

ETC5510: Introduction to Data Analysis

Week 5, part A

Missing Data

Lecturer: *Nicholas Tierney and Stuart Lee*

Department of Econometrics and Business Statistics

✉ ETC5510.Clayton-x@monash.edu

April 2020



Recap

- Joins
- advanced data vis

Example

San Francisco weather data

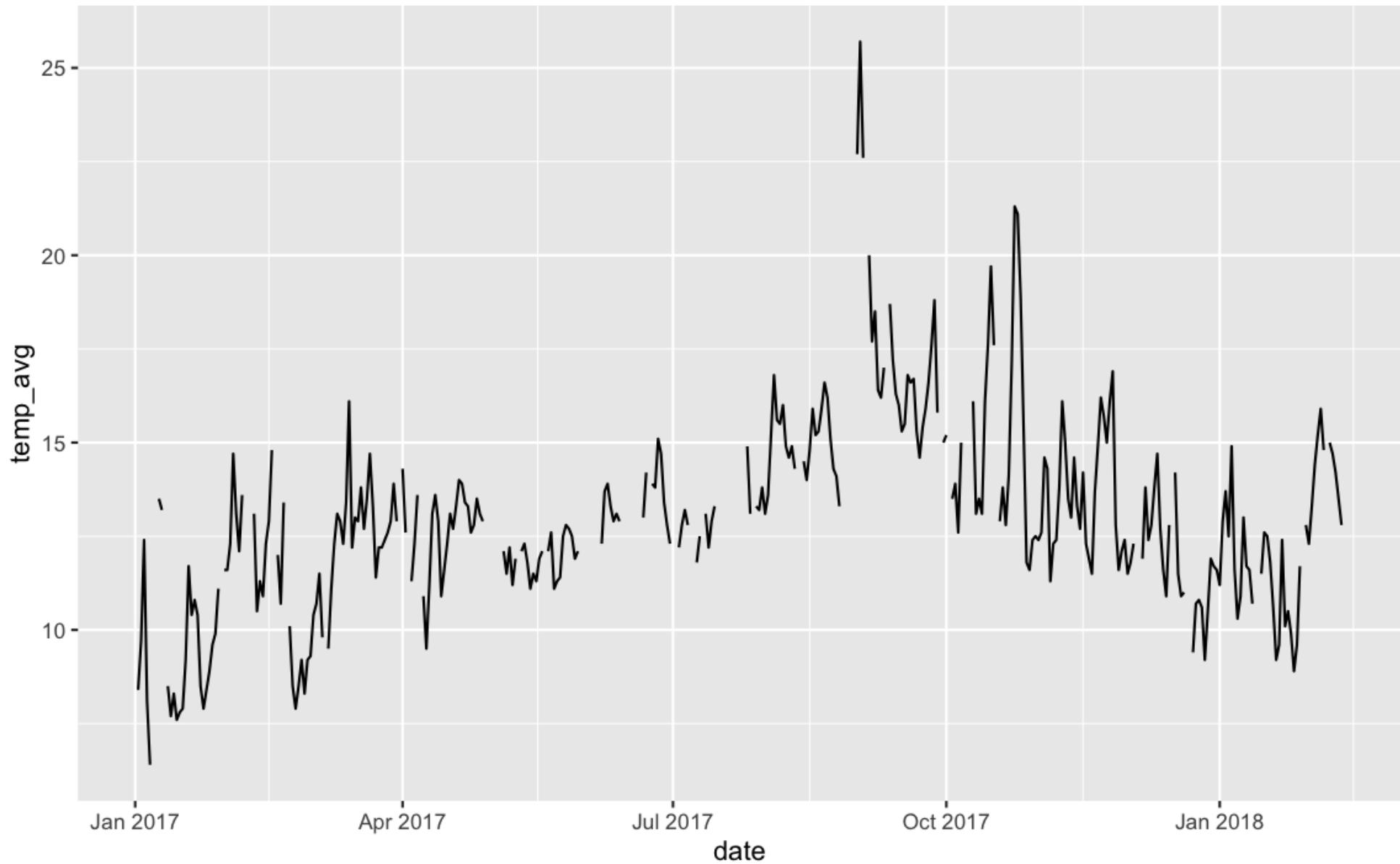
|| Date | Wind | Temp ||

Using the R package: [GSODR](#)

([Global Surface Summary of the Day](#)).

