

Filtering water quality dataset

Amber Lee

6/15/2021

Unfinished

Line 177 in the code is where we should filter for the appropriate QF values. Note how `waterQF_var` is defined (line 71). Use the tutorial (very last link) to use `filter_at()`. Based on how I did the filtering, which was to remove *ALL ROWS* for which the QF code is 8, I only get 53,000 rows.

Summary

One important goal before we interpolate missing values is to filter the data, to decide on the scope of our analysis. Filtering the data to what is necessary will also reduce the amount of interpolation needed. We will see that removing duplicate rows will greatly reduce the missing values (that aren't actually missing), thus making the interpolation step significantly easier.

In cleaning the LTRM Water Quality dataset, we encounter the following questions:

- Which variables (columns) are the most important for us to keep (out of the 133 total variables)?
- Why do duplicate rows happen, and how do we deal with them?
- Which samples (rows) are high enough quality for us to keep? (This involves the QF codes.)

The LTRM Water Quality dataset has duplicate rows for the same SHEETBAR, which is problematic because SHEETBAR is a unique identifier for a water data sheet (sample at a date, time, and location).

```
library(tidyverse)
library(ggplot2)
library(lubridate)
library(corrplot)
library(RColorBrewer)
```

```
water20 <- read.csv(file = "pools EDA/pool data/ltrm_water_data_lat_long.csv")
```

Important variables

```
water_var <- c('TN', 'TP', 'TEMP', 'DO', 'TURB',
              'COND', 'VEL', 'SS', 'WDP', 'CHLcal', 'SECCHI')

waterQF_var <- paste(water_var, "QF", sep = "")

identifier_var <- c('SHEETBAR', 'DATE', 'LATITUDE', 'LONGITUDE')
```

We decided that the 11 continuous variables of importance were: total nitrogen, total phosphorous, temperature, dissolved oxygen, turbidity, water condition, velocity, suspended solids, water depth, chlorophyll-a, and Secchi distance.

In addition, we will want to include the QA/QC codes, along with identifier variables like SHEETBAR and date.

Lastly, we manually edit these variable strings because:

- The water depth variable is WDP, but the corresponding quality factor is ZMAXQF rather than WDPQF.
- CHLcal, calibrated fluorometric chlorophyll a, does not have a corresponding quality factor code. According to the metadata: "CHLcal is generated by calibration of fluorometric chlorophyll readings (CHLF) to season and year specific measurements of spectrophotometric chlorophyll (CHLS). Data from sites where CHLS and CHLF are both collected are used to build river-specific calibration curves for these data. Values are corrected for pheophytin. Units are micrograms per liter."

```
waterQF_var <- waterQF_var[waterQF_var != "WDPQF" &
                           waterQF_var != "CHLcalQF"]

waterQF_var <- c(waterQF_var, "ZMAXQF")
```

Duplicate rows

There are 204305 total rows in the LTRM water quality dataset. Of these rows, there are 156474 distinct SHEETBAR codes.

