

Assignment 1

Problem 3

Import necessary libraries

```
In [1]: 1 import pandas as pd
        2 import numpy as np
        3 import matplotlib.pyplot as plt
        4 import seaborn as sns
        5
        6 %matplotlib inline
        7 from plotly import __version__
        8 from plotly.offline import download_plotlyjs, init_notebook_mode, plot, iplot
        9
       10 print(__version__) # requires version >= 1.9.0
       11 import cufflinks as cf
       12 # For Notebooks
       13 init_notebook_mode(connected=True)
       14
       15 # For offline use
       16 cf.go_offline()
       17
       18 import warnings
       19
       20 warnings.filterwarnings("ignore")
```

5.18.0

(a) Read the text file, how many classes exist in the dataset? Which class has the most amount of data?

Read data and plot it based on type of the drink

```
In [2]: 1 df = pd.read_csv('Dataset_III.csv')
```

```
In [3]: 1 df['type'].unique()
```

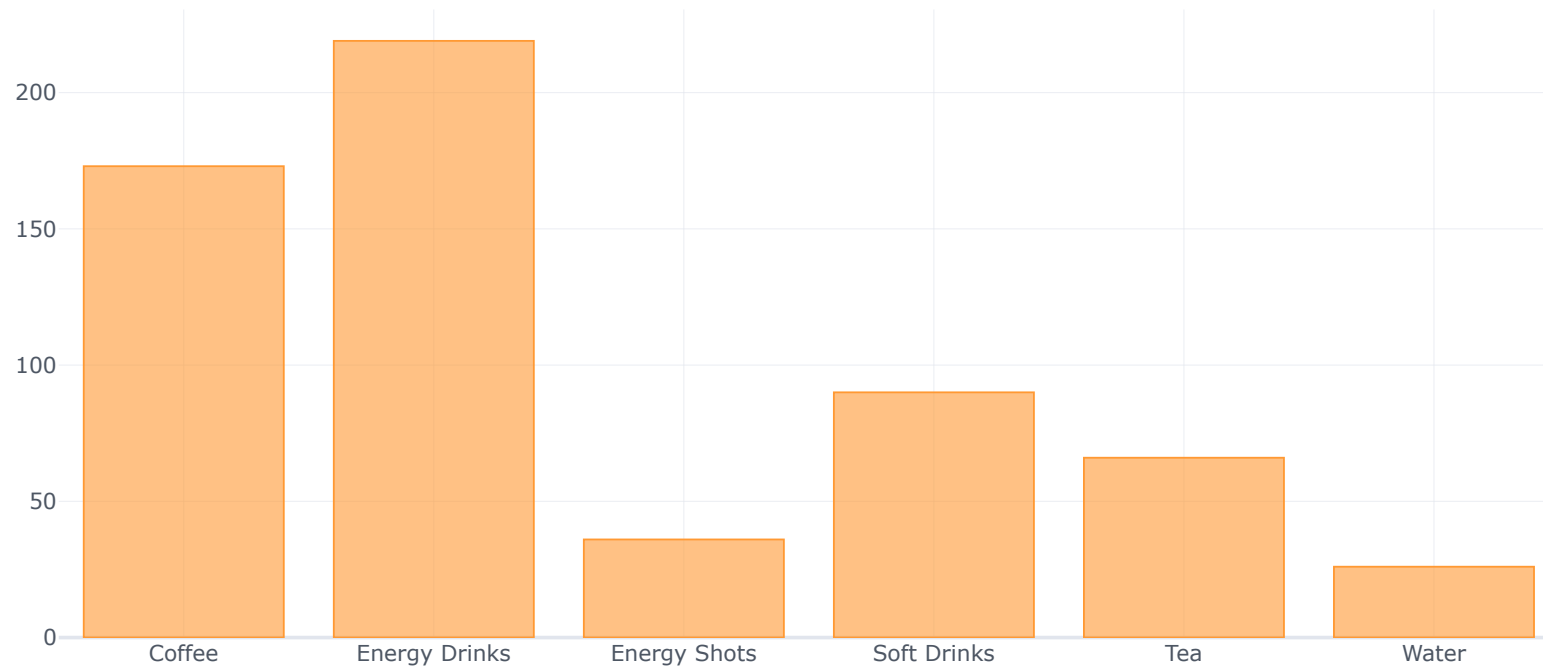
```
Out[3]: array(['Coffee', 'Energy Drinks', 'Energy Shots', 'Soft Drinks', 'Tea',  
              'Water'], dtype=object)
```

```
In [4]: 1 df.groupby('type')['drink'].count()
```

```
Out[4]: type  
Coffee          173  
Energy Drinks   219  
Energy Shots    36  
Soft Drinks     90  
Tea             66  
Water          26  
Name: drink, dtype: int64
```

