# R Brown Bag session: tidyverse overview

JMKlein

04-19-2022

```
install.packages("tidyverse", repos = "http://cran.us.r project.org")

install.packages("readxl", repos = "http://cran.us.r project.org")

library(tidyverse)
library(readxl)
```

## Part 1: Exploring your data

### Load 2020 Census Population dataset
```
Census2020 <-  read_excel("2020 Census File.xlsx")
```

### Investigate with glimpse
```
glimpse(Census2020)

## Rows: 51
## Columns: 10
## $ Area                                              <chr> "Alabama", "Alas~
## $ Region                                            <chr> "South", "West",~
## $ `2020 Census Resident Population`                 <dbl> 5024279, 733391,~
## $ `2010 Census Resident Population`                 <dbl> 4779736, 710231,~
## $ `Numeric Change`                                  <dbl> 244543, 23160, 7~
## $ `Percent Change`                                  <dbl> 5.1, 3.3, 11.9, ~
## $ `State Rank Based on 2020 Census Resident Population` <chr> "24", "48", "14"~
## $ `State Rank Based on 2010 Census Resident Population` <chr> "23", "47", "16"~
## $ `State Rank Based on Numeric Change`               <chr> "24", "45", "8",~
## $ `State Rank Based on Percent Change`               <chr> "27", "36", "9",~
```

### Explore the dimensions
```
dim(Census2020)

## [1] 51 10
```

### Display column and row names
```
colnames(Census2020)

##  [1] "Area"
##  [2] "Region"
##  [3] "2020 Census Resident Population"
##  [4] "2010 Census Resident Population"
##  [5] "Numeric Change"
##  [6] "Percent Change"
##  [7] "State Rank Based on 2020 Census Resident Population"
##  [8] "State Rank Based on 2010 Census Resident Population"
##  [9] "State Rank Based on Numeric Change"
## [10] "State Rank Based on Percent Change"
```

```
rownames(Census2020)
```

```
##  [1] "1"  "2"  "3"  "4"  "5"  "6"  "7"  "8"  "9"  "10" "11" "12" "13" "14" "15"
## [16] "16" "17" "18" "19" "20" "21" "22" "23" "24" "25" "26" "27" "28" "29" "30"
## [31] "31" "32" "33" "34" "35" "36" "37" "38" "39" "40" "41" "42" "43" "44" "45"
## [46] "46" "47" "48" "49" "50" "51"
```

## View top and bottom observations

```
head(Census2020)
```

```
## # A tibble: 6 x 10
##   Area     Region `2020 Census Resident ~ `2010 Census Resident~ `Numeric Change`
##   <chr>    <chr>                    <dbl>                  <dbl>            <dbl>
## 1 Alabama  South                  5024279                4779736           244543
## 2 Alaska   West                    733391                 710231            23160
## 3 Arizona  West                   7151502                6392017           759485
## 4 Arkans~  South                  3011524                2915918            95606
## 5 Califo~  West                  39538223               37253956          2284267
## 6 Colora~  West                   5773714                5029196           744518
## # ... with 5 more variables: Percent Change <dbl>,
## #   State Rank Based on 2020 Census Resident Population <chr>,
## #   State Rank Based on 2010 Census Resident Population <chr>,
## #   State Rank Based on Numeric Change <chr>,
## #   State Rank Based on Percent Change <chr>
```

```
tail(Census2020)
```

```
## # A tibble: 6 x 10
##   Area     Region  `2020 Census Resident~ `2010 Census Residen~ `Numeric Change`
##   <chr>    <chr>                    <dbl>                 <dbl>            <dbl>
## 1 Vermont  North                   643077                625741            17336
## 2 Virginia South                  8631393               8001024           630369
## 3 Washing~ West                   7705281               6724540           980741
## 4 West Vi~ South                  1793716               1852994           -59278
## 5 Wiscons~ Midwest                5893718               5686986           206732
## 6 Wyoming  West                    576851                563626            13225
## # ... with 5 more variables: Percent Change <dbl>,
## #   State Rank Based on 2020 Census Resident Population <chr>,
## #   State Rank Based on 2010 Census Resident Population <chr>,
## #   State Rank Based on Numeric Change <chr>,
## #   State Rank Based on Percent Change <chr>
```

## Explore largest and smallest values in a column

```
max(Census2020$`2020 Census Resident Population`)
```

```
## [1] 39538223
```

```
min(Census2020$`2020 Census Resident Population`)
```

```
## [1] 576851
```

## Display summary stats

```
summary(Census2020)
```

```
##      Area              Region          2020 Census Resident Population
##  Length:51          Length:51          Min.   :  576851
##  Class :character   Class :character   1st Qu.: 1816411
```

```
##   Mode  :character   Mode  :character   Median : 4505836
##                                         Mean   : 6499006
##                                         3rd Qu.: 7428392
##                                         Max.   :39538223
##
##   2010 Census Resident Population Numeric Change     Percent Change
##   Min.   :  563626                Min.   : -59278   Min.   :-3.200
##   1st Qu.: 1696962                1st Qu.:  86292   1st Qu.: 2.900
##   Median : 4339367                Median : 206732   Median : 5.700
##   Mean   : 6053834                Mean   : 445171   Mean   : 7.024
##   3rd Qu.: 6636084                3rd Qu.: 495080   3rd Qu.:10.400
##   Max.   :37253956                Max.   :3999944   Max.   :18.400
##
##   State Rank Based on 2020 Census Resident Population
##   Length:51
##   Class :character
##   Mode  :character
##
##
##
##   State Rank Based on 2010 Census Resident Population
##   Length:51
##   Class :character
##   Mode  :character
##
##
##
##   State Rank Based on Numeric Change State Rank Based on Percent Change
##   Length:51                          Length:51
##   Class :character                   Class :character
##   Mode  :character                   Mode  :character
##
##
##
```

### Open and explore the dataset in a new pane- with filtering options

```
View(Census2020)
```

### Identify a column

```
Census2020$`2020 Census Resident Population`
```

```
##  [1]  5024279   733391  7151502  3011524 39538223  5773714  3605944   989948
##  [9]   689545 21538187 10711908  1455271  1839106 12812508  6785528  3190369
## [17]  2937880  4505836  4657757  1362359  6177224  7029917 10077331  5706494
## [25]  2961279  6154913  1084225  1961504  3104614  1377529  9288994  2117522
## [33] 20201249 10439388   779094 11799448  3959353  4237256 13002700  1097379
## [41]  5118425   886667  6910840 29145505  3271616   643077  8631393  7705281
## [49]  1793716  5893718   576851
```

```
Census2020$Region
```

```
##  [1] "South"   "West"    "West"    "South"   "West"    "West"    "North"
##  [8] "South"   "South"   "South"   "South"   "West"    "West"    "Midwest"
## [15] "Midwest" "Midwest" "Midwest" "South"   "South"   "North"   "South"
## [22] "North"   "Midwest" "Midwest" "South"   "Midwest" "West"    "Midwest"
## [29] "West"    "North"   "North"   "West"    "North"   "South"   "Midwest"
```

```
## [36] "Midwest" "South"    "West"     "North"    "North"    "South"    "Midwest"
## [43] "South"   "South"    "West"     "North"    "South"    "West"     "South"
## [50] "Midwest" "West"
```

