

《机器学习》上机实践

031802326-苏炜斌

实验题:

Iris 数据集（鸢尾花数据集）是常用的分类实验数据集，由 Fisher 于 1936 收集整理。数据集包含 150 个数据样本，分为 3 类，每类 50 个数据，每个数据包含 4 个属性。4 个属性分别为花萼长度，花萼宽度，花瓣长度，花瓣宽度，单位是 cm。3 个类别分别为 Setosa（山鸢尾），Versicolour（杂色鸢尾），Virginica（维吉尼亚鸢尾）。

1. Iris 数据集已与常见的机器学习工具集成，请查阅资料找出 MATLAB 平台或 Python 平台加载内置 Iris 数据集方法，并简要描述该数据集结构。

```
{'data': array([[5.1, 3.5, 1.4, 0.2],  
                [4.9, 3. , 1.4, 0.2],  
                [4.7, 3.2, 1.3, 0.2],  
                [4.6, 3.1, 1.5, 0.2],  
                [5. , 3.6, 1.4, 0.2],  
                [5.4, 3.9, 1.7, 0.4],  
                [4.6, 3.4, 1.4, 0.3],  
                [5. , 3.4, 1.5, 0.2],  
                [4.4, 2.9, 1.4, 0.2],  
                [4.9, 3.1, 1.5, 0.1],  
                [5. , 3.5, 1.5, 0.2])
```

[illegible]

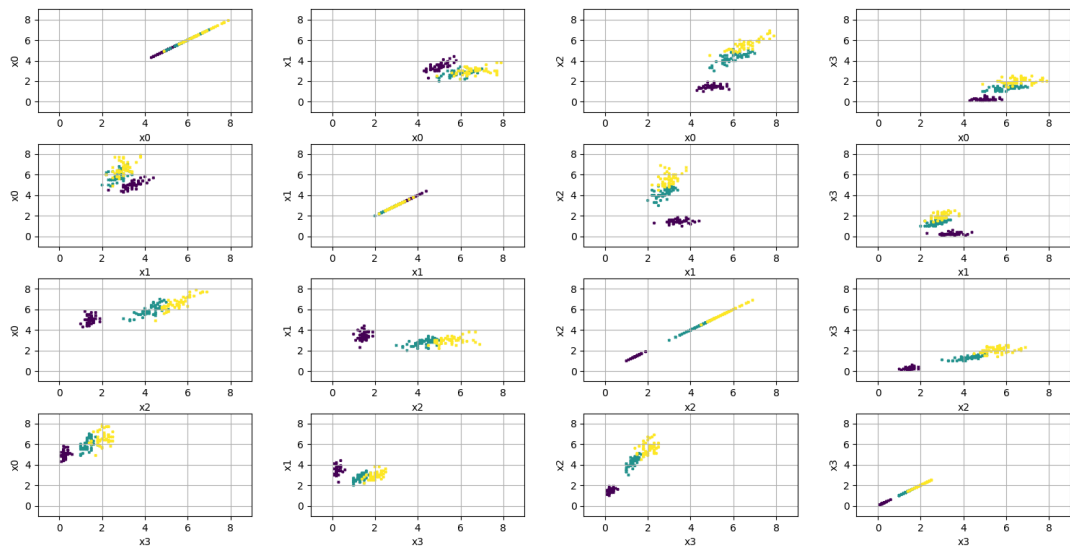
该数据集是一个 map，其中 data 中含有 150 组数据，每组数据有四种属性值；target 中为 data 中数据对应的分类，有 0.1.2 共三种；target names 中为分类值所对应的花名

核心代码

1. `iris = datasets.load_iris()`
2. `iris_data = iris.data`
3. `iris_target = iris.target`
4. `iris_target_names=iris.target_names`
5. `print(iris)`
6. `visualization(iris_data,iris_target,4,4) #可视化`

def visualization(feature, label, m, n)见附录

2. Iris 数据集中有一个种类与另外两个类是线性可分的，其余两个类是线性不可分的。请你通过数据可视化的方法找出该线性可分类并给出判断依据。



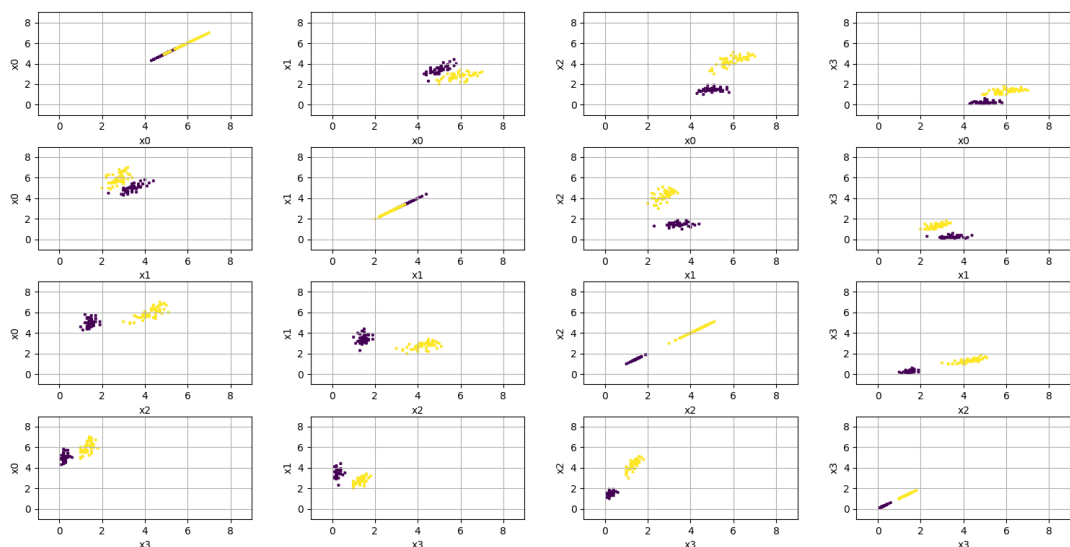
如上图所示 0 是线性可分的，1,2 是线性不可分的，可以通过一个线性函数将其与另外两个类分开