The distribution of Drinks based on their type is presented in the following plot

```
In [5]: 1 df.groupby('type')['drink'].count().plot(kind = 'bar')
```



[Export to plot.ly »](#)

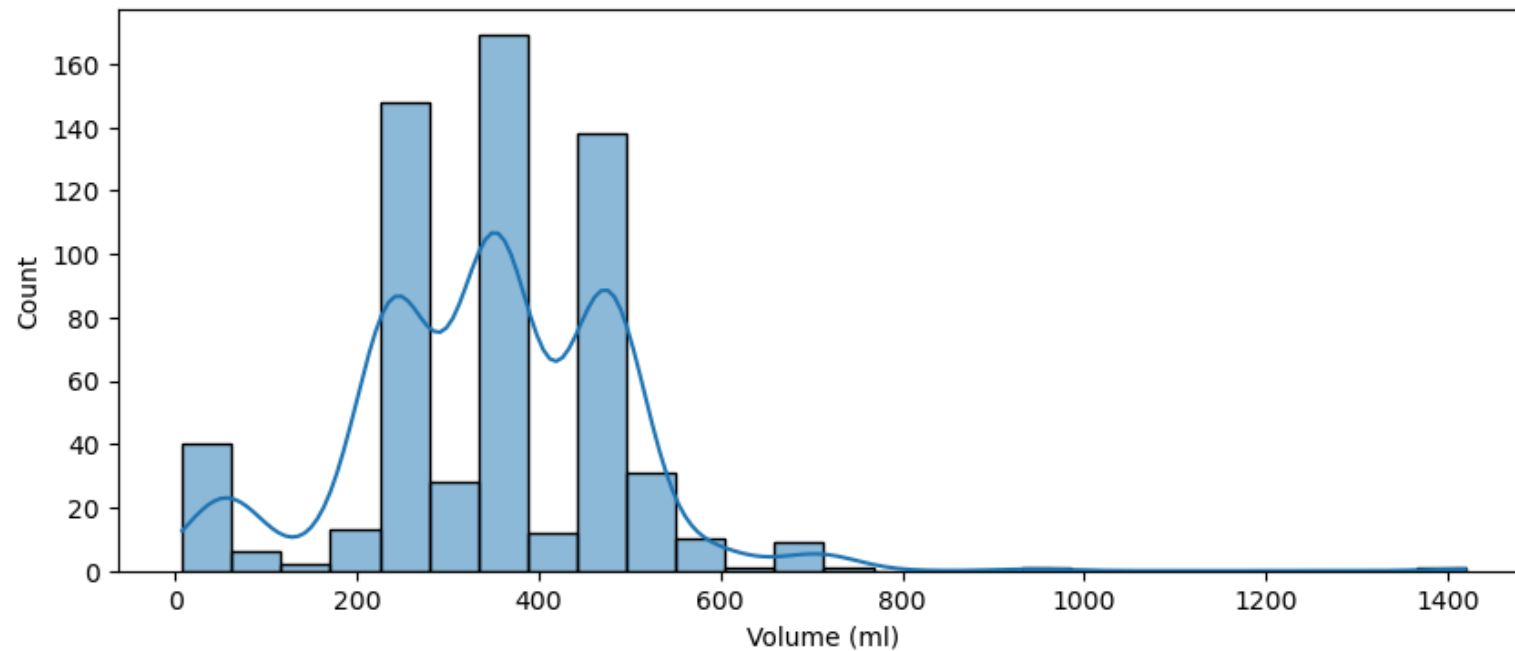
We can see that there are 6 types of drinks among our dataset. Also, this Coffee and Energy Drinks have the most amount of data

