# Unsupervised Learning - K-mean Clustering

Abhirup Sen

11/06/2021

## R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see http://rmarkdown.rstudio.com.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
set.seed(256)
# Load iris dataset into a new variable iris2
data(iris)
str(iris)
```

```
## 'data.frame':    150 obs. of  5 variables:
##  $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
##  $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
##  $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
##  $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
##  $ Species     : Factor w/ 3 levels "setosa","versicolor",..: 1 1 1 1 1 1 1 1 1 1 ...
```

```
#Remove the initial label of Species from original dataset
iris2 <- iris[,-5]
iris2
```

```
##     Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1            5.1         3.5          1.4         0.2
## 2            4.9         3.0          1.4         0.2
## 3            4.7         3.2          1.3         0.2
## 4            4.6         3.1          1.5         0.2
## 5            5.0         3.6          1.4         0.2
## 6            5.4         3.9          1.7         0.4
## 7            4.6         3.4          1.4         0.3
## 8            5.0         3.4          1.5         0.2
## 9            4.4         2.9          1.4         0.2
## 10           4.9         3.1          1.5         0.1
## 11           5.4         3.7          1.5         0.2
## 12           4.8         3.4          1.6         0.2
## 13           4.8         3.0          1.4         0.1
## 14           4.3         3.0          1.1         0.1
## 15           5.8         4.0          1.2         0.2
## 16           5.7         4.4          1.5         0.4
```

```
## 17           5.4         3.9         1.3         0.4
## 18           5.1         3.5         1.4         0.3
## 19           5.7         3.8         1.7         0.3
## 20           5.1         3.8         1.5         0.3
## 21           5.4         3.4         1.7         0.2
## 22           5.1         3.7         1.5         0.4
## 23           4.6         3.6         1.0         0.2
## 24           5.1         3.3         1.7         0.5
## 25           4.8         3.4         1.9         0.2
## 26           5.0         3.0         1.6         0.2
## 27           5.0         3.4         1.6         0.4
## 28           5.2         3.5         1.5         0.2
## 29           5.2         3.4         1.4         0.2
## 30           4.7         3.2         1.6         0.2
## 31           4.8         3.1         1.6         0.2
## 32           5.4         3.4         1.5         0.4
## 33           5.2         4.1         1.5         0.1
## 34           5.5         4.2         1.4         0.2
## 35           4.9         3.1         1.5         0.2
## 36           5.0         3.2         1.2         0.2
## 37           5.5         3.5         1.3         0.2
## 38           4.9         3.6         1.4         0.1
## 39           4.4         3.0         1.3         0.2
## 40           5.1         3.4         1.5         0.2
## 41           5.0         3.5         1.3         0.3
## 42           4.5         2.3         1.3         0.3
## 43           4.4         3.2         1.3         0.2
## 44           5.0         3.5         1.6         0.6
## 45           5.1         3.8         1.9         0.4
## 46           4.8         3.0         1.4         0.3
## 47           5.1         3.8         1.6         0.2
## 48           4.6         3.2         1.4         0.2
## 49           5.3         3.7         1.5         0.2
## 50           5.0         3.3         1.4         0.2
## 51           7.0         3.2         4.7         1.4
## 52           6.4         3.2         4.5         1.5
## 53           6.9         3.1         4.9         1.5
## 54           5.5         2.3         4.0         1.3
## 55           6.5         2.8         4.6         1.5
## 56           5.7         2.8         4.5         1.3
## 57           6.3         3.3         4.7         1.6
## 58           4.9         2.4         3.3         1.0
## 59           6.6         2.9         4.6         1.3
## 60           5.2         2.7         3.9         1.4
## 61           5.0         2.0         3.5         1.0
## 62           5.9         3.0         4.2         1.5
## 63           6.0         2.2         4.0         1.0
## 64           6.1         2.9         4.7         1.4
## 65           5.6         2.9         3.6         1.3
## 66           6.7         3.1         4.4         1.4
## 67           5.6         3.0         4.5         1.5
## 68           5.8         2.7         4.1         1.0
## 69           6.2         2.2         4.5         1.5
## 70           5.6         2.5         3.9         1.1
```

