

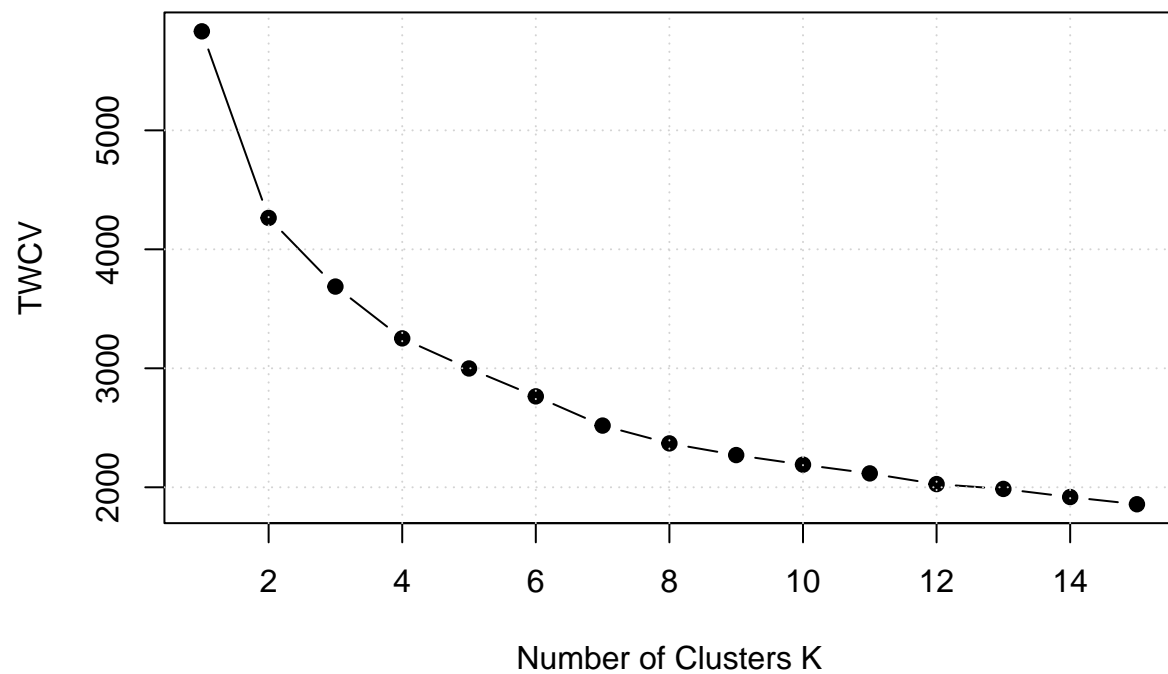
# Siyu 535 hw4

Siyu Wu | USC ID: 2350-0175-97

```
library(readxl)
library(tibble)
df0 = read_excel("cities1(1).xlsx",sheet="Data")
df0$Crime_Trend = NULL
df0$Unemployment_Threat = NULL
df = column_to_rownames(df0,"Metropolitan_Area")
df = scale(df)
```

```
df_prep = df0
df_prep = column_to_rownames(df_prep,"Metropolitan_Area")
df_mm = data.frame(apply(df_prep,2,function(x) (x-min(x))/(max(x)-min(x))))
```

```
set.seed(123)
twcv = function(k) kmeans(df,k,nstart=10)$tot.withinss
k = 1:15
twcv_values = sapply(k,twcv)
plot(k,twcv_values,type="b",pch=19,
      xlab="Number of Clusters K", ylab="TWCV")
grid()
```



```
set.seed(123)
library(cluster)
library(factoextra)
```

```
## Loading required package: ggplot2
```

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

```
final = kmeans(df,centers=4,nstart=25)
fviz_cluster(final,data=df,geom = "point")
```



```
cluster_number = as.factor(final$cluster)
m1 = prcomp(df,scale=T)
fviz_pca_biplot(m1,habillage=cluster_number,geom="point",labelsize=3)
```

The biplot displays the following variables as vectors (blue arrows):

