

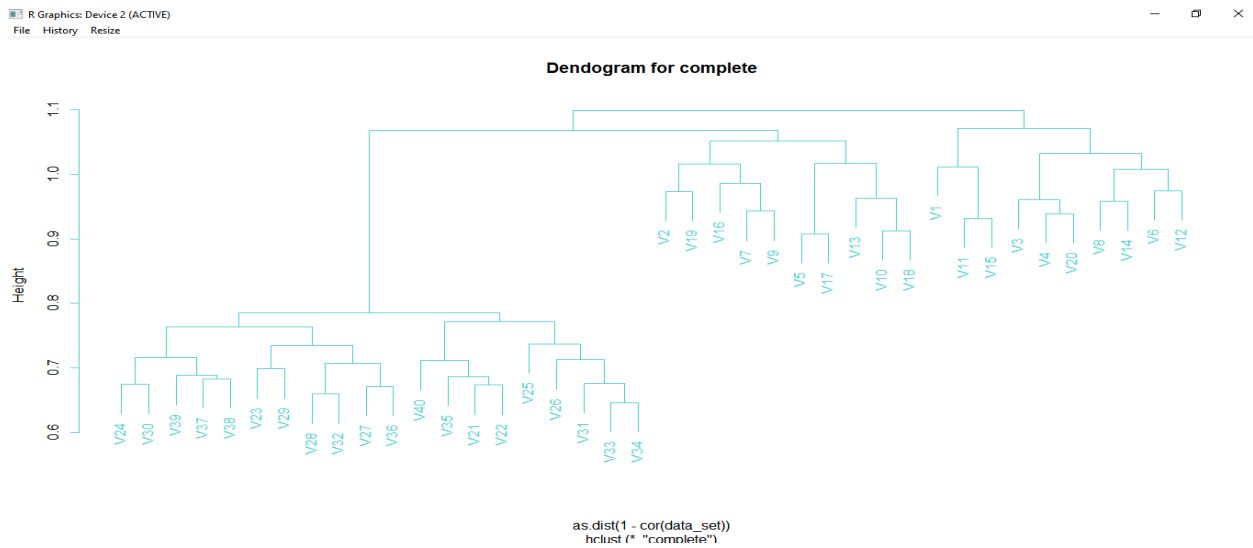
Question 2:

a)

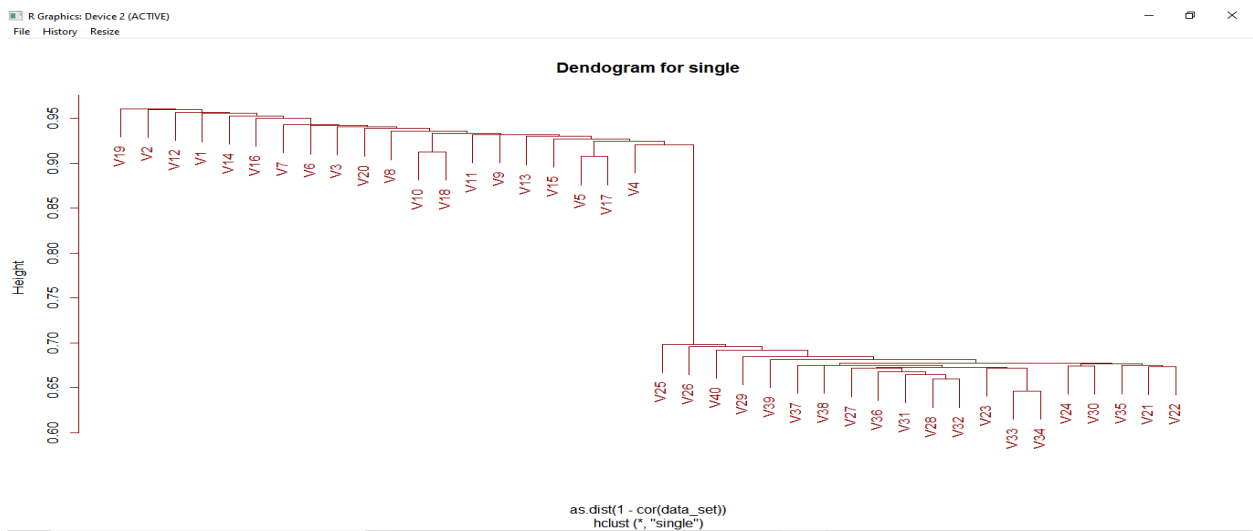
```
> data_set = read.csv("c:/shreyas/document/books/statistical_data_mining/ch10Ex11.csv", header = FALSE)
> |
```

b) Hierarchical clustering using correlation based distance.

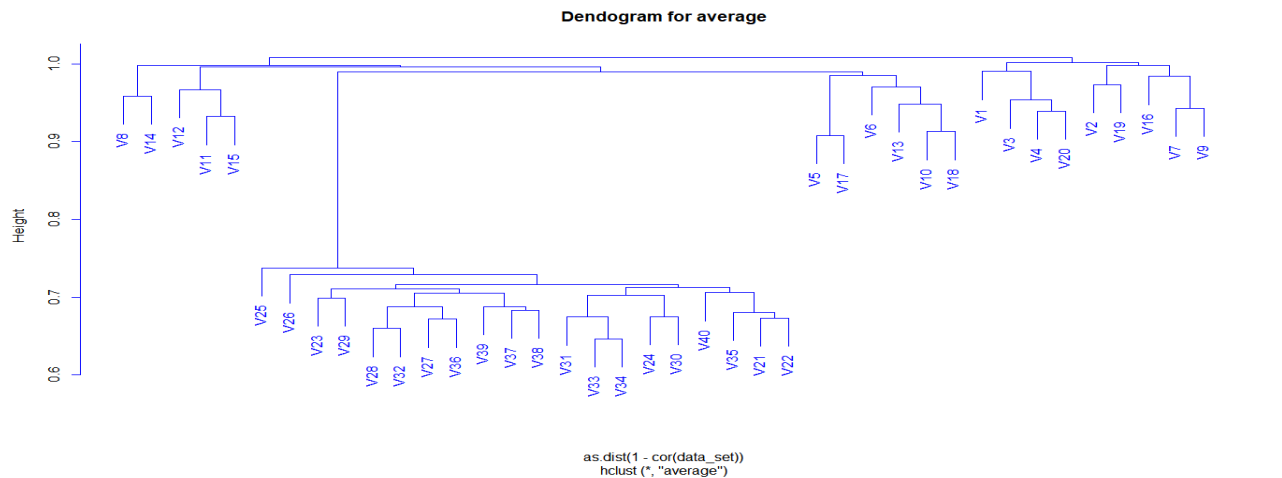
Dendrogram for complete linkage



Dendrogram of single linkage



Dendrogram for Average Linkage



Dendrogram for centroid Linkage

- Yes the result does depend on the type linkage used for clustering. In the above case, We got two clusters for complete and single linkage and three clusters for average linkage. Therefore , the result depends upon the type of linkage used.

c)

- There fore when tried with K-means as follows:

```
> kmean = kmeans(t(data_set), centers=2)
> kmean$cluster
```

	V1	V2	V3	V4	V5	V6	V7	V8	V9	V10	V11	V12	V13	V14	V15	V16	V17	V18	V19	V20	V21	V22	V23	V24	V25	V26	V27	V28	V29	V30	V31	V32	V33	V34	V35	V36	V37	V38	V39	V40
Cluster	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		

- K-means separates the data into two different groups as shown above.

Running K-means with 2 clusters for different expression values.

[illegible]

- The K-means with 2 clusters result shows that the genes 11-20 and 501-600 differ most between 2 groups as shown above figure.