```
## 71          5.9          3.2          4.8          1.8
## 72          6.1          2.8          4.0          1.3
## 73          6.3          2.5          4.9          1.5
## 74          6.1          2.8          4.7          1.2
## 75          6.4          2.9          4.3          1.3
## 76          6.6          3.0          4.4          1.4
## 77          6.8          2.8          4.8          1.4
## 78          6.7          3.0          5.0          1.7
## 79          6.0          2.9          4.5          1.5
## 80          5.7          2.6          3.5          1.0
## 81          5.5          2.4          3.8          1.1
## 82          5.5          2.4          3.7          1.0
## 83          5.8          2.7          3.9          1.2
## 84          6.0          2.7          5.1          1.6
## 85          5.4          3.0          4.5          1.5
## 86          6.0          3.4          4.5          1.6
## 87          6.7          3.1          4.7          1.5
## 88          6.3          2.3          4.4          1.3
## 89          5.6          3.0          4.1          1.3
## 90          5.5          2.5          4.0          1.3
## 91          5.5          2.6          4.4          1.2
## 92          6.1          3.0          4.6          1.4
## 93          5.8          2.6          4.0          1.2
## 94          5.0          2.3          3.3          1.0
## 95          5.6          2.7          4.2          1.3
## 96          5.7          3.0          4.2          1.2
## 97          5.7          2.9          4.2          1.3
## 98          6.2          2.9          4.3          1.3
## 99          5.1          2.5          3.0          1.1
## 100         5.7          2.8          4.1          1.3
## 101         6.3          3.3          6.0          2.5
## 102         5.8          2.7          5.1          1.9
## 103         7.1          3.0          5.9          2.1
## 104         6.3          2.9          5.6          1.8
## 105         6.5          3.0          5.8          2.2
## 106         7.6          3.0          6.6          2.1
## 107         4.9          2.5          4.5          1.7
## 108         7.3          2.9          6.3          1.8
## 109         6.7          2.5          5.8          1.8
## 110         7.2          3.6          6.1          2.5
## 111         6.5          3.2          5.1          2.0
## 112         6.4          2.7          5.3          1.9
## 113         6.8          3.0          5.5          2.1
## 114         5.7          2.5          5.0          2.0
## 115         5.8          2.8          5.1          2.4
## 116         6.4          3.2          5.3          2.3
## 117         6.5          3.0          5.5          1.8
## 118         7.7          3.8          6.7          2.2
## 119         7.7          2.6          6.9          2.3
## 120         6.0          2.2          5.0          1.5
## 121         6.9          3.2          5.7          2.3
## 122         5.6          2.8          4.9          2.0
## 123         7.7          2.8          6.7          2.0
## 124         6.3          2.7          4.9          1.8
```

```
## 125            6.7            3.3            5.7            2.1
## 126            7.2            3.2            6.0            1.8
## 127            6.2            2.8            4.8            1.8
## 128            6.1            3.0            4.9            1.8
## 129            6.4            2.8            5.6            2.1
## 130            7.2            3.0            5.8            1.6
## 131            7.4            2.8            6.1            1.9
## 132            7.9            3.8            6.4            2.0
## 133            6.4            2.8            5.6            2.2
## 134            6.3            2.8            5.1            1.5
## 135            6.1            2.6            5.6            1.4
## 136            7.7            3.0            6.1            2.3
## 137            6.3            3.4            5.6            2.4
## 138            6.4            3.1            5.5            1.8
## 139            6.0            3.0            4.8            1.8
## 140            6.9            3.1            5.4            2.1
## 141            6.7            3.1            5.6            2.4
## 142            6.9            3.1            5.1            2.3
## 143            5.8            2.7            5.1            1.9
## 144            6.8            3.2            5.9            2.3
## 145            6.7            3.3            5.7            2.5
## 146            6.7            3.0            5.2            2.3
## 147            6.3            2.5            5.0            1.9
## 148            6.5            3.0            5.2            2.0
## 149            6.2            3.4            5.4            2.3
## 150            5.9            3.0            5.1            1.8
```

```r
# Apply K-mean clustering to understand the Species from other attributes
?kmeans
```

