

Problem set 4

2025-10-05

In the next problem set, we plan to explore the relationship between COVID-19 death rates and vaccination rates across US states by visually examining their correlation. This analysis will involve gathering COVID-19 related data from the CDC's API and then extensively processing it to merge the various datasets. Since the population sizes of states vary significantly, we will focus on comparing rates rather than absolute numbers. To facilitate this, we will also source population data from the US Census to accurately calculate these rates.

In this problem set we will learn how to extract and wrangle data from the data US Census and CDC APIs.

1. Get an API key from the US Census at https://api.census.gov/data/key_signup.html. You can't share this public key. But your code has to run on a TFs computer. Assume the TF will have a file in their working directory named `census-key.R` with the following one line of code:

```
census_key <- ""
```

Write a first line of code for your problem set that defines `census_key` by running the code in the file `census-key.R`.

```
source("census-key.R")
```

2. The [US Census API User Guide](#) provides details on how to leverage this valuable resource. We are interested in vintage population estimates for years 2021 and 2022. From the documentation we find that the *endpoint* is:

```
url <- "https://api.census.gov/data/2021/pep/population"
```

Use the `httr2` package to construct the following GET request.

`https://api.census.gov/data/2021/pep/population?get=POP_2020,POP_2021,NAME&for=state:*&key=Y`

Create an object called `request` of class `httr2_request` with this URL as an endpoint. Hint: Print out `request` to check that the URL matches what we want.

```
library(httr2)
request <- request(url) |>
  req_url_query(
    get = "POP_2020,POP_2021,NAME",
    `for` = "state:",
    key = census_key
  )
```

3. Make a request to the US Census API using the `request` object. Save the response to an object named `response`. Check the response status of your request and make sure it was successful. You can learn about *status codes* [here](#).

```
response <- request |> req_perform()
resp_status(response)
```

```
[1] 200
```

```
resp_status_desc(response)
```

```
[1] "OK"
```

4. Use a function from the `httr2` package to determine the content type of your response.

```
resp_content_type(response)
```

```
[1] "application/json"
```

5. Use just one line of code and one function to extract the data into a matrix. Hints: 1) Use the `resp_body_json` function. 2) The first row of the matrix will be the variable names and this OK as we will fix in the next exercise.

```
population <- resp_body_json(response) |> do.call(rbind, args = _)
```

6. Examine the `population` matrix you just created. Notice that 1) it is not tidy, 2) the column types are not what we want, and 3) the first row is a header. Convert `population` to a tidy dataset. Remove the state ID column and change the name of the column with state names to `state_name`. Add a column with state abbreviations called `state`. Make sure you assign the abbreviations for DC and PR correctly. Hint: Use the **janitor** package to make the first row the header.

```
library(tidyverse)
library(janitor)
#population <- population |> ## Use janitor row to names function
# convert to tibble
# remove stat column
# rename state column to state_name
# use pivot_longer to tidy
# remove POP_ from year
# parse all relevant columns to numeric
# add state abbreviations using state.abb variable
# use case_when to add abbreviations for DC and PR
library(purrr)

population <- population |>
  as_tibble(.name_repair = "minimal") |>
  row_to_names(1) |>
  rename(state_name = NAME, state_fips = state) |>
  pivot_longer(starts_with("POP_"),
               names_to = "year", values_to = "population") |>
  mutate(
    year = sub("^POP_", "", year),
    population = as.numeric(population)
  ) |>
  select(-state_fips) |>
  mutate(
    state_name = map_chr(state_name, 1),
    state = case_when(
      state_name == "District of Columbia" ~ "DC",
      state_name == "Puerto Rico" ~ "PR",
      TRUE ~ state.abb[match(state_name, state.name)]
    )
  )
```

```
head(population)
```

```
# A tibble: 6 x 4
  state_name year  population state
  <chr>      <chr>      <dbl> <chr>
1 Oklahoma   2020      3962031 OK
2 Oklahoma   2021      3986639 OK
3 Nebraska   2020      1961455 NE
4 Nebraska   2021      1963692 NE
5 Hawaii     2020      1451911 HI
6 Hawaii     2021      1441553 HI
```

