

EE583 Pattern Recognition HW4

Kutay Uğurlu 2232841

December 5, 2021

1 Question 1

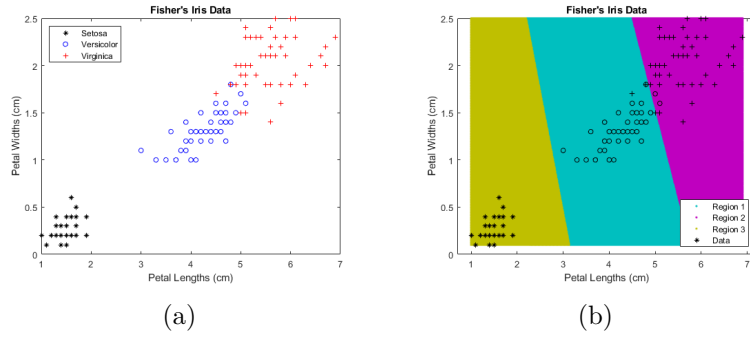


Figure 1: The data distribution and partitioning

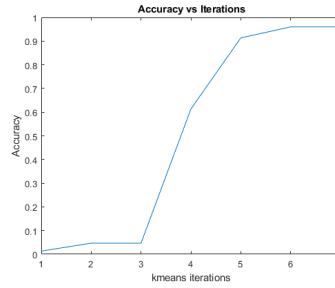


Figure 2: classification success for different centroids

The centroids are manually set according to the approximate mean deduced from the data distribution plot. Initial centroids can be found in 5.1. The accuracy is calculated in two steps. First, the confusion matrix is calculated, then the diagonal entries of it are added to obtain the true positive classifications. Finally, $Accuracy \triangleq \frac{\#TP}{\#Samples}$. Figure 2 shows the accuracies of the iterations initialized with different centroids. Through careful centroid selection, the accuracy can be increased from 0.1 to 0.9647.

2 Question 2

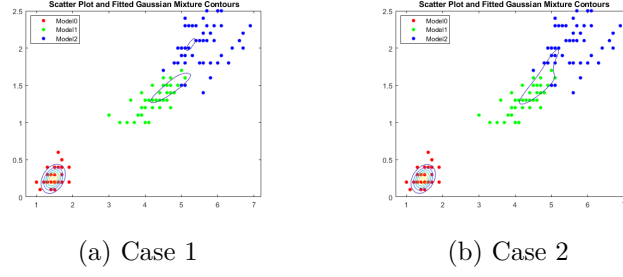


Figure 3: Gaussian Fits

I have included 2 different versions of Gaussian fit to the feature vectors, due to the randomness of the algorithm. In Figure 3a, I have observed 3 different Gaussian fits. However, the second trial resulted in algorithm fitting 2 different features to the same Gaussian distribution in Figure 3b.

3 Question 3

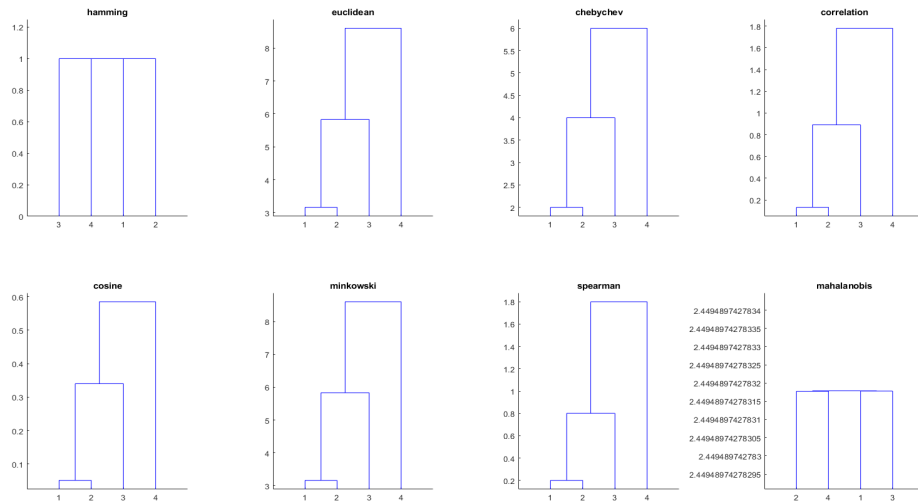


Figure 4: Dendrograms for different metrics

All the metrics beside Hamming and Mahalanobis, resulted in similar dendrograms with easily noticeable difference in the distances in the relevant scales. However, distance calculated using Mahalanobis resulted in close numbers, hence even in order of 10^{-8} the dendrogram presentation is not clear. Hamming distance calculate the percentage of different coordinates in the data matrix X . Since all rows of X carries unique elements, Hamming metric resulted in 100% for all off the feature vectors.

