

**Question 4. (30 points)** Implement K-Means clustering algorithm for the dataset (**buddymove.csv**) provided with this assignment. You can run the algorithm three times for the values of  $K=3$ ,  $K=4$ , and  $K=5$  for each run (here  $K$  refers to # of clusters). A sample code file (**KMeans Clustering.py**) has been attached with this assignment. The description of data set is given below:

1. **Data Set Information:**

This dataset was populated from destination reviews published by 249 reviewers of holidayiq.com till October 2014. Reviews falling in 6 categories among destinations across South India were considered and the count of reviews in each category for every reviewer (traveler) is captured. More information about this dataset is available at the following link:

<https://archive.ics.uci.edu/ml/datasets/BuddyMove+Data+Set>

2. **Attribute Information:**

Attribute 1 : Unique user id

Attribute 2 : Number of reviews on stadiums, sports complex, etc.

Attribute 3 : Number of reviews on religious institutions

Attribute 4 : Number of reviews on beach, lake, river, etc.

Attribute 5 : Number of reviews on theatres, exhibitions, etc.

Attribute 6 : Number of reviews on malls, shopping places, etc.

Attribute 7 : Number of reviews on parks, picnic spots, etc.

You might only need to consider Attributes 2 through 7, because User\_id can be dropped. The idea is to cluster the Users that have similar characteristics of posting a # of reviews. You can implement the algorithm in any programming language you want. However, I highly recommend you use Python, R, or Matlab. Your program for question 4 should have the following:

1. Attach a snapshot of each run of your program (for different values of  $K$ ) that prints clustering labels for each run.

```
Windows PowerShell
Copyright (C) Microsoft Corporation. All rights reserved.

Try the new cross-platform PowerShell https://aka.ms/pscore6

PS C:\Users\Jacob\Desktop\hw4> conda activate base
PS C:\Users\Jacob\Desktop\hw4> & C:/Users/Jacob/anaconda3/python.exe "c:/Users/Jacob/Desktop/hw4/KMean Clustering.py"
(249, 7)
  User Id  Sports  Religious  Nature  Theatre  Shopping  Picnic
0  User 1      2        77        79        69        68        95
1  User 2      2        62        76        76        69        68
2  User 3      2        50        97        87        50        75
3  User 4      2        68        77        95        76        61
4  User 5      2        98        54        59        95        86

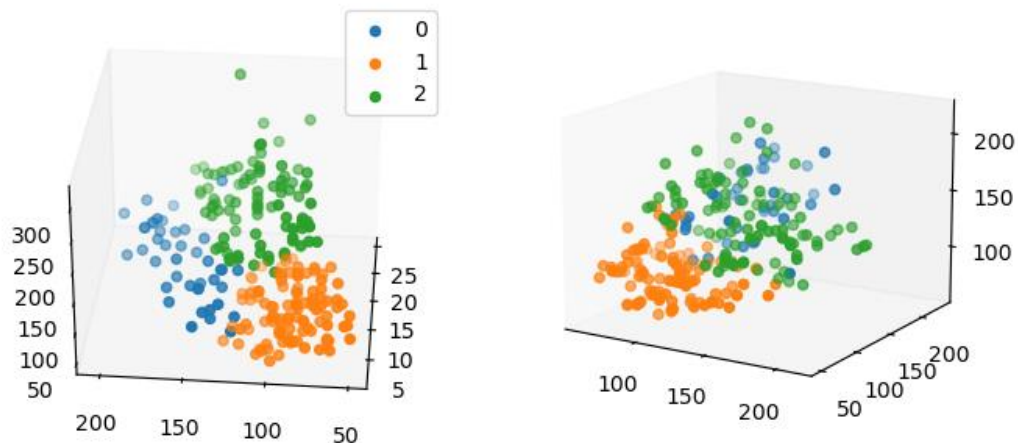
Clustering buddymove.csv with K-Means(k=3)
```

[illegible]