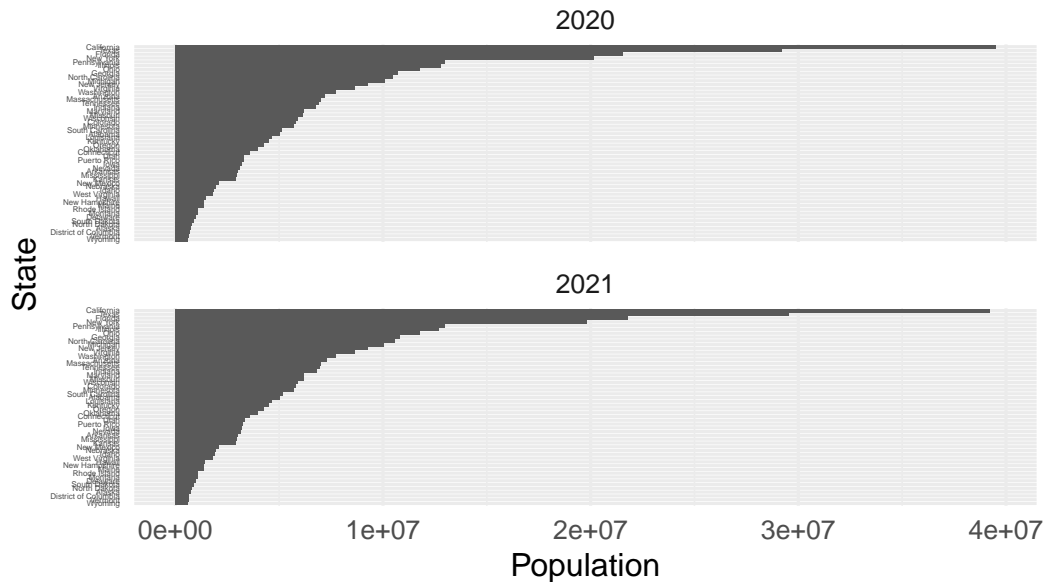
```
nrow(population)
```

```
[1] 104
```

7. As a check, make a barplot of states' 2021 and 2022 populations. Show the state names in the y-axis ordered by population size. Hint: You will need to use **reorder** and use **facet_wrap**.

```
# population |>
# reorder state
# assign aesthetic mapping
# use geom_col to plot barplot
# flip coordinates
# facet by year
population |>
  filter(year %in% c("2020", "2021")) |>
  mutate(state_order = reorder(state_name, population)) |>
  ggplot(aes(x = state_order, y = population)) +
  geom_col() +
  coord_flip() +
  facet_wrap(~ year, ncol = 1, scales = "free_y") +
  labs(x = "State", y = "Population", title = "US State Populations (2020 vs
    ↪ 2021)") +
  theme_minimal(base_size = 12)+
  theme(axis.text.y = element_text(size = 3))
```

US State Populations (2020 vs 2021)



