

Sanity Checks on the GWL Dataset

The first course of action is to analyse the dimensions of `gwl_2023_24.feather`.

```
library(dplyr)
library(feather)
library(ggplot2)
library(rmarkdown)

gwl_2023_24 <- read_feather('../data/gwl_2023_24.feather')
glimpse(gwl_2023_24)

## Rows: 4,308,184
## Columns: 14
## $ concat          <chr> "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54", "#AAXIO54"
## $ station_code    <chr> "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54", "AAXIO54"
## $ name            <chr> "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria", "Himmatpur Andheria"
## $ latitude        <dbl> 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303, 29.3303
## $ longitude       <dbl> 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670, 79.05670
## $ agency         <chr> "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB"
## $ state          <chr> "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand", "Uttarakhand"
## $ district       <chr> "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital", "Nainital"
## $ tahsil         <chr> "-", "-", "-", "-", "-", "-", "-", "-", "-", "-", "-", "-", "-", "-", "-"
## $ datatype_code  <chr> "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ"
## $ description    <chr> "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level"
## $ unit_code      <chr> "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m"
## $ data_time      <dtm> 2023-06-17 18:00:00, 2023-06-18 00:00:00, 2023-06-18 06:00:00, 2023-06-18 12:00:00, 2023-06-18 18:00:00, 2023-06-19 00:00:00, 2023-06-19 06:00:00, 2023-06-19 12:00:00, 2023-06-19 18:00:00, 2023-06-20 00:00:00, 2023-06-20 06:00:00, 2023-06-20 12:00:00, 2023-06-20 18:00:00, 2023-06-21 00:00:00
## $ data_value     <dbl> 162.975, 162.945, 163.095, 162.990, 162.863, 162.736, 162.654, 162.348, 162.105, 162.000, 162.000, 162.000, 162.000, 162.000
```

There are **4,308,184** rows of data for one year! This does not align with our ex-ante knowledge that the order of count of stations is in the *ten thousands*, and around *4 observations* are taken over one year.

```
observations_per_station <- gwl_2023_24 %>%
  count(latitude, longitude, sort = TRUE)
glimpse(observations_per_station)

## Rows: 18,849
## Columns: 3
## $ latitude <dbl> 30.60668, 30.48083, 31.34666, 28.70780, 30.72694, 30.20917, 29.93250, 29.84444, 31
## $ longitude <dbl> 75.88028, 75.93638, 74.69117, 77.24900, 76.80307, 75.86305, 76.09778, 75.93889, 75
## $ n <int> 12924, 12526, 7173, 7089, 6612, 5813, 5809, 5775, 5726, 4716, 4411, 4403, 4400, 43

summary(observations_per_station$n)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       1.0     1.0     3.0   228.6     7.0 12924.0
```

One half of our prior knowledge has been confirmed-there are *18,849 unique lat-long pairs*, each referring to a groundwater level measuring station. All upto the 3rd quartile, every latitude-longitude pair has a single-digit observation count. What happens beyond this quartile is open to question!

```
q3 <- quantile(observations_per_station$n, 0.75)
stations_beyond_75_percentile <- observations_per_station %>%
  filter(n > q3)
glimpse(stations_beyond_75_percentile)
```

```
## Rows: 4,691
## Columns: 3
## $ latitude <dbl> 30.60668, 30.48083, 31.34666, 28.70780, 30.72694, 30.20917, 29.93250, 29.84444, 31
## $ longitude <dbl> 75.88028, 75.93638, 74.69117, 77.24900, 76.80307, 75.86305, 76.09778, 75.93889, 75
## $ n <int> 12924, 12526, 7173, 7089, 6612, 5813, 5809, 5775, 5726, 4716, 4411, 4403, 4400, 43
```

As a test to make things simpler for ourselves, we will consider only the first station (30.60668, 75.88028) and see what is going on with the datapoints from this station.

```
most_observed_station <- gwl_2023_24 %>%
  filter(latitude == 30.60668, longitude == 75.88028) %>%
  arrange(data_time)
glimpse(most_observed_station)
```

```
## Rows: 12,924
## Columns: 14
## $ concat <chr> "#CGWBHR125", "#CGWBPB08", "#CGWBPB09", "#CGWBPB10", "#CGWBPB11", "#CGWBPB185"
## $ station_code <chr> "CGWBHR125", "CGWBPB08", "CGWBPB09", "CGWBPB10", "CGWBPB11", "CGWBPB185", "CGW
## $ name <chr> "Punjab BARNALA MEHLA KALAN CHANANWAL", "Punjab BARNALA BARNALA Barnala (M)",
## $ latitude <dbl> 30.60668, 30.60668, 30.60668, 30.60668, 30.60668, 30.60668, 30.60668, 30.60668
## $ longitude <dbl> 75.88028, 75.88028, 75.88028, 75.88028, 75.88028, 75.88028, 75.88028, 75.88028
## $ agency <chr> "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB", "CGWB"
## $ state <chr> "Punjab", "Punjab", "Punjab", "Punjab", "Punjab", "Punjab", "Punjab", "Punjab"
## $ district <chr> "BARNALA", "BARNALA", "BARNALA", "BARNALA", "BARNALA", "SBS NAGAR", "SBS NAGAR
## $ tahsil <chr> "-", "BARNALA", "BARNALA", "-", "-", "-", "-", "NAWANSHAHR", "-", "-", "BARNAL
## $ datatype_code <chr> "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "GGZ", "
## $ description <chr> "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level", "GPRS-Water Level"
## $ unit_code <chr> "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m", "m"
## $ data_time <dtm> 2023-05-01 00:00:00, 2023-05-01 00:00:00, 2023-05-01 00:00:00, 2023-05-01 00:
## $ data_value <dbl> -36.519, -44.887, -43.131, -40.021, -35.517, -18.710, -19.918, -6.693, -18.693
```

