**Week 9**

CS 273a - Introduction to Machine Learning (Winter '15)*

Prof. Alex Ihler

**Varad Meru**

Student # 26648958

Due Date: 03/10/2015

# Homework 5[†‡]

## Problem 1: Basics of Clustering

(a) The dataset is loaded in Matlab, as shown in Listing 1 and Figure 1

Listing 1: Loading Data

```matlab
% Problem a
iris=load('data/iris.txt');      % load the text file
X = iris(:,1:2);          % features are other columns
features = char('Sepal length','Sepal width','Petal length','Petal ...
    width','Species');
features_short = char('SL','SW','PL','PW','SP');
whos

f = figure;
scatter(X(:,1), X(:,2), 'filled');
saveas(f,'scatter.png','png');
```

(b) The Listing 2 shows the code listing for problem of k-means for different values of k and initialization methods. The Figure 2 shows the plot of the created clusters.

Listing 2: K-Means on iris dataset for different values of k and initialization methods.

```matlab
%% Problem b
k = 5;
[z,c,sumd] = kmeans(X,k);
[z1,c1,sumd1] = kmeans(X,k,'k++');

f = figure;
plotClassify2D([],X,z);
saveas(f,'kmeans_k_5_simple.png', 'png');

f = figure;
plotClassify2D([],X,z1);
saveas(f,'kmeans_k_5_kpp.png', 'png');

k = 20;
[z,c,sumd] = kmeans(X,k);
[z1,c1,sumd1] = kmeans(X,k,'k++');

f = figure;
plotClassify2D([],X,z);
saveas(f,'kmeans_k_20_simple.png', 'png');

f = figure;
plotClassify2D([],X,z1);
saveas(f,'kmeans_k_20_kpp.png', 'png');
```

(c) The Listing 3 shows the code listing for problem of k-means for different values of k and initialization methods. The Figure 3 shows the plot of the created clusters.

---

Listing 3: Agglomerative clustering on iris dataset for different values of k and linkage methods.

```
1  %% Problem c
2  k = 5;
3  Z = linkage(X,'single');
4  c = cluster(Z,'maxclust',k);
5  f = figure;
6  plotClassify2D([],X,c);
7  saveas(f,'linkage_single_5.png', 'png');
8  f = figure;
9  dendrogram(Z)
10 saveas(f,'dendogram_5.png', 'png');
11
12 Z = linkage(X,'complete');
13 c = cluster(Z,'maxclust',k);
14 f = figure;
15 plotClassify2D([],X,c);
16 saveas(f,'linkage_complete_5.png', 'png');
17
18 k = 20;
19 Z = linkage(X,'single');
20 c = cluster(Z,'maxclust',k);
21 f = figure;
22 plotClassify2D([],X,c);
23 saveas(f,'linkage_single_20.png', 'png');
24 f = figure;
25 dendrogram(Z)
26 saveas(f,'dendogram_20.png', 'png');
27
28 Z = linkage(X,'complete');
29 c = cluster(Z,'maxclust',k);
30 f = figure;
31 plotClassify2D([],X,c);
32 saveas(f,'linkage_complete_20.png', 'png');
```

(d) The EM Gaussian mixture model is run with 5 and 20 components, as shown in Listing 4. The generated clusters can be seen in Figure 4.

Listing 4: The EM Gaussian mixture model for 5 and 20 components.

```
1  %% Problem d
2  k = 5;
3  [zx,Tx,softx,llx] = emCluster(X,k);
4  f = figure;
5  plotClassify2D([],X,zx);
6  saveas(f,'emgm_5.png', 'png');
7
8  k = 20;
9  [zx,Tx,softx,llx] = emCluster(X,k);
10 f = figure;
11 plotClassify2D([],X,zx);
12 saveas(f,'emgm_20.png', 'png');
```