8. The following URL:

```
url <- "https://github.com/datasciencelabs/2025/raw/refs/heads/main/data/regions.json"
```

points to a JSON file that lists the states in the 10 Public Health Service (PHS) defined by CDC. We want to add these regions to the `population` dataset. To facilitate this create a data frame called `regions` that has two columns `state_name`, `region`, `region_name`. One of the regions has a long name. Change it to something shorter.

```
jsonlite::fromJSON("https://github.com/datasciencelabs/2025/raw/refs/heads/main/data/regions.json") |> str()
```

```
'data.frame':  10 obs. of  3 variables:
 $ region      :List of 10
 ..$ : int 1
 ..$ : int 2
 ..$ : int 3
 ..$ : int 4
 ..$ : int 5
 ..$ : int 6
 ..$ : int 7
```

```

..$ : int 8
..$ : int 9
..$ : int 10
$ region_name: chr "New England" "New York and New Jersey, Puerto Rico, Virgin Islands" "M
$ states      :List of 10
..$ : chr "Connecticut" "Maine" "Massachusetts" "New Hampshire" ...
..$ : chr "New Jersey" "New York" "Puerto Rico" "Virgin Islands"
..$ : chr "Delaware" "District of Columbia" "Maryland" "Pennsylvania" ...
..$ : chr "Alabama" "Florida" "Georgia" "Kentucky" ...
..$ : chr "Illinois" "Indiana" "Michigan" "Minnesota" ...
..$ : chr "Arkansas" "Louisiana" "New Mexico" "Oklahoma" ...
..$ : chr "Iowa" "Kansas" "Missouri" "Nebraska"
..$ : chr "Colorado" "Montana" "North Dakota" "South Dakota" ...
..$ : chr "Arizona" "California" "Hawaii" "Nevada" ...
..$ : chr "Alaska" "Idaho" "Oregon" "Washington"

```

```

library(jsonlite)
library(purrr)
url <- "https://github.com/datasciencelabs/2025/raw/refs/heads/main/data/regi
↪ ions.json"
# regions <- use jsonlit JSON parser
# regions <- convert list to data frame. You can use map_df in purrr package

regions <- fromJSON(url) |>
  mutate(region = vapply(region, function(x) as.integer(x[1]), integer(1)))
↪ |>
  unnest_longer(states) |>
  rename(state_name = states) |>
  mutate(
    region_name = if_else(
      region_name == "New York and New Jersey, Puerto Rico, Virgin Islands",
      "NY+NJ+PR",
      region_name
    )
  ) |>
  filter(state_name %in% c(state.name, "District of Columbia", "Puerto
↪ Rico")) |>
  select(state_name, region, region_name)

```

```
head(regions)
```

```
# A tibble: 6 x 3
```

| | state_name | region | region_name |
|---|---------------|--------|-------------|
| | <chr> | <int> | <chr> |
| 1 | Connecticut | 1 | New England |
| 2 | Maine | 1 | New England |
| 3 | Massachusetts | 1 | New England |
| 4 | New Hampshire | 1 | New England |
| 5 | Rhode Island | 1 | New England |
| 6 | Vermont | 1 | New England |

9. Add a region and region name columns to the `population` data frame.

```
population <- population |>
  left_join(regions, by = "state_name")
```

```
#population |> filter(is.na(region)) |> distinct(state_name)
```

10. From reading <https://data.cdc.gov/> we learn the endpoint <https://data.cdc.gov/resource/pwn4-m3yp> provides state level data from SARS-COV2 cases. Use the **httr2** tools you have learned to download this into a data frame. Is all the data there? If not, comment on why.

```
api <- "https://data.cdc.gov/resource/pwn4-m3yp.json"
# cases_raw <-
cases_raw <- request(api) |>
  req_perform() |>
  resp_body_json(simplifyDataFrame = TRUE)
```

We see exactly 1,000 rows. We should be seeing over 52×3 rows per state.