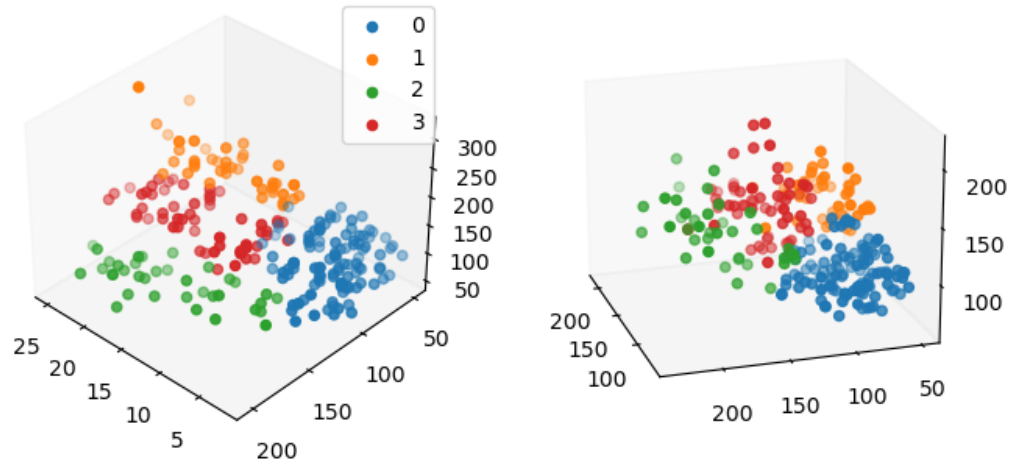
```
Clustering buddymove.csv with K-Means(k=5)
Unique labels:
[0 1 2 3 4]
Labels:
[4 4 4 4 2 4 4 4 4 4 2 2 4 2 4 4 4 4 2 4 2 4 4 4 4 2 4 4 4 4 2 4 4 4 2 4 2 4 2 4 4
 4 4 2 4 4 4 4 4 4 2 2 4 4 3 4 4 4 4 2 4 4 2 4 2 4 4 4 2 4 4 4 2 4 4 4 2 4 2 4 4 4
 4 2 4 2 4 2 4 2 2 4 4 4 4 2 4 4 2 4 4 4 4 2 2 4 4 2 4 2 4 4 4 2 4 4 4 2 4 1 1 1
 1 0 1 0 1 3 0 1 1 1 0 0 1 4 1 1 0 4 3 1 1 0 4 0 1 3 1 1 4 1 4 1 3 4 3 0 4
 2 0 0 0 0 1 1 4 1 0 3 1 3 1 3 0 0 3 1 0 3 1 0 1 1 0 1 3 0 1 0 3 0 1 1 3 0
 3 0 0 0 1 3 1 0 0 3 0 1 1 0 3 1 1 3 0 3 1 1 0 1 3 1 0 3 3 3 0 1 3 0 0 1 1
 1 1 3 0 1 3 1 1 1 3 1 0 1 0 3 4 1 1 3 0 0 0 1 0 1 1 1]
Centers:
[[ 17.6744186  96.86046512 195.72093023 132.58139535  90.30232558
 147.44186047]
 [ 17.06779661 127.25423729 140.94915254 144.74576271 134.69491525
 138.05084746]
 [  6.09677419 116.          73.48387097  82.61290323 116.93548387
 106.09677419]
 [ 17.3         169.4         85.63333333 112.83333333 192.4
 143.16666667]
 [  5.93023256 81.20930233 109.60465116 102.22093023  79.30232558
 91.98837209]]
```

