

Pset 1 - Water usage

425/625

Spring 2024

Introduction

Water scarcity is a major issue in many parts of the world. According to the United Nations, “About two billion people worldwide don’t have access to safe drinking water today (SDG Report 2022), and roughly half of the world’s population is experiencing severe water scarcity for at least part of the year (IPCC). These numbers are expected to increase, exacerbated by climate change and population growth (WMO).”

In this problem set, we will investigate water usage estimates by crop in the United States. The `.csv` for this data set comes from here (by checking Select All and clicking Get Custom Zip) and the associated academic journal article is here. See this thread on X for a summary.

Read the academic article to familiarize yourself with the basics of the water usage data. You don’t need to know how these water usage levels were estimated, so you can skip over those parts. We are going to focus on visualizing the water levels using the estimates that they generated.

Data preparation

The `.zip` file `rawdata/DOI-10-13012-b2idb-4607538_v1.zip` contains one `.csv` file per source (SWW, GWW, GWD) per year from 2008 to 2020. There are also a couple of `.txt` files in the folder. We can use `unzip` with `list = TRUE` to see what’s in the `.zip` file.

```
unzip(zipfile = 'rawdata/DOI-10-13012-b2idb-4607538_v1.zip',  
      list = TRUE) ## this lists the filename, but does not unzip the file
```

##	Name	Length	Date
## 1	DOI-10-13012-b2idb-4607538_v1/readme.txt	1053	2023-10-29 14:08:00
## 2	DOI-10-13012-b2idb-4607538_v1/gwa_2008.csv	2274812	2023-10-29 14:08:00
## 3	DOI-10-13012-b2idb-4607538_v1/gwa_2009.csv	2274812	2023-10-29 14:08:00
## 4	DOI-10-13012-b2idb-4607538_v1/gwa_2010.csv	2200859	2023-10-29 14:08:00
## 5	DOI-10-13012-b2idb-4607538_v1/gwa_2011.csv	2274812	2023-10-29 14:08:00
## 6	DOI-10-13012-b2idb-4607538_v1/gwa_2012.csv	2274812	2023-10-29 14:08:00
## 7	DOI-10-13012-b2idb-4607538_v1/gwa_2013.csv	2274812	2023-10-29 14:08:00
## 8	DOI-10-13012-b2idb-4607538_v1/gwa_2014.csv	2274812	2023-10-29 14:08:00
## 9	DOI-10-13012-b2idb-4607538_v1/gwa_2015.csv	2200859	2023-10-29 14:08:00
## 10	DOI-10-13012-b2idb-4607538_v1/gwa_2016.csv	2275517	2023-10-29 14:08:00
## 11	DOI-10-13012-b2idb-4607538_v1/gwa_2017.csv	2275517	2023-10-29 14:08:00
## 12	DOI-10-13012-b2idb-4607538_v1/gwa_2018.csv	2275517	2023-10-29 14:08:00
## 13	DOI-10-13012-b2idb-4607538_v1/gwa_2019.csv	2275517	2023-10-29 14:08:00
## 14	DOI-10-13012-b2idb-4607538_v1/gwa_2020.csv	2275517	2023-10-29 14:08:00
## 15	DOI-10-13012-b2idb-4607538_v1/gwd_2008.csv	211884	2023-10-29 14:08:00
## 16	DOI-10-13012-b2idb-4607538_v1/gwd_2009.csv	208249	2023-10-29 14:08:00

```
## 17 DOI-10-13012-b2idb-4607538_v1/gwd_2010.csv 214546 2023-10-29 14:08:00
## 18 DOI-10-13012-b2idb-4607538_v1/gwd_2011.csv 213608 2023-10-29 14:08:00
## 19 DOI-10-13012-b2idb-4607538_v1/gwd_2012.csv 210157 2023-10-29 14:08:00
## 20 DOI-10-13012-b2idb-4607538_v1/gwd_2013.csv 207564 2023-10-29 14:08:00
## 21 DOI-10-13012-b2idb-4607538_v1/gwd_2014.csv 209619 2023-10-29 14:08:00
## 22 DOI-10-13012-b2idb-4607538_v1/gwd_2015.csv 208683 2023-10-29 14:08:00
## 23 DOI-10-13012-b2idb-4607538_v1/gwd_2016.csv 206644 2023-10-29 14:08:00
## 24 DOI-10-13012-b2idb-4607538_v1/gwd_2017.csv 206188 2023-10-29 14:08:00
## 25 DOI-10-13012-b2idb-4607538_v1/gwd_2018.csv 206429 2023-10-29 14:08:00
## 26 DOI-10-13012-b2idb-4607538_v1/gwd_2019.csv 208246 2023-10-29 14:08:00
## 27 DOI-10-13012-b2idb-4607538_v1/gwd_2020.csv 208252 2023-10-29 14:08:00
## 28 DOI-10-13012-b2idb-4607538_v1/sw_2008.csv 2274792 2023-10-29 14:08:00
## 29 DOI-10-13012-b2idb-4607538_v1/sw_2009.csv 2274792 2023-10-29 14:08:00
## 30 DOI-10-13012-b2idb-4607538_v1/sw_2010.csv 2200839 2023-10-29 14:08:00
## 31 DOI-10-13012-b2idb-4607538_v1/sw_2011.csv 2274792 2023-10-29 14:08:00
## 32 DOI-10-13012-b2idb-4607538_v1/sw_2012.csv 2274792 2023-10-29 14:08:00
## 33 DOI-10-13012-b2idb-4607538_v1/sw_2013.csv 2274792 2023-10-29 14:08:00
## 34 DOI-10-13012-b2idb-4607538_v1/sw_2014.csv 2274792 2023-10-29 14:08:00
## 35 DOI-10-13012-b2idb-4607538_v1/sw_2015.csv 2200839 2023-10-29 14:08:00
## 36 DOI-10-13012-b2idb-4607538_v1/sw_2016.csv 2275497 2023-10-29 14:08:00
## 37 DOI-10-13012-b2idb-4607538_v1/sw_2017.csv 2275497 2023-10-29 14:08:00
## 38 DOI-10-13012-b2idb-4607538_v1/sw_2018.csv 2275497 2023-10-29 14:08:00
## 39 DOI-10-13012-b2idb-4607538_v1/sw_2019.csv 2275497 2023-10-29 14:08:00
## 40 DOI-10-13012-b2idb-4607538_v1/sw_2020.csv 2275497 2023-10-29 14:08:00
## 41 DOI-10-13012-b2idb-4607538_v1/dataset_info.txt 3894 2023-10-29 14:08:00
```