11. The reason you see exactly 1,000 rows is because CDC has a default limit. You can change this limit by adding `$limit=10000000000` to the request. Rewrite the previous request to ensure that you receive all the data. Then wrangle the resulting data frame to produce a data frame with columns `state`, `date` (should be the end date) and `cases`. Make sure the cases are numeric and the dates are in Date ISO-8601 format.

```
api <- "https://data.cdc.gov/resource/pwn4-m3yp.json"
cases_raw <- request(api) |>
  req_url_query(`$limit` = 10000000000) |>
  req_perform() |>
  resp_body_json(simplifyDataFrame = TRUE)

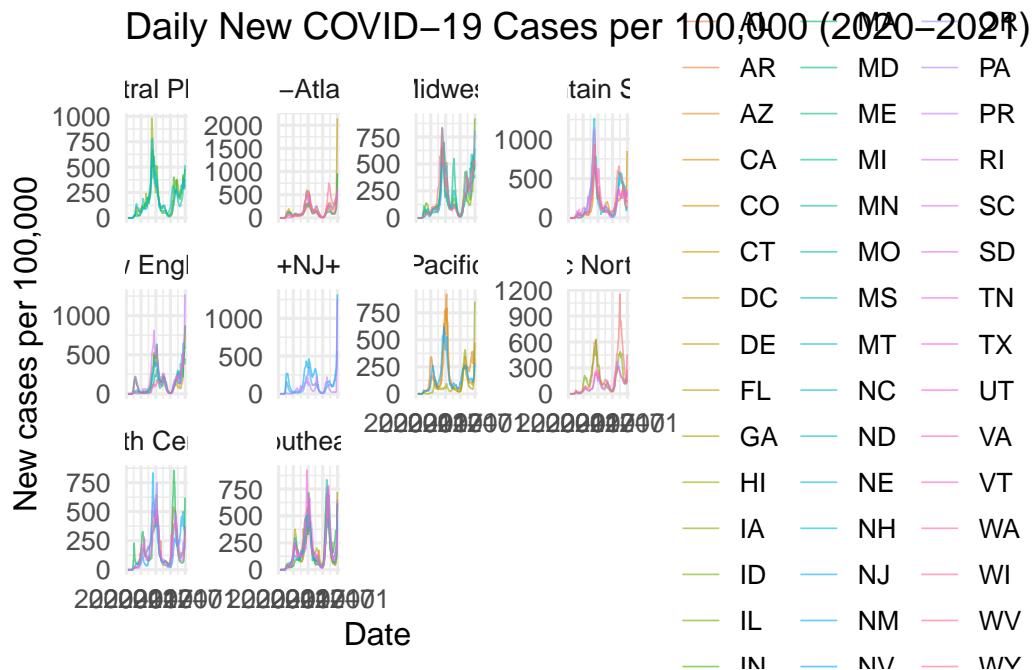
cases <- cases_raw |>
```

```
transmute(
  state = state,
  date = as.Date(end_date),
  cases = as.numeric(new_cases)
)|>
filter(state %in% c(state.abb, "DC", "PR"), !is.na(cases), !is.na(date))
```

```
#names(cases_raw)
```

12. For 2020 and 2021, make a time series plot of cases per 100,000 versus time for each state. Stratify the plot by region name. Make sure to label you graph appropriately.

```
cases |>
  filter(format(date, "%Y") %in% c("2020", "2021")) |>
  mutate(year = as.integer(format(date, "%Y"))) |>
  left_join(
    population |>
      filter(year %in% c("2020", "2021")) |>
      mutate(year = as.integer(year)) |>
      select(state, region_name, year, population),
    by = c("state", "year")
  ) |>
  mutate(cases_per_100k = cases / population * 1e5) |>
  ggplot(aes(x = date, y = cases_per_100k, group = state, col = state)) +
  geom_line(alpha = 0.6, linewidth = 0.3) +
  facet_wrap(~ region_name, scales = "free_y") +
  labs(
    title = "Daily New COVID-19 Cases per 100,000 (2020-2021)",
    x = "Date", y = "New cases per 100,000"
  ) +
  theme_minimal(base_size = 12)
```

13. The dates in the `cases` dataset are stored as character strings. Use the **lubridate** package to properly parse the `date` column, then create a summary table showing the total COVID-19 cases by month and year for 2020 and 2021. The table should have columns for year, month (as month name), and total cases across all states. Order by year and month. Use the **knitr** package and `kable()` function to display the results as a formatted table.

```
library(knitr)
#cases |>
cases |>
  mutate(
    date = date,
    year = year(date),
    month_num = month(date),
    month = month(date, label = TRUE, abbr = FALSE)
  ) |>
  filter(year %in% c(2020, 2021)) |>
  group_by(year, month_num, month) |>
  summarise(total_cases = sum(cases, na.rm = TRUE), .groups = "drop") |>
  arrange(year, month_num) |>
  select(year, month, total_cases) |>
  kable(
    caption = "Total New COVID-19 Cases by Month (All States), 2020-2021",
```

```
align = "lcr"
)
```

