

机器学习基础 实验二 实验报告

2021级软件5班 马贵亮 202122202214

实验目的

了解线性判别函数与参数、非参数估计的相关知识，完成如下四个任务。

实验内容

将在后续实验过程中依次阐述对应的内容。

实验过程

1 线性判别函数

1.1 实验内容

采用 `exp2_1.mat` 中的数据，实现线性判别函数分类算法，其 `x1`、`x2` 为二维自变量，`y` 为样本类别。编程实现线性判别函数分类，并做出分类结果可视化。

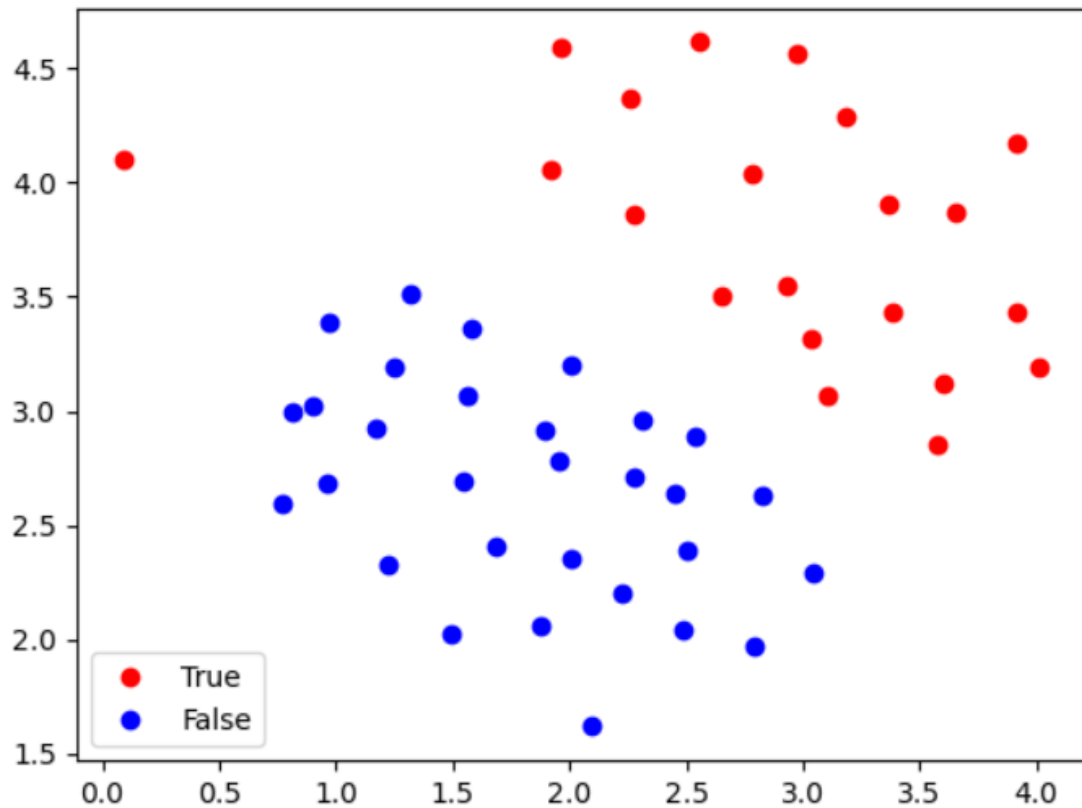
1.2 实验过程

首先我需要对数据进行读取，由于原始数据存在在一个 `.mat` 文件中，因此需要借助 `scipy` 库对 `.mat` 文件进行读入，读入后为了后续处理方便，我在该代码中将其转换为 `dataFrame` 来便于我后续的数据处理。

有关代码如下：

```
1 import scipy
2 import pandas as pd
3 data = scipy.io.loadmat('../data/e2data1.mat')
4 x = pd.DataFrame(data['x'], columns=['x1', 'x2'])
5 y = pd.DataFrame(data['y'], columns=['y'])
```

我们执行完这段代码后我们便有了一组数据，首先根据 `y` 对应的值将所有的 `x` 分成两类，由于所有 `x` 均为二维散点，因此可以根据不同的颜色绘制散点图，来实现对数据的基本可视化。



通过对散点分布的大致观察，该组数据大致线性可分，因此采用简单的线性分类器来对其进行线性判别。针对线性判别，我设计了一个 `LinearClassifier` 类来实现线性判别模型的构建。

在原理方面，我采用单个感知机的模型进行实现，该感知机的激活函数设计为 `sigmoid` 函数。即： $Sigmoid(x) = \frac{1}{1+e^{-x}}$ 。损失函数采用交叉熵即： $Loss(z) = y \ln z + (1 - y) \ln(1 - z)$ 。设该线性模型的权重为 ω ，偏置为 b 。则整体模型为 $z = Sigmoid(X\omega + b)$ 。而最终判别函数以 0.5 为界，大于者为正例，小于者为负例。整体损失为 $Cost = \frac{1}{n} \sum Loss(z)$

我们的目的是使得损失值尽可能小，由梯度下降则

$$d\omega = \frac{dCost}{dz} \cdot Sigmoid'(X\omega + b) \cdot X^T = \frac{1}{n} X^T (z - y)$$

则下降方向为 $-d\omega = -\frac{1}{n} X^T (z - y)$

设学习率为 α 则每次更新为 $\omega = \omega - \alpha d\omega$

再设计一个迭代轮次，则可以实现整个线性分类器的设计，在具体实现的过程中将 X ，转换成 [1 X]，整体代码如下：

`Sigmoid` 函数实现，以及数据转换代码：

```

1  '''
2  DataProcessor.py
3  '''
4  import pandas as pd
5  import numpy as np
6  import matplotlib.pyplot as plt
7
8
9  def PrepareForTraining(df):
10     df.insert(0, '', 1)
11     return df

```

```

12
13 def sigmoid(x):
14     return 1 / (1 + np.exp(-x))
15

```

线性分类设计:

```

1  '''
2  LinearClassifier.py
3  '''
4  import numpy as np
5  from DataProcessor import PrepareForTraining, Sigmoid
6
7
8  class LinearClassifier:
9      def __init__(self, alpha=0.01, iterations=1000):
10         self.alpha = alpha
11         self.iterations = iterations
12         self.weights = None
13
14     def gradientDescent(self, X, y):
15         n = X.shape[0]
16         pred = Sigmoid(np.dot(X, self.weights))
17         delta = pred - y
18         self.weights -= self.alpha * (1 / n) * np.dot(X.T, delta) #交叉熵
19         # self.weights -= self.alpha * (1 / n) * np.dot(X.T, delta * pred *
20         # (1 - pred)) # $L(w)=\frac{1}{2n}[\text{Sigmoid}(Xw)-y]^2$
21
22     def fit(self, X, y):
23         X = PrepareForTraining(X)
24         n, m = X.shape
25         self.weights = np.zeros([m, 1])
26         for _ in range(self.iterations):
27             self.gradientDescent(X, y)
28         return self.weights
29
30     def predict(self, X):
31         X = PrepareForTraining(X)
32         return Sigmoid(np.dot(X, self.weights))>=0.5

```

随后我对我自己的代码进行调用和训练, 可以获得对应的权重

```

1  from LinearClassifier import LinearClassifier
2  from DataProcessor import PrepareForTraining, LinearBoundary
3
4  LC = LinearClassifier(alpha=0.1, iterations=10000)
5  weight = LC.fit(X, y)

```

最终通过简单的数据处理, 在该区间内绘制出线性分类直线, 绘制对应的分类染色, 基于如下代码实现。最终在执行后可以获得如下图。

```

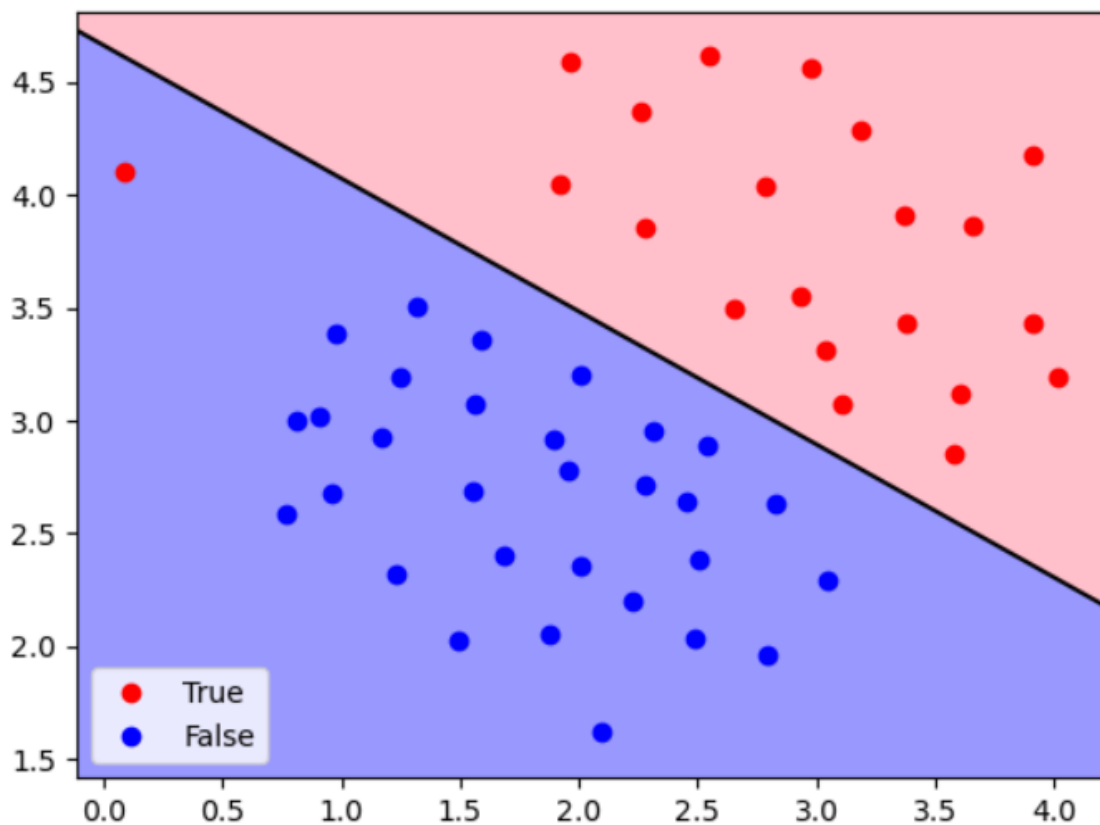
1  import matplotlib.pyplot as plt
2  import numpy as np
3  import pandas as pd
4  from matplotlib.colors import ListedColormap

