

OvR 的优势: OvR 只需要训练 N 个分类器.

OvR 的劣势: 每次训练时都用到了所有的数据, 在类别很多的时候, 训练时间开销更大.

这两种多分类推广方式都可能存在难以处理的情况.

例如对 OvR 来说, 如果类别特别多的时候, 每次都要使用全部的数据来处理, 时间开销就接近 $O(n^2)$ 了, 几乎是不可接受的.

即使是对于 OvO 来说, 我们也能构造出一个类似的例子, 例如有一半的样例都是属于第 0 类, 但是剩下的一半几乎每个样例都对应一个类别, 这样时间开销和空间开销也都到达了 $O(n^2)$ 了, 也几乎是不可接受的.

2. 对 OvR 来说, 由于对每一个类进行了相同的处理, 其拆解出来的二分类任务中类别不平衡的影响会相互抵消, 因此通常不需要显式考虑正负类别不平衡带来的影响.