

Clustering.R

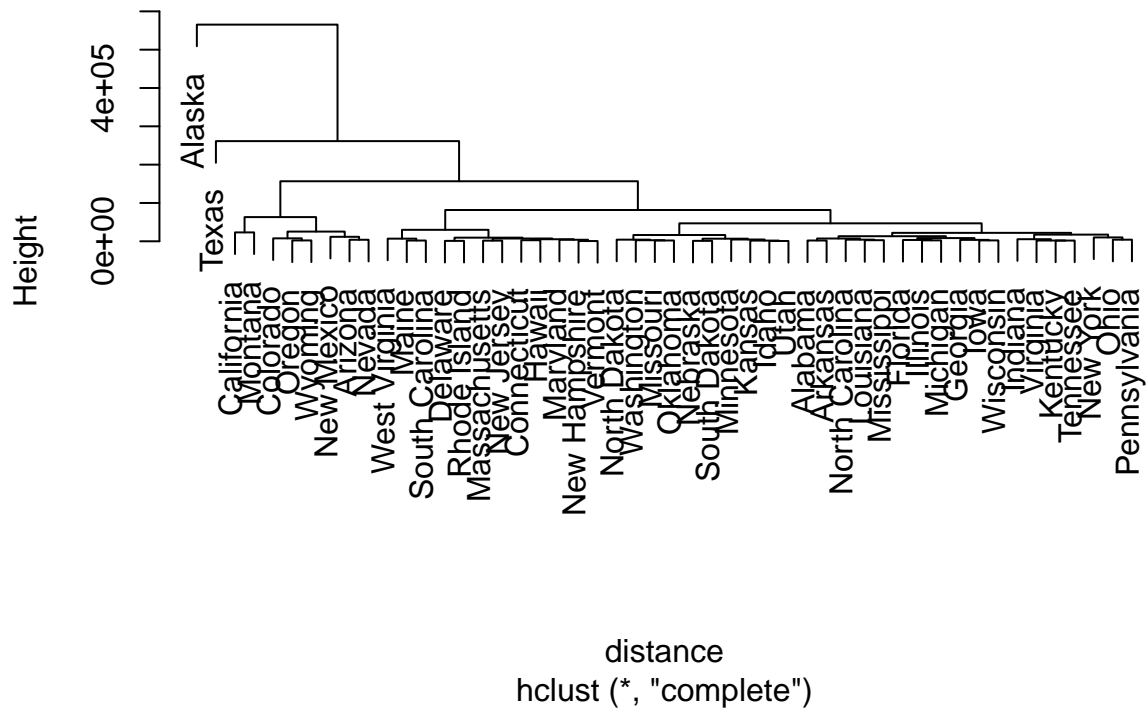
jasho

2020-03-14

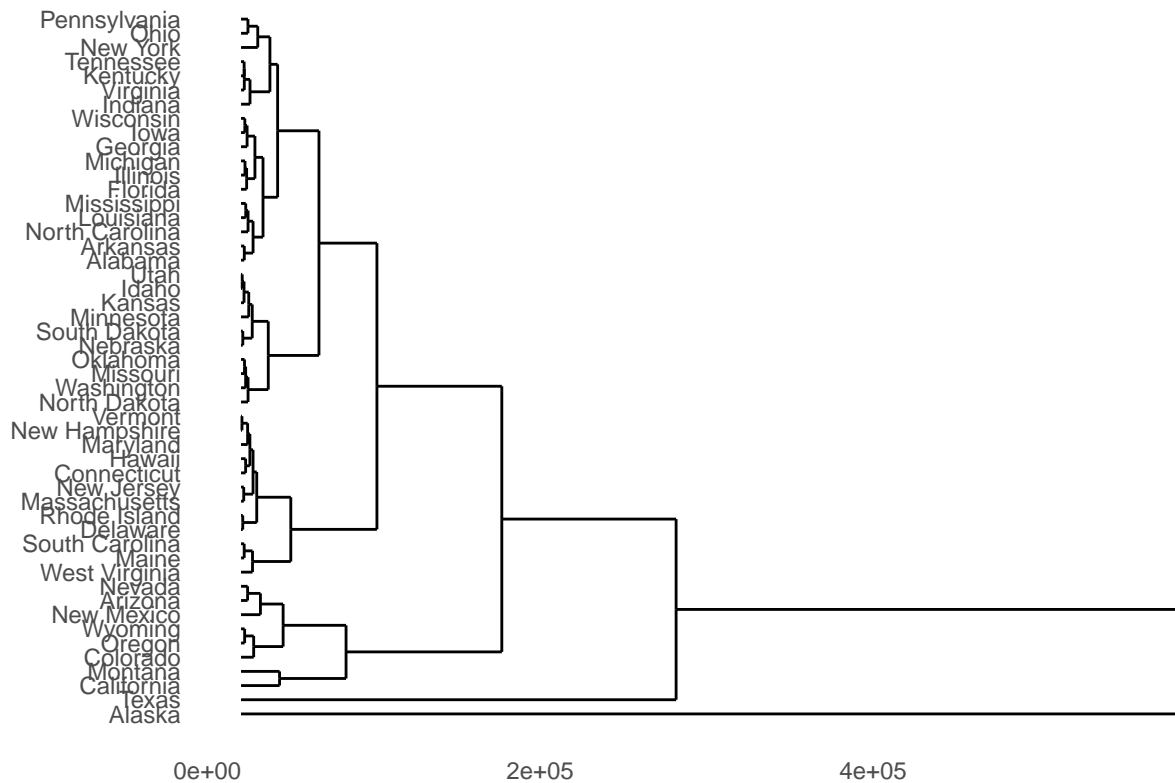
```
pacman::p_load(datasets, tidyverse, ggdendro, pander, cluster)
dat = state.x77

distance <- dist(as.matrix(dat))
hc <- hclust(distance)
plot(hc)
```

Cluster Dendrogram

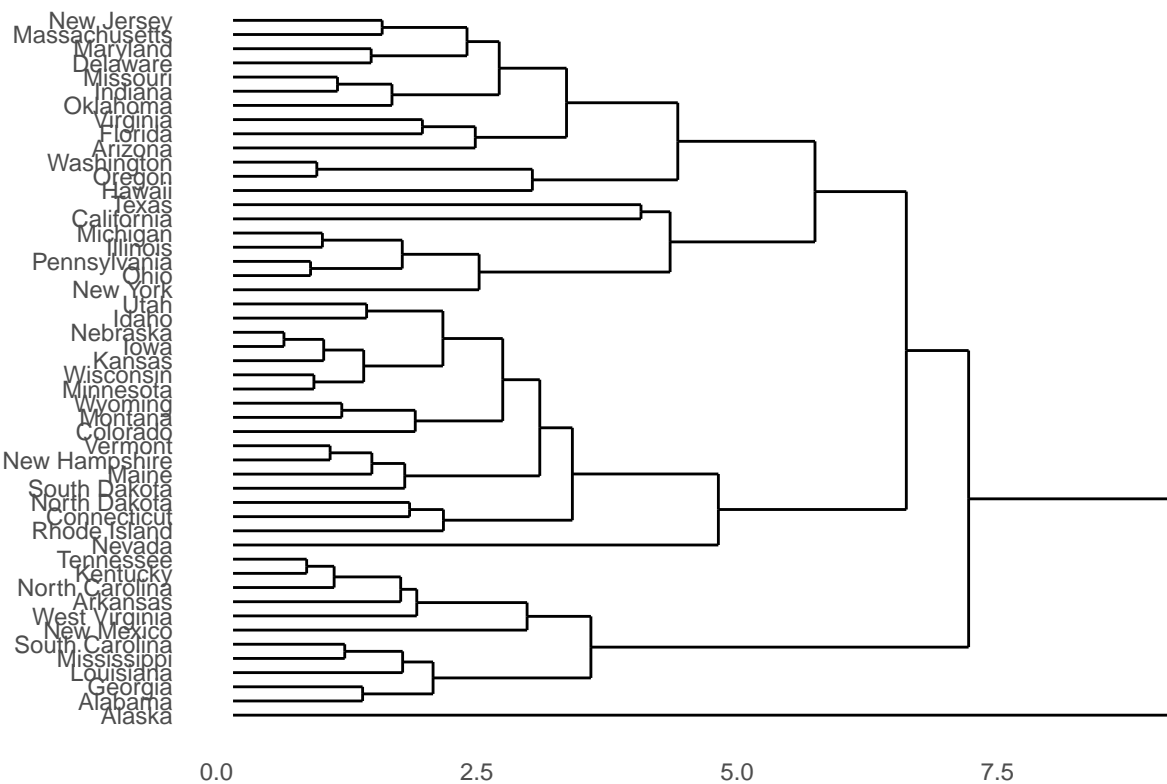


```
ggdendro::ggdendrogram(hc, rotate = 1)
```



*#Scaling eliminated the effects of population
#and area that overpowered the other variables.
#this allowed the population to have a similar
#effect to area. i.e cali and texas being the
#closest to each other.*

```
datasc <- as.data.frame(scale(dat))
scdist <- dist(as.matrix(datasc))
hcsc <- hclust(scdist)
ggdendro::ggdendrogram(hcsc, rotate = 1)
```