## Display contents of column as a table

```
table(Census2020$Region)

##
## Midwest   North   South    West
##      12       9      17      13

table(Census2020$Area, Census2020$Region)

##
##                        Midwest North South West
##    Alabama                   0     0     1    0
##    Alaska                    0     0     0    1
##    Arizona                   0     0     0    1
##    Arkansas                  0     0     1    0
##    California                0     0     0    1
##    Colorado                  0     0     0    1
##    Connecticut               0     1     0    0
##    Delaware                  0     0     1    0
##    District of Columbia      0     0     1    0
##    Florida                   0     0     1    0
##    Georgia                   0     0     1    0
##    Hawaii                    0     0     0    1
##    Idaho                     0     0     0    1
##    Illinois                  1     0     0    0
##    Indiana                   1     0     0    0
##    Iowa                      1     0     0    0
##    Kansas                    1     0     0    0
##    Kentucky                  0     0     1    0
##    Louisiana                 0     0     1    0
##    Maine                     0     1     0    0
##    Maryland                  0     0     1    0
##    Massachusetts             0     1     0    0
##    Michigan                  1     0     0    0
##    Minnesota                 1     0     0    0
##    Mississippi               0     0     1    0
##    Missouri                  1     0     0    0
##    Montana                   0     0     0    1
##    Nebraska                  1     0     0    0
##    Nevada                    0     0     0    1
##    New Hampshire             0     1     0    0
##    New Jersey                0     1     0    0
##    New Mexico                0     0     0    1
##    New York                  0     1     0    0
##    North Carolina            0     0     1    0
##    North Dakota              1     0     0    0
##    Ohio                      1     0     0    0
##    Oklahoma                  0     0     1    0
##    Oregon                    0     0     0    1
##    Pennsylvania              0     1     0    0
##    Rhode Island              0     1     0    0
##    South Carolina            0     0     1    0
```

```
##    South Dakota                    1      0      0      0
##    Tennessee                       0      0      1      0
##    Texas                           0      0      1      0
##    Utah                            0      0      0      1
##    Vermont                         0      1      0      0
##    Virginia                        0      0      1      0
##    Washington                      0      0      0      1
##    West Virginia                   0      0      1      0
##    Wisconsin                       1      0      0      0
##    Wyoming                         0      0      0      1
```

## Identify an exact position, [rows, columns]

```
Census2020[,1]

## # A tibble: 51 x 1
##    Area
##    <chr>
##  1 Alabama
##  2 Alaska
##  3 Arizona
##  4 Arkansas
##  5 California
##  6 Colorado
##  7 Connecticut
##  8 Delaware
##  9 District of Columbia
## 10 Florida
## # ... with 41 more rows

Census2020[1,]

## # A tibble: 1 x 10
##   Area    Region `2020 Census Resident ~ `2010 Census Resident ~ `Numeric Change`
##   <chr>  <chr>                    <dbl>                    <dbl>            <dbl>
## 1 Alaba~ South                  5024279                  4779736           244543
## # ... with 5 more variables: Percent Change <dbl>,
## #   State Rank Based on 2020 Census Resident Population <chr>,
## #   State Rank Based on 2010 Census Resident Population <chr>,
## #   State Rank Based on Numeric Change <chr>,
## #   State Rank Based on Percent Change <chr>

Census2020[1,1]

## # A tibble: 1 x 1
##   Area
##   <chr>
## 1 Alabama
```

## Export to csv

```
write.csv(Census2020, "Census2020.csv")
```

# Part 2: Manipulate and transform with Tidyverse: intro to dplyr commands using select, rename, filter, arrange, mutate, summarize

## Read-in two ACS files: 2019 population and 2019 poverty rate

```
Census2019 <-  read_csv("2019Pop.csv")

##
## -- Column specification ---------------------------------------------------
## cols(
##   State = col_character(),
##   Estimate = col_double()
## )

Poverty2019 <-  read_csv("2019Poverty.csv")

##
## -- Column specification ---------------------------------------------------
## cols(
##   State = col_character(),
##   PovertyStatus = col_double(),
##   BelowPoverty = col_double(),
##   AbovePoverty = col_double()
## )
```

## Use the select function to keep/select the columns: state name, region, 2020 population, numeric change, percent change, and state rank

```
Census2020Sub1 <-  Census2020 %>%
  select(`Area`,
         `Region`,
         `2020 Census Resident Population`,
         `Numeric Change`,
         `Percent Change`,
         `State Rank Based on 2020 Census Resident Population`)
```

## View the subsetted object

```
Census2020Sub1

## # A tibble: 51 x 6
##    Area          Region `2020 Census Resident ~ `Numeric Change` `Percent Change`
##    <chr>         <chr>                    <dbl>            <dbl>            <dbl>
##  1 Alabama       South                  5024279           244543              5.1
##  2 Alaska        West                    733391            23160              3.3
##  3 Arizona       West                   7151502           759485             11.9
##  4 Arkansas      South                  3011524            95606              3.3
##  5 California    West                  39538223          2284267              6.1
##  6 Colorado      West                   5773714           744518             14.8
##  7 Connecticut   North                  3605944            31847              0.9
##  8 Delaware      South                   989948            92014             10.2
##  9 District of~  South                   689545            87822             14.6
## 10 Florida       South                 21538187          2736877             14.6
## # ... with 41 more rows, and 1 more variable:
## #   State Rank Based on 2020 Census Resident Population <chr>
```

### Use the rename function to rename columns to easy to work with names

```
Census2020Sub1 <-  Census2020Sub1 %>%
  rename(State = Area,
         Pop2020 = `2020 Census Resident Population`,
         NumChange2020 = `Numeric Change`,
         PercentChange2020 = `Percent Change`,
         StateRank = `State Rank Based on 2020 Census Resident Population`)
```

### View new column names

```
str(Census2020Sub1)

## tibble [51 x 6] (S3: tbl_df/tbl/data.frame)
##  $ State            : chr [1:51] "Alabama" "Alaska" "Arizona" "Arkansas" ...
##  $ Region           : chr [1:51] "South" "West" "West" "South" ...
##  $ Pop2020          : num [1:51] 5024279 733391 7151502 3011524 39538223 ...
##  $ NumChange2020    : num [1:51] 244543 23160 759485 95606 2284267 ...
##  $ PercentChange2020: num [1:51] 5.1 3.3 11.9 3.3 6.1 14.8 0.9 10.2 14.6 14.6 ...
##  $ StateRank        : chr [1:51] "24" "48" "14" "33" ...
```

### Use the filter function to subset rows by pop size, using 9999999 as the limit

```
PopAboveLimit <-  Census2020Sub1 %>%
  filter(Pop2020 > 9999999)


PopBelowLimit <-  Census2020Sub1 %>%
  filter(Pop2020 <= 9999999)
```

