

1.

```
1 from cProfile import label
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.decomposition import PCA
4 import seaborn as sns
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import pandas as pd
8 from sklearn.neighbors import KNeighborsClassifier
9 from sklearn.metrics import accuracy_score, confusion_matrix
10
11 """Given the dataset: recipes_muffins_cupcakes_scones.csv,
12 1. print the variance ratio and plot the cumulative sum of the variance ratio for all 8 features. /*DONE*/
13 2. In addition, using as a reference the plots in slides 110-115, create a scatter plot between PC1 and PC2. /*DONE*/
14 3. In addition, create a histogram of features and plot a heatmap with the features with the largest variation in PC1 and PC2. /*DONE*/
15 4. Furthermore, find and print the features with the highest (max) and lowest (min) variation both in PC1 and PC2
16 (positively and negatively correlated). /*DONE*/
17 5. Finally, plot a correlation heatmap. /*DONE*/
18 """
19 Note: You will need to standardize your data points.
20 """
21
22 df = pd.read_csv('recipes_muffins_cupcakes_scones.csv')
23
24 df.loc[df['Type'] == 'Muffin', 'Type'] = 0
25 df.loc[df['Type'] == 'Cupcake', 'Type'] = 1
26 df.loc[df['Type'] == 'Scone', 'Type'] = 2
27
28 X = np.array(df.loc[:, 'Flour':'Salt'])
29 y = np.array(df.loc[:, 'Type'])
30
31
32 Xscaler = StandardScaler()
33 X = Xscaler.fit_transform(X)
34
35
36 pca = PCA(n_components=8)
37 PC = pca.fit_transform(X)
38 explained_variance = pca.explained_variance_ratio_
39 print(f'\nVariance Ratio: {explained_variance}')
40 print()
41
42 #cumulative sum
43 x = np.arange(1,9)
44 plt.plot(x, np.cumsum(explained_variance))
45 plt.xlabel('Principal Components')
46 plt.ylabel('Cumulative Ratio')
47 plt.title('PC= 1-8')
48 plt.xticks(range(1,11))
49 plt.show()
50
51 # scatter plot
52 plt.scatter(PC[:, 0], PC[:, 1], c=y, cmap='plasma')
53 plt.xlabel('PC1')
54 plt.ylabel('PC2')
55 plt.title(f'PC= 2, Variance: {explained_variance[0:1]}')
56 plt.show()
57
58 # histogram
59 Muffin = X[y==0]
60 Cupcake = X[y==1]
61 Scone = X[y==2]
62
63 fig, axes = plt.subplots(2, 4, figsize=(12, 9)) # 3 columns each containing 10 figures, total 30 features
64 ax = axes.ravel() # flat axes with numpy ravel
65 feature_names = [*df]
66 feature_names.pop(0)
```

```

68     for i in range(8):
69         _, bins = np.histogram(X[:, i], bins=25)
70         ax[i].hist(Muffin[:, i], bins=bins, color='r', alpha=.5) # red color for malignant class
71         ax[i].hist(Cupcake[:, i], bins=bins, color='g', alpha=.3) # alpha is for transparency in the overlapped region
72         ax[i].hist(Scone[:, i], bins=bins, color='b', alpha=.2) # alpha is for transparency in the overlapped region
73         ax[i].set_title(feature_names[i], fontsize=9)
74         ax[i].axes.get_xaxis().set_visible(False) # the x-axis co-ordinates are not so useful, as we just want to look how well separated the histograms are
75         ax[i].set_yticks(())
76
77     ax[0].legend(['Muffin', 'Cupcake', 'Scone'], loc='best', fontsize=8)
78     plt.tight_layout() # let's make good plots
79     plt.show()
80
81     #highest (max) and lowest (min) variation
82
83     PCAcomp = pca.components_
84     PCAcomp = PCAcomp[0:2,:]
85
86
87     print(f"PC1 highest variation feature {feature_names[int(np.where(PCAcomp == max(PCAcomp[0,:]))[1][0])]}: {max(PCAcomp[0,:])} ")
88     print(f"PC1 lowest variation feature {feature_names[int(np.where(PCAcomp == min(PCAcomp[0,:]))[1][0])]}: {min(PCAcomp[0,:])} ")
89
90     print(f"PC2 highest variation feature {feature_names[int(np.where(PCAcomp == max(PCAcomp[1,:]))[1][0])]}: {max(PCAcomp[1,:])} ")
91     print(f"PC2 lowest variation feature {feature_names[int(np.where(PCAcomp == min(PCAcomp[1,:]))[1][0])]}: {min(PCAcomp[1,:])} ")
92
93

