

Kwapong - Assignment 5

Due at 11:59pm on November 26.

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 assignment repo: <https://github.com/TKwapong/surv-assgn5>

```
library(tidyverse)
library(magrittr)
library(gtrendsR)
library(censusapi)
library(ggplot2)
library(factoextra)
library(ggmap)
```

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

```
#Make sure to obtain your own api key and read into the cs_key object
cs_key <- read_file("census-key.txt")
```

```
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 %<>% mutate(county = tolower(gsub(" County, Illinois", "", NAME)))
head(acs_il_c)

```

	state	county	NAME	pop	hh_income	income
1	17	hancock	Hancock County, Illinois	18633	50077	25647
2	17	grundy	Grundy County, Illinois	50338	67162	30232
3	17	kankakee	Kankakee County, Illinois	111493	54697	25111
4	17	dupage	DuPage County, Illinois	930514	81521	40547
5	17	alexander	Alexander County, Illinois	7051	29071	16067
6	17	menard	Menard County, Illinois	12576	60420	31323

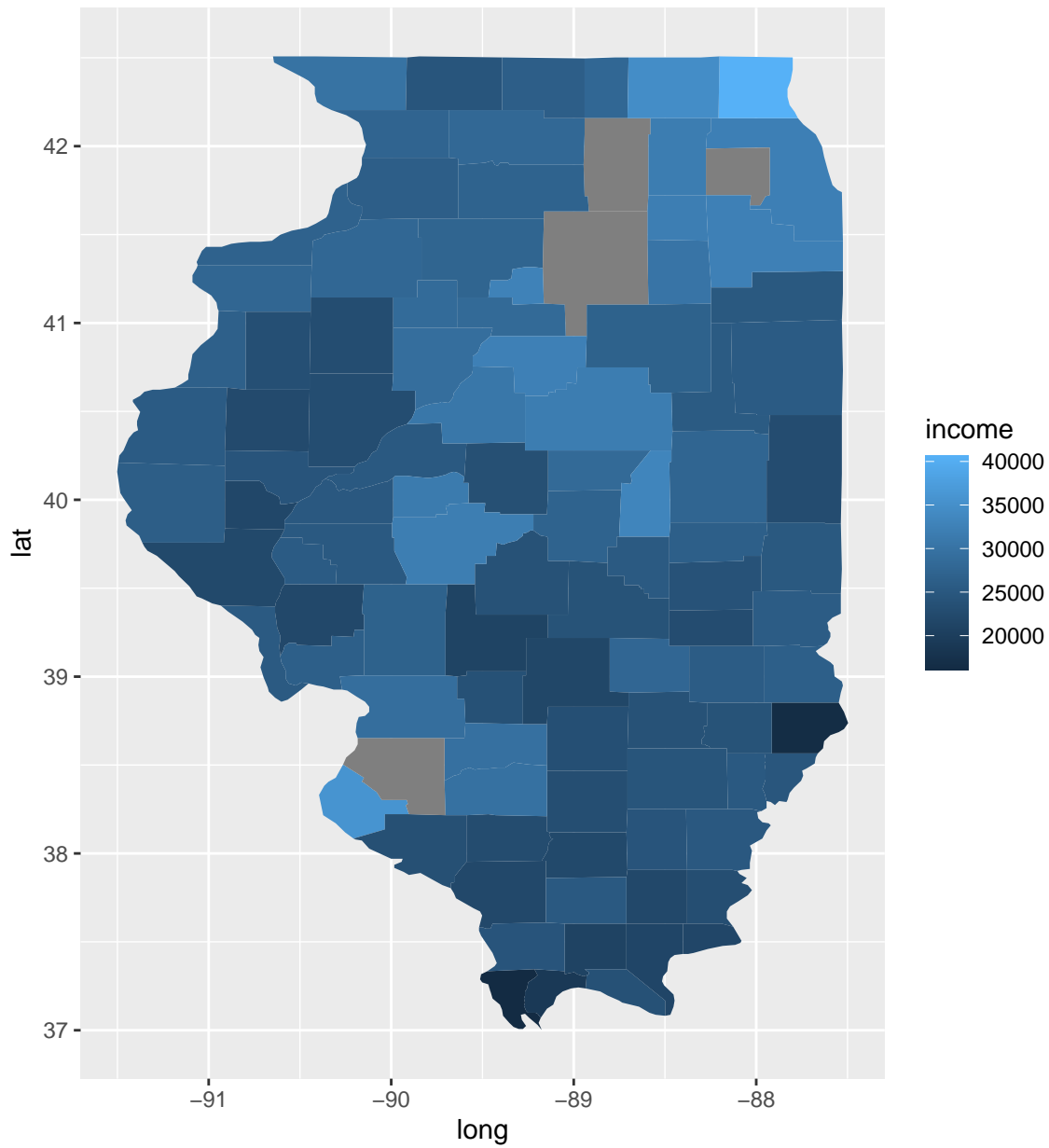
```
acs_map <- il_map %>%
  left_join(acs_il_c, by = c("subregion" = "county"))
head(acs_map)
```

