



西南财经大学
SOUTHWESTERN UNIVERSITY OF FINANCE AND ECONOMICS

Name: 兰文婕

Mission objectives: iris

Course name: 多元统计

Date: 2022 年 5 月 17 日

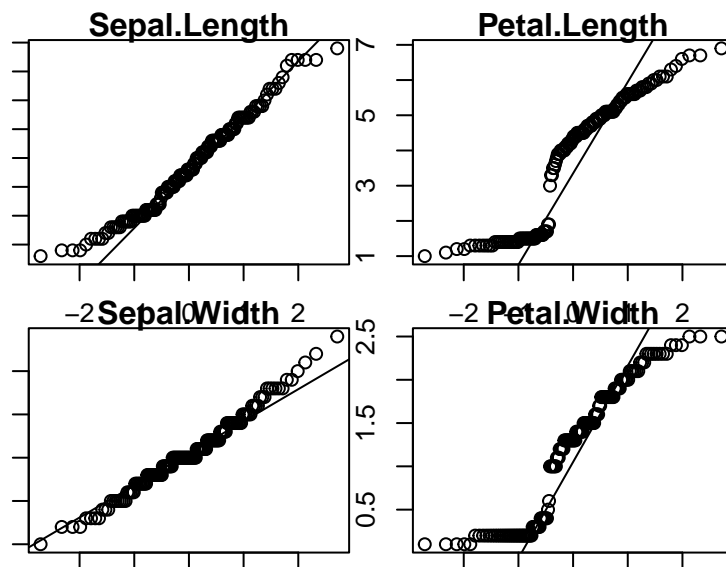
1 正态性检验

1.1 单变量的多元正态性检验

单因变量正态是多因变量多元正态的必要非充分条件。

```
1 iris_d=iris[,1:4]
2
3 #检验每一个变量的qqplot
4 layout(matrix(1:8,nc=4))
5 sapply(colnames(iris_d),function(x)
6 {qqnorm(iris_d[[x]],main=x)
7 qqline(iris_d[[x]])
8 })
9
```

code 1: 单变量多元性检验



```
## $Sepal.Length
## NULL
##
## $Sepal.Width
## NULL
##
## $Petal.Length
## NULL
##
## $Petal.Width
## NULL
```

Figure 1: 1.1

由上图可知,sepal.length 与 sepal.width 的概率图大致均匀分布在 $y=x$ 线两侧,petal.length 与 petal.width 的图明显偏离 $y=x$ 线,左边朝直线下方弯曲,右边朝直线上方弯曲,具有长尾分布的特征。

1.2 二维散点图

判定指标: 通过变量间的两两散点图,如果存在不服从线性关系的两个变量,也就是散点图不成直线,则说明数据不服从多元正态分布。

```
1 pairs(iris_d)
2
```