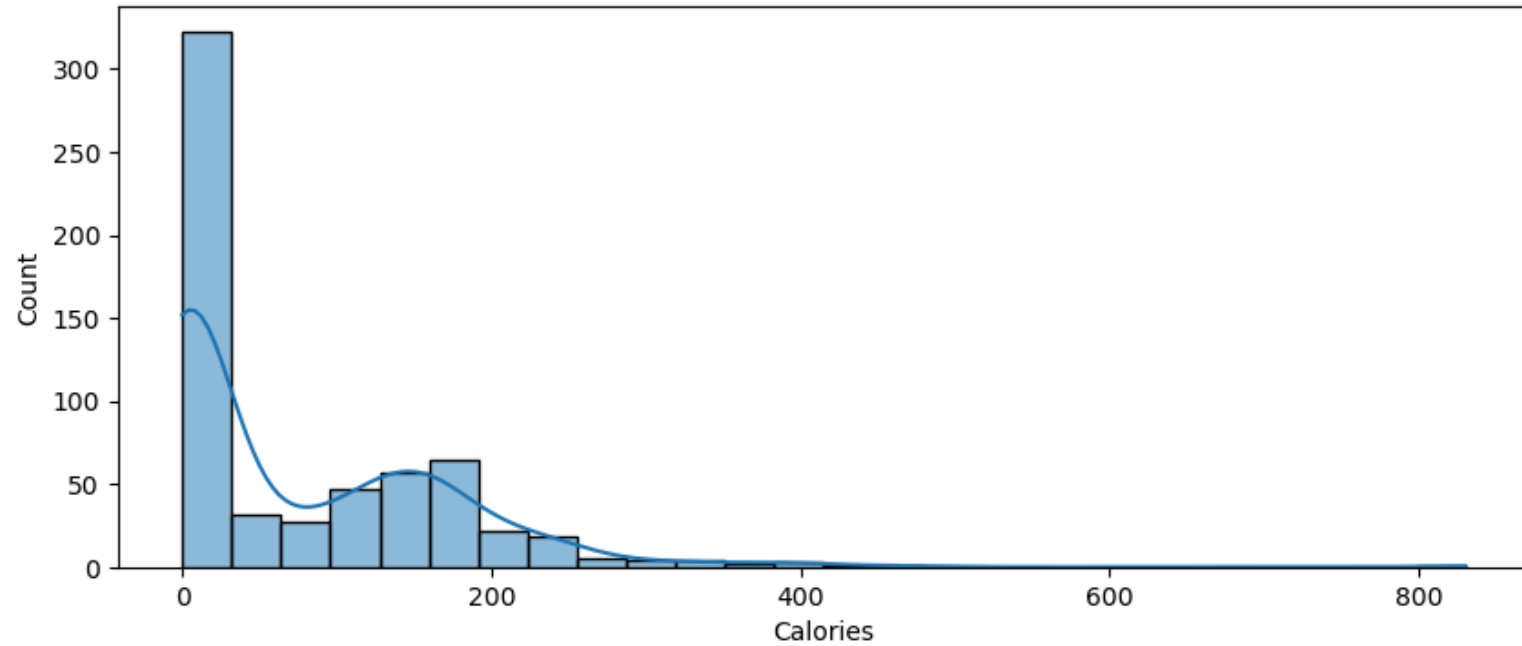
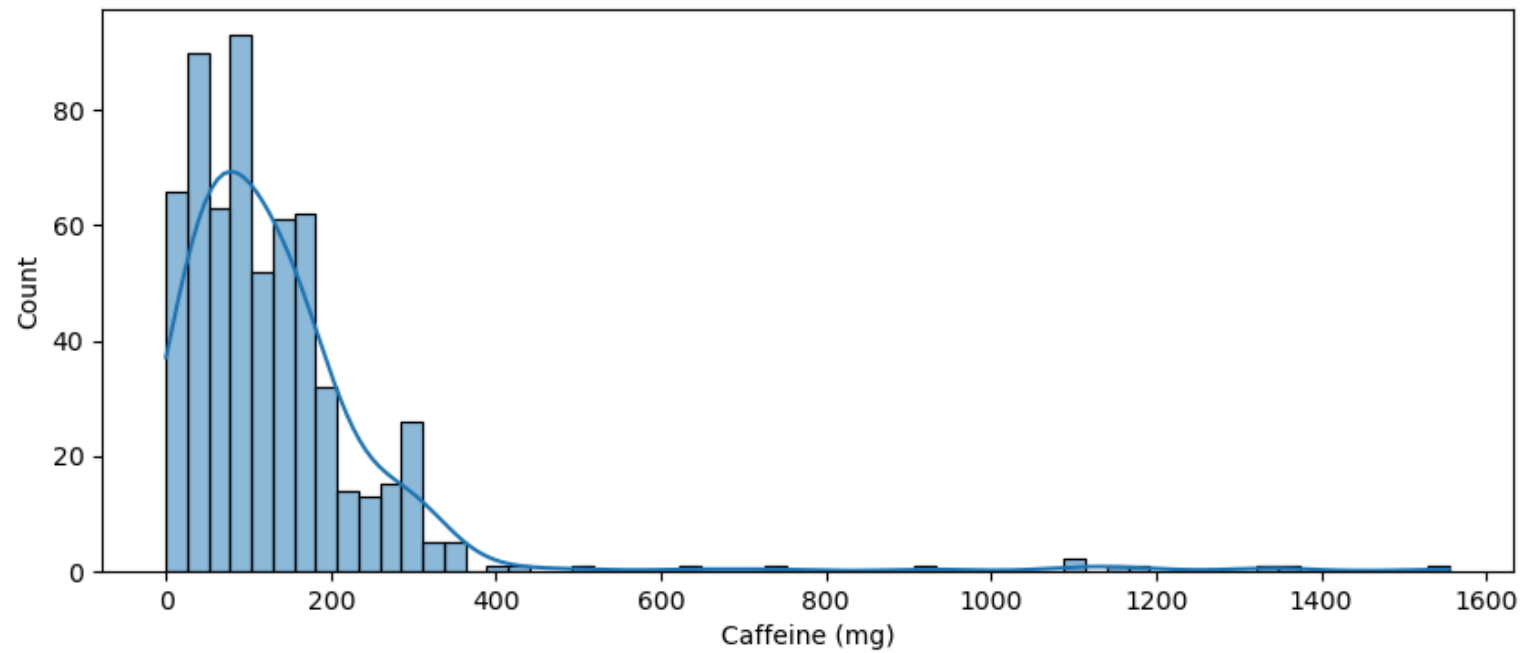
(b) Plot the histogram of the numerical features. Explain what you understand.