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

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

```
hclust_data <- acs_map %>%
  select(pop, hh_income, income) %>%
  na.omit() %>%
  mutate_all(scale)
```

```
hclust_dist <- dist(hclust_data)
#hc_complete <- hclust(hclust_dist, method = "complete")
#hc_average <- hclust(hclust_dist, method = "average")
hc_ward <- hclust(hclust_dist, 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 (Ward's Method)", xlab = "", sub = "")
rect.hclust(hc_ward, k = 4, border = "red")
```

Hierarchical Clustering (Ward's Method)



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`.

```
cluster_assignments <- cutree(hc_ward, k = 5)

acs_map_cluster <- acs_map %>%
  filter(complete.cases(pop, hh_income, income)) %>%
  mutate(cluster = factor(cluster_assignments))

cluster_means <- acs_map_cluster %>%
  group_by(cluster) %>%
  summarise(
    mean_pop = mean(pop, na.rm = TRUE),
```

```

    mean_hh_income = mean(hh_income, na.rm = TRUE),
    mean_income = mean(income, na.rm = TRUE)
  )

head(cluster_means)

```

```

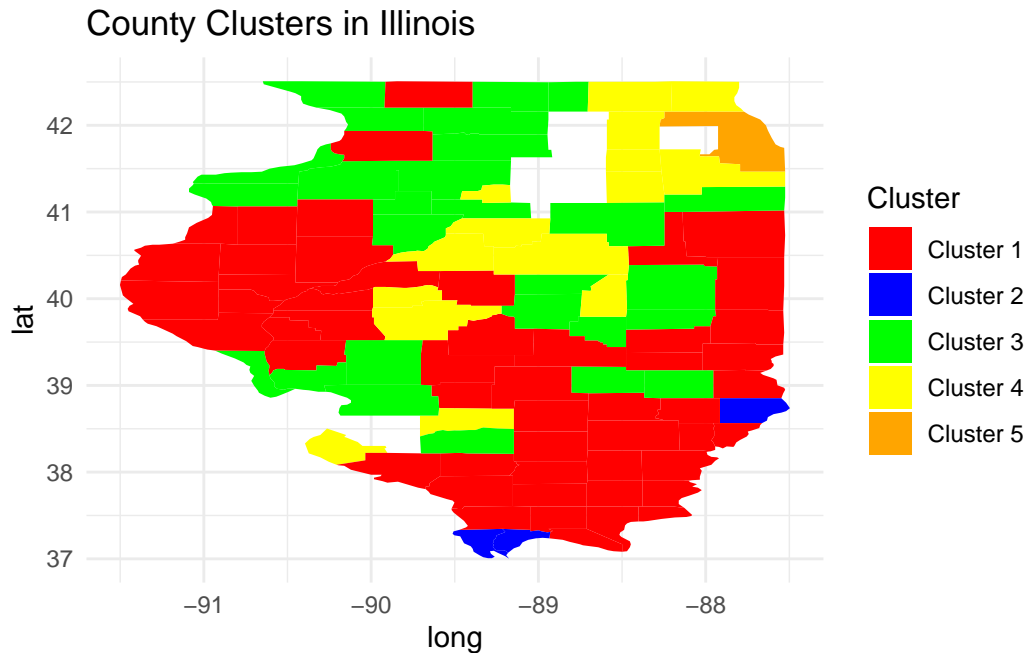
# A tibble: 5 x 4
  cluster mean_pop mean_hh_income mean_income
  <fct>      <dbl>         <dbl>         <dbl>
1 1         22974.         45393.         23888.
2 2          7908.         32334.         17226.
3 3         71660.         52901.         27342.
4 4        173989.         66562.         32456.
5 5        5227575         56902          32179