1. # 核心代码
2. `iris_data_linear, iris_target_linear = remove_from_data(iris_data, iris_target, 2)`
3. `visualization(iris_data_linear, iris_target_linear, 4, 4)`

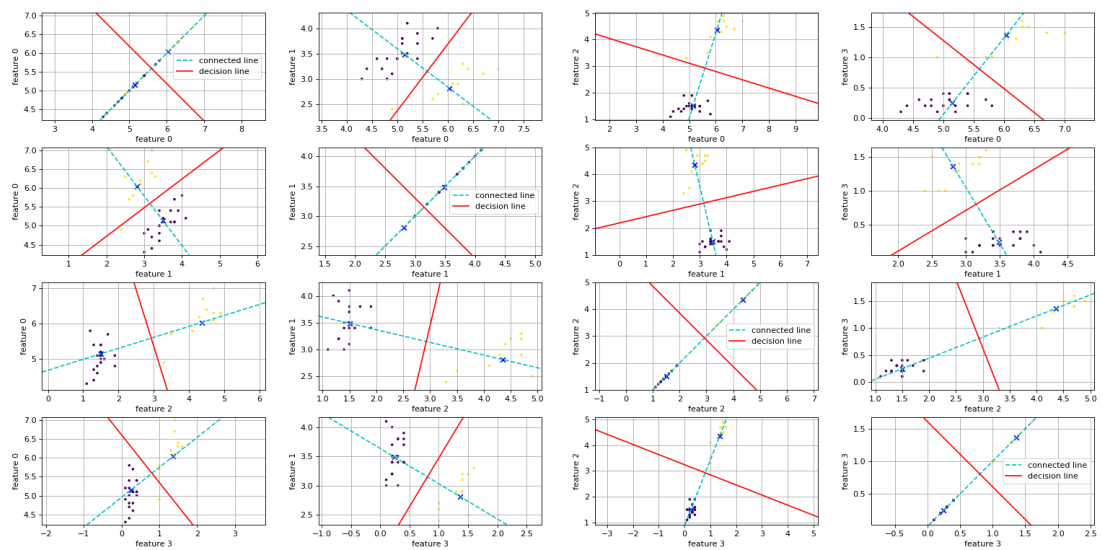
`def remove_from_data(feature, label, num):`见附录

3. 去除 Iris 数据集中线性不可分的类中最后一个，余下的两个线性可分的类构成的数据集命令为 `Iris_linear`，请使用留出法将 `Iris_linear` 数据集按 7:3 分为训练集与测试集，并使用训练集训练一个 MED 分类器，在测试集上测试训练好的分类器的性能，给出《模式识别与机器学习-评估方法与性能指标》中所有量化指标并可视化分类结果。

去除 2 后的散点图



分类效果图



指标:

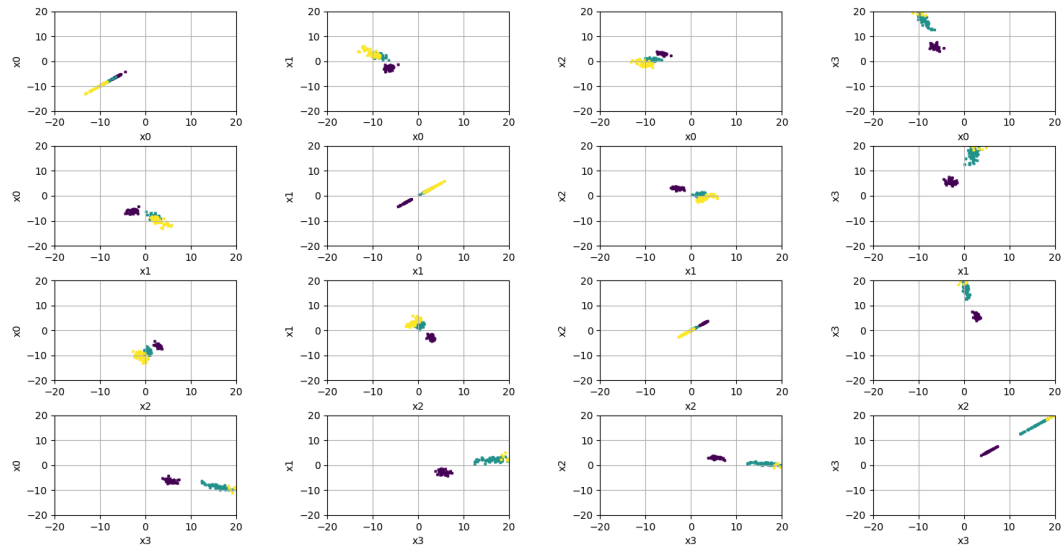
```
tp= 20 tn= 10 fn= 0 fp= 0
accuracy: 1.0 precision: 1.0 recall: 1.0 specificity: 1.0 F1_Score: 1.0
```