All the three clustering mechanisms bring different properties and have different usecases. The Agglomerative clustering has a simple architecture and works well for simple and small dataset. EM Gaussian mixture models are well suited for fuzzy nature of dataset where a data might not be strictly associated with any cluster. Traditional means as a low computational complexity compared to others, $O(nkt)$, where $n$ is the number of datapoints, $k$ is the number of clusters and $t$ is the number of iterations. And thus, it can be computed
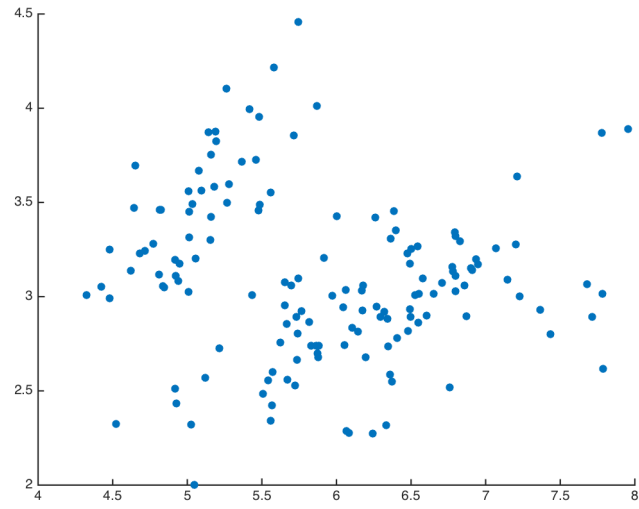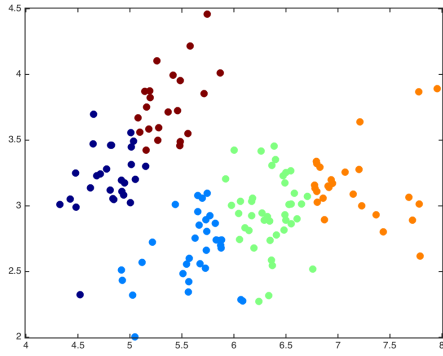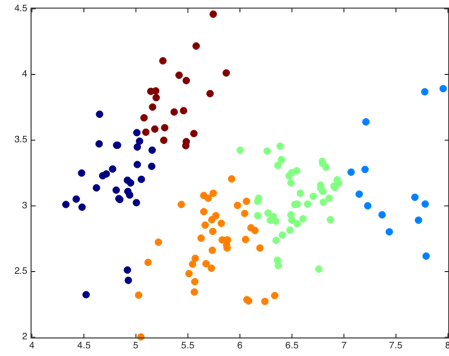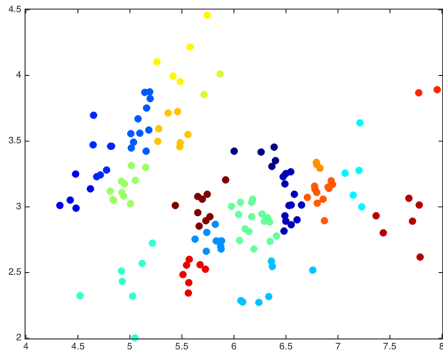
Figure 1: Data Loaded

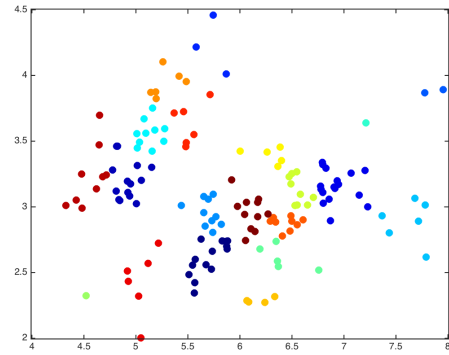repeatedly to search for a good configuration. With all the methods, 'k-means ++' gave good results.

