

```
In [24]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix
from sklearn.metrics import f1_score
import statsmodels.api as sm
import sys
import warnings

if not sys.warnoptions:
    warnings.simplefilter("ignore")

In [25]: df = pd.read_csv('glass.csv')
df.head()
```

Out[25]:

	RI	Na	Mg	Al	Si	K	Ca	Ba	Fe	Type
0	1.52101	13.64	4.49	1.10	71.78	0.06	8.75	0.0	0.0	1
1	1.51761	13.89	3.60	1.36	72.73	0.48	7.83	0.0	0.0	1
2	1.51618	13.53	3.55	1.54	72.99	0.39	7.78	0.0	0.0	1
3	1.51766	13.21	3.69	1.29	72.61	0.57	8.22	0.0	0.0	1
4	1.51742	13.27	3.62	1.24	73.08	0.55	8.07	0.0	0.0	1

```
In [26]: df.isnull().sum()
```

Out[26]:

```
RI      0
Na      0
Mg      0
Al      0
Si      0
K        0
Ca       0
Ba       0
Fe       0
Type     0
dtype: int64
```

```
In [27]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 214 entries, 0 to 213
Data columns (total 10 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0    RI      214 non-null      float64
 1    Na      214 non-null      float64
 2    Mg      214 non-null      float64
 3    Al      214 non-null      float64
 4    Si      214 non-null      float64
 5    K        214 non-null      float64
 6    Ca      214 non-null      float64
 7    Ba      214 non-null      float64
 8    Fe      214 non-null      float64
 9    Type    214 non-null      int64
dtypes: float64(9), int64(1)
memory usage: 16.8 KB
```

They all consist of 214 rows, the x columns are float type, the y column is integer type. Data is clean so far.

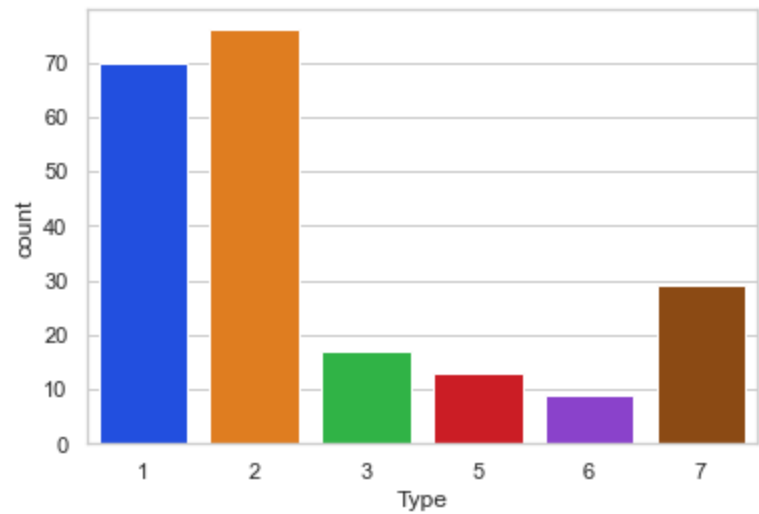
```
In [28]: df['Type'].value_counts().sort_values(ascending=False)
```

Out[28]:

```
2    76
1    70
7    29
3    17
5    13
6     9
Name: Type, dtype: int64
```

```
In [29]: sns.set(style='whitegrid')
sns.countplot(x='Type', data=df, palette='bright')
```

Out[29]: <AxesSubplot:xlabel='Type', ylabel='count'>



```
In [30]: df.describe()
```

Out[30]:

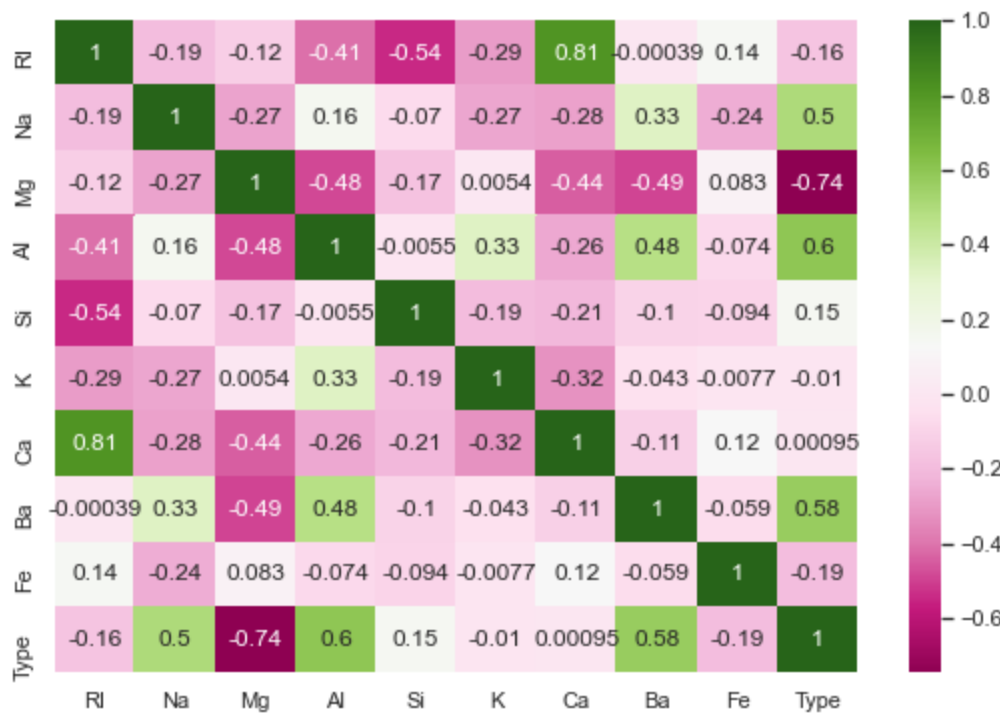
	RI	Na	Mg	Al	Si	K	Ca	Ba	Fe	Type
count	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000	214.000000
mean	1.518365	13.407850	2.684533	1.444907	72.650935	0.497056	8.956963	0.175047	0.057009	2.780374
std	0.003037	0.816604	1.442408	0.499270	0.774546	0.652192	1.423153	0.497219	0.097439	2.103739
min	1.511150	10.730000	0.000000	0.290000	69.810000	0.000000	5.430000	0.000000	0.000000	1.000000
25%	1.516523	12.907500	2.115000	1.190000	72.280000	0.122500	8.240000	0.000000	0.000000	1.000000
50%	1.517680	13.300000	3.480000	1.360000	72.790000	0.555000	8.600000	0.000000	0.000000	2.000000
75%	1.519157	13.825000	3.600000	1.630000	73.087500	0.610000	9.172500	0.000000	0.100000	3.000000
max	1.533930	17.380000	4.490000	3.500000	75.410000	6.210000	16.190000	3.150000	0.510000	7.000000

```
In [31]: df.corr()
```

Out[31]:

	RI	Na	Mg	Al	Si	K	Ca	Ba	Fe	Type
RI	1.000000	-0.191885	-0.122274	-0.407326	-0.542052	-0.289833	0.810403	-0.000386	0.143010	-0.164237
Na	-0.191885	1.000000	-0.273732	0.156794	-0.069809	-0.266087	-0.275442	0.326603	-0.241346	0.502898
Mg	-0.122274	-0.273732	1.000000	-0.481799	-0.165927	0.005396	-0.443750	-0.492262	0.083060	-0.744993
Al	-0.407326	0.156794	-0.481799	1.000000	-0.005524	0.325958	-0.259592	0.479404	-0.074402	0.598829
Si	-0.542052	-0.069809	-0.165927	-0.005524	1.000000	-0.193331	-0.208732	-0.102151	-0.094201	0.151565
K	-0.289833	-0.266087	0.005396	0.325958	-0.193331	1.000000	-0.317836	-0.042618	-0.007719	-0.010054
Ca	0.810403	-0.275442	-0.443750	-0.259592	-0.208732	-0.317836	1.000000	-0.112841	0.124968	0.000952
Ba	-0.000386	0.326603	-0.492262	0.479404	-0.102151	-0.042618	-0.112841	1.000000	-0.058692	0.575161
Fe	0.143010	-0.241346	0.083060	-0.074402	-0.094201	-0.007719	0.124968	-0.058692	1.000000	-0.188278
Type	-0.164237	0.502898	-0.744993	0.598829	0.151565	-0.010054	0.000952	0.575161	-0.188278	1.000000

```
In [32]: cor = df.corr()
plt.figure(figsize=(9,6))
sns.heatmap(data = cor, annot = True, cmap = 'PiYG')
plt.show()
```



