

day-66-pca

December 14, 2023

Day 66 Principal Component Analysis

By: Loga Aswin

```
[1]: #import libraries
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
```

```
[3]: #load dataset
df = pd.read_csv("/content/winequality-red.csv")
```

```
[4]: df.head()
```

```
[4]:   fixed acidity  volatile acidity  citric acid  residual sugar  chlorides \
0           7.4             0.70         0.00           1.9        0.076
1           7.8             0.88         0.00           2.6        0.098
2           7.8             0.76         0.04           2.3        0.092
3          11.2             0.28         0.56           1.9        0.075
4           7.4             0.70         0.00           1.9        0.076

      free sulfur dioxide  total sulfur dioxide  density    pH  sulphates \
0              11.0             34.0    0.9978  3.51        0.56
1              25.0             67.0    0.9968  3.20        0.68
2              15.0             54.0    0.9970  3.26        0.65
3              17.0             60.0    0.9980  3.16        0.58
4              11.0             34.0    0.9978  3.51        0.56

      alcohol  quality
0         9.4         5
1         9.8         5
2         9.8         5
3         9.8         6
4         9.4         5
```

```
[19]: df.isnull().sum()
```

```
[19]: fixed acidity      0
      volatile acidity  0
      citric acid       0
      residual sugar    0
      chlorides         0
      free sulfur dioxide 0
      total sulfur dioxide 0
      density          0
      pH              0
      sulphates        0
      alcohol          0
      quality          0
      dtype: int64
```

```
[20]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  -
0   fixed acidity         1599 non-null   float64
1   volatile acidity      1599 non-null   float64
2   citric acid           1599 non-null   float64
3   residual sugar        1599 non-null   float64
4   chlorides             1599 non-null   float64
5   free sulfur dioxide    1599 non-null   float64
6   total sulfur dioxide   1599 non-null   float64
7   density               1599 non-null   float64
8   pH                   1599 non-null   float64
9   sulphates             1599 non-null   float64
10  alcohol               1599 non-null   float64
11  quality               1599 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
```

Splitting into features and target sets

```
[6]: X = df.iloc[:, :-1].values
      y = df.iloc[:, -1].values
```

splitting the data into the training and test set

```
[7]: from sklearn.model_selection import train_test_split
      X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2,
      random_state = 0)
```

```
[9]: # Feature scaling
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

We are importing the PCA class from the decomposition module in sklearn

```
[10]: from sklearn.decomposition import PCA

pca = PCA(n_components = 2)
X_train = pca.fit_transform(X_train)
```

```
[11]: X_test = pca.transform(X_test)
```

Training logistic model on new training dataset

```
[13]: from sklearn.linear_model import LogisticRegression

clf = LogisticRegression(random_state = 0)
clf.fit(X_train, y_train)
```

```
[13]: LogisticRegression(random_state=0)
```

Model Evaluation Metrics

```
[14]: from sklearn.metrics import confusion_matrix, accuracy_score
y_pred = clf.predict(X_test)
```

```
[15]: matrix = confusion_matrix(y_test, y_pred)
print(matrix)
```

```
[[ 0  0  0  2  0  0]
 [ 0  0  4  7  0  0]
 [ 0  0 89 45  1  0]
 [ 0  0 55 81  6  0]
 [ 0  0  4 21  2  0]
 [ 0  0  0  2  1  0]]
```

```
[18]: print("Accuracy:", accuracy_score(y_test, y_pred)*100)
```

Accuracy: 53.75

```
[17]: from sklearn.metrics import classification_report
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
3	0.00	0.00	0.00	2

4	0.00	0.00	0.00	11
5	0.59	0.66	0.62	135
6	0.51	0.57	0.54	142
7	0.20	0.07	0.11	27
8	0.00	0.00	0.00	3
accuracy			0.54	320
macro avg	0.22	0.22	0.21	320
weighted avg	0.49	0.54	0.51	320

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344:
UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
0.0 in labels with no predicted samples. Use `zero_division` parameter to
control this behavior.
```

```
_warn_prf(average, modifier, msg_start, len(result))
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344:
UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
0.0 in labels with no predicted samples. Use `zero_division` parameter to
control this behavior.
```

```
_warn_prf(average, modifier, msg_start, len(result))
```

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344:
UndefinedMetricWarning: Precision and F-score are ill-defined and being set to
0.0 in labels with no predicted samples. Use `zero_division` parameter to
control this behavior.
```

```
_warn_prf(average, modifier, msg_start, len(result))
```