```

```

94
95     #heatmap
96     df2 = df.drop('Type', axis='columns')
97     s = sns.heatmap(df2.corr(), cmap='coolwarm')
98     s.set_yticklabels(s.get_yticklabels(), rotation=30, fontsize=7)
99     s.set_xticklabels(s.get_xticklabels(), rotation=30, fontsize=7)
100     plt.show()
101
102     #KNN
103     pca = PCA(n_components=2)
104     X_t = pca.fit_transform(X)
105     y = y.astype('int')
106
107     knn = KNeighborsClassifier(n_neighbors=3)
108     knn.fit(X_t, y)
109
110     data = np.array([38, 18, 23, 20, 9, 3, 1, 0]).reshape(1,-1)
111     data = Xscaler.transform(data)
112     data = pca.transform(data)
113
114
115
116     pred = knn.predict(data)
117
118     prediction = ''
119     if pred == 0:
120         prediction = 'Muffin'
121     elif pred == 1:
122         prediction = 'Cupcake'
123     else:
124         prediction = 'Scone'
125
126     print()
127     print('the data is predicted to be a ', prediction)
128

```

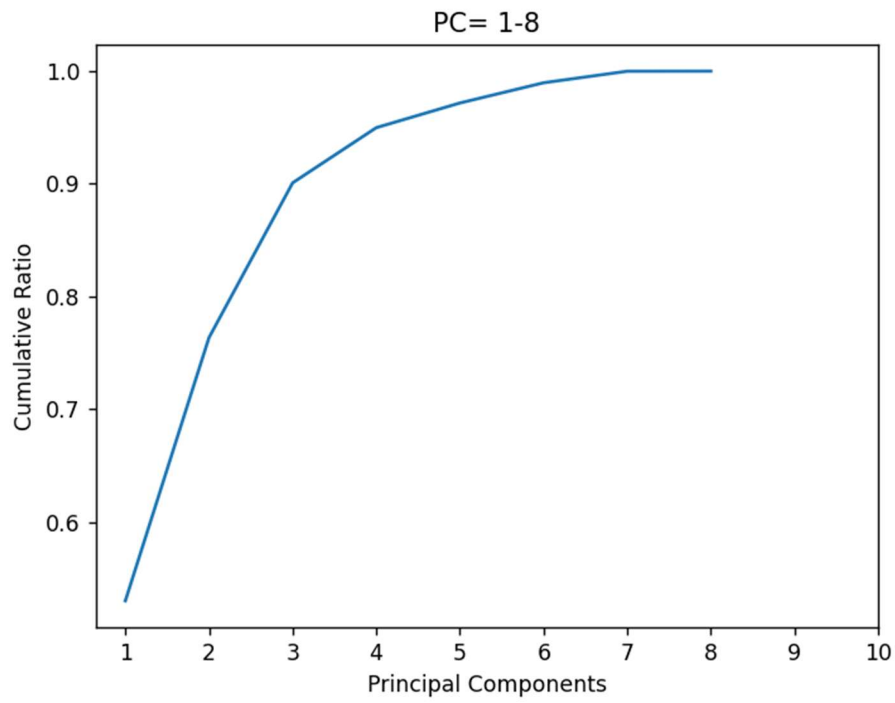
```

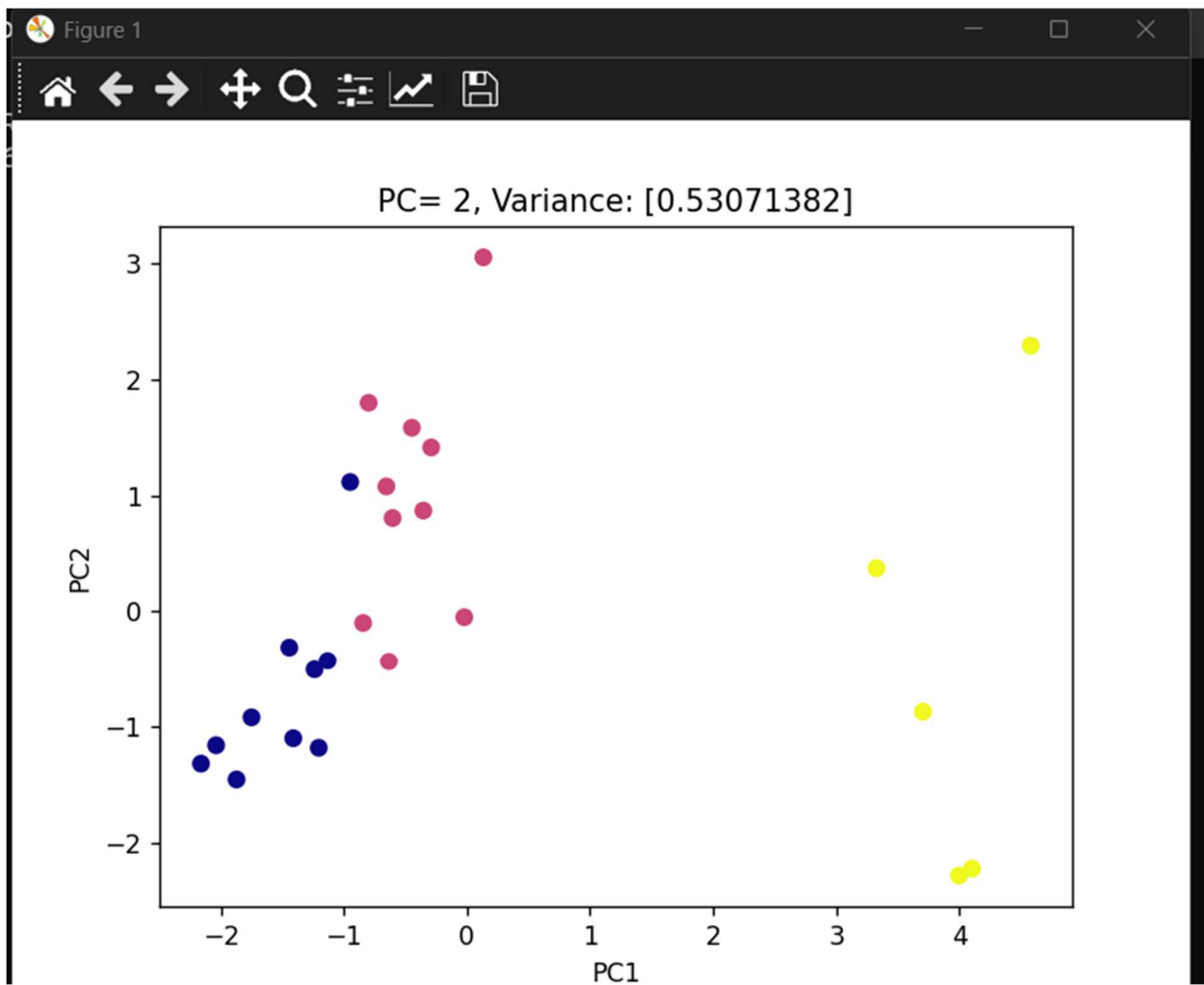
128
129     plt.scatter(X_t[:, 0], X_t[:, 1], c=y, cmap='plasma')
130     plt.scatter(data[0,0], data[0,1], c='r', cmap='plasma')
131     plt.xlabel('PC1')
132     plt.ylabel('PC2')
133     plt.title(f'PC= 2, Variance: {explained_variance[0:1]}')
134
135     # Displaying the legend
136     plt.show()
137

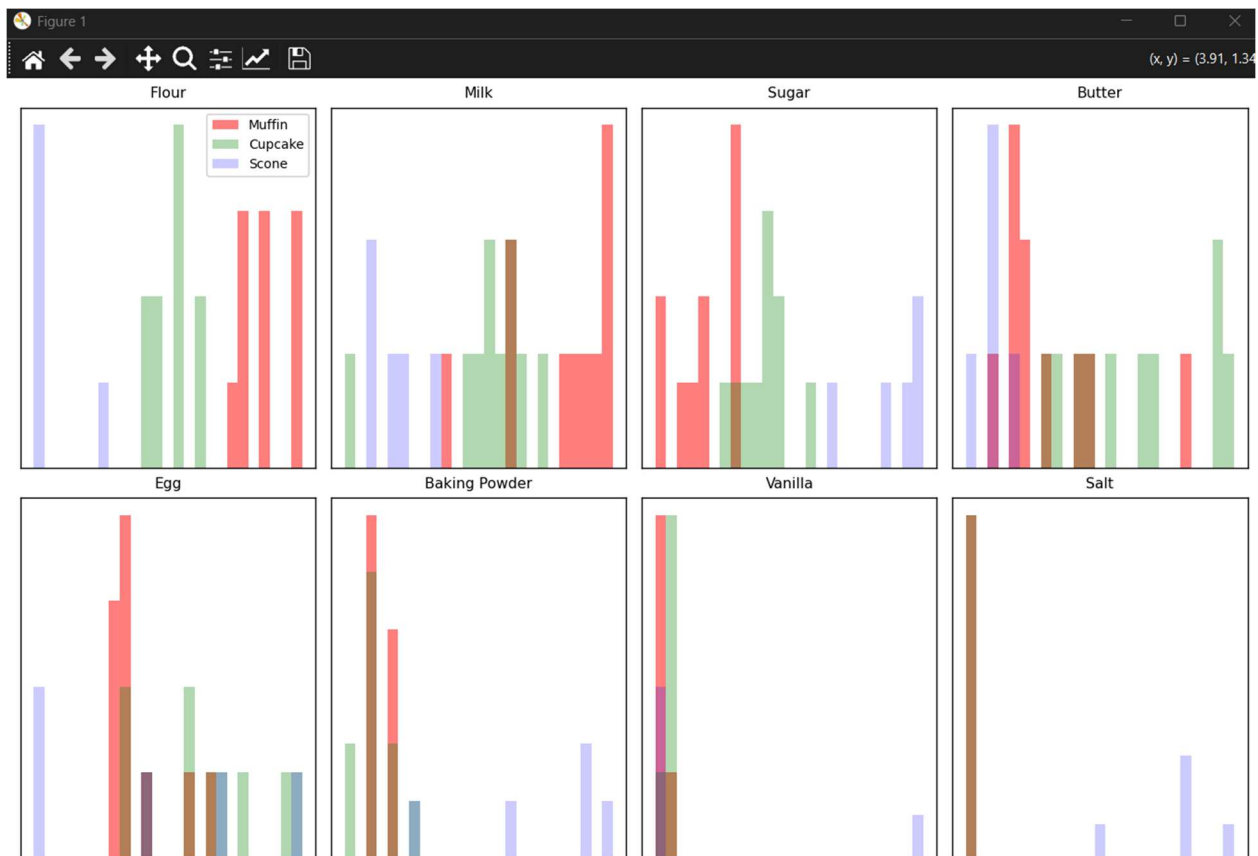
```