```
## starting httpd help server ... done
```

```r
kmeans.result <- kmeans(iris2, centers = 3, nstart = 20)
kmeans.result
```

```
## K-means clustering with 3 clusters of sizes 50, 38, 62
##
## Cluster means:
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1     5.006000    3.428000     1.462000    0.246000
## 2     6.850000    3.073684     5.742105    2.071053
## 3     5.901613    2.748387     4.393548    1.433871
##
## Clustering vector:
##    [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##   [38] 1 1 1 1 1 1 1 1 1 1 1 1 1 3 3 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
##   [75] 3 3 3 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 2 3 2 2 2 2 3 2 2 2 2
##  [112] 2 2 3 3 2 2 2 2 3 2 3 2 3 2 2 3 3 2 2 2 2 2 3 2 2 2 2 3 2 2 2 3 2 2 2 3 2
##  [149] 2 3
##
## Within cluster sum of squares by cluster:
## [1] 15.15100 23.87947 39.82097
##  (between_SS / total_SS =  88.4 %)
```
```

```
##
## Available components:
##
## [1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss"
## [6] "betweenss"    "size"         "iter"         "ifault"
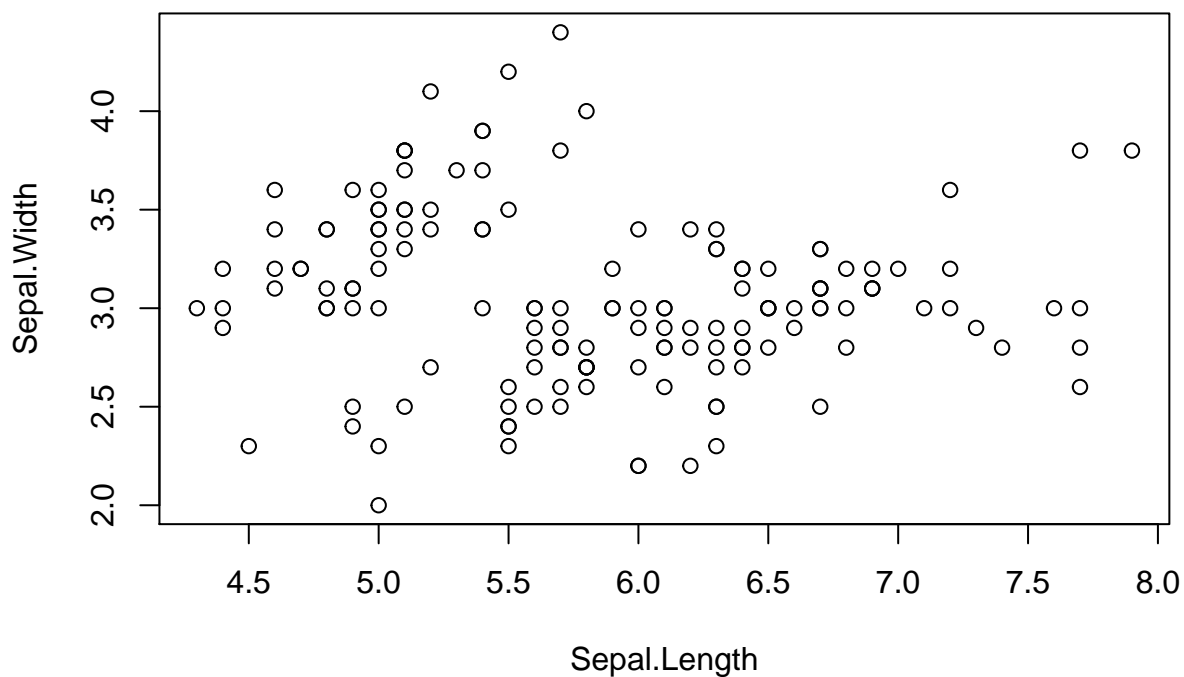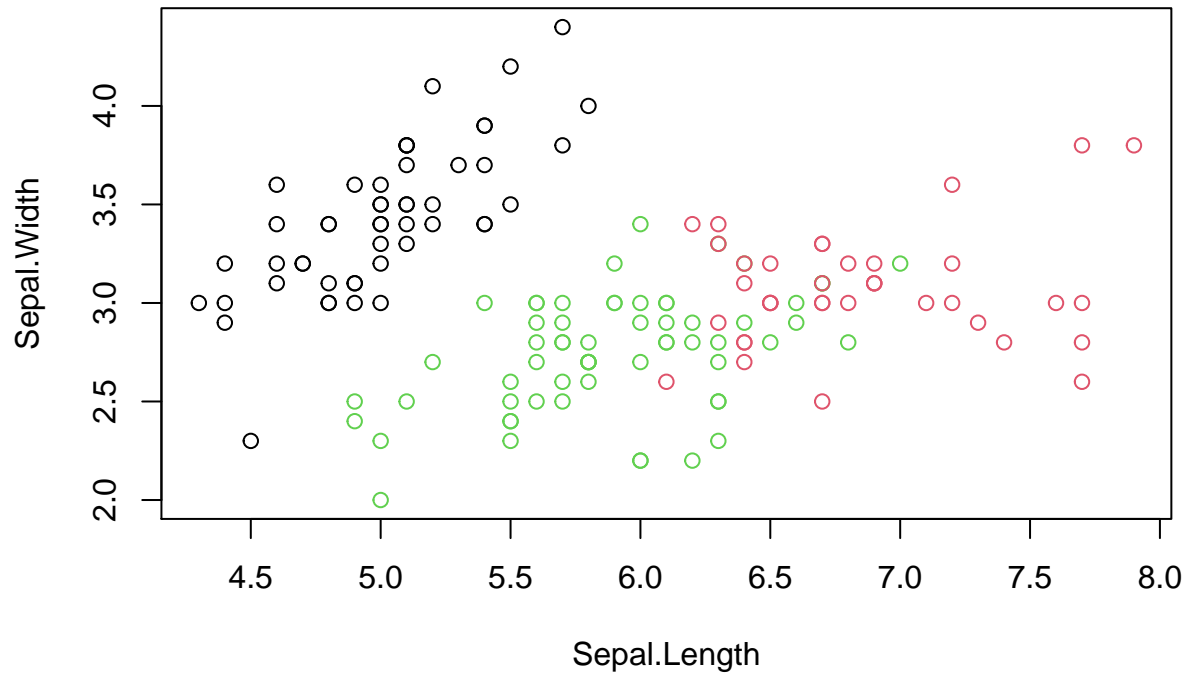```

