



中国科学院大学

University of Chinese Academy of Sciences

模式识别与机器学习

081203M04004H

Chap 6 课程作业解答

2022 年 11 月 3 号

Professor: 黄庆明



学生: 周胤昌

学号: 202228018670052

学院: 网络安全学院

所属专业: 网络安全

方向: 安全协议理论与技术

Problem 1

给定训练数据集 $X = \begin{pmatrix} 1 & 2 & 5 & 4 \\ 2 & 5 & 1 & 2 \end{pmatrix}$, $y = (19 \ 26 \ 19 \ 20)^T$, 令 $\alpha = 0.001$, $w_0 = (1 \ 1)^T$. 编程实现 SGD 和 GD 算法, 求解 w .

Solution: 最优化问题 (即代价函数) 为

$$\min_w J(w) = \frac{1}{2N} \sum_{i=1}^N (w^T x^{(i)} - y^{(i)})^2$$

于是批梯度下降 (BGD) 和梯度迭代更新规则分别为 (具体的 python 训练代码见如下)

$$\frac{\partial J(w)}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N x_j^{(i)} (w^T x^{(i)} - y^{(i)}), \quad w_j \leftarrow w_j - \frac{\alpha}{N} \sum_{i=1}^N (w^T x^{(i)} - y^{(i)}) x_j^{(i)}, \alpha > 0$$

```

1 # 批量梯度下降 BGD
2 # 拟合函数为: y = theta * x
3 # 代价函数为: J = 1 / (2 * m) * ((theta * x) - y) * ((theta * x) - y).T;
4 # 梯度迭代为: theta = theta - alpha / m * (x * (theta * x - y).T);
5 import time
6 import numpy as np
7 # 1、多元线性回归的 BGD 程序
8 def bgd_multi():
9     # 训练集, 每个样本有 2 个分量
10    x = np.array([(1, 2), (2, 5), (5, 1), (4, 2)])
11    y = np.array([19, 26, 19, 20])
12    # 初始化
13    m, dim = x.shape
14    theta = np.ones(dim) # 参数
15    alpha = 0.001 # 学习率
16    threshold = 0.0001 # 停止迭代的错误阈值
17    iterations = 1500 # 迭代次数
18    error = 0 # 初始错误为 0
19    # 迭代开始
20    for i in range(iterations):
21        error = 1 / (2 * m) * np.dot((np.dot(x, theta) - y).T, (np.dot(x, theta) - y))
22        # 迭代停止
23        if abs(error) <= threshold:
24            break
25        theta -= alpha / m * (np.dot(x.T, (np.dot(x, theta) - y)))
26    print('BGD 的迭代次数为%d,' % (i + 1), 'theta:', theta, ', error: %f' % error)
27 if __name__ == '__main__':
28     start = time.time()
29     bgd_multi()
30     end = time.time()
31     print('运行时间为: ', (end - start) * 1000, 'ms')

```

上述代码的输出结果为 $w = (2.868 \ 4.565)^T$, 循环迭代次数为 1500, 循环终止时的线性拟合误差为 $\epsilon = 6.995$, BGD 算法运行时间为 13.84ms.

而对于 SGD 算法, 最优化问题 (即代价函数) 仍为

$$\min_w J(w) = \frac{1}{2N} \sum_{i=1}^N (w^T x^{(i)} - y^{(i)})^2$$

于是随机梯度下降 (SGD) 和梯度迭代更新规则分别为 (具体的 python 训练代码见如下)

$$\frac{\partial J(\mathbf{w})}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N x_j^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)}), \quad w_j \leftarrow w_j - \alpha (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)}) x_j^{(i)}, \alpha > 0$$

```

1  # 随机梯度下降 SGD
2  import time
3  import numpy as np
4  # 2、多元线性回归的 SGD 程序
5  def sgd():
6      # 训练集, 每个样本有 2 个分量
7      x = np.array([(1, 2), (2, 5), (5, 1), (4, 2)])
8      y = np.array([19, 26, 19, 20])
9      # 初始化
10     m, dim = x.shape
11     theta = np.ones(dim) # 参数
12     alpha = 0.001 # 学习率
13     threshold = 0.0001 # 停止迭代的错误阈值
14     iterations = 1500 # 迭代次数
15     error = 0 # 初始错误为 0
16     # 迭代开始
17     for i in range(iterations):
18         error = 1 / (2 * m) * np.dot((np.dot(x, theta) - y).T, (np.dot(x,
19         ↪ theta) - y))
20         # 迭代停止
21         if abs(error) <= threshold:
22             break
23         j = np.random.randint(0, m)
24         theta -= alpha * (x[j] * (np.dot(x[j], theta) - y[j]))
25     print('SGD 的迭代次数为%d,' % (i + 1), 'theta:', theta, ', error: %f' %
26     ↪ error)
27 if __name__ == '__main__':
28     start = time.time()
29     sgd()
30     end = time.time()
31     print('运行时间为: ', (end - start) * 1000, 'ms')