Before summarizing/visualizing this data, we'll want to join these data sets. We could certainly unzip the file manually. We can also do this in R using `unzip`.

```
unzip(zipfile = 'rawdata/DOI-10-13012-b2idb-4607538_v1.zip',
      junkpaths = TRUE,
      exdir = 'rawdata') ## gets rid of paths, keeps only filenames
```

1. Join data First, let's create a data set with all years/crops together in one data frame. Below is some code to help you get started. Add comments to each place there is `##` to explain what the chunk of code is doing. Then add code to the **Transforming data** Section to transform the data into a data frame with 5 columns: `GEOID`, `crop`, `source`, `year`, and `value` (indicating km^3 of water).

Note that `eval = F` at the start of the chunk will prevent this chunk from evaluating when you knit the document. You can temporarily remove it if you'd like, but you'll want to add it back before knitting the document so that knitting takes less time.

```
sources = c('gwd', 'sw', 'gwa')
years = 2008:2020
d = NULL

for(s in sources){
  cat(s, ' ') ## show progress

  for(year in years){
    cat(year, ' ') ## show progress

    ## load raw data file named s_year.csv
```

```

filename = paste0('rawdata/', s, '_', year, '.csv')
df = read.csv(filename)
head(df)

## Tranform data #####
## Use `pivot_longer`, `separate`, and/or other functions to transform this
## data frame into a data frame with 5 columns:
## GEOID, crop, source, year, and value (indicating km^3 of water)
df <- df %>% pivot_longer(cols = !GEOID,
                        names_to = c("src", "crop", "year"),
                        names_sep = "[.]",
                        values_to = "value")

## end of transforming data #####

## concatenate transformed data
d = rbind(d, df)
}

cat('\n') ## start a new line before showing progress for the next source
}
d <- d[, c("GEOID", "crop", "src", "year", "value")]
head(d)
tail(d)