### View dimenstions of the new objects

```
dim(PopAboveLimit)

## [1] 10  6

dim(PopBelowLimit)

## [1] 41  6
```

### Use filter to subset rows by two conditions, using population and state rank

- Use a population limit of 9999999 and state rank limits to narrow down data

```
PopAboveLimitAND <-  Census2020Sub1 %>%
  filter(Pop2020 > 9999999 & StateRank >= 9)


PopAboveLimitOR <-  Census2020Sub1 %>%
  filter(Pop2020 > 9999999 | StateRank >= 9)
```

### View the contents of the new object

```
glimpse(PopAboveLimitAND)

## Rows: 1
## Columns: 6
## $ State             <chr> "North Carolina"
## $ Region            <chr> "South"
## $ Pop2020           <dbl> 10439388
## $ NumChange2020     <dbl> 903905
## $ PercentChange2020 <dbl> 9.5
## $ StateRank         <chr> "9"
```

```
glimpse(PopAboveLimitOR)

## Rows: 11
## Columns: 6
## $ State           <chr> "California", "District of Columbia", "Florida", "Ge~
## $ Region          <chr> "West", "South", "South", "South", "Midwest", "Midwe~
## $ Pop2020         <dbl> 39538223, 689545, 21538187, 10711908, 12812508, 1007~
## $ NumChange2020   <dbl> 2284267, 87822, 2736877, 1024255, -18124, 193691, 82~
## $ PercentChange2020 <dbl> 6.1, 14.6, 14.6, 10.6, -0.1, 2.0, 4.2, 9.5, 2.3, 2.4~
## $ StateRank       <chr> "1", "X", "3", "8", "6", "10", "4", "9", "7", "5", "~
```

## Convert state rank from integer to numeric

```
str(Census2020Sub1$StateRank)

##  chr [1:51] "24" "48" "14" "33" "1" "21" "29" "45" "X" "3" "8" "40" "38" ...

Census2020Sub1$StateRank <-  as.numeric(Census2020Sub1$StateRank, na.rm = TRUE)

## Warning: NAs introduced by coercion
```

## Use the arrange function to sort the two population objects by state rank

- Order the filtered objects by ascending

```
TopPopAsce <-  PopAboveLimit %>%
  arrange(StateRank)


LowPopAsce <-  PopBelowLimit %>%
  arrange(StateRank)
```

## View new object containing large states arranged by state rank- ascending

```
head(TopPopAsce)

## # A tibble: 6 x 6
##    State         Region   Pop2020 NumChange2020 PercentChange2020 StateRank
##    <chr>         <chr>     <dbl>        <dbl>              <dbl> <chr>
## 1 California    West     39538223      2284267               6.1 1
## 2 Michigan      Midwest 10077331       193691               2   10
## 3 Texas         South    29145505      3999944              15.9 2
## 4 Florida       South    21538187      2736877              14.6 3
## 5 New York      North    20201249       823147               4.2 4
## 6 Pennsylvania North     13002700       300321               2.4 5
```

## View new object containing small states arranged by state rank- ascending

```
head(LowPopAsce)

## # A tibble: 6 x 6
##    State          Region Pop2020 NumChange2020 PercentChange2020 StateRank
##    <chr>          <chr>   <dbl>        <dbl>              <dbl> <chr>
## 1 New Jersey     North  9288994       497100               5.7 11
## 2 Virginia       South  8631393       630369               7.9 12
## 3 Washington     West   7705281       980741              14.6 13
## 4 Arizona        West   7151502       759485              11.9 14
## 5 Massachusetts  North  7029917       482288               7.4 15
## 6 Tennessee      South  6910840       564735               8.9 16
```

## Use the arrange function to sort the two population objects by state rank

- Order the filtered objects by descending

```
TopPopDesc <-  PopAboveLimit %>%
  arrange(desc(StateRank))

LowPopDesc <-  PopBelowLimit %>%
  arrange(desc(StateRank))
```

## View new object with large states arranged by state rank- descending

```
head(TopPopDesc)

## # A tibble: 6 x 6
##   State          Region  Pop2020 NumChange2020 PercentChange2020 StateRank
##   <chr>          <chr>     <dbl>         <dbl>             <dbl> <chr>
## 1 North Carolina South   10439388        903905               9.5 9
## 2 Georgia        South   10711908       1024255              10.6 8
## 3 Ohio           Midwest 11799448        262944               2.3 7
## 4 Illinois       Midwest 12812508        -18124              -0.1 6
## 5 Pennsylvania   North   13002700        300321               2.4 5
## 6 New York       North   20201249        823147               4.2 4
```

## View new object with small states arranged by state rank- descending

```
head(LowPopDesc)

## # A tibble: 6 x 6
##   State                Region  Pop2020 NumChange2020 PercentChange2020 StateRank
##   <chr>                <chr>     <dbl>         <dbl>             <dbl> <chr>
## 1 District of Columbia South   689545         87822              14.6 X
## 2 Wyoming              West    576851         13225               2.3 50
## 3 Vermont              North   643077         17336               2.8 49
## 4 Alaska               West    733391         23160               3.3 48
## 5 North Dakota         Midwest 779094        106503              15.8 47
## 6 South Dakota         Midwest 886667         72487               8.9 46
```

## Use the mutate function to add a new column

- Calculate the 2010 pop using the 2020 pop and numeric change columns

```
Census2020Mutate <-  Census2020Sub1 %>%
  mutate(Pop2010 = Pop2020 - NumChange2020)
```

## View top observations of new object

```
head(Census2020Mutate)

## # A tibble: 6 x 7
##   State       Region Pop2020 NumChange2020 PercentChange2020 StateRank  Pop2010
##   <chr>       <chr>    <dbl>         <dbl>             <dbl>    <dbl>    <dbl>
## 1 Alabama     South  5024279        244543               5.1       24  4779736
## 2 Alaska      West    733391         23160               3.3       48   710231
## 3 Arizona     West   7151502        759485              11.9       14  6392017
## 4 Arkansas    South  3011524         95606               3.3       33  2915918
## 5 California  West  39538223       2284267               6.1        1 37253956
## 6 Colorado    West   5773714        744518              14.8       21  5029196
```

## Use the summarise function to determine the total population in the US across all states, for 2020 and 2010

- 2020

```
Census2020PopSum <-  Census2020Mutate %>%
  summarise(Total2020 = sum(Pop2020))
```

  • 2010

```
Census2010PopSum <-  Census2020Mutate %>%
  summarise(Total2010 = sum(Pop2010))
```

## View new objects with totals of 2020 and 2010 population size

  • 2020

```
Census2020PopSum
```

```
## # A tibble: 1 x 1
##   Total2020
##       <dbl>
## 1 331449281
```

  • 2010

```
Census2010PopSum
```

```
## # A tibble: 1 x 1
##   Total2010
##       <dbl>
## 1 308745538
```

## Use the summarise function to determine the total population in the US across all states, for 2020 and 2010. Include group_by region

  • 2020

```
Census2020PopbyRegion <-  Census2020Mutate %>%
  group_by(Region) %>%
  summarise(Total2020 = sum(Pop2020))
```