```

```

5
6 x_min, x_max = x['x1'].min() - 0.2, x['x1'].max() + 0.2
7 y_min, y_max = x['x2'].min() - 0.2, x['x2'].max() + 0.2
8 x = np.linspace(x_min, x_max, 100)
9 liny = LinearBoundary(weight, x)
10 plt.plot(x, liny, 'k')
11
12 custom_cmap = ListedColormap(['#9898ff', # 浅红
13                                '#FFC0CB',])
14
15 x0, x1 = np.meshgrid(np.linspace(x_min, x_max, 1000).reshape(-1, 1),
16                      np.linspace(y_min, y_max, 1000).reshape(-1, 1))
17 X_new = np.c_[x0.ravel(), x1.ravel()]
18 X_new = pd.DataFrame(X_new)
19 Y_new = LC.predict(X_new)
20 z = Y_new.reshape(x0.shape)
21 plt.contourf(x0, x1, z, cmap=custom_cmap)
22
23 plt.axis([x_min, x_max, y_min, y_max])
24 plt.legend(loc='lower left')
25 plt.show()

```



即可实现对实验该部分的实现。

2 最大似然估计

2.1 实验内容

掌握用最大似然估计进行参数估计的原理；当训练样本服从多元正态分布时，计算不同高斯情况下的均值和方差。

使用上面给出的三维数据或者使用 `exp2_2.xlsx` 中的数据：

- (1) 编写程序，对类 1 和类 2 中的三个特征 x_i 分别求解最大似然估计的均值 $\hat{\mu}$ 和方差 σ^2 。
- (2) 编写程序，处理二维数据的情形 $p(x) \sim N(\mu, \Sigma)$ 。对类 1 和类 2 中任意两个特征的组合分别求解最大似然估计的均值 $\hat{\mu}$ 和方差 $\hat{\Sigma}$ （每个类有 3 种可能）。
- (3) 编写程序，处理三维数据的情形 $p(x) \sim N(\mu, \Sigma)$ 。对类 1 和类 2 中三个特征求解最大似然估计的均值 $\hat{\mu}$ 和方差 $\hat{\Sigma}$ 。
- (4) 假设该三维高斯模型是可分离的，即 $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \sigma_3^2)$ ，编写程序估计类 1 和类 2 中的均值和协方差矩阵中的参数。
- (5) 比较前 4 种方法计算出来的每一个特征的均值 μ_i 的异同，并加以解释。
- (6) 比较前 4 种方法计算出来的每一个特征的方差 σ_i 的异同，并加以解释。

2.2 实验过程

首先是基本的数据读入，由于该组数据为 `.xlsx` 文件，因此采用 `pandas` 可以对其进行读入，读入后再做剩余的处理，采用最简单的 `pd.read_excel` 来实现，并且利用 `y` 对应的值对其进行分类以便后边所有题的处理。

```
1 df = pd.read_excel('../data/exp2-2.xlsx')
2 y = df['y']
3 x = df.drop('y', axis=1)
4 x1 = x[y == 1]
5 x2 = x[y == 0]
```

对于 (1)，对于每个类别的每个特征 x_i

$$\mu = \frac{1}{n} x_i$$
$$\sigma^2 = \frac{1}{n} (x_i - \mu)^2$$

同样在python的 `numpy` 包中，有着 `mean` 和 `var` 两个函数可以用来计算 μ 和 σ^2 ，对此我设计了如下代码来进行计算和比较

```
1 def Part1(data):
2     '''
3         编写程序，对类 1 和类 2 中的三个特征  $x_i$  分别求解最大似然估计的均值  $\hat{\mu}$  和方差  $\hat{\sigma}^2$ 。
4         当前要处理的工作即将类中三个特征各自看作独立，求其均值和方差
5         即此时  $\mu$  和  $\Sigma$  均不知
6     '''
7     # 手动实现计算均值和方差
8     data = data.to_numpy()
9     n, m = data.shape
10    miu = np.zeros(m)
11    sigma2 = np.zeros(m)
12    for x in data:
13        miu += x / n
14    for x in data:
```

```

15     sigma2 += (x - miu) ** 2 / n
16     print(f'当对三个特征分别手算求最大似然估计
    时:\nmiu=\n{miu}\nsigma2=\n{sigma2}')
17
18     # 也可以采用pandas中的包来计算均值和方差
19     miu = np.mean(data, axis=0)
20     sigma2 = np.var(data, axis=0, ddof=0)
21     print(f'当对三个特征分别利用numpy求最大似然估计
    时:\nmiu=\n{miu}\nsigma2=\n{sigma2}\n')
22
23 print('第一部分:')
24 print("For Class 1:")
25 Part1(X1)
26 print("For Class 2:")
27 Part1(X2)

```

得到的运行结果如下

```

1  第一部分:
2  For Class 1:
3  当对三个特征分别手算求最大似然估计时:
4  miu=
5  [-0.0709 -0.6047 -0.911 ]
6  sigma2=
7  [0.90617729 4.20071481 4.541949 ]
8  当对三个特征分别利用numpy求最大似然估计时:
9  miu=
10 [-0.0709 -0.6047 -0.911 ]
11 sigma2=
12 [0.90617729 4.20071481 4.541949 ]
13
14 For Class 2:
15 当对三个特征分别手算求最大似然估计时:
16 miu=
17 [-0.1216  0.4299  0.00372]
18 sigma2=
19 [0.05820804 0.04597009 0.00726551]
20 当对三个特征分别利用numpy求最大似然估计时:
21 miu=
22 [-0.1216  0.4299  0.00372]
23 sigma2=
24 [0.05820804 0.04597009 0.00726551]

```

通过比对发现，两者的计算完全一致。

对于 (2) 的部分，我依次枚举两种组合，并以此计算其组合后各自的 μ 和 \sum 。

对于一个两两组合的数据，其 μ 的求法与上述基本一致，只不过换成了矩阵运算。而对于 $\sum = \frac{1}{n} \sum (x_i - \mu) \cdot (x_i - \mu)^T$ 。并且也尝试采用python中的 mean 方法和 cov 方法来进行求解。代码和输出如下

```

1  def DealForMatrix(data):
2      data = data.to_numpy()
3      n, m = data.shape
4      miu = np.zeros(m)
5      sigma2 = np.zeros([m, m])
6      for x in data:

```

```

7         miu += x
8     miu = miu / n
9     for x in data:
10         delta = (x - miu).reshape(-1, 1)
11         sigma2 += np.dot(delta, delta.T)
12     sigma2 = sigma2 / n
13     print(f'当手算求最大似然估计时:\nmiu=\n{miu}\nsigma2=\n{sigma2}')
14
15     # 也可以采用pandas中的包来计算均值和方差
16     miu = np.mean(data, axis=0)
17     sigma2 = np.cov(data, rowvar=False, ddof=0)
18     print(f'当利用numpy求最大似然估计时:\nmiu=\n{miu}\nsigma2=\n{sigma2}\n')
19
20
21 def Part2(data):
22     n, m = data.shape
23     for i in range(m):
24         for j in range(i + 1, m):
25             print(f'采用x{i + 1}和x{j + 1}进行计算所得:')
26             tmp = pd.concat([data.iloc[:, i], data.iloc[:, j]],
axis=1).copy()
27             DealForMatrix(tmp)
28
29 print("第二部分:")
30 print("For Class 1:")
31 Part2(x1)
32 print("For Class 2:")
33 Part2(x2)

```

```

1  第二部分:
2  For Class 1:
3  采用x1和x2进行计算所得:
4  当手算求最大似然估计时:
5  miu=
6  [-0.0709 -0.6047]
7  sigma2=
8  [[0.90617729 0.56778177]
9   [0.56778177 4.20071481]]
10 当利用numpy求最大似然估计时:
11  miu=
12  [-0.0709 -0.6047]
13  sigma2=
14  [[0.90617729 0.56778177]
15   [0.56778177 4.20071481]]
16
17 采用x1和x3进行计算所得:
18 当手算求最大似然估计时:
19  miu=
20  [-0.0709 -0.911 ]
21  sigma2=
22  [[0.90617729 0.3940801 ]
23   [0.3940801 4.541949 ]]
24 当利用numpy求最大似然估计时:
25  miu=
26  [-0.0709 -0.911 ]
27  sigma2=
28  [[0.90617729 0.3940801 ]

```

```
29     [0.3940801  4.541949  ]]
30
31 采用x2和x3进行计算所得：
32 当手算求最大似然估计时：
33  miu=
34  [-0.6047 -0.911  ]
35  sigma2=
36  [[4.20071481 0.7337023  ]
37   [0.7337023  4.541949  ]]
38 当利用numpy求最大似然估计时：
39  miu=
40  [-0.6047 -0.911  ]
41  sigma2=
42  [[4.20071481 0.7337023  ]
43   [0.7337023  4.541949  ]]
44
45  For Class 2:
46  采用x1和x2进行计算所得：
47  当手算求最大似然估计时：
48  miu=
49  [-0.1216  0.4299]
50  sigma2=
51  [[ 0.05820804 -0.01321216]
52   [-0.01321216  0.04597009]]
53 当利用numpy求最大似然估计时：
54  miu=
55  [-0.1216  0.4299]
56  sigma2=
57  [[ 0.05820804 -0.01321216]
58   [-0.01321216  0.04597009]]
59
60 采用x1和x3进行计算所得：
61 当手算求最大似然估计时：
62  miu=
63  [-0.1216  0.00372]
64  sigma2=
65  [[ 0.05820804 -0.00478645]
66   [-0.00478645  0.00726551]]
67 当利用numpy求最大似然估计时：
68  miu=
69  [-0.1216  0.00372]
70  sigma2=
71  [[ 0.05820804 -0.00478645]
72   [-0.00478645  0.00726551]]
73
74 采用x2和x3进行计算所得：
75 当手算求最大似然估计时：
76  miu=
77  [0.4299  0.00372]
78  sigma2=
79  [[0.04597009 0.00850987]
80   [0.00850987 0.00726551]]
81 当利用numpy求最大似然估计时：
82  miu=
83  [0.4299  0.00372]
84  sigma2=
85  [[0.04597009 0.00850987]
86   [0.00850987 0.00726551]]
```