Table 1: Total New COVID-19 Cases by Month (All States), 2020–2021

| year | month | total_cases |
|------|-----------|-------------|
| 2020 | January | 11 |
| 2020 | February | 68 |
| 2020 | March | 50335 |
| 2020 | April | 822648 |
| 2020 | May | 616691 |
| 2020 | June | 642552 |
| 2020 | July | 1977016 |
| 2020 | August | 1452393 |
| 2020 | September | 1401917 |
| 2020 | October | 1608932 |
| 2020 | November | 3887222 |
| 2020 | December | 6907540 |
| 2021 | January | 5649115 |
| 2021 | February | 2543964 |
| 2021 | March | 1928749 |
| 2021 | April | 1694189 |
| 2021 | May | 948953 |
| 2021 | June | 484817 |
| 2021 | July | 1120939 |
| 2021 | August | 3519407 |
| 2021 | September | 4960807 |
| 2021 | October | 2317854 |
| 2021 | November | 2289118 |
| 2021 | December | 5293391 |

14. The following URL provides additional COVID-19 data from the CDC in JSON format:

```
deaths_url <- "https://data.cdc.gov/resource/9bhg-hcku.json"
```

Use **httr2** to download COVID-19 death data from this endpoint. Make sure to remove the default limit to get all available data. Create a clean dataset called **deaths** with columns **state**, **date**, and **deaths** (renamed from the original column name). Ensure dates are in proper Date format and deaths are numeric.

```

deaths_raw <- request(deaths_url) |>
  req_url_query(`$limit` = 1000000000) |>
  req_perform() |>
  resp_body_json(simplifyDataFrame = TRUE)

deaths <- deaths_raw |>
  transmute(
    state,
    date = as.Date(substr(end_date, 1, 10)),
    deaths = as.numeric(covid_19_deaths)
  ) |>
  filter(state %in% population$state_name, !is.na(deaths), !is.na(date))

```

```
names(deaths_raw)
```

```

[1] "data_as_of"           "start_date"
[3] "end_date"             "group"
[5] "state"                "sex"
[7] "age_group"            "covid_19_deaths"
[9] "total_deaths"         "pneumonia_deaths"
[11] "pneumonia_and_covid_19_deaths" "influenza_deaths"
[13] "pneumonia_influenza_or_covid" "footnote"
[15] "year"                 "month"

```

```
str(deaths)
```

```

'data.frame':  93933 obs. of  3 variables:
 $ state : chr  "Alabama" "Alabama" "Alabama" "Alabama" ...
 $ date  : Date, format: "2023-09-23" "2023-09-23" ...
 $ deaths: num  21520 19 46 142 267 ...

```

15. Using the `deaths` dataset you created, make a bar plot showing the total COVID-19 deaths by state. Show only the top 10 states with the highest death counts. Order the bars from highest to lowest and use appropriate labels and title.

```

deaths |>
  filter(state != "United States", state != "New York City") |>
  group_by(state) |>
  summarise(total_deaths = sum(deaths, na.rm = TRUE), .groups = "drop") |>
  slice_max(total_deaths, n = 10) |>

```

```

ggplot(aes(x = reorder(state, total_deaths), y = total_deaths)) +
  geom_col(fill = "#A51C30") +
  coord_flip() +
  labs(
    title = "Top 10 States by Total COVID-19 Deaths",
    x = "State",
    y = "Total Deaths"
  ) +
  theme_minimal(base_size = 12)

```