  • 2010

```
Census2010PopbyRegion <-  Census2020Mutate %>%
  group_by(Region) %>%
  summarise(Total2010 = sum(Pop2010))
```

## View new objects with totals of 2020 and 2010 population size, grouped by region

  • 2020

```
Census2020PopbyRegion
```

```
## # A tibble: 4 x 2
##   Region   Total2020
##   <chr>        <dbl>
## 1 Midwest  68985454
## 2 North    57609148
## 3 South    126266107
## 4 West     78588572
```

  • 2010

```
Census2010PopbyRegion
```

```
## # A tibble: 4 x 2
##   Region   Total2010
##   <chr>        <dbl>
## 1 Midwest  66927001
```

```
## 2 North      55317240
## 3 South     114555744
## 4 West       71945553
```

## Calculate the average national population for 2020 and 2010, include group_by region

- 2020

```
Census2020PopbyRegion <-  Census2020Mutate %>%
  group_by(Region) %>%
  summarize(Total2020 = mean(Pop2020))
```

- 2010

```
Census2010PopbyRegion <-  Census2020Mutate %>%
  group_by(Region) %>%
  summarize(Total2010 = mean(Pop2010))
```

## View new objects with averages of 2020 and 2010 population size, grouped by region

- 2020

```
Census2020PopbyRegion

## # A tibble: 4 x 2
##    Region  Total2020
##    <chr>        <dbl>
## 1 Midwest  5748788.
## 2 North     6401016.
## 3 South     7427418.
## 4 West      6045275.
```

- 2010

```
Census2010PopbyRegion

## # A tibble: 4 x 2
##    Region  Total2010
##    <chr>        <dbl>
## 1 Midwest  5577250.
## 2 North     6146360
## 3 South     6738573.
## 4 West      5534273.
```

## Calculate the sum of large states, include group_by region

```
PopAboveLimitbyRegion <-  PopAboveLimit %>%
  group_by(Region) %>%
  summarize(TotalLarge2020 = sum(Pop2020))
```

## View new object with total population of large states, grouped by region

```
PopAboveLimitbyRegion

## # A tibble: 4 x 2
##    Region  TotalLarge2020
##    <chr>              <dbl>
## 1 Midwest         34689287
## 2 North            33203949
## 3 South            71834988
## 4 West             39538223
```

## Calculate the sum of small states, include group_by region

- Use the object PopBelowLimit

```
PopBelowLimitbyRegion <- PopBelowLimit %>%
  group_by(Region) %>%
  summarize(TotalSmall2020 = sum(Pop2020))
```

## View new object with total population of small states, grouped by region

```
PopBelowLimitbyRegion

## # A tibble: 4 x 2
##   Region  TotalSmall2020
##   <chr>            <dbl>
## 1 Midwest       34296167
## 2 North         24405199
## 3 South         54431119
## 4 West          39050349
```

## Examples of combining multiple dplyr verbs in one workflow - You can use all of the verbs chained together in logical order to achieve complex results

## Utilize select and rename functions in one workflow

```
Census2020Bonus <-  Census2020 %>%
  select(`Area`,
         `2020 Census Resident Population`,
         `2010 Census Resident Population`,
         `State Rank Based on 2020 Census Resident Population`) %>%
  rename(State = Area,
         Pop2020 = `2020 Census Resident Population`,
         Pop2010 = `2010 Census Resident Population`,
         StateRank = `State Rank Based on 2020 Census Resident Population`)
```

## View top observations of new object

```
head(Census2020Bonus)

## # A tibble: 6 x 4
##   State         Pop2020  Pop2010 StateRank
##   <chr>           <dbl>    <dbl> <chr>
## 1 Alabama       5024279  4779736 24
## 2 Alaska         733391   710231 48
## 3 Arizona       7151502  6392017 14
## 4 Arkansas      3011524  2915918 33
## 5 California   39538223 37253956 1
## 6 Colorado      5773714  5029196 21
```

## Utilize filter and arrange in one workflow

```
Census2020Bonus1 <-  Census2020Bonus %>%
  filter(StateRank >= 2 & StateRank <= 50) %>%
  arrange(desc(Pop2020))
```

## View glimpse of new object

```
glimpse(Census2020Bonus1)

## Rows: 35
## Columns: 4
```

```
## $ State      <chr> "Texas", "Florida", "New York", "Pennsylvania", "Wisconsin",~
## $ Pop2020    <dbl> 29145505, 21538187, 20201249, 13002700, 5893718, 5773714, 57~
## $ Pop2010    <dbl> 25145561, 18801310, 19378102, 12702379, 5686986, 5029196, 53~
## $ StateRank <chr> "2", "3", "4", "5", "20", "21", "22", "23", "24", "25", "26"~
```

## Combine the mutate and summarize functions in one workflow

- Sum the population of top largest and smallest states using prior object

```
Census2020Bonus2 <-  Census2020Bonus1 %>%
  mutate(size = case_when(Pop2020 > 9999999 ~ 'Big',
                          Pop2020 <= 9999999 ~ 'Small')) %>%
  group_by(size) %>%
  summarize(Total2020 = sum(Pop2020))
```

## View glimpse of new object

```
glimpse(Census2020Bonus2)

## Rows: 2
## Columns: 2
## $ size      <chr> "Big", "Small"
## $ Total2020 <dbl> 83887641, 85657697
```

## Put it all together

```
Census2020Workflow <-  Census2020 %>%
  select(`Area`,
         `2020 Census Resident Population`,
         `2010 Census Resident Population`,
         `State Rank Based on 2020 Census Resident Population`) %>%
  rename(State = Area,
         Pop2020 = `2020 Census Resident Population`,
         Pop2010 = `2010 Census Resident Population`,
         StateRank = `State Rank Based on 2020 Census Resident Population`) %>%
  filter(StateRank >= 2 & StateRank <= 50) %>%
  arrange(desc(Pop2020)) %>%
  mutate(size = case_when(Pop2020 > 9999999 ~ 'Big',
                          Pop2020 <= 9999999 ~ 'Small')) %>%
  group_by(size) %>%
  summarize(Total2020 = sum(Pop2020))
```

## View outcome, it is the same as the workflow seen prior

```
Census2020Workflow

## # A tibble: 2 x 2
##   size   Total2020
##   <chr>      <dbl>
## 1 Big     83887641
## 2 Small   85657697
```

## Join 2020 Census with 2019 ACS Population, by state

```
CensusData1 <-  left_join(Census2020Sub1, Census2019, by = "State")
```

## View new joined object

```
head(CensusData1)

## # A tibble: 6 x 7
##   State     Region  Pop2020 NumChange2020 PercentChange2020 StateRank Estimate
```

```
##    <chr>     <chr>        <dbl>           <dbl>             <dbl>    <dbl>     <dbl>
## 1 Alabama   South     5024279          244543               5.1       24   4876250
## 2 Alaska    West       733391           23160               3.3       48    737068
## 3 Arizona   West      7151502          759485              11.9       14   7050299
## 4 Arkansas  South     3011524           95606               3.3       33   2999370
## 5 California West     39538223         2284267               6.1        1  39283497
## 6 Colorado  West      5773714          744518              14.8       21   5610349
```