4. # 核心代码
5. #划分训练集、测试集
6. `x_train, x_test, y_train, y_test = train_test_split(iris_data_linear, iris_target_linear, test_size=0.3)`
7. `meds=Medclass()`
8. `meds.train(x_train,y_train)`
9. `meds.performance(x_test, y_test, 0)`
10. # 展示每个特征两两对比图，显示决策线
11. `show_decision_line(x_test, y_test, meds, class_1=0, class_2=1, n=4)`

Medclass()详情见附录

4. 将 Iris 数据集白化，可视化白化结果并于原始可视化结果比较，讨论白化的作用。

白化效果图：



白化后的数据

```
[ -6.42454729e+00 -2.96247867e+00  3.11009965e+00  5.13162769e+00]
[ -6.96976967e+00 -3.19637173e+00  2.84250029e+00  6.80915085e+00]
[ -6.02392527e+00 -2.85668753e+00  2.56226112e+00  5.51944956e+00]
[ -6.36609597e+00 -2.45223688e+00  3.06697143e+00  5.29615980e+00]
[ -5.59699408e+00 -1.90026774e+00  2.61377134e+00  4.91982490e+00]
[ -6.18816203e+00 -1.79636987e+00  3.23053159e+00  4.74460434e+00]
[ -6.69565468e+00 -2.91334224e+00  3.33919987e+00  5.41182419e+00]
[ -6.60936652e+00 -2.23021643e+00  3.01090892e+00  5.45648839e+00]
[ -5.86433003e+00 -1.85783773e+00  3.12641560e+00  4.51668055e+00]
[ -5.32784196e+00 -2.47868380e+00  2.83632763e+00  3.84402816e+00]
[ -6.40085096e+00 -4.01584608e+00  3.52964738e+00  4.94195738e+00]
[ -7.16723375e+00 -4.37740790e+00  3.08893523e+00  6.54768183e+00]
[ -6.13728644e+00 -4.05156969e+00  2.73345905e+00  6.02844259e+00]
[ -6.06933698e+00 -3.02546108e+00  2.80576243e+00  5.62189323e+00]
[ -6.98810837e+00 -2.85338894e+00  3.24309462e+00  6.35793730e+00]
[ -6.67731598e+00 -3.25632504e+00  2.93860554e+00  5.86303774e+00]
```

特征白化的目的：将原始特征映射到新的一个特征空间，使得在新空间中特征的协方差为单位矩阵，从而去除特征变化的不同及特征之间的相关性。

特征白化后的数据的特征之间的相关性去除

```
12. # 核心代码
13. iris_data_white = whitening(iris_data)
14. print(iris_data_white)
15. visualization_white(iris_data_white,iris_target,4,4)
16. # 核心代码
17. def whitening(data):
18.     Ex=np.cov(data,rowvar=False) #Ex 为 data 的协方差矩阵
19.     print(Ex.shape)
```

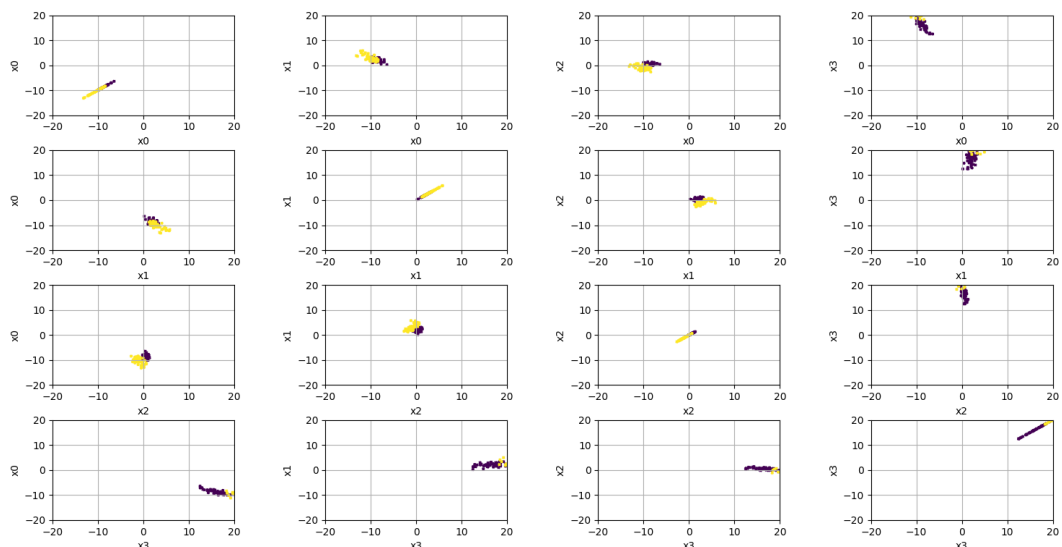
```

20. a, b = np.linalg.eig(Ex) #原始特征协方差矩阵 Ex 的特征值和特征向量
21. #特征向量单位化
22. modulus=[]
23. b=np.real(b)
24. for i in range(b.shape[1]):
25.     sum=0
26.     for j in range(b.shape[0]):
27.         sum+=b[i][j]**2
28.     modulus.append(sum)
29. modulus=np.asarray(modulus,dtype="float64")
30. b=b/modulus
31. #对角矩阵 A
32. a=np.real(a)
33. A=np.diag(a**(-0.5))
34. W=np.dot(A,b.transpose())
35. X=np.dot(W,np.dot(Ex,W.transpose()))
36. for i in range(W.shape[0]):
37.     for j in range(W.shape[1]):
38.         if np.isnan(W[i][j]):
39.             W[i][j]=0
40. print(W)
41. return np.dot(data,W)print(iris_data_white)

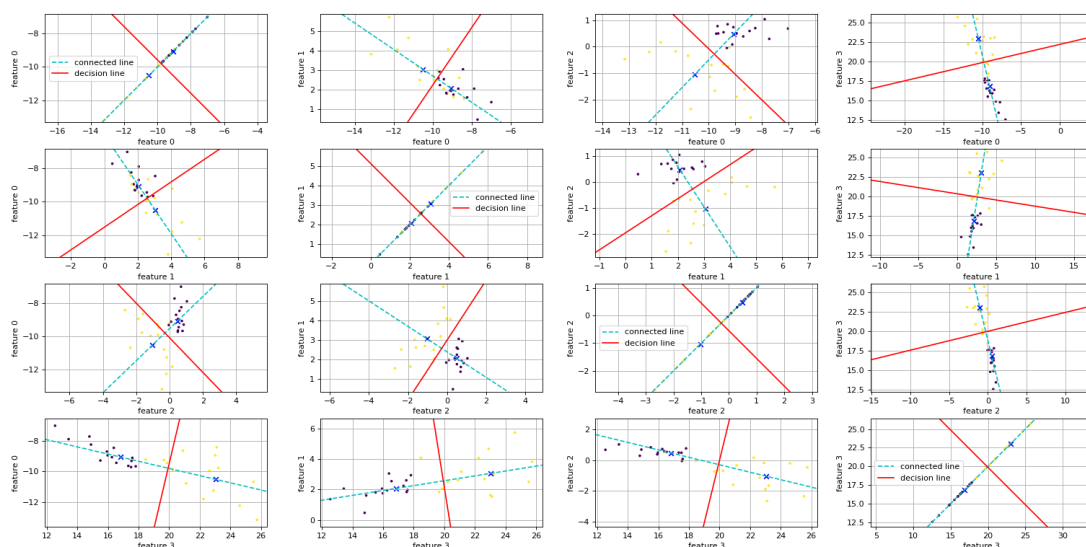
```