code 2: 变量之间相关性检验

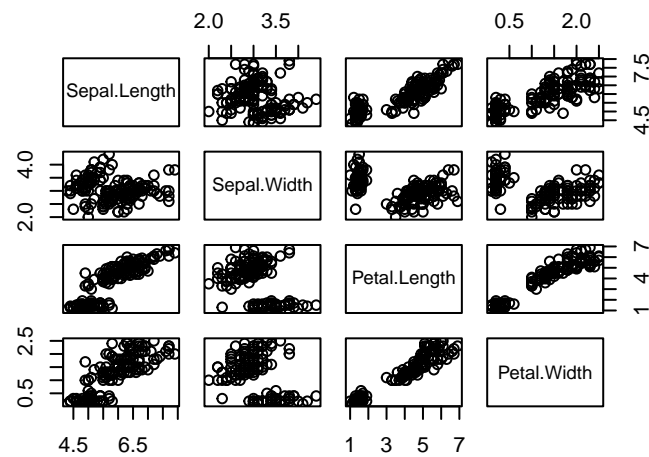


Figure 2: 1.2

1.3 马氏距离是否服从卡方分布

```

1  y <- iris_d
2  cm <- colMeans(y)
3  S <- cov(y)
4  d <- apply(y, 1, function(y) t(y - cm) %*% solve(S) %*% (y - cm))
5  plot(qc <- qchisq((1:nrow(y) - 1/2) / nrow(y), df = 7), sd <- sort(d),
6       xlab = expression(paste(chi[7]^2, " Quantile")),
7       ylab = "Ordered distances", xlim = range(qc) * c(1, 1.1))
8  oups <- which(rank(abs(qc - sd), ties = "random") > nrow(y) - 3)
9  text(qc[oups], sd[oups] - 1.5, names(oups))
10 abline(a = 0, b = 1)
11

```

code 3: 马氏距离卡方分布检验

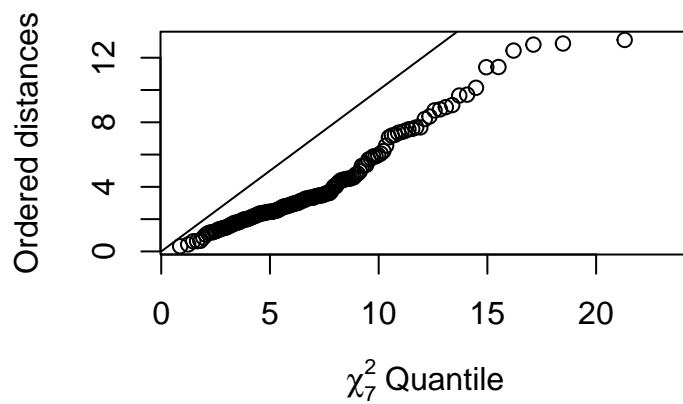


Figure 3: 1.3

从 qq 图可以看出，点距离标准线的偏离程度较大，说明可以拒绝数据服从正态分布的假设。

1.4 夏皮罗-威尔克检验 Shapiro-Wilk test

```
1 library(mvnormtest)
2 iris_d=t(iris[,1:4])
3 mshapiro.test(iris_d)
4
```

code 4: Shapiro-Wilk test

```
##
## Shapiro-Wilk normality test
##
## data:  Z
## W = 0.97935, p-value = 0.02342
```