### Join 2020 and 2019 population object with 2019 ACS Poverty, by state

- Use rename function to change generic "estimate" column to something specific before join

```
CensusData1 <-  CensusData1 %>%
  rename(PopEstimate2019 = Estimate)


CensusData2 <-  left_join(CensusData1, Poverty2019, by = "State")
```

### View top observations of the new object

```
head(CensusData2)
```

```
## # A tibble: 6 x 10
##    State  Region Pop2020 NumChange2020 PercentChange20~ StateRank PopEstimate2019
##    <chr>  <chr>    <dbl>         <dbl>            <dbl>     <dbl>           <dbl>
## 1 Alaba~ South    5.02e6        244543              5.1        24         4876250
## 2 Alaska West     7.33e5         23160              3.3        48          737068
## 3 Arizo~ West     7.15e6        759485             11.9        14         7050299
## 4 Arkan~ South    3.01e6         95606              3.3        33         2999370
## 5 Calif~ West     3.95e7       2284267              6.1         1        39283497
## 6 Color~ West     5.77e6        744518             14.8        21         5610349
## # ... with 3 more variables: PovertyStatus <dbl>, BelowPoverty <dbl>,
## #   AbovePoverty <dbl>
```

### Use filter and mutate functions to add a ranking variable for states based on below poverty variable

```
CensusDataRanked <-  CensusData2 %>%
  mutate(PovertyRank = dense_rank(desc(BelowPoverty))) %>%
  filter(PovertyRank <= 10)
```

### View a glimpse of new object

```
glimpse(CensusDataRanked)
```

```
## Rows: 10
## Columns: 11
## $ State            <chr> "California", "Florida", "Georgia", "Illinois", "Mic~
## $ Region           <chr> "West", "South", "South", "Midwest", "Midwest", "Nor~
## $ Pop2020          <dbl> 39538223, 21538187, 10711908, 12812508, 10077331, 20~
## $ NumChange2020    <dbl> 2284267, 2736877, 1024255, -18124, 193691, 823147, 9~
## $ PercentChange2020 <dbl> 6.1, 14.6, 10.6, -0.1, 2.0, 4.2, 9.5, 2.3, 2.4, 15.9
## $ StateRank        <dbl> 1, 3, 8, 6, 10, 4, 9, 7, 5, 2
## $ PopEstimate2019  <dbl> 39283497, 20901636, 10403847, 12770631, 9965265, 195~
## $ PovertyStatus    <dbl> 38733295, 21048884, 10332523, 12373209, 9772151, 189~
## $ BelowPoverty     <dbl> 4552837, 2664772, 1373909, 1420542, 1269062, 2467006~
## $ AbovePoverty     <dbl> 34180458, 18384112, 8958614, 10952667, 8503089, 1646~
## $ PovertyRank      <int> 1, 3, 9, 7, 10, 4, 8, 6, 5, 2
```

```
glimpse(CensusDataRanked$PovertyRank)
```

```
##  int [1:10] 1 3 9 7 10 4 8 6 5 2
```

**Visualize using ggplot**

```
ggplot(CensusDataRanked) +
  geom_bar(mapping = aes(x = reorder(State,  BelowPoverty),
                         y = BelowPoverty,
                         fill = PercentChange2020),
           stat = 'identity') +
  labs(title = "Top 10 Most Populated States in 2020",
       x = "State",
       y = "Population Below Poverty") +
  coord_flip()
```

# Part 3: Explore with Tidycensus and API

## API Key and load Tidycensus package

```
library(tidycensus)

census_api_key("INSERT YOUR API KEY HERE")

## To install your API key for use in future sessions, run this function with `install =
TRUE`.
```

