

Machine Learning and Data Mining II

LABWORK 2: Clustering

Groupwork:

Nguyễn Thị Quỳnh Anh
Vũ Yến Linh

BI11-029
BI11-152

I. AHC

Given the data $X = \{1, 2, 9, 12, 20\}$ in 1-D space.

1. Apply the AHC clustering using Single Linkage / Complete Linkage for the X dataset.

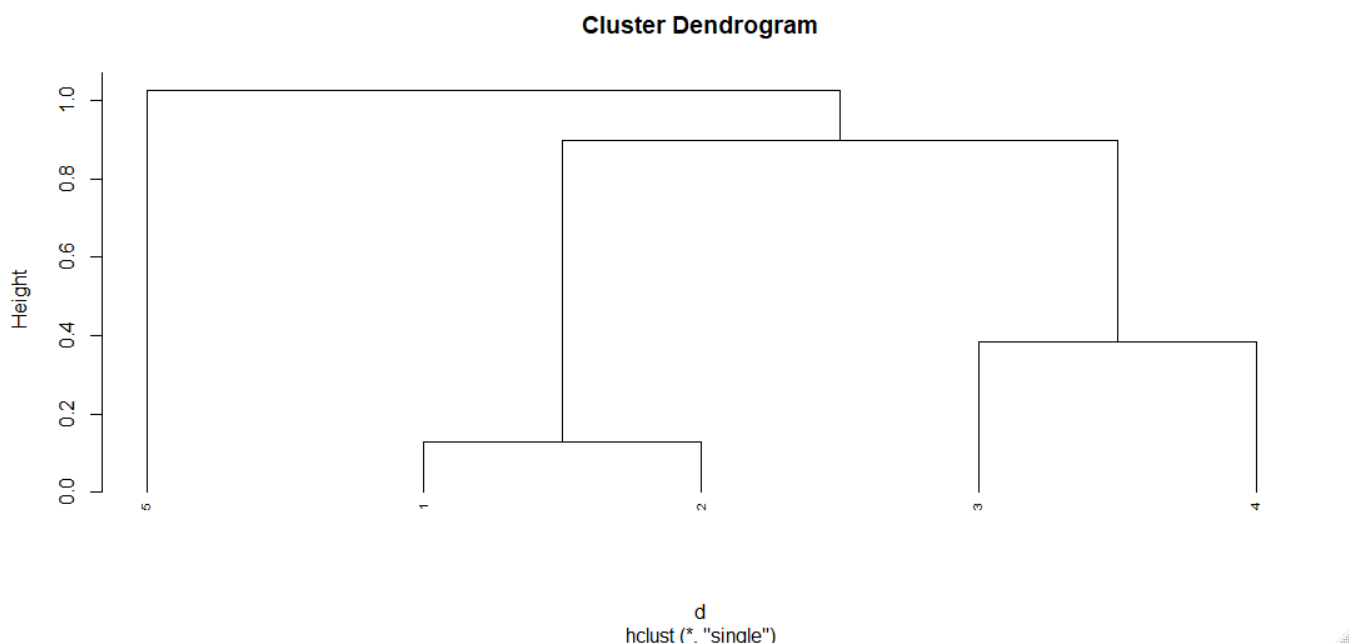
- Using R studio

```
> X <- c(1,2,9,12,20)
> X <- na.omit(X)
> X <- scale(X)
> d <- dist(X, method = "euclidean")
> hc1 <- hclust(d, method = "single")
> hc2 <- hclust(d, method = "complete")
> plot(hc1, cex = 0.6, hang = -1)
> |
```

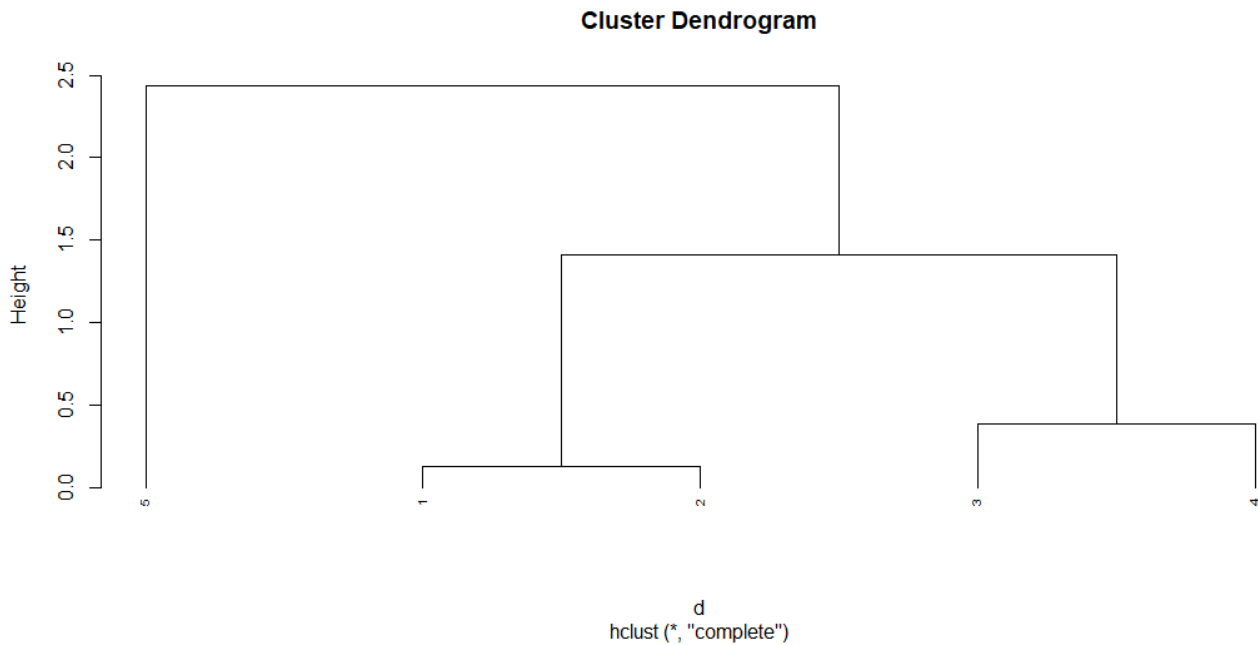
	V1
1	-1.00115255
2	-0.87279966
3	0.02567058
4	0.41072925
5	1.43755238

