

# Assignment 5

Kangrui Liu

2023-12-01

You may work in pairs or individually for this assignment. Make sure you join a group in Canvas if you are working in pairs. Turn in this assignment as an HTML or PDF file to ELMS. Make sure to include the R Markdown or Quarto file that was used to generate it. Include the GitHub link for the repository containing these files.

## Github Link

- [https://github.com/krlui67/Assignment\\_SURV727/tree/main/a5](https://github.com/krlui67/Assignment_SURV727/tree/main/a5)

```
library(censusapi)
library(tidyverse)
library(magrittr)
library(factoextra)
library(dplyr)
library(stringr)
library(ggplot2)
library(MASS)
```

## Exploring ACS Data

In this notebook, we use the Census API to gather data from the American Community Survey (ACS). This requires an access key, which can be obtained here:

- [https://api.census.gov/data/key\\_signup.html](https://api.census.gov/data/key_signup.html)

```
cs_key <- read.table("cs_key.txt")[1,1]
```

```
acs_il_c <- getCensus(name = "acs/acs5",
  vintage = 2016,
  vars = c("NAME", "B01003_001E", "B19013_001E", "B19301_001E"),
  region = "county:*",
  regionin = "state:17",
  key = cs_key)
```

```
acs_il_c <- acs_il_c %>% dplyr::rename(pop = B01003_001E, hh_income = B19013_001E, income = B19301_001E)
head(acs_il_c)
```

```
##   state county                NAME    pop hh_income income
## 1    17    067 Hancock County, Illinois 18633    50077  25647
```

```
## 2    17    063 Grundy County, Illinois 50338    67162 30232
## 3    17    091 Kankakee County, Illinois 111493    54697 25111
## 4    17    043 DuPage County, Illinois 930514    81521 40547
## 5    17    003 Alexander County, Illinois 7051    29071 16067
## 6    17    129 Menard County, Illinois 12576    60420 31323
```

Pull map data for California into a data frame.

```
il_map <- map_data("county", region = "illinois")
head(il_map)
```

```
##      long      lat group order  region subregion
## 1 -91.49563 40.21018     1     1 illinois      adams
## 2 -90.91121 40.19299     1     2 illinois      adams
## 3 -90.91121 40.19299     1     3 illinois      adams
## 4 -90.91121 40.10704     1     4 illinois      adams
## 5 -90.91121 39.83775     1     5 illinois      adams
## 6 -90.91694 39.75754     1     6 illinois      adams
```

Join the ACS data with the map data. Note that `il_map` has a column `subregion` which includes county names. We need a corresponding variable in the ACS data to join both data sets. This needs some transformations, among which the function `tolower()` might be useful. Call the joined data `acs_map`.