```

Data exploration and summaries

Let's load the data we'll use for the rest of the assignment. This is the data set created in #1, so if you were unable to finish #1, you can still do the rest of the assignment.

```

d = readRDS('data/water.usage.rds')
head(d)

```

```

## # A tibble: 6 x 5
##   GEOID crop      src  year  value
##   <int> <chr>    <chr> <chr> <dbl>
## 1  1001 barley   gwd   2008     0
## 2  1001 corn    gwd   2008     0
## 3  1001 cotton  gwd   2008     0
## 4  1001 millet  gwd   2008     0
## 5  1001 oats    gwd   2008     0
## 6  1001 other_sctg2 gwd   2008     0

```

2. Summaries of data Find the mean, the change from 2008 to 2020, and the percent change from 2008 to 2020, for each crop and each source (SWW, GWW, GWD).

```

## calculate mean for each crop and source
d_mean <- d %>%
  group_by(crop, src) %>%
  summarise(mean = mean(value))

```

```

## 'summarise()' has grouped output by 'crop'. You can override using the
## '.groups' argument.

```

```
# head(d_mean)

## calculate change and % change for each crop and source
d_chg <- d %>%
  filter(year == "2008" | year == "2020") %>%
  group_by(crop, src, year) %>%
  summarise(total = sum(value)) %>%
  mutate(change = c(0, diff(total))) %>%
  mutate(pct_change = change/lag(total) * 100) %>%
  filter(year == "2020") %>%
  select(crop, src, change, pct_change)
```

'summarise()' has grouped output by 'crop', 'src'. You can override using the
'.groups' argument.

```
# head(d_chg)

d_summ <- full_join(d_mean, d_chg, by = c("crop" = "crop", "src" = "src"))
head(d_summ)
```