```

上述代码的输出结果为 $\mathbf{w} = (2.880 \ 4.613)^T$, 循环迭代次数为 1500, 循环终止时的线性拟合误差为 $\epsilon = 7.01$, SGD 算法的运行时间为 13.51ms. 可以看出 SGD 算法的速度是略快于 BGD 算法的.

Problem 2

利用下面表格 1 中的训练数据训练朴素贝叶斯分类器. 给定测试样本 $\mathbf{x} = (2, S)^T$ 和 $\mathbf{x} = (1, N)^T$, 请预测他们的标签.

表 1: Problem 2 的训练数据

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
x_1	1	1	1	1	1	2	2	2	2	2	3	3	3	3	3
x_2	S	M	M	S	S	S	M	M	L	L	L	M	M	L	L
y	-1	-1	1	1	-1	-1	-1	1	1	1	1	1	1	1	-1

由于测试样本中出现了训练样本中未出现的特征分量, 所以需要对朴素贝叶斯分类器参数的极大似然估计做 *Laplace* 平滑 (其中平滑系数 λ 选取为 1):

$$p(y = 1) = \frac{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + 1}{N + K} = \frac{9 + 1}{15 + 2} = \frac{10}{17} \Rightarrow p(y = -1) = 1 - \frac{10}{17} = \frac{7}{17}$$

第 1 个特征分量的似然函数为

$$p(x_1 = 1|y = 1) = \frac{\sum_{i=1}^N I_{\{x_1^{(i)}=1 \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_1} = \frac{2 + 1}{9 + 3} = \frac{1}{4},$$

$$p(x_1 = 2|y = 1) = \frac{\sum_{i=1}^N I_{\{x_1^{(i)}=2 \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_1} = \frac{3 + 1}{9 + 3} = \frac{1}{3},$$

$$p(x_1 = 3|y = 1) = \frac{\sum_{i=1}^N I_{\{x_1^{(i)}=3 \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_1} = \frac{4 + 1}{9 + 3} = \frac{5}{12}$$

第 2 个特征分量的似然函数为:

$$p(x_2 = S|y = 1) = \frac{\sum_{i=1}^N I_{\{x_2^{(i)}=S \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_2} = \frac{1 + 1}{9 + 4} = \frac{2}{13},$$

$$p(x_2 = M|y = 1) = \frac{\sum_{i=1}^N I_{\{x_2^{(i)}=M \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_2} = \frac{4 + 1}{9 + 4} = \frac{5}{13},$$

$$p(x_2 = L|y = 1) = \frac{\sum_{i=1}^N I_{\{x_2^{(i)}=L \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_2} = \frac{4+1}{9+4} = \frac{5}{13},$$

$$p(x_2 = N|y = 1) = \frac{\sum_{i=1}^N I_{\{x_2^{(i)}=N \wedge y^{(i)}=1\}} + 1}{\sum_{i=1}^N I_{\{y^{(i)}=1\}} + S_2} = \frac{0+1}{9+4} = \frac{1}{13}$$

预测表达式为

$$p(y = 1|\mathbf{x}) = \frac{p(\mathbf{x}|y = 1)p(y = 1)}{p(\mathbf{x})} = \frac{\prod_{j=1}^D p(x_j|y = 1)p(y = 1)}{\prod_{j=1}^D p(x_j|y = 1)p(y = 1) + \prod_{j=1}^D p(x_j|y = -1)p(y = -1)}$$

于是对于样本 $\mathbf{x}^{(1)} = (2, S)^T$, 预测概率为

$$p(y = 1|\mathbf{x}^{(1)}) = \frac{\frac{4}{12} \times \frac{2}{13} \times \frac{10}{17}}{\frac{4}{12} \times \frac{2}{13} \times \frac{10}{17} + \frac{2+1}{6+3} \times \frac{3+1}{6+4} \times \frac{7}{17}} = 0.35461$$

$$p(y = -1|\mathbf{x}^{(1)}) = \frac{\frac{2+1}{6+3} \times \frac{3+1}{6+4} \times \frac{7}{17}}{\frac{4}{12} \times \frac{2}{13} \times \frac{10}{17} + \frac{2+1}{6+3} \times \frac{3+1}{6+4} \times \frac{7}{17}} = 0.64539$$

所以样本 $\mathbf{x}^{(1)} = (2, S)^T$ 的标签预测为 $y = -1$.

