

Assignment 5

Due at 11:59pm on November 26.

Feiran Ge

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/IvyG-a/727_Assignment5.git

```
library(censusapi)
```

Attaching package: 'censusapi'

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

getFunction

```
library(tidyverse)
```

```
-- 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    3.5.1      v tibble     3.2.1
v lubridate  1.9.3      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>

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

```
cs_key <- "a966f6817d43c2f9870cf68a8dc7979c68a8c1e2"

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`.

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

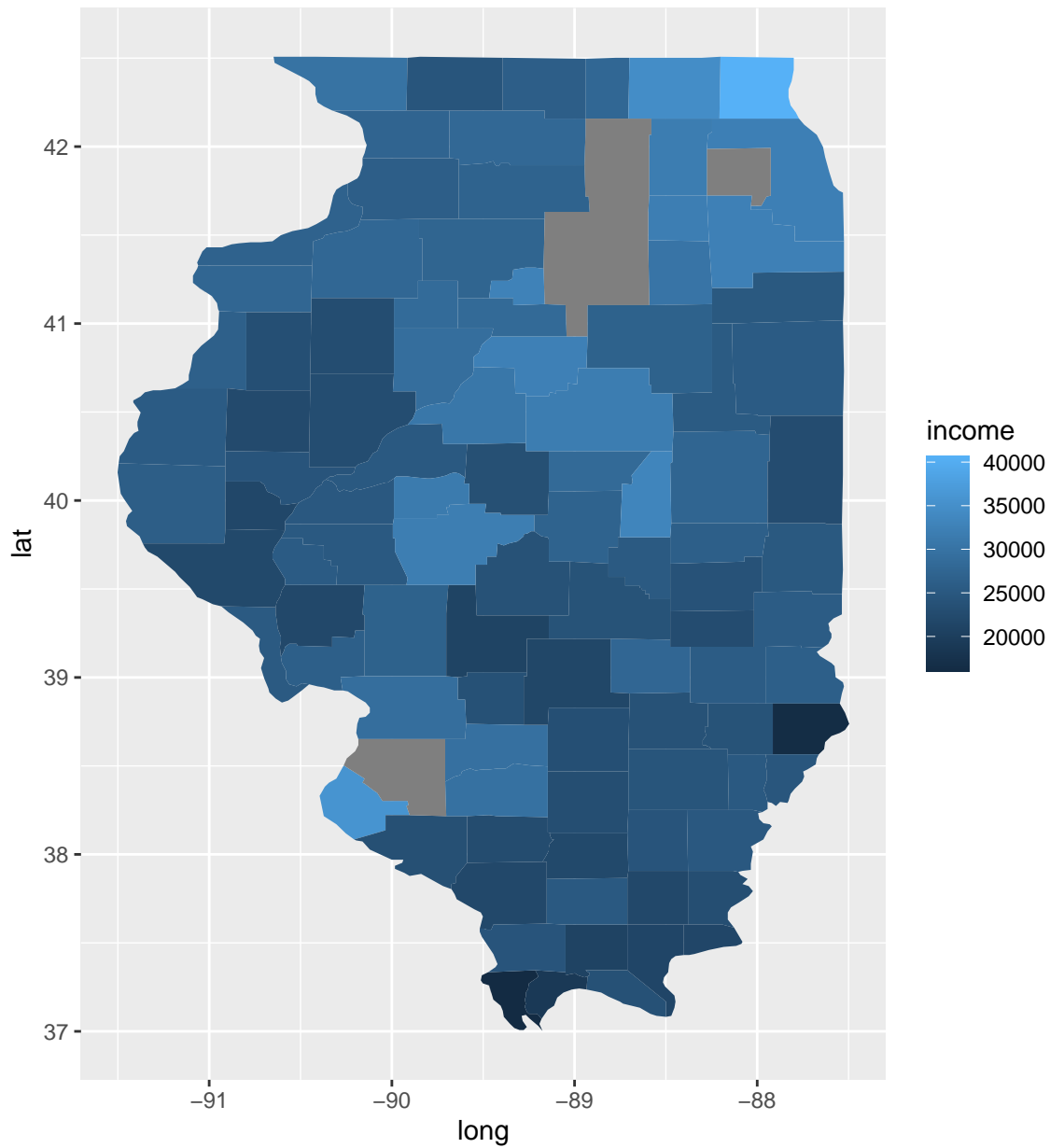
```
acs_il_c <- acs_il_c %>%
  mutate(subregion = tolower(gsub(" County, Illinois", "", NAME)))
acs_map <- il_map %>%
  left_join(acs_il_c, by = "subregion")
head(acs_map)
```

	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

```
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.

