

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn import metrics
from sklearn.neighbors import KNeighborsClassifier

In [2]: df = pd.read_csv('D:\\24 - Machine_Learning\\download files\\winequality-red.csv', sep=";")
df
```

Out[2]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.700	0.00	1.9	0.076	11.0	34.0	0.99780	3.51	0.56	9.4	5
1	7.8	0.880	0.00	2.6	0.098	25.0	67.0	0.99680	3.20	0.68	9.8	5
2	7.8	0.760	0.04	2.3	0.092	15.0	54.0	0.99700	3.26	0.65	9.8	5
3	11.2	0.280	0.56	1.9	0.075	17.0	60.0	0.99800	3.16	0.58	9.8	6
4	7.4	0.700	0.00	1.9	0.076	11.0	34.0	0.99780	3.51	0.56	9.4	5
...
1594	6.2	0.600	0.08	2.0	0.090	32.0	44.0	0.99490	3.45	0.58	10.5	5
1595	5.9	0.550	0.10	2.2	0.062	39.0	51.0	0.99512	3.52	0.76	11.2	6
1596	6.3	0.510	0.13	2.3	0.076	29.0	40.0	0.99574	3.42	0.75	11.0	6
1597	5.9	0.645	0.12	2.0	0.075	32.0	44.0	0.99547	3.57	0.71	10.2	5
1598	6.0	0.310	0.47	3.6	0.067	18.0	42.0	0.99549	3.39	0.66	11.0	6

1599 rows × 12 columns

```
In [3]: 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
```

In [4]:

Out[4]:

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
count	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000	1599.000000
mean	8.319637	0.527821	0.270976	2.538806	0.087467	15.874922	46.467792	0.996747	3.311113	0.658149	10.422983	5.636023
std	1.741096	0.179060	0.194801	1.409928	0.047065	10.460157	32.895324	0.001887	0.154386	0.169507	1.065668	0.807569
min	4.600000	0.120000	0.000000	0.900000	0.012000	1.000000	6.000000	0.990070	2.740000	0.330000	8.400000	3.000000
25%	7.100000	0.390000	0.090000	1.900000	0.070000	7.000000	22.000000	0.995600	3.210000	0.550000	9.500000	5.000000
50%	7.900000	0.520000	0.260000	2.200000	0.079000	14.000000	38.000000	0.996750	3.310000	0.620000	10.200000	6.000000
75%	9.200000	0.640000	0.420000	2.600000	0.090000	21.000000	62.000000	0.997835	3.400000	0.730000	11.100000	6.000000
max	15.900000	1.580000	1.000000	15.500000	0.611000	72.000000	289.000000	1.003690	4.010000	2.000000	14.900000	8.000000

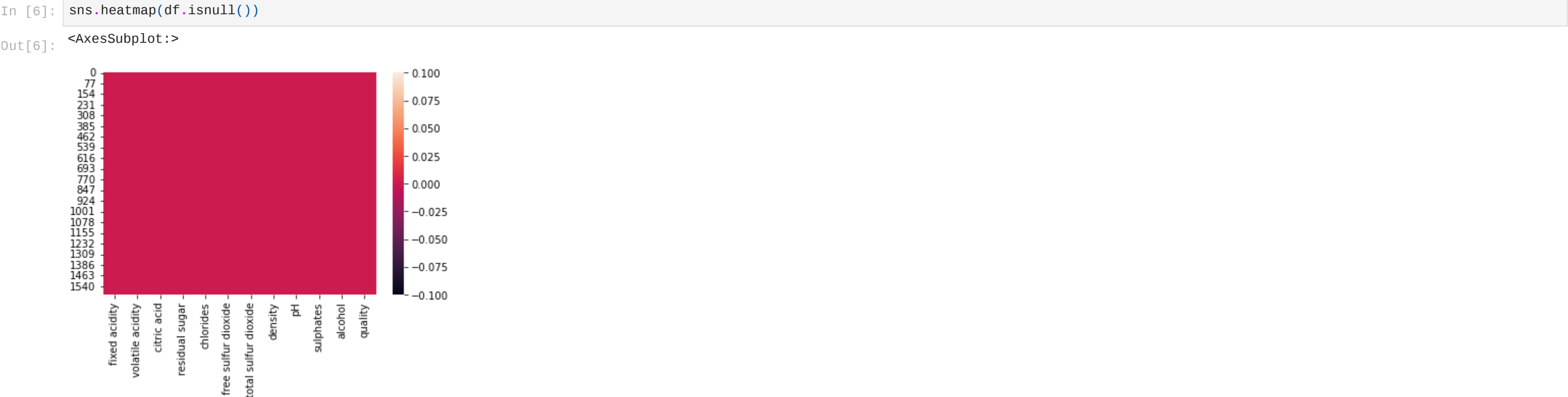
```
In [5]: df.isnull