而对于样本 $\mathbf{x}^{(2)} = (1, N)^T$, 预测概率为

$$p(y = 1|\mathbf{x}^{(2)}) = \frac{\frac{2+1}{9+3} \times \frac{0+1}{9+4} \times \frac{10}{17}}{\frac{2+1}{9+3} \times \frac{0+1}{9+4} \times \frac{10}{17} + \frac{3+1}{6+3} \times \frac{0+1}{6+4} \times \frac{7}{17}} = 0.382003$$

$$p(y = -1|\mathbf{x}^{(2)}) = \frac{\frac{3+1}{6+3} \times \frac{0+1}{6+4} \times \frac{7}{17}}{\frac{2+1}{9+3} \times \frac{0+1}{9+4} \times \frac{10}{17} + \frac{3+1}{6+3} \times \frac{0+1}{6+4} \times \frac{7}{17}} = 0.617997$$

所以样本 $\mathbf{x}^{(2)} = (1, N)^T$ 的标签预测为 $y = -1$. 上述算法可以用 Python 代码实现, 不妨将 S, M, L, N 分别记为 1, 2, 3, 4, 具体代码如下:

```
1 import numpy as np
2 from sklearn.naive_bayes import MultinomialNB
3 X = np.array([[1, 1], [1, 2], [1, 2], [1, 1], [1, 1], [2, 1], [2, 2], [2, 2],
4               [2, 3], [2, 3], [3, 3], [3, 2], [3, 2], [3, 3], [3, 3]])
5 # S, M, L, N 分别用数字 1, 2, 3, 4 替代
6 Y = np.array([-1, -1, 1, 1, -1, -1, -1, 1, 1, 1, 1, 1, 1, 1, -1])
7 clf = MultinomialNB(alpha = 1.0, fit_prior = True, class_prior = None)
8 clf.fit(X, Y)
9 clf.predict_proba([[2, 1], [1, 4]]) # 输出 (2, S) 和 (1, N) 划分到各个类别的概率值
```

经过上述代码可得: 对于样本 $\mathbf{x}^{(1)} = (2, S)^T$, 预测概率为 $p(y = 1|\mathbf{x}^{(1)}) = 0.406$, $p(y = -1|\mathbf{x}^{(1)}) = 0.594$, 因此标签为 $y = -1$; 对于样本 $\mathbf{x}^{(2)} = (1, N)^T$, 预测概率为 $p(y = 1|\mathbf{x}^{(2)}) = 0.384$, $p(y = -1|\mathbf{x}^{(2)}) = 0.616$, 故标签为 $y = -1$. 且经过 MultinomialNB 后所习得的 (平滑) 先验概率的对数为 $-0.511, -0.916$.

上述算法也可以用 Matlab 代码实现, 具体代码见下所示:

```

1  n=input('请输入训练集的个数:n\n');
2  k=2; A=zeros(n,3);% 存储训练集, 用来学习
3  A(:,1)=input('请输入所有样本的第一个特征 (用列向量表示): \n');
4  A(:,2)=input('请输入所有样本的第二个特征 (用列向量表示): \n');
5  A(:,3)=input('请输入所有样本所属的类 (用列向量表示): \n');
6  x1=input('请输入需要预测的样本的第一个特征: 1 或 2 或 3 \n');
7  x2=input('请输入需要预测的样本的第二个特征: S 或 M 或 L 或 N\n','s');
8  % 计算其属于 1 类的概率
9  F1=zeros(1,3);
10 for i=1:n
11     if A(i,3)==1
12         F1(1,1)=F1(1,1)+1;
13         if A(i,1)==x1
14             F1(1,2)=F1(1,2)+1;
15         end
16         if A(i,2)==x2
17             F1(1,3)=F1(1,3)+1;
18         end
19     end
20 end
21 C1=(F1(1,1)+1)/(n+k)*(F1(1,2)+1)/(F1(1,1)+3)*(F1(1,3)+1)/(F1(1,1)+4);
22 % 属于 2 类的概率
23 F2=zeros(1,3);
24 for i=1:n
25     if A(i,3)==-1
26         F2(1,1)=F2(1,1)+1;
27         if A(i,1)==x1
28             F2(1,2)=F2(1,2)+1;
29         end
30         if A(i,2)==x2
31             F2(1,3)=F2(1,3)+1;
32         end
33     end
34 end
35 C2=(F2(1,1)+1)/(n+k)*(F2(1,2)+1)/(F2(1,1)+3)*(F2(1,3)+1)/(F2(1,1)+4);
36 if C1>C2
37     disp('该样本的类为 1');
38 else if C1==C2
39     disp('该样本的类为-1 或 1')
40 else
41     disp('该样本的类为-1')
42 end
43 end