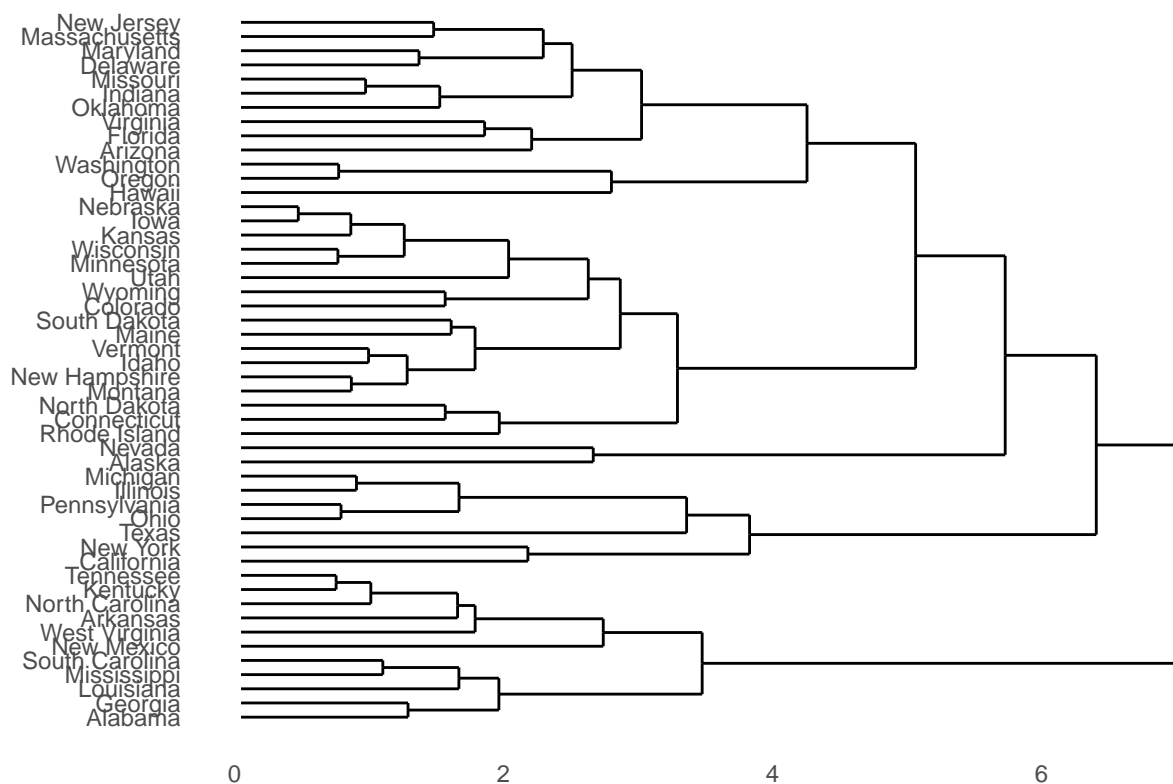


*#Taking out area really starts to show the
#important stuff like HS grad and literacy,
#these
#things were not nearly as important when area
#was involved*

```

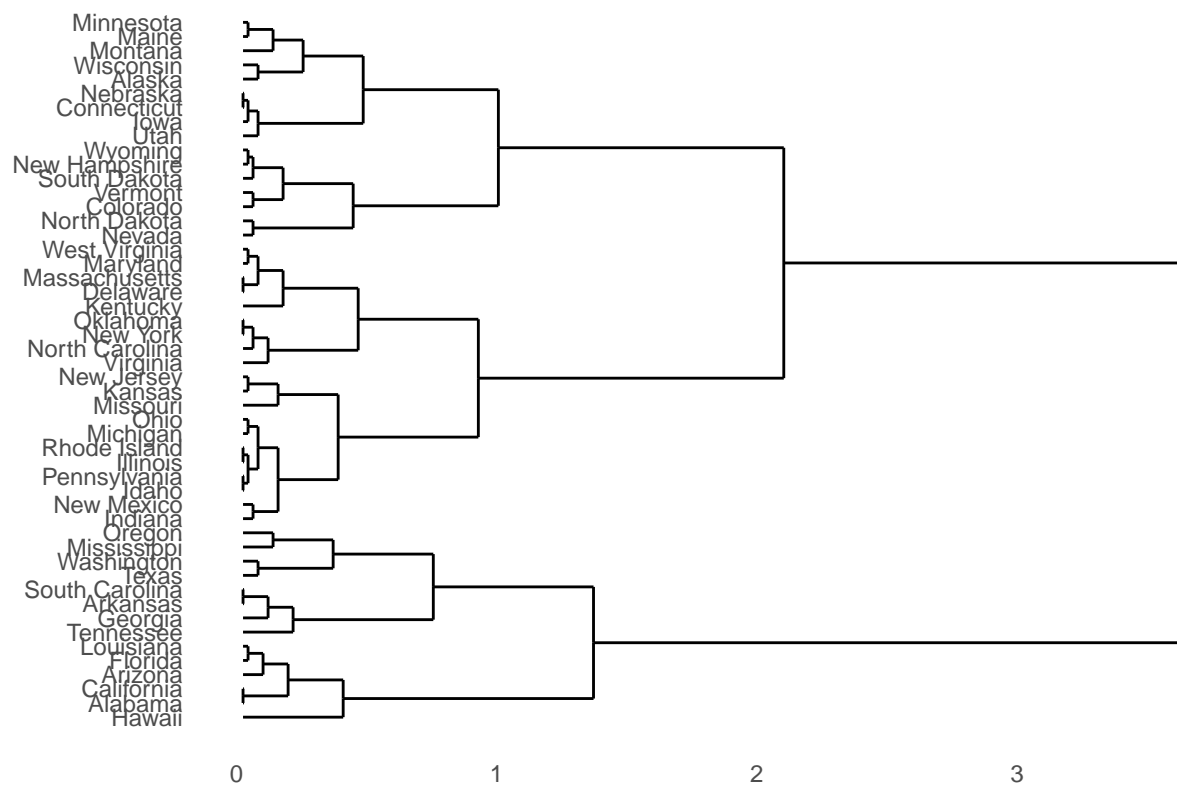
namos <- dat[,0]
datAsc <- datsc %>%
  select(-Area)
Ascdist <- dist(as.matrix(datAsc))
Ahc <- hclust(Ascdist)
ggdendro::ggdendrogram(Ahc, rotate = 1)

```



*#Frost only is super cool, and since I lived in
#Texas and Washington can confirm that Frost*

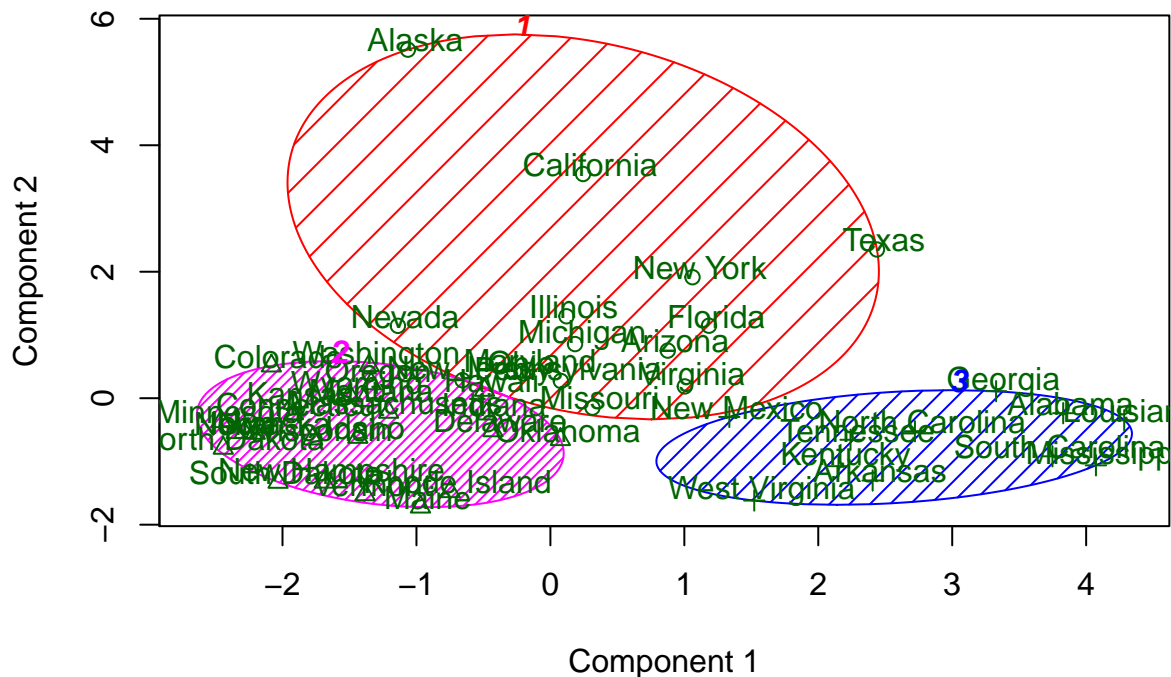
```
Fdat <- datsc %>%
  select(Frost)
Fdist <- dist(as.matrix(Fdat))
Fhc <- hclust(Fdist)
ggdendrogram(Fhc, rotate = 1)
```



#K-means

```
clus1 <- kmeans(datasc, 3)
clusplot(datasc, clus1$cluster, color = TRUE, shade = TRUE, labels = 2, lines = 0)
```

CLUSPLOT(datsc)