Figure 4: 1.4

p-value<0.05, 说明可以在 0.05 的水平上拒绝数据服从多元正态分布的假设。

1.5 总结

正态性检验在样本量不大时是有必要的，当样本量足够大时，由中心极限定理，不必太过纠结正态性的问题。

2 分类

使用 python。

2.1 LDA

2.1.1 降维至三维

```
1  #导入各种程序包
2  import numpy as np
3  import pandas as pd
4  import seaborn as sns
5  import matplotlib.pyplot as plt
6  from sklearn.metrics import classification_report
7  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
8  clf = LinearDiscriminantAnalysis()
9  from sklearn.model_selection import train_test_split
10 ##划分训练集与测试集
11 from sklearn.datasets import load_iris
12 iris = load_iris()
13 X_train=iris.data
14 y_train=iris.target
15 X_train,X_test,y_train,y_test=train_test_split(X_train,y_train,
16         test_size=0.4,random_state=0,stratify=y_train)
17 #三维 (数据可视化)
18 def plot_LDA(converted_X,y):
19     '''
20     绘制经过 LDA 转换后的数据
21     '''
22     from mpl_toolkits.mplot3d import Axes3D
23     fig=plt.figure()
24     ax=Axes3D(fig)
25     colors=['dodgerblue','lime','darkorange']
26     target_names = iris.target_names
27     markers='o*s'
28
29     for target,color,marker,target_name in zip([0,1,2],colors,markers,target_names):
30         pos=(y==target).ravel()
31         X=converted_X[pos,:]
32
33         ax.scatter(X[:,0], X[:,1], X[:,2],color=color,marker=marker,
34                 label=target_name)
35     ax.legend(loc="best")
36     fig.suptitle("After LDA")
37 def plot_LDA1(converted_X,y):
38     '''
39     绘制未经 LDA 转换后的数据
40
```

```

41     '''
42     from mpl_toolkits.mplot3d import Axes3D
43     fig=plt.figure()
44     ax=Axes3D(fig)
45     colors=['dodgerblue','lime','darkorange']
46     target_names = iris.target_names
47     markers='o*s'
48
49     for target,color,marker,target_name in zip([0,1,2],colors,markers,target_names):
50         pos=(y==target).ravel()
51         X=converted_X[pos,:]
52
53         ax.scatter(X[:,0], X[:,1], X[:,2],color=color,marker=marker,
54                 label=target_name)
55     ax.legend(loc="best")
56     fig.suptitle("Before LDA")
57     clf.fit(X_train,y_train)
58 # converted_X(150,3) 对数据降维了      权值lda.coef_(3, 4)  lda.intercept_(3,)
59 converted_X=np.dot(X_train,np.transpose(clf.coef_))+clf.intercept_ # X*权值+偏置b 就是输出值
60 plot_LDA(converted_X,y_train)
61 #最初未转化的三维的情况
62 plot_LDA1(X_train,y_train)
63 y_pred = clf.predict(X_test)
64 print(y_pred)
65 print(y_test)
66 #评估
67 # 分类报告: precision/recall/f1-score/均值/分类个数
68 from sklearn.metrics import classification_report
69 target_names = ['class 0', 'class 1', 'class 2']
70 print(classification_report(y_test, y_pred, target_names=target_names))
71 #准确率
72 from sklearn.metrics import accuracy_score
73 accuracy_score(y_test, y_pred)
74 # 混淆矩阵
75 from sklearn.metrics import confusion_matrix
76 confusion_matrix(y_test, y_pred)
77

```

code 5: 3—dimension-LDA

运行结果:

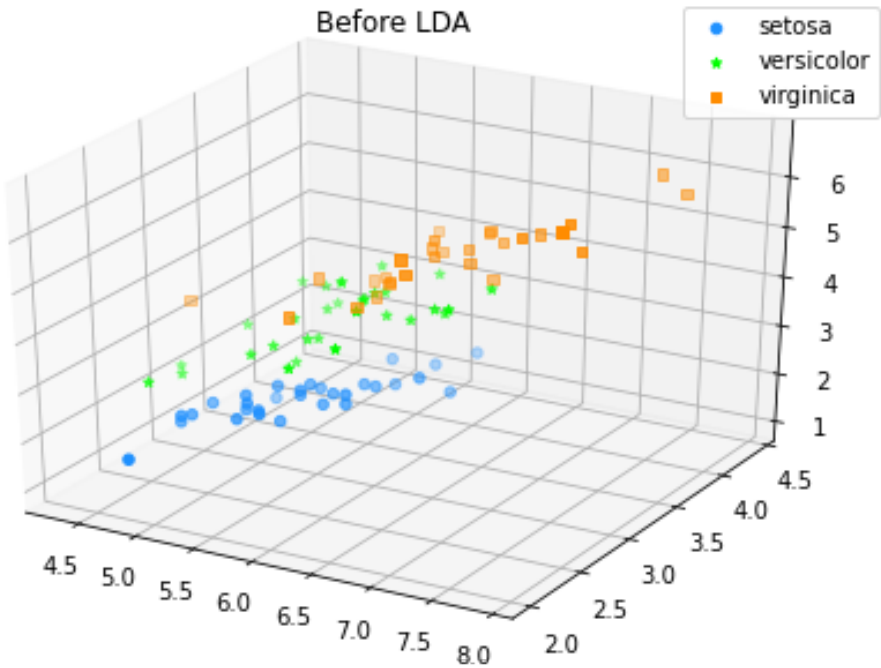


Figure 5: Before 3-dimension-LDA

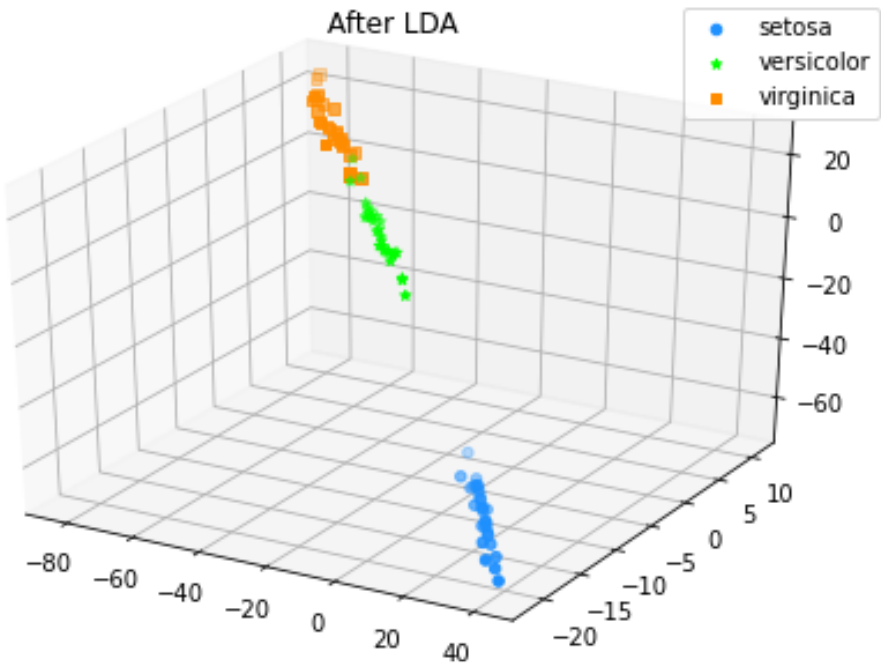


Figure 6: After 2-dimension-LDA

2.1.2 降维至二维

大致过程与上述相同，下面只给出了关键部分的代码。

```
1 #二维
```

```
2 #数据训练
3 lda = LinearDiscriminantAnalysis(n_components=2)
4 X_transform = lda.fit(X_train, y_train).transform(X_test)
5 y_pred = lda.predict(X_test)
6 #绘图
7 colors=['dodgerblue','lime','darkorange']
8 plt.figure()
9 target_names = iris.target_names
10 lw=2
11 for color, i, target_name in zip(colors, [0, 1, 2], target_names):
12     plt.scatter(X_transform[y_pred == i, 0], X_transform[y_pred == i, 1], alpha=.8, color=color,
13                label=target_name)
14 plt.legend(loc='best', shadow=False, scatterpoints=1)
15 plt.title('After 2-dimension-LDA')
16 plt.show()
17
```

code 6: 2—dimension-LDA

2.1.3 运行结果

Tabel 1: 评估分类报告 LDA

	precision	recall	f1-score	support
class 0	1.00	1.00	1.00	20
class 1	1.00	1.00	1.00	20
class 2	1.00	1.00	1.00	20
accuracy	1.00	60		
macro avg	1.00	1.00	1.00	60
weighted avg	1.00	1.00	1.00	60

准确率: 1.0

混淆矩阵:

20 0 0
0 20 0
0 0 20

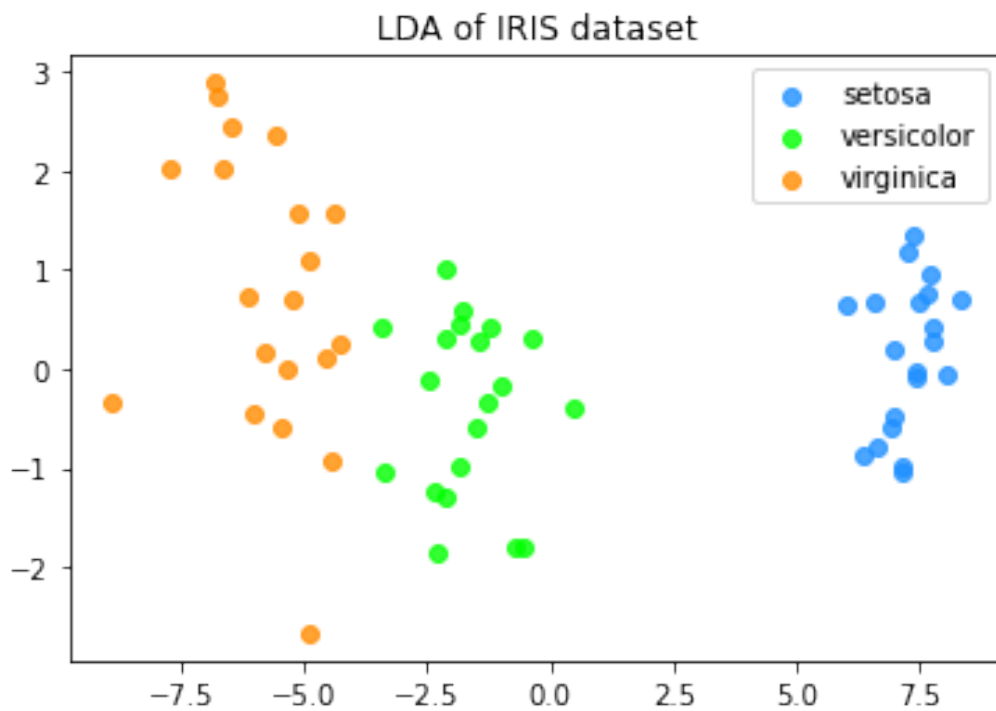


Figure 7: After 2-dimension-LDA

2.2 SVM

```
1  #导入各种程序包
2  import numpy as np
3  import pandas as pd
4  import seaborn as sns
5  import matplotlib.pyplot as plt
6  from sklearn.metrics import classification_report
7  from sklearn.model_selection import train_test_split
8  from sklearn.model_selection import GridSearchCV
9  from sklearn.metrics import classification_report
10 from sklearn.svm import SVC
11 #网格搜索调参
12 tuned_parameters = [{'kernel': ['rbf'], 'gamma': [1e-3, 1e-4],
13                      'C': [1, 10, 100, 1000]},
14                      {'kernel': ['linear'], 'C': [1, 10, 100, 1000]}]
15 #评分方法定义
16 scores = ['precision', 'recall']
17 for score in scores:
18     print("# Tuning hyper-parameters for %s" % score)
19     print()
20
21     # 调用 GridSearchCV, 将 SVC(), tuned_parameters, cv=5, 还有 scoring 传递进去,
22     clf = GridSearchCV(SVC(), tuned_parameters, cv=5, scoring='%s_macro' % score)
23     # 用训练集训练这个学习器 clf
```

```

24     clf.fit(X_train, y_train)
25
26     print("Best parameters set found on development set:")
27     print()
28
29     # 再调用 clf.best_params_ 就能直接得到最好的参数搭配结果
30     print(clf.best_params_)
31
32     print()
33     print("Grid scores on development set:")
34     print()
35     means = clf.cv_results_['mean_test_score']
36     stds = clf.cv_results_['std_test_score']
37
38     # 看一下具体的参数间不同数值的组合后得到的分数是多少
39     for mean, std, params in zip(means, stds, clf.cv_results_['params']):
40         print("%0.3f (+/-%0.03f) for %r"
41               % (mean, std * 2, params))
42     y_pred=clf.predict(X_test)
43

```

code 7: SVM

2.2.1 运行结果

Tabel 2: 评估分类报告 SVM

	precision	recall	f1-score	support
class 0	1.00	1.00	1.00	20
class 1	1.00	0.90	0.95	20
class 2	0.91	1.00	0.95	20
accuracy	0.97	60		
macro avg	0.97	0.97	0.97	60
weighted avg	0.97	0.97	0.97	60

准确率: 0.9666666666666667

混淆矩阵:

```

20 0 0
0 18 2
0 0 20

```

2.3 总结

从上述评估结果对比可知, LDA 在 iris 数据集上的分类效果优于 SVM 的分类效果。