Using Seaborn package, we plot histograms of desired features. also, a kde line is plotted above the bar charts to make it more convenient for analysis

In [6]:

```
1 # Plot for 'Volume (ml)'  
2 plt.figure(figsize=(10, 4))  
3 sns.histplot(df['Volume (ml)'], kde=True)  
4 plt.show()  
5  
6 # Plot for 'Caffeine (mg)'  
7 plt.figure(figsize=(10, 4))  
8 sns.histplot(df['Caffeine (mg)'], kde=True)  
9 plt.show()  
10  
11 # Plot for 'Calories'  
12 plt.figure(figsize=(10, 4))  
13 sns.histplot(df['Calories'], kde=True)  
14 plt.show()  
15  
16
```



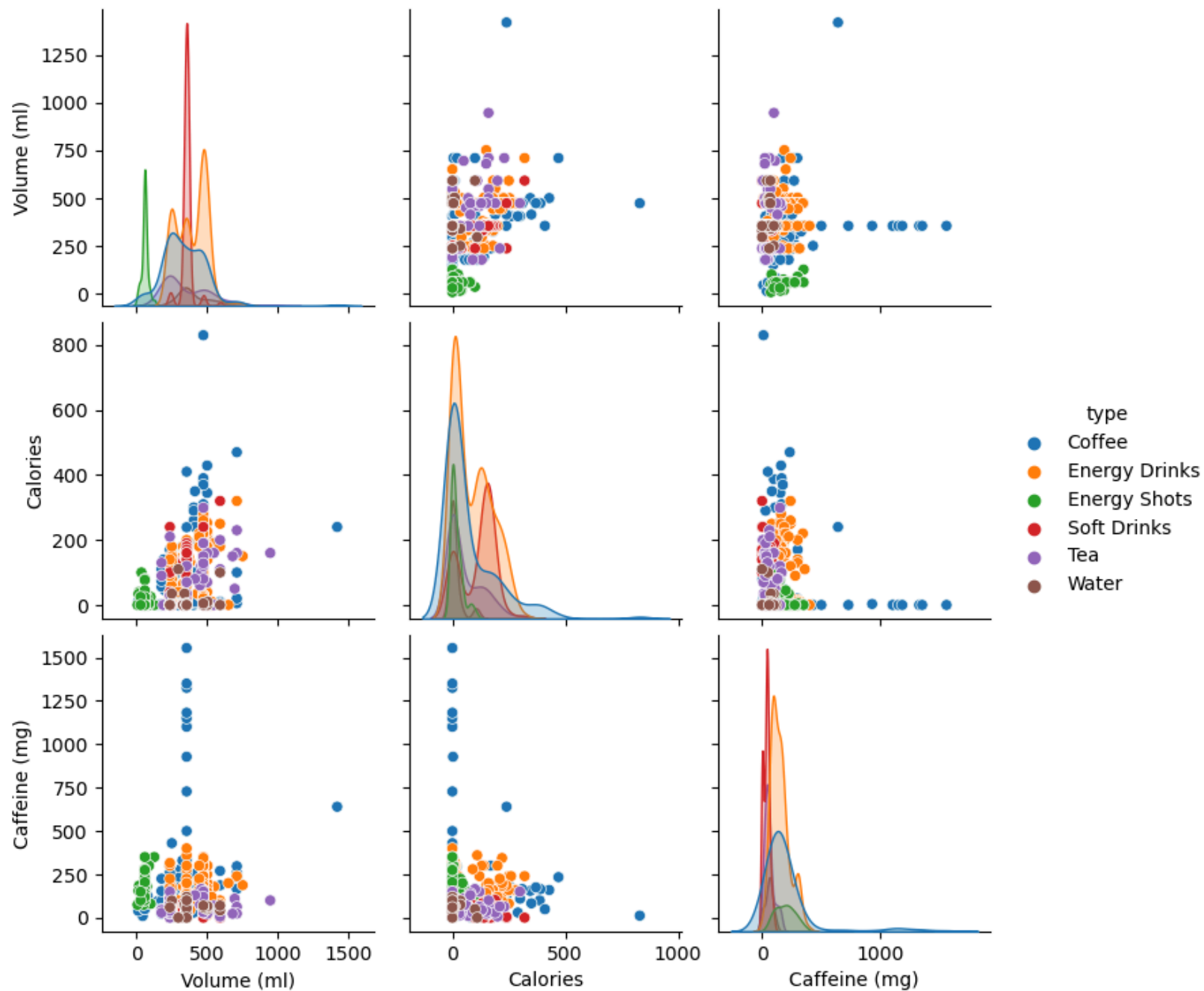


By analyzing above diagrams,the following results can be concluded:

- Most of drinks in data set, have a volume in range of 200-450 ml. this indicates that there are fewer samples related to small sized drinks, like shots and similar cups.
- Also, the distribution of drinks based on Caffeine shows that there are very few drinks with high amount of Caffeine, and most of the drinks have a normal amount of 0-200 mg. Meanwhile, a significant number of drinks have very little amount of Caffeine
- Most of the drinks present in dataset have very little amount of calories. this can be due to the fact in the description of dataset which indicates that the value of calories are not defined, because of the variable amount of sugar which is used besides the drinks

```
In [7]: 1 sns.pairplot(df , hue = 'type')
```

```
Out[7]: <seaborn.axisgrid.PairGrid at 0x17410a0fe20>
```



Scaling feature values to a normalized range

```
In [8]: 1 from sklearn.preprocessing import StandardScaler
        2 scaler = StandardScaler()
        3 scaler.fit(df.drop(['type', 'drink'],axis=1))
```

```
Out[8]: ▾ StandardScaler
        StandardScaler()
```

```
In [9]: 1 scaled_features = scaler.transform(df.drop(['type', 'drink'],axis=1))
        2 df_feat = pd.DataFrame(scaled_features,columns=df.columns[:-2])
        3 df_feat.head()
```

```
Out[9]:
```

	drink	Volume (ml)	Calories
0	-0.623477	-0.797362	0.916714
1	-0.670834	-0.797362	0.066393
2	-0.670834	0.786216	-0.223489
3	-0.670834	-0.797362	1.902314
4	-0.670834	-0.797362	-0.442511

(c) Use KNN to fit the appropriate model. For validation, use the F1 score metric.

Using sklearn, first data is divided into test and train and then a KNN model is trained with initial number of neighbors equals to 1. This value will be modified later

```
In [10]: 1 from sklearn.model_selection import train_test_split
        2 X_train, X_test, y_train, y_test = train_test_split(scaled_features,df['type'], test_size=0.2)
```

```
In [11]: 1 from sklearn.neighbors import KNeighborsClassifier
2 knn = KNeighborsClassifier(n_neighbors=1)
3 knn.fit(X_train,y_train)
```

```
Out[11]: KNeighborsClassifier
KNeighborsClassifier(n_neighbors=1)
```

Using the trained model, output values of test data is predicted. Next, we can see the confusion matrix which shows True and False predictions based on each type finally, some metrics are calculated

```
In [12]: 1 pred = knn.predict(X_test)
```

```
In [13]: 1 from sklearn.metrics import classification_report,confusion_matrix
2 print(confusion_matrix(y_test,pred))
```

```
[[19  6  0  0  5  0]
 [ 7 30  0  0  3  2]
 [ 0  0  5  0  0  0]
 [ 1  3  0 18  2  0]
 [ 4  3  0  0  8  0]
 [ 1  0  0  1  3  1]]
```

```
In [14]: 1 print(classification_report(y_test,pred))
```

	precision	recall	f1-score	support
Coffee	0.59	0.63	0.61	30
Energy Drinks	0.71	0.71	0.71	42
Energy Shots	1.00	1.00	1.00	5
Soft Drinks	0.95	0.75	0.84	24
Tea	0.38	0.53	0.44	15
Water	0.33	0.17	0.22	6
accuracy			0.66	122
macro avg	0.66	0.63	0.64	122
weighted avg	0.68	0.66	0.67	122

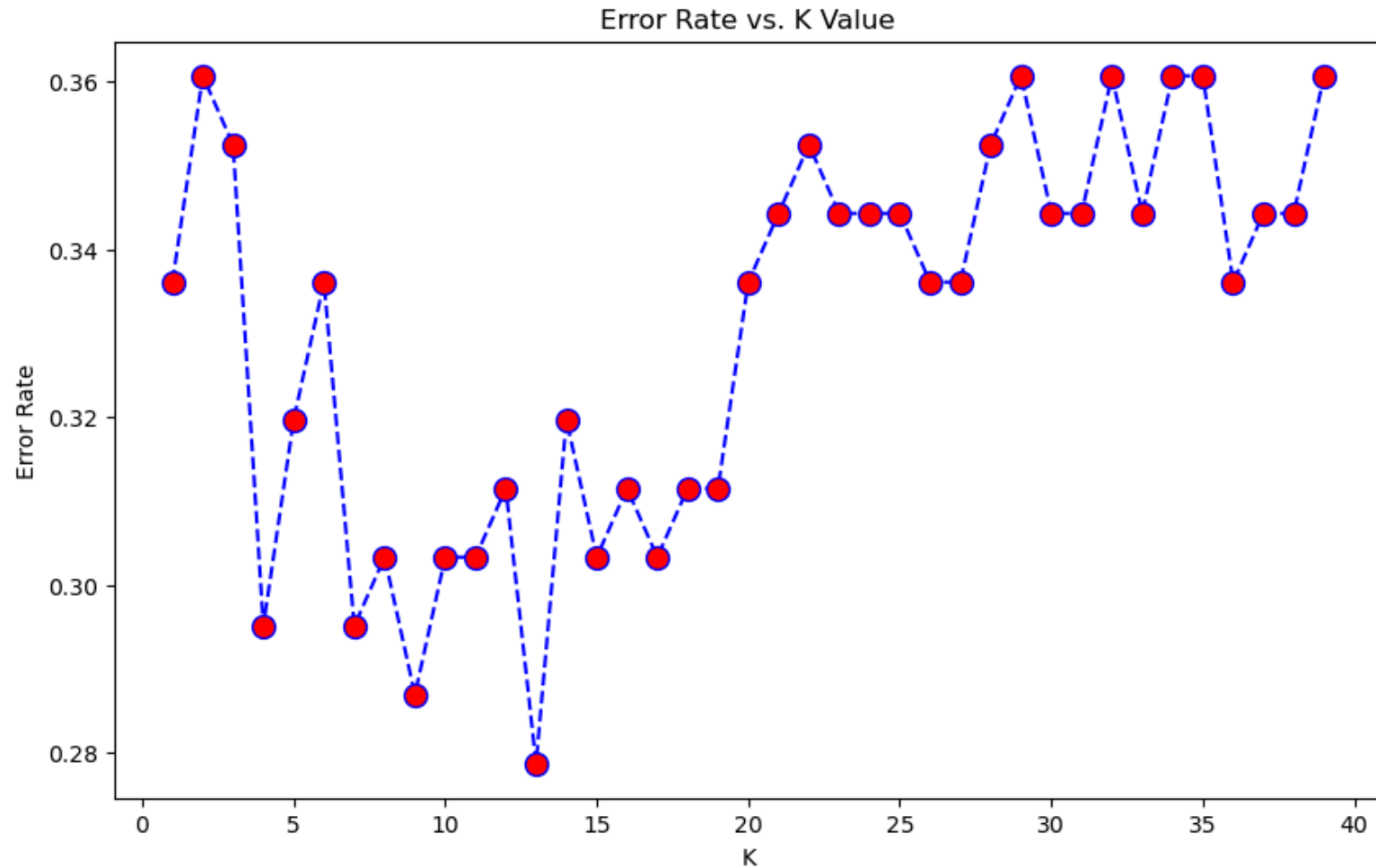
(d) Find the best K.

To optimize the model, we must define the number of neighbors. a convenient method for this is utilizing Elbow Method. in this method, an error rate is calculated for different number of neighbors. later using this plot, we can choose the best value of K

```
In [15]: 1 error_rate = []
          2
          3 # Will take some time
          4 for i in range(1,40):
          5
          6     knn = KNeighborsClassifier(n_neighbors=i)
          7     knn.fit(X_train,y_train)
          8     pred_i = knn.predict(X_test)
          9     error_rate.append(np.mean(pred_i != y_test))
```

```
In [16]: 1 plt.figure(figsize=(10,6))
2 plt.plot(range(1,40),error_rate,color='blue', linestyle='dashed', marker='o', markerfacecolor='red', markersize=10)
3 plt.title('Error Rate vs. K Value')
4 plt.xlabel('K')
5 plt.ylabel('Error Rate')
```

Out[16]: Text(0, 0.5, 'Error Rate')



As shown on the diagram, the minimum Error Rate is 15. so we set number 15 as our number of neighbors

```
In [17]: 1 # NOW WITH K=15
2 knn = KNeighborsClassifier(n_neighbors=15)
3
4 knn.fit(X_train,y_train)
5 pred = knn.predict(X_test)
6
7 print('WITH K=15')
8 print('\n')
9 confusion_matrix(y_test,pred)
10 print('\n')
11 print(classification_report(y_test,pred))
```

WITH K=15

	precision	recall	f1-score	support
Coffee	0.71	0.57	0.63	30
Energy Drinks	0.63	0.86	0.73	42
Energy Shots	1.00	0.80	0.89	5
Soft Drinks	0.87	0.83	0.85	24
Tea	0.62	0.53	0.57	15
Water	0.00	0.00	0.00	6
accuracy			0.70	122
macro avg	0.64	0.60	0.61	122
weighted avg	0.68	0.70	0.68	122

Thank you