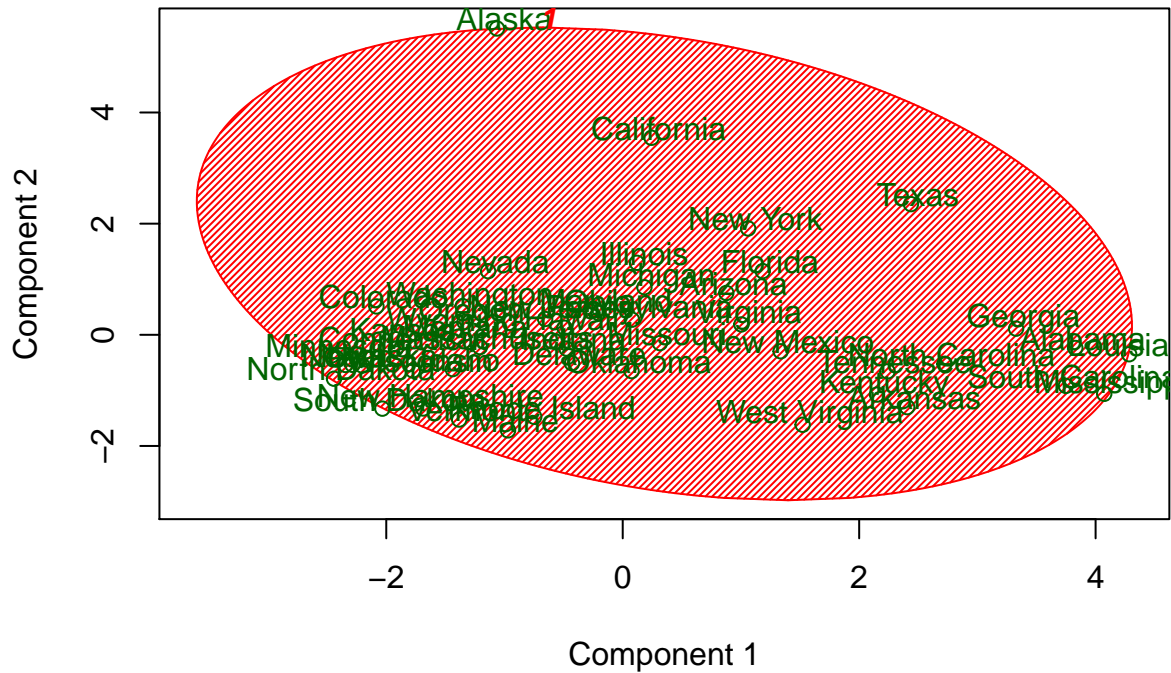


These two components explain 65.39 % of the point variability.

```
#large population/large area is the dominating
#factor for the clusters

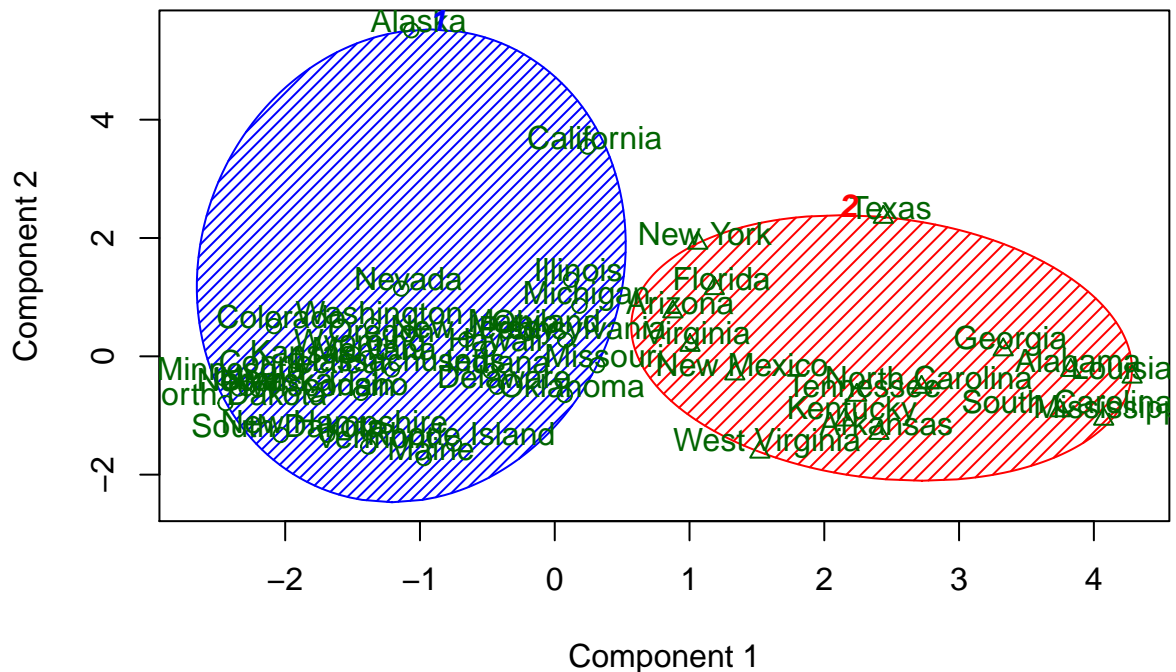
for (i in 1:25) {
  thisclus <- kmeans(datasc, i)
  clusplot(datasc, thisclus$cluster, color = TRUE, shade = TRUE, labels = 2, lines = 0)
}
```

CLUSPLOT(datsc)



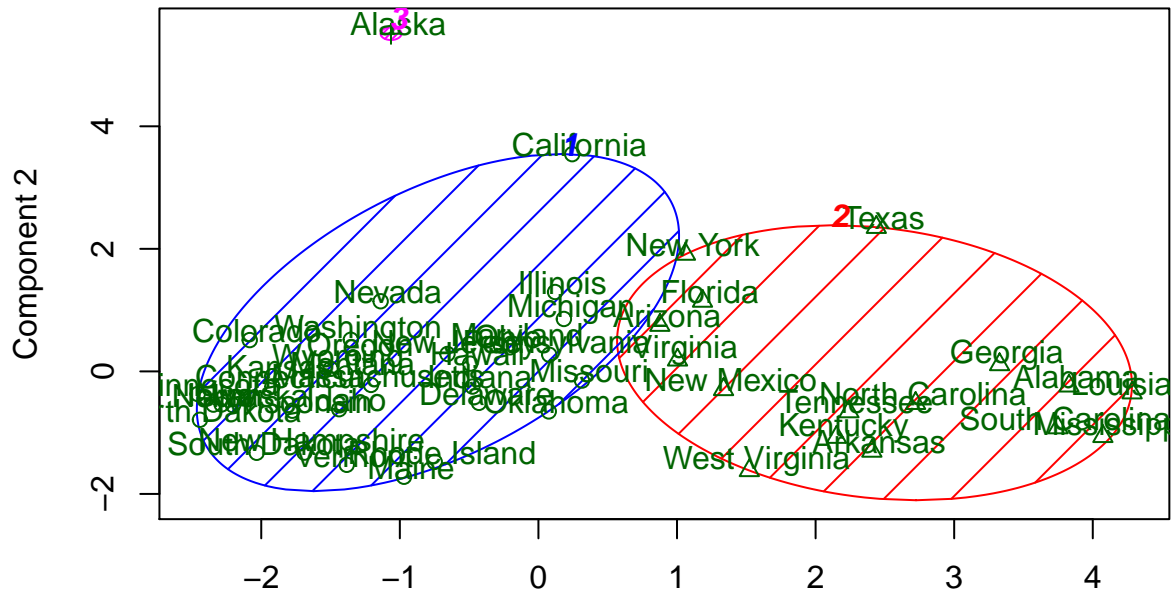
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



These two components explain 65.39 % of the point variability.

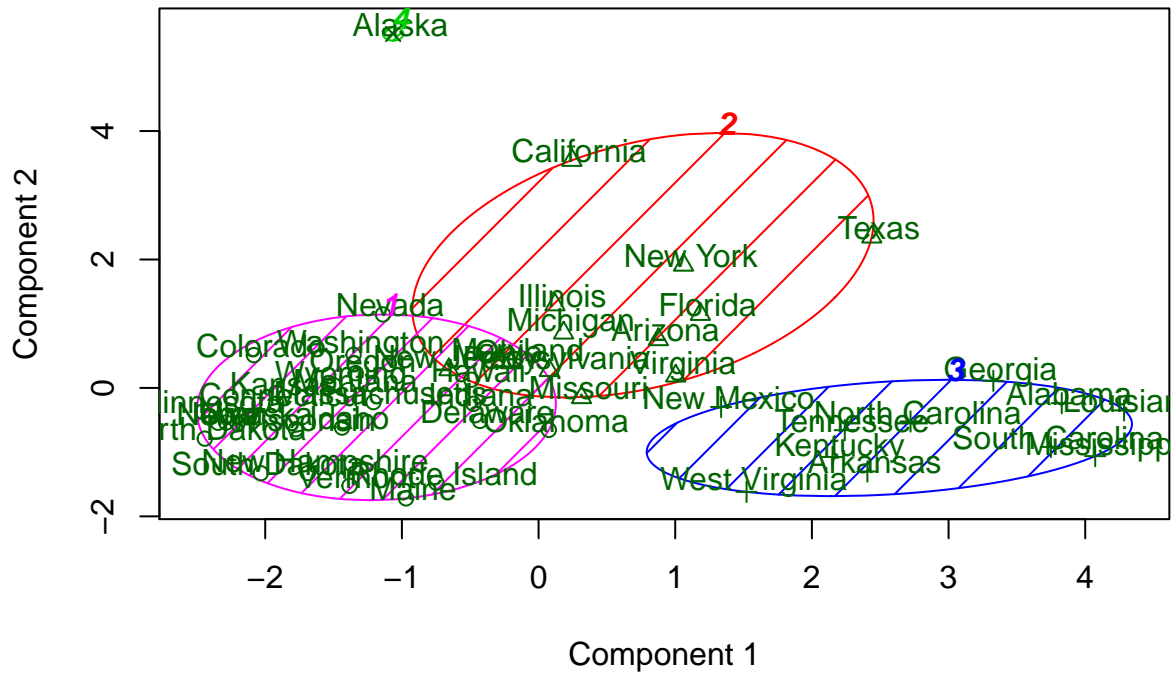
CLUSPLOT(datsc)



Component 1

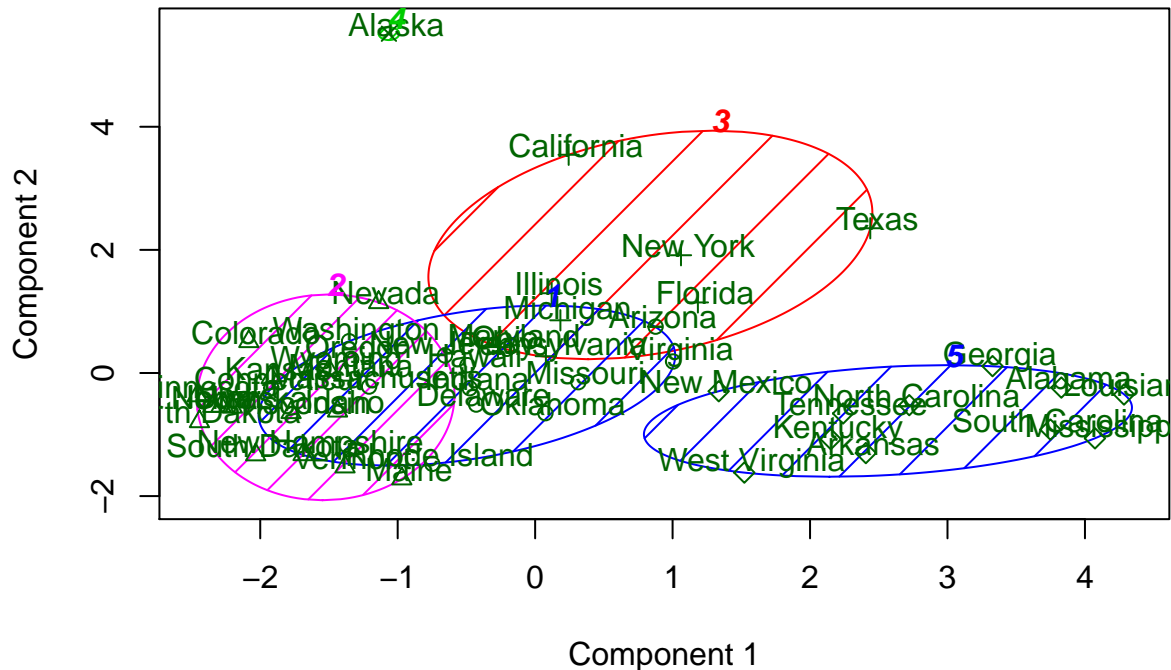
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