```

```

ggplot(acs_map_clust) +
  geom_polygon(aes(x = long, y = lat, group = group, fill = cluster)) +
  labs(title = "County Clusters in Illinois") +
  scale_fill_manual(values = c("red", "blue", "green", "yellow", "orange"),
                    name = "Cluster",
                    breaks = c(1, 2, 3, 4, 5),
                    labels = c("Cluster 1", "Cluster 2", "Cluster 3",
                              "Cluster 4", "Cluster 5")) +
  theme_minimal()

```



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(.== -666666666, 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

```
# Clean tract data
acs_il_t <- acs_il_t %>%
  mutate(
    county = str_remove(NAME, "Census Tract [0-9.]+, ") %>%
      str_remove(", Illinois") %>%
      str_remove(" County") %>%
      str_trim() %>%
      tolower(),
  )
```

```
head(acs_il_t)
```

	state	county	tract	NAME	pop
1	17	cook	806002	Census Tract 8060.02, Cook County, Illinois	7304
2	17	cook	806003	Census Tract 8060.03, Cook County, Illinois	7577
3	17	cook	806400	Census Tract 8064, Cook County, Illinois	2684
4	17	cook	806501	Census Tract 8065.01, Cook County, Illinois	2590
5	17	cook	750600	Census Tract 7506, Cook County, Illinois	3594
6	17	cook	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.

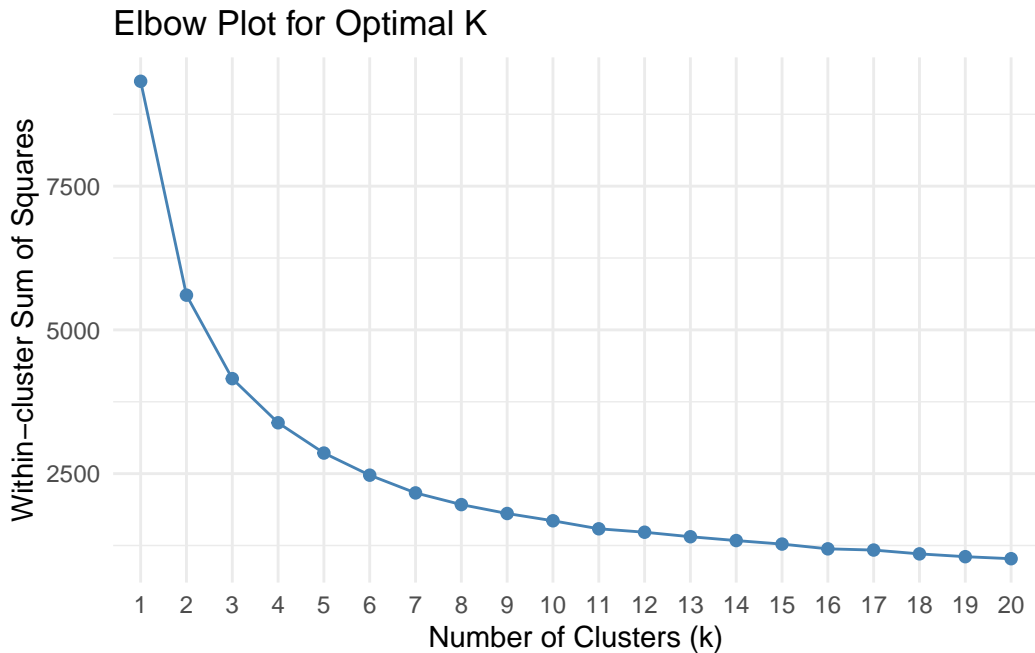
```
kmeans_data <- acs_il_t %>%
  select(pop, hh_income, income) %>%
  na.omit() %>%
  mutate_all(scale)

head(kmeans_data)
```

	pop	hh_income	income
1	1.6189842	-0.14115103	-0.43466339
2	1.7582445	-0.24892639	-0.35470216
3	-0.7377284	0.05298581	-0.03018336
4	-0.7856788	-0.25517911	-0.65370411
5	-0.2735274	-0.70759359	-0.77591974
6	-1.3309874	0.06979420	0.04870414