- Cost\_Living
- Crime
- Past\_Job\_Growth
- Future\_Job\_Growth
- Total\_Property
- Climate
- Total\_Violent
- Jobs
- Blue\_Collar\_Jobs
- High\_Jobs
- White\_Collar\_Jobs
- Average\_Jobs
- Transportation
- Population\_2000
- Health\_Care
- Education
- Arts

The legend indicates four groups of data points:

- Group 1: Red circles
- Group 2: Green triangles
- Group 3: Cyan squares
- Group 4: Purple pluses

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:data.table':
##
##     between, first, last

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
df_avg = data.frame(apply(df_avg %>% select(2:19),2,function(x) x / sum(x)))
df_avg['cluster'] = c(1,2,3,4)
```

```
library(tidyr)
df_td = gather(df_avg,key='type',value='value',1:18)
```

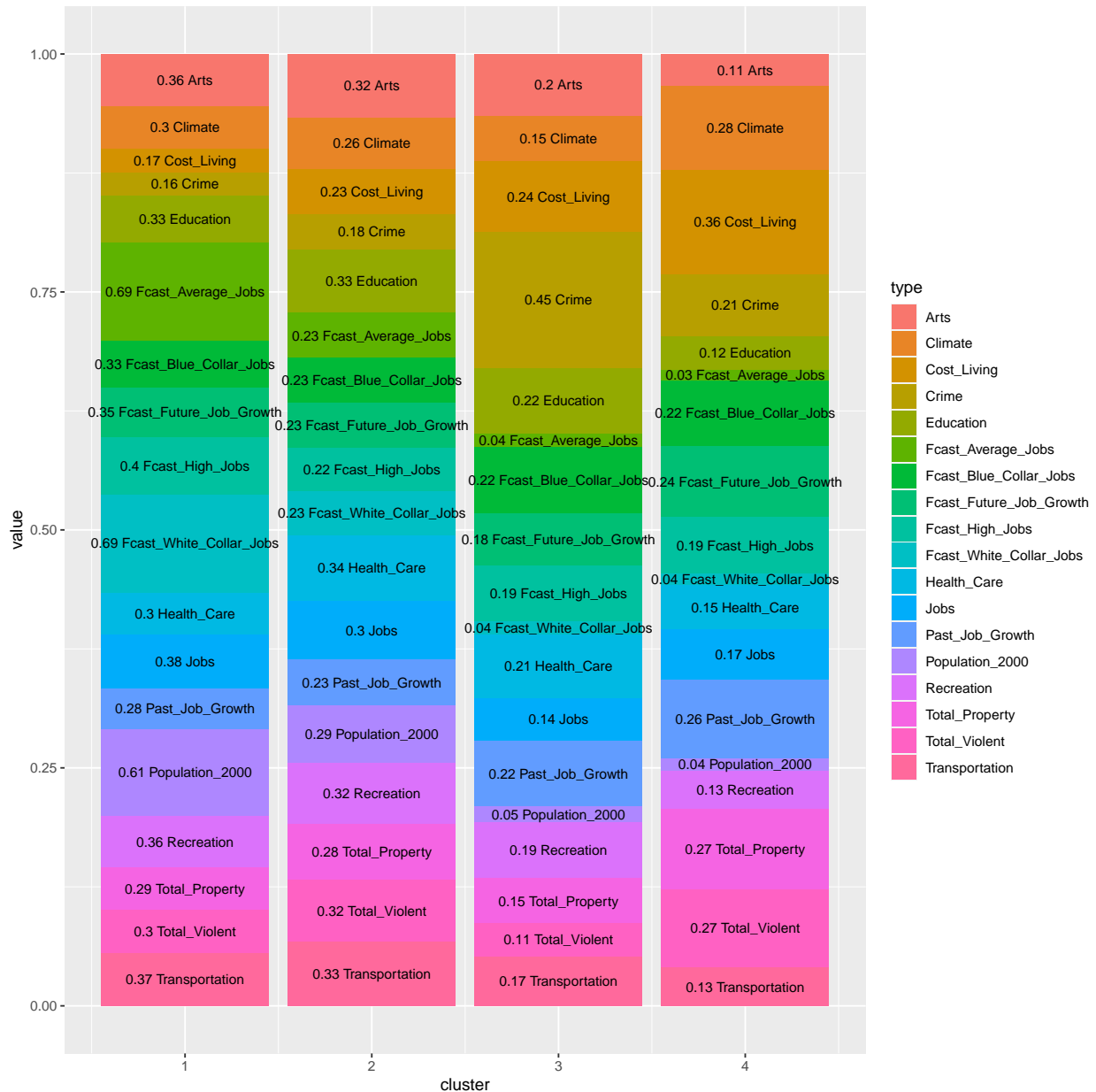
```
df_min = df_td %>% group_by(type) %>% slice(which.min(value))
for (i in 1:4){
  cat('-----')
  cat('kCluster',i,'Minimum Attribute\n')
  print(filter(df_min, cluster==i)$type)
}
```

