

Individual Assignment 11

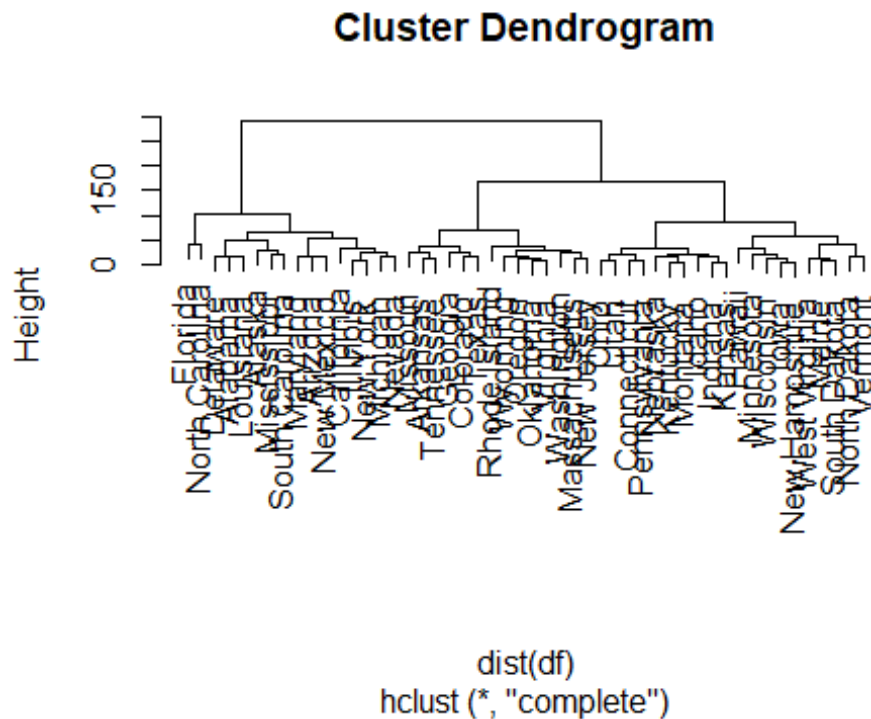
Charlie Ling

2021/12/2

Exercise 10.7: Problem 9 ISLR p.417

9. Consider the USArrests data. We will now perform hierarchical clustering on the states.
 - (a) Using hierarchical clustering with complete linkage and Euclidean distance, cluster the states.

```
library(ISLR)
df=na.omit(USArrests)
hc.complete = hclust(dist(df), method = "complete")
plot(hc.complete)
```



- (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.

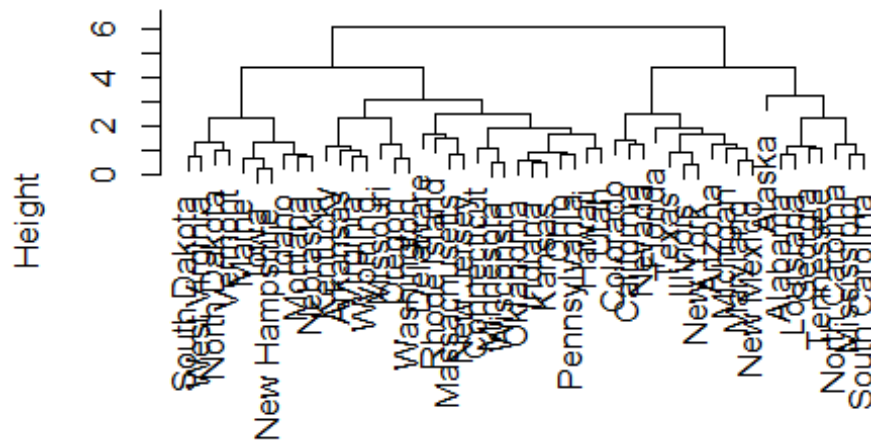
```
library(caret)

## 载入需要的程辑包: ggplot2

## 载入需要的程辑包: lattice

norm.values = preProcess(df, method = c("center", "scale"))
df.norm = predict(norm.values, df)
hc.complete = hclust(dist(df.norm), method = "complete")
plot(hc.complete)
```

Cluster Dendrogram



```
dist(df.norm)
hclust (*, "complete")
```

```
cutree(hc.complete,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

- (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.

#Scaling makes the cluster more evenly distributed and into more groups. As for this case, I do not think the variables should be scaled, because the scaled data are clustered into more than 3 groups. And the absolute value of crimes instead of scaled value is more useful to define the crime situation of a state.