<bound method DataFrame.isnull of
0      fixed acidity  volatile acidity  citric acid  residual sugar  chlorides \
1      7.4          0.700            0.00         0.00         0.00
2      7.8          0.880            0.00         0.00         0.098
3      7.8          0.760            0.04         0.04         0.092
4     11.2          0.280            0.56         0.56         0.075
...     7.4          0.700            0.00         0.00         0.076
1594    6.2          0.600            0.08         0.08         0.090
1595    5.9          0.550            0.10         0.10         0.062
1596    6.3          0.510            0.13         0.13         0.076
1597    5.9          0.645            0.12         0.12         0.075
1598    6.0          0.310            0.47         0.47         0.067

      free sulfur dioxide  total sulfur dioxide  density  pH  sulphates \
0                11.0                34.0  0.99780  3.51         0.56
1                25.0                67.0  0.99680  3.20         0.68
2                15.0                54.0  0.99700  3.26         0.65
3                17.0                60.0  0.99800  3.16         0.58
4                11.0                34.0  0.99780  3.51         0.56
...                ...                ...     ...     ...     ...
1594             32.0                44.0  0.99490  3.45         0.58
1595             39.0                51.0  0.99512  3.52         0.76
1596             29.0                40.0  0.99574  3.42         0.75
1597             32.0                44.0  0.99547  3.57         0.71
1598             18.0                42.0  0.99549  3.39         0.66

      alcohol  quality
0          9.4         5
1          9.8         5
2          9.8         5
3          9.8         6
4          9.4         5
...         ...     ...
1594       10.5         5
1595       11.2         6
1596       11.0         6
1597       10.2         5
1598       11.0         6

[1599 rows x 12 columns]>
```



```
In [7]: df.shape

Out[7]: (1599, 12)
```

```
In [8]: features = df.drop('quality', axis=1).values
features

Out[8]: array([[ 7.4 ,  0.7 ,  0. , ...,  3.51 ,  0.56 ,  9.4 ],
       [ 7.8 ,  0.88 ,  0. , ...,  3.2 ,  0.68 ,  9.8 ],
       [ 7.8 ,  0.76 ,  0.04 , ...,  3.26 ,  0.65 ,  9.8 ],
       ...,
       [ 6.3 ,  0.51 ,  0.13 , ...,  3.42 ,  0.75 , 11. ],
       [ 5.9 ,  0.645 ,  0.12 , ...,  3.57 ,  0.71 , 10.2 ],
       [ 6. ,  0.31 ,  0.47 , ...,  3.39 ,  0.66 , 11. ]])
```

```
In [9]: classes = df['quality'].values
classes

Out[9]: array([5, 5, 5, ..., 6, 5, 6], dtype=int64)
```

```
In [10]: (train_feat, test_feat, train_classes, test_classes) = train_test_split(features, classes, train_size = 0.8, random_state=40)
```

```
In [11]: knn = KNeighborsClassifier(n_neighbors=2)
```

```
In [12]: knn.fit(train_feat, train_classes)
```

```
Out[12]: KNeighborsClassifier(n_neighbors=2)
```

```
In [13]: pred = knn.predict(test_feat)
```

```
In [14]: print("Accuracy: ", metrics.accuracy_score(test_classes, pred))

Accuracy:  0.496875
```

```
In [15]: neighbors = np.arange(1,9)
```

```
In [16]: train_accuracy = np.empty(len(neighbors))
```

```
In [17]: test_accuracy = np.empty(len(neighbors))
```

```
In [18]: for i,k in enumerate(neighbors):
#Setup a knn classifier with k neighbors
knn = KNeighborsClassifier(n_neighbors=k)
#Fit the model
knn.fit(train_feat, train_classes)
#compute accuracy on training set
train_accuracy[i] = knn.score(train_feat, train_classes)
#compute accuracy on test set
test_accuracy[i] = knn.score(test_feat, test_classes)
```

```
In [21]: plt.title('k-NN Varying number of neighbors')
plt.plot(neighbors, test_accuracy, label='Testing Accuracy')
plt.plot(neighbors, train_accuracy, label='Training Accuracy')
plt.legend()
plt.xlabel('Number of neighbors')
plt.ylabel('Accuracy')
plt.show()
```

