

Assignment 5

Due at 11:59pm on November 25.

Effy Tang & Weiqi Yang

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 Repository: <https://github.com/effytang/Tang-Yang-a5>

```
library(censusapi)
library(tidyverse)
```

Warning: package 'ggplot2' was built under R version 4.4.3

```
-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.4      v readr      2.1.5
v forcats    1.0.0      v stringr    1.5.1
v ggplot2     4.0.1      v tibble     3.2.1
v lubridate  1.9.4      v tidyr      1.3.1
v purrr       1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
```

```
library(magrittr)
```

Attaching package: 'magrittr'

The following object is masked from 'package:purrr':

```
set_names
```

The following object is masked from 'package:tidyr':

```
extract
```

```
library(factoextra)
```

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

```
library(maps)
```

Attaching package: 'maps'

The following object is masked from 'package:purrr':

```
map
```

```
cs_key <- readLines("census-key.txt")
```

Warning in readLines("census-key.txt"): incomplete final line found on
'census-key.txt'

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

```
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) %>%
  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 Illinois 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. Not 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`.

```

acs_il_c <- acs_il_c %>%
  mutate(subregion = str_replace(NAME, " County, Illinois", ""),
         subregion = tolower(subregion))

acs_map <- il_map %>%
  left_join(acs_il_c, by = "subregion")

head(acs_map)

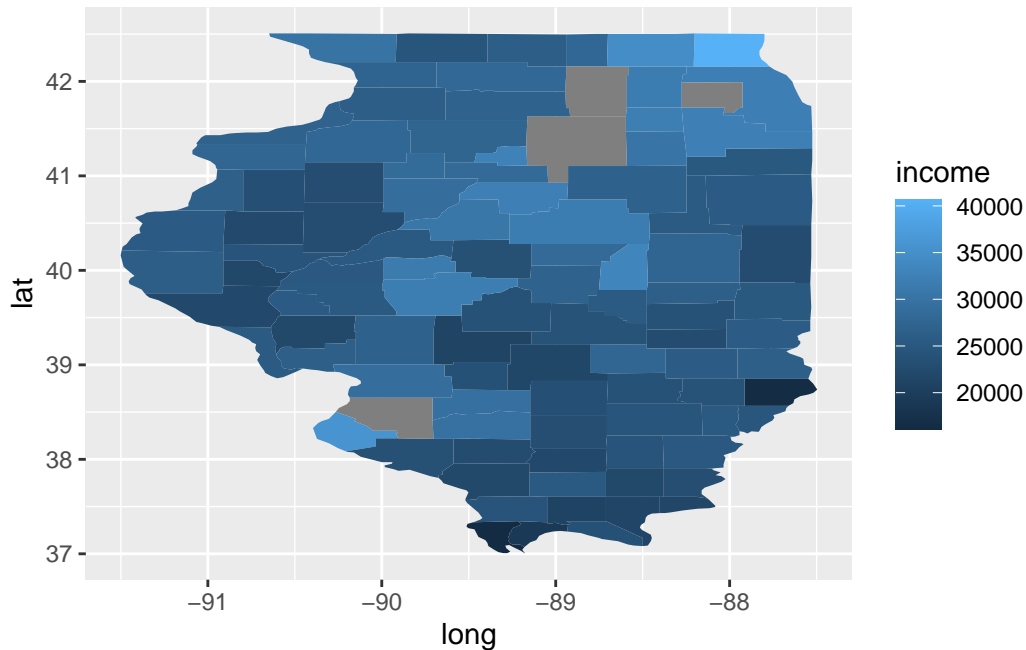
```

	long	lat	group	order	region	subregion	state	county
1	-91.49563	40.21018	1	1	illinois	adams	17	001
2	-90.91121	40.19299	1	2	illinois	adams	17	001
3	-90.91121	40.19299	1	3	illinois	adams	17	001
4	-90.91121	40.10704	1	4	illinois	adams	17	001
5	-90.91121	39.83775	1	5	illinois	adams	17	001
6	-90.91694	39.75754	1	6	illinois	adams	17	001

	NAME	pop	hh_income	income
1	Adams County, Illinois	66949	48065	26053
2	Adams County, Illinois	66949	48065	26053
3	Adams County, Illinois	66949	48065	26053
4	Adams County, Illinois	66949	48065	26053
5	Adams County, Illinois	66949	48065	26053
6	Adams County, Illinois	66949	48065	26053

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