可以发现其计算结果保持一致，即代码设计正确。

对于 (3) 部分，其关键步骤于 (2) 保持一致，直接调用 `DealForMatrix` 即可，其具体代码和结果如下：

```
1 def Part3(data):
2     DealForMatrix(data)
3
4     print("第三部分:")
5     print("For Class 1:")
6     Part3(X1)
7     print("For Class 2:")
8     Part3(X2)
```

```
1 第三部分:
2 For Class 1:
3 当手算求最大似然估计时:
4 miu=
5 [-0.0709 -0.6047 -0.911 ]
6 sigma2=
7 [[0.90617729 0.56778177 0.3940801 ]
8  [0.56778177 4.20071481 0.7337023 ]
9  [0.3940801  0.7337023  4.541949  ]]
10 当利用numpy求最大似然估计时:
11 miu=
12 [-0.0709 -0.6047 -0.911 ]
13 sigma2=
14 [[0.90617729 0.56778177 0.3940801 ]
15  [0.56778177 4.20071481 0.7337023 ]
16  [0.3940801  0.7337023  4.541949  ]]
17
18 For Class 2:
19 当手算求最大似然估计时:
20 miu=
21 [-0.1216  0.4299  0.00372]
22 sigma2=
23 [[ 0.05820804 -0.01321216 -0.00478645]
24  [-0.01321216  0.04597009  0.00850987]
25  [-0.00478645  0.00850987  0.00726551]]
26 当利用numpy求最大似然估计时:
27 miu=
28 [-0.1216  0.4299  0.00372]
29 sigma2=
30 [[ 0.05820804 -0.01321216 -0.00478645]
31  [-0.01321216  0.04597009  0.00850987]
32  [-0.00478645  0.00850987  0.00726551]]
```

对于第四部分，当 $\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \sigma_3^2)$ 时，其对应的 σ 和第一部分计算保持一致，因此不再过多赘述。

```
1 def Part4(data):
2     '''
3     绝大多数代码在Part1中已经呈现，因此此处直接使用numpy实现
4     '''
5     data = data.to_numpy()
```

```

6      miu = np.mean(data,axis=0)
7      sigma2 = np.var(data,axis=0,ddof=0)
8      cov = np.diag(sigma2)
9      print(f'当利用numpy求最大似然估计时:\nmiu=\n{miu}\ncov=\n{cov}\n')
10
11     print("第四部分:")
12     print("For Class 1:")
13     Part4(x1)
14
15     print("For Class 2:")
16     Part4(x2)

```

```

1  第四部分:
2  For Class 1:
3  当利用numpy求最大似然估计时:
4  miu=
5  [-0.0709 -0.6047 -0.911 ]
6  cov=
7  [[0.90617729 0.          0.          ]
8   [0.          4.20071481 0.          ]
9   [0.          0.          4.541949  ]]
10
11  For Class 2:
12  当利用numpy求最大似然估计时:
13  miu=
14  [-0.1216  0.4299  0.00372]
15  cov=
16  [[0.05820804 0.          0.          ]
17   [0.          0.04597009 0.          ]
18   [0.          0.          0.00726551]]

```

对于 (5) 和(6)在求解均值和协方差的过程中，计算最大似然估计，协方差计算中除以 n 是合适的。即采用有偏估计。整体的值均保持一致。

- 对比均值：各种方法计算的均值应该一致，因为均值的估计与协方差矩阵的形式无关。
- 对比方差：直接方差估计与在协方差矩阵中提取的方差应一致。然而，对角协方差矩阵只考虑单个变量的方差，忽略变量之间的相关性，可能在某些应用中提供不同的视角。

3-Parzen窗

3-Parzen.1 实验内容

使用上面表格中的数据或者使用 exp2_3.xlsx 中的数据进行 Parzen 窗估计和设计分类器。窗函数为一个球形的高斯函数如公式2-1所示：

$$\varphi\left(\frac{x-x_i}{h}\right) \propto \exp\left[-\frac{(x-x_i)^T(x-x_i)}{2h^2}\right]$$

编写程序，使用 Parzen 窗估计方法对任意一个的测试样本点 xx 进行分类。对分类器的训练则使用表2-2中的三维数据。令 $h = 1$ ，分类样本点为 $(0.5, 1.0, 0.0)^T$ ， $(0.31, 1.51, -0.50)^T$ ， $(-0.3, 0.44, -0.1)^T$ 。

3-Parzen.2 实验过程

Parzen窗的基本原理为以当前点为中心，绘制一个h的高斯窗口，对其中进行数点，即可返回类条件概率密度，在理想情况下也可根据点数的多少直接进行分类。那么基于如上基础进行了简单代码实现如下：

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5
6 def ParzenWindow(X, h, test):
7     X = X.to_numpy()
8     sigma = 0
9     for x in X:
10         sigma += np.exp(-np.dot((test - x), (test - x).T) / 2 * h ** 2)
11     return sigma
12
13
14 def Parzen(X, y, h, test):
15     test = np.array(test)
16     yList = y.drop_duplicates().tolist()
17     res = np.zeros([len(yList), 1])
18     for index, testX in enumerate(test):
19         K = np.zeros(len(yList))
20         for i, yEle in enumerate(yList):
21             XList = X[y == yEle]
22             k = ParzenWindow(XList, h, testX)
23             K[i] = k / (XList.shape[0] * h ** 2)
24         res[index] = yList[np.argmax(K)]
25     return res
26
27
28 df = pd.read_excel('../data/exp2-3.xlsx')
29 df.to_csv('../data/exp2-3.csv', index=False)
30 y = df['y']
31 X = df.drop('y', axis=1)
32 Xtest = [[0.5, 1.0, 0.0],
33          [0.31, 1.51, -0.50],
34          [-0.3, 0.44, -0.1]]
35 yPred = Parzen(X, y, 1, Xtest)
36 print(yPred)
```

在进行计算后，可以获得如下结果：

```
1 [[2.]
2  [2.]
3  [2.]]
```

则将该过程实现。

3-KNN

3-KNN.1 实验内容

k-近邻概率密度估计：

对上面表格中的数据使用k-近邻方法进行概率密度估计：

1. 编写程序，对于一维的情况，当有 n 个数据样本点时，进行 k -近邻概率密度估计。对表格中的类3的特征 x_1 ，用程序画出当 $k = 1, 3, 5$ 时的概率密度估计结果。
2. 编写程序，对于二维的情况，当有 n 个数据样本点时，进行k-近邻概率密度估计。对表格中的类2的特征 $(x_1, x_2)^T$ ，用程序画出当 $k = 1, 3, 5$ 时的概率密度估计结果。
3. 编写程序，对表格中的3个类别的三维特征，使用 k -近邻概率密度估计方法。并且对下列点处的概率密度进行估计：
 $(-0.41, 0.82, 0.88)^T$ ， $(0.14, 0.72, 4.1)^T$ ， $(-0.81, 0.61, -0.38)^T$ 。

3-KNN.2 实验过程

对于整个实验我大致有两种不同的理解方案，其实际的预测结果差值不大，但其实际运算的类条件概率密度具有一定的差异性，而且在后段的预测过程中使用了贝叶斯的思想，在预测不仅仅只考虑类条件概率密度，也考虑进先验概率。

我的第一种理解时仅考虑题面所提及的类3（题1），和类2（题2）。仅使用这一个类，利用knn求得这一类得类条件概率密度（与其他类无关）。即k个点只针对一个样本，考虑到如果这个点附近的不够密，那么其距离会很远，则类条件概率密度很小。

使用多个类，但是只展示类3（题1），类2（题2）的概率密度。即k个点是所有样本。那么在这k个点中最多的那类的类条件概率密度应当最高。

再利用求得的类条件概率密度与先验概率相乘再比较大小进行预测。

在计算其中的距离后对应的体积时，我采用了高维球体的体积计算公式进行

超球体体积 $V_n(r) = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2}+1)} r^n$ ，令 $\alpha = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2}+1)}$ 为系数，则 $V_n(r) = \alpha r^n$

并且基于我个人对于KNN的理解，构建了一个KNN类来实现该实验以及4 KNN实战的部分的KNN代码，如下：

```
1  '''
2  该代码为在实验4内容修改后的knn代码
3  '''
4
5  import numpy as np
6  import matplotlib.pyplot as plt
7  import pandas as pd
8  from numpy.linalg import inv
9  import math
10
11
12  class kNNClassifier:
13
14      def KNN_Vn(self, x, k, t):
15          '''
16          基于  $p_{\{n\}}(x) = k_{\{n\}}/(V_{\{n\}}*n)$  来计算概率密度
17          :param x: 原始数据集的某一类别
18          :param k: KNN 的k
19          :param t: 测试集
20          :return: 一个概率密度函数
```

```

21     ...
22     res = np.zeros((t.shape[0],))
23     for i, test in enumerate(t):
24         dist = np.linalg.norm(X - test, axis=1)
25         dist = np.sort(dist)
26         if dist[k - 1] != 0:
27             res[i] = k / (X.shape[0] * dist[k - 1] ** X.shape[1])
28         else:
29             res[i] = k / (X.shape[0] * (1e-5) ** X.shape[1])
30     return res
31
32 def knn_Vn(self, X_train, y_train, X_test, k):
33     res = [{_ for _ in range(X_test.shape[0])}]
34     yvalueList = y_train.unique().tolist()
35     for y in yvalueList:
36         trainSet = X_train[y_train == y].copy()
37         yRes = self.knn_Vn(trainSet, k, X_test)
38         for i in range(len(yRes)):
39             # 为 res 中的每一行的字典添加键值对，键为 y，值为 yRes 中对应的元素
40             res[i][y] = yRes[i]
41     return res
42
43 def gamma(self, x):
44     # 递归利用已知的Gamma(1/2) = sqrt(pi)
45     if abs(x - 0.5) < 1e-6:
46         return math.sqrt(math.pi)
47     elif abs(x - 1) < 1e-6:
48         return 1
49     else:
50         return (x - 1) * self.gamma(x - 1)
51
52 def knn_Euler_Count(self, X_train, y_train, X_test, k):
53     ...
54     基于欧拉距离的计数方法来计算类条件概率密度
55     :param X_train: 训练集
56     :param y_train: 训练标签
57     :param X_test: 测试集
58     :param k: k近邻的k
59     :return: 测试集类条件概率密度
60     ...
61     res = []
62     # 超球体体积系数  $V_n(r) = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
63     令  $\alpha = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
64     pi_power = math.pi ** (X_train.shape[1] / 2)
65     gamma_value = self.gamma((X_train.shape[1] / 2) + 1)
66     alpha = ((pi_power) / gamma_value) if gamma_value != 0 else
67     float('inf')
68     for i, test in enumerate(X_test):
69         dist = np.linalg.norm(X_train - test, axis=1)
70         knn_indices = np.argsort(dist)[:k]
71         knn_labels = y_train[knn_indices]
72         d = dist[knn_indices[-1]]
73         if d == 0:
74             d = 1e-5
75         class_probs = {
76             cls: np.sum(knn_labels == cls) / (X_train.shape[0] * alpha
77             * (d ** X_train.shape[1]))
78             for cls in np.unique(y_train)}

```