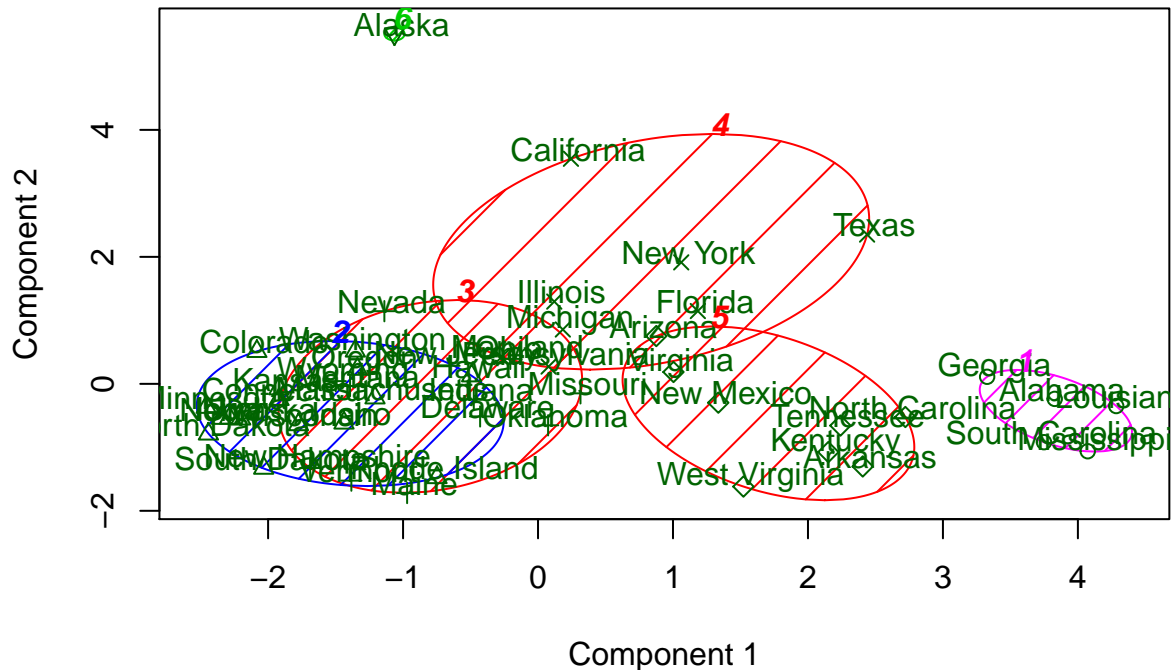


These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

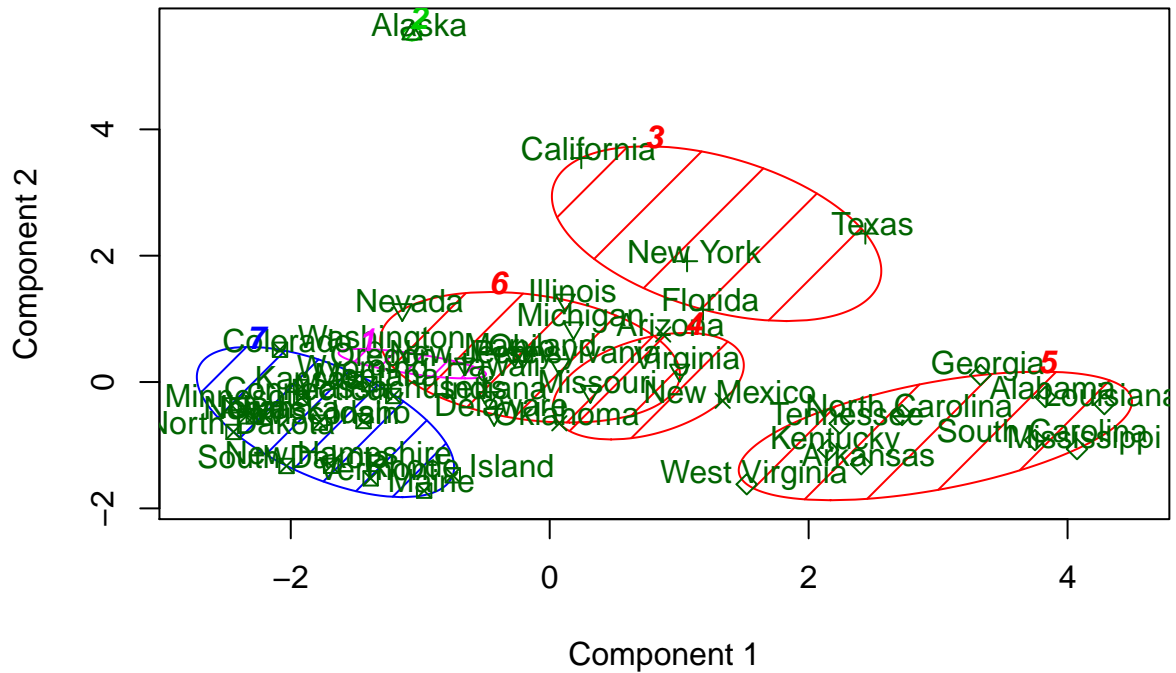


CLUSPLOT(datsc)



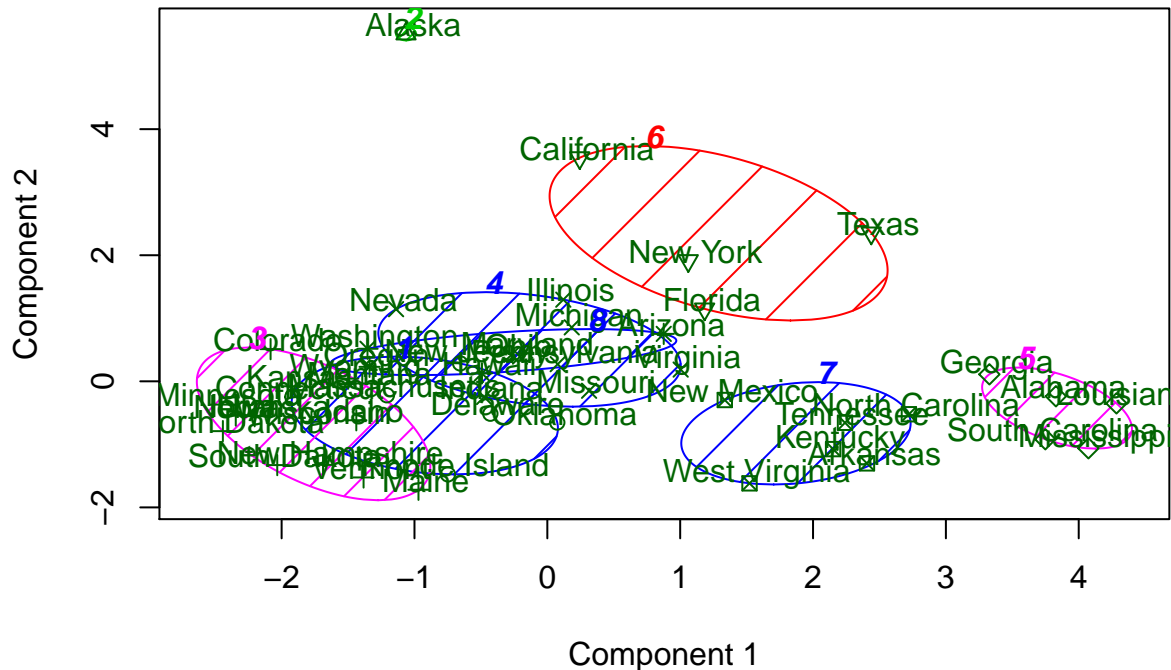
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



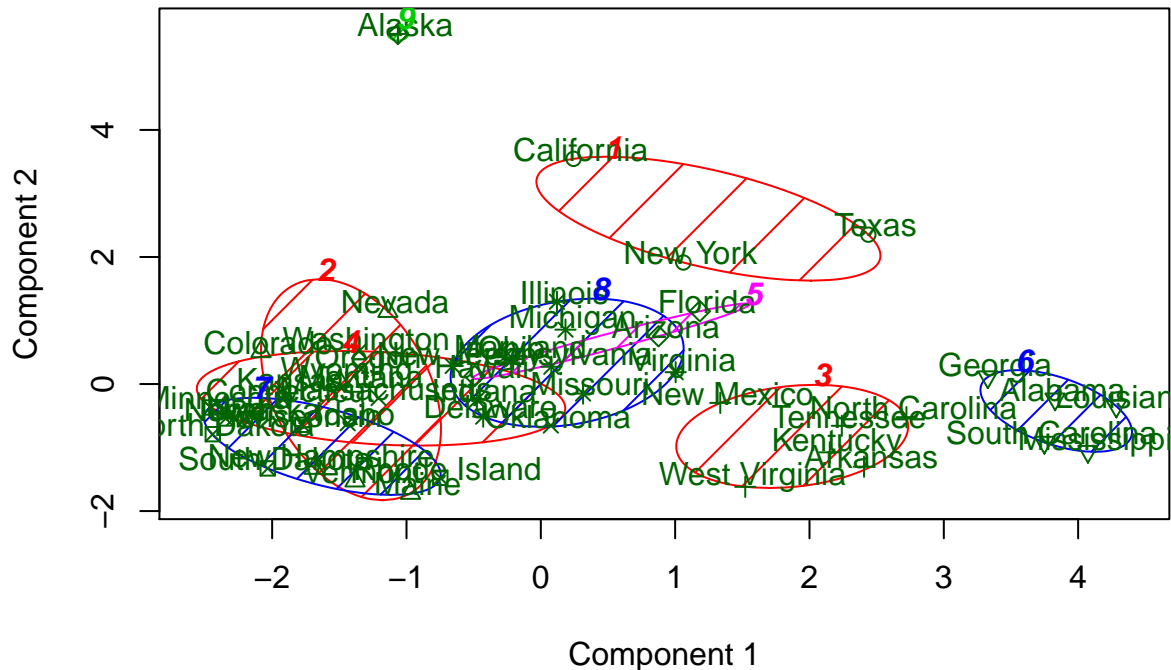
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



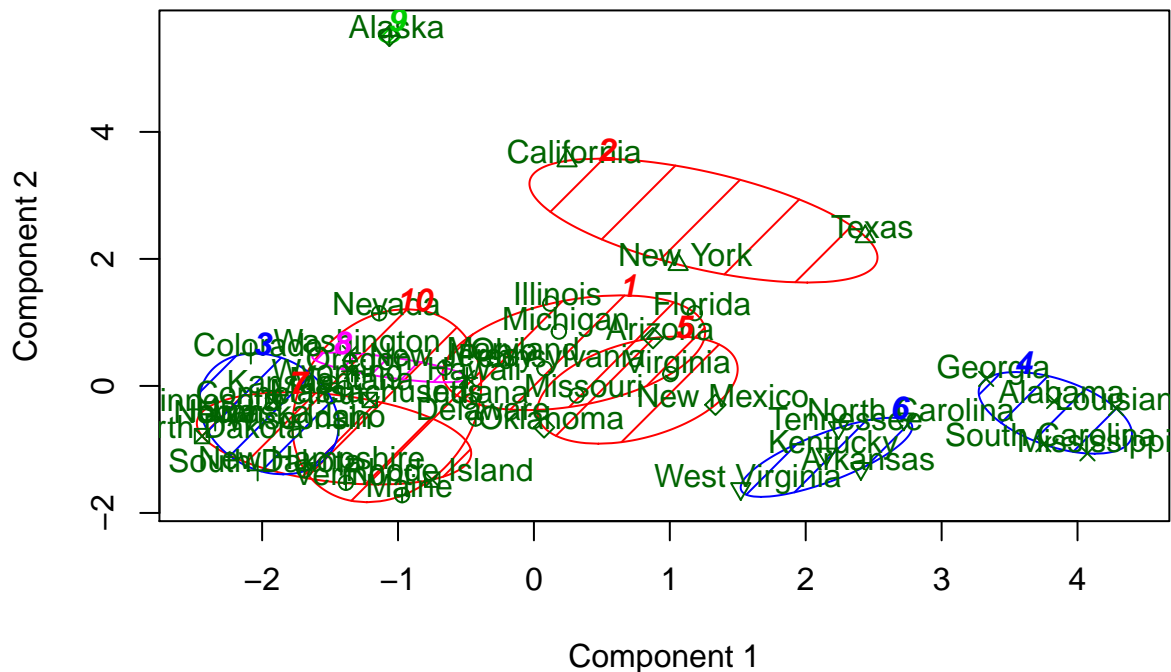
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