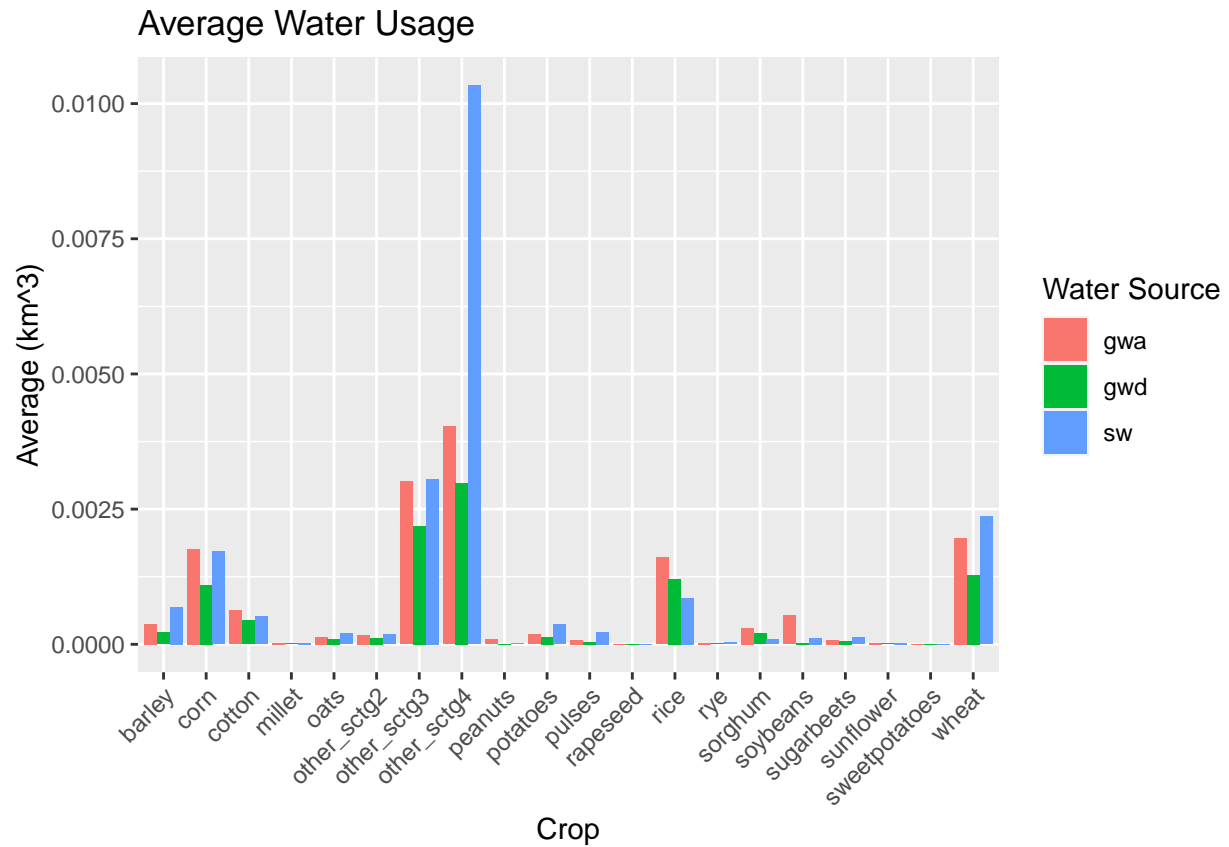
```
## # A tibble: 6 x 5
## # Groups:   crop [2]
##   crop   src      mean  change pct_change
##   <chr> <chr>    <dbl>  <dbl>    <dbl>
## 1 barley gwa  0.000372  0.0631     5.21
## 2 barley gwd  0.000222 -0.118    -17.4
## 3 barley sw   0.000688 -0.508    -21.4
## 4 corn   gwa  0.00176   0.617     11.5
## 5 corn   gwd  0.00110  -0.167     -4.61
## 6 corn   sw   0.00172  -2.19    -30.7
```