```

76         res.append(class_probs)
77     return res
78
79     def mahalanobis_distance(self, x, dataset, invCov):
80         '''
81         计算单个样本与数据集所有样本的马氏距离
82         :param x:
83         :param dataset:
84         :param invCov:
85         :return:
86         '''
87         dist = np.zeros(dataset.shape[0])
88         dataset = dataset.to_numpy()
89         for i, data in enumerate(dataset):
90             delta = x - data
91             dist[i] = np.sqrt(delta.dot(invCov).dot(delta.T))
92         return dist
93
94     def knn_mahalanobis_Count(self, X_train, y_train, X_test, k):
95         '''
96         基于马氏距离的计数方法来计算类条件概率密度
97         :param X_train: 训练集
98         :param y_train: 训练标签
99         :param X_test: 测试集
100        :param k: k近邻的k
101        :return: 测试集的条件概率密度
102        '''
103        # 超球体体积系数  $V_n(r) = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
104        r^n$ 令  $\alpha = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
105        pi_power = math.pi ** (X_train.shape[1] / 2)
106        gamma_value = self.gamma((X_train.shape[1] / 2) + 1)
107        alpha = ((pi_power) / gamma_value) if gamma_value != 0 else
108        float('inf')
109        cov = np.cov(X_train.T)
110        invCov = inv(cov)
111        res = []
112        for i, test in enumerate(X_test):
113            dist = self.mahalanobis_distance(test, X_train, invCov)
114            KNN_indices = np.argsort(dist)[:k]
115            KNN_labels = y_train[KNN_indices]
116            d = dist[KNN_indices[-1]]
117            if d == 0:
118                d = 1e-5
119            class_probs = {
120                cls: np.sum(KNN_labels == cls) / (X_train.shape[0] * alpha
121                * (d ** X_train.shape[1]))
122                for cls in np.unique(y_train)}
123            res.append(class_probs)
124        return res
125
126    def density(self, X_train, y_train, X_test, k, typ=0):
127        yvalueList = y_train.unique().tolist()
128        prior_prob = {}
129        if not isinstance(X_test, np.ndarray):
130            # 如果不是, 转换它为numpy数组
131            X_test = np.array(X_test)
132        for y in yvalueList:
133            prior_prob[y] = X_train[y_train == y].shape[0]

```

```

131         if typ == 0:
132             probs = self.kNN_Euler_Count(X_train, y_train, X_test, k)
133         elif typ == 1:
134             probs = self.kNN_Mahalanobis_Count(X_train, y_train, X_test, k)
135         elif typ == 2:
136             probs = self.kNN_Vn(X_train, y_train, X_test, k)
137         return probs

```

如上述实验一致，我们需要先对数据进行读入，依旧采用 `pandas` 包中的 `read_excel` 方法将其读入为一个 `dataFrame`，并将 `x` 和 `y` 提取出来，并且导入(3)中需要预测的数据

```

1 df = pd.read_excel('../data/exp2-3.xlsx')
2 df.to_csv('../data/exp2-3.csv', index=False)
3 y = df['y']
4 x = df.drop('y', axis=1)
5 xtest = [[-0.41, 0.82, 0.88],
6           [0.14, 0.72, 4.1],
7           [-0.81, 0.61, -0.38]]
8 xtest = np.array(xtest)

```

以下为两种不同的理解对应的各自结果和图像。

首先对于(1):

对于第一种理解，需要将数据根据类完全摘离，数据处理即为：

```

1 x1 = x['x1'].copy()
2 x_train = x1[y == 3].copy().reset_index(drop=True)
3 y_train = y[y == 3].copy().reset_index(drop=True)
4 DealPart1(x_train, y_train)

```

对于第二种理解，不需要单独摘出数据，数据预处理为：

```

1 x1 = x['x1'].copy()
2 x_train = x1.copy().reset_index(drop=True)
3 y_train = y.copy().reset_index(drop=True)
4 DealPart1(x_train, y_train)

```

而两者调用KNN的方法保持一致，即如下函数 `DealPart1`：

```

1 def DealPart1(X_train, y_train):
2     x_min = X_train.min()
3     x_max = X_train.max()
4     X_train = X_train.to_numpy().reshape(-1, 1)
5     X_test = np.linspace(x_min, x_max, 1000, endpoint=True).reshape(-1, 1)
6     knn = kNNClassifier()
7
8     # 获得概率密度函数
9     y1_density = knn.density(X_train, y_train, X_test, 1)
10    y3_density = knn.density(X_train, y_train, X_test, 3)
11    y5_density = knn.density(X_train, y_train, X_test, 5)
12
13    y1_test = np.array([d[3] for d in y1_density]).reshape(-1, 1)
14    y3_test = np.array([d[3] for d in y3_density]).reshape(-1, 1)
15    y5_test = np.array([d[3] for d in y5_density]).reshape(-1, 1)

```

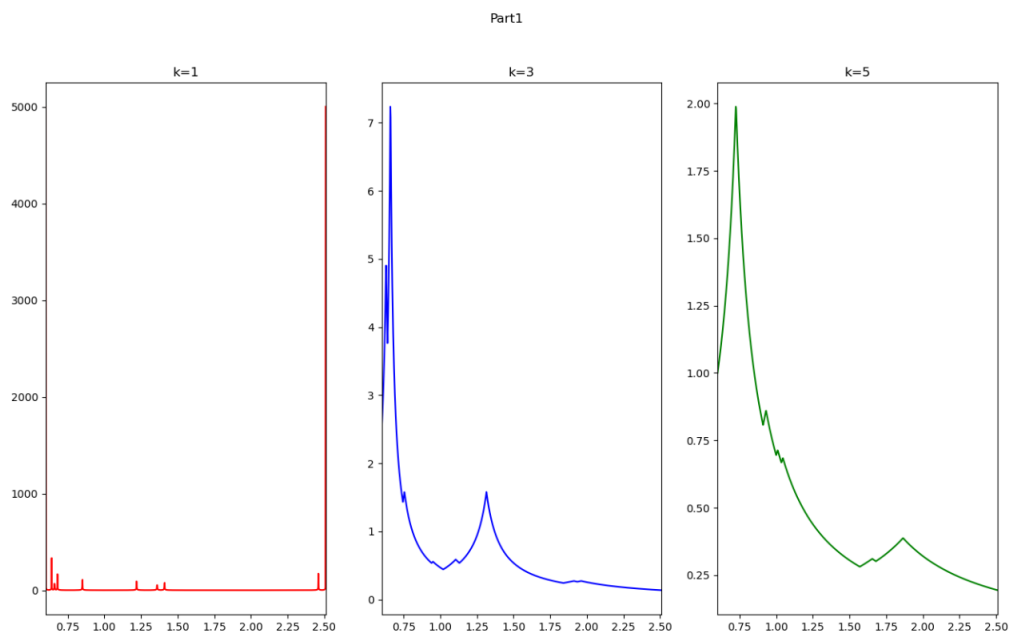
```

16
17 plt.figure(figsize=(16, 9))
18 plt.subplot(131)
19 plt.plot(X_test, y1_test, 'r')
20 plt.xlim([x_min, x_max])
21 plt.title('k=1')
22
23 plt.subplot(132)
24 plt.plot(X_test, y3_test, 'b')
25 plt.xlim([x_min, x_max])
26 plt.title('k=3')
27
28 plt.subplot(133)
29 plt.plot(X_test, y5_test, 'g')
30 plt.xlim([x_min, x_max])
31 plt.title('k=5')
32
33 plt.suptitle('Part1')
34 plt.show()

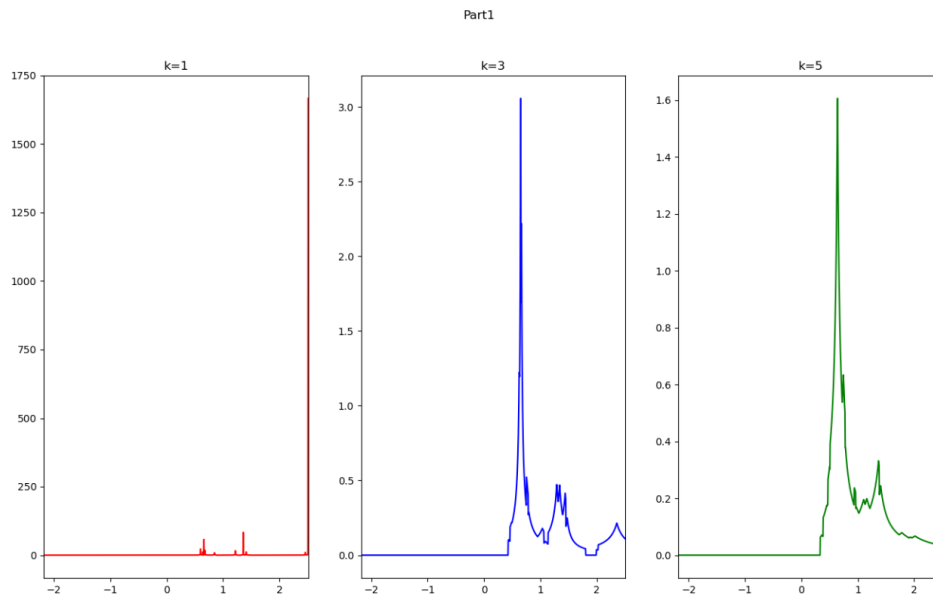
```

分别执行后：

对于第一种理解不同k下的类条件概率密度如下：



对于第二种理解不同k下的类条件概率密度如下：



如果将其对应的x轴对其，可以发现理解1的图像更加光滑，因为其在数点的过程中新增一个点增加距离均保持一定的一致性，因此整体图像比较光滑。而对于第二种理解的情况下由于考虑到k个点的范围内不一定存在该类的点，因此会出现数据为0的情况。

而对于k=1的时候，由于会导致距离变为0，因此需要进行设置一个极小的数值来保证其数据输出的合理性。

同样相仿上述对于(1)的操作，对于(2)的操作我们仍需先进行数据处理是否摘除，再执行KNN的操作。

对于理解1的数据处理如下：