```
ggplot(acs_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 capita income. First, clean the data so that you have the appropriate variables to use for clustering.

```
acs_il_c_clean <- acs_il_c %>%
  select(NAME, pop, hh_income, income) %>%
  na.omit()

acs_clustering <- acs_il_c_clean %>%
  select(pop, hh_income, income) %>%
  scale()

head(acs_clustering)
```

	pop	hh_income	income
1	-0.20225946	-0.1129887	-0.1265936
2	-0.14253141	1.5457905	0.9823871
3	-0.02732344	0.3355661	-0.2562368
4	1.51560431	2.9399029	3.4772915
5	-0.22407842	-2.1524570	-2.4437225
6	-0.21367005	0.8912111	1.2462689

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

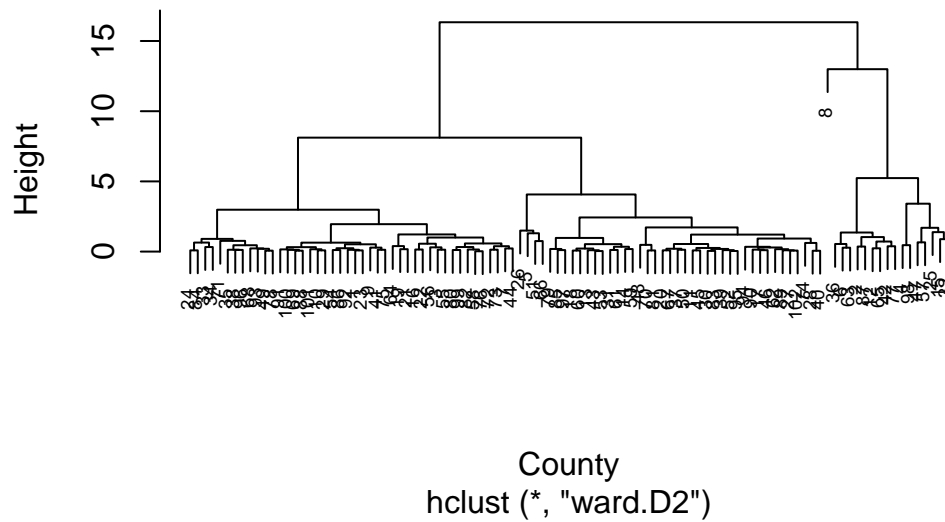
```
dist_matrix <- dist(acs_clustering, method = "euclidean")

hc_ward <- hclust(dist_matrix, method = "ward.D2")
```

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

```
plot(hc_ward,
     main = "Hierarchical Clustering of Illinois Counties",
     xlab = "County",
     ylab = "Height",
     cex = 0.6)
```

Hierarchical Clustering of Illinois Counties



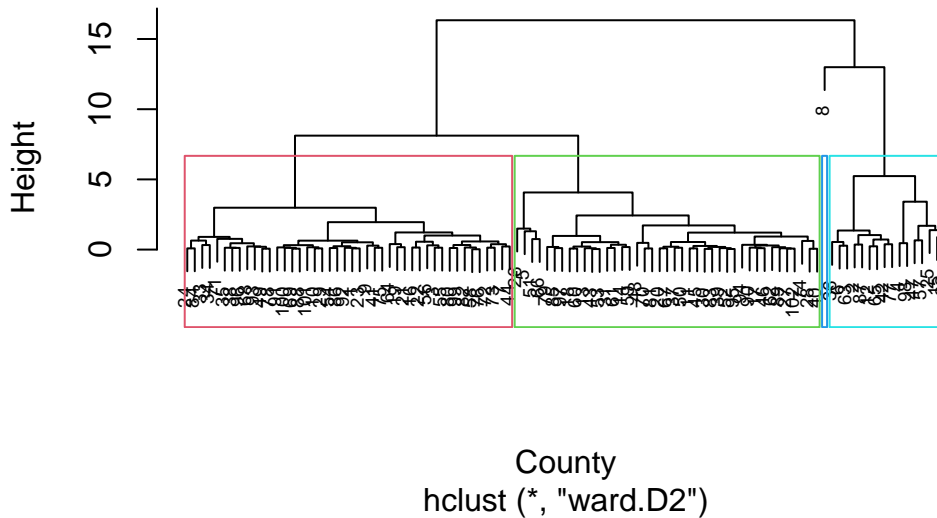
```
k <- 4

clusters <- cutree(hc_ward, k = k)

plot(hc_ward,
     main = "Hierarchical Clustering of Illinois Counties",
     xlab = "County",
     ylab = "Height",
     cex = 0.6)

rect.hclust(hc_ward, k = k, border = 2:5)
```