## Model using KNN

```
In [33]: # x and y assignment
x = df.iloc[:, [0,1,2,3,4,5,6,7,8]]
y = df.iloc[:,9:]
```

```
In [34]: #Let's see how x and y look
x.head(3)
```

Out[34]:

	RI	Na	Mg	Al	Si	K	Ca	Ba	Fe
0	1.52101	13.64	4.49	1.10	71.78	0.06	8.75	0.0	0.0
1	1.51761	13.89	3.60	1.36	72.73	0.48	7.83	0.0	0.0
2	1.51618	13.53	3.55	1.54	72.99	0.39	7.78	0.0	0.0

```
In [35]: y.head(3)
```

Out[35]:

	Type
0	1
1	1
2	1

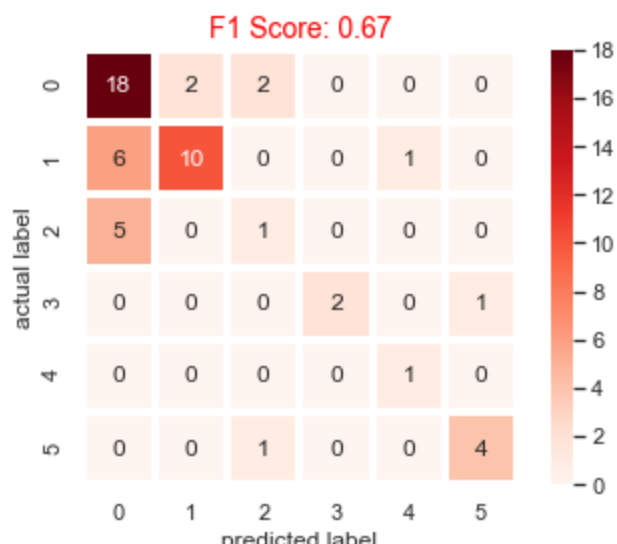
```
In [36]: # Splitting the Data Set into Independent Variables and Dependent Variables
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(x,y, test_size=0.25, random_state=18)
```

```
In [37]: from sklearn.neighbors import KNeighborsClassifier

knn = KNeighborsClassifier(n_neighbors=1, metric = 'minkowski')
knn.fit(x_train,y_train)

y_pred = knn.predict(x_test)

# confusion matrix and f1 score
f1_score_knn = f1_score(y_test,y_pred, average='micro')
cm_knn = confusion_matrix(y_test,y_pred)
sns.heatmap(cm_knn, annot=True,fmt=".0f",linewidths=3,square=True, cmap='Reds', color="#cd1076")
plt.ylabel('actual label')
plt.xlabel('predicted label')
plt.title(f'F1 Score: {f1_score_knn:.2f}',size=14,color='red')
plt.show()
```



```
In [38]: # finding optimum k
error_rate = []
for i in range(1,40):
    knn_test = KNeighborsClassifier(n_neighbors=i, metric = 'minkowski')
    knn_test.fit(x_train,y_train)
    pred_i = knn_test.predict(x_test)
    pred_i=pred_i.reshape(54,1)
    error_rate.append(np.mean(pred_i != y_test))

plt.figure(figsize=(10,6))
plt.plot(range(1,40),error_rate,color='blue', linestyle='dashed',
        marker='o',markerfacecolor='red', markersize=10)
plt.title('Error Rate vs. K Value')
plt.xlabel('K')
plt.ylabel('Error Rate')

error_rate = np.array(error_rate)
a = error_rate.tolist().index(error_rate.min())
print(f'Minimum error: {error_rate.min():.2f} at K: {a}')
```

Minimum error: 0.30 at K: 1

