

# INF552: Programming Assignment 3

Yi-Li Chen

## Part 1: Implementation

### PCA

- All imported library

```
1. import csv
2. import numpy as np
```

- Read text file into a list and convert to array

```
1. data = []
2. with open('pca-data.txt') as file:
3.     reader = csv.reader(file, delimiter='\t')
4.     for line in reader:
5.         data.append([float(line[0]),float(line[1]),float(line[2])])
6. data = np.array(data)
```

- Step0: Centralize the data point by deducting mean value

```
1. def data_centralize(data):
2.     mu = np.mean(data,axis=0)
3.     data_std = data-mu
4.     return data_std
```

- Step1: Calculate the covariance matrix for the features in the dataset.

```
1. def calculate_covariance(data_std):
2.     data_covar = np.cov(data_std, rowvar=False, bias=True)
3.     return data_covar
```

- Step2: Calculate the eigenvalues and eigenvectors for the covariance matrix.

```
1. def calculate_eigen(data_covar):
2.     eigenvalues,eigenvectors = np.linalg.eig(data_covar)
3.     return eigenvalues,eigenvectors
```

- Step3: Pick k eigenvalues and form a matrix of eigenvectors.

```
1. def build_k_dim_eigenvectorMatrix(eigenvectors, k):
2.     k_eigen = eigenvectors[:,0:k]
3.     return k_eigen
```

- After defining each step of algorithm, do the matrix multiplication of original data and matrix of eigenvectors with Nxk dimension.

```

1. def PCA(data, k):
2.     data_std = data_centralize(data)
3.     data_covar = calculate_covariance(data_std)
4.     eigenvalues,eigenvectors = calculate_eigen(data_covar)
5.     k_eigen = build_k_dim_eigenvectorMatrix(eigenvectors, k)
6.     print(k_eigen)
7.     new_data = np.matmul(data_std, k_eigen)
8.     return new_data

```

- The final result of PCA:

```

1. new_data = PCA(data, k=2)

```

The directions of the first two principal components:

	1 <sup>st</sup>	2 <sup>nd</sup>	
[ [	0.86667137	-0.4962773	]
[	-0.23276482	-0.4924792	]
[	0.44124968	0.71496368	]]

## FastMap

- All imported library

```

1. import csv
2. import random
3. import numpy as np
4. import matplotlib.pyplot as plt

```

- Read the fastmap\_data into a list and convert to array

```

1. fastmap_data = []
2. with open('fastmap-data.txt') as file:
3.     reader = csv.reader(file, delimiter='\t')
4.     for line in reader:
5.         fastmap_data.append([float(line[0]),float(line[1]),float(line[2])])
6. fastmap_data = np.array(fastmap_data)

```

- Read the fastmap\_wordlist into a list

```
1. fastmap_wordlist = []
2. with open('fastmap-wordlist.txt') as file:
3.     reader = csv.reader(file)
4.     for line in reader:
5.         fastmap_wordlist+=line
```

- Build a function that we can find out the corresponding farthest object when the one object and the distance are given

```
1. def find_farthest_obj(obj1_idx, fastmap_data):
2.     # The row related to obj1_idx
3.     obj1_distances = np.concatenate((fastmap_data[fastmap_data[:,0]==obj1_idx],
4.                                       fastmap_data[fastmap_data[:,1]==obj1_idx]))
5.     # Find out Oa in the above rows with max distance
6.     max_dist_idx = np.argmax(obj1_distances[:,2])
7.     obj2 = list(obj1_distances[max_dist_idx][0:2])
8.     obj2.remove(obj1_idx)
9.     obj2_idx = int(obj2[0])
10.    return obj2_idx
```

- Build a function that we can heuristically find out the pair of farthest objects:

- (1) Initialize an object arbitrarily and declare it to be Ob
- (2) Set Oa to be the object that is farthest apart from Ob according to the above function
- (3) Set Ob to be the object that is farthest apart from Oa again
- (4) Return the pair if Ob are same as the previous iteration

```
1. def find_farthest_pair(fastmap_wordlist, fastmap_data):
2.     Ob = random.sample(fastmap_wordlist,1)[0]
3.     Ob_idx = fastmap_wordlist.index(Ob)+1
4.     old_ob_idx = Ob_idx
5.     while 1:
6.         Oa_idx = find_farthest_obj(Ob_idx, fastmap_data)
7.         Ob_idx = find_farthest_obj(Oa_idx, fastmap_data)
8.         if old_ob_idx == Ob_idx:
9.             return Oa_idx, Ob_idx
10.        else:
11.            old_ob_idx = Ob_idx
```

- Build a distance function so that we can gain the distance when two objects and fastmap\_data are given. This is because we can calculate the coordinate of object easily and more clearly.

```

1. def distance_function(obj1_idx, obj2_idx, fastmap_data):
2.     idx = []
3.     idx.append(obj1_idx)
4.     idx.append(obj2_idx)
5.     idx.sort()
6.     for idx_objs in range(len(fastmap_data)):
7.         if fastmap_data[:,0][idx_objs]==idx[0] and fastmap_data[:,1][idx_objs]==idx[1]:
8.             d = fastmap_data[:,2][idx_objs]
9.             return d
10.        else:
11.            d = 0
12.    return d