```
library(dplyr)

clustering_data <- acs_il_c %>%
  dplyr::select(pop, hh_income, income) %>%
  drop_na() %>%
  dplyr::mutate_all(scale)

head(clustering_data)
```

```
      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
```

```
hclust_d <- dist(clustering_data)
dim(as.matrix(hclust_d))
```

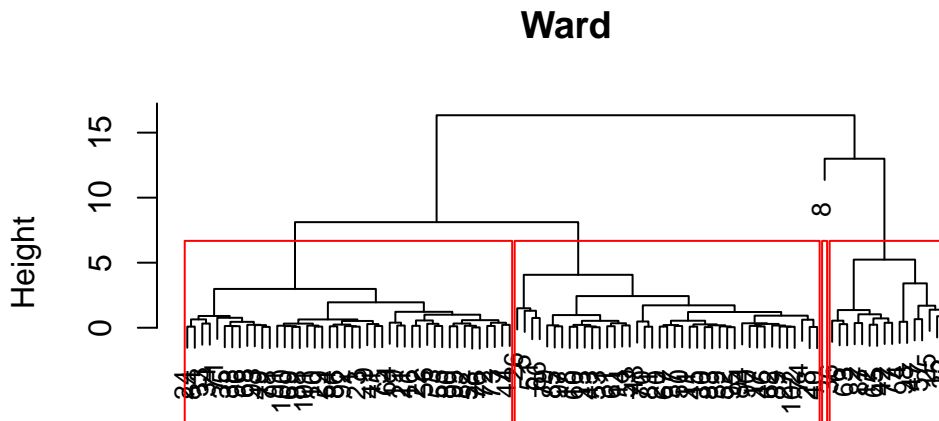
```
[1] 102 102
```

```
as.matrix(hclust_d)[1:5, 1:5]
```

```
      1      2      3      4      5
1 0.0000000 1.996235 0.4986095 5.025852 3.086907
2 1.9962350 0.000000 1.7355417 3.304166 5.042015
3 0.4986095 1.735542 0.0000000 4.806499 3.318745
4 5.0258523 3.304166 4.8064987 0.000000 8.001064
5 3.0869067 5.042015 3.3187446 8.001064 0.000000
```

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

```
hc_ward <- hclust(hclust_d, method = "ward.D2")
plot(hc_ward, main = "Ward", xlab = "", sub = "", cex = 0.8)
rect.hclust(hc_ward, k = 4, border = "red")
```



4 is a reasonable number of clusters.