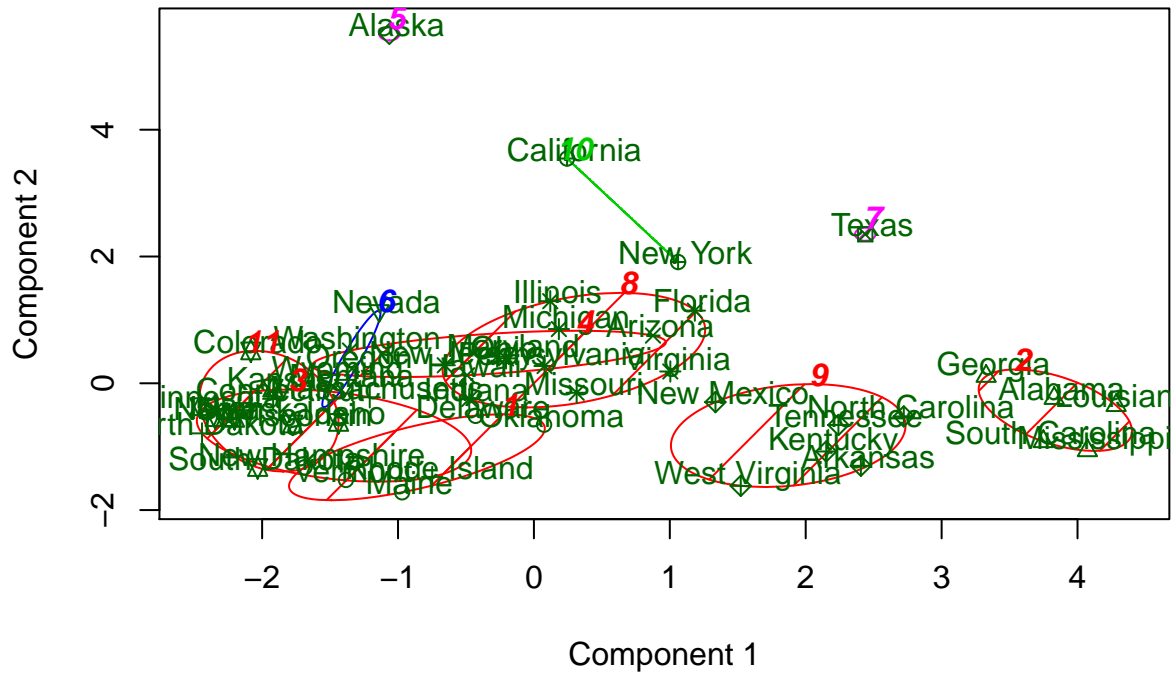
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



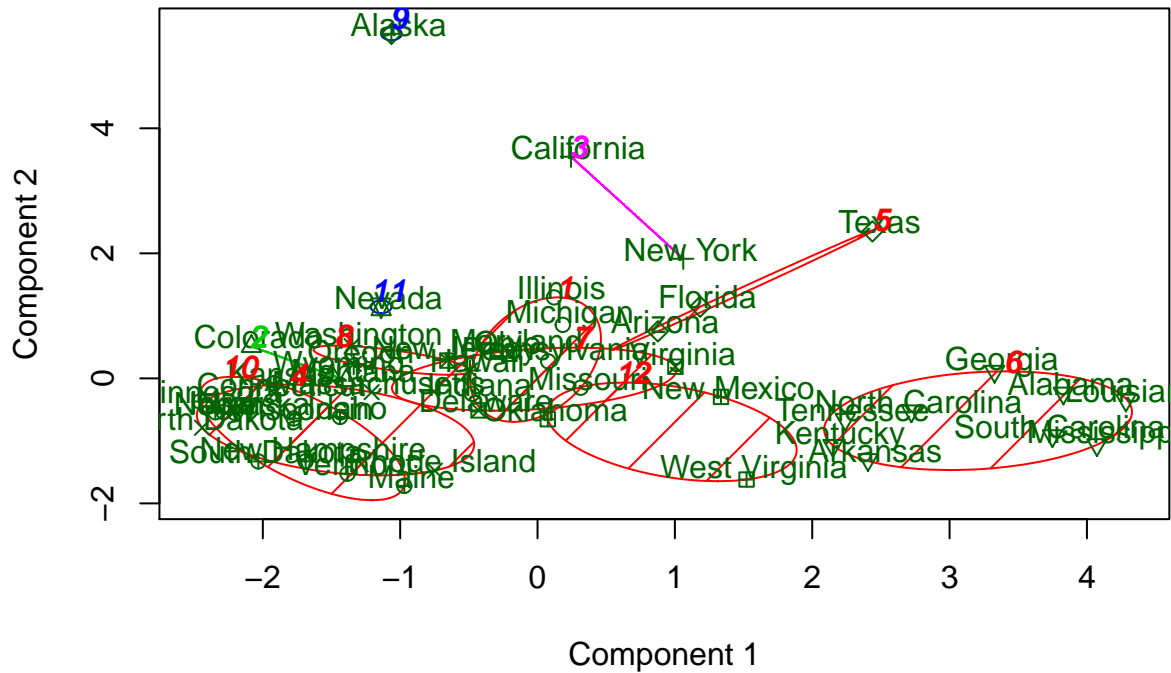
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



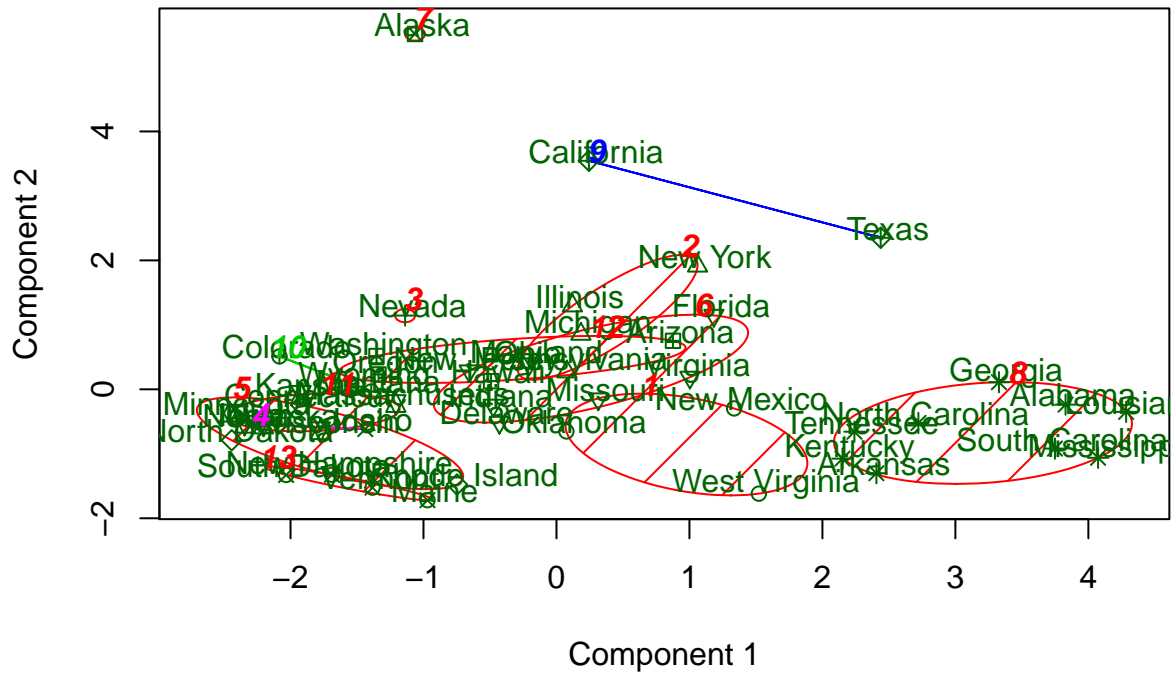
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



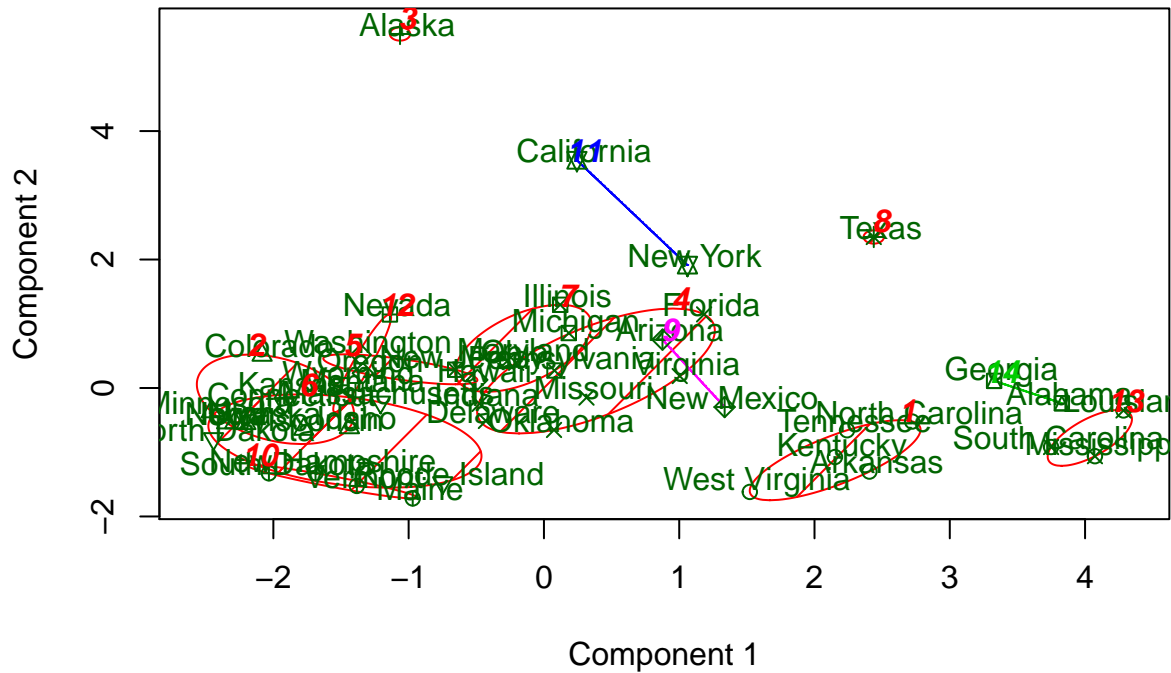
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