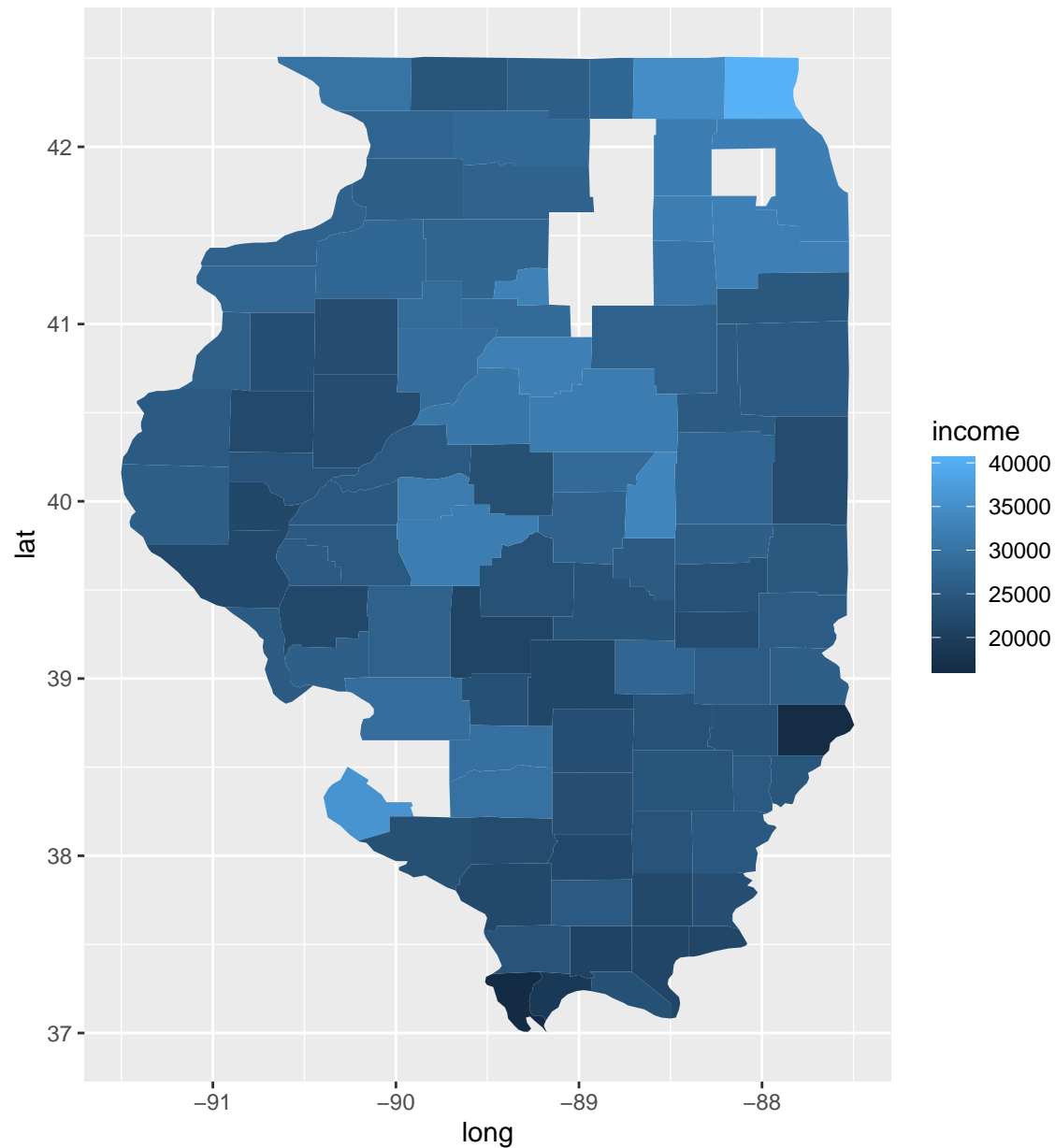
```
# Clean data in acs_il_c
acs_il_c %<>%
  separate(NAME, c("location","state_name"), sep = ",")

for (i in 1:dim(acs_il_c)[1]) {
  acs_il_c$location[i] <- gsub(" County","",acs_il_c$location[i])
  acs_il_c$location[i] <- tolower(acs_il_c$location[i])
}

# Join the ACS data with the map data.
acs_il_map <- inner_join(acs_il_c,il_map,by = join_by(location == subregion))
```

After you do this, plot a map of Illinois with Counties colored by per capita income.

```
ggplot(acs_il_map) + geom_polygon(aes(x = long, y = lat, group = group, fill = income))
```



## Hierarchical Clustering

We want to find clusters of counties that are similar in their population, average household income and per capital income.

First, clean the data so that you have the appropriate variables to use for clustering.