```
# See the cluster identification for each observation
kmeans.result$cluster
```

```
##   [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##  [38] 1 1 1 1 1 1 1 1 1 1 1 1 1 3 3 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
##  [75] 3 3 3 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 2 3 2 2 2 2 3 2 2 2 2
## [112] 2 2 3 3 2 2 2 2 3 2 3 2 3 2 2 3 3 2 2 2 2 2 3 2 2 2 2 3 2 2 2 2 3 2 2 2 3 2
## [149] 2 3
```

```
# Compare with original label
table(iris$Species, kmeans.result$cluster)
```

```
##
##               1  2  3
##   setosa     50  0  0
##   versicolor  0  2 48
##   virginica   0 36 14
```

```
# Visualizing and interpreting results of k-means()
plot(iris2[c("Sepal.Length", "Sepal.Width")])
```

```
plot(iris2[c("Sepal.Length", "Sepal.Width")], col = kmeans.result$cluster)
```



```
plot(iris2[c("Sepal.Length", "Sepal.Width")], col = kmeans.result$cluster, main = "K-Means with 3 cluste
```

## K–Means with 3 clusters



```r
# plot cluster centers
kmeans.result$centers
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width
## 1     5.006000    3.428000     1.462000    0.246000
## 2     6.850000    3.073684     5.742105    2.071053
## 3     5.901613    2.748387     4.393548    1.433871
```

```r
kmeans.result$centers[,c("Sepal.Length", "Sepal.Width")]
```

```
##   Sepal.Length Sepal.Width
## 1     5.006000    3.428000
## 2     6.850000    3.073684
## 3     5.901613    2.748387
```

```r
#points(kmeans.result$centers[,c("Sepal.Length", "Sepal.Width")], col = 1:3,pch = 8, cex=3)
```

```r
# Visualising the clusters
library(cluster)
y_kmeans <- kmeans.result$cluster
?clusplot()
clusplot(iris2[,c("Sepal.Length", "Sepal.Width")],
         y_kmeans,
         lines = 0,
```

```
        shade = TRUE,
        color = TRUE,
        labels = 2,
        plotchar = FALSE,
        span = TRUE,
        main = paste('Clusters of iris'),
        xlab = 'Sepal.Length',
        ylab = 'Sepal.Width')
```

## Clusters of iris



Sepal.Length

These two components explain 100 % of the point variability.