4 Question 4

4.1 Precomputed clustering

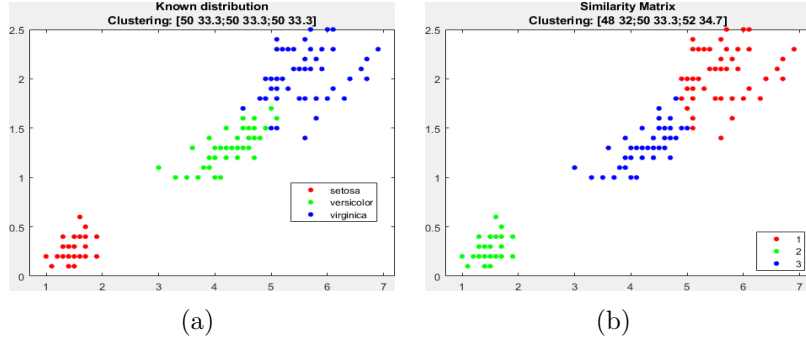


Figure 5: Clustering for different cases

4.2 Laplacian Matrix Normalization

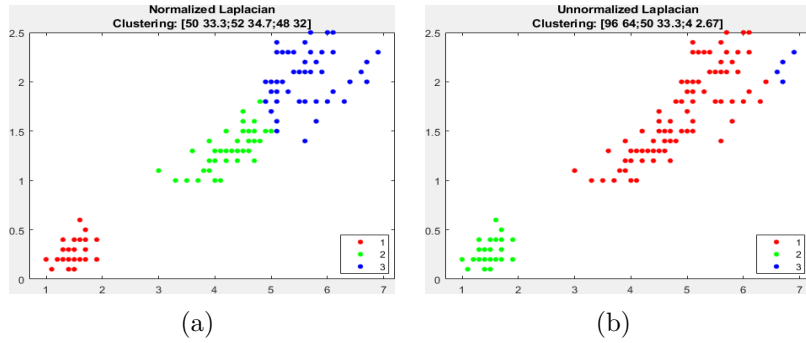


Figure 6: Clustering for different cases

Using unnormalized Laplacian matrix increased the number of misclassifications.

4.3 Distance Metrics

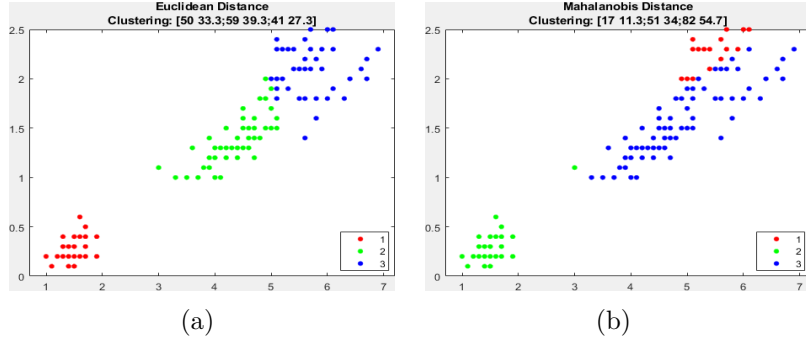


Figure 7: Clustering for different cases

When the clustering is based on Mahalanobis distance metric, the accuracy gets bad. Since, it takes into distances to the mean of normalized distribution. On the other hand, Euclidean metrics finds the neighborhood relationship simply calculating the distance between data point on 2 dimensional feature space. Hence, the result is closer to the real distribution.

4.4 Kernel Scale

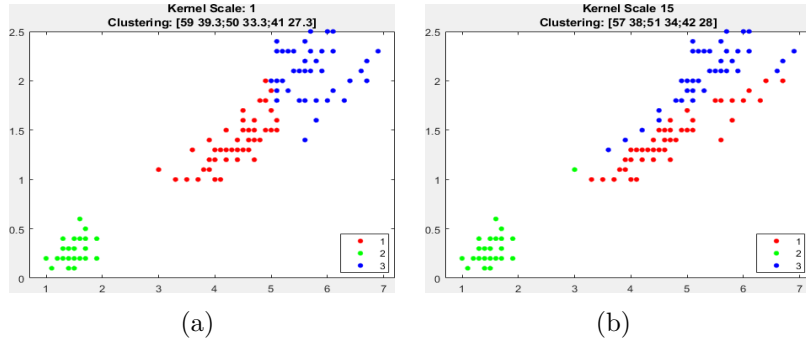


Figure 8: Clustering for different cases

Kernel scale parameter is used to scale the distance in the pairwise similarity definition:

$$S_{i,j} = e^{-\left(\frac{D_{i,j}}{\sigma}\right)^2} \quad (1)$$

According to Eqn. 1, the similarity of 2 data points is exponentially proportional to the similarity. Hence, using a larger kernel scale will result in increased similarity in further data points. This is exactly what is observed in Figure 8. Kernel scale of 15 resulted in further points to be classified in the same cluster, hence caused misclassifications.