2. Draw the clustering result using available functions (dendrogram() in Matlab or hclust() in R, etc.).

- Single Linkage



- Complete Linkage



3. Apply on 2 more datasets from UCI. Make a study of data features. Observe the dendrogram and comment on results.

- Dataset from UCI :

Synchronous Machine Data Set:

<http://archive.ics.uci.edu/ml/datasets/Synchronous+Machine+Data+Set>

Energy efficiency Data Set:

<http://archive.ics.uci.edu/ml/datasets/Energy+efficiency>

- For Synchronous Machine Dataset:

Since we have already made a data analysis for the Synchronous dataset in the previous labwork, so now we will just make a summary for it.

```

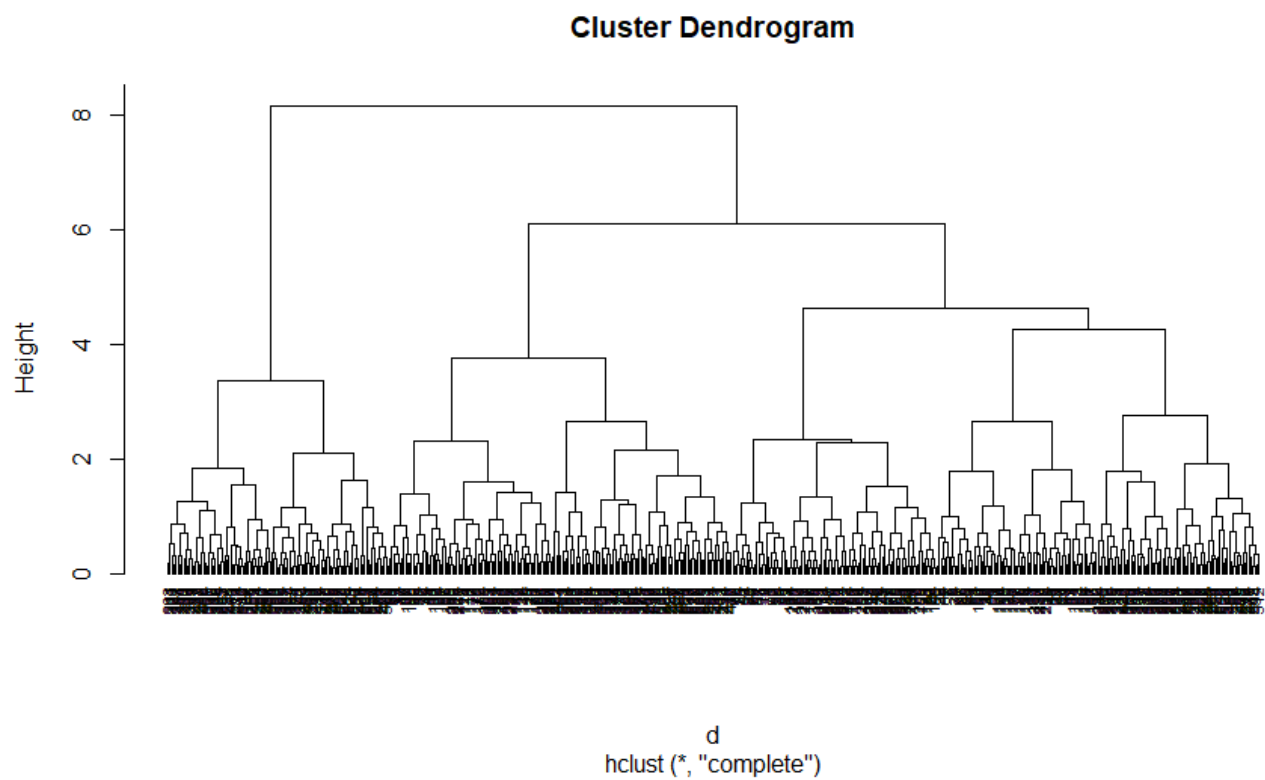
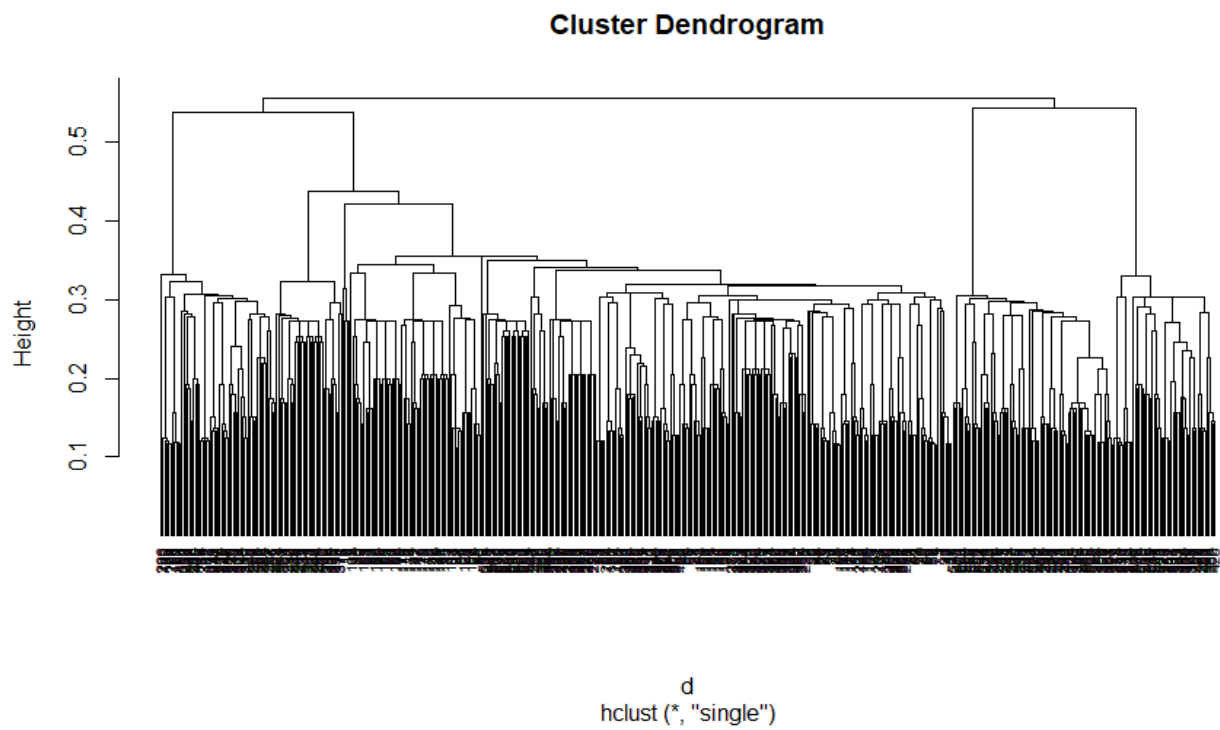
> view(synchronous.machine)
> syn <- synchronous.machine
> summary(syn)
      Iy      PF      e      dIf      If
Min.   :3.0   Min.   :0.6500   Min.   :0.0000   Min.   :0.0370   Min.   :1.217
1st Qu.:3.7   1st Qu.:0.7400   1st Qu.:0.0800   1st Qu.:0.1890   1st Qu.:1.369
Median :4.5   Median :0.8200   Median :0.1800   Median :0.3450   Median :1.525
Mean   :4.5   Mean   :0.8253   Mean   :0.1747   Mean   :0.3507   Mean   :1.531
3rd Qu.:5.3   3rd Qu.:0.9200   3rd Qu.:0.2600   3rd Qu.:0.4860   3rd Qu.:1.666
Max.   :6.0   Max.   :1.0000   Max.   :0.3500   Max.   :0.7690   Max.   :1.949
> cor(syn)
      Iy      PF      e      dIf      If
Iy  1.00000000 -0.04157389  0.04157389  0.4249449  0.4249449
PF -0.04157389  1.00000000 -1.00000000 -0.8610135 -0.8610135
e   0.04157389 -1.00000000  1.00000000  0.8610135  0.8610135
dIf 0.42494491 -0.86101347  0.86101347  1.0000000  1.0000000
If  0.42494491 -0.86101347  0.86101347  1.0000000  1.0000000
> |