```
## -----kCluster 1 Minimum Attribute
## [1] "Cost_Living" "Crime"
## -----kCluster 2 Minimum Attribute
## character(0)
## -----kCluster 3 Minimum Attribute
## [1] "Climate" "Fcast_Blue_Collar_Jobs"
## [3] "Fcast_Future_Job_Growth" "Fcast_High_Jobs"
## [5] "Jobs" "Past_Job_Growth"
## [7] "Total_Property" "Total_Violent"
## -----kCluster 4 Minimum Attribute
## [1] "Arts" "Education"
## [3] "Fcast_Average_Jobs" "Fcast_White_Collar_Jobs"
## [5] "Health_Care" "Population_2000"
## [7] "Recreation" "Transportation"
```

```
df_max = df_td %>% group_by(type) %>% slice(which.max(value))
for (i in 1:4){
  cat('-----')
  cat('kCluster',i,'Maximum Attribute\n')
  print(filter(df_max, cluster==i)$type)
}
```

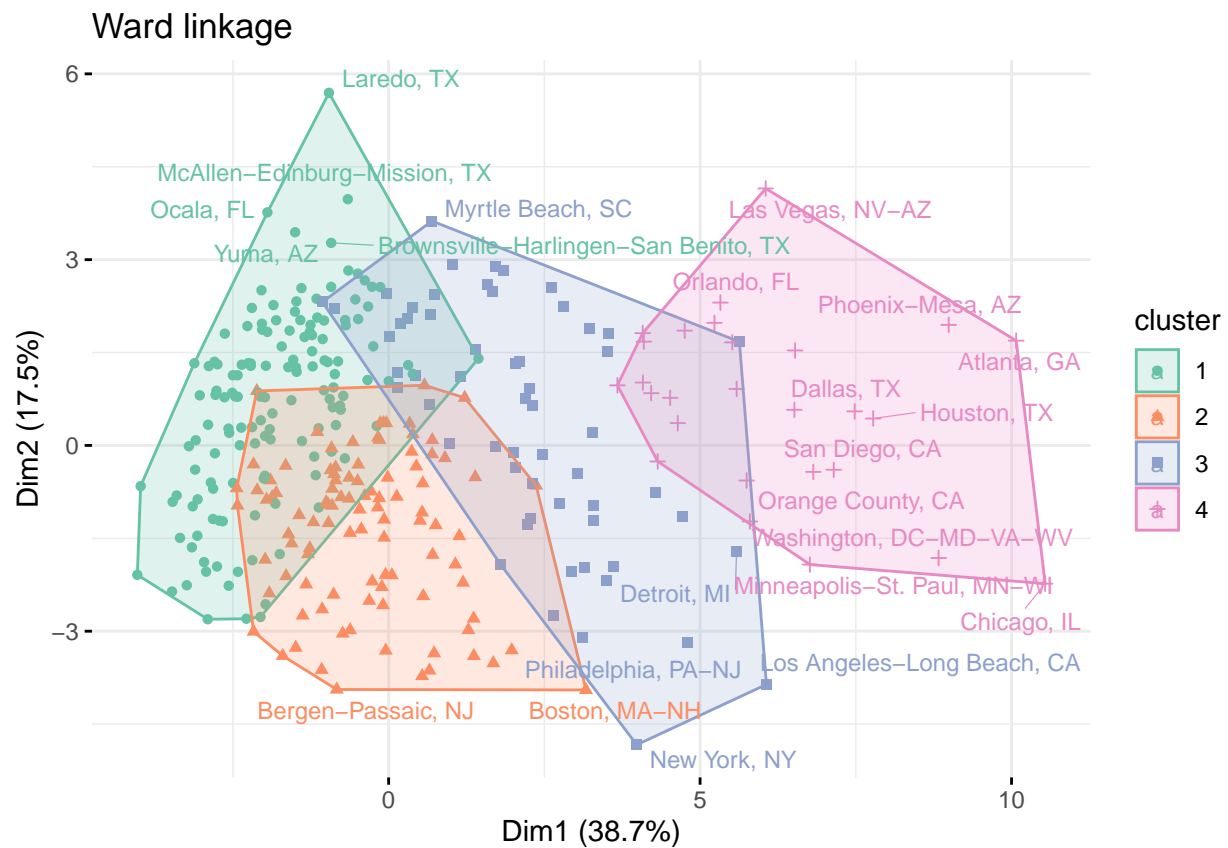
```
## -----kCluster 1 Maximum Attribute
## [1] "Arts" "Climate"
## [3] "Education" "Fcast_Average_Jobs"
## [5] "Fcast_Blue_Collar_Jobs" "Fcast_Future_Job_Growth"
## [7] "Fcast_High_Jobs" "Fcast_White_Collar_Jobs"
## [9] "Jobs" "Past_Job_Growth"
## [11] "Population_2000" "Recreation"
## [13] "Total_Property" "Transportation"
## -----kCluster 2 Maximum Attribute
## [1] "Health_Care" "Total_Violent"
## -----kCluster 3 Maximum Attribute
## [1] "Crime"
## -----kCluster 4 Maximum Attribute
## [1] "Cost_Living"
```