```

- In sum, firstly, initialize the Ob randomly, and find out the farthest pair of the fastmap data. Secondly, with the farthest pair of Oa and Ob, we can calculate the coordinate for each object. Then, update the fastmap with new distance function. Afterwards, we call FastMap function again with the new fastmap data and k-1. In addition, a stopping criteria which is when k=0 need to be set so as to return the and finish calculating all the coordinate for each object. One thing needs to be noticed is that the k will be subtracted 1 and the coordinate\_idx will be added 1 for every next iteration.

```

1. def FastMap(k, fastmap_data, fastmap_wordlist, coordinate_idx):
2.     # Stopping criteria
3.     if k<=0:
4.         return
5.     else:
6.         coordinate_idx+=1
7.         # Find each farthest pair of fastmap_data
8.         Oa_idx, Ob_idx = find_farthest_pair(fastmap_wordlist, fastmap_data)
9.         # Calculate the coordinate for each iteration
10.        for obj_idx in range(1,len(fastmap_wordlist)+1):
11.            d_ai = distance_function(Oa_idx, obj_idx, fastmap_data)
12.            d_ab = distance_function(Oa_idx, Ob_idx, fastmap_data)
13.            d_ib = distance_function(Ob_idx, obj_idx, fastmap_data)
14.            x = (d_ai**2+d_ab**2-d_ib**2)/(2*d_ab)
15.            output_k_coordinate[obj_idx-1][coordinate_idx]=x

```

```

16.     # Create new distance function which is new fastmap data
17.     new_fastmap_data = fastmap_data.copy()
18.     for row in range(len(new_fastmap_data)):
19.         D_old = distance_function(new_fastmap_data[row][0],
20.                                   new_fastmap_data[row][1],
21.                                   fastmap_data)
22.         xi = output_k_coordinate[int(new_fastmap_data[row][0]-1)][coordinate_idx]
23.         xj = output_k_coordinate[int(new_fastmap_data[row][1]-1)][coordinate_idx]
24.         new_fastmap_data[row][2] = (D_old**2-(xi-xj)**2)**0.5
25.     # Call FastMap function again with the input of new fastmap
26.     FastMap(k-1, new_fastmap_data, fastmap_wordlist, coordinate_idx)
27.     return output_k_coordinate

```

- The final result of FastMap algorithm:

```

1. output_k_coordinate = np.zeros((len(fastmap_wordlist),2)) # k=2
2. result = FastMap(2, fastmap_data, fastmap_wordlist, -1)

```

#### → The result of output coordinate

```

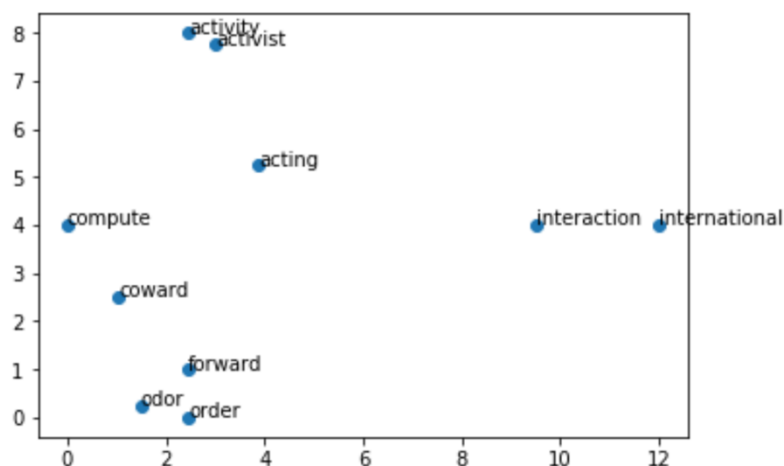
array([[ 3.875      ,  5.25      ],
       [ 3.        ,  7.75      ],
       [ 0.        ,  4.        ],
       [ 1.04166667,  2.5       ],
       [ 2.45833333,  1.        ],
       [ 9.5       ,  4.        ],
       [ 2.45833333,  8.        ],
       [ 1.5       ,  0.25      ],
       [ 2.45833333,  0.        ],
       [12.        ,  4.        ]])

```

```

1. x = result[:,0]
2. y = result[:,1]
3. plt.plot(x,y, 'o')
4. for xx,yy,word in zip(x,y,fastmap_wordlist):
5.     plt.annotate(word,(xx,yy))

```



In this assignment, the most challenging problem is to build the recursive algorithm for FastMap function. Since we have to make it to be a linear, we can't use two for-loop to run all the objects. We have to use one for loop to calculate the new distance function and the new fastmap data. Moreover, setting the stopping criteria is always the hardest things for me. I have to know clearly when it should stop and return the result. Last but not least, the relationship between the index of objects in fastmap-wordlist and the IDs in the fastmap-data need to be careful to deal with.