```
1 x2 = pd.concat([X['x1'], X['x2']], axis=1)
2 x_train = x2[y == 2].copy().reset_index(drop=True)
3 y_train = y[y == 2].copy().reset_index(drop=True)
4 DealPart2(x_train, y_train)
```

对于理解2的数据处理如下：

```
1 x2 = pd.concat([X['x1'], X['x2']], axis=1)
2 x_train = x2.copy().reset_index(drop=True)
3 y_train = y.copy().reset_index(drop=True)
4 DealPart2(x_train, y_train)
```

同样两者计算类条件概率密度的函数相同，为 DealPart2 如下

```
1 def DrawPart2(x0, x1, z, k):
2     # matplotlib.use('TkAgg')
3     # 创建一个图形和两个子图（一个2D，一个3D）
4     fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 6))
5
6     # 第一个子图为3D图
7     ax1 = fig.add_subplot(121, projection='3d')
8     surf = ax1.plot_surface(x0, x1, z, cmap='viridis')
9     ax1.set_xlabel('x1')
10    ax1.set_ylabel('x2')
11    ax1.set_zlabel('p')
12    ax1.set_title(f'k={k}, 3D')
```

```

13     ax1.patch.set_visible(False)
14     ax1.grid(False)
15
16     # 第二个子图为2D等高线图
17     contour = ax2.contourf(x0, x1, z, cmap='viridis')
18     ax2.set_xlabel('x1')
19     ax2.set_ylabel('x2')
20     ax2.set_title(f'k={k}, 2D')
21     # 调整布局
22     plt.tight_layout()
23     # 显示图形
24     plt.show()
25
26 def DealPart2(X_train, y_train):
27     x1_min = X_train['x1'].min()
28     x1_max = X_train['x1'].max()
29     x2_min = X_train['x2'].min()
30     x2_max = X_train['x2'].max()
31     x0, x1 = np.meshgrid(np.linspace(x1_min, x1_max, 100).reshape(-1, 1),
32                           np.linspace(x2_min, x2_max, 100).reshape(-1, 1))
33     X_test = np.c_[x0.ravel(), x1.ravel()]
34     X_train = X_train.to_numpy()
35     knn = KNNClassifier()
36     y1_post = knn.density(X_train, y_train, X_test, 1)
37     y3_post = knn.density(X_train, y_train, X_test, 3)
38     y5_post = knn.density(X_train, y_train, X_test, 5)
39
40     # 获得概率密度函数
41     y1_density = knn.density(X_train, y_train, X_test, 1)
42
43     y3_density = knn.density(X_train, y_train, X_test, 3)
44     y5_density = knn.density(X_train, y_train, X_test, 5)
45
46     y1_test = np.array([d[2] for d in y1_density]).reshape(-1, 1)
47     y3_test = np.array([d[2] for d in y3_density]).reshape(-1, 1)
48     y5_test = np.array([d[2] for d in y5_density]).reshape(-1, 1)
49
50     DrawPart2(x0, x1, y1_test.reshape(x0.shape), 1)
51     DrawPart2(x0, x1, y3_test.reshape(x0.shape), 3)
52     DrawPart2(x0, x1, y5_test.reshape(x0.shape), 5)

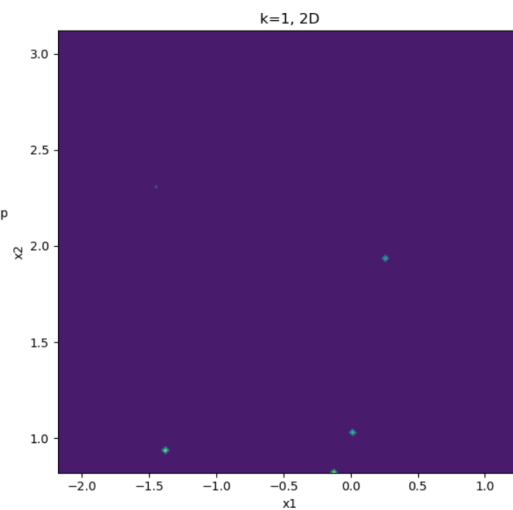
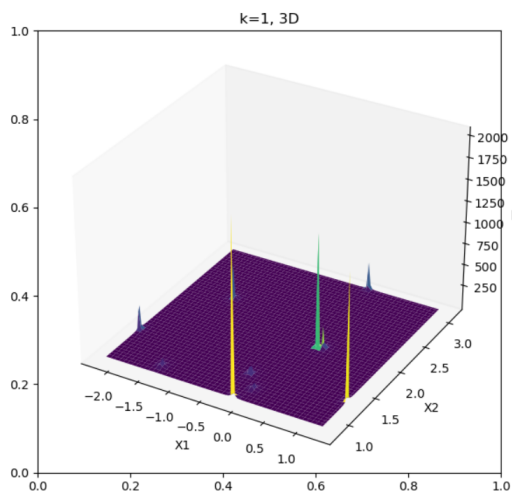
```

其运行的结果如下：

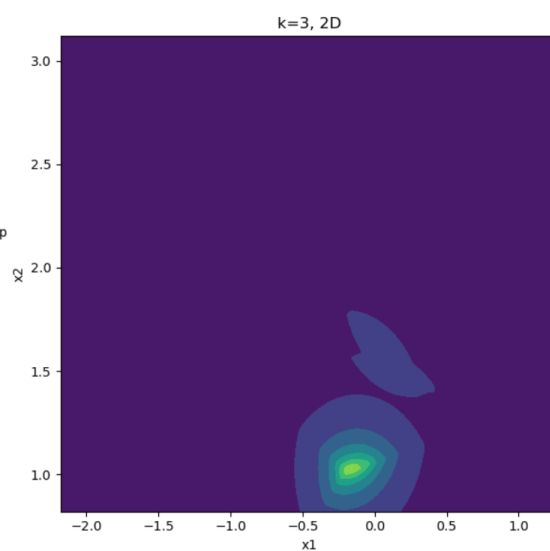
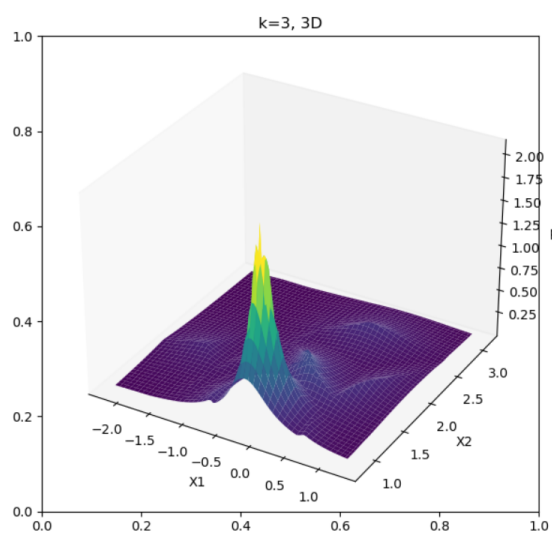
左侧为二维曲面图，右侧为对应的等高线

第一种理解：

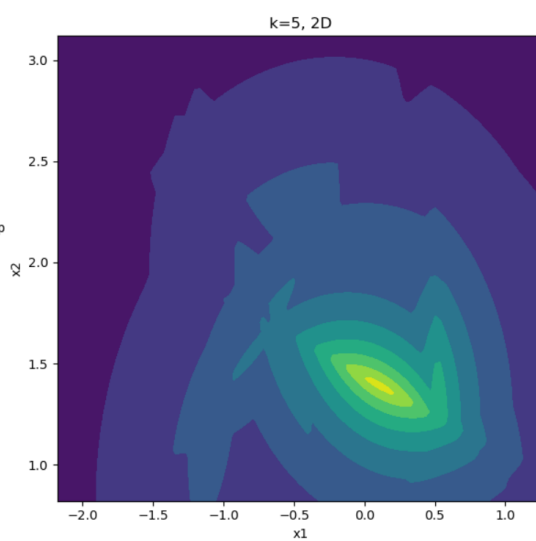
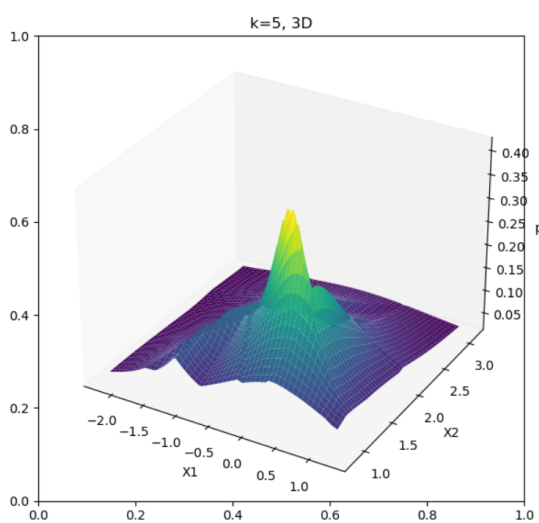
k=1



$k=3$

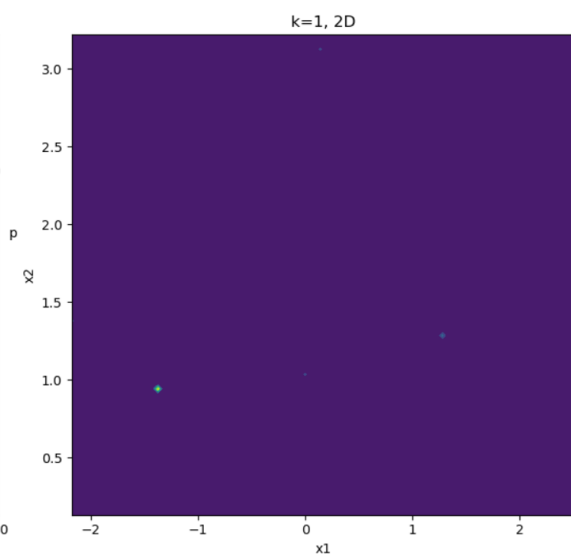
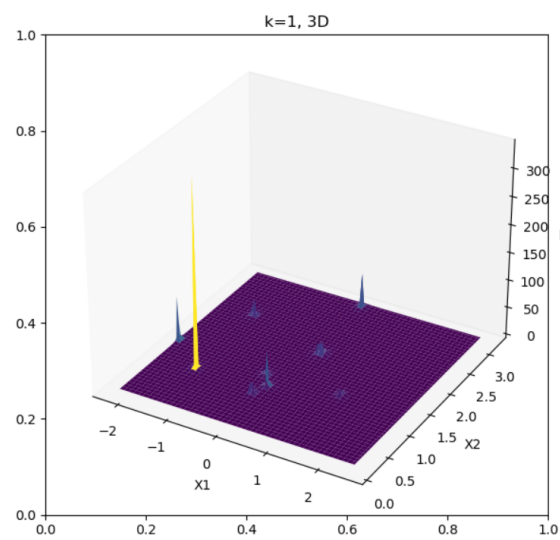


$k=5$

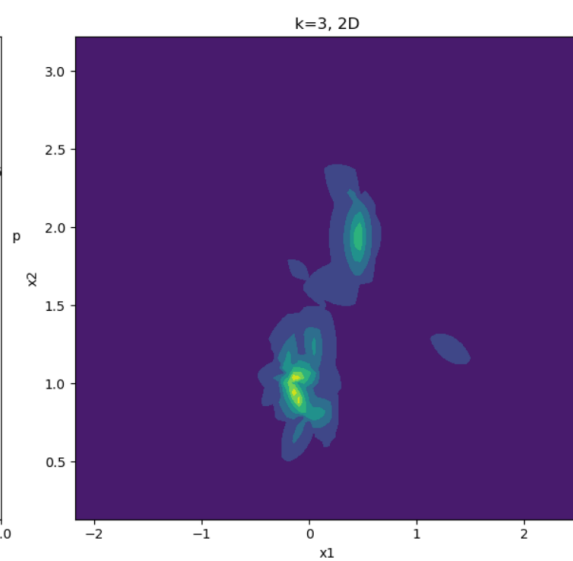
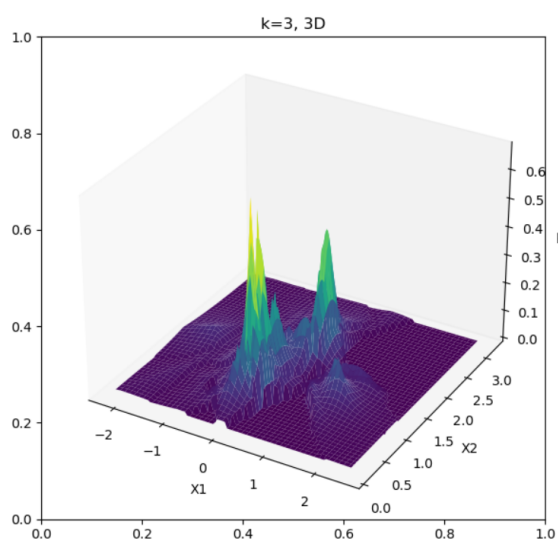


第二种理解:

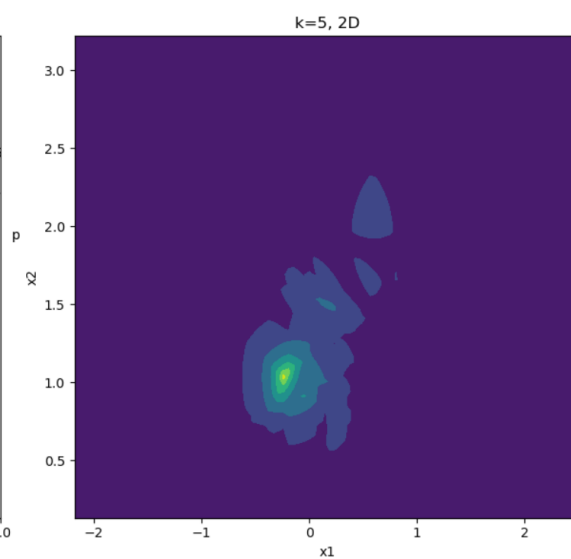
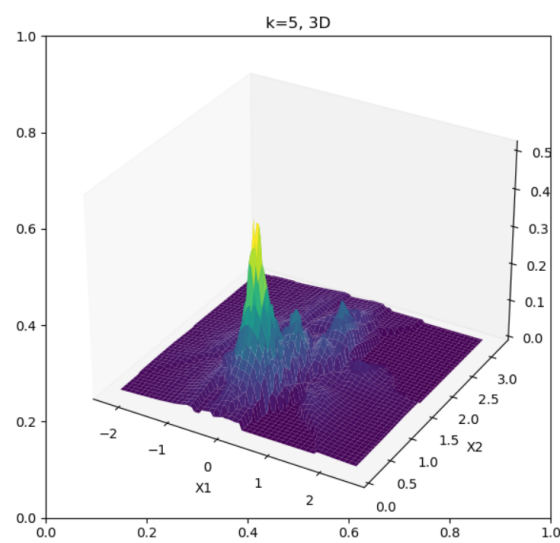
$k=1$



k=3



k=5



而对于(3) 采用不同的k进行执行，获得其对应的类条件概率密度即可，执行如下代码，得到如下结果