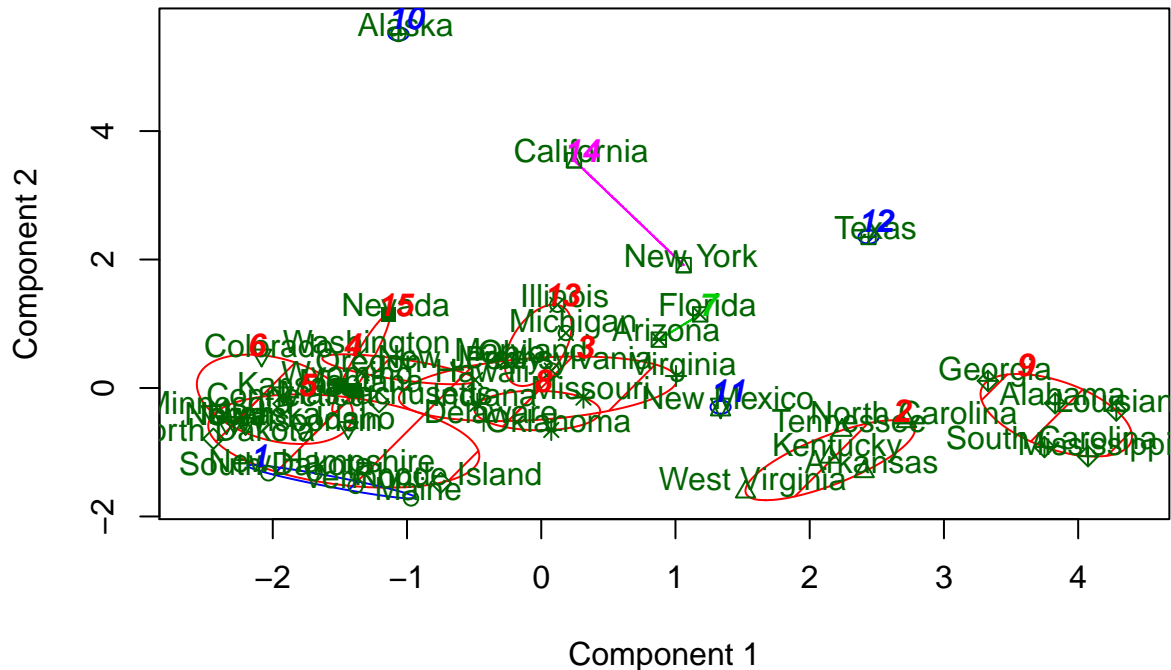
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



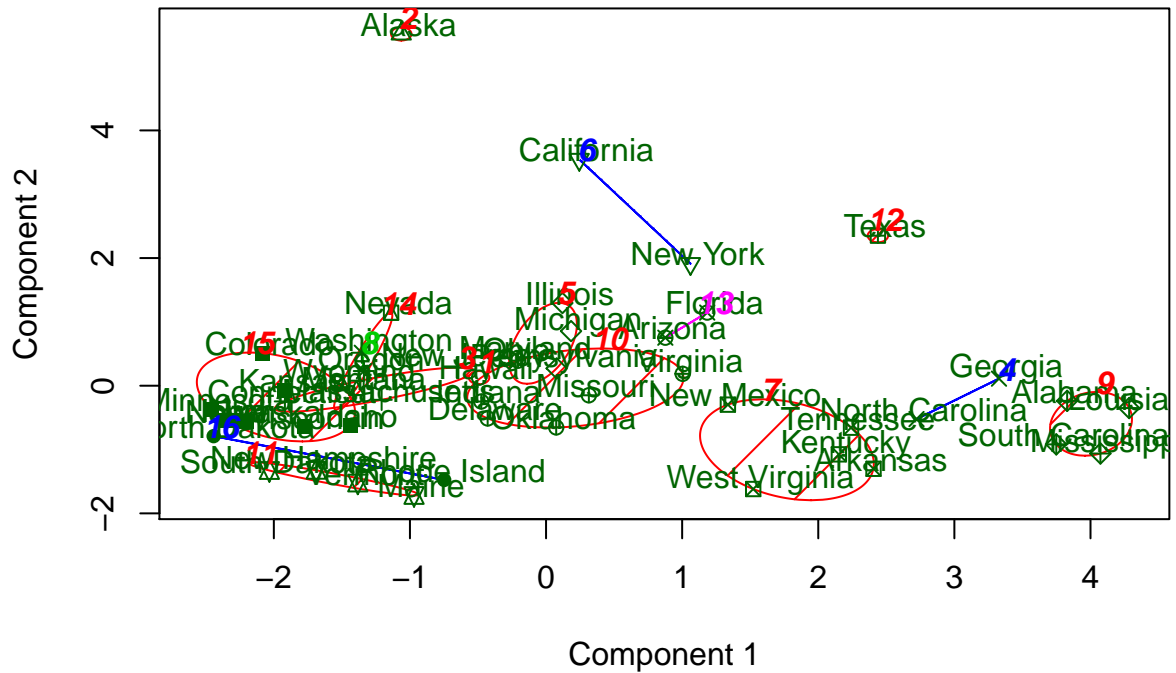
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



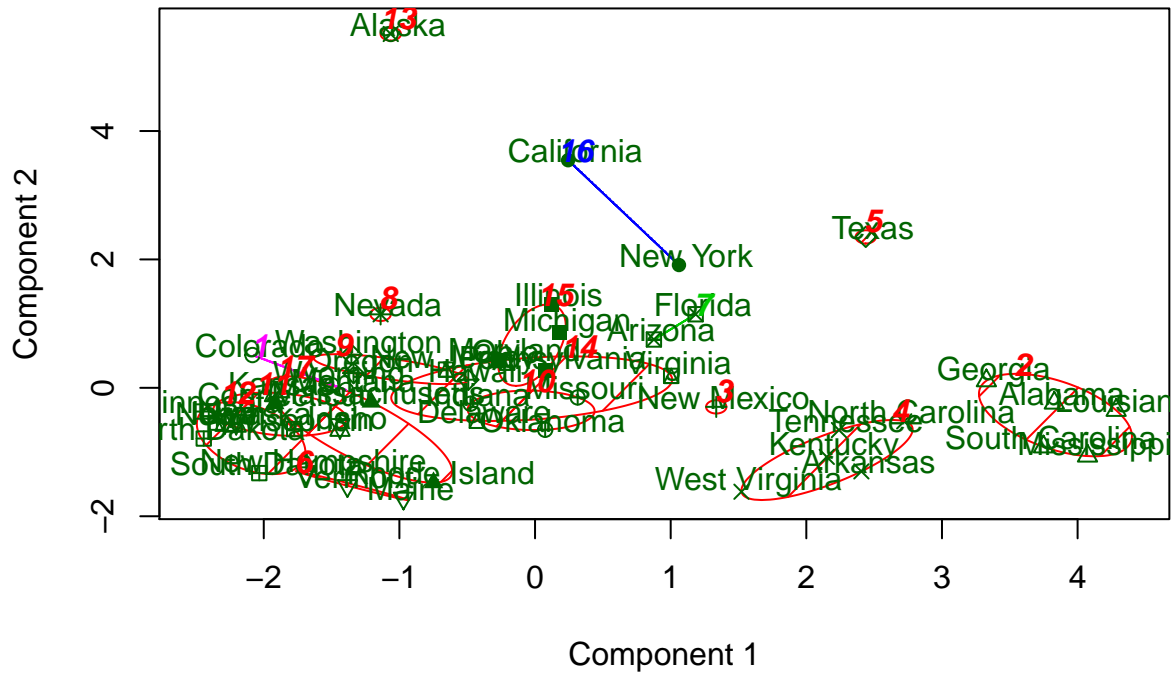
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