Hierarchical Clustering of Illinois Counties



```
acs_il_c_clean <- acs_il_c_clean %>%
  mutate(cluster = as.factor(clusters))

table(acsil_c_clean$cluster)
```

```
1  2  3  4
44 16 41  1
```

```
head(acsil_c_clean)
```

	NAME	pop	hh_income	income	cluster
1	Hancock County, Illinois	18633	50077	25647	1
2	Grundy County, Illinois	50338	67162	30232	2
3	Kankakee County, Illinois	111493	54697	25111	1
4	DuPage County, Illinois	930514	81521	40547	2
5	Alexander County, Illinois	7051	29071	16067	3
6	Menard County, Illinois	12576	60420	31323	2

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_il_c_clean <- acs_il_c_clean %>%
  mutate(subregion = str_replace(NAME, " County, Illinois", ""),
         subregion = tolower(subregion))

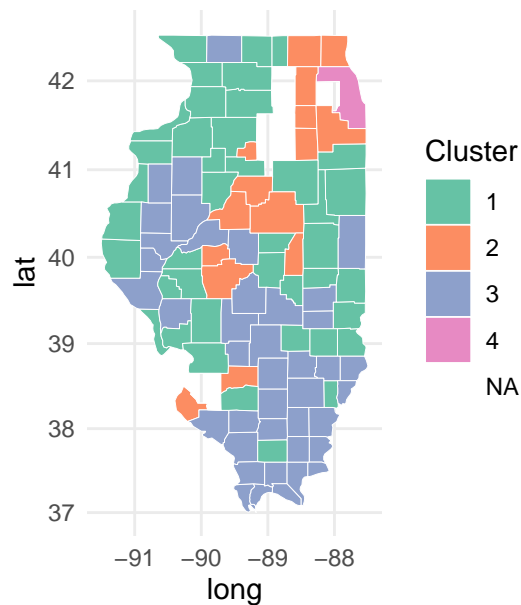
acs_map_clusters <- il_map %>%
  left_join(acs_il_c_clean, by = "subregion")

ggplot(acs_map_clusters) +
  geom_polygon(aes(x = long,
                  y = lat,
                  group = group,
                  fill = cluster),
             color = "white", size = 0.1) +
  coord_map() +
  scale_fill_brewer(palette = "Set2") +
  labs(fill = "Cluster",
       title = "County Clusters in Illinois based on Population and Income") +
  theme_minimal()

```

Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.
 i Please use `linewidth` instead.

County Clusters in Illinois based on Population



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) %>%
  mutate_all(funs(ifelse(.==66666666, NA, .))) %>%
  rename(pop = B01003_001E,
         hh_income = B19013_001E,
         income = B19301_001E)
```

Warning: `funs()` was deprecated in dplyr 0.8.0.

i Please use a list of either functions or lambdas:

```
# Simple named list: list(mean = mean, median = median)
```

```
# Auto named with `tibble::lst()`: tibble::lst(mean, median)
```

```
# Using lambdas list(~ mean(., trim = .2), ~ median(., na.rm = TRUE))
```

```
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

k-Means

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

```
acs_il_t_clean <- acs_il_t %>%  
  select(pop, hh_income, income) %>%  
  na.omit() %>%  
  scale()  
  