```

```

> syn <- na.omit(syn)
> syn <- scale(syn)
> d <- dist(syn, method = "euclidean")
> hc1 <- hclust(d, method = "single")
> hc2 <- hclust(d, method = "complete")
> plot(hc1, cex = 0.6, hang = -1)
> |

```



- For Energy efficiency Dataset:

```

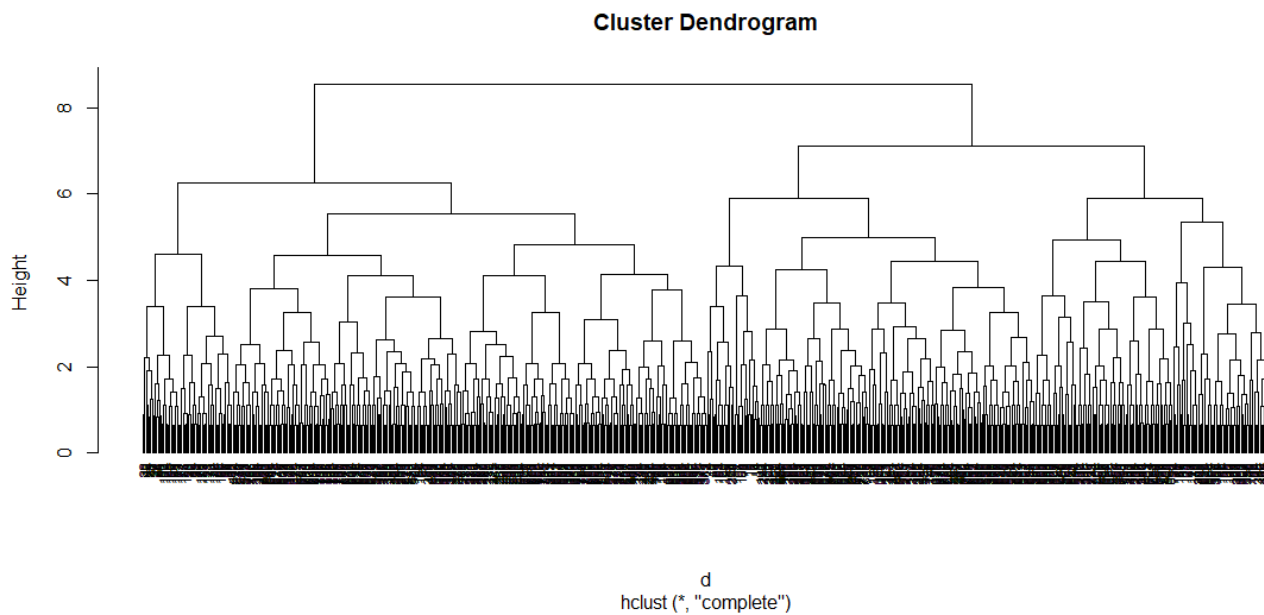
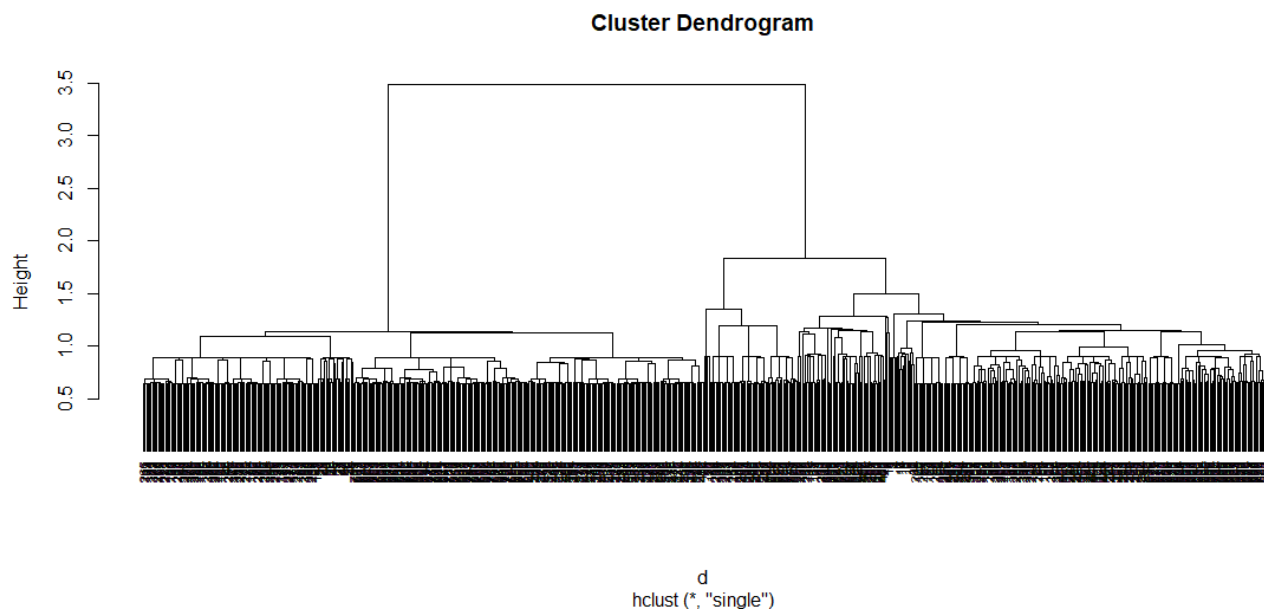
> enb <- ENB2012_data
> summary(enb)
      X1          X2          X3          X4          X5          X6          X7          X8          Y1          Y2
Min.   :0.6200   Min.   :514.5   Min.   :245.0   Min.   :110.2   Min.   :3.50   Min.   :2.00   Min.   :0.0000   Min.   :0.0000   Min.   : 6.01   Min.   :10.90
1st Qu.:0.6825   1st Qu.:606.4   1st Qu.:294.0   1st Qu.:140.9   1st Qu.:3.50   1st Qu.:2.75   1st Qu.:0.1000   1st Qu.:1.750   1st Qu.:12.99   1st Qu.:15.62
Median :0.7500   Median :673.8   Median :318.5   Median :183.8   Median :5.25   Median :3.50   Median :0.2500   Median :3.000   Median :18.95   Median :22.08
Mean   :0.7642   Mean   :671.7   Mean   :318.5   Mean   :176.6   Mean   :5.25   Mean   :3.50   Mean   :0.2344   Mean   :2.812   Mean   :22.31   Mean   :24.59
3rd Qu.:0.8300   3rd Qu.:741.1   3rd Qu.:343.0   3rd Qu.:220.5   3rd Qu.:7.00   3rd Qu.:4.25   3rd Qu.:0.4000   3rd Qu.:4.000   3rd Qu.:31.67   3rd Qu.:33.13
Max.   :0.9800   Max.   :808.5   Max.   :416.5   Max.   :220.5   Max.   :7.00   Max.   :5.00   Max.   :0.4000   Max.   :5.000   Max.   :43.10   Max.   :48.03

> var(enb)
      X1          X2          X3          X4          X5          X6          X7          X8          Y1          Y2
X1 1.118887e-02   -9.242069e+00   -0.9403911   -4.150839e+00   0.1533246   0.00000000   1.073424e-21   0.00000000   0.66416098   0.6383312
X2 -9.242069e+00   7.759164e+03   751.2907432   3.503937e+03   -132.3702738   0.00000000   5.473313e-19   0.00000000   -584.94150880   -563.9664689
X3 -9.403911e-01   7.512907e+02   1903.2698827   -5.759896e+02   21.4654498   0.00000000   0.000000e+00   0.00000000   200.58657888   177.2672425
X4 -4.150839e+00   3.503937e+03   -575.9895698   2.039963e+03   -76.9178618   0.00000000   -7.203513e-19   0.00000000   -392.76404384   -370.6168557
X5 1.533246e-01   -1.323703e+02   21.4654498   -7.691786e+01   3.0664928   0.00000000   0.000000e+00   0.00000000   15.71567080   14.9230052
X6 0.000000e+00   0.000000e+00   0.0000000   0.000000e+00   0.0000000   1.25162973   0.000000e+00   0.00000000   -0.02920078   0.1520860
X7 1.073424e-21   5.473313e-19   0.0000000   -7.203513e-19   0.0000000   0.00000000   1.774772e-02   0.04400261   0.36272731   0.2629852
X8 0.000000e+00   0.000000e+00   0.0000000   0.000000e+00   0.0000000   0.00000000   4.400261e-02   2.40547588   1.36727265   0.7454857
Y1 6.641610e-01   -5.849415e+02   200.5865789   -3.927640e+02   15.7156708   -0.02920078   3.627273e-01   1.36727265   101.81221616   93.6741331
Y2 6.383312e-01   -5.639665e+02   177.2672425   -3.706169e+02   14.9230052   0.15208605   2.629852e-01   0.74548566   93.67413308   90.5029827

> cov(enb)