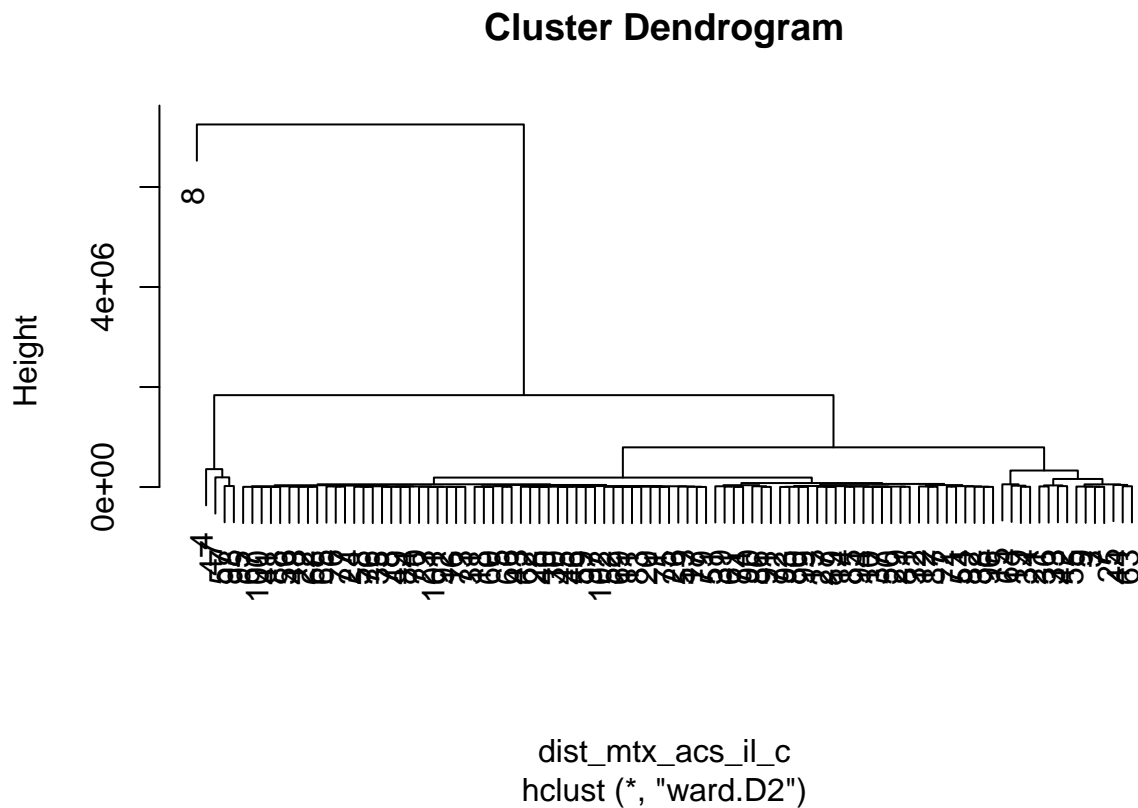
```
cleaned_acs_il_c <- acs_il_c %>% dplyr::select(location, pop, hh_income, income) %>% na.omit()
```

Next, create the distance matrix of the cleaned data. This distance matrix can be used to cluster counties, e.g. using the ward method.

```
# Create the distance matrix
dist_mtx_acs_il_c <- dist(cleaned_acs_il_c[, -1]) # excluding the county identifier
# Perform hierarchical clustering
hc_c <- hclust(dist_mtx_acs_il_c, method = "ward.D2")
```

Plot the dendrogram to find a reasonable number of clusters. Draw boxes around the clusters of your cluster solution.