5. 去除 Iris 数据集中线性可分的类，余下的两个线性不可分的类构成的数据集命令为 `Iris_nonlinear`，请使用留出法将 `Iris_nonlinear` 数据集按 7:3 分为训练集与测试集，并使用训练集训练一个 MED 分类器，在测试集上测试训练好的分类器的性能，给出《模式识别与机器学习-评估方法与性能指标》中所有量化指标并可视化分类结果。讨论本题结果与 3 题结果的差异。

去除 0 的效果图



分类效果图



指标:

```
tp= 15 tn= 12 fn= 0 fp= 3
accuracy: 0.9 precision: 0.8333333333333334 recall: 1.0 specificity: 0.8 F1_Score: 0.9090909090909091
```

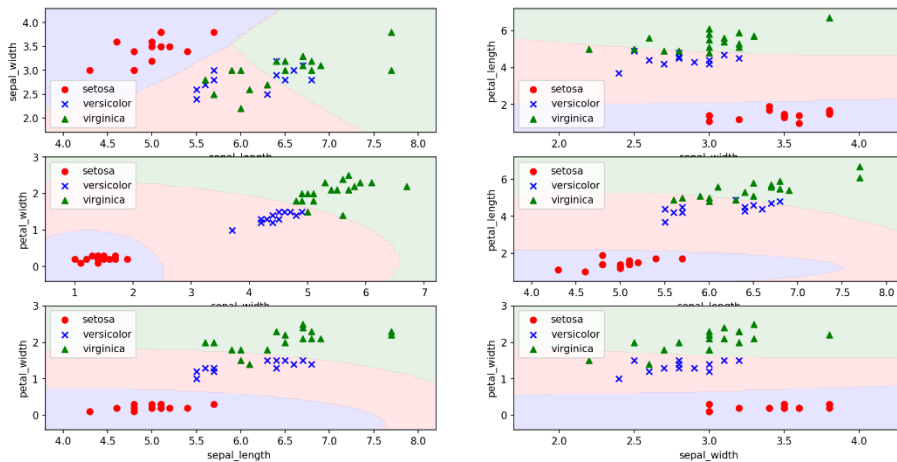
第五题的分类效果与第三题的分类效果较差，正确率下降

42. # 核心代码
43. iris_data_nolinear, iris_target_nolinear = remove_from_data(iris_data_white, iris_target, 0)#白化
44. visualization_white(iris_data_nolinear, iris_target_nolinear, 4, 4)
45. #划分训练集、测试集
46. x_train, x_test, y_train, y_test = train_test_split(iris_data_nolinear, iris_target_nolinear, test_size=0.3)
47. meds2=Medclass()
48. meds2.train(x_train, y_train)
49. meds2.performance(x_test, y_test, 1)
50. # 展示每个特征两两对比图，显示决策线
51. show_decision_line(x_test, y_test, meds2, class_1=1, class_2=2, n=4)

6. 请使用 5 折交叉验证为 Iris 数据集训练一个多分类的贝叶斯分类器。给出平均 Accuracy, 并可视化实验结果。与第 3 题和第 5 题结果做比较, 讨论贝叶斯分类器的优劣。

平均 accuracy: 0.9933333333333334

分类效果图



贝叶斯分类器与其他方法相比最大的优势或许就在于, 它在接受大数据量训练和查询时所具备的高速度。贝叶斯分类器的另一大优势是, 对分类器实际学习状况的解释还是相对简单的。贝叶斯分类器的最大缺陷就是, 它无法处理基于特征组合所产生的变化结果。