```
ggplot(df_td,aes(x = cluster, y = value, fill = type)) +
  geom_col(position = "fill") +
  geom_text(aes(label = paste(round(value,2),type)),
            position = position_fill(vjust=0.5),
            size=3)
```



```
distance=dist(df)
h1 = hclust(distance,method="ward.D")
cut1 = cutree(h1,k=4)
fviz_cluster(list(data=df,cluster=cut1),main="Ward linkage",
              palette = "Set2", show.clust.cent=F,labels=10,
              repel=T,
              ggtheme=theme_minimal())
```

```
## Warning: ggrepel: 302 unlabeled data points (too many overlaps). Consider
## increasing max.overlaps
```

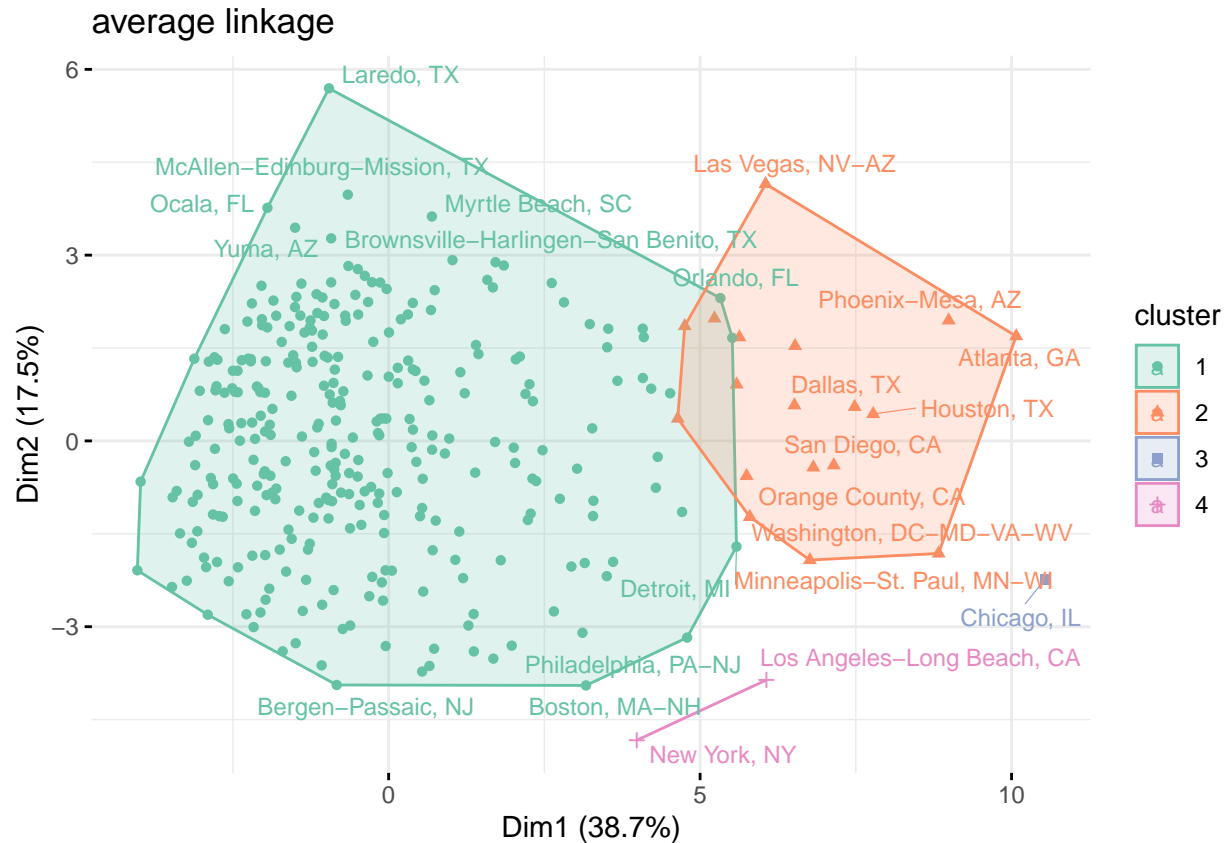