```

在命令行窗口的输出过程为:

```

1 >> MultinomialNB
2 请输入训练集的个数:n
3 15
4 请输入所有样本的第一个特征 (用列向量表示):
5 [1,1,1,1,1,2,2,2,2,2,3,3,3,3,3]
6 请输入所有样本的第二个特征(用列向量表示):
7 ['S','M','M','S','S','S','M','M','L','L','L','M','M','L','L']
8 请输入所有样本所属的类(用列向量表示):
9 [-1,-1,1,1,-1,-1,-1,1,1,1,1,1,1,1,-1]
10 请输入需要预测的样本的第一个特征: 1 或 2 或 3
11 1
12 请输入需要预测的样本的第二个特征: S 或 M 或 L 或 N
13 N
14 该样本的类为-1

```

```

1 >> MultinomialNB
2 请输入训练集的个数:n
3 15
4 请输入所有样本的第一个特征 (用列向量表示):
5 [1,1,1,1,1,2,2,2,2,2,3,3,3,3,3]
6 请输入所有样本的第二个特征(用列向量表示):
7 ['S','M','M','S','S','S','M','M','L','L','L','M','M','L','L']
8 请输入所有样本所属的类(用列向量表示):
9 [-1,-1,1,1,-1,-1,-1,1,1,1,1,1,1,1,-1]
10 请输入需要预测的样本的第一个特征: 1 或 2 或 3
11 2
12 请输入需要预测的样本的第二个特征: S 或 M 或 L 或 N
13 S
14 该样本的类为-1

```

通过 Matlab 右侧的变量值栏目可以计算出: 测试样本 $\mathbf{x}^{(1)} = (2, S)^T$ 的预测概率为 $p(y = 1|\mathbf{x}^{(1)}) = 0.3546$, $p(y = -1|\mathbf{x}^{(1)}) = 0.6454$, 因此标签为 $y = -1$. $\mathbf{x}^{(2)} = (1, N)^T$ 的预测概率为 $p(y = 1|\mathbf{x}^{(2)}) = 0.3820$, $p(y = -1|\mathbf{x}^{(2)}) = 0.6180$, 因此标签为 $y = -1$.

Problem 3

生成式判别模型: 高斯贝叶斯分类器和逻辑回归

Part A

考虑一类特定的高斯朴素贝叶斯分类器, 其中

- Y 是服从伯努利分布的布尔变量, 其中参数 $\pi = P(Y = 1)$, $P(Y = 0) = 1 - \pi$;
- $X = [x_1, \dots, x_D]^T$, 其中特征分量 x_i 是连续的随机变量. 对于每个 x_i , $P(x_i|Y = k)$ 是高斯分布 $\mathcal{N}(\mu_{ik}, \sigma_i)$. 注意到 σ_i 是高斯分布的标准差 (且不依赖于 k);
- 给定 Y 时, $\forall i \neq j$, x_i 和 x_j 都是条件独立的 (因此称之为朴素分类器).

问题: 请证明判别式分类器 (如逻辑回归) 与上述特定类别的高斯朴素贝叶斯分类器之间的关系正是逻辑回归所使用的形式.