3. Convert Table 2 to a visualization

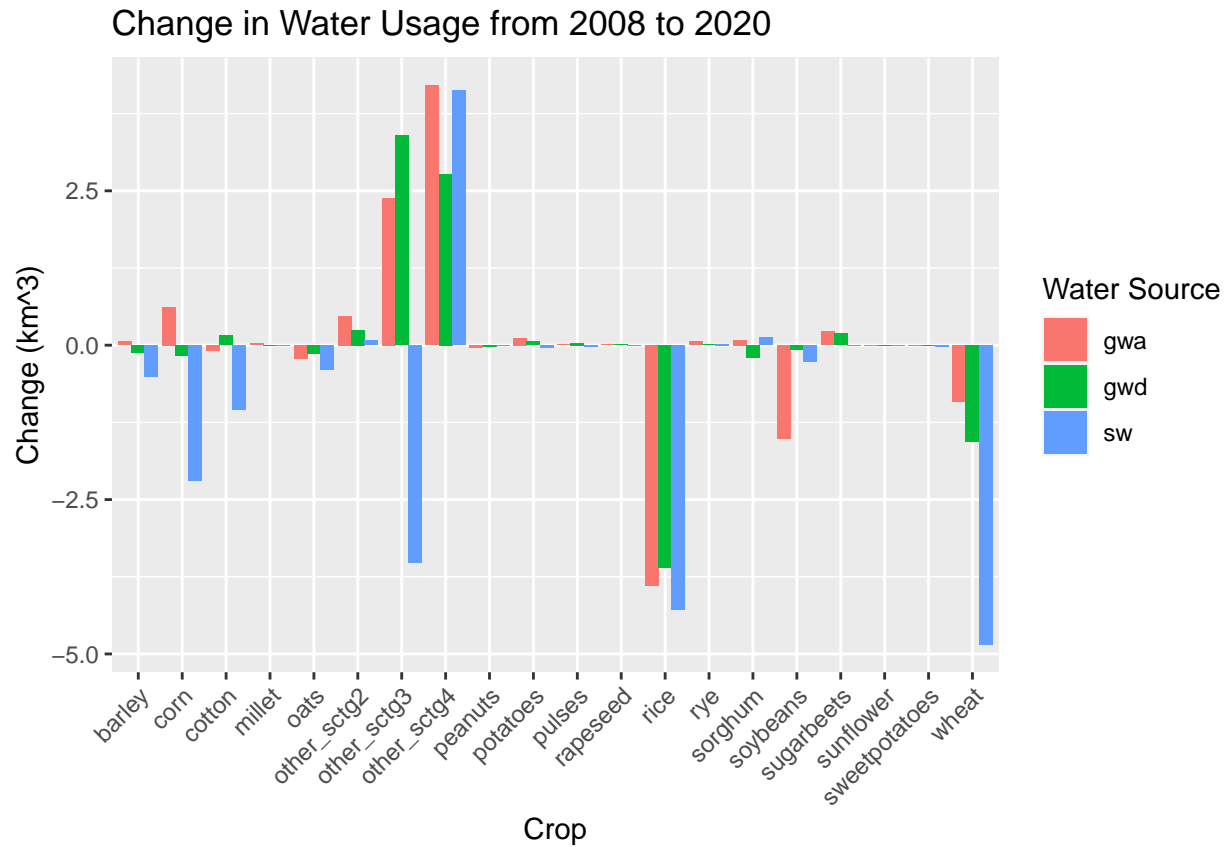
Create a visual representation of the information in Table 2. Create a visualization (or visualizations) that contains mean, change, and percent change in water usage from each crop and source.

```
library(ggplot2)

## plot mean
ggplot(d_summ) +
  geom_bar(aes(fill = src, x = crop, y = mean), stat = "identity",
           position = "dodge") +
  theme(axis.text.x = element_text(angle = 45, hjust=1)) +
  labs(title = "Average Water Usage",
       x = "Crop", y = "Average (km^3)") +
  scale_fill_discrete(name = "Water Source")
```



```
## plot change
ggplot(d_summ) +
  geom_bar(aes(fill = src, x = crop, y = change), stat = "identity",
    position = "dodge") +
  theme(axis.text.x = element_text(angle = 45, hjust=1)) +
  labs(title = "Change in Water Usage from 2008 to 2020",
    x = "Crop", y = "Change (km^3)") +
  scale_fill_discrete(name = "Water Source")
```



```
## plot percent change
ggplot(d_summ) +
  geom_bar(aes(fill = src, x = crop, y = pct_change), stat = "identity",
    position = "dodge") +
  theme(axis.text.x = element_text(angle = 45, hjust=1)) +
  labs(title = "Percent Change in Water Usage from 2008 to 2020",
    x = "Crop", y = "Percent") +
  scale_fill_discrete(name = "Water Source")
```