```

1 knn = KNNClassifier()
2 for k in range(5):
3     density = knn.density(X, y, Xtest, k + 1, 2)
4     print(f'k={k + 1} 时的概率密度')
5     for index in range(len(density)):
6         print(density[index])

```

```

1 k=1 时的概率密度
2 {1: 0.01883296229251353, 2: 0.22167283203251173, 3: 0.07712096076645929}
3 {1: 0.04131910919281311, 2: 0.0016351496179347956, 3: 0.0030716239012411896}
4 {1: 0.1771613350916689, 2: 0.13534049221543457, 3: 0.013723919694781556}
5 k=2 时的概率密度
6 {1: 0.03178787537795644, 2: 0.16427006941590264, 3: 0.14466794867681504}
7 {1: 0.026792930905745034, 2: 0.003220749147742275, 3: 0.005967511370200231}
8 {1: 0.22110756108533755, 2: 0.26802033748210696, 3: 0.026840350461291625}
9 k=3 时的概率密度
10 {1: 0.00879596401261926, 2: 0.23162580018270137, 3: 0.15009942277353658}
11 {1: 0.017964386176866163, 2: 0.0040412635615593925, 3: 0.008591866296713964}
12 {1: 0.01102829558991302, 2: 0.3638189115792198, 3: 0.0354162032665907}
13 k=4 时的概率密度
14 {1: 0.010231864737074553, 2: 0.23905049106054166, 3: 0.1874814465698455}
15 {1: 0.013391533141938477, 2: 0.005348468289115052, 3: 0.011381430132195194}
16 {1: 0.00919145633055199, 2: 0.3362516370263627, 3: 0.040106087762924514}
17 k=5 时的概率密度
18 {1: 0.010953452694776046, 2: 0.1393888918855546, 3: 0.11134270190658971}
19 {1: 0.004295766694553839, 2: 0.006489497138077377, 3: 0.010832627307478103}
20 {1: 0.007823205569486222, 2: 0.12488416706442222, 3: 0.039732878044842136}

```

4 KNN实战

4-1 实验目的

掌握KNN算法的使用。

一、数据预处理

- 1.将e2.txt中的数据处理成可以输入给模型的格式
- 2.是否还需要对特征值进行归一化处理？目的是什么？

二、数据可视化分析

将预处理好的数据以散点图的形式进行可视化，通过直观感觉总结规律，感受 KNN 模型思想与人类经验的相似之处。

三、构建 KNN 模型并测试

- 1.输出测试集各样本的预测标签和真实标签，并计算模型准确率。
- 2.选择哪种距离更好？欧氏还是马氏？
- 3.改变数据集的划分以及 k 的值，观察模型准确率随之的变化情况。

注意：选择训练集与测试集的随机性

四、使用模型构建可用系统

利用构建好的 KNN 模型实现系统，输入为新的数据的三个特征，输出为预测的类别。

4-2 实验过程

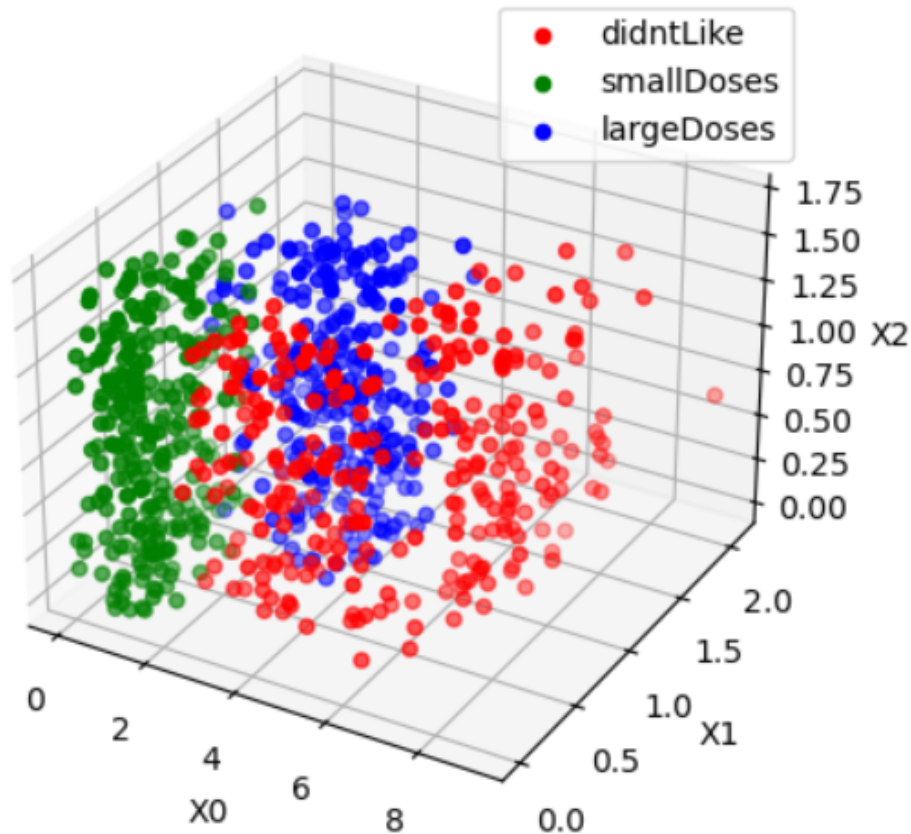
首先获得该数据，利用 pandas 的 `read_csv` 函数进行读入，通过对数据的观察，对数据进行一个暴力的归一化处理以及标签的数据化。如下：