## Search for Variables

```
vars <-  load_variables(2020, "pl")

print(tbl_df(vars), n=301)

## Warning: `tbl_df()` was deprecated in dplyr 1.0.0.
## Please use `tibble::as_tibble()` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_warnings()` to see where this warning was generated.

## # A tibble: 301 x 3
##      name   label                              concept
##      <chr>  <chr>                              <chr>
##    1 H1_00~ " !!Total:"                        OCCUPANCY STATUS
##    2 H1_00~ " !!Total:!!Occupied"              OCCUPANCY STATUS
##    3 H1_00~ " !!Total:!!Vacant"                OCCUPANCY STATUS
##    4 P1_00~ " !!Total:"                        RACE
##    5 P1_00~ " !!Total:!!Population of one rac~ RACE
##    6 P1_00~ " !!Total:!!Population of one rac~ RACE
##    7 P1_00~ " !!Total:!!Population of one rac~ RACE
##    8 P1_00~ " !!Total:!!Population of one rac~ RACE
##    9 P1_00~ " !!Total:!!Population of one rac~ RACE
##   10 P1_00~ " !!Total:!!Population of one rac~ RACE
##   11 P1_00~ " !!Total:!!Population of one rac~ RACE
##   12 P1_00~ " !!Total:!!Population of two or ~ RACE
##   13 P1_01~ " !!Total:!!Population of two or ~ RACE
##   14 P1_01~ " !!Total:!!Population of two or ~ RACE
##   15 P1_01~ " !!Total:!!Population of two or ~ RACE
##   16 P1_01~ " !!Total:!!Population of two or ~ RACE
##   17 P1_01~ " !!Total:!!Population of two or ~ RACE
##   18 P1_01~ " !!Total:!!Population of two or ~ RACE
##   19 P1_01~ " !!Total:!!Population of two or ~ RACE
##   20 P1_01~ " !!Total:!!Population of two or ~ RACE
##   21 P1_01~ " !!Total:!!Population of two or ~ RACE
##   22 P1_01~ " !!Total:!!Population of two or ~ RACE
##   23 P1_02~ " !!Total:!!Population of two or ~ RACE
##   24 P1_02~ " !!Total:!!Population of two or ~ RACE
##   25 P1_02~ " !!Total:!!Population of two or ~ RACE
##   26 P1_02~ " !!Total:!!Population of two or ~ RACE
##   27 P1_02~ " !!Total:!!Population of two or ~ RACE
##   28 P1_02~ " !!Total:!!Population of two or ~ RACE
##   29 P1_02~ " !!Total:!!Population of two or ~ RACE
##   30 P1_02~ " !!Total:!!Population of two or ~ RACE
##   31 P1_02~ " !!Total:!!Population of two or ~ RACE
##   32 P1_02~ " !!Total:!!Population of two or ~ RACE
```

```
## 33 P1_03~ " !!Total:!!Population of two or ~ RACE
## 34 P1_03~ " !!Total:!!Population of two or ~ RACE
## 35 P1_03~ " !!Total:!!Population of two or ~ RACE
## 36 P1_03~ " !!Total:!!Population of two or ~ RACE
## 37 P1_03~ " !!Total:!!Population of two or ~ RACE
## 38 P1_03~ " !!Total:!!Population of two or ~ RACE
## 39 P1_03~ " !!Total:!!Population of two or ~ RACE
## 40 P1_03~ " !!Total:!!Population of two or ~ RACE
## 41 P1_03~ " !!Total:!!Population of two or ~ RACE
## 42 P1_03~ " !!Total:!!Population of two or ~ RACE
## 43 P1_04~ " !!Total:!!Population of two or ~ RACE
## 44 P1_04~ " !!Total:!!Population of two or ~ RACE
## 45 P1_04~ " !!Total:!!Population of two or ~ RACE
## 46 P1_04~ " !!Total:!!Population of two or ~ RACE
## 47 P1_04~ " !!Total:!!Population of two or ~ RACE
## 48 P1_04~ " !!Total:!!Population of two or ~ RACE
## 49 P1_04~ " !!Total:!!Population of two or ~ RACE
## 50 P1_04~ " !!Total:!!Population of two or ~ RACE
## 51 P1_04~ " !!Total:!!Population of two or ~ RACE
## 52 P1_04~ " !!Total:!!Population of two or ~ RACE
## 53 P1_05~ " !!Total:!!Population of two or ~ RACE
## 54 P1_05~ " !!Total:!!Population of two or ~ RACE
## 55 P1_05~ " !!Total:!!Population of two or ~ RACE
## 56 P1_05~ " !!Total:!!Population of two or ~ RACE
## 57 P1_05~ " !!Total:!!Population of two or ~ RACE
## 58 P1_05~ " !!Total:!!Population of two or ~ RACE
## 59 P1_05~ " !!Total:!!Population of two or ~ RACE
## 60 P1_05~ " !!Total:!!Population of two or ~ RACE
## 61 P1_05~ " !!Total:!!Population of two or ~ RACE
## 62 P1_05~ " !!Total:!!Population of two or ~ RACE
## 63 P1_06~ " !!Total:!!Population of two or ~ RACE
## 64 P1_06~ " !!Total:!!Population of two or ~ RACE
## 65 P1_06~ " !!Total:!!Population of two or ~ RACE
## 66 P1_06~ " !!Total:!!Population of two or ~ RACE
## 67 P1_06~ " !!Total:!!Population of two or ~ RACE
## 68 P1_06~ " !!Total:!!Population of two or ~ RACE
## 69 P1_06~ " !!Total:!!Population of two or ~ RACE
## 70 P1_06~ " !!Total:!!Population of two or ~ RACE
## 71 P1_06~ " !!Total:!!Population of two or ~ RACE
## 72 P1_06~ " !!Total:!!Population of two or ~ RACE
## 73 P1_07~ " !!Total:!!Population of two or ~ RACE
## 74 P1_07~ " !!Total:!!Population of two or ~ RACE
## 75 P2_00~ " !!Total:"                     HISPANIC OR LATINO, AND NOT HISPA~
## 76 P2_00~ " !!Total:!!Hispanic or Latino"    HISPANIC OR LATINO, AND NOT HISPA~
## 77 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 78 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 79 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 80 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 81 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 82 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 83 P2_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 84 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 85 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 86 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 87 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
```

```
##  88 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  89 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  90 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  91 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  92 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  93 P2_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  94 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  95 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  96 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  97 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  98 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
##  99 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 100 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 101 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 102 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 103 P2_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 104 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 105 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 106 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 107 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 108 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 109 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 110 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 111 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 112 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 113 P2_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 114 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 115 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 116 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 117 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 118 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 119 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 120 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 121 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 122 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 123 P2_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 124 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 125 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 126 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 127 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 128 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 129 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 130 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 131 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 132 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 133 P2_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 134 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 135 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 136 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 137 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 138 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 139 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 140 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 141 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 142 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
```

```
## 143 P2_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 144 P2_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 145 P2_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 146 P2_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 147 P2_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 148 P3_00~ " !!Total:"                        RACE FOR THE POPULATION 18 YEARS ~
## 149 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 150 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 151 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 152 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 153 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 154 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 155 P3_00~ " !!Total:!!Population of one rac~ RACE FOR THE POPULATION 18 YEARS ~
## 156 P3_00~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 157 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 158 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 159 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 160 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 161 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 162 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 163 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 164 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 165 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 166 P3_01~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 167 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 168 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 169 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 170 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 171 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 172 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 173 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 174 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 175 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 176 P3_02~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 177 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 178 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 179 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 180 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 181 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 182 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 183 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 184 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 185 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 186 P3_03~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 187 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 188 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 189 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 190 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 191 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 192 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 193 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 194 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 195 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 196 P3_04~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 197 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
```

```
## 198 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 199 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 200 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 201 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 202 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 203 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 204 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 205 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 206 P3_05~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 207 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 208 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 209 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 210 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 211 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 212 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 213 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 214 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 215 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 216 P3_06~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 217 P3_07~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 218 P3_07~ " !!Total:!!Population of two or ~ RACE FOR THE POPULATION 18 YEARS ~
## 219 P4_00~ " !!Total:"                      HISPANIC OR LATINO, AND NOT HISPA~
## 220 P4_00~ " !!Total:!!Hispanic or Latino"   HISPANIC OR LATINO, AND NOT HISPA~
## 221 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 222 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 223 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 224 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 225 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 226 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 227 P4_00~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 228 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 229 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 230 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 231 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 232 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 233 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 234 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 235 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 236 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 237 P4_01~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 238 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 239 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 240 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 241 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 242 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 243 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 244 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 245 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 246 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 247 P4_02~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 248 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 249 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 250 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 251 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 252 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
```

```
## 253 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 254 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 255 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 256 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 257 P4_03~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 258 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 259 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 260 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 261 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 262 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 263 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 264 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 265 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 266 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 267 P4_04~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 268 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 269 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 270 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 271 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 272 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 273 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 274 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 275 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 276 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 277 P4_05~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 278 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 279 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 280 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 281 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 282 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 283 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 284 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 285 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 286 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 287 P4_06~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 288 P4_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 289 P4_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 290 P4_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 291 P4_07~ " !!Total:!!Not Hispanic or Latin~ HISPANIC OR LATINO, AND NOT HISPA~
## 292 P5_00~ " !!Total:"                        GROUP QUARTERS POPULATION BY MAJO~
## 293 P5_00~ " !!Total:!!Institutionalized pop~ GROUP QUARTERS POPULATION BY MAJO~
## 294 P5_00~ " !!Total:!!Institutionalized pop~ GROUP QUARTERS POPULATION BY MAJO~
## 295 P5_00~ " !!Total:!!Institutionalized pop~ GROUP QUARTERS POPULATION BY MAJO~
## 296 P5_00~ " !!Total:!!Institutionalized pop~ GROUP QUARTERS POPULATION BY MAJO~
## 297 P5_00~ " !!Total:!!Institutionalized pop~ GROUP QUARTERS POPULATION BY MAJO~
## 298 P5_00~ " !!Total:!!Noninstitutionalized ~ GROUP QUARTERS POPULATION BY MAJO~
## 299 P5_00~ " !!Total:!!Noninstitutionalized ~ GROUP QUARTERS POPULATION BY MAJO~
## 300 P5_00~ " !!Total:!!Noninstitutionalized ~ GROUP QUARTERS POPULATION BY MAJO~
## 301 P5_01~ " !!Total:!!Noninstitutionalized ~ GROUP QUARTERS POPULATION BY MAJO~
```