nrow(acs_il_t_clean)
```

```
[1] 3109
```

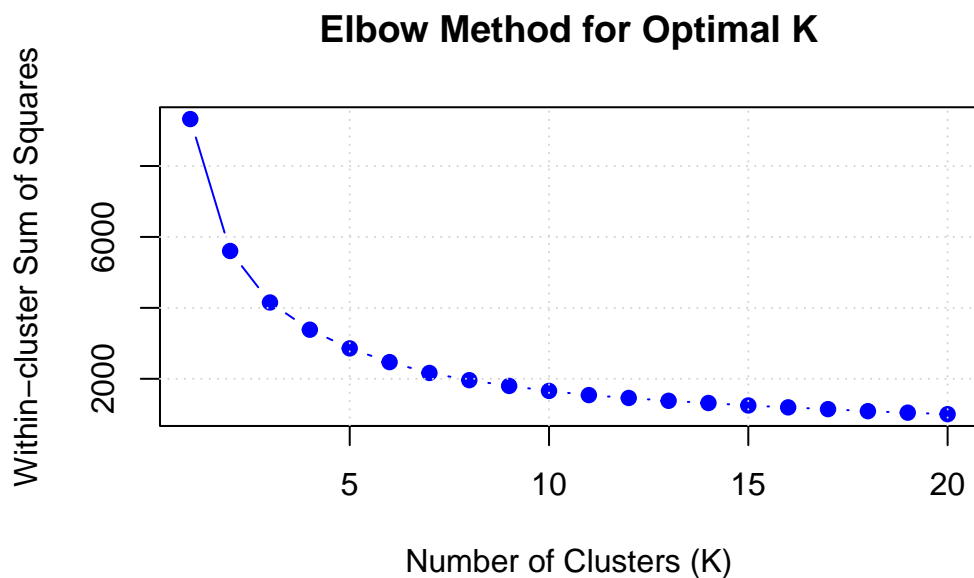
Since we want to use K Means in this section, we start by determining the optimal number of K that results in Clusters with low within but high between variation. Plot within cluster sums of squares for a range of K (e.g. up to 20).

```
wss <- sapply(1:20, function(k) {  
  kmeans(acs_il_t_clean, centers = k, nstart = 25)$tot.withinss  
})
```

```
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations  
Warning: did not converge in 10 iterations
```

Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations
Warning: did not converge in 10 iterations

```
plot(1:20, wss,  
     type = "b",  
     xlab = "Number of Clusters (K)",  
     ylab = "Within-cluster Sum of Squares",  
     main = "Elbow Method for Optimal K",  
     pch = 19,  
     col = "blue")  
  
grid()
```



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

```

set.seed(123)
optimal_k <- 5

km_result <- kmeans(acs_il_t_clean,
                    centers = optimal_k,
                    nstart = 25)

acs_il_t <- acs_il_t %>%
  filter(!is.na(pop) & !is.na(hh_income) & !is.na(income)) %>%
  mutate(cluster = as.factor(km_result$cluster))

table(acs_il_t$cluster)

```

```

  1    2    3    4    5
276 1016 274 744 799

```

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.

```

cluster_means <- acs_il_t %>%
  group_by(cluster) %>%
  summarise(
    mean_pop = mean(pop),
    mean_hh_income = mean(hh_income),
    mean_income = mean(income),
    n_tracts = n()
  )

print(cluster_means)

```

```

# A tibble: 5 x 5
  cluster mean_pop mean_hh_income mean_income n_tracts
  <fct>    <dbl>         <dbl>         <dbl>    <int>
1 1        3896.         122368.         67665.     276
2 2        2686.         37123.         19778.    1016
3 3        7838.         86010.         38154.     274
4 4        5381.         49260.         23275.     744
5 5        3610.         73195.         35913.     799

```

```

most_frequent_county <- acs_il_t %>%
  group_by(cluster, county) %>%
  summarise(n = n(), .groups = "drop") %>%
  group_by(cluster) %>%
  slice_max(n, n = 1) %>%
  select(cluster, county, n)

print(most_frequent_county)

```

```

# A tibble: 5 x 3
# Groups:   cluster [5]
  cluster county    n
  <fct>   <chr> <int>
1 1      031    156
2 2      031    432
3 3      031     87
4 4      031   353
5 5      031   284

```

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.

We want to utilize this function to iterate over multiple Ks (e.g., $K = 2, \dots, 10$) 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.

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

```

run_kmeans <- function(k) {
  set.seed(123)
  km <- kmeans(acs_il_t_clean, centers = k, nstart = 25)
  return(km$cluster)
}

for (k in 2:10) {
  col_name <- paste0("cluster_", k)

  acs_il_t <- acs_il_t %>%
    mutate(!!col_name := as.factor(run_kmeans(k)))
}

```

```
}
```

```
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	cluster	cluster_2	cluster_3	cluster_4	cluster_5	cluster_6
1	56975	23750	4	2	3	4	4	3
2	53769	25016	4	2	3	4	4	3
3	62750	30154	5	2	2	1	5	5
4	53583	20282	2	2	2	2	2	5
5	40125	18347	2	2	2	2	2	4
6	63250	31403	5	2	2	1	5	5

	cluster_7	cluster_8	cluster_9	cluster_10
1	6	3	7	8
2	6	3	7	8
3	3	2	6	7
4	2	2	6	7
5	2	6	4	10
6	3	2	6	7