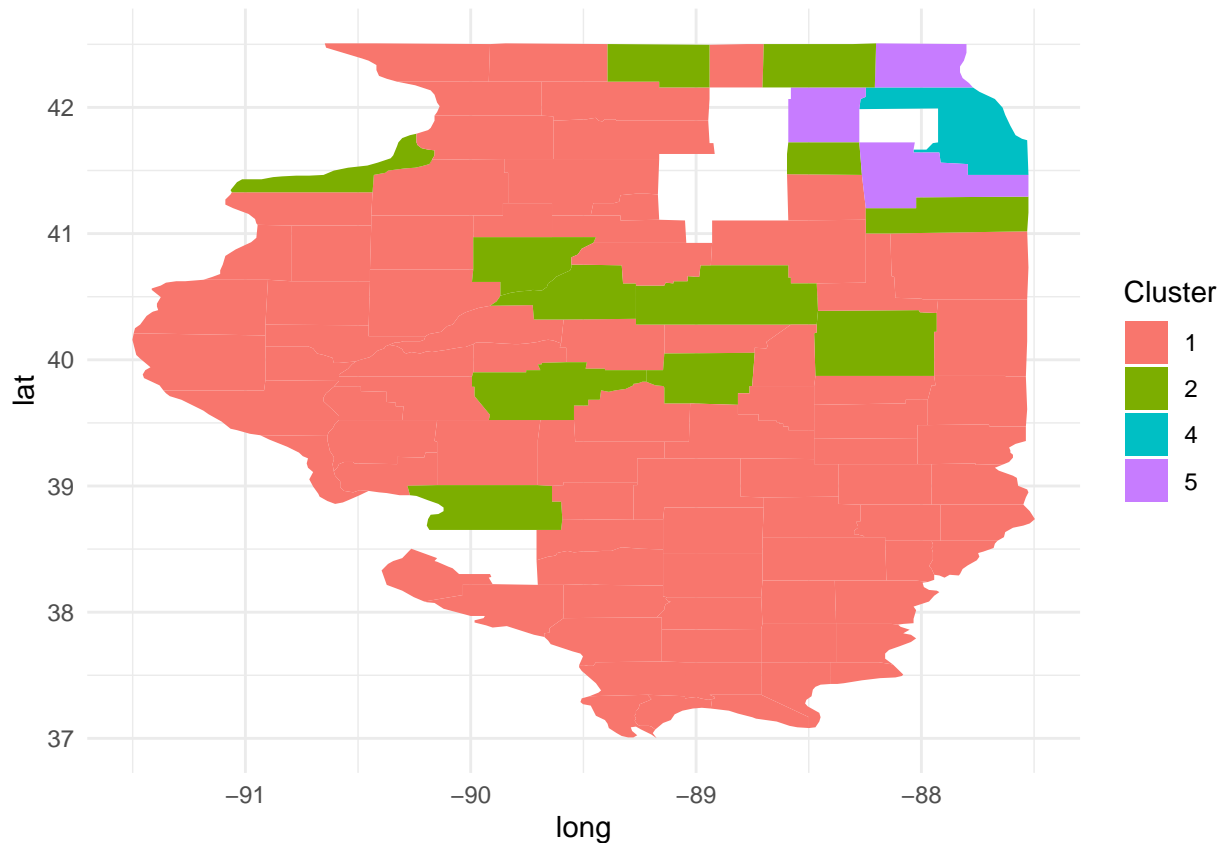
```
# Plot the dendrogram
plot(hc_c)
```



```
#+ rect.hclust(hc_c,k = 3) # Here 'k' is the number of clusters you choose
# Cut the dendrogram to form clusters
clusters <- cutree(hc_c, k = 5)
# Merge the cluster memberships with the original data
cleaned_acs_il_c$cluster <- clusters
```

Visualize the county clusters on a map. For this task, create a new `acs_map` object that now also includes cluster membership as a new column. This column should be called `cluster`.

```
acs_map <- inner_join(cleaned_acs_il_c, il_map, by = join_by(location == subregion))
ggplot(acs_map) + geom_polygon(aes(x = long, y = lat, group = group, fill = factor(cluster))) + theme_m...
```



## Census Tracts

For the next section we need ACS data on a census tract level. We use the same variables as before.

```
acs_il_t <- getCensus(name = "acs/acs5",
  vintage = 2016,
  vars = c("NAME", "B01003_001E", "B19013_001E", "B19301_001E"),
  region = "tract:*",
  regionin = "state:17",
  key = cs_key) %>% dplyr::rename(pop = B01003_001E,
  hh_income = B19013_001E,
  income = B19301_001E)
acs_il_t[acs_il_t == -66666666] <- NA
head(acs_il_t)
```

```
##   state county  tract                                NAME  pop
## 1    17     031 806002 Census Tract 8060.02, Cook County, Illinois 7304
## 2    17     031 806003 Census Tract 8060.03, Cook County, Illinois 7577
## 3    17     031 806400 Census Tract 8064, Cook County, Illinois 2684
## 4    17     031 806501 Census Tract 8065.01, Cook County, Illinois 2590
## 5    17     031 750600 Census Tract 7506, Cook County, Illinois 3594
## 6    17     031 310200 Census Tract 3102, Cook County, Illinois 1521
##   hh_income income
## 1    56975  23750
```

```
## 2      53769  25016
## 3      62750  30154
## 4      53583  20282
## 5      40125  18347
## 6      63250  31403
```

```
# Clean NAME in acs_il_t
acs_il_t %>%
  separate(NAME, c("ct", "location", "state_name"), sep = ",")
for (i in 1:dim(acs_il_t)[1]) {
  acs_il_t$location[i] <- gsub(" County", "", acs_il_t$location[i])
  acs_il_t$location[i] <- tolower(acs_il_t$location[i])
}
acs_il_t$ct <- trimws(acs_il_t$ct)
acs_il_t$location <- trimws(acs_il_t$location)
acs_il_t$state_name <- trimws(acs_il_t$state_name)
```