与第 3 题和第 5 题结果做比较贝叶斯分类器的分类效果比较好, 准确率高

代码见附录 Bayess.py

```
52. class BayesClassifier(): #贝叶斯分类器,高斯分布概率估计
53.
54. def __init__(self):
55.     self.parameters=[]
56.
57. def train(self,X_data,Y_data):
58.
59.     for categorys in set(Y_data):#遍历每一种类别
60.         selected= Y_data==categorys #选中对应类别的数据
61.         X_newData= X_data[selected] #得到新数据
62.         mean=np.mean(X_newData,axis=0) #得到均值
63.         cov = np.cov(X_newData.transpose())
64.         self.parameters.append(BayesParameter(mean,cov,categorys))
65.
66. def predict(self,data):
67.     res=-1
68.     probability=0
69.     for parameter in self.parameters:
70.         if stats.multivariate_normal.pdf(data, mean=parameter.mean,
cov=parameter.cov)>probability:
```

```
71.             res=parameter.category
72.             probability=stats.multivariate_normal.pdf(data, mean=parameter.mean,
cov=parameter.cov)
73.         return res
74.
75. def    K_Folds_Cross_Validation(data,tar,k):
76.     Set=[]
77.     Tar=[]
78.     for i in range(k):
79.         tempSet=[]
80.         tempTar=[]
81.         tempSet.extend(data[i*10:(i+1)*10])
82.         tempTar.extend(tar[i*10:(i+1)*10])
83.         tempSet.extend(data[(i+5) * 10:(i + 6) * 10])
84.         tempTar.extend(tar[(i+5) * 10:(i + 6) * 10])
85.         tempSet.extend(data[(i+10) * 10:(i + 11) * 10])
86.         tempTar.extend(tar[(i+10) * 10:(i + 11) * 10])
87.         Set.append(tempSet)
88.         Tar.append(tempTar)
89.     return np.asarray(Set),np.asarray(Tar)
```


附录:

Main.py(1-5 题)

```
90. import numpy as np
91. from sklearn import datasets
92. from sklearn.model_selection import train_test_split
93. import matplotlib.pyplot as plt
94. import numpy as np
95. from sklearn.model_selection import train_test_split
96.
97. # MED 分类器
98. class Medclass:
99.     def __init__(self):
100.         self.center_dict = {} # 分类中心点, 以类别标签为键    label: center_point(list)
101.         self.feature_number = 0 # 特征维度
102.         self.train_state = False # 训练状态, True 为训练完成, False 表示还没训练过
103.
104.     def train(self, feature_set, label_set):
105.         new_label_set = {key: value for key, value in enumerate(label_set)} # 将标签集合转换
            为以下标为键的字典    index: label
106.         self.feature_number = len(feature_set[0])
107.         sample_num = len(label_set) # 样本个数
108.         count = {} # 计算每个类别的样本个数    label: count(int)
109.         # 计算每个类别的分类中心点
110.         for index in range(sample_num):
111.             if new_label_set[index] not in count.keys():
112.                 count[new_label_set[index]] = 0
113.             else:
114.                 count[new_label_set[index]] += 1 # 计算对应标签的样本数
115.             if new_label_set[index] not in self.center_dict.keys():
116.                 self.center_dict[new_label_set[index]] = feature_set[index]
117.             else:
118.                 self.center_dict[new_label_set[index]] += feature_set[index]
119.         for _key_ in self.center_dict.keys():
120.             for _feature_ in range(self.feature_number):
121.                 self.center_dict[_key_][_feature_] /= count[_key_]
122.         self.train_state = True
123.
124.         # 根据输入来进行分类预测, 输出以 下标—预测分类 为键值对的字典
125.     def predict(self, feature_set):
126.         # 先判断此分类器是否经过训练
127.         if not self.train_state:
128.             return {}
129.         sample_num = len(feature_set)
```

```

130.         distance_to = {} # 计算某个样本到各分类中心点距离的平方 label: float
131.         result = {} # 保存分类结果 index: label
132.         for _sample_ in range(sample_num):
133.             for _key_ in self.center_dict.keys():
134.                 delta = feature_set[_sample_] - self.center_dict[_key_]
135.                 distance_to[_key_] = np.dot(delta.T, delta)
136.                 result[_sample_] = min(distance_to, key=distance_to.get) # 返回最小值的键（即
label)
137.         return result
138.
139. # 判断预测准确率
140. def accuracy(self, feature_set, label_set):
141.     if not self.train_state:
142.         return 0.0
143.     correct_num = 0
144.     total_num = len(label_set)
145.     predict = self.predict(feature_set)
146.     for _sample_ in range(total_num):
147.         if predict[_sample_] == label_set[_sample_]:
148.             correct_num += 1
149.     return correct_num / total_num
150.
151. # 根据指定的阳性类别，计算分类器的性能指标（准确率 accuracy，精度 precision，召回
率 recall，特异性 specificity，F1_Score）
152. def performance(self, feature_set, label_set, positive):
153.     if not self.train_state:
154.         return {}
155.     total_num = len(label_set)
156.     predict = self.predict(feature_set)
157.     true_positive, false_positive, true_negative, false_negative = 0, 0, 0, 0
158.     for _sample_ in range(total_num):
159.         if predict[_sample_] == label_set[_sample_]:
160.             if label_set[_sample_] == positive:
161.                 true_positive += 1
162.             else:
163.                 true_negative += 1
164.         else:
165.             if label_set[_sample_] == positive:
166.                 false_negative += 1
167.             else:
168.                 false_positive += 1
169.     print("tp=",true_positive,"tn=",true_negative,"fn=",false_negative,"fp=",false_positive)
170.     accuracy = (true_positive + true_negative) / total_num # 准确率（预测正确的样本与
总样本数之比）

