

Clustering

K-means Clustering

The following document shows k-means clustering of the 'iris' data. At first, we remove species from the data to cluster. After that, we apply function `kmeans()` to `iris2`, and store the clustering result in `kmeans.result`. The cluster number is set to 3 in the code below.

```
> iris2<-iris
> iris2$Species<-NULL
```

[illegible]

The clustering result is then compared with the class label (species) to check whether similar objects are grouped together.

```
> table(iris$Species, kmeans.result$cluster)

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

The above result shows that cluster “setosa” can be easily separated from the other clusters, and that clusters “versicolor” and “virginica” are to a small degree overlapped with each other.

Next, the clusters and their centres are plotted.

Note that there are four dimensions in the data and that only the first two dimensions are used to draw the plot below.

Some black points close to the green centre (asterisk) are actually closer to the black center in the four dimensional space. We also need to be aware that the results of k-means clustering may vary from run to run, due to random selection of initial cluster centers.

```
> library(tm)
Warning message:
package 'tm' was built under R version 2.15.3
> library(SnowballC)
Warning message:
package 'SnowballC' was built under R version 2.15.3
> |
```

```
> plot(iris2[c("Sepal.Length", "Sepal.Width")], col=kmeans.result$cluster)
> #plot cluster centres
> points(kmeans.result$centres[,c("Sepal.Length", "Sepal.Width")], col=1:3,
+ pch=8, cex=2)
```

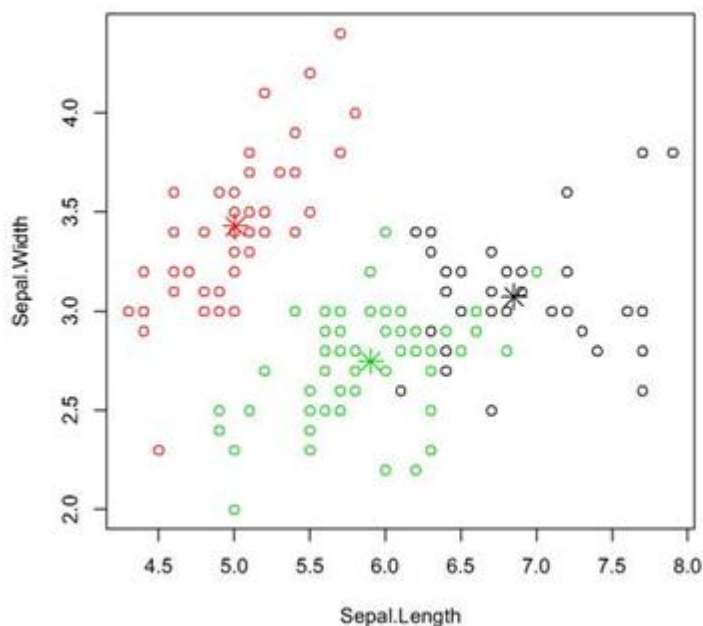


Fig: Results of K-Means Clustering

Text Clustering using K-means clustering:



Text Clustering.R