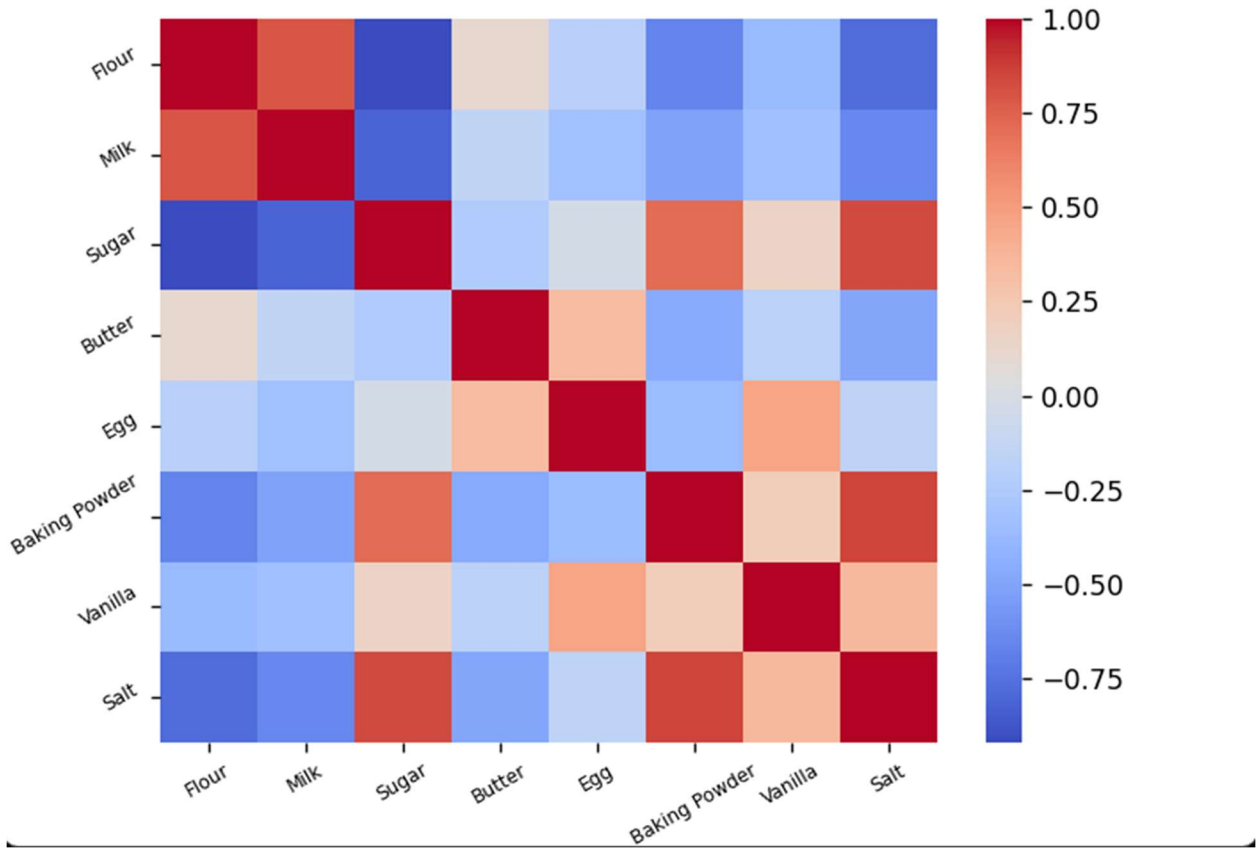
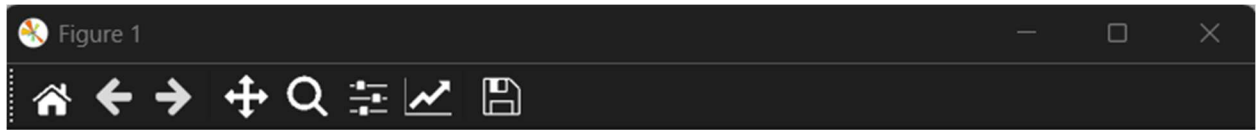
Variance Ratio: [5.30713818e-01 2.32926199e-01 1.37271639e-01 4.89399359e-02
2.18813363e-02 1.79333825e-02 1.02410755e-02 9.26139804e-05]