We visualize duplicates as follows. To identify the duplicate rows, we count the number of occurrences of each unique SHEETBAR value in the dataset. Then, we can calculate and plot the distribution of SHEETBAR duplicates.

```
duplicates <- water20 %>%
  select(SHEETBAR) %>%
  group_by(SHEETBAR) %>%
  summarize(count = n())

duplicates %>% head()

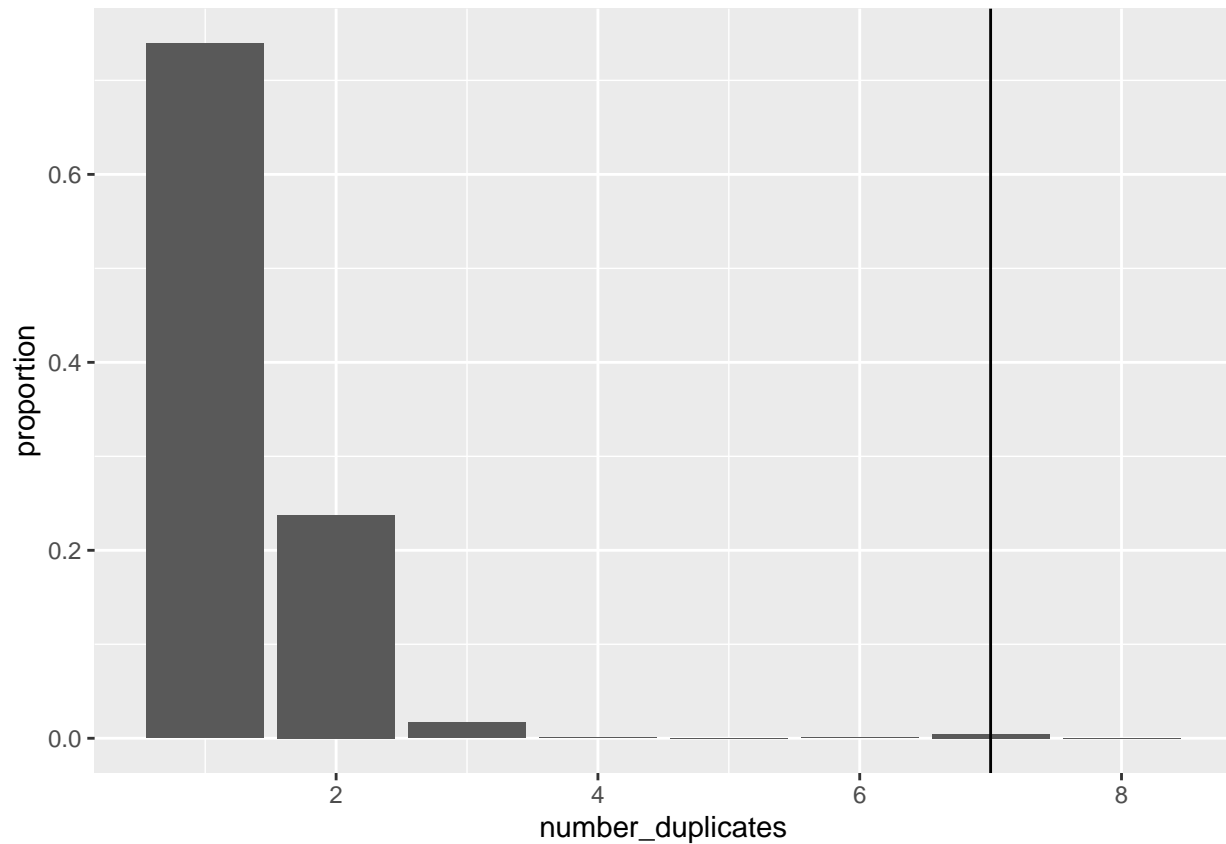
## # A tibble: 6 x 2
##   SHEETBAR count
##   <int> <int>
## 1 -4604348     1
## 2 -4604347     2
## 3 -4604346     2
## 4 -4604345     2
## 5 -4604344     1
## 6 -4604343     2

count_n_duplicates <- function(n, df) {
  return((df %>% filter(count == n) %>% dim())[1]/156474) #156k distinct sheetbars
}

count_duplicates <- data.frame(proportion = sapply(1:8, count_n_duplicates,
                                                    duplicates),
                              number_duplicates = 1:8)

ggplot(count_duplicates, aes(x = number_duplicates, y = proportion)) +
```

```
geom_bar(stat = "identity") +
geom_vline(xintercept = 7)
```



```
count_duplicates # use kable to make nice table output
```

```
##      proportion number_duplicates
## 1 0.7393496683                1
## 2 0.2377072229                2
## 3 0.0166481332                3
## 4 0.0006199113                4
## 5 0.0002364610                5
## 6 0.0007605097                6
## 7 0.0046780935                7
## 8 0.0000000000                8
```

The proportion of SHEETBARs with at least one duplicated row is 0.2606503 (representing about 47,000 rows). When do duplicated rows occur?

We look at two SHEETBARs with duplicated rows.

```
water20 %>%
  filter(SHEETBAR == -4604347) %>%
  select(SHEETBAR, Z, CALCZCD, DO, TP, TN)
```

```
##   SHEETBAR  Z CALCZCD DO   TP   TN
## 1 -4604347 0.2      SF 7.4 0.265 7.49
## 2 -4604347 4.8      BT 7.5   NA   NA
```

```
water20 %>%
  filter(SHEETBAR == 41015929 ) %>%
  select(SHEETBAR, Z, CALCZCD, DO, TP, TN)
```

```
##   SHEETBAR   Z CALCZCD   DO    TP    TN
## 1 41015929 0.2      SF 14.4 0.075 2.557
## 2 41015929 1.0      OT 14.8    NA    NA
## 3 41015929 2.0      OT 14.9    NA    NA
## 4 41015929 3.0      OT 15.0    NA    NA
## 5 41015929 4.0      OT 15.2    NA    NA
## 6 41015929 4.6      BT 15.2    NA    NA
```

Here, we see that TP and TN, total phosphorous and total nitrogen, are measured only at the surface level (when CALCZCD == "SF"). The variable CALCZCD is a categorical variable with levels surface, middle, bottom, and other. It is calculated with the sample depth and the total water depth (of the river site).