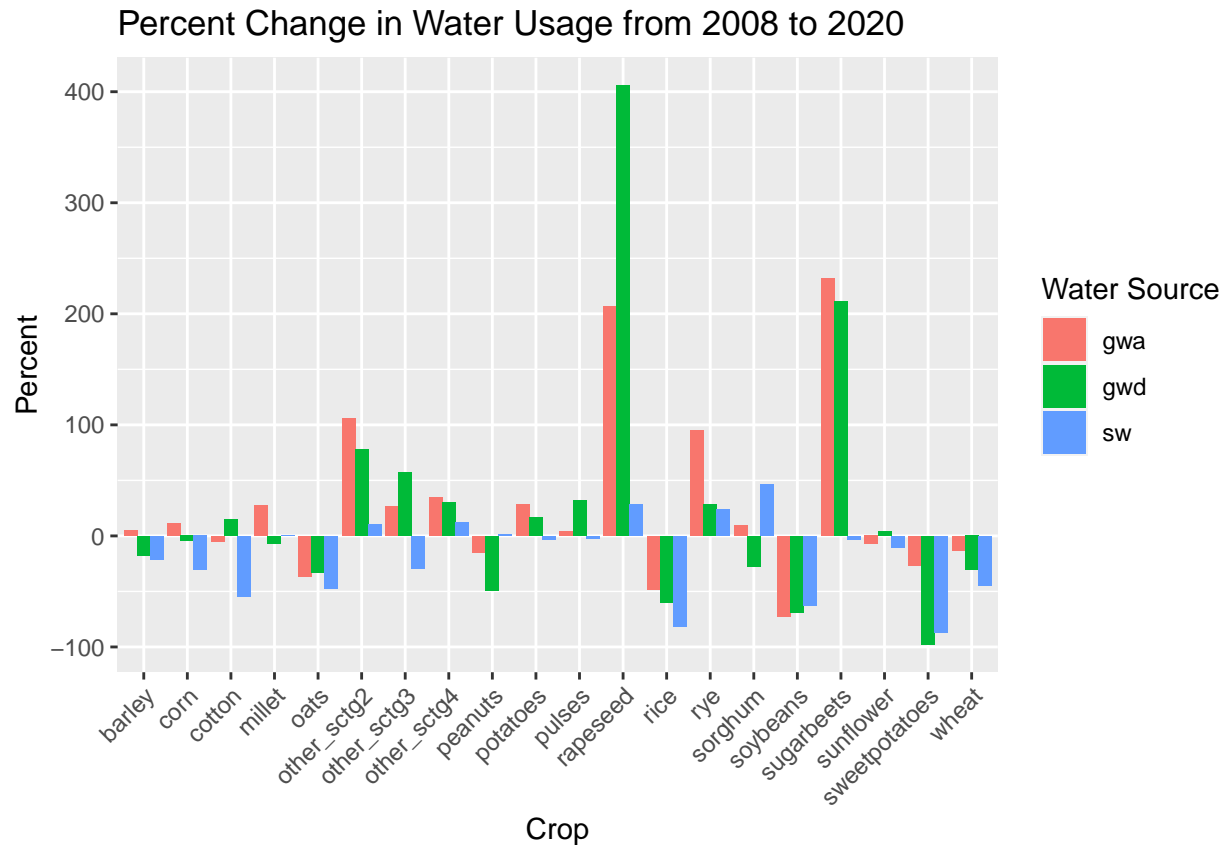


Figure 4

Figure 4 shows the average water usage by crop and source.

- A. average irrigation water usage by source, colored by crop,
- B. average irrigation water usage by crop, colored by source

Two other options for visualizing a numeric variable broken down by two different categorical variable would be a tile plot/grid plot (e.g. <https://github.com/bmacGTPM/pubtheme?tab=readme-ov-file#grid-plot>) and a mosaic plot (<https://haleyjeppson.github.io/ggmosaic/>).

4. Create a tile plot/grid plot of the data in Figure 4.

5. Create a mosaic plot of the data in Figure 4.

6. What are the benefits (other than it fits on one plot) and drawbacks of these two plots?

7. Figure 6

Figure 6 uses a different color scale for each plot. Discuss the benefits and drawbacks of this choice. What was the main purposes of this figure? Given the main purpose, would you recommend using the same color scale, or different color scales, for each plot?

8. Figure 8

Figure 8 also uses a different color scale for each plot. Discuss the benefits and drawbacks of this choice. What was the main purposes of this figure? Given the main purpose, would you recommend using the same color scale, or different color scales, for each plot?

9. Breakdown of GWW

The paper notes in Section 3.1 that $GWW = GWW_{sustainable} + GWW_{unsustainable}$, and that $GWD = GWW_{unsustainable}$. Create a visualization showing the percent of GWW that is GWD for each crop. Use the mean values for water usage.

10. Custom visualization

What is another question you have about this data? Create a visualization that attempt to answer your question.