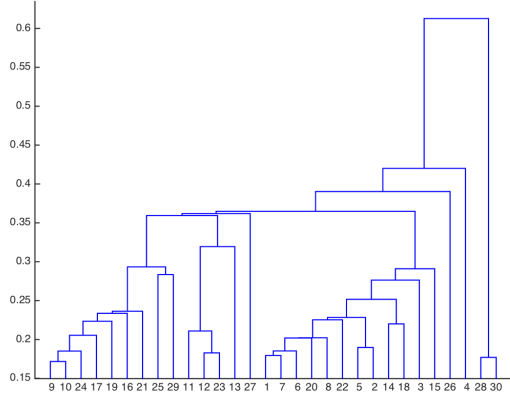(a) k=5 Random Initialisation

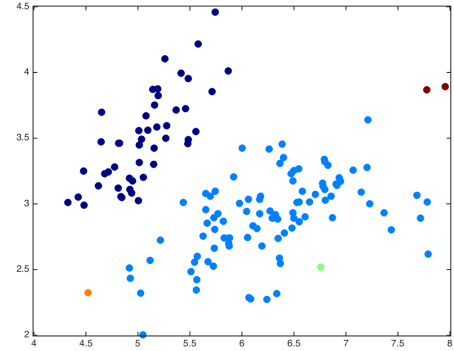(b) k=5 Initialisation using K-Means++

(c) k=20 Random Initialisation
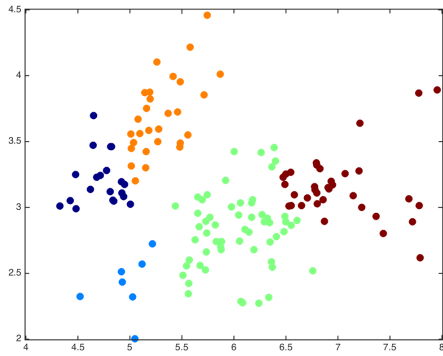
(d) k=20 Initialisation using K-Means++

Figure 2: Various runs of K-Means different values of k and initialization methods
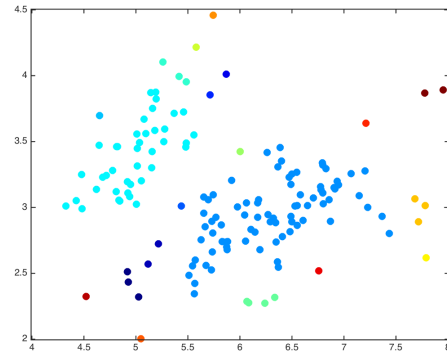
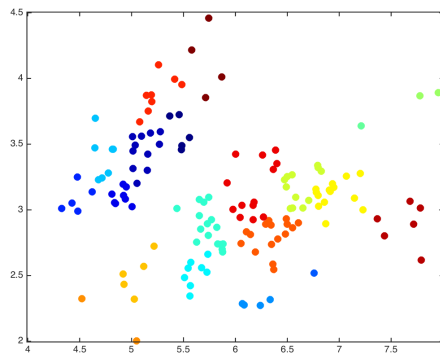(a) Complete Dendrogram of the created clusters

(b) k=5 Linkage method = single

(c) k=5 Linkage method = complete

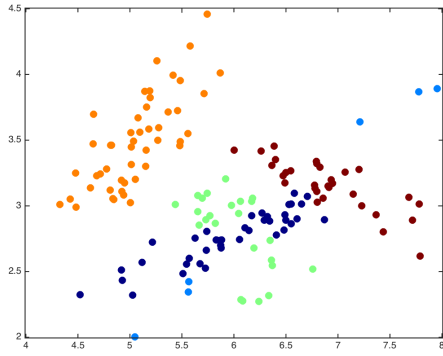(d) k=20 Linkage method = single
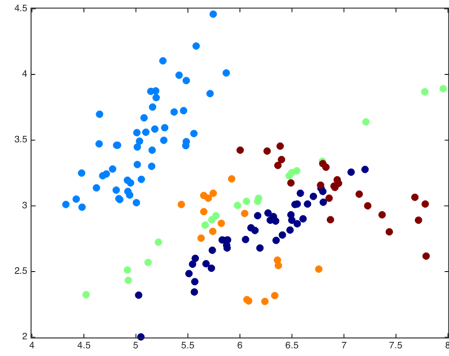
(e) k=20 Linkage method = complete
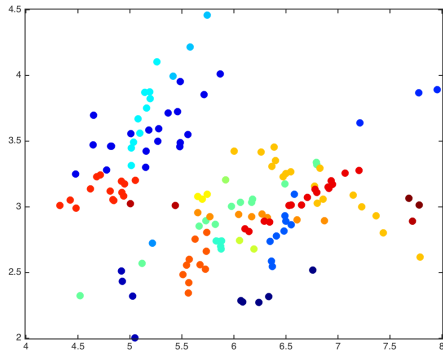
Figure 3: Various runs of Agglomerative clustering with different values of k and linkage methods
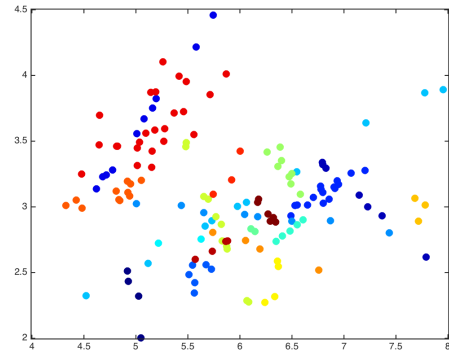
(a) k=5 Random Initialisation

(b) k=5 Initialisation using K-Means++

(c) k=20 Initialisation using K-Means++

(d) k=20 Initialisation using K-Means++

Figure 4: EM Gaussian mixture model with different values of k and initialization methods

## Problem 2: K Means Clustering on Text

(a) The clusters are computed using `kmeans()` method provided by Matlab and the code that produces can be seen in Listing 5.