The `glimpse()` function is showing us at a glance that the Punjab Barnala Mehla Kalan Chananwal station has been taking readings at a much higher frequency than was expected! For good measure, we will confirm if the time stamps are unique throughout those 12,924 readings:

```
timestamp_count <- most_observed_station %>%
  summarize(count = n_distinct(data_time)) %>%
  pull(count)
timestamp_count
```

```
## [1] 1493
```

A fraction of those datapoints have unique timestamps! Are these duplicates?

```
## Rows: 1,493  
## Columns: 4  
## $ data_time <dtm> 2023-05-01 00:00:00, 2023-05-01 06:00:00, 2023-05-01 12:00:00, 2023-05-01 18:00:  
## $ n <int> 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9,  
## $ n_distinct <int> 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9, 9,  
## $ mean <dbl> -29.34322, -29.45856, -29.52244, -29.43967, -29.35944, -29.43967, -29.43200, -29.43200,
```

```
unequal_timestamps <- readings_per_time %>%
  filter(n != n_distinct) %>%
  arrange(desc(n - n_distinct))
glimpse(unequal_timestamps)
```

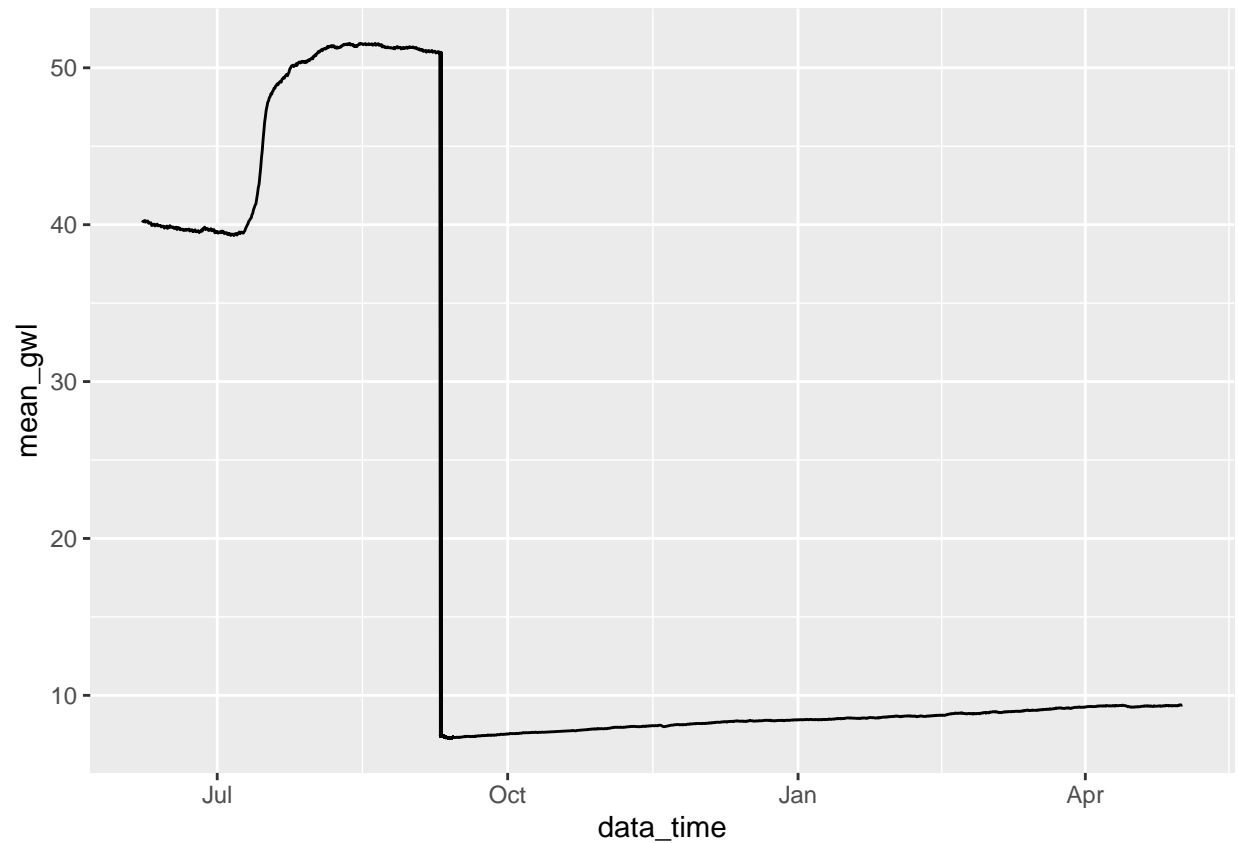
```
## Rows: 2
## Columns: 4
## $ data_time  <dtm> 2023-09-09 18:00:00, 2024-04-20 18:00:00
## $ n          <int> 9, 9
## $ n_distinct <int> 8, 8
## $ mean       <dbl> -31.56222, -30.44656
```

```
random_datetime <- most_observed_station %>%
  filter(data_time < as.POSIXct('2023-05-01 06:00:00'))
summary(random_datetime$data_value)
```

That is surely a large range of observations for a measurement that had been allegedly taken at the same location and time! This is a conclusion on the quality of data collected at the Punjab Barnala Mehla Kalan Chananwal station. Such an exercise can be extended to other stations as well, which also report an abnormally high amount of data.

Averaging over spatio-temporal locations

3



Initially thought to be a pointless exercise, the data visualisation has brought us a peculiar graph, and brings us to ask stranger questions: what happened to the readings in September?

In any case, the time-stamps are at the heart of the true reason as to why the dataset is this massive.

We will conclude here briefly until better ideas come up to manage this dataset.

- Aggregate data on time-stamps?