      X1          X2          X3          X4          X5          X6          X7          X8          Y1          Y2
X1 1.118887e-02   -9.242069e+00   -0.9403911   -4.150839e+00   0.1533246   0.00000000   1.073424e-21   0.00000000   0.66416098   0.6383312
X2 -9.242069e+00   7.759164e+03   751.2907432   3.503937e+03   -132.3702738   0.00000000   5.473313e-19   0.00000000   -584.94150880   -563.9664689
X3 -9.403911e-01   7.512907e+02   1903.2698827   -5.759896e+02   21.4654498   0.00000000   0.000000e+00   0.00000000   200.58657888   177.2672425
X4 -4.150839e+00   3.503937e+03   -575.9895698   2.039963e+03   -76.9178618   0.00000000   -7.203513e-19   0.00000000   -392.76404384   -370.6168557
X5 1.533246e-01   -1.323703e+02   21.4654498   -7.691786e+01   3.0664928   0.00000000   0.000000e+00   0.00000000   15.71567080   14.9230052
X6 0.000000e+00   0.000000e+00   0.0000000   0.000000e+00   0.0000000   1.25162973   0.000000e+00   0.00000000   -0.02920078   0.1520860
X7 1.073424e-21   5.473313e-19   0.0000000   -7.203513e-19   0.0000000   0.00000000   1.774772e-02   0.04400261   0.36272731   0.2629852
X8 0.000000e+00   0.000000e+00   0.0000000   0.000000e+00   0.0000000   0.00000000   4.400261e-02   2.40547588   1.36727265   0.7454857
Y1 6.641610e-01   -5.849415e+02   200.5865789   -3.927640e+02   15.7156708   -0.02920078   3.627273e-01   1.36727265   101.81221616   93.6741331
Y2 6.383312e-01   -5.639665e+02   177.2672425   -3.706169e+02   14.9230052   0.15208605   2.629852e-01   0.74548566   93.67413308   90.5029827

> cor(enb)
      X1          X2          X3          X4          X5          X6          X7          X8          Y1          Y2
X1 1.000000e+00   -9.919015e-01   -0.2037817   -8.688234e-01   0.8277473   0.000000000   7.617400e-20   0.00000000   0.622271936   0.63433907
X2 -9.919015e-01   1.000000e+00   0.1955016   8.807195e-01   -0.8581477   0.000000000   4.664140e-20   0.00000000   -0.658119917   -0.67299893
X3 -0.2037817e-01   0.1955016e-01   1.0000000   -2.923165e-01   0.2809757   0.000000000   0.000000e+00   0.00000000   0.455671365   0.42711700
X4 -8.688234e-01   8.807195e-01   -0.2923165   1.000000e+00   -0.9725122   0.000000000   -1.197187e-19   0.00000000   -0.861828052   -0.86254660
X5 8.277473e-01   -8.581477e-01   0.2809757   -9.725122e-01   1.0000000   0.000000000   0.000000e+00   0.00000000   0.889430464   0.89578517
X6 0.000000e+00   0.000000e+00   0.0000000   0.000000e+00   0.0000000   1.000000000   0.000000e+00   0.00000000   -0.002586763   0.01428960
X7 7.617400e-20   4.664140e-20   0.0000000   -1.197187e-19   0.0000000   0.000000000   1.000000e+00   0.21296422   0.269841685   0.20750499
X8 0.000000e+00   0.000000e+00   0.0000000   0.000000e+00   0.0000000   0.000000000   2.129642e-01   1.00000000   0.087368460   0.05052512
Y1 6.222719e-01   -6.581199e-01   0.4556714   -8.618281e-01   0.8894305   -0.002586763   2.698417e-01   0.08736846   1.000000000   0.97586174
Y2 6.343391e-01   -6.729989e-01   0.4271170   -8.625466e-01   0.8957852   0.014289598   2.075050e-01   0.05052512   0.975861739   1.000000000