```

```

171.         precision = true_positive / (true_positive + false_positive) # 精度（所有 预测 为阳性的
           样本中， 真值 为阳性的比例）
172.         recall = true_positive / (true_positive + false_negative) # 召回率（所有 真值 为阳性的
           样本中， 预测 为阳性的比例）
173.         specificity = true_negative / (true_negative + false_positive) # 特异性（所有 真值 为
           阴性的样本中， 预测 为阴性的比例）
174.         F1_Score = (2 * precision * recall) / (precision + recall) # 精度与召回率的加权平均
175.         print("accuracy:", accuracy, "precision:", precision, "recall:", recall,
           "specificity:", specificity, "F1_Score:", F1_Score)
176.
177.
178.     # 获取某一类的样本中心点
179.     def get_center(self, key):
180.         if key in self.center_dict.keys():
181.             return self.center_dict[key]
182.         else:
183.             return []
184.
185.     def get_center_dict(self):
186.         return self.center_dict
187. #end
188.
189. #画分割线
190.
191. # 展示二维平面上，二分类问题的决策线（class_1 和 class_2）
192. # feature 是样本特征集合，label 是对应的标签集合，对每一维特征进行两两比较，n 表示特征
           维数
193. def show_decision_line(feature, label, med_classifier, class_1=0, class_2=0, n=0):
194.     plt.figure(figsize=(16, 12), dpi=80) # 整张画布大小与分辨率
195.     img = [[] for i in range(n * n)]
196.     for i in range(n):
197.         for j in range(n):
198.             img[i * n + j] = plt.subplot(n, n, i * n + j + 1)
199.             center_1 = med_classifier.get_center(class_1)
200.             center_2 = med_classifier.get_center(class_2)
201.             c_1 = [center_1[i], center_1[j]] # class_1 类中心点的 i, j 两维的分量
202.             c_2 = [center_2[i], center_2[j]] # class_2 类中心点的 i, j 两维的分量
203.             center_3 = [(c_1[0] + c_2[0]) / 2, (c_1[1] + c_2[1]) / 2] # 两点连线的中点
204.             k2, b2 = calculate_vertical_line(c_1, c_2) # 两点中垂线的斜率和截距
205.             plt.scatter(feature[:, i], feature[:, j], c=label, s=20, marker='.') # 整个样本集在特
           征 0 和 2 上的散点图
206.             plt.scatter(c_1[0], c_1[1], c='b', marker='x') # 显示 med 分类器计算的样本中心点
207.             plt.scatter(c_2[0], c_2[1], c='b', marker='x')
208.             plt.grid(True) # 显示网格线

```

```

209.         plt.axis('equal') # 横纵坐标间隔大小相同
210.         plt.axline(c_1, c_2, color='c', linestyle="--", label="connected line")
211.         plt.axline(center_3, slope=k2, color='r', label="decision line")
212.         if i == j:
213.             plt.legend() # 对角线上的子图显示出图例
214.             plt.xlabel("feature " + str(i))
215.             plt.ylabel("feature " + str(j))
216.             plt.tight_layout() # 自动调整子图大小，减少相互遮挡的问题
217.     plt.show()
218.
219.
220. # 计算两点连线，返回斜率和纵截距（假设是二维平面上的点，并且用列表表示）
221. def calculate_connected_line(point_1, point_2):
222.     if len(point_1) != 2 or len(point_2) != 2:
223.         return None
224.     k = (point_1[1] - point_2[1]) / (point_1[0] - point_2[0])
225.     b = (point_1[0] * point_2[1] - point_2[0] * point_1[1]) / (point_1[0] - point_2[0])
226.     return k, b
227.
228.
229. # 计算两点中垂线，返回斜率和纵截距（假设是二维平面上的点，并且用列表表示）
230. def calculate_vertical_line(point_1, point_2):
231.     if len(point_1) != 2 or len(point_2) != 2:
232.         return None
233.     k = -(point_1[0] - point_2[0]) / (point_1[1] - point_2[1])
234.     b = (point_1[1] + point_2[1] + (point_1[0] + point_2[0]) * (point_1[0] - point_2[0]) /
        (point_1[1] - point_2[1]))/2
235.     return k, b
236. #画分割线 end
237.
238. # feature 表示样本特征，label 表示对应的标签,m 行 n 列共计 m*n 个子图
239. def visualization(feature, label, m, n):
240.     plt.figure(figsize=(10, 10), dpi=100)
241.     img = [[] for i in range(m*n)]
242.     for i in range(m):
243.         for j in range(n):
244.             img[i*n+j] = plt.subplot(m, n, i*n+j+1)
245.             plt.xlabel("x"+str(i))
246.             plt.ylabel("x"+str(j))
247.             plt.xlim(-1, 9)
248.             plt.ylim(-1, 9)
249.             plt.scatter(feature[:, i], feature[:, j], s=5, c=label, marker='x')
250.             plt.grid(True) # 显示网格线
251.             plt.tight_layout() # 自动调整子图大小，减少相互遮挡的问题

```

```

252.     plt.show()
253.
254. # feature 表示样本特征, label 表示对应的标签,m 行 n 列共计 m*n 个子图
255. def visualization_white(feature, label, m, n):
256.     plt.figure(figsize=(10, 10), dpi=100)
257.     img = [[] for i in range(m*n)]
258.     for i in range(m):
259.         for j in range(n):
260.             img[i*n+j] = plt.subplot(m, n, i*n+j+1)
261.             plt.xlabel("x"+str(i))
262.             plt.ylabel("x"+str(j))
263.             plt.xlim(-20, 20)
264.             plt.ylim(-20, 20)
265.             plt.scatter(feature[:, i], feature[:, j], s=5, c=label, marker='x')
266.             plt.grid(True) # 显示网格线
267.             plt.tight_layout() # 自动调整子图大小, 减少相互遮挡的问题
268.     plt.show()
269.
270. # 去除某个类别的样本, 返回两个 numpy 数组
271. def remove_from_data(feature, label, num):
272.     new_feature = []
273.     new_label = []
274.     for index in range(len(label)):
275.         if label[index] != num:
276.             new_feature.append(feature[index])
277.             new_label.append(label[index])
278.     return np.asarray(new_feature), np.asarray(new_label)
279.
280.
281. # 特征白化, 返回白化后的矩阵 (numpy 数组格式)
282. # 参数为 numpy 格式的数组, 其格式为数学上的矩阵的转置
283. def whitening(data):
284.     Ex=np.cov(data,rowvar=False) #Ex 为 data 的协方差矩阵
285.     print(Ex.shape)
286.     a, b = np.linalg.eig(Ex) #原始特征协方差矩阵 Ex 的特征值和特征向量
287.     #特征向量单位化
288.     modulus=[]
289.     b=np.real(b)
290.     for i in range(b.shape[1]):
291.         sum=0
292.         for j in range(b.shape[0]):
293.             sum+=b[i][j]**2
294.         modulus.append(sum)
295.     modulus=np.asarray(modulus,dtype="float64")

```

```

296. b=b/modulus
297. #对角矩阵 A
298. a=np.real(a)
299. A=np.diag(a**(-0.5))
300. W=np.dot(A,b.transpose())
301. X=np.dot(W,np.dot(Ex,W.transpose()))
302. for i in range(W.shape[0]):
303.     for j in range(W.shape[1]):
304.         if np.isnan(W[i][j]):
305.             W[i][j]=0
306. print(W)
307. return np.dot(data,W)
308.
309.
310. if __name__ == '__main__':
311.     iris = datasets.load_iris()
312.     iris_data = iris.data
313.     iris_target = iris.target
314.     iris_target_names=iris.target_names
315.     print(iris)
316.     #可视化
317.     visualization(iris_data,iris_target,4,4)
318.
319.     #去除线性不可分的最后一个
320.     iris_data_linear, iris_target_linear = remove_from_data(iris_data, iris_target, 2)
321.     visualization(iris_data_linear,iris_target_linear,4,4)
322.     #划分训练集、测试集
323.     x_train, x_test, y_train, y_test = train_test_split(iris_data_linear, iris_target_linear,
        test_size=0.3)
324.     meds=Medclass()
325.     meds.train(x_train,y_train)
326.     meds.performance(x_test, y_test, 0)
327.     # 展示每个特征两两对比图，显示决策线
328.     show_decision_line(x_test, y_test, meds, class_1=0, class_2=1, n=4)
329.
330.     #特征白化
331.     iris_data_white = whitening(iris_data)
332.     print(iris_data_white)
333.     visualization_white(iris_data_white,iris_target,4,4)
334.
335.     #去除线性可分的类
336.     #iris_data_nolinear, iris_target_nolinear = remove_from_data(iris_data, iris_target, 0) #无白
        化
337.     #visualization(iris_data_nolinear,iris_target_nolinear,4,4)

```

```

338.     iris_data_nolinear, iris_target_nolinear = remove_from_data(iris_data_white, iris_target, 0)#
        白化
339.     visualization_white(iris_data_nolinear,iris_target_nolinear,4,4)
340.     #划分训练集、测试集
341.     x_train, x_test, y_train, y_test = train_test_split(iris_data_nolinear, iris_target_nolinear,
        test_size=0.3)
342.     meds2=Medclass()
343.     meds2.train(x_train,y_train)
344.     meds2.performance(x_test, y_test, 1)
345.     # 展示每个特征两两对比图，显示决策线
346.     show_decision_line(x_test, y_test, meds2, class_1=1, class_2=2, n=4)

```

Bayess.py(第 6 题)

```

1.  import numpy as np
2.  from sklearn import datasets
3.  from sklearn.naive_bayes import GaussianNB
4.  from sklearn.model_selection import train_test_split
5.  import matplotlib.pyplot as plt
6.  import numpy as np
7.  from sklearn.model_selection import train_test_split
8.  import matplotlib.pyplot as plt
9.  import numpy as np
10. from scipy import stats
11. plt.rcParams['savefig.dpi'] = 150 #图片像素
12. plt.rcParams['figure.dpi'] = 100 #分辨率
13. #贝叶斯分类器
14. class BayesParameter(): #存储贝叶斯分类器参数
15.
16.     def __init__(self,mean,cov,category):
17.         self.mean=mean
18.         self.cov=cov
19.         self.category=category
20.
21. class BayesClassifier(): #贝叶斯分类器,高斯分布概率估计
22.
23.     def __init__(self):
24.         self.parameters=[]
25.
26.     def train(self,X_data,Y_data):
27.
28.         for categorys in set(Y_data):#遍历每一种类别
29.             selected= Y_data==categorys #选中对应该类别的数据
30.             X_newData= X_data[selected] #得到新数据
31.             mean=np.mean(X_newData,axis=0) #得到均值