### Look at Decennial Population Numbers

```r
pop20 <-  get_decennial(
  geography = "state",
  variables = "P1_001N",
  year = 2020)
```

```
## Getting data from the 2020 decennial Census

## Using the PL 94-171 Redistricting Data summary file

## Note: 2020 decennial Census data use differential privacy, a technique that
## introduces errors into data to preserve respondent confidentiality.
## i Small counts should be interpreted with caution.
## i See https://www.census.gov/library/fact-sheets/2021/protecting-the-confidentiality-o
## f-the-2020-census-redistricting-data.html for additional guidance.
## This message is displayed once per session.
```

### View table of decennial counts

```
print(tbl_df(pop20), n=52)

## # A tibble: 52 x 4
##    GEOID NAME                 variable    value
##    <chr> <chr>                <chr>       <dbl>
##  1 01    Alabama              P1_001N   5024279
##  2 02    Alaska               P1_001N    733391
##  3 04    Arizona              P1_001N   7151502
##  4 05    Arkansas             P1_001N   3011524
##  5 06    California           P1_001N  39538223
##  6 08    Colorado             P1_001N   5773714
##  7 09    Connecticut          P1_001N   3605944
##  8 10    Delaware             P1_001N    989948
##  9 11    District of Columbia P1_001N    689545
## 10 16    Idaho                P1_001N   1839106
## 11 12    Florida              P1_001N  21538187
## 12 13    Georgia              P1_001N  10711908
## 13 15    Hawaii               P1_001N   1455271
## 14 17    Illinois             P1_001N  12812508
## 15 18    Indiana              P1_001N   6785528
## 16 19    Iowa                 P1_001N   3190369
## 17 20    Kansas               P1_001N   2937880
## 18 21    Kentucky             P1_001N   4505836
## 19 22    Louisiana            P1_001N   4657757
## 20 23    Maine                P1_001N   1362359
## 21 24    Maryland             P1_001N   6177224
## 22 25    Massachusetts        P1_001N   7029917
## 23 26    Michigan             P1_001N  10077331
## 24 27    Minnesota            P1_001N   5706494
## 25 28    Mississippi          P1_001N   2961279
## 26 29    Missouri             P1_001N   6154913
## 27 30    Montana              P1_001N   1084225
## 28 31    Nebraska             P1_001N   1961504
## 29 32    Nevada               P1_001N   3104614
## 30 33    New Hampshire        P1_001N   1377529
## 31 34    New Jersey           P1_001N   9288994
## 32 35    New Mexico           P1_001N   2117522
## 33 36    New York             P1_001N  20201249
## 34 37    North Carolina       P1_001N  10439388
## 35 38    North Dakota         P1_001N    779094
## 36 39    Ohio                 P1_001N  11799448
## 37 40    Oklahoma             P1_001N   3959353
## 38 41    Oregon               P1_001N   4237256
## 39 42    Pennsylvania         P1_001N  13002700
```

```
## 40 44    Rhode Island       P1_001N    1097379
## 41 45    South Carolina     P1_001N    5118425
## 42 46    South Dakota       P1_001N     886667
## 43 47    Tennessee          P1_001N    6910840
## 44 48    Texas              P1_001N   29145505
## 45 49    Utah               P1_001N    3271616
## 46 50    Vermont            P1_001N     643077
## 47 51    Virginia           P1_001N    8631393
## 48 53    Washington         P1_001N    7705281
## 49 54    West Virginia      P1_001N    1793716
## 50 55    Wisconsin          P1_001N    5893718
## 51 56    Wyoming            P1_001N     576851
## 52 72    Puerto Rico        P1_001N    3285874
```

## View DMV population from Census provided data

- District of Columbia

```
pop20 %>% filter(GEOID == 11)
```

```
## # A tibble: 1 x 4
##    GEOID NAME                 variable  value
##    <chr> <chr>                <chr>     <dbl>
## 1 11    District of Columbia P1_001N  689545
```

- Maryland

```
pop20 %>% filter(GEOID == 24)
```

```
## # A tibble: 1 x 4
##    GEOID NAME     variable   value
##    <chr> <chr>    <chr>      <dbl>
## 1 24    Maryland P1_001N  6177224
```

- Virginia

```
pop20 %>% filter(GEOID == 51)
```

```
## # A tibble: 1 x 4
##    GEOID NAME     variable   value
##    <chr> <chr>    <chr>      <dbl>
## 1 51    Virginia P1_001N  8631393
```

## View DMV population from outside source provided data

- District of Columbia

```
Census2020 %>% filter(Area == "District of Columbia")
```

```
## # A tibble: 1 x 10
##    Area        Region `2020 Census Residen~ `2010 Census Residen~ `Numeric Change`
##    <chr>       <chr>                  <dbl>                 <dbl>            <dbl>
## 1 District ~ South                 689545                601723            87822
## # ... with 5 more variables: Percent Change <dbl>,
## #   State Rank Based on 2020 Census Resident Population <chr>,
## #   State Rank Based on 2010 Census Resident Population <chr>,
## #   State Rank Based on Numeric Change <chr>,
## #   State Rank Based on Percent Change <chr>
```

- Maryland

```
Census2020 %>% filter(Area == "Maryland")
```

```
## # A tibble: 1 x 10
##   Area     Region `2020 Census Resident ~ `2010 Census Resident~ `Numeric Change`
##   <chr>    <chr>                   <dbl>                  <dbl>            <dbl>
## 1 Maryla~ South                 6177224                5773552           403672
## # ... with 5 more variables: Percent Change <dbl>,
## #   State Rank Based on 2020 Census Resident Population <chr>,
## #   State Rank Based on 2010 Census Resident Population <chr>,
## #   State Rank Based on Numeric Change <chr>,
## #   State Rank Based on Percent Change <chr>
```

- Virginia

```
Census2020 %>% filter(Area == "Virginia")
```

```
## # A tibble: 1 x 10
##   Area     Region `2020 Census Resident ~ `2010 Census Resident~ `Numeric Change`
##   <chr>    <chr>                   <dbl>                  <dbl>            <dbl>
## 1 Virgin~ South                 8631393                8001024           630369
## # ... with 5 more variables: Percent Change <dbl>,
## #   State Rank Based on 2020 Census Resident Population <chr>,
## #   State Rank Based on 2010 Census Resident Population <chr>,
## #   State Rank Based on Numeric Change <chr>,
## #   State Rank Based on Percent Change <chr>
```

## Compare the two sources of data, create new objects for each

- District of Columbia

```
API_DC <-  pop20 %>%
  filter(GEOID == 11) %>%
  select(value)

ACS_DC <-  Census2020 %>%
  filter(Area == "District of Columbia") %>%
  select(`2020 Census Resident Population`)
```