```



4. Conclude on the advantages and drawbacks of AHC.

- Advantage:

- AHC outputs a hierarchy → easier to decide on the number of clusters (by looking at the dendrogram).
- Easy to implement, doesn't require any input.
- Drawbacks:
 - For a n -element dataset, AHC must calculate the distance from one element to $(n-1)$ remaining elements so the AHC has the complexity of $O(n^2)$ → time complexity: not suitable for large databases.
 - Initial seeds have a strong impact on the final results.
 - The order of the data has an impact on the final results.
 - Very sensitive to noise/outliers.

II. K-means:

Select two datasets from UCI:

- 1.Iris Dataset: <http://archive.ics.uci.edu/ml/datasets/Iris>
- 2.Wine Data Set: <http://archive.ics.uci.edu/ml/datasets/Wine>

1. Run experiments with k-means. Explain the experimental protocol.

For Iris Dataset:

- **Model kmeans_re:**

[illegible]

The 3 clusters are made which are of 50, 62, and 38 sizes respectively. Within the cluster, the sum of squares is 88.4%.

- **Cluster identification and confusion matrix:**

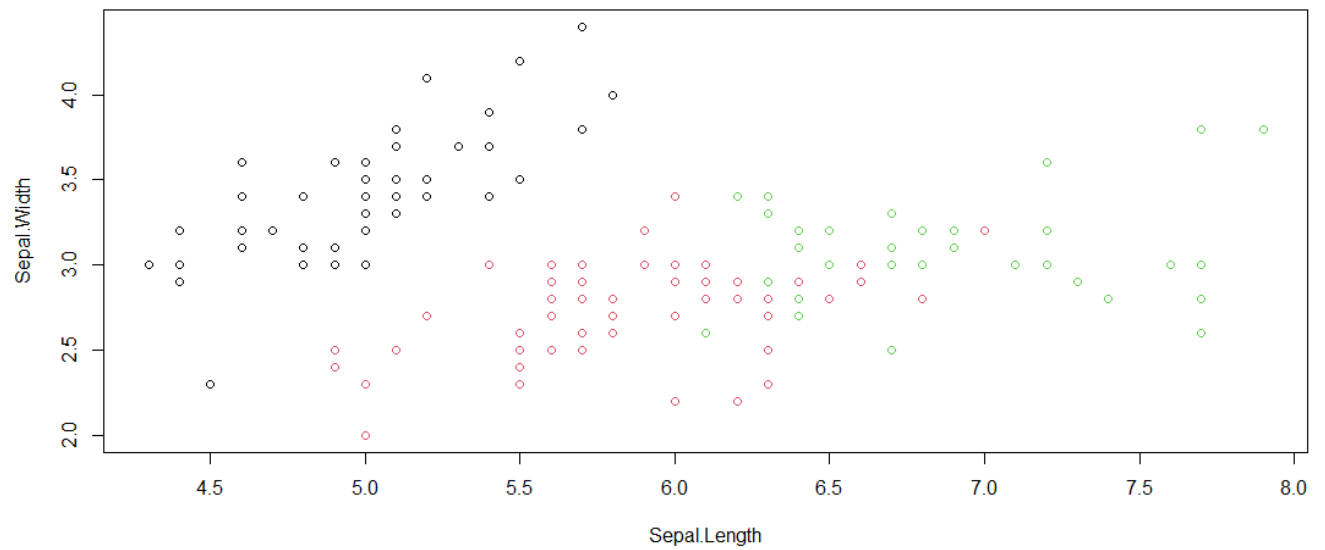
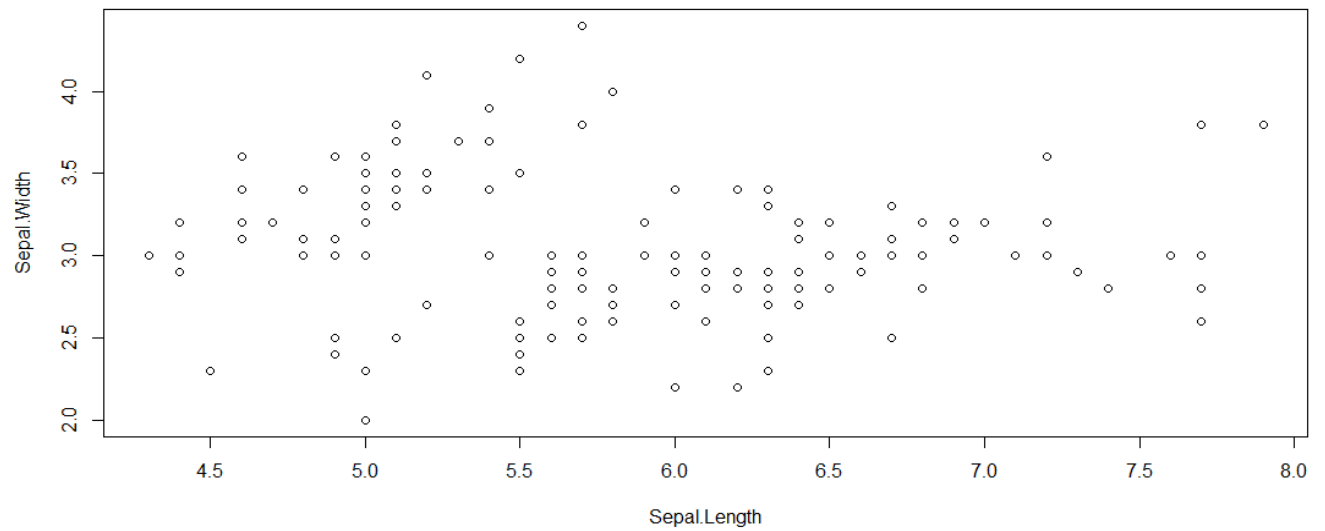
[illegible]

50 Setosa are so appropriately categorized as Setosa. 48 of the 62 Versicolor are appropriately categorized as Versicolor, and 14 are categorized as virginica. 19 of the 36 virginica have been reliably identified as virginica, while two have been identified as Versicolor.