```

```

32.         cov = np.cov(X_newData.transpose())
33.         self.parameters.append(BayesParameter(mean,cov,categrys))
34.
35.     def predict(self,data):
36.         res=-1
37.         probability=0
38.         for parameter in self.parameters:
39.             if stats.multivariate_normal.pdf(data, mean=parameter.mean,
cov=parameter.cov)>probability:
40.                 res=parameter.category
41.                 probability=stats.multivariate_normal.pdf(data, mean=parameter.mean,
cov=parameter.cov)
42.         return res
43.
44.
45. def    K_Folds_Cross_Validation(data,tar,k):
46.     Set=[]
47.     Tar=[]
48.     for i in range(k):
49.         tempSet=[]
50.         tempTar=[]
51.         tempSet.extend(data[i*10:(i+1)*10])
52.         tempTar.extend(tar[i*10:(i+1)*10])
53.         tempSet.extend(data[(i+5) * 10:(i + 6) * 10])
54.         tempTar.extend(tar[(i+5) * 10:(i + 6) * 10])
55.         tempSet.extend(data[(i+10) * 10:(i + 11) * 10])
56.         tempTar.extend(tar[(i+10) * 10:(i + 11) * 10])
57.         Set.append(tempSet)
58.         Tar.append(tempTar)
59.     return np.asarray(Set),np.asarray(Tar)
60.
61. def data_visualization(data,tar):
62.     trainSet,testSet, trainTar,testTar = train_test_split(data, tar, test_size=0.3)
63.     bc = BayesClassifier()
64.     bc.train(trainSet, trainTar)
65.     testPredict = np.array([bc.predict(x) for x in testSet],dtype="int")
66.     # 画图部分
67.     fig = plt.figure(figsize=(8, 8))
68.     xx = [[0, 1], [1, 2], [2, 3], [0,2],[0,3],[1,3]]
69.     yy = ["sepal_length", "sepal_width",
70.         ["sepal_width", "petal_length"],
71.         ["sepal_width", "petal_width"],
72.         ["sepal_length", "petal_length"],
73.         ["sepal_length ", "petal_width"],

```



```

74.         ["sepal_width","petal_width"])
75.     for i in range(6):
76.         ax = fig.add_subplot(321 + i)
77.         x_max,x_min=testSet.max(axis=0)[xx[i][0]]+0.5,testSet.min(axis=0)[xx[i][0]]-0.5
78.         y_max,y_min=testSet.max(axis=0)[xx[i][1]]+0.5,testSet.min(axis=0)[xx[i][1]]-0.5
79.         xlist = np.linspace(x_min, x_max, 80)
80.         ylist = np.linspace(y_min, y_max, 100)
81.         XX, YY = np.meshgrid(xlist, ylist)
82.         bc = GaussianNB()
83.         bc.fit(trainSet[:, xx[i]],trainTar)
84.         xys = [np.array([xx, yy]).reshape(1, -1) for xx, yy in zip(np.ravel(XX), np.ravel(YY))]
85.         zz = np.array([bc.predict(x) for x in xys])
86.         Z = zz.reshape(XX.shape)
87.         plt.contourf(XX, YY, Z, 2, alpha=.1, colors=('blue', 'red', 'green'))
88.         ax.scatter(testSet[testPredict == 0, xx[i][0]], testSet[testPredict == 0, xx[i][1]],
89.                    c='r', marker='o',
90.                    label="setosa")
91.         ax.scatter(testSet[testPredict==1, xx[i][0]], testSet[testPredict==1, xx[i][1]], c='b',
92.                    marker='x',
93.                    label="versicolor")
94.         ax.scatter(testSet[testPredict==2, xx[i][0]], testSet[testPredict==2, xx[i][1]], c='g',
95.                    marker='^',
96.                    label="virginica")
97.         ax.set_xlabel(yy[i][0])
98.         ax.set_ylabel(yy[i][1])
99.         ax.legend(loc=0)
100. plt.show()
101.
102. datas=datasets.load_iris()
103. data=datas.data
104. tar=datas.target
105. data_visualization(data,tar)
106.
107. if __name__=="__main__":
108.     sets,tar=K_Folds_Cross_Validation(data,tar,5)
109.     accuracy=0
110.     print(tar[0].shape)
111.     for i in range(5): #第 i 个子集作为测试集
112.         x,y=0,0
113.         X_data,Y_data=None,None
114.         for j in range(5):
115.             if i!=j:
116.                 if x*y==0:
117.                     X0_data=sets[i]

```

```
116.             Y0_data=tar[i]
117.         else:
118.             X0_data=np.concatenate((X0_data,sets[i]),axis=0)
119.             Y0_data = np.concatenate((Y0_data, tar[i]), axis=0)
120.             x+=1
121.             y+=1
122.
123.         bc= BayesClassifier()
124.         bc.train(X0_data,Y0_data)
125.
126.         y_predict=[bc.predit(x) for x in sets[i]]
127.         tempAccuracy=np.sum(y_predict==tar[i])/tar[i].shape[0]
128.         accuracy+=tempAccuracy
129.
130.     accuracy=accuracy/5
131. print("accuracy:",accuracy)
```