- Maryland

```
API_MD <- pop20 %>% filter(GEOID == 24) %>%
  select(value)

ACS_MD <-  Census2020 %>%
  filter(Area == "Maryland") %>%
  select(`2020 Census Resident Population`)
```

- Virginia

```
API_VA <-  pop20 %>% filter(GEOID == 51) %>%
  select(value)

ACS_VA <-  Census2020 %>%
  filter(Area == "Virginia") %>%
  select(`2020 Census Resident Population`)
```

## Do the two sources of population data match?

- District of Columbia

```
all(API_DC == ACS_DC)
```

```
## [1] TRUE
```

- Maryland

```
all(API_MD == ACS_MD)

## [1] TRUE
```

- Virginia

```
all(API_VA == ACS_VA)

## [1] TRUE
```

### Group quarters data

```
group_quarters <-  get_decennial(
  geography = "state",
  table = "P5",
  year = 2020,
  output = "wide")

## Getting data from the 2020 decennial Census

## Loading PL variables for 2020 from table P5. To cache this dataset for faster access t
o Census tables in the future, run this function with `cache_table = TRUE`. You only need
to do this once per Census dataset.

## Using the PL 94-171 Redistricting Data summary file
```

### Show top observations of group quarters data

```
head(group_quarters)

## # A tibble: 6 x 12
##    GEOID NAME      P5_001N P5_002N P5_003N P5_004N P5_005N P5_006N P5_007N P5_008N
##    <chr> <chr>       <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 01    Alabama    127934   70648   39749    1479   27869    1551   57286   45489
## 2 02    Alaska      30291    7177    4842     457    1781      97   23114    1472
## 3 04    Arizona    160269   89904   64154    2331   21938    1481   70365   38945
## 4 05    Arkansas    82518   48001   27079    1248   19266     408   34517   26887
## 5 06    Califor~   917932  344896  201570    8966  124804    9556  573036  230361
## 6 08    Colorado   126848   55851   32307    1525   21379     640   70997   38819
## # ... with 2 more variables: P5_009N <dbl>, P5_010N <dbl>
```

### Group quarters DMV data

- District of Columbia

```
dc_group_quarters <-  get_decennial(
  geography = "state",
  table = "P5",
  state = "DC",
  year = 2020,
  output = "wide")

## Getting data from the 2020 decennial Census

## Loading PL variables for 2020 from table P5. To cache this dataset for faster access t
o Census tables in the future, run this function with `cache_table = TRUE`. You only need
to do this once per Census dataset.

## Using the PL 94-171 Redistricting Data summary file
```

- Maryland

```r
md_group_quarters <- get_decennial(
  geography = "state",
  table = "P5",
  state = "MD",
  year = 2020,
  output = "wide")
```

```
## Getting data from the 2020 decennial Census

## Loading PL variables for 2020 from table P5. To cache this dataset for faster access t
o Census tables in the future, run this function with `cache_table = TRUE`. You only need
to do this once per Census dataset.

## Using the PL 94-171 Redistricting Data summary file
```

- Virginia

```r
va_group_quarters <- get_decennial(
  geography = "state",
  table = "P5",
  state = "VA",
  year = 2020,
  output = "wide")
```

```
## Getting data from the 2020 decennial Census

## Loading PL variables for 2020 from table P5. To cache this dataset for faster access t
o Census tables in the future, run this function with `cache_table = TRUE`. You only need
to do this once per Census dataset.

## Using the PL 94-171 Redistricting Data summary file
```

## Use rbind to concatenate rows

```r
dmv_group_quarters <- rbind(va_group_quarters,
                            md_group_quarters,
                            dc_group_quarters)
```

## View DMV group quarters object

```r
dmv_group_quarters
```

```
## # A tibble: 3 x 12
##    GEOID NAME      P5_001N P5_002N P5_003N P5_004N P5_005N P5_006N P5_007N P5_008N
##    <chr> <chr>       <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 51    Virginia   236646   96832   57014    2038   36195    1585  139814   92450
## 2 24    Maryland   125505   58693   27040    1008   29252    1393   66812   46179
## 3 11    Distric~    40682    5606    2278     315    2727     286   35076   23802
## # ... with 2 more variables: P5_009N <dbl>, P5_010N <dbl>
```

## Show hispanic DMV data

- District of Columbia

```r
dc_hispanic <- get_decennial(
  geography = "county",
  variables = "P2_002N",
  state = "DC",
  year = 2020)
```

```
## Getting data from the 2020 decennial Census

## Using the PL 94-171 Redistricting Data summary file
```

- Maryland

```
md_hispanic <-  get_decennial(
  geography = "county",
  variables = "P2_002N",
  state = "MD",
  year = 2020)

## Getting data from the 2020 decennial Census

## Using the PL 94-171 Redistricting Data summary file
```

- Virginia

```
va_hispanic <-  get_decennial(
  geography = "county",
  variables = "P2_002N",
  state = "VA",
  year = 2020)

## Getting data from the 2020 decennial Census

## Using the PL 94-171 Redistricting Data summary file
```

## Show DMV Hispanic data

- District of Columbia

```
dc_hispanic

## # A tibble: 1 x 4
##   GEOID NAME                                        variable value
##   <chr> <chr>                                       <chr>    <dbl>
## 1 11001 District of Columbia, District of Columbia P2_002N   77652
```

- Maryland
```
md_hispanic

## # A tibble: 24 x 4
##    GEOID NAME                                  variable  value
##    <chr> <chr>                                 <chr>     <dbl>
##  1 24003 Anne Arundel County, Maryland         P2_002N   56796
##  2 24005 Baltimore County, Maryland            P2_002N   61492
##  3 24011 Caroline County, Maryland             P2_002N    2820
##  4 24013 Carroll County, Maryland              P2_002N    7745
##  5 24017 Charles County, Maryland              P2_002N   11677
##  6 24019 Dorchester County, Maryland           P2_002N    1777
##  7 24023 Garrett County, Maryland              P2_002N     321
##  8 24025 Harford County, Maryland              P2_002N   14007
##  9 24029 Kent County, Maryland                 P2_002N    1061
## 10 24033 Prince George's County, Maryland P2_002N   205463
## # ... with 14 more rows
```

- Virginia

```
va_hispanic

## # A tibble: 133 x 4
##    GEOID NAME                       variable value
##    <chr> <chr>                      <chr>    <dbl>
##  1 51003 Albemarle County, Virginia P2_002N   8453
##  2 51005 Alleghany County, Virginia P2_002N    178
##  3 51009 Amherst County, Virginia   P2_002N    838
##  4 51011 Appomattox County, Virginia P2_002N   344
##  5 51015 Augusta County, Virginia   P2_002N   2728
##  6 51017 Bath County, Virginia      P2_002N     73
##  7 51021 Bland County, Virginia     P2_002N     60
##  8 51023 Botetourt County, Virginia P2_002N    776
##  9 51027 Buchanan County, Virginia  P2_002N    177
## 10 51029 Buckingham County, Virginia P2_002N   413
## # ... with 123 more row
```