

```
In [1]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
```

## Importing the dataset

```
In [2]: 1 df = pd.read_csv("Wine.csv")
2 df.head()
```

```
Out[2]:
   Alcohol  Malic_Acid  Ash  Ash_Alcinity  Magnesium  Total_Phenols  Flavanoids  Nonflavanoid_Phenols  Proanthocyanins  Color_Intensity  Hue  OD280
0      14.23       1.71    2.43        15.6       127         2.80      3.06        0.28          2.29        5.64     1.04     3.92
1      13.20       1.78    2.14        11.2       100         2.65      2.76        0.26          1.28        4.38     1.05     3.40
2      13.16       2.36    2.67        18.6       101         2.80      3.24        0.30          2.81        5.68     1.03     3.17
3      14.37       1.95    2.50        16.8       113         3.85      3.49        0.24          2.18        7.80     0.86     3.45
4      13.24       2.59    2.87        21.0       118         2.80      2.69        0.39          1.82        4.32     1.04     2.93
```

```
In [3]: 1 df.shape
```

```
Out[3]: (178, 14)
```

```
In [4]: 1 df.isna().sum()
```

```
Out[4]:
Alcohol          0
Malic_Acid       0
Ash              0
Ash_Alcinity     0
Magnesium        0
Total_Phenols    0
Flavanoids       0
Nonflavanoid_Phenols 0
Proanthocyanins 0
Color_Intensity   0
Hue              0
OD280            0
Proline           0
Customer_Segment 0
dtype: int64
```

## Principal Component Analysis (PCA)

常用于高维数据的降维，可用于提取数据的主要特征分量

```
In [5]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

```
In [6]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.decomposition import PCA
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
8
9 sc = StandardScaler()
10 X_train = sc.fit_transform(X_train)
11 X_test = sc.transform(X_test)
12
13 pca = PCA(n_components=2) # creating 2 principal features (X)
14 X_train = pca.fit_transform(X_train)
15 X_test = pca.transform(X_test)
16
17 lr = LogisticRegression(random_state=0)
18 lr.fit(X_train, y_train)
19
20 predictions = lr.predict(X_test)
21
22 print(confusion_matrix(y_test, predictions))
23 print(classification_report(y_test, predictions))
24 print(accuracy_score(y_test, predictions))
```

```
[14  0  0]
[ 1 15  0]
[ 0  0  6]]
```

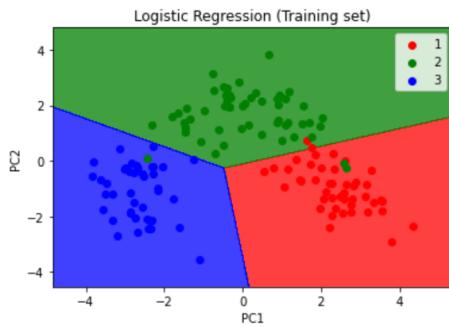
	precision	recall	f1-score	support
1	0.93	1.00	0.97	14
2	1.00	0.94	0.97	16
3	1.00	1.00	1.00	6

	precision	recall	f1-score	support
accuracy			0.97	36
macro avg	0.98	0.98	0.98	36
weighted avg	0.97	0.97	0.97	36

0.9722222222222222

```
In [7]: 1 from matplotlib.colors import ListedColormap
2 X_set, y_set = X_train, y_train
3 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
4                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
5 plt.contourf(X1, X2, lr.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75, cmap = ListedColormap(('red', 'green', 'blue')))
7 plt.xlim(X1.min(), X1.max())
8 plt.ylim(X2.min(), X2.max())
9 for i, j in enumerate(np.unique(y_set)):
10     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
11                 c = ListedColormap(('red', 'green', 'blue'))(i), label = j)
12 plt.title('Logistic Regression (Training set)')
13 plt.xlabel('PC1')
14 plt.ylabel('PC2')
15 plt.legend()
16 plt.show()

*c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence
in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row
if you intend to specify the same RGB or RGBA value for all points.
*c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence
in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row
if you intend to specify the same RGB or RGBA value for all points.
*c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence
in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row
if you intend to specify the same RGB or RGBA value for all points.
```