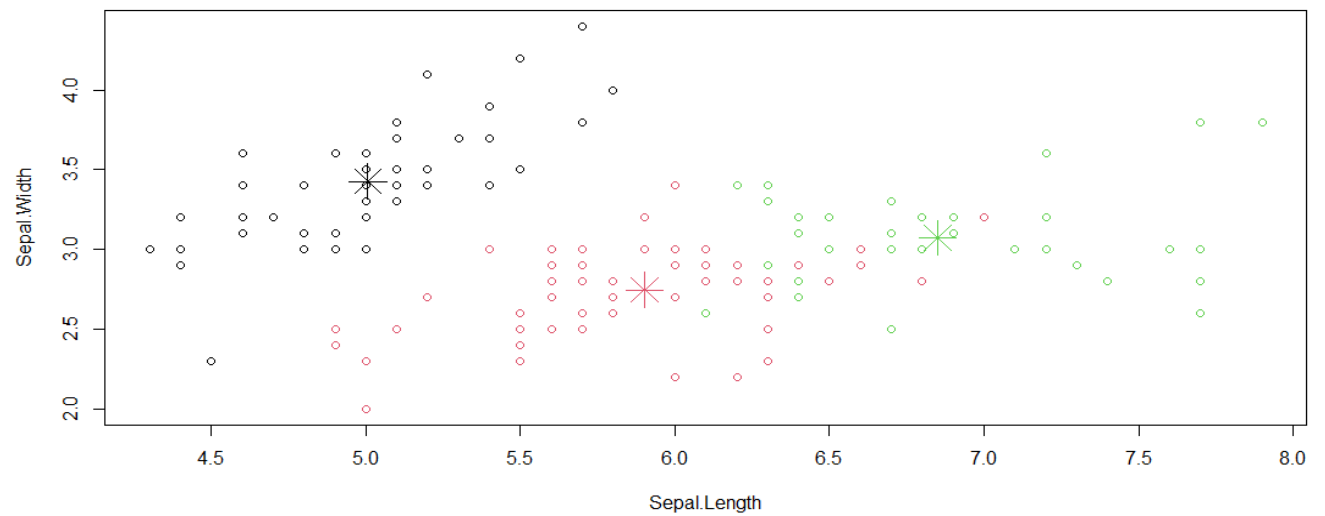
- **K-means with 3 clusters plot:**

```
# Model Evaluation and visualization
plot(iris_1[c("Sepal.Length", "Sepal.width")])
plot(iris_1[c("Sepal.Length", "Sepal.width")],
     col = kmeans.re$cluster)
plot(iris_1[c("Sepal.Length", "Sepal.width")],
     col = kmeans.re$cluster,