Solution: 先利用贝叶斯公式:

$$\begin{aligned} P(Y = 1|X) &= \frac{P(X|Y = 1) \cdot P(Y = 1)}{P(X|Y = 0) \cdot P(Y = 0) + P(X|Y = 1) \cdot P(Y = 1)} \\ &= \frac{1}{1 + \frac{P(X|Y = 0) \cdot P(Y = 0)}{P(X|Y = 1) \cdot P(Y = 1)}} \\ &= \frac{1}{1 + \exp\left(\ln \frac{P(X|Y = 0) \cdot P(Y = 0)}{P(X|Y = 1) \cdot P(Y = 1)}\right)} \end{aligned}$$

再代入 $P(Y = 1) = \pi$, $P(Y = 0) = 1 - \pi$, 同时将 \ln 中的乘法变为加法:

$$P(Y = 1|X) = \frac{1}{1 + \exp\left(\ln \frac{P(X|Y = 0)}{P(X|Y = 1)} + \ln \frac{P(Y = 0)}{P(Y = 1)}\right)} = \frac{1}{1 + \exp\left(\ln \frac{P(X|Y = 0)}{P(X|Y = 1)} + \ln \frac{1 - \pi}{\pi}\right)}$$

因为 X 是 D 维的, 同时每一个特征都是相互条件独立的, 因此 $P(X|Y) = \prod_{i=1}^D P(x_i|Y)$, 于是有

$$P(Y = 1|X) = \frac{1}{1 + \exp\left(\ln \frac{\prod_{i=1}^D P(x_i|Y = 0)}{\prod_{i=1}^D P(x_i|Y = 1)} + \ln \frac{1 - \pi}{\pi}\right)}$$

由于 $\forall x_i$, $P(x_i|Y = k)$ 都服从高斯分布 $\mathcal{N}(\mu_{ik}, \sigma_i)$, 即 $P(x_i|Y = k) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(x_i - \mu_{ik})^2}{2\sigma_i^2}\right)$. 我

我们先看 $\ln \frac{\prod_{i=1}^D P(x_i|Y=0)}{\prod_{i=1}^D P(x_i|Y=1)}$:

$$\begin{aligned}
\ln \frac{\prod_{i=1}^D P(x_i|Y=0)}{\prod_{i=1}^D P(x_i|Y=1)} &= \ln \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2}\right) - \ln \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}\right) \\
&= \sum_{i=1}^D \left(\ln \frac{1}{\sqrt{2\pi}\sigma_i} - \frac{(x_i - \mu_{i0})^2}{2\sigma_i^2} \right) - \sum_{i=1}^D \left(\ln \frac{1}{\sqrt{2\pi}\sigma_i} - \frac{(x_i - \mu_{i1})^2}{2\sigma_i^2} \right) \\
&= \sum_{i=1}^D \left(\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2} - \frac{(x_i - \mu_{i0})^2}{2\sigma_i^2} \right) = \sum_{i=1}^D \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} x_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right)
\end{aligned}$$

再将其代入到 $P(Y=1|X)$ 的表达式中:

$$\begin{aligned}
P(Y=1|X) &= \frac{1}{1 + \exp \left\{ \sum_{i=1}^D \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} x_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right) + \ln \frac{1 - \pi}{\pi} \right\}} \\
&= \frac{1}{1 + \exp \left\{ - \left[\sum_{i=1}^D \underbrace{\frac{\mu_{i1} - \mu_{i0}}{\sigma_i^2} x_i}_{\text{即 } w_i} + \underbrace{\left(\sum_{i=1}^D \frac{\mu_{i0}^2 - \mu_{i1}^2}{2\sigma_i^2} + \ln \frac{\pi}{1 - \pi} \right)}_{\text{即 } b} \right] \right\}} \\
&= \frac{1}{1 + \exp \left[- \left(\sum_{i=1}^D w_i x_i + b \right) \right]}
\end{aligned}$$

同理可得:

$$P(Y=0|X=1) = 1 - P(Y=1|X) = \frac{\exp \left[- \left(\sum_{i=1}^D w_i x_i + b \right) \right]}{1 + \exp \left[- \left(\sum_{i=1}^D w_i x_i + b \right) \right]}$$

也就是说, 上述特定类别的高斯朴素贝叶斯分类器其实正好就是逻辑回归的形式! 逻辑回归的参数 w 通常采用 SGD 算法去学习, 而高斯朴素贝叶斯分类器根据假设直接给出了这些参数, 这也说明了逻辑回归是一种更强的模型, 因为它的假设很弱, 而朴素贝叶斯的假设非常的强.