```
1 yMapping = {'didn'tLike': 0, 'smallDoses': 1, 'largeDoses': 2}
2 yList = ['didn'tLike', 'smallDoses', 'largeDoses']
3
4
5 def dealForData(data):
6     data['x0'] = data['x0'] / 10000
7     data['x1'] = data['x1'] / 10
8     data['y'] = data['y'].replace(yMapping)
9     return data
10
11 data = pd.read_csv('../data/e2.txt', sep="\t", names=['x0', 'x1', 'x2',
12     'y'])
12 data = dealForData(data)
```

随后简单对数据集进行一个4: 1的划分，并对训练集绘制散点图来可视化

```
1 def DrawScatterPlot(x, y):
2     fig = plt.figure()
3     ax = fig.add_subplot(111, projection='3d')
4
5     # 为每个分类的点设置不同的颜色和标签
6     colors = ['red', 'green', 'blue']
7     labels = ['didn'tLike', 'smallDoses', 'largeDoses']
8     for i in range(3):
9         subset = x[y == i]
10        ax.scatter(subset['x0'], subset['x1'], subset['x2'], c=colors[i],
11            label=labels[i])
12        # 设置图例和坐标轴标签
13        ax.legend()
14        ax.set_xlabel('x0')
15        ax.set_ylabel('x1')
16        ax.set_zlabel('x2')
17        plt.show()
18
19 def splitForData(data, test_size, random_state):
20     np.random.seed(random_state)
21     data_shuffled = data.sample(frac=1).reset_index(drop=True)
22     train_size = int((1 - test_size) * len(data))
23     train_data = data_shuffled[:train_size]
24     test_data = data_shuffled[train_size:]
25     x_train = train_data.drop('y', axis=1)
26     y_train = train_data['y']
27     x_test = test_data.drop('y', axis=1)
28     y_test = test_data['y'].copy()
29     return x_train, y_train, x_test, y_test
30
31 x_train, y_train, x_test, y_test = splitForData(data, 0.2, 7)
32 y_test = y_test.to_numpy()
33 DrawScatterPlot(x_train, y_train)
```

即可获得如右图的散点图：



随后即可执行KNN算法，在上述给定的代码中，我采用可调节的参数来进行选择不同的距离计算方式以及不同的计算概率密度的方式，并在预测中再乘统计的先验概率来进行预测，比单纯采用类条件概率密度的预测在数据不稳定打乱的情况下更好。KNN类的具体描述如下。

```

1  class knnClassifier:
2
3      # 该方法基于上述第三问题的KNN中的第一种理解，即将类别分为不同的类进行，不同类之间
      # 互不干涉，但是该方法存在某一类的点少于k个，则会出现概念问题，并且不易解决，因此采用该方法
      # 时k不应该选取太大。
4      def KNN_Vn(self, x, k, t):
5          '''
6              基于  $p_{\{n\}}(x) = k_{\{n\}} / (V_{\{n\}} * n)$  来计算概率密度
7              :param x: 原始数据集的某一类别
8              :param k: KNN 的k
9              :param t: 测试集
10             :return: 一个概率密度函数
11             '''
12             res = np.zeros((t.shape[0],))
13             for i, test in enumerate(t):
14                 dist = np.linalg.norm(x - test, axis=1)
15                 dist = np.sort(dist)
16                 if dist[k - 1] != 0:
17                     res[i] = k / (x.shape[0] * dist[k - 1] ** x.shape[1])
18                 else:
19                     res[i] = k / (x.shape[0] * (1e-5) ** x.shape[1])
20             return res
21

```

```

22 # 该方法与上述的方法为配套使用的方法，来实现上述描述的理解。
23 def knn_vn(self, X_train, y_train, X_test, k):
24     res = [{ } for _ in range(X_test.shape[0])]
25     yvalueList = y_train.unique().tolist()
26     for y in yvalueList:
27         trainSet = X_train[y_train == y].copy()
28         yRes = self.knn_vn(trainSet, k, X_test)
29         for i in range(len(yRes)):
30             # 为 res 中的每一行的字典添加键值对，键为 y，值为 yRes 中对应的元素
31             res[i][y] = yRes[i]
32     return res
33
34 # 计算超球体系数中的gamma函数
35 def gamma(self, x):
36     # 递归利用已知的Gamma(1/2) = sqrt(pi)
37     if abs(x - 0.5) < 1e-6:
38         return math.sqrt(math.pi)
39     elif abs(x - 1) < 1e-6:
40         return 1
41     else:
42         return (x - 1) * self.gamma(x - 1)
43
44 # 常规基于欧氏距离数点的计算方式
45 def knn_euler_count(self, X_train, y_train, X_test, k):
46     '''
47     基于欧拉距离的计数方法来计算类条件概率密度
48     :param X_train: 训练集
49     :param y_train: 训练标签
50     :param X_test: 测试集
51     :param k: k近邻的k
52     :return: 测试集类条件概率密度
53     '''
54     res = []
55     # 超球体体积系数  $V_n(r) = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
56     # 令  $\alpha = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
57     pi_power = math.pi ** (X_train.shape[1] / 2)
58     gamma_value = self.gamma((X_train.shape[1] / 2) + 1)
59     alpha = ((pi_power) / gamma_value) if gamma_value != 0 else float('inf')
60     for i, test in enumerate(X_test):
61         dist = np.linalg.norm(X_train - test, axis=1)
62         knn_indices = np.argsort(dist)[:k]
63         knn_labels = y_train[knn_indices]
64         d = dist[knn_indices[-1]]
65         if d == 0:
66             d = 1e-5
67         class_probs = {
68             cls: np.sum(knn_labels == cls) / (X_train.shape[0] * alpha
69 * (d ** X_train.shape[1]))
70             for cls in np.unique(y_train)}
71         res.append(class_probs)
72     return res
73
74 # 马氏距离的计算
75 def mahalanobis_distance(self, x, dataset, invCov):
76     '''
77     计算单个样本与数据集所有样本的马氏距离
78     :param x:

```



```

77         :param dataset:
78         :param invCov:
79         :return:
80         '''
81         dist = np.zeros(dataset.shape[0])
82         dataset = dataset.to_numpy()
83         for i, data in enumerate(dataset):
84             delta = x - data
85             dist[i] = np.sqrt(delta.dot(invCov).dot(delta.T))
86         return dist
87
88     # 基于马氏距离的KNN
89     def KNN_Mahalanobis_Count(self, X_train, y_train, X_test, k):
90         '''
91         基于马氏距离的计数方法来计算类条件概率密度
92         :param X_train: 训练集
93         :param y_train: 训练标签
94         :param X_test: 测试集
95         :param k: k近邻的k
96         :return: 测试集的条件概率密度
97         '''
98         # 超球体体积系数  $V_n(r) = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
99         # 令  $\alpha = \frac{\pi^{n/2}}{\Gamma(\frac{n}{2} + 1)}$ 
100         pi_power = math.pi ** (X_train.shape[1] / 2)
101         gamma_value = self.gamma((X_train.shape[1] / 2) + 1)
102         alpha = ((pi_power) / gamma_value) if gamma_value != 0 else
103         float('inf')
104         cov = np.cov(X_train.T)
105         invCov = inv(cov)
106         res = []
107         for i, test in enumerate(X_test):
108             dist = self.mahalanobis_distance(test, X_train, invCov)
109             KNN_indices = np.argsort(dist)[:k]
110             KNN_labels = y_train[KNN_indices]
111             d = dist[KNN_indices[-1]]
112             if d == 0:
113                 d = 1e-5
114             class_probs = {
115                 cls: np.sum(KNN_labels == cls) / (X_train.shape[0] * alpha
116                 * (d ** X_train.shape[1]))
117                 for cls in np.unique(y_train)}
118             res.append(class_probs)
119         return res
120
121     # predict函数
122     def execute(self, X_train, y_train, X_test, k, typ=0):
123         yvalueList = y_train.unique().tolist()
124         prior_prob = {}
125         for y in yvalueList:
126             prior_prob[y] = X_train[y_train == y].shape[0]
127         if typ == 0:
128             probs = self.KNN_Euler_Count(X_train, y_train,
129             X_test.to_numpy(), k)
130         elif typ == 1:
131             probs = self.KNN_Mahalanobis_Count(X_train, y_train,
132             X_test.to_numpy(), k)
133         elif typ == 2:
134             probs = self.KNN_Vn(X_train, y_train, X_test.to_numpy(), k)

```

```

130
131     predict = np.zeros(X_test.shape[0])
132     # print(probs)
133     for i, class_prob in enumerate(probs):
134         max_prob = -1
135         for y in class_prob:
136             current_prob = class_prob[y] * prior_prob[y]
137             if current_prob > max_prob:
138                 max_prob = current_prob
139                 predict[i] = y
140     return predict
141
142     # 单纯计算类条件概率密度函数
143     def density(self, X_train, y_train, X_test, k, typ=0):
144         yvalueList = y_train.unique().tolist()
145         prior_prob = {}
146         if not isinstance(X_test, np.ndarray):
147             # 如果不是, 转换它为numpy数组
148             X_test = np.array(X_test)
149         for y in yvalueList:
150             prior_prob[y] = X_train[y_train == y].shape[0]
151         if typ == 0:
152             probs = self.knn_Euler_Count(X_train, y_train, X_test, k)
153         elif typ == 1:
154             probs = self.knn_Mahalanobis_Count(X_train, y_train, X_test, k)
155         elif typ == 2:
156             probs = self.knn_Vn(X_train, y_train, X_test, k)
157         return probs

```

随后在随机种子, k保持不变的情况下, 在 `execute` 方法中执行不同的 `typ` 即可返回不同的预测结果。

```

1  # 基于一般欧氏距离方法
2  y_pred = knn.execute(X_train, y_train, X_test,5,0)
3  accuracy = np.mean(y_pred == y_test)
4  print(f'一般欧式-Accuracy: {accuracy}')
5
6  # 基于马氏距离方法
7  y_pred = knn.execute(X_train, y_train, X_test,5,1)
8  accuracy = np.mean(y_pred == y_test)
9  print(f'一般马氏-Accuracy: {accuracy}')
10
11 # 基于单类欧式距离的方法
12 y_pred = knn.execute(X_train, y_train, X_test,5,2)
13 accuracy = np.mean(y_pred == y_test)
14 print(f'独类欧式-Accuracy: {accuracy}')

```

在数据集划分为随机种子为7时, 结果如下:

```

1  一般欧式-Accuracy: 0.995
2  一般马氏-Accuracy: 0.965
3  独类欧式-Accuracy: 0.98

```

可以发现此时采用欧式距离的准确率大于马氏距离。

再采用不同的随机种子, 例如21, 42, 56

当randstate=21时

```
1 一般欧式-Accuracy: 0.98
2 一般马氏-Accuracy: 0.95
3 独类欧式-Accuracy: 0.985
```

当randstate=42时

```
1 一般欧式-Accuracy: 0.97
2 一般马氏-Accuracy: 0.96
3 独类欧式-Accuracy: 0.975
```

当randstate=56时

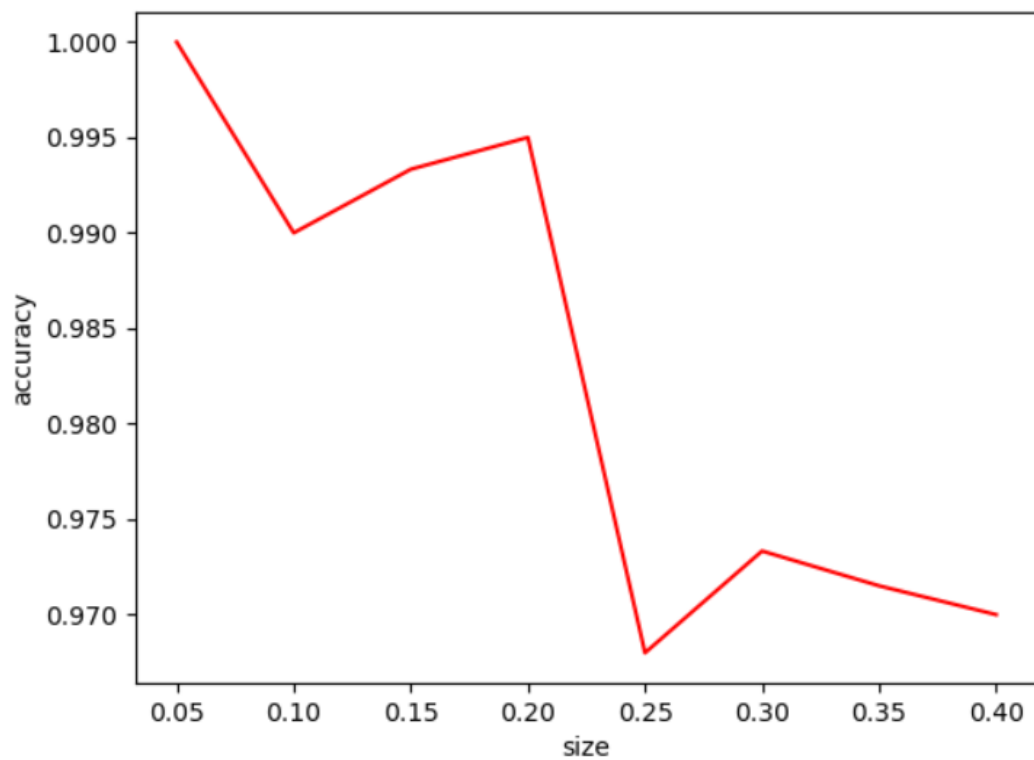
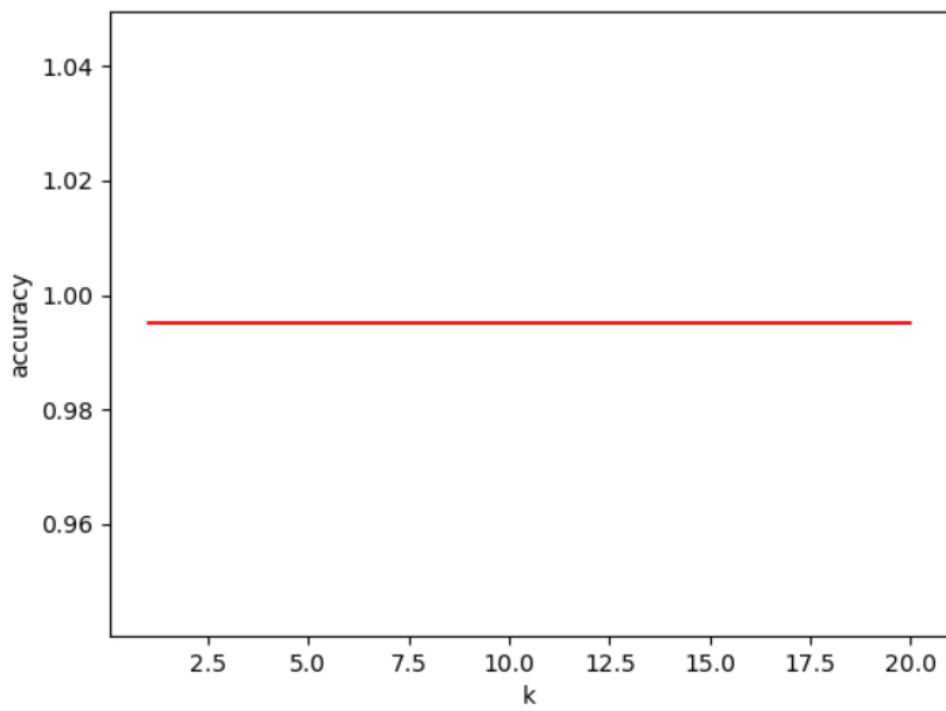
```
1 一般欧式-Accuracy: 0.945
2 一般马氏-Accuracy: 0.94
3 独类欧式-Accuracy: 0.945
```

总的来看，欧式距离应当优于马氏距离在该问题上。

随后执行这两个函数来绘制准确率根据不同的划分（测试集占比）和不同k的变化成都

```
1 def plotAccuracyK(X_train, y_train, X_test, y_test):
2     x = np.linspace(1,20,20)
3     accuracy = np.zeros(20)
4     for k in range(20):
5         knn = kNNClassifier()
6         y_pred = knn.execute(X_train, y_train, X_test,5, 0)
7         accuracy[k]=np.mean(y_pred == y_test)
8     plt.plot(x, accuracy, 'r-')
9     plt.xlabel('k')
10    plt.ylabel('accuracy')
11    plt.show()
12
13 def plotAccuracySize(data, random_state=42):
14     x = np.linspace(0.05,0.4,8)
15     accuracy = np.zeros(8)
16     for i,size in enumerate(x):
17         X_train, y_train, X_test, y_test = splitForData(data, size, 7)
18         y_test = y_test.to_numpy()
19         knn = kNNClassifier()
20         y_pred = knn.execute(X_train, y_train, X_test,5, 0)
21         accuracy[i] = np.mean(y_pred == y_test)
22     plt.plot(x, accuracy, 'r-')
23     plt.xlabel('size')
24     plt.ylabel('accuracy')
25     plt.show()
26
27 plotAccuracyK(X_train, y_train, X_test, y_test)
28 plotAccuracySize(data)
```

得到如下结果：



可以发现随着k的增长，KNN的准确率变化不大在该数据集下。

但随着测试集占比越来越大，KNN的准确率逐步下降。

随后输入这样的数据，查看其是否可以预测

```
1 Temp = [40000,8,0.9]
2 Temp = pd.DataFrame([Temp], columns=['x0', 'x1', 'x2'])
3 y_pred = knn.execute(X_train, y_train, Temp,5,1)
4 index = int(y_pred[0]) # 转换为整数
5 print(yList[index])
```

其输出 `didntLike` 即可以正常进行预测。

实验到此结束。

代码结构

```
1 /Exp2/
2 |----- code/
3 |         |----- Part1 第一部分的代码
4 |         |         |----- DataProcessor.py 数据处理的工具类
5 |         |         |----- exp2-1.py 线性分类主函数代码
6 |         |         |----- LinearClassifier.py 线性分类类
7 |         |----- Part2 第二部分的代码
8 |         |         |----- exp2-2.py 最大似然主函数代码
9 |         |----- Part3 第三部分的代码
10 |        |         |----- exp2-3-kNearestNeighbor.py KNN执行主代码
11 |        |         |----- exp2-3-ParzenWindow.py Parzen窗执行主代码
12 |        |----- Part4 第四部分代码
13 |        |         |----- exp2-4.py KNN实战执行主代码
14 |        |         |----- KNN.py part3和part4有关KNN的KNN类
15 |
16 |----- data/
17 |         |----- e2data1.mat part1数据
18 |         |----- exp2-2.xlsx part2数据
19 |         |----- exp2-3.xlsx part3数据
20 |         |----- e2.txt part4数据
21 |
22 |----- images/实验报告有关图片
23 |
24 |----- Exp2.md 实验报告Markdown
25 |
26 |----- Exp2.pdf 实验报告pdf
```

心得体会

在这次机器学习的基础实验中，我有幸深入探讨了四个核心部分：简单线性分类模型、有参估计、多维统计分析及无参估计，每个部分都让我获得了宝贵的学术及实践经验。

简单线性分类模型：通过实现简单的线性分类模型，我不仅复习了线性代数的基本知识，还学习了如何通过编程将理论应用到实际数据分析中。这一部分的挑战在于选择合适的模型参数和理解模型的决策边界。实验过程中，我通过可视化手段直观地观察了分类效果，这极大地增强了我的直观理解和对模型调优的实践能力。

有参估计中的极大似然估计：在这一部分，我学习了如何在单维和多维情境下使用极大似然估计来计算数据的均值和方差。通过对比单维和多维的计算方法，我更深入地理解了这些统计量在不同维度下的行为和意义，尤其是在处理实际数据集时如何应用这些理论来提取有用的统计信息。

无参估计的Parzen窗和KNN算法：探索无参估计的方法开阔了我的视野，特别是在使用Parzen窗和KNN算法进行数据分类和回归分析方面。我实践了如何根据数据的分布选择合适的窗口大小和邻居数，这对于优化模型性能至关重要。通过实际操作，我学会了调整这些参数以适应不同的数据集，进而优化分类和预测的准确性。

整体反思：这四个实验部分使我认识到，无论是有参还是无参估计，理解其背后的数学原理和如何将这些原理应用于实际问题都是至关重要的。此外，实验不仅提高了我的编程技能，还增强了我对机器学习模型如何在现实世界中应用的认识。我对未来在更复杂数据集上应用这些技术感到兴奋，并期待在未来的学习和研究中继续探索更多机器学习的领域。