2. Your code should also print visualization of each cluster. Research on how to create a visualization for clusters.

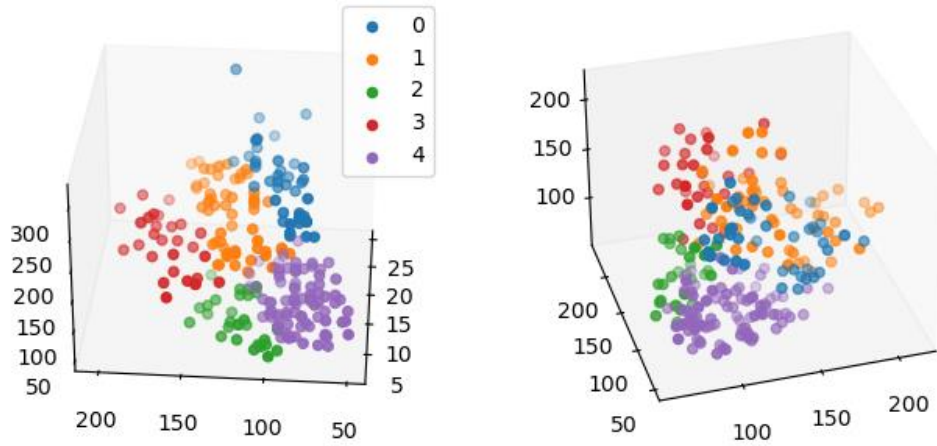
Kmeans(k=3) Clustering of buddymove.csv



Kmeans(k=4) Clustering of buddymove.csv



Kmeans(k=5) Clustering of buddymove.csv



3. Submit the code file along with your submission.

```
4. import numpy as np
5. import pandas as pd
6. import matplotlib.pyplot as plt
7. from mpl_toolkits import mplot3d
8. from sklearn.cluster import KMeans
9.
10.##@Author: Jacob Hopkins
11.## 3d plot refernce https://likegeeks.com/3d-plotting-in-python/#Putting\_legends
12.
13.## read from csv file
14.dataset_filename = 'buddymove.csv'
15.data = pd.read_csv(dataset_filename)
16.
17.##The following code displays the shape of data set
18.print(data.shape)
19.print(data.head())
20.
21.#names of columns
22.f1label = 'Sports'
23.f2label = 'Religious'
24.f3label = 'Nature'
25.f4label = 'Theatre'
26.f5label = 'Shopping'
27.f6label = 'Picnic'
28.
29.##Extracting each column
30.f1 = data[f1label]
31.f2 = data[f2label]
32.f3 = data[f3label]
33.f4 = data[f4label]
34.f5 = data[f5label]
35.f6 = data[f6label]
36.
37.## Creating an array of data points
38.X = np.array(list(zip(f1,f2,f3,f4,f5,f6)))
39.
40.## K-Means clustering algorithm with different parameters, 3-5 clusters
41.
42.for k in range(3,6):
43.    print(f'\n\nClustering {dataset_filename} with K-Means(k={k})')
44.
45.    kmeans = KMeans(n_clusters=k, random_state=0).fit(X)
46.
```

```
47.     lables = kmeans.labels_  
48.     unique_labels = np.unique(lables)  
49.  
50.     ## print cluster labels  
51.     print("Unique labels: ")  
52.     print(unique_labels)  
53.  
54.     print("Labels: ")  
55.     print(lables)  
56.  
57.     print("Centers: ")  
58.     print(kmeans.cluster_centers_)  
59.  
60.     ## visualize the clusters  
61.     fig = plt.figure(figsize=(8,4))  
62.  
63.     ## plot of features 0 1 2  
64.     ax = fig.add_subplot(121, projection='3d')  
65.  
66.     for l in unique_labels:  
67.         ax.scatter(X[lables == l , 0] , X[lables == l , 1] , X[lables == l  
68. , 2], label = l)  
69.     ax.set_title(f'Kmeans(k={k}) Clustering of {dataset_filename}')  
70.     ax.grid(False)  
71.     ax.legend(loc="best")  
72.  
73.     ## plot of features 3, 4, 5  
74.     ax2 = fig.add_subplot(122, projection='3d')  
75.  
76.     for l in unique_labels:  
77.         ax2.scatter(X[lables == l , 3] , X[lables == l , 4] , X[lables ==  
78. l , 5], label = l)  
79.     ax2.grid(False)  
80.     plt.show()
```