```
In [8]: 1 from matplotlib.colors import ListedColormap
2 X_set, y_set = X_test, y_test
3 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
4                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
5 plt.contourf(X1, X2, lr.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75, cmap = ListedColormap(('red', 'green', 'blue')))
7 plt.xlim(X1.min(), X1.max())
8 plt.ylim(X2.min(), X2.max())
9 for i, j in enumerate(np.unique(y_set)):
10     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
11                 c = ListedColormap(('red', 'green', 'blue'))(i), label = j)
12 plt.title('Logistic Regression (Test set)')
13 plt.xlabel('PC1')
14 plt.ylabel('PC2')
15 plt.legend()
16 plt.show()

*c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence
in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row
if you intend to specify the same RGB or RGBA value for all points.
*c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence
in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row
if you intend to specify the same RGB or RGBA value for all points.
*c argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence
in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row
if you intend to specify the same RGB or RGBA value for all points.
```



## Linear Discriminant Analysis (LDA)

线性判别分析

```
In [9]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

```

In [10]: 
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
8
9 sc = StandardScaler()
10 X_train = sc.fit_transform(X_train)
11 X_test = sc.transform(X_test)
12
13 lda = LDA(n_components=2)
14 X_train = lda.fit_transform(X_train, y_train)
15 X_test = lda.transform(X_test)
16
17 lr = LogisticRegression(random_state=0)
18 lr.fit(X_train, y_train)
19
20 predictions = lr.predict(X_test)
21
22 print(confusion_matrix(y_test, predictions))
23 print(classification_report(y_test, predictions))
24 print(accuracy_score(y_test, predictions))

```

```

[[14  0  0]
 [ 0 16  0]
 [ 0  0  6]]
      precision    recall   f1-score   support
          1       1.00     1.00     1.00      14
          2       1.00     1.00     1.00      16
          3       1.00     1.00     1.00       6

   accuracy                           1.00      36
  macro avg       1.00     1.00     1.00      36
weighted avg       1.00     1.00     1.00      36

```

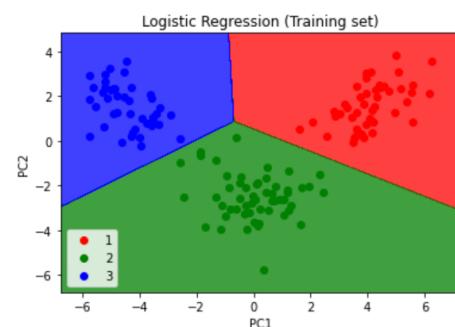
1.0

```

In [11]: 
1 from matplotlib.colors import ListedColormap
2 X_set, y_set = X_train, y_train
3 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
4                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
5 plt.contourf(X1, X2, lr.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75, cmap = ListedColormap(['red', 'green', 'blue']))
7 plt.xlim(X1.min(), X1.max())
8 plt.ylim(X2.min(), X2.max())
9 for i, j in enumerate(np.unique(y_set)):
10     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
11                 c = ListedColormap(['red', 'green', 'blue'])(i), label = j)
12 plt.title('Logistic Regression (Training set)')
13 plt.xlabel('PC1')
14 plt.ylabel('PC2')
15 plt.legend()
16 plt.show()

```

\*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.  
 \*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.  
 \*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.

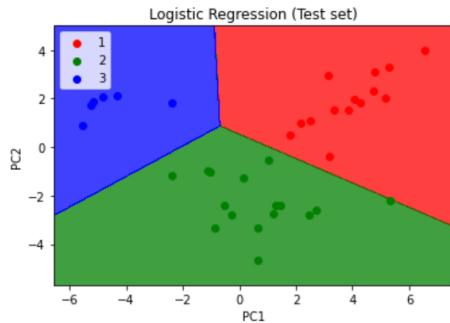


```

In [12]: 
1 from matplotlib.colors import ListedColormap
2 X_set, y_set = X_test, y_test
3 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
4                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
5 plt.contourf(X1, X2, lr.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75, cmap = ListedColormap(['red', 'green', 'blue']))
7 plt.xlim(X1.min(), X1.max())
8 plt.ylim(X2.min(), X2.max())
9 for i, j in enumerate(np.unique(y_set)):
10     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
11                 c = ListedColormap(['red', 'green', 'blue'])(i), label = j)
12 plt.title('Logistic Regression (Test set)')
13 plt.xlabel('PC1')
14 plt.ylabel('PC2')
15 plt.legend()
16 plt.show()