```
cutree(hc_ward, 4)
```

```
[1] 1 2 1 2 3 2 1 4 1 1 3 2 3 3 2 1 3 3 1 3 1 1 1 1 2 3 3 3 1 3 3 2 1 1 1 2 3
[38] 1 3 3 1 2 3 1 3 3 2 3 1 3 3 3 3 1 1 1 2 1 3 3 3 3 2 1 2 3 3 1 1 3 1 1 1 3
[75] 1 1 2 1 3 1 3 1 3 1 3 1 2 1 3 3 1 1 1 3 3 1 3 2 1 1 1 3
```

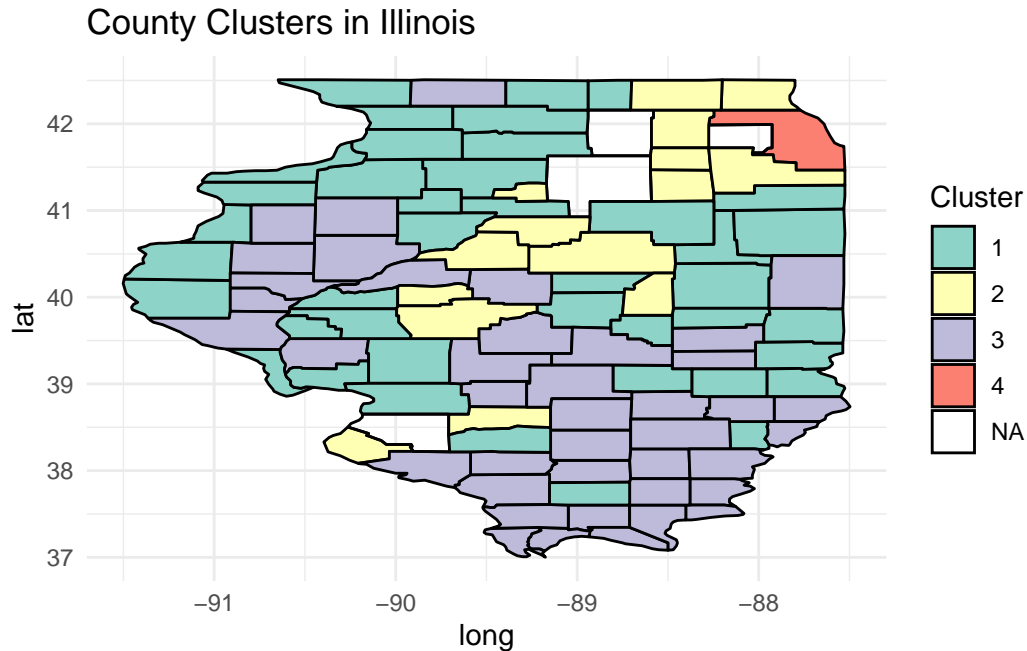
Visualize the county clusters on a map. For this task, create a new `acs_map` object that now also includes cluster membership as a new column. This column should be called `cluster`.

```
acs_il_c <- acs_il_c %>%
  mutate(cluster = cutree(hc_ward, k = 4))
acs_map <- il_map %>%
  left_join(acs_il_c, by = "subregion")
```

```
library(ggplot2)

ggplot(acs_map, aes(long, lat, group = group, fill = factor(cluster))) +
  geom_polygon(color = "black") +
```

```
scale_fill_brewer(palette = "Set3") +
labs(title = "County Clusters in Illinois", fill = "Cluster") +
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(~ ifelse(. == -666666666, NA, .)) %>%
rename(pop = B01003_001E,
  hh_income = B19013_001E,
  income = B19301_001E)

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.

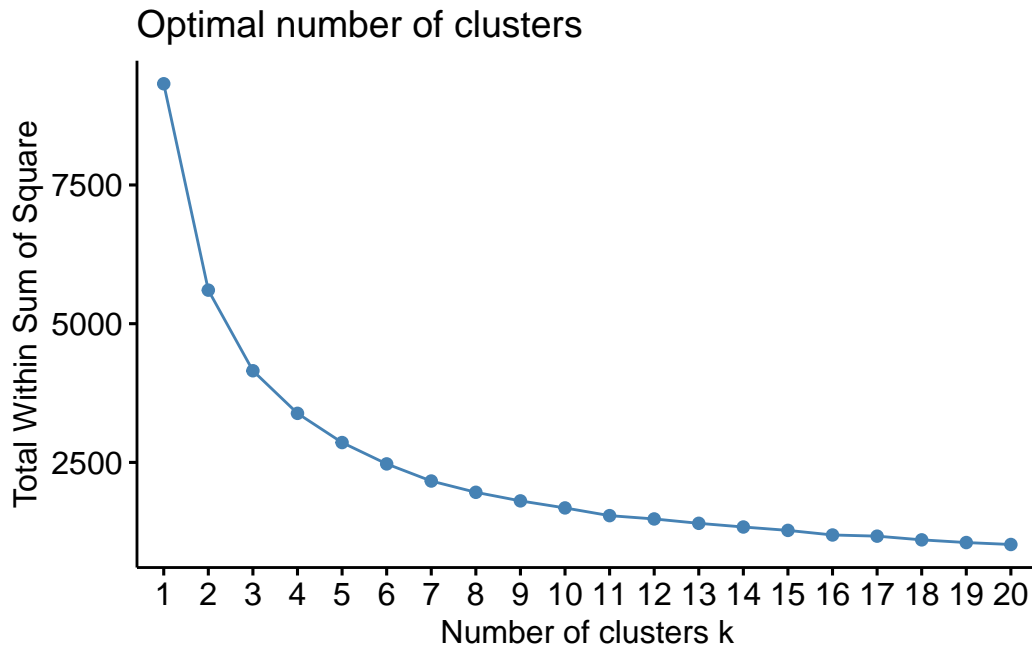
```
kmeans_data <- acs_il_t %>%
  dplyr::select(pop, hh_income, income) %>%
  drop_na(.) %>%
  dplyr::mutate_all(scale)

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)
```