## k-Means

As before, clean our data for clustering census tracts based on population, average household income and per capita income.

```
cleaned_acs_il_t <- acs_il_t %>%
  dplyr::select(location, pop, hh_income, income) %>%
  na.omit()
```

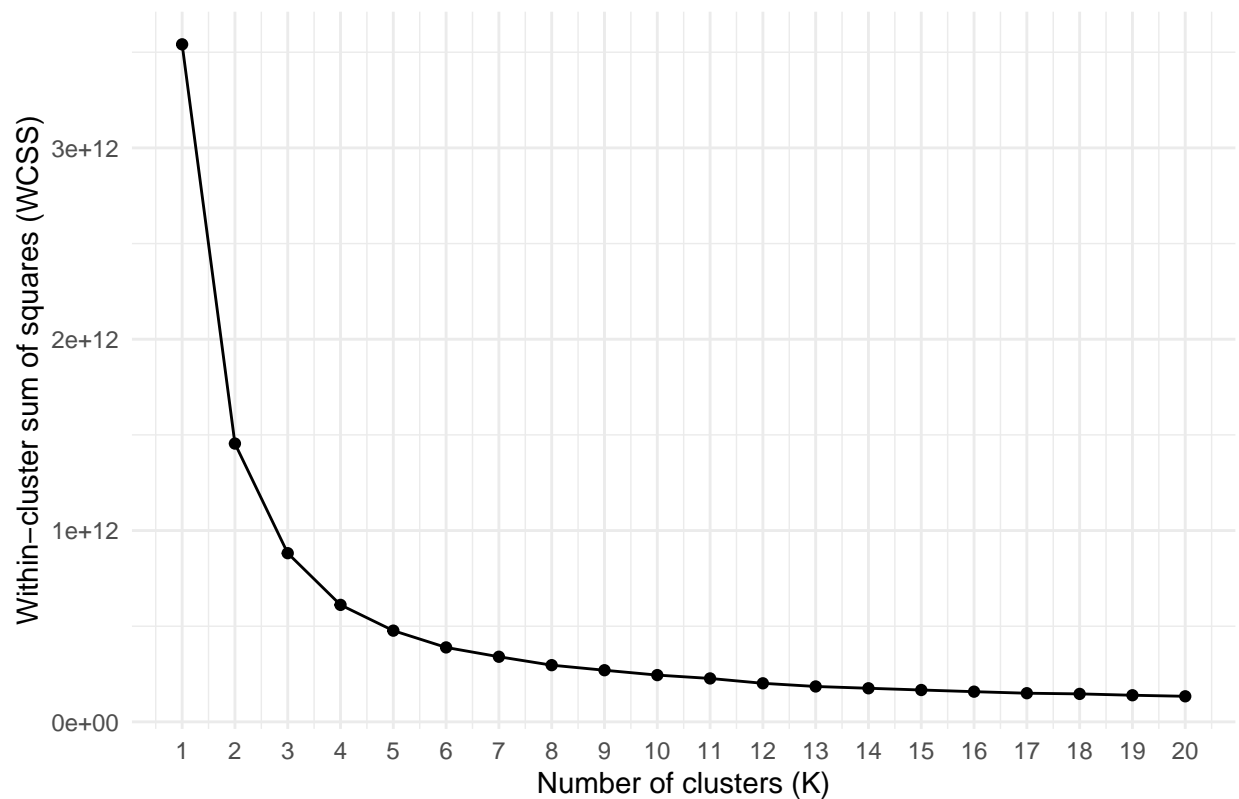
Since we want to use K Means in this section, we start by determining the optimal number of K that results in Clusters with low within but high between variation. Plot within cluster sums of squares for a range of K (e.g. up to 20).

```
# Compute WCSS for a range of K values
set.seed(123) # Set seed for reproducibility
wcss <- map_dbl(1:20, function(k) {
  kmeans(cleaned_acs_il_t[, -1], centers = k, nstart = 20)$tot.withinss
})

# Create a dataframe for plotting
elbow_data <- tibble(k = 1:20, wcss = wcss)

# Plot the elbow plot
ggplot(elbow_data, aes(x = k, y = wcss)) +
  geom_line() +
  geom_point() +
  scale_x_continuous(breaks = 1:20) +
  labs(x = "Number of clusters (K)", y = "Within-cluster sum of squares (WCSS)",
       title = "Elbow Method for Optimal K") +
  theme_minimal()
```