```

```
*c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.
*c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.
** argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.
```



## Kernel PCA

核主成分分析, KernelPCA是PCA的一个改进版. 主成分分析(Principal Components Analysis, PCA)适用于数据的线性降维。而核主成分分析(Kernel PCA, KPCA)可实现数据的非线性降维

```
In [13]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

```
In [14]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.decomposition import KernelPCA
4 from sklearn.linear_model import LogisticRegression
5 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
8
9 sc = StandardScaler()
10 X_train = sc.fit_transform(X_train)
11 X_test = sc.transform(X_test)
12
13 kpca = KernelPCA(n_components=2, kernel='rbf')
14 X_train = kpca.fit_transform(X_train)
15 X_test = kpca.transform(X_test)
16
17 lr = LogisticRegression(random_state=0)
18 lr.fit(X_train, y_train)
19
20 predictions = lr.predict(X_test)
21
22 print(confusion_matrix(y_test, predictions))
23 print(classification_report(y_test, predictions))
24 print(accuracy_score(y_test, predictions))
```

```
[14]   0   0
[ 0 16  0]
[ 0   0  6]]
      precision    recall  f1-score   support
1         1.00     1.00    1.00      14
2         1.00     1.00    1.00      16
3         1.00     1.00    1.00       6

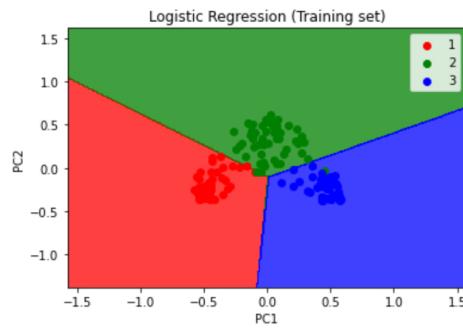
accuracy                           1.00      36
macro avg      1.00     1.00    1.00      36
weighted avg    1.00     1.00    1.00      36
```

1.0

```
In [15]: 1 from matplotlib.colors import ListedColormap
2 X_set, y_set = X_train, y_train
3 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
4                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
5 plt.contourf(X1, X2, lr.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75, cmap = ListedColormap(['red', 'green', 'blue']))
7 plt.xlim(X1.min(), X1.max())
8 plt.ylim(X2.min(), X2.max())
9 for i, j in enumerate(np.unique(y_set)):
10     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
11                 c = ListedColormap(['red', 'green', 'blue'])(i), label = j)
12 plt.title('Logistic Regression (Training set)')
13 plt.xlabel('PC1')
14 plt.ylabel('PC2')
15 plt.legend()
16 plt.show()
```

```
*c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.
*c* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with *x* & *y*. Please use the *color* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.
```

if you intend to specify the same RGB or RGBA value for all points.  
\*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.



```
In [16]: 1 from matplotlib.colors import ListedColormap
2 X_set, y_set = X_test, y_test
3 X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
4                      np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
5 plt.contourf(X1, X2, lr.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75, cmap = ListedColormap(['red', 'green', 'blue']))
7 plt.xlim(X1.min(), X1.max())
8 plt.ylim(X2.min(), X2.max())
9 for i, j in enumerate(np.unique(y_set)):
10     plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
11                 c = ListedColormap(['red', 'green', 'blue'])(i), label = j)
12 plt.title('Logistic Regression (Test set)')
13 plt.xlabel('PC1')
14 plt.ylabel('PC2')
15 plt.legend()
16 plt.show()
```

\*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.

\*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.

\*c\* argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with \*x\* & \*y\*. Please use the \*color\* keyword-argument or provide a 2-D array with a single row if you intend to specify the same RGB or RGBA value for all points.