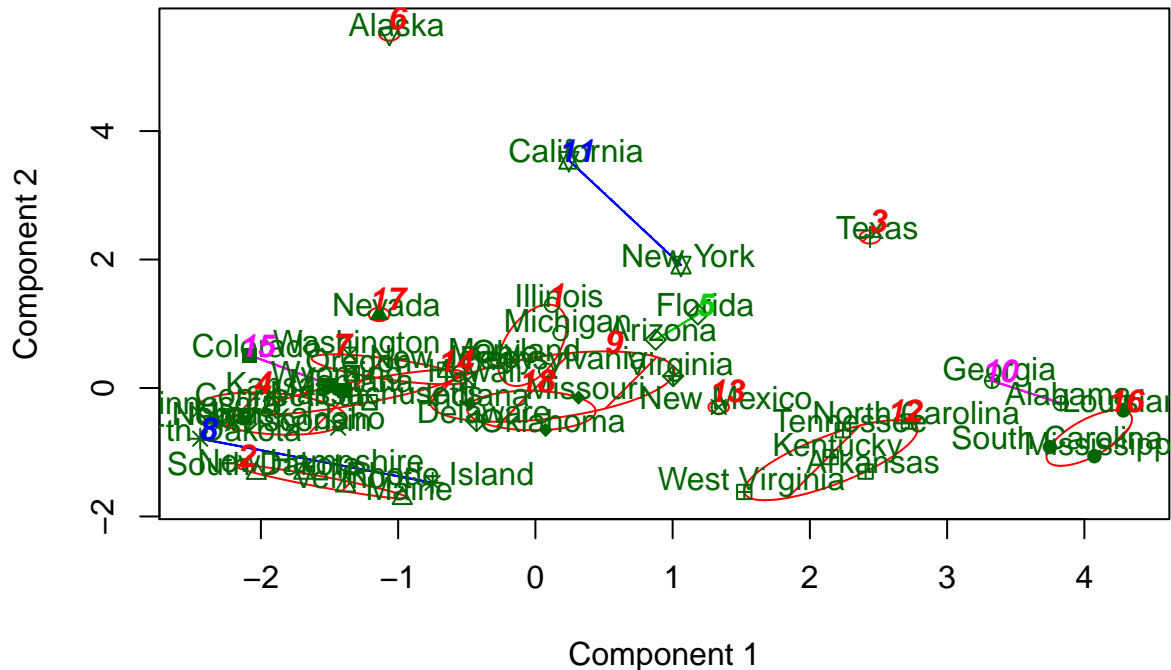


These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

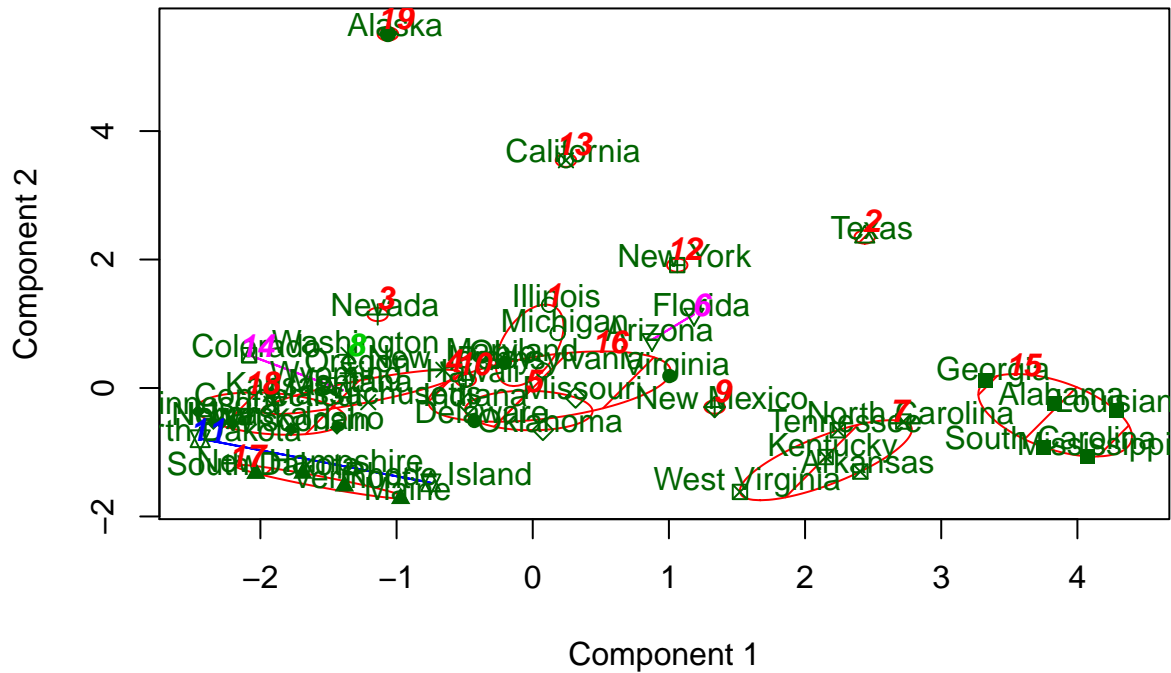


CLUSPLOT(datsc)

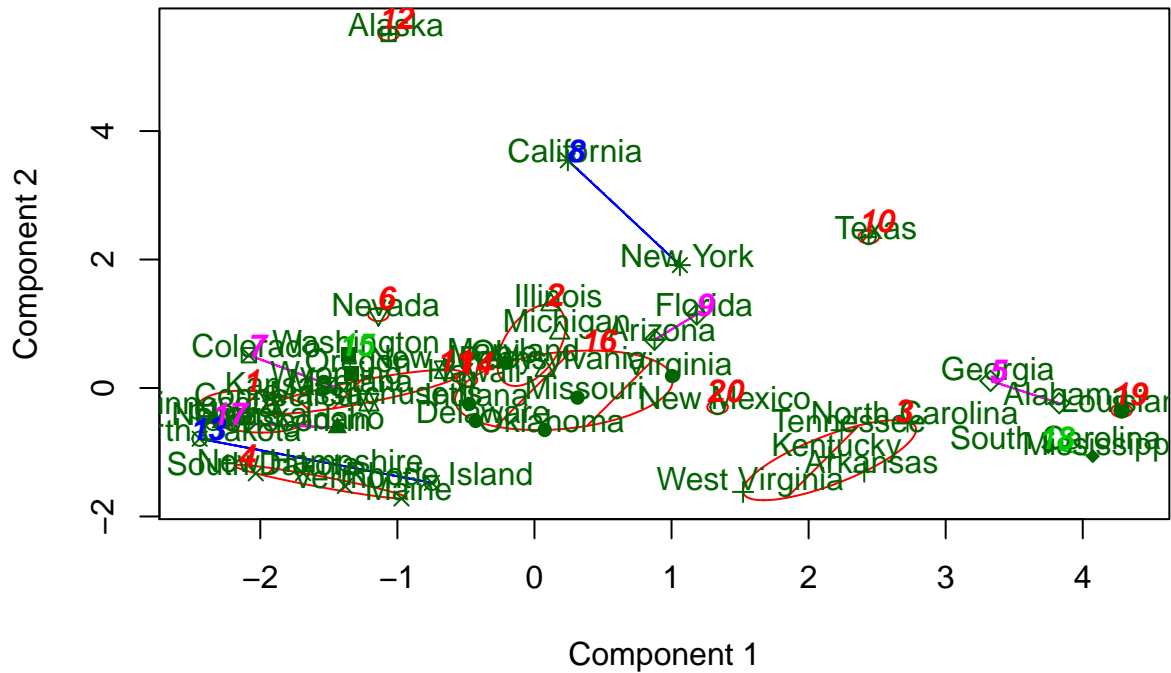


These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

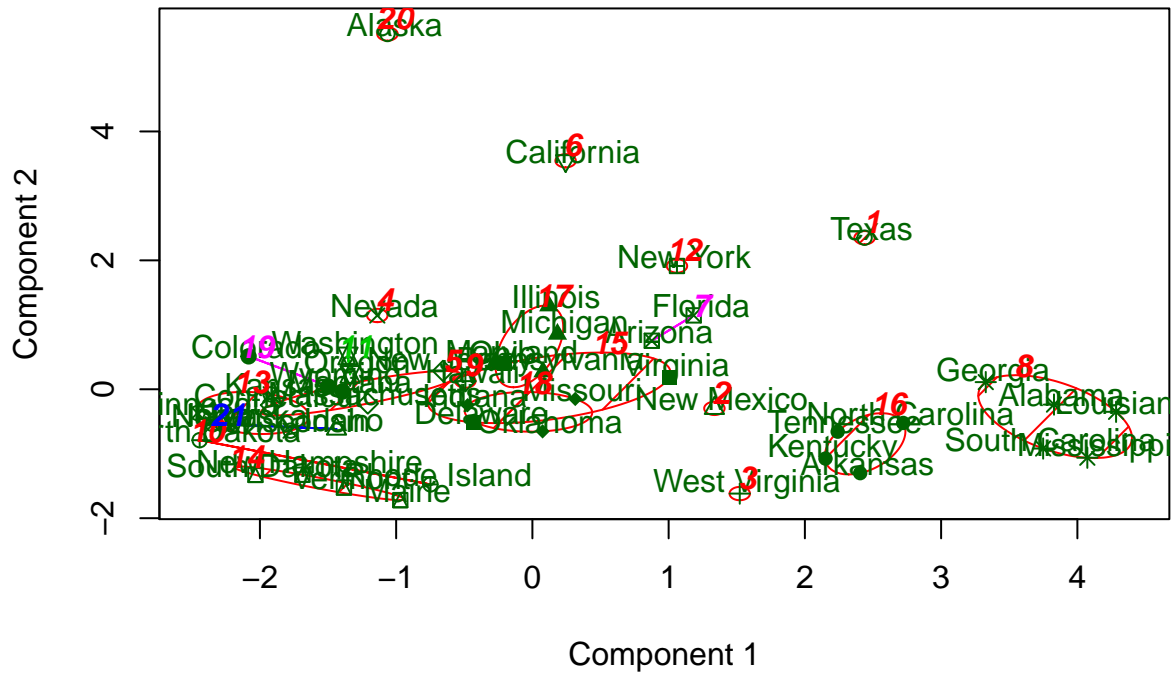


CLUSPLOT(datsc)

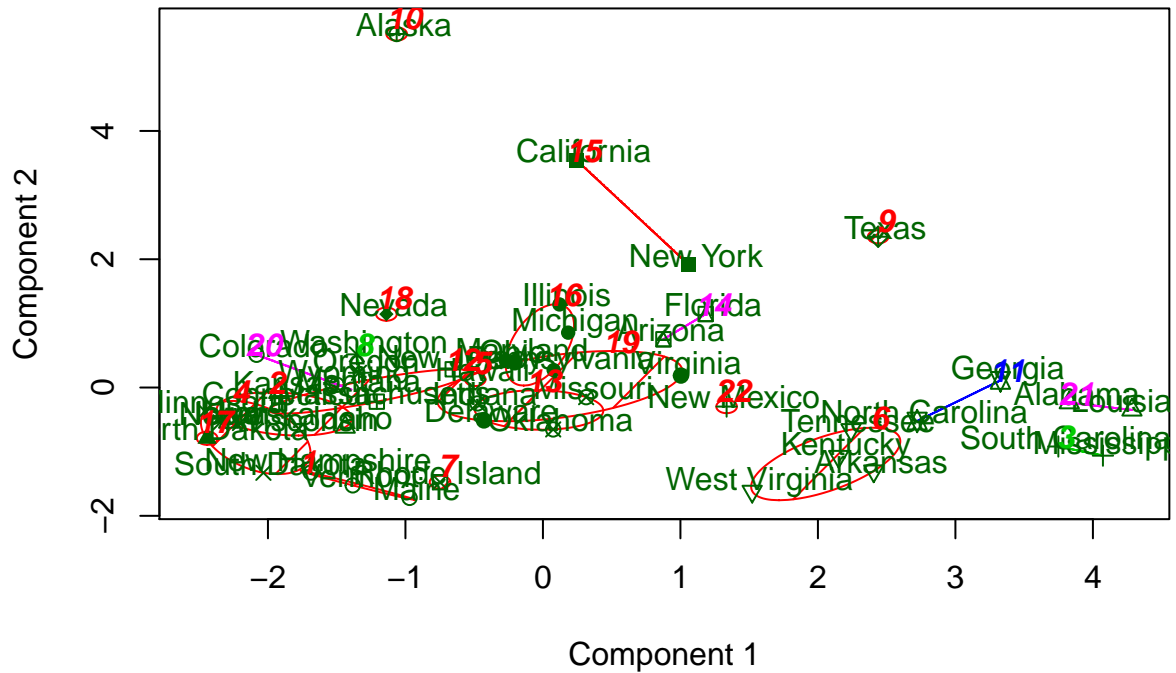


These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

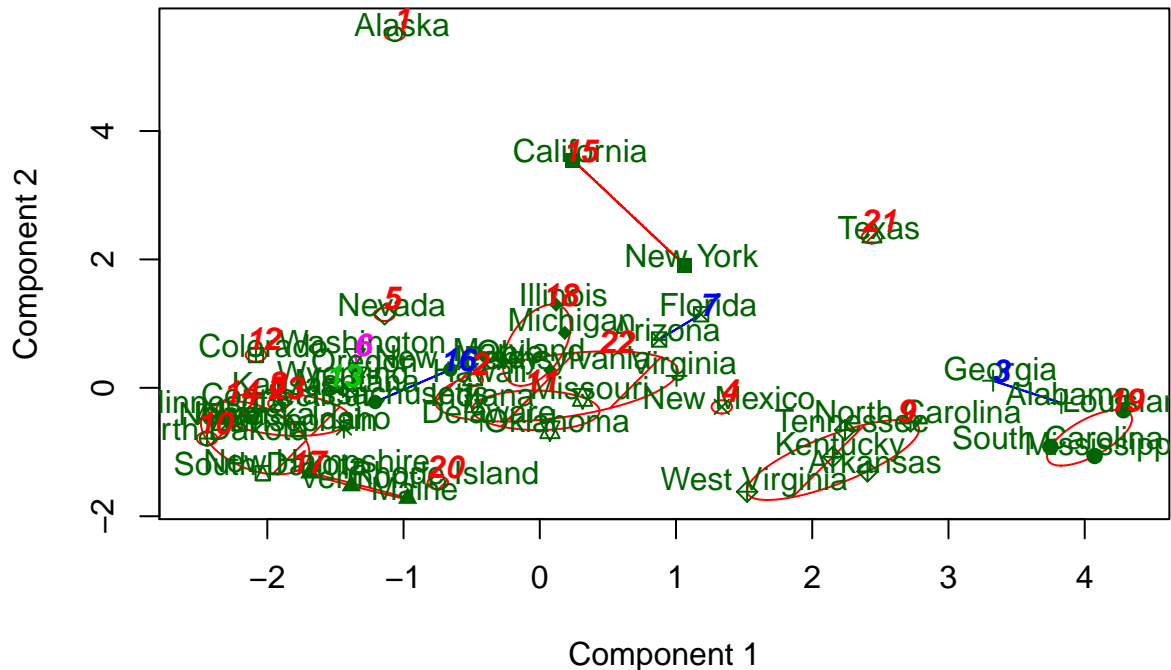


CLUSPLOT(datsc)