```
# Determining number of clusters
# See Total within sum of square error
kmeans.result$tot.withinss
```

```
## [1] 78.85144
```

```
# Initialize total within cluster sum of squares error: wcss
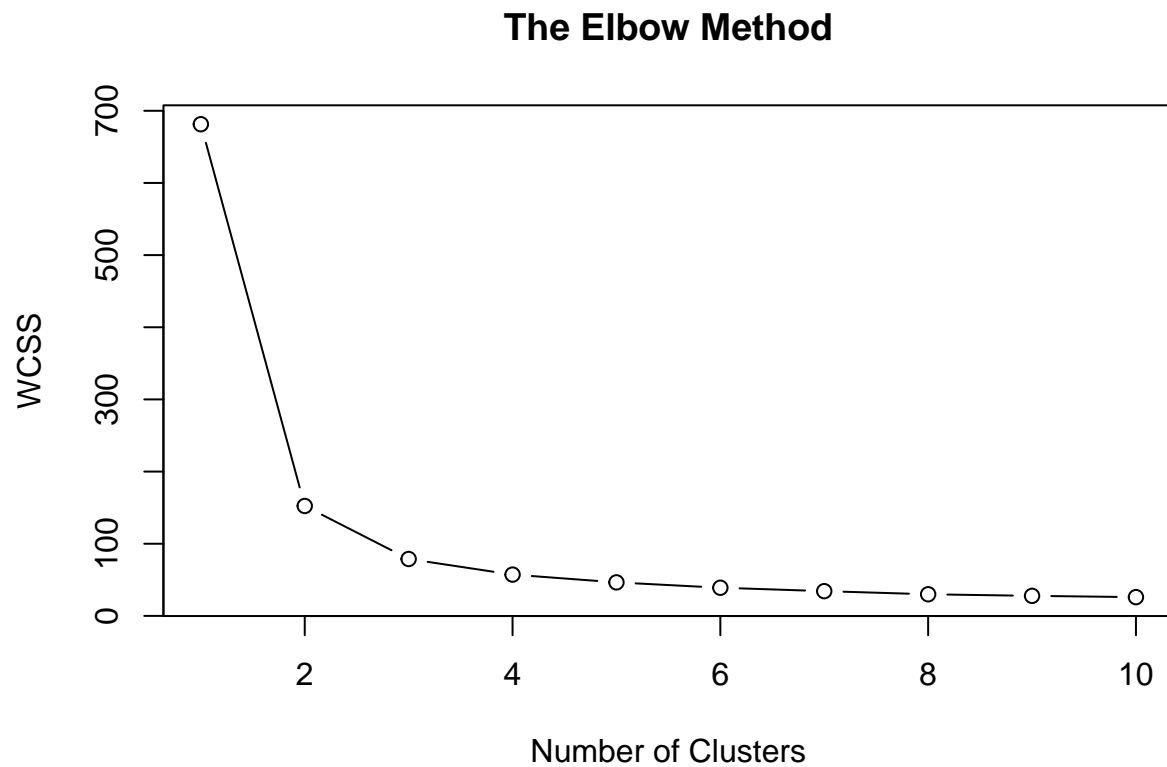set.seed(6)
wcss = vector()
```

```
# For 1 to 15 cluster centers check the WCSS
for (i in 1:10) wcss[i] = sum(kmeans(iris2, centers = i,nstart = 20)$withinss)
# Plot WSS vs. number of clusters
plot(1:10,
     wcss,
```

```
      type = "b",
      main = paste('The Elbow Method'),
      xlab = "Number of Clusters",
      ylab = "WCSS")
```

## The Elbow Method



```
# Set k equal to the number of clusters corresponding to the elbow location
k <-3

# Fitting K-Means to the dataset
set.seed(29)
kmeans.result <- kmeans(iris2, centers = 3, nstart = 20)
y_kmeans = kmeans.result$cluster
```