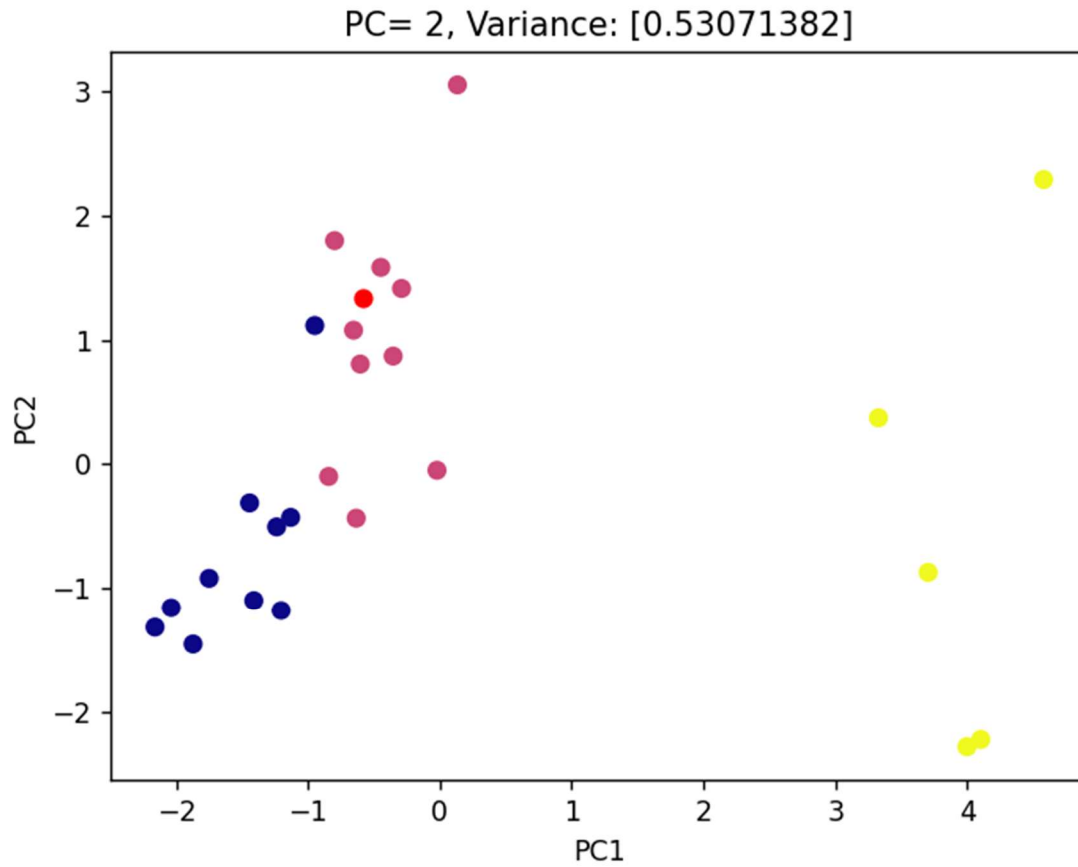
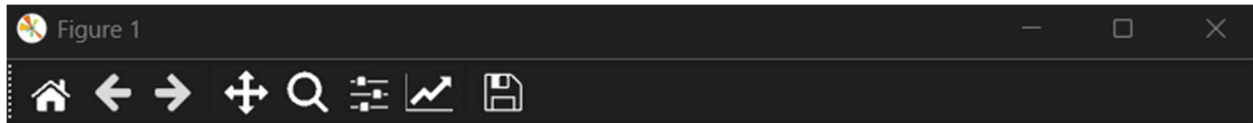
Figure 1











Variance Ratio: [5.30713818e-01 2.32926199e-01 1.37271639e-01 4.89399359e-02 2.18813363e-02 1.79333825e-02 1.02410755e-02 9.26139804e-05]

PC1 highest variation feature Salt: 0.45560967418642023

PC1 lowest variation feature Flour: -0.4452797934996038

PC2 highest variation feature Egg: 0.6628137748507001

PC2 lowest variation feature Milk: -0.3233860340000608

the data is predicted to be a Cupcake