Listing 5: K-means on textual data with `k=20` and for part (b), the number of runs $= 20$

```
1   %% Problem (a, b)
2   k = 20;
3   [z,c,sumd] = kmeans(Xn,k);
4   disp(sumd)
5   disp('-------------')
6   for i=1:20;
7       [z1,c1,sumd1] = kmeans(Xn,k);
8       disp(sumd1)
9       if sumd1 < sumd
10          z = z1;
11          c = c1;
12          sumd = sumd1;
13      end;
14  end;
15  display('Minimum-sumd')
16  disp(sumd)
17  %{
18      2.4211
19  --------------
20      2.0789
21      2.1009
22      2.0453
23      2.0625
24      2.0861
25      2.4494
26      2.0615
27      2.0709
28      2.0646
29      2.1144
30      2.0098
31      2.0530
32      2.4806
33      2.0479
34      2.0587
35      2.0811
36      2.4695
37      1.9865
38      2.4891
39      2.0573
40
41  Minimum-sumd
42      1.9865
43  %}
```

(b) Listing 5 shows 20 runs of K-Means on the textual data and I got the best value of the cost function to be `sumd = 1.9865`.

(c) For counting the document accusation with clusters, I use a function called `count_unique()`, and it gives the output as given in Listing 6 as well as the bar graph of the distribution can be seen at Figure 5. For the second part of the question, I used a slight modification of the original code given in the homework to populate the document terms. The code and the output can be seen in Listing 7.

Listing 6: The number of documents for each cluster

```
 1  >> [uniques,numUnique] = count_unique(z)
 2  >> [uniques,numUnique]
 3  ans =
 4        1    21
 5        2    45
 6        3     4
 7        4     2
 8        5     7
 9        6     8
10        7    69
11        8     3
12        9    10
13       10     2
14       11     2
15       12     1
16       13     1
17       14     3
18       15     1
19       16     1
20       17    15
21       18     2
22       19     3
23       20     2
```

Listing 7: The "Most-Likely" terms in the 20 clusters

```
 1  %% Problem (c) - 2
 2  for i =1:size(c, 1);
 3      [sorted,order] = sort( c(i,:), 2, 'descend');
 4      fprintf('Doc %d: ',i);
 5      fprintf('%s ', vocab{order(1:10)});
 6      fprintf('\n');
 7  end;
 8  %{
 9  Doc 1: times square millennium city 2000 night 000 eve york midnight
10  Doc 2: team game season coach games players league play going win
11  Doc 3: archbishop york bishop cardinal sports church began column ...
        american close
12  Doc 4: america boat team zealand cup nippon gilmour challengers round true
13  Doc 5: book century marks amp war week finds lives school boy
14  Doc 6: yeltsin putin russia russian president power political kremlin ...
        chechnya russians
15  Doc 7: city american national president 000 home millennium end political ...
        going
16  Doc 8: fireworks island city midnight celebration lot millennium hour ...
        calls celebrations
17  Doc 9: y2k koskinen system problems saturday 2000 reported computers ...
        friday officials
18  Doc 10: tutsi hutu rwanda burundi ethnic country experts africa van 1994
19  Doc 11: texas arkansas yards line offensive game season sacks games defensive
20  Doc 12: test end houston 000 0101 0102 100 1900 1900s 1968
21  Doc 13: cats beijing owners police association called carry chinese eat ...
        eating
22  Doc 14: hijackers hostages pakistan burger told government indian india ...
        passengers killed
23  Doc 15: sports angeles began brooklyn column los seen young eye game
24  Doc 16: economy government putin system america businesses country ...
        economic president russia
25  Doc 17: 2000 computer internet government systems york problem problems ...
        news city
26  Doc 18: lakers jackson game star phil players record conference practice ...
        monday
```
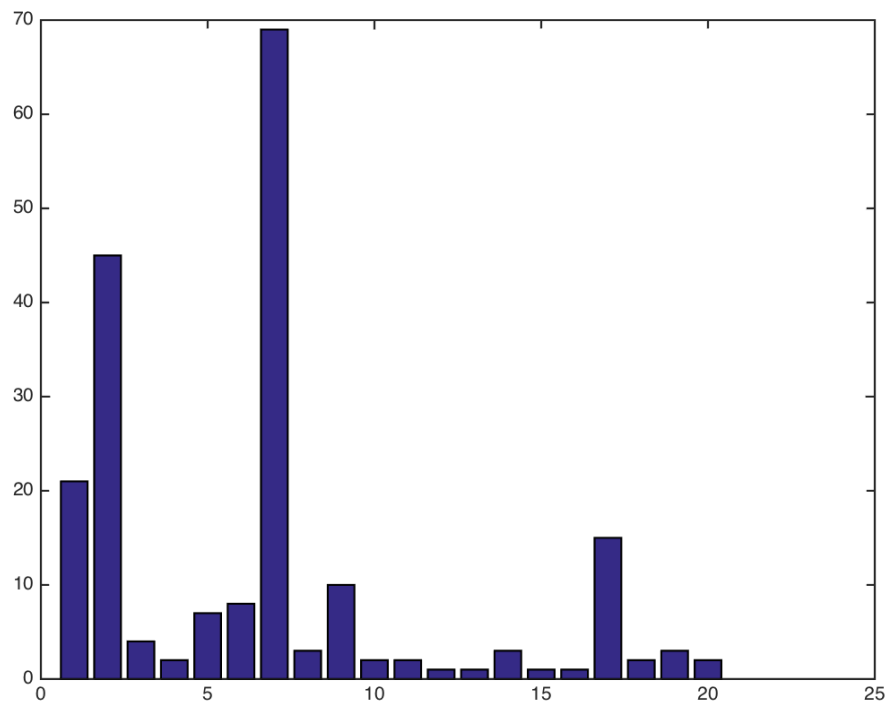
Figure 5: Distribution of the documents spread across k=20 clusters.

```
27  Doc 19: news atlanta constitution journal service moved cox y2k cnn ...
       millennium
28  Doc 20: buses authority diesel natural gas plan mta city york hybrid
29  %}
```

(d) A

# Problem 3: Eigen Faces

(a) Loading Eigen Faces Data, as can be seen in Listing 9. Also see Figure 6

Listing 8: Loading Face Dataset

```
1  clc;close all;clear all;
2  rand('seed',0);
3  X = load('data/faces.txt'); % load face dataset
```

(b) Normalisng the dataset.

Listing 9: Normalizing

```
1  %% Faces A
2  mu = mean(X);
3  X0 = bsxfun(minus,X,mu);size(mu);ans =1 576
```
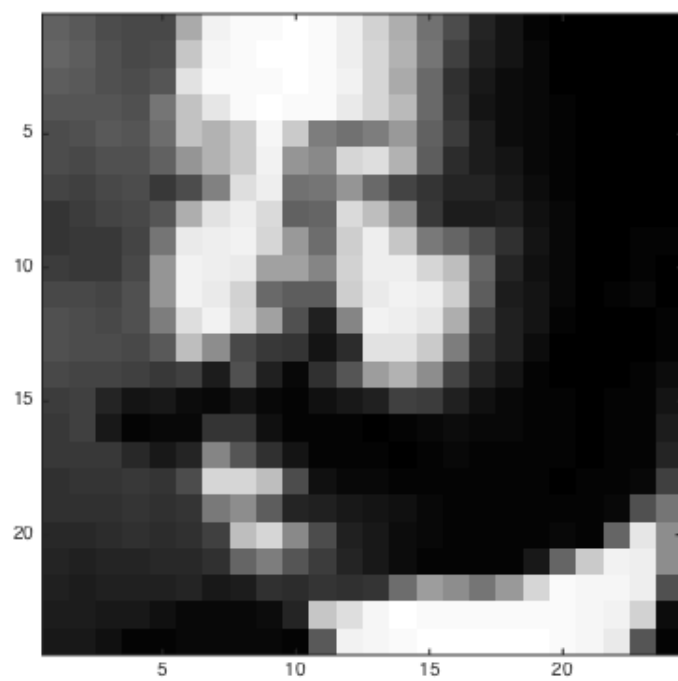
Figure 6: Face sample.