5 APPENDIX

The code given in this section is shared @[Q](#).

5.1 Q1

```

1 load fisheriris
2 X = meas(:,3:4);
3 setosa_idx = strcmp(species,"setosa");
4 versicolor_idx = strcmp(species,"versicolor");
5 virginica_idx = strcmp(species,"virginica");
6
7 SETOSA = X(setosa_idx,:);
8 VERSICOLOR = X(versicolor_idx,:);
9 VIRGINICA = X(virginica_idx,:);
10
11 figure;
12 plot(SETOSA(:,1),SETOSA(:,2),'k*','MarkerSize',5);
13 hold on
14 plot(VERSICOLOR(:,1),VERSICOLOR(:,2),'bo','MarkerSize',5);
15 hold on
16 plot(VIRGINICA(:,1),VIRGINICA(:,2),'r+','MarkerSize',5);
17 title 'Fisher's Iris Data';
18 xlabel 'Petal Lengths (cm)';
19 ylabel 'Petal Widths (cm)';
20 legend('Setosa','Versicolor','Virginica','Location','northwest')
21
22 %%
23 % rng(1); % For reproducibility
24
25 Cs = {};
26 Cs{1} = [6 2 ; 1.5 0.3 ; 4 1];
27 Cs{2} = [5.2 2.5; 4.4 1.25; 3.8 5.5];
28 Cs{3} = [5 2.4; 4.6 1.45; 3.3 3.5];
29 Cs{4} = [4.5 2.1; 4.8 1.55; 3.2 3.2];
30 Cs{5} = [4.2 1.8; 5.2 1.78; 2.5 2.1];
31 Cs{6} = [4.1 1.4; 5.5 1.89; 1.8 1.05];
32 Cs{7} = [4 1 ; 6 2 ; 1.5 0.3];
33 gnd_truth = repelem([3;1;2],50,1);
34 accs = [];
35 CONFMATS = {};
36 for i = 1:7
37 [idx,C_new] = kmeans(X,3,'MaxIter',1,'Start',Cs{i});
38 confusion_matrix = confusionmat(gnd_truth,idx);
39 accuracy = sum(diag(confusion_matrix)) / size(X,1);
40 C_old = C_new;
41 accs(end+1) = accuracy;
42 CONFMATS{end+1} = confusion_matrix;
43 end
44 figure
45 plot(accs)

```

```

46 title({'\bf Accuracy vs Iterations'})
47 ylabel('Accuracy')
48 xlabel('kmeans iterations')
49 %%
50 x1 = min(X(:,1)):0.01:max(X(:,1));
51 x2 = min(X(:,2)):0.01:max(X(:,2));
52 [x1G,x2G] = meshgrid(x1,x2);
53 XGrid = [x1G(:),x2G(:)]; % Defines a fine grid on the plot
54
55 idx2Region = kmeans(XGrid,3,'MaxIter',1,'Start',C_new);
56 % Assigns each node in the grid to the closest centroid
57
58 figure;
59 gscatter(XGrid(:,1),XGrid(:,2),idx2Region,...
60         [0,0.75,0.75;0.75,0,0.75;0.75,0.75,0],'..');
61 hold on;
62 plot(X(setosa_idx,1),X(setosa_idx,2),'k*','MarkerSize',5);
63 hold on
64 plot(X(versicolor_idx,1),X(versicolor_idx,2),'ko','MarkerSize',5);
65 hold on
66 plot(X(virginica_idx,1),X(virginica_idx,2),'k+','MarkerSize',5);
67 title('Fisher''s Iris Data');
68 xlabel('Petal Lengths (cm)');
69 ylabel('Petal Widths (cm)');
70 legend('Region 1','Region 2','Region 3','Data','Location','SouthEast
        ');
71 hold off;