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).

```
fviz_nbclust(kmeans_data,
             kmeans,
             method = "wss",
             k.max = 20) +
  labs(title = "Elbow Plot for Optimal K",
       x = "Number of Clusters (k)",
       y = "Within-cluster Sum of Squares") +
  theme_minimal()
```



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

```
set.seed(123)
kmeans_result <- kmeans(kmeans_data, centers = 6, nstart = 25)

kmeans_data$cluster <- kmeans_result$cluster
acs_il_t$cluster <- NA
acs_il_t$cluster[complete.cases(acs_il_t[c("pop", "hh_income", "income")])] <- kmeans_result$cluster
```

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.

```
# Calculate cluster summaries
cluster_summary <- acs_il_t %>%
  group_by(cluster) %>%
  summarise(
    mean_pop = round(mean(pop, na.rm = TRUE), 2),
    mean_hh_income = round(mean(hh_income, na.rm = TRUE), 2),
    mean_income = round(mean(income, na.rm = TRUE), 2),
    n = n()
  ) %>%
  filter(!is.na(cluster))

# Find most frequent county
```

```

most_frequent_county <- acs_il_t %>%
  filter(!is.na(cluster)) %>%
  group_by(cluster, county) %>%
  summarise(count = n(), .groups = 'drop') %>%
  group_by(cluster) %>%
  slice_max(order_by = count, n = 1) %>%
  arrange(cluster)

print("Cluster Summary Statistics:")

```

```
[1] "Cluster Summary Statistics:"
```

```
print(as.data.frame(cluster_summary), row.names = FALSE)
```

cluster	mean_pop	mean_hh_income	mean_income	n
1	4518.59	92962.88	45054.83	527
2	3812.25	135623.34	77010.10	154
3	5965.38	53871.63	24940.38	690
4	2689.25	32061.32	17260.47	764
5	3305.63	58016.03	29402.11	914
6	11339.97	93651.17	39361.42	60

```
print("\nMost Frequent County per Cluster:")
```

```
[1] "\nMost Frequent County per Cluster:"
```

```
print(as.data.frame(most_frequent_county), row.names = FALSE)
```

cluster	county	count
1	cook	220
2	cook	97
3	cook	326
4	cook	379
5	cook	282
6	will	12

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.

```

means_cluster <- function(data, K) {
  set.seed(123)
  data_numeric <- data[sapply(data, is.numeric)]
  data_numeric <- na.omit(data_numeric)
  # Run K-means clustering for the given number of clusters (K)
  kmeans_result <- kmeans(data_numeric, centers = K, nstart = 25)

  # Add the resulting cluster membership to the original dataset
  # Ensure to match rows back to the original data
  data$cluster_K <- kmeans_result$cluster[match(rownames(data_numeric), rownames(data))]

  # Return the modified dataset with clusters
  return(data)
}

```

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.

```

#K values to iterate over
K_values <- 2:10

# Initialize the dataset
acs_il_t_numeric <- acs_il_t %>%
  select(county, pop, hh_income, income) %>%
  na.omit() %>%
  mutate_at(vars(pop, hh_income, income), scale)

# Iterate over K_values and apply means_cluster for each value of K,
# storing the results in a list
clustered_data_list <- lapply(K_values, function(K) {
  means_cluster(acs_il_t_numeric, K) # Apply the function for each K
})

# Combine the results into a single dataset with each cluster column
final_data <- clustered_data_list[[1]] # Start with the first result

# Combine the cluster columns for each value of K (2 to 10)
for (i in 2:length(clustered_data_list)) {
  final_data <- cbind(final_data, clustered_data_list[[i]]$cluster_K)
}

```

```
# Rename the columns to indicate the cluster for each K
colnames(final_data)[(ncol(final_data) - length(K_values) + 1)
                     :ncol(final_data)] <- paste0("cluster_K", K_values)
```

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

```
head(final_data, 1)
```

	county	pop	hh_income	income	cluster_K2	cluster_K3	cluster_K4
1	cook	1.618984	-0.141151	-0.4346634	2	3	1
	cluster_K5	cluster_K6	cluster_K7	cluster_K8	cluster_K9	cluster_K10	
1	4	3	6	8	7	8	