     main = "K-means with 3 clusters")

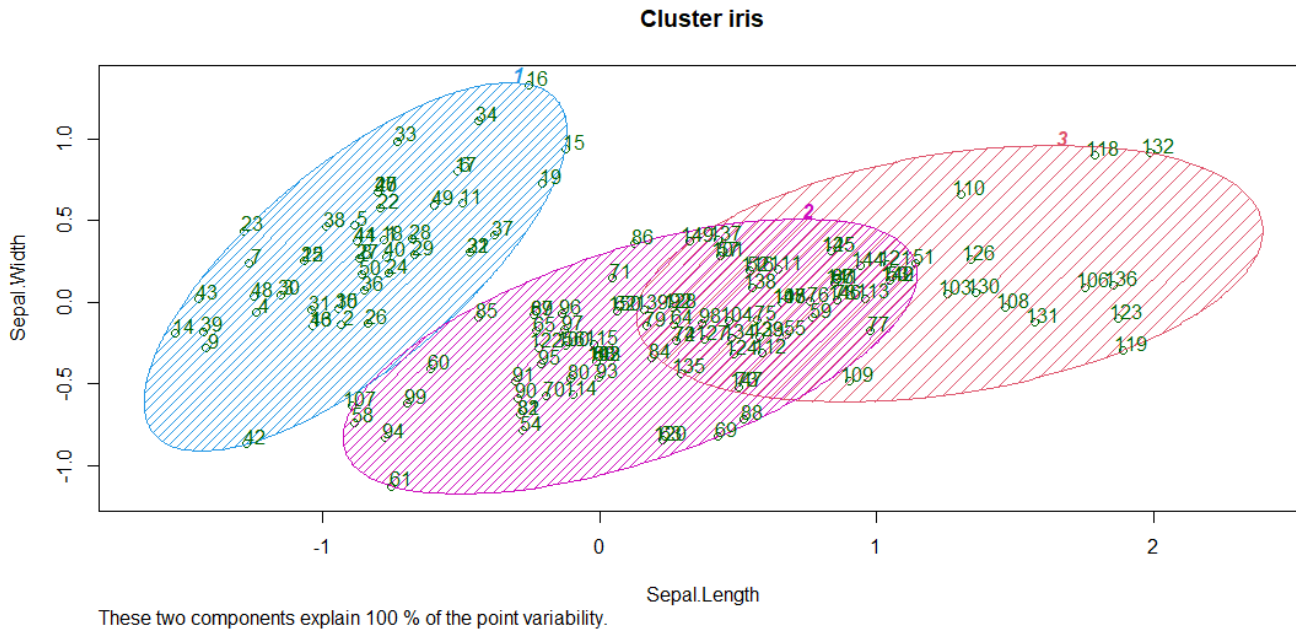
## Plotting cluster centers
kmeans.re$centers
kmeans.re$centers[, c("Sepal.Length", "Sepal.width")]
```



K-means with 3 clusters



- **Plot of clusters:**

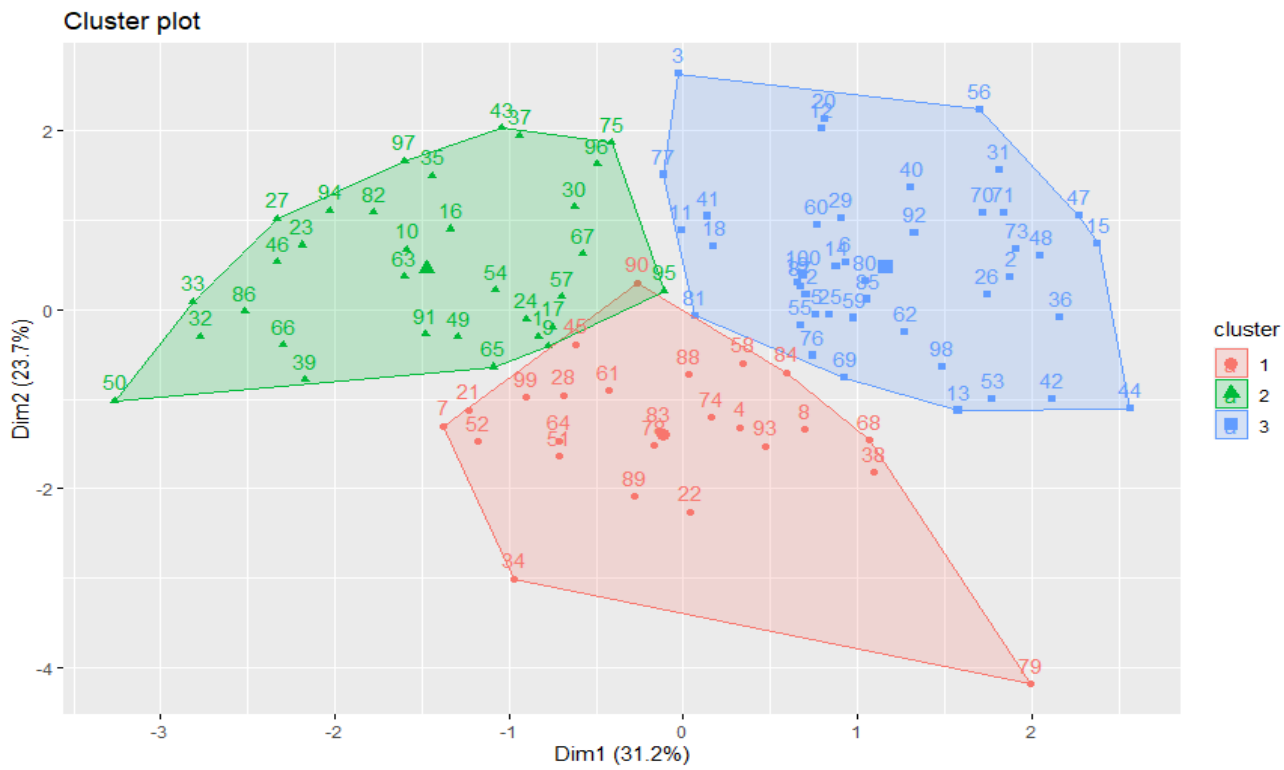


For Wine dataset:

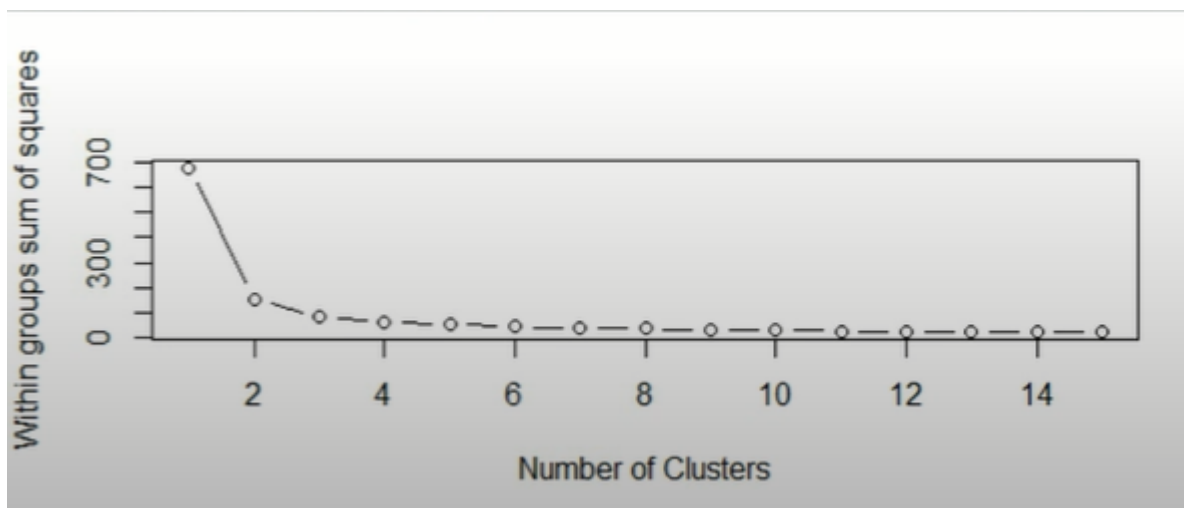
- **Model kmeans_re:**

[illegible]

- Plot of clusters:



2. Explain the centroid initialization of the k-means used in the first question.



We base on the elbow method to choose $k = 3$ for the Iris Dataset

For Wine dataset: $k = 4$

[illegible]

4. Calculate the clustering quality (any criteria in slides).

```
> confusionMatrix(cm)
Confusion Matrix and Statistics

          setosa versicolor virginica
setosa      20           0           0
versicolor   0          20           0
virginica    0           0          20

Overall Statistics

          Accuracy : 1
          95% CI : (0.9404, 1)
    No Information Rate : 0.3333
    P-Value [Acc > NIR] : < 2.2e-16

          Kappa : 1

McNemar's Test P-Value : NA

Statistics by Class:

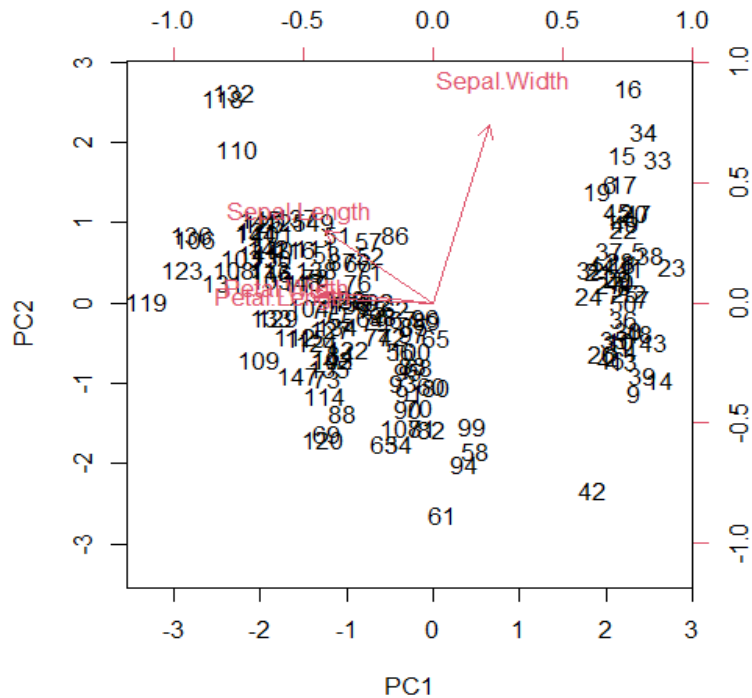
                Class: setosa Class: versicolor Class: virginica
Sensitivity      1.0000          1.0000          1.0000
Specificity      1.0000          1.0000          1.0000
Pos Pred Value   1.0000          1.0000          1.0000
Neg Pred Value   1.0000          1.0000          1.0000
Prevalence       0.3333          0.3333          0.3333
Detection Rate   0.3333          0.3333          0.3333
Detection Prevalence 0.3333          0.3333          0.3333
Balanced Accuracy 1.0000          1.0000          1.0000
```

The model's accuracy was 100% and its p-value was less than 1.

This suggests that the model is good.

5. Use PCA or SVD to visualize the data distribution in 2D/3D.

- PCA for Iris dataset :



- PCA for Wine dataset:

