

Individual_assignment11

Yuhan_Xu_474154

2019/11/21

Prefix

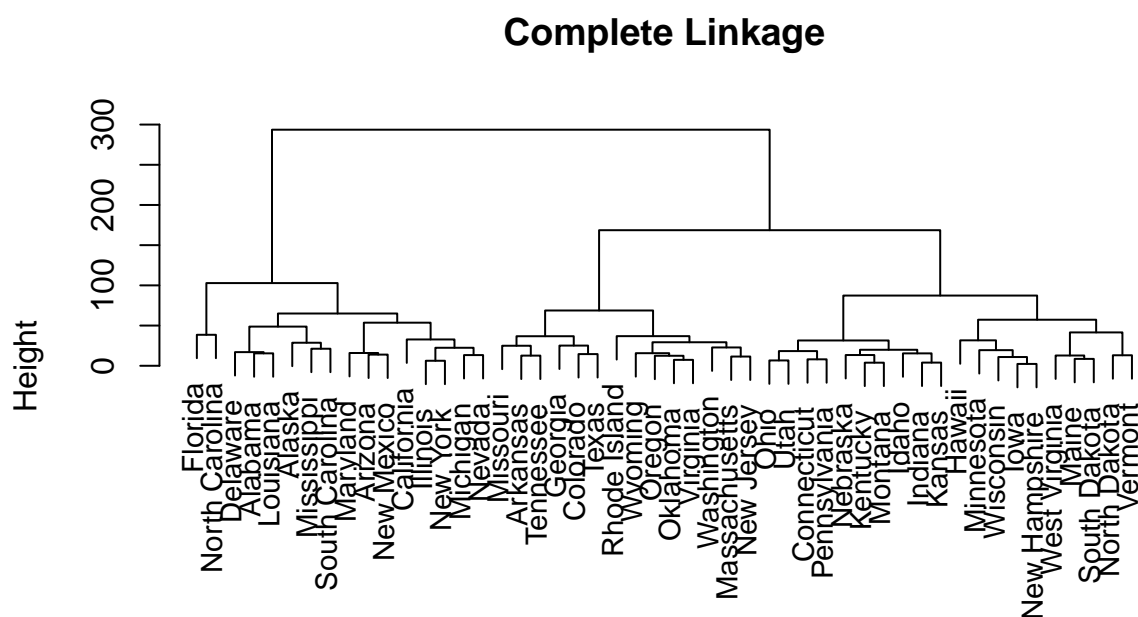
Consider the USArrests data. We will now perform hierarchical clustering on the states.

```
fix(USArrests)
```

(a)

Using hierarchical clustering with complete linkage and Euclidean distance, cluster the states.

```
hc.complete = hclust(dist(USArrests), method = "complete")  
plot(hc.complete, main = "Complete Linkage", xlab = "", sub = "", cex = 0.9)
```



(b)

Cut the dendrogram at a height that results in three distinct clusters. Which states belong to which clusters?

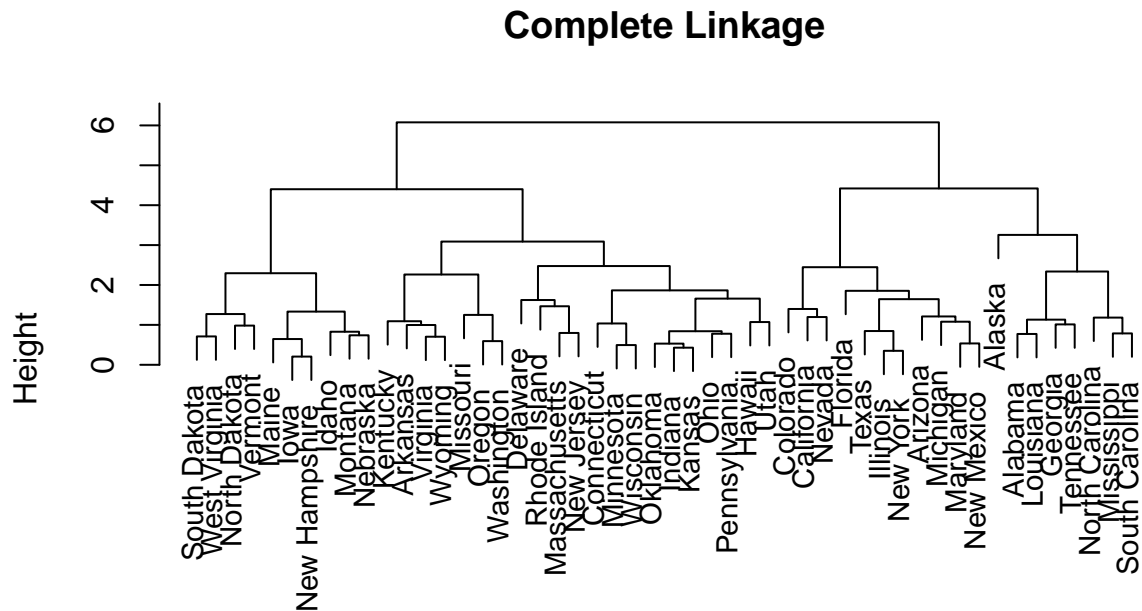
```
cutree(hc.complete,3)
```

##	Alabama	Alaska	Arizona	Arkansas	California
##	1	1	1	2	1
##	Colorado	Connecticut	Delaware	Florida	Georgia
##	2	3	1	1	2
##	Hawaii	Idaho	Illinois	Indiana	Iowa
##	3	3	1	3	3
##	Kansas	Kentucky	Louisiana	Maine	Maryland
##	3	3	1	3	1
##	Massachusetts	Michigan	Minnesota	Mississippi	Missouri
##	2	1	3	1	2
##	Montana	Nebraska	Nevada	New Hampshire	New Jersey
##	3	3	1	3	2
##	New Mexico	New York	North Carolina	North Dakota	Ohio
##	1	1	1	3	3
##	Oklahoma	Oregon	Pennsylvania	Rhode Island	South Carolina
##	2	2	3	2	1
##	South Dakota	Tennessee	Texas	Utah	Vermont
##	3	2	2	3	3
##	Virginia	Washington	West Virginia	Wisconsin	Wyoming
##	2	2	3	3	2

(c)

Hierarchically cluster the states using complete linkage and Euclidean distance, after scaling the variables to have standard deviation one.

```
USArrests.sc = scale(USArrests)
hc.complete1 = hclust(dist(USArrests.sc), method = "complete")
plot(hc.complete1, main="Complete Linkage", xlab = "", sub="", cex=0.9)
```



(d)

What effect does scaling the variables have on the hierarchical clustering obtained? In your opinion, should the variables be scaled before the inter-observation dissimilarities are computed? Provide a justification for your answer.

```
cutree(hc.complete1,3)
```

##	Alabama	Alaska	Arizona	Arkansas	California
##	1	1	2	3	2
##	Colorado	Connecticut	Delaware	Florida	Georgia
##	2	3	3	2	1
##	Hawaii	Idaho	Illinois	Indiana	Iowa
##	3	3	2	3	3
##	Kansas	Kentucky	Louisiana	Maine	Maryland
##	3	3	1	3	2
##	Massachusetts	Michigan	Minnesota	Mississippi	Missouri
##	3	2	3	1	3
##	Montana	Nebraska	Nevada	New Hampshire	New Jersey
##	3	3	2	3	3
##	New Mexico	New York	North Carolina	North Dakota	Ohio
##	2	2	1	3	3
##	Oklahoma	Oregon	Pennsylvania	Rhode Island	South Carolina
##	3	3	3	3	1

```
##      South Dakota      Tennessee      Texas      Utah      Vermont
##           3           1           2           3           3
##      Virginia      Washington      West Virginia      Wisconsin      Wyoming
##           3           3           3           3           3
```

```
table(cutree(hc.complete,3))
```

```
##
##  1  2  3
## 16 14 20
```

```
table(cutree(hc.complete1,3))
```

```
##
##  1  2  3
##  8 11 31
```

Scaling the variables will result in different clusters.

From the summary of this dataset, we know that different variables have different magnitudes. So, I think the variables should be scaled to avoid those variables with bigger magnitude dominating in clustering.

```
summary(USArrests)
```

```
##      Murder      Assault      UrbanPop      Rape
##  Min.   : 0.800   Min.    : 45.0   Min.    :32.00   Min.    : 7.30
##  1st Qu.: 4.075   1st Qu.:109.0   1st Qu.:54.50   1st Qu.:15.07
##  Median : 7.250   Median :159.0   Median :66.00   Median :20.10
##  Mean   : 7.788   Mean    :170.8   Mean    :65.54   Mean    :21.23
##  3rd Qu.:11.250   3rd Qu.:249.0   3rd Qu.:77.75   3rd Qu.:26.18
##  Max.   :17.400   Max.    :337.0   Max.    :91.00   Max.    :46.00
```