## Elbow Method for Optimal K



Run `kmeans()` for the optimal number of clusters based on the plot above.

Find the mean population, household income and per capita income grouped by clusters. In addition, display the most frequent county that can be observed within each cluster.

```
k <- 3
# Create the distance matrix
dist_mtx_acs_il_t <- dist(cleaned_acs_il_t[, -1]) # excluding the county identifier
# Perform hierarchical clustering
hc_t <- hclust(dist_mtx_acs_il_t, method = "ward.D2")
# Cut the dendrogram to form clusters
clusters <- cutree(hc_t, k = k)
# Merge the cluster memberships with the original data
cleaned_acs_il_t$cluster <- clusters
```

```
cleaned_acs_il_t %>% group_by(cluster) %>% summarise(m_pop = mean(pop), m_hh=mean(hh_income), m_inc=mean(inc))
```

```
## # A tibble: 3 x 4
##   cluster m_pop    m_hh    m_inc
##   <int> <dbl>    <dbl> <dbl>
## 1       1 3810.   45051. 22649.
## 2       2 4482.  117231. 59573.
## 3       3 4958.   77859. 38121.
```

```
cleaned_acs_il_t %>% group_by(cluster) %>% count(location) %>% top_n(1, wt = n)
```

```
## # A tibble: 3 x 3
## # Groups:   cluster [3]
##   cluster location      n
##   <int> <chr>      <int>
## 1      1 cook      864
## 2      2 cook      193
## 3      3 cook      255
```

```
# top_n Select top (or bottom) n rows (by value)
```

As you might have seen earlier, it's not always clear which number of clusters is the optimal choice. To automate K Means clustering, program a function based on `kmeans()` that takes K as an argument. You can fix the other arguments, e.g. such that a specific dataset is always used when calling the function.

```
optimize_kmeans <- function(k){
  scaled_data <- scale(cleaned_acs_il_t[, 2:4])
  set.seed(123) # Set seed for reproducibility
  kmeans_result <- kmeans(scaled_data, centers = k, nstart = 20)$tot.withinss
  cluster_var_name <- paste("cluster_k", k, sep = "")
  cleaned_acs_il_t[[cluster_var_name]] <- kmeans_result
  return(cleaned_acs_il_t)
}

cleaned_acs_il_t <- optimize_kmeans(2)
```

We want to utilize this function to iterate over multiple Ks (e.g., K = 2, ..., 10) and -- each time -- add the resulting cluster membership as a new variable to our (cleaned) original data frame (`acs_il_t`). There are multiple solutions for this task, e.g. think about the `apply` family or `for` loops.

```
k_values <- 2:10 # range of K

for (k in k_values) {
  cleaned_acs_il_t <- optimize_kmeans(k)
}
```

Finally, display the first rows of the updated data set (with multiple cluster columns).

```
head(cleaned_acs_il_t,1)
```

```
##   location pop hh_income income cluster cluster_k2 cluster_k3 cluster_k4
## 1   cook 7304   56975  23750      1   5604.163  4151.776  3384.043
##   cluster_k5 cluster_k6 cluster_k7 cluster_k8 cluster_k9 cluster_k10
## 1  2857.685  2470.315  2164.351  1961.135  1798.685  1658.639
```