```
c1 = cophenetic(h1)
cor(distance,c1)
```

```
## [1] 0.5079247
```

```
h2 = hclust(distance, method = 'average')
cut2 = cutree(h2,k=4)
fviz_cluster(list(data = df, cluster = cut2),main="average linkage",
  palette = "Set2",show.clust.cent = F, labels = 10,
  repel = T, # Avoid label overlap (slow)
  ggtheme = theme_minimal()
)
```

```
## Warning: ggrepel: 302 unlabeled data points (too many overlaps). Consider
## increasing max.overlaps
```



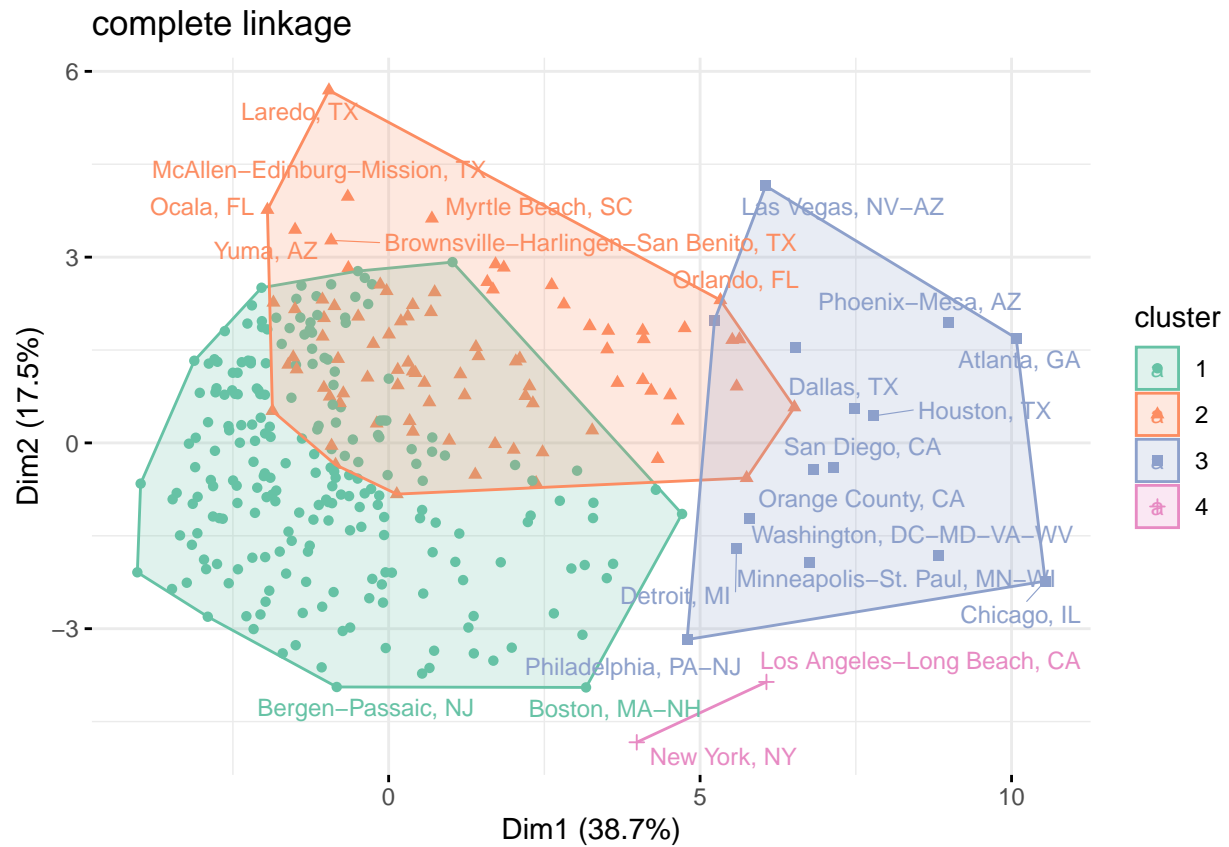
```
c2 = cophenetic(h2)
cor(distance,c2)
```