```

5.2 Q2

```
1 load fisheriris.mat
2 X = meas(:,3:4);
3 GMMModel = fitgmdist(X,3);
4 figure
5 y = [zeros(50,1);ones(50,1);2*ones(50,1)];
6 h = gscatter(X(:,1),X(:,2),y);
7 hold on
8 gmPDF = @(x,y) arrayfun(@(x0,y0) pdf(GMMModel,[x0 y0]),x,y);
9 g = gca;
10 fcontour(gmPDF,[g.XLim g.YLim])
11 title('\bf Scatter Plot and Fitted Gaussian Mixture Contours')
12 legend(h,'Model0','Model1','Model2')
13 hold off
```


5.3 Q3

```
1 X = [0 1 2 3; 1 0 4 5; 2 4 0 6; 3 5 6 0];
2 y = squareform(X);
3
4 distances = ["hamming","euclidean","chebychev","correlation",...
5             "cosine","minkowski","spearman","mahalanobis"];
6 figure('units','normalized','outerposition',[0 0 1 1])
7 for distance = distances
8     idx = find(distances == distance);
9     subplot(2,4,idx)
10    Z = linkage(X,'complete',distance);
11    dendrogram(Z)
12    title(distance)
13 end
```

5.4 Q4

```

1 %%
2 chdir('..')
3 addpath('export_fig')
4 chdir('HW4')
5 %%
6 load fisheriris
7 X = meas(:,3:4);
8 gscatter(X(:,1),X(:,2),species);
9 dist_temp = pdist(X);
10 dist = squareform(dist_temp);
11 S = exp(-dist.^2);
12 k = 3; % Number of clusters
13 rng('default') % For reproducibility
14 mat = tabulate(species);
15 title({'Known distribution',[ 'Clustering: ',mat2str(repmat
    ([50,33.3],3,1),3) ]})
16 %% Use similarity mat(:,2:end)rix
17 figure
18 idx = spectralcluster(S,k,'Distance','precomputed','
    LaplacianNormalization','symmetric');
19 gscatter(X(:,1),X(:,2),idx);
20 mat = tabulate(idx);
21 title({'Similarity Matrix',[ 'Clustering: ',mat2str(mat(:,2:end),3)
    ]})
22 %% Laplacian normalized
23 figure
24 idx2 = spectralcluster(X,k,'NumNeighbors',size(X,1),'
    LaplacianNormalization','symmetric');
25 gscatter(X(:,1),X(:,2),idx2);
26 mat = tabulate(idx2);
27 title({'Normalized Laplacian',[ 'Clustering: ',mat2str(mat(:,2:end)
    ,3) ]})
28 %% Laplacian unnormalized
29 figure
30 idx2 = spectralcluster(X,k,'NumNeighbors',size(X,1),'
    LaplacianNormalization','none');
31 gscatter(X(:,1),X(:,2),idx2);
32 mat = tabulate(idx2);
33 title({'Unnormalized Laplacian',[ 'Clustering: ',mat2str(mat(:,2:end)
    ,3) ]})
34
35 %% Distance Euclidean
36 figure
37 idx2 = spectralcluster(X,k,'NumNeighbors',size(X,1),'Distance','
    euclidean');
38 gscatter(X(:,1),X(:,2),idx2);
39 mat = tabulate(idx2);
40 title({'Euclidean Distance',[ 'Clustering: ',mat2str(mat(:,2:end),3)
    ]})

```

```

41 %% Distance mahalanobis
42 figure
43 idx2 = spectralcluster(X,k, 'NumNeighbors', size(X,1), 'Distance', '
    mahalanobis');
44 gscatter(X(:,1),X(:,2),idx2);
45 mat = tabulate(idx2);
46 title({'Mahalanobis Distance', ['Clustering: ', mat2str(mat(:,2:end)
    ,3)]})
47
48 %% Kernel Scale 1
49 figure
50 idx2 = spectralcluster(X,k, 'NumNeighbors', size(X,1), 'KernelScale', 1)
    ;
51 gscatter(X(:,1),X(:,2),idx2);
52 mat = tabulate(idx2);
53 title({'Kernel Scale: 1', ['Clustering: ', mat2str(mat(:,2:end),3)]})
54
55 %% Kernel Scale 15
56 figure
57 idx2 = spectralcluster(X,k, 'NumNeighbors', size(X,1), 'KernelScale',
    ,15);
58 gscatter(X(:,1),X(:,2),idx2);
59 mat = tabulate(idx2);
60 title({'Kernel Scale 15', ['Clustering: ', mat2str(mat(:,2:end),3)]})
61
62
63 %%
64 figHandles = findall(0, 'Type', 'figure');
65
66 for i = 1:numel(figHandles)
67     export_fig(['Q4_', num2str(i)], '-png', figHandles(i), '-append')
68 end

```