In contrast, dissolved oxygen DO is measured at various depths (denoted by Z) because different parts of the water column have different levels of DO. It would be inappropriate to average the dissolved oxygen levels because they were taken at different sample depths.

Thus, this missing values of TP and TN are occurring at different sample depths at the same sampling site. These missing values aren't *really* missing values; they would be redundant to interpolate.

We can reasonably keep only the samples taken at the surface level

We decided to filter for rows that were labelled as surface level, CALCZCD == "SF". Implicitly, this filtering step removes samples for which the sample depth is missing.

```
table(water20$CALCZCD)
```

```
##
##           BT           MD           OT           SF
## 8436 34991 2971 10736 147171
```

More than 70% of the samples were taken on the surface level. Of these measurements taken at the surface level, the eleven important continuous variables (in `water_var`) were recorded with a recording rate of at least 50%. This is a good sanity check because TN and TP are never recorded in the middle and bottom water depths. we checked to see the recording rate of the eleven important variables and found recording rates greater than 50%.

```
filterwater <- water20 %>%
  filter(CALCZCD == "SF") %>%
  select(all_of(water_var))

sapply(filterwater, function(x) sum(is.na(x))/147171))
```

```
##           TN           TP           TEMP           DO           TURB           COND           VEL
## 0.49246115 0.50078480 0.01157837 0.01390220 0.01440501 0.01510488 0.39597475
##           SS           WDP           CHLcal           SECCHI
## 0.17570717 0.07648925 0.26196058 0.09751921
```

What if CALCZCD is missing?

```
(water20 %>%
  filter(CALCZCD == "") %>%
  dim())[1]
```

```
## [1] 8436
```

There are about 8000 samples with missing CALCZCD. We are

Water column variables

WDP,

Filter for QA/QC

Combine

```
filterwater <- water20 %>%
  filter(CALCZCD == "SF") %>%
  select(all_of(c(identifier_var, water_var, waterQF_var))) %>%
  filter_at(vars(contains("QF")), all_vars(. != 8)) # forrest, edit this

filterwater %>% head()
```

```
##  SHEETBAR      DATE LATITUDE LONGITUDE  TN    TP TEMP  DO TURB COND  VEL
## 1 41004753 12/30/1996 44.44746 -92.22680 2.897 0.121 0.1 11.6 4 490 NA
## 2 41004755 12/30/1996 44.41210 -92.10046 3.054 0.112 0.1 11.9 4 484 NA
## 3 41004757 12/30/1996 44.41015 -92.08419 1.282 0.066 0.1 11.0 4 175 NA
## 4 41004761 12/30/1996 44.42327 -92.13400 3.266 0.115 0.1 11.7 3 478 NA
## 5 41004763 12/30/1996 44.32699 -91.93292 2.908 0.105 0.1 10.6 3 446 0.07
## 6 41004764 12/30/1996 44.37975 -91.96503 1.899 0.087 0.1 9.3 6 284 0.06
##  SS WDP CHLcal SECCHI TNQF TPQF TEMPQF DOQF TURBQF CONDQF VELQF SSQF
## 1 2.9 7.82 NA 119 0 0
## 2 2.4 5.50 NA 187 0 0
## 3 2.4 0.71 NA 89 0 0
## 4 2.2 5.33 NA 123 0 0
## 5 1.8 1.10 NA 119 0 0
## 6 3.5 3.54 NA 73 0 0
##  SECCHIQF ZMAXQF
## 1
## 2
## 3 2
## 4
## 5
## 6
```

```
dim(filterwater)
```

```
## [1] 53255 25
```

use this tutorial about `filter_at` <https://suzan.rbind.io/2018/02/dplyr-tutorial-3/>