```
## [1] 0.8047003
```

```
h3 = hclust(distance, method = 'complete')
cut3 = cutree(h3,k=4)
fviz_cluster(list(data = df, cluster = cut3),main="complete linkage",
  palette = "Set2",show.clust.cent = F, labels = 10,
  repel = T,
  ggtheme = theme_minimal()
)
```

```
## Warning: ggrepel: 302 unlabeled data points (too many overlaps). Consider
## increasing max.overlaps
```

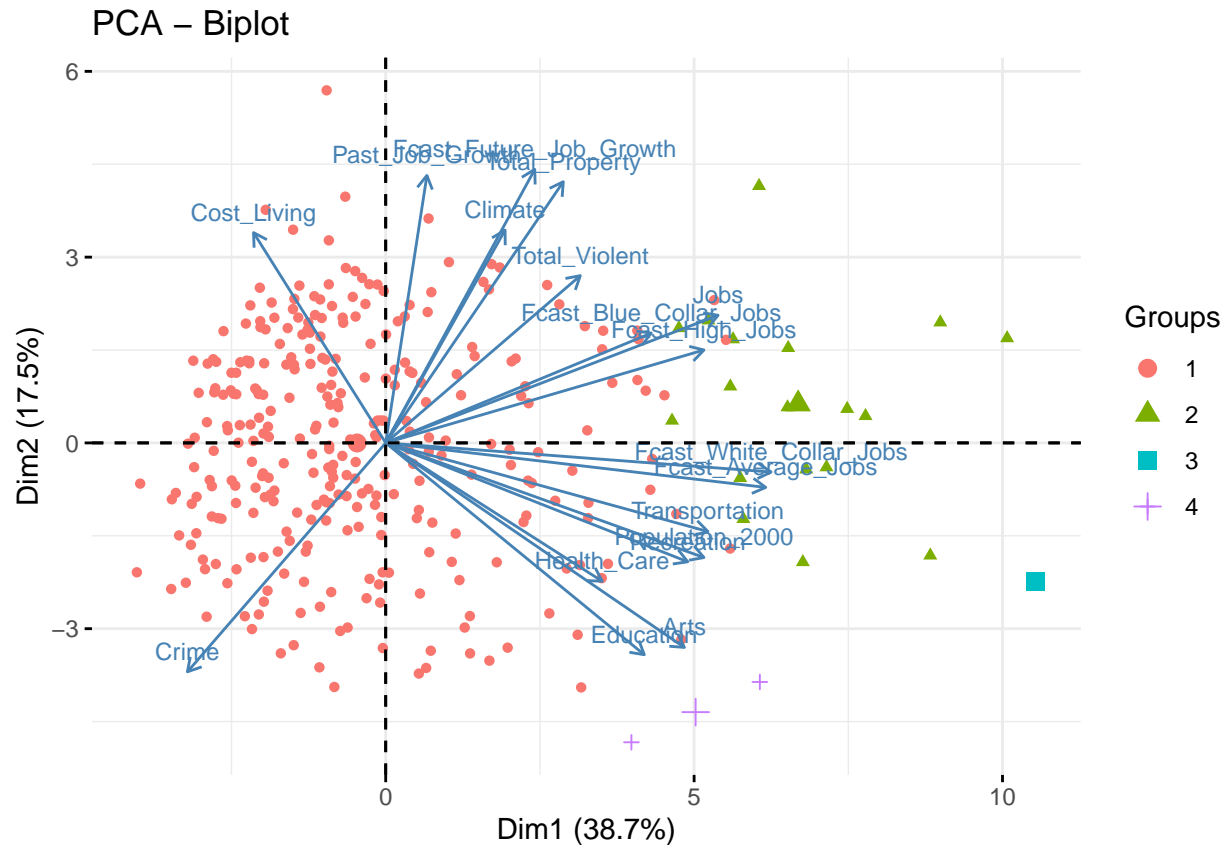