k=4 is likely a good choice for the number of clusters.

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

```
km <- kmeans(kmeans_data, 4, nstart = 20)
```

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.

```
kmeans_data <- kmeans_data %>%
  mutate(cluster = km$cluster)

head(kmeans_data)
```

```
pop    hh_income    income cluster
```

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

```
str(kmeans_data)
```

```
'data.frame':  3109 obs. of  4 variables:
 $ pop      : num [1:3109, 1] 1.619 1.758 -0.738 -0.786 -0.274 ...
 ..- attr(*, "scaled:center")= num 4130
 ..- attr(*, "scaled:scale")= num 1960
 $ hh_income: num [1:3109, 1] -0.141 -0.249 0.053 -0.255 -0.708 ...
 ..- attr(*, "scaled:center")= num 61174
 ..- attr(*, "scaled:scale")= num 29747
 $ income   : num [1:3109, 1] -0.4347 -0.3547 -0.0302 -0.6537 -0.7759 ...
 ..- attr(*, "scaled:center")= num 30632
 ..- attr(*, "scaled:scale")= num 15833
 $ cluster  : int  4 4 3 2 2 3 2 3 2 2 ...
```

```
acs_il_t_re <- acs_il_t %>%
  drop_na() %>%
  mutate(
    pop = scale(pop),
    hh_income = scale(hh_income),
    income = scale(income)
  )
```

```
head(acs_il_t_re)
```

	state	county	tract	NAME	pop
1	17	031	806002	Census Tract 8060.02, Cook County, Illinois	1.6189842
2	17	031	806003	Census Tract 8060.03, Cook County, Illinois	1.7582445
3	17	031	806400	Census Tract 8064, Cook County, Illinois	-0.7377284
4	17	031	806501	Census Tract 8065.01, Cook County, Illinois	-0.7856788
5	17	031	750600	Census Tract 7506, Cook County, Illinois	-0.2735274
6	17	031	310200	Census Tract 3102, Cook County, Illinois	-1.3309874

	hh_income	income
1	-0.14115103	-0.43466339
2	-0.24892639	-0.35470216

```

3  0.05298581 -0.03018336
4 -0.25517911 -0.65370411
5 -0.70759359 -0.77591974
6  0.06979420  0.04870414

```

```
str(acs_il_t_re)
```

```

'data.frame':  3109 obs. of  7 variables:
 $ state      : chr  "17" "17" "17" "17" ...
 $ county     : chr  "031" "031" "031" "031" ...
 $ tract      : chr  "806002" "806003" "806400" "806501" ...
 $ NAME       : chr  "Census Tract 8060.02, Cook County, Illinois" "Census Tract 8060.03, Cook
 $ pop        : num [1:3109, 1] 1.619 1.758 -0.738 -0.786 -0.274 ...
 ..- attr(*, "scaled:center")= num 4130
 ..- attr(*, "scaled:scale")= num 1960
 $ hh_income  : num [1:3109, 1] -0.141 -0.249 0.053 -0.255 -0.708 ...
 ..- attr(*, "scaled:center")= num 61174
 ..- attr(*, "scaled:scale")= num 29747
 $ income     : num [1:3109, 1] -0.4347 -0.3547 -0.0302 -0.6537 -0.7759 ...
 ..- attr(*, "scaled:center")= num 30632
 ..- attr(*, "scaled:scale")= num 15833

```

```

acs_combined <- acs_il_t_re %>%
  left_join(kmeans_data %>%
    dplyr::select(pop, hh_income, income, cluster),
    by = c("pop", "hh_income", "income"))
head(acs_combined)

```

	state	county	tract	NAME	pop
1	17	031	806002	Census Tract 8060.02, Cook County, Illinois	1.6189842
2	17	031	806003	Census Tract 8060.03, Cook County, Illinois	1.7582445
3	17	031	806400	Census Tract 8064, Cook County, Illinois	-0.7377284
4	17	031	806501	Census Tract 8065.01, Cook County, Illinois	-0.7856788
5	17	031	750600	Census Tract 7506, Cook County, Illinois	-0.2735274
6	17	031	310200	Census Tract 3102, Cook County, Illinois	-1.3309874

	hh_income	income	cluster
1	-0.14115103	-0.43466339	4
2	-0.24892639	-0.35470216	4
3	0.05298581	-0.03018336	3
4	-0.25517911	-0.65370411	2
5	-0.70759359	-0.77591974	2
6	0.06979420	0.04870414	3

```
cluster_summary <- acs_combined %>%
  group_by(cluster) %>%
  summarize(
    mean_pop = mean(pop, na.rm = TRUE),
    mean_hh_income = mean(hh_income, na.rm = TRUE),
    mean_income = mean(income, na.rm = TRUE),
    most_frequent_county = names(which.max(table(county)))
  )
cluster_summary
```

```
# A tibble: 4 x 5
  cluster mean_pop mean_hh_income mean_income most_frequent_county
  <int>     <dbl>         <dbl>         <dbl> <chr>
1     1  0.00202         1.99          2.20  031
2     2 -0.507          -0.787        -0.689  031
3     3 -0.187           0.308         0.243  031
4     4  1.47            0.135        -0.0951 031
```

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.

```
kmeans_clustering <- function(k, data) {
  km <- kmeans(data, centers = k, nstart = 25)
  return(km$cluster)
}
```

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.

```
str(acs_il_t %>%
  dplyr::select(pop, hh_income, income))
```

```
'data.frame':  3123 obs. of  3 variables:
 $ pop      : num  7304 7577 2684 2590 3594 ...
 $ hh_income: num  56975 53769 62750 53583 40125 ...
 $ income   : num  23750 25016 30154 20282 18347 ...
```

```

acs_il_t <- acs_il_t %>%
  dplyr::mutate(
    pop = as.numeric(pop),
    hh_income = as.numeric(hh_income),
    income = as.numeric(income)
  )

acs_il_t <- acs_il_t %>%
  filter(!is.na(pop) & !is.na(hh_income) & !is.na(income))

acs_il_t <- acs_il_t %>%
  dplyr::mutate(
    hh_income = ifelse(is.na(hh_income), mean(hh_income, na.rm = TRUE), hh_income),
    income = ifelse(is.na(income), mean(income, na.rm = TRUE), income)
  )

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

```

acs_clustered <- acs_il_t

for (k in 2:10) {
  cluster_col <- kmeans_clustering(k, acs_il_t %>%
    dplyr::select(pop, hh_income, income))
  acs_clustered <- acs_clustered %>%
    dplyr::mutate(!!paste0("cluster_k", k) := cluster_col)
}

```

Warning: Quick-TRANSfer stage steps exceeded maximum (= 155450)

Warning: did not converge in 10 iterations

Warning: did not converge in 10 iterations

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

```
head(acs_clustered)
```

	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_k2	cluster_k3	cluster_k4	cluster_k5	cluster_k6
1	56975	23750	2	2	3	3	4
2	53769	25016	2	2	3	3	6
3	62750	30154	2	3	3	3	4
4	53583	20282	2	2	3	3	6
5	40125	18347	2	2	2	4	6
6	63250	31403	2	3	3	3	4
	cluster_k7	cluster_k8	cluster_k9	cluster_k10			
1	7	7	5	3			
2	3	7	2	3			
3	7	7	5	4			
4	3	7	2	3			
5	3	8	3	7			
6	7	7	5	4			