Part B

一般的高斯朴素贝叶斯分类器和逻辑回归：将“ $P(x_i|Y = k)$ 的标准差 σ_i 不依赖于 k ”的假设删除掉。即 $\forall x_i, P(x_i|Y = k)$ 是高斯分布 $\mathcal{N}(\mu_{ik}, \sigma_{ik})$, 其中 $i = 1, \dots, D$ 且 $k = 0, 1$ 。

问题：更一般的高斯朴素贝叶斯分类器中 $P(Y|X)$ 的表达式是否仍然形如逻辑回归的形式。写出 $P(Y|X)$ 的新格式来证明你的答案。

Solution：其实就是去掉高斯分布中的方差与其所属类别无关的假设，即 $P(x_i|Y = k)$ 服从高斯分布 $\mathcal{N}(\mu_{ik}, \sigma_{ik})$ 。前面的推导与上面一致，即

$$P(Y = 1|X) = \frac{1}{1 + \exp \left(\ln \frac{\prod_{i=1}^D P(x_i|Y = 0)}{\prod_{i=1}^D P(x_i|Y = 1)} + \ln \frac{1 - \pi}{\pi} \right)}$$

由于高斯分布假设更改，即 $P(x_i|Y = k) = \frac{1}{\sqrt{2\pi}\sigma_{ik}} \exp \left(-\frac{(x_i - \mu_{ik})^2}{2\sigma_{ik}^2} \right)$ ，将其代入 $\ln \frac{\prod_{i=1}^D P(x_i|Y = 0)}{\prod_{i=1}^D P(x_i|Y = 1)}$

得到：

$$\begin{aligned} \ln \frac{\prod_{i=1}^D P(x_i|Y = 0)}{\prod_{i=1}^D P(x_i|Y = 1)} &= \ln \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_{i0}} \exp \left(-\frac{(x_i - \mu_{i0})^2}{2\sigma_{i0}^2} \right) - \ln \prod_{i=1}^D \frac{1}{\sqrt{2\pi}\sigma_{i1}} \exp \left(-\frac{(x_i - \mu_{i1})^2}{2\sigma_{i1}^2} \right) \\ &= \sum_{i=1}^D \left[\ln \frac{1}{\sqrt{2\pi}\sigma_{i0}} - \frac{(x_i - \mu_{i0})^2}{2\sigma_{i0}^2} \right] - \sum_{i=1}^D \left[\ln \frac{1}{\sqrt{2\pi}\sigma_{i1}} - \frac{(x_i - \mu_{i1})^2}{2\sigma_{i1}^2} \right] \\ &= \sum_{i=1}^D \left[\left(\frac{1}{2\sigma_{i1}^2} - \frac{1}{2\sigma_{i0}^2} \right) x_i^2 + \left(\frac{\mu_{i0}}{\sigma_{i0}^2} - \frac{\mu_{i1}}{\sigma_{i1}^2} \right) x_i + \left(\frac{\mu_{i1}^2}{2\sigma_{i1}^2} - \frac{\mu_{i0}^2}{2\sigma_{i0}^2} \right) + \ln \frac{\sigma_{i1}}{\sigma_{i0}} \right] \end{aligned}$$

将上式结果代入 $P(Y = 1|X)$ 得到：

$$P(Y = 1|X) = \frac{1}{1 + \exp \left(\sum_{i=1}^D \left[\underbrace{\left(\frac{1}{2\sigma_{i1}^2} - \frac{1}{2\sigma_{i0}^2} \right) x_i^2}_{\text{二次项非零}} + \left(\frac{\mu_{i0}}{\sigma_{i0}^2} - \frac{\mu_{i1}}{\sigma_{i1}^2} \right) x_i + \left(\frac{\mu_{i1}^2}{2\sigma_{i1}^2} - \frac{\mu_{i0}^2}{2\sigma_{i0}^2} \right) + \ln \frac{\sigma_{i1}}{\sigma_{i0}} \right] + \ln \frac{1 - \pi}{\pi} \right)}$$

显然上述形式中的二次项系数非零，即去掉了高斯分布中方差与类别无关的假设后朴素贝叶斯变成了非线性的方式。但是逻辑回归中没有二次项，所以这种假设条件下的高斯朴素贝叶斯分类器无法直接用 **Logistic 回归** 表示。