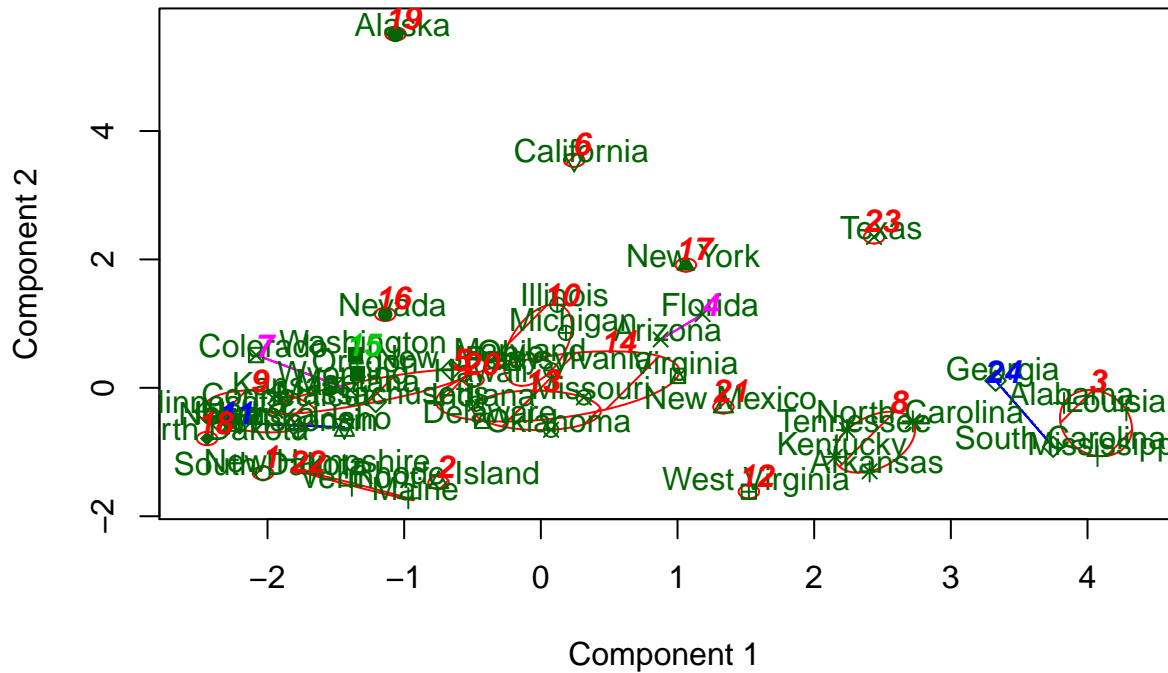
These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)

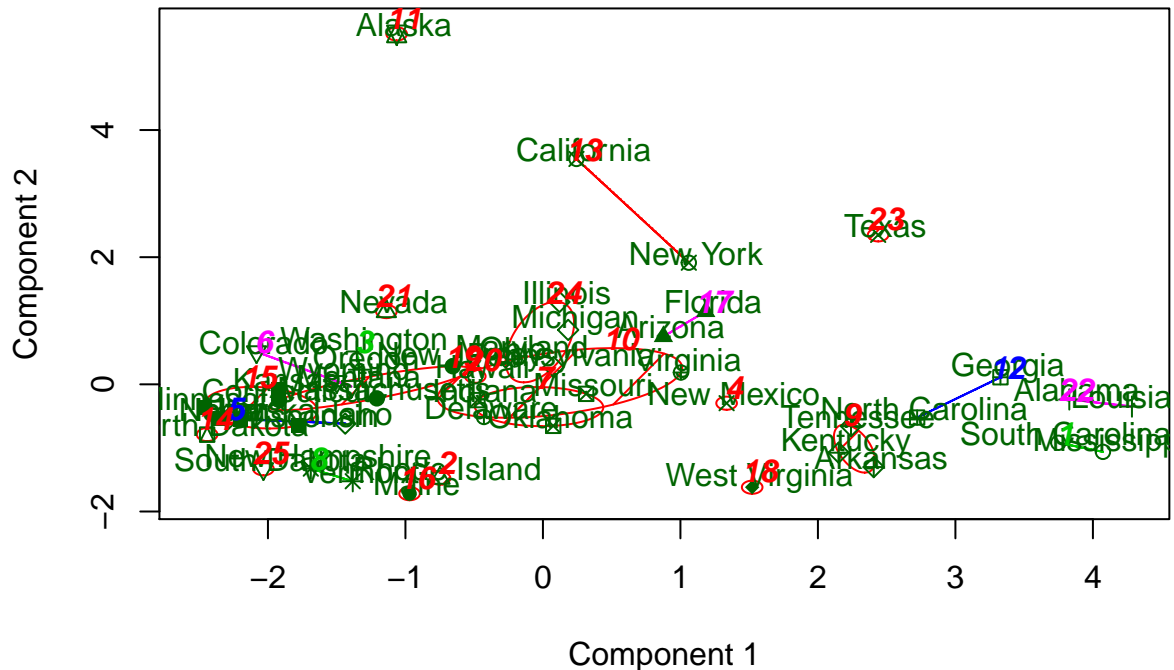


These two components explain 65.39 % of the point variability.

CLUSPLOT(datsc)



CLUSPLOT(datsc)

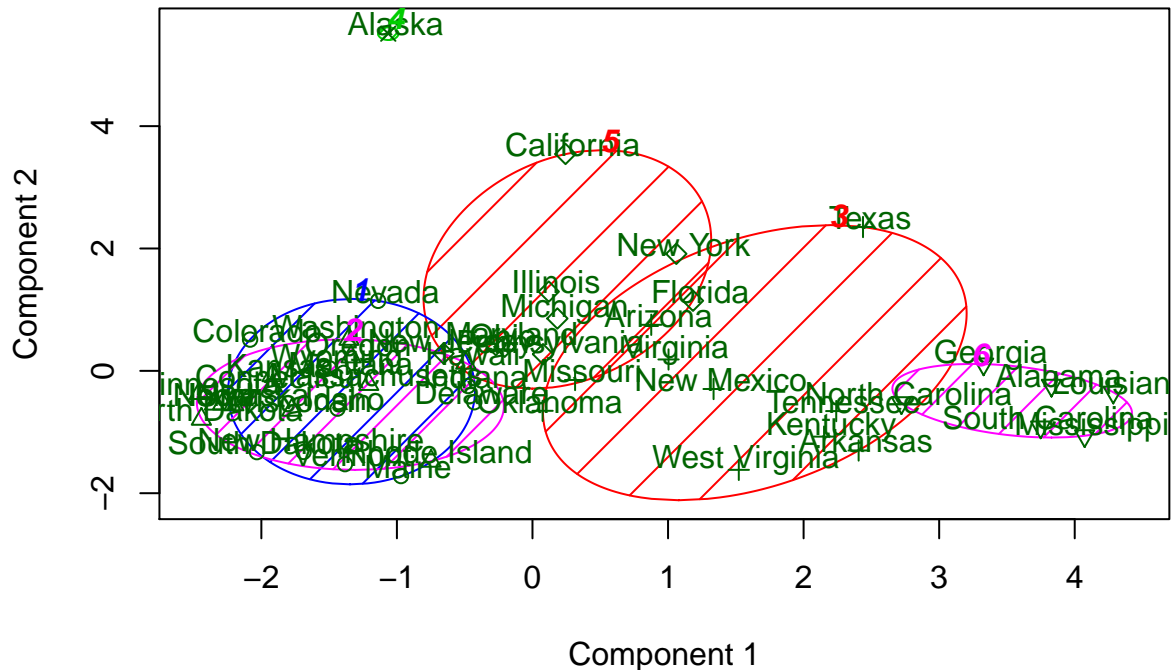


These two components explain 65.39 % of the point variability.

#I like 6

```
clus6 <- kmeans(datasc, 6)
clusplot(datasc, clus6$cluster, color = TRUE, shade = TRUE, labels = 2, lines = 0)
```

CLUSPLOT(datsc)



```
pander(clus6$cluster)
```

Table 1: Table continues below

Alabama	Alaska	Arizona	Arkansas	California	Colorado	Connecticut
6	4	3	3	5	1	2

Table 2: Table continues below

Delaware	Florida	Georgia	Hawaii	Idaho	Illinois	Indiana	Iowa
1	5	6	2	1	5	1	2

Table 3: Table continues below

Kansas	Kentucky	Louisiana	Maine	Maryland	Massachusetts	Michigan
2	3	6	1	5	2	5

Table 4: Table continues below

Minnesota	Mississippi	Missouri	Montana	Nebraska	Nevada
2	6	3	1	2	1

Table 5: Table continues below

New Hampshire	New Jersey	New Mexico	New York	North Carolina
1	5	3	5	6

Table 6: Table continues below

North Dakota	Ohio	Oklahoma	Oregon	Pennsylvania	Rhode Island
2	5	3	2	5	2

Table 7: Table continues below

South Carolina	South Dakota	Tennessee	Texas	Utah	Vermont	Virginia
6	1	3	3	2	1	3

Washington	West Virginia	Wisconsin	Wyoming
2	3	2	1

*#New york texas and cali are clearly grouped
 #due to poulation. Alaska is alone due to area,
 #A lot of the middle sized states are grouoed
 #together in the bottom right, and these ultr
 #small states are all in the bottom left,*