```
c3 = cophenetic(h3)
cor(distance, c3)
```

```
## [1] 0.6848473
```

```
fviz_pca_biplot(m1, habillage=cut2, geom="point", labelsize=3)
```



```
dfh = df_mm
dfh$cluster = cut2
dfh$Metropolitan_Area = NULL
dfh_avg = setDT(dfh)[, lapply(.SD, mean), keyby = cluster]

dfh_avg = data.frame(apply(dfh_avg %>% select(2:19), 2, function(x) x / sum(x)))
dfh_avg['cluster'] = c(1,2,3,4)

dfh_td = gather(dfh_avg, key='type', value='value', 1:18)
```

```
dfh_min = dfh_td %>% group_by(type) %>% slice(which.min(value))
for (i in 1:4){
  cat('-----')
  cat('hCluster', i, 'Minimum Attribute\n')
  print(filter(dfh_min, cluster==i)$type)
}
```

```
## -----hCluster 1 Minimum Attribute
## [1] "Arts" "Education"
## [3] "Fcast_Average_Jobs" "Fcast_White_Collar_Jobs"
## [5] "Health_Care" "Population_2000"
## [7] "Recreation" "Total_Property"
## [9] "Total_Violent" "Transportation"
## -----hCluster 2 Minimum Attribute
## character(0)
```

```
## -----hCluster 3 Minimum Attribute
## [1] "Climate"
## -----hCluster 4 Minimum Attribute
## [1] "Cost_Living"      "Crime"
## [3] "Fcast_Blue_Collar_Jobs" "Fcast_Future_Job_Growth"
## [5] "Fcast_High_Jobs"    "Jobs"
## [7] "Past_Job_Growth"

dfh_max = dfh_td %>% group_by(type) %>% slice(which.max(value))
for (i in 1:4){
  cat('-----')
  cat('hCluster',i,'Maximum Attribute\n')
  print(filter(dfh_max, cluster==i)$type)
}
```

```
## -----hCluster 1 Maximum Attribute
## [1] "Cost_Living" "Crime"
## -----hCluster 2 Maximum Attribute
## [1] "Fcast_Future_Job_Growth" "Fcast_High_Jobs"
## [3] "Jobs"                  "Past_Job_Growth"
## [5] "Total_Property"
## -----hCluster 3 Maximum Attribute
## [1] "Education"      "Fcast_Average_Jobs"
## [3] "Fcast_Blue_Collar_Jobs" "Fcast_White_Collar_Jobs"
## [5] "Health_Care"     "Recreation"
## [7] "Transportation"
## -----hCluster 4 Maximum Attribute
## [1] "Arts"          "Climate"          "Population_2000" "Total_Violent"
```

```
ggplot(dfh_td,aes(x = cluster, y = value, fill = type)) +
  geom_col(position = "fill") +
  geom_text(aes(label = paste(round(value,2),type)),
            position = position_fill(vjust=0.5),
            size=3)
```