Part C

高斯贝叶斯分类器和逻辑回归：现在考虑下述假设中的高斯贝叶斯分类器 (不带有“朴素”)：

- Y 是服从伯努利分布的布尔变量, 其中参数 $\pi = P(Y = 1)$, $P(Y = 0) = 1 - \pi$;
- $X = [x_1, x_2]^T$, 即我们只考虑样本的特征维数是 2 的情况, 并且每个特征分量都是一个连续的随机变量. 给定 y 时, x_1, x_2 不再是条件独立的了. 我们假设 $P(x_1, x_2|Y = k)$ 是双变量高斯分布 $\mathcal{N}(\mu_{1k}, \mu_{2k}, \sigma_1, \sigma_2, \rho)$, 其中 μ_{1k}, μ_{2k} 分别是 x_1, x_2 的均值, σ_1, σ_2 分别是 x_1, x_2 的标准差, 并且 ρ 是 x_1, x_2 的相关系数. 因此双变量高斯分布的密度函数为

$$P(x_1, x_2|Y = k) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left\{ \frac{1}{1-\rho^2} \left[-\frac{(x_1 - \mu_{1k})^2}{2\sigma_1^2} + \frac{\rho}{\sigma_1\sigma_2} (x_1 - \mu_{1k})(x_2 - \mu_{2k}) - \frac{(x_2 - \mu_{2k})^2}{2\sigma_2^2} \right] \right\}$$

问题：不朴素的高斯贝叶斯分类器中 $P(Y|X)$ 的表达式是否仍然形如逻辑回归的形式. 写出 $P(Y|X)$ 的新格式来证明你的答案.

Solution：由于 $X = [x_1, x_2]^T$, 直接使用 **Part A** 中推导的公式

$$P(Y = 1|X) = \frac{1}{1 + \exp \left(\ln \frac{P(X|Y = 0)}{P(X|Y = 1)} + \ln \frac{1 - \pi}{\pi} \right)}$$

将题中的二元正态分布的密度函数表达式代入 $P(Y = 1|X)$, 此时我们不妨先看 $\ln P(X|Y = k)$:

$$\ln P(X|Y = k) = \ln \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} + \frac{1}{1-\rho^2} \left[-\frac{(x_1 - \mu_{1k})^2}{2\sigma_1^2} + \frac{\rho}{\sigma_1\sigma_2} (x_1 - \mu_{1k})(x_2 - \mu_{2k}) - \frac{(x_2 - \mu_{2k})^2}{2\sigma_2^2} \right]$$

于是经过**精心整理**可得:

$$\begin{aligned} \ln \frac{P(X|Y = 0)}{P(X|Y = 1)} &= \ln P(X|Y = 0) - \ln P(X|Y = 1) \\ &= \frac{1}{1-\rho^2} \left\{ \left[\frac{\mu_{10} - \mu_{11}}{\sigma_1^2} + \frac{\rho(\mu_{21} - \mu_{20})}{\sigma_1\sigma_2} \right] x_1 + \left[\frac{\mu_{20} - \mu_{21}}{\sigma_2^2} + \frac{\rho(\mu_{11} - \mu_{10})}{\sigma_1\sigma_2} \right] x_2 \right. \\ &\quad \left. + \left[\frac{\mu_{11}^2 - \mu_{10}^2}{2\sigma_1^2} + \frac{\mu_{21}^2 - \mu_{20}^2}{2\sigma_2^2} + \frac{\rho(\mu_{10}\mu_{20} - \mu_{11}\mu_{21})}{\sigma_1\sigma_2} \right] \right\} \end{aligned}$$

将上述结果代入到 $P(Y = 1|X)$, 可将 $P(Y = 1|X)$ 表示为逻辑回归的形式:

$$P(Y = 1|X) = \frac{1}{1 + \exp \left[- \left(\sum_{i=1}^2 w_i x_i + b \right) \right]}$$

其中

$$\begin{aligned} w_1 &= \frac{-1}{1-\rho^2} \left[\frac{\mu_{10} - \mu_{11}}{\sigma_1^2} + \frac{\rho(\mu_{21} - \mu_{20})}{\sigma_1\sigma_2} \right] \\ w_2 &= \frac{-1}{1-\rho^2} \left[\frac{\mu_{20} - \mu_{21}}{\sigma_2^2} + \frac{\rho(\mu_{11} - \mu_{10})}{\sigma_1\sigma_2} \right] \\ b &= \frac{-1}{1-\rho^2} \left[\frac{\mu_{11}^2 - \mu_{10}^2}{2\sigma_1^2} + \frac{\mu_{21}^2 - \mu_{20}^2}{2\sigma_2^2} + \frac{\rho(\mu_{10}\mu_{20} - \mu_{11}\mu_{21})}{\sigma_1\sigma_2} \right] + \ln \frac{1 - \pi}{\pi} \end{aligned}$$

因此可以知道去掉“naive”的高斯贝叶斯分类器依然是线性的, 可以用 Logistic 回归表示.

通过上述推导，分别去掉了两个不同的假设来揭开高斯朴素贝叶斯分类器的面纱，最后简单总结为：

- 传统的高斯朴素贝叶斯等价于通过假设直接从训练数据中算出参数的逻辑回归；
- 去掉高斯分布中方差与类别无关的假设后，高斯朴素贝叶斯将变成非线性分类器，其中多了二次项特征；
- 移除特征之间的条件独立性假设后，高斯贝叶斯依然是线性的，但是其联合概率密度计算会比较复杂。

Problem 4

Logistic 回归 (LR) 的 MLE 参数估计：如果有参数的正则项 (用以抑制过拟合)，那么该如何估计参数？其中

$$l(\mathbf{w}) = \log L(\mathbf{w}) - \lambda \|\mathbf{w}\|_2^2 = \sum_{i=1}^N y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) + (1 - y^{(i)}) \log (1 - f(\mathbf{x}^{(i)}, \mathbf{w})) - \lambda \|\mathbf{w}\|_2^2$$

Solution: 易知

$$l(\mathbf{w}) = \sum_{i=1}^N y^{(i)} \log f(\mathbf{x}^{(i)}, \mathbf{w}) + (1 - y^{(i)}) \log (1 - f(\mathbf{x}^{(i)}, \mathbf{w})) - \lambda \mathbf{w}^T \mathbf{w}$$

故可以如下求得梯度：

$$\frac{\partial l(\mathbf{w})}{\partial w_j} = \sum_{i=1}^N \left[y^{(i)} \frac{1}{f(\mathbf{x}^{(i)}, \mathbf{w})} \cdot \frac{\partial f(\mathbf{x}^{(i)}, \mathbf{w})}{\partial w_j} + (1 - y^{(i)}) \frac{1}{1 - f(\mathbf{x}^{(i)}, \mathbf{w})} \cdot \frac{-\partial (f(\mathbf{x}^{(i)}, \mathbf{w}))}{\partial w_j} \right] - 2\lambda w_j$$

又因为

$$\frac{\partial f(\mathbf{x}^{(i)}, \mathbf{w})}{\partial w_j} = \underbrace{\frac{\partial g(\mathbf{w}^T \mathbf{x}^{(i)})}{\partial w_j}}_{g(z) = \frac{1}{1 + \exp(-z)} \Rightarrow g'(z) = g(z)(1 - g(z))} = g(\mathbf{w}^T \mathbf{x}^{(i)}) (1 - g(\mathbf{w}^T \mathbf{x}^{(i)})) x_j^{(i)} = f(\mathbf{x}^{(i)}, \mathbf{w}) (1 - f(\mathbf{x}^{(i)}, \mathbf{w})) x_j^{(i)}$$

于是得到梯度表达式：

$$\begin{aligned} \frac{\partial l(\mathbf{w})}{\partial w_j} &= \sum_{i=1}^N \left\{ y^{(i)} \cdot \left[(1 - f(\mathbf{x}^{(i)}, \mathbf{w})) x_j^{(i)} \right] + (1 - y^{(i)}) \cdot (-1) \cdot \left[f(\mathbf{x}^{(i)}, \mathbf{w}) x_j^{(i)} \right] \right\} - 2\lambda w_j \\ &= \sum_{i=1}^N \left[y^{(i)} \cdot (1 - f(\mathbf{x}^{(i)}, \mathbf{w})) \cdot x_j^{(i)} + (y^{(i)} - 1) \cdot f(\mathbf{x}^{(i)}, \mathbf{w}) \cdot x_j^{(i)} \right] - 2\lambda w_j \\ &= \sum_{i=1}^N \left[y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}) \right] x_j^{(i)} - 2\lambda w_j \end{aligned}$$

为了得到 MLE 的参数估计，我们需要使用 SGD 算法来得到 \mathbf{w} ，因此 SGD 的梯度更新规则如下：

$$w_j \leftarrow w_j + \alpha \left\{ \left[y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}) \right] x_j^{(i)} - 2\lambda w_j \right\} = (1 - 2\lambda\alpha) w_j + \alpha \left[y^{(i)} - f(\mathbf{x}^{(i)}, \mathbf{w}) \right] x_j^{(i)}$$

其中 α 为步长； λ 为超参数，需要手动调整来观察过拟合的抑制情况。