Written by Adam Sparks

github.com/ropensci/GSODR

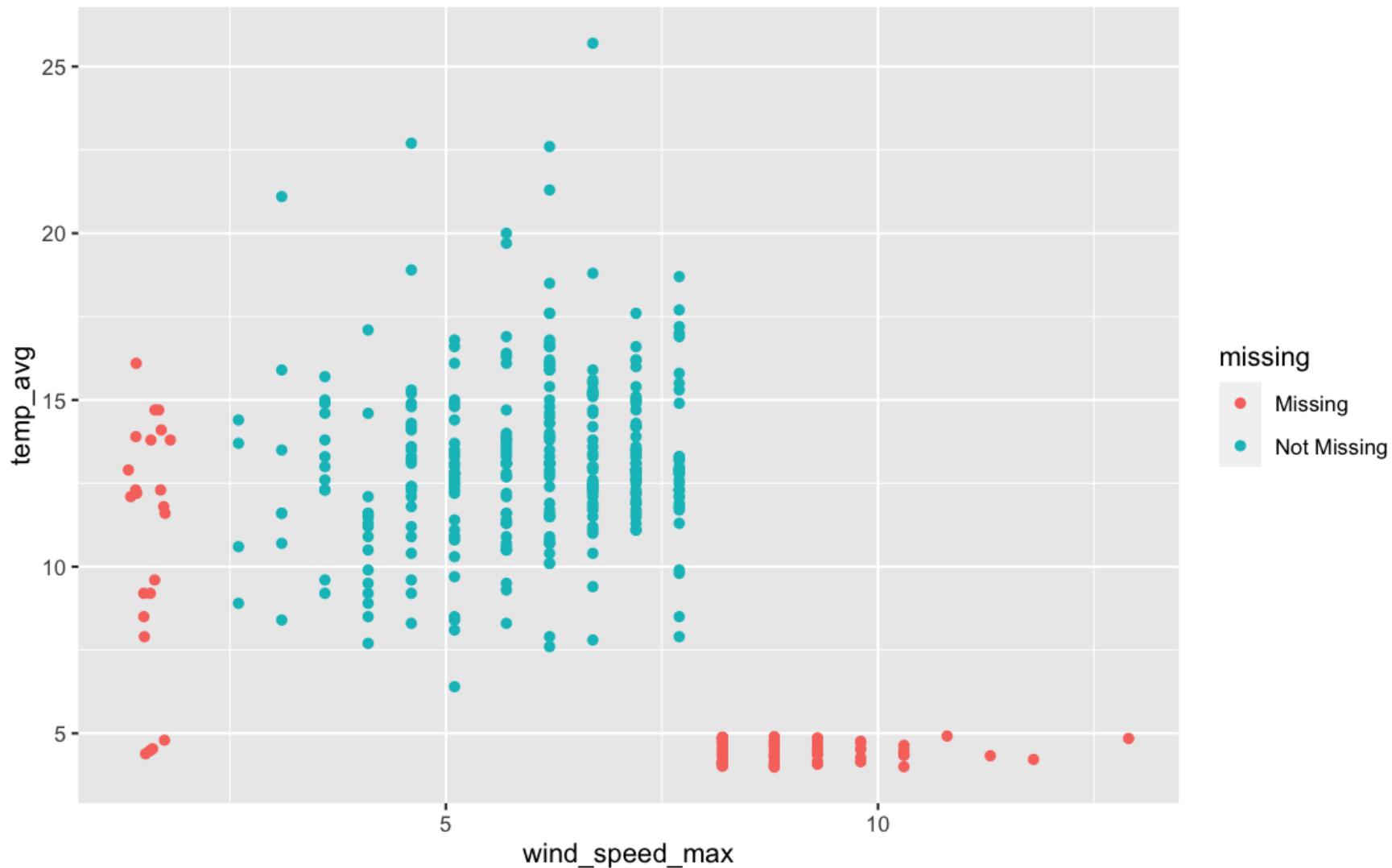


Your Turn: These gaps are missing values! What are some reasons this might be a problem?

Some thoughts

- What is missing?
- Why are they missing?
- How can we summarise and explore this?

One way to show missing data



Wait, What?

What people think dealing
with missing data looks
like



What dealing with missing data actually looks like

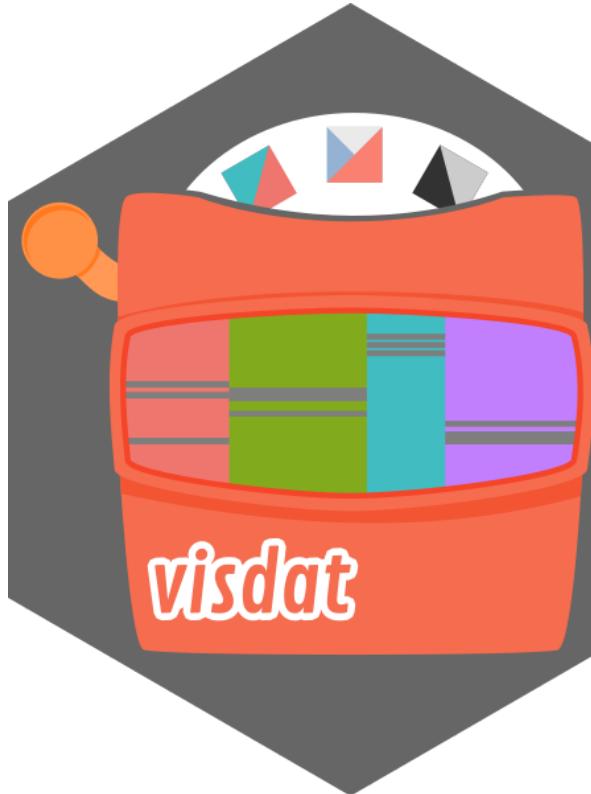


What I want dealing with missing data to be like



@草根足球记录员

Learn more



visdat.njtierney.com



naniar.njtierney.com

Overview

1. What even are missing values
2. How to start looking at missing data
3. How to start exploring missing data
4. How to impute (fill in) Missing values

What are missing values?

Missing values are values that should have been recorded but were not.

NA = Not Available.

How do I check if I have missing values?

```
x <- c(1, NA, 3, NA, NA, 5)

library(naniar)
any_na(x)

[1] TRUE

are_na(x)

[1] FALSE  TRUE FALSE  TRUE  TRUE FALSE

n_miss(x)

[1] 3

prop_miss(x)

[1] 0.5
```

Working with missing data

NA + [anything] = NA

```
heights
```

```
Sophie      Dan      Fred  
  165       177      NA
```

```
sum(heights)
```

```
[1] NA
```

Working with missing data

`na.rm = TRUE` will removes missings

```
sum(heights, na.rm = TRUE)
```

```
[1] 342
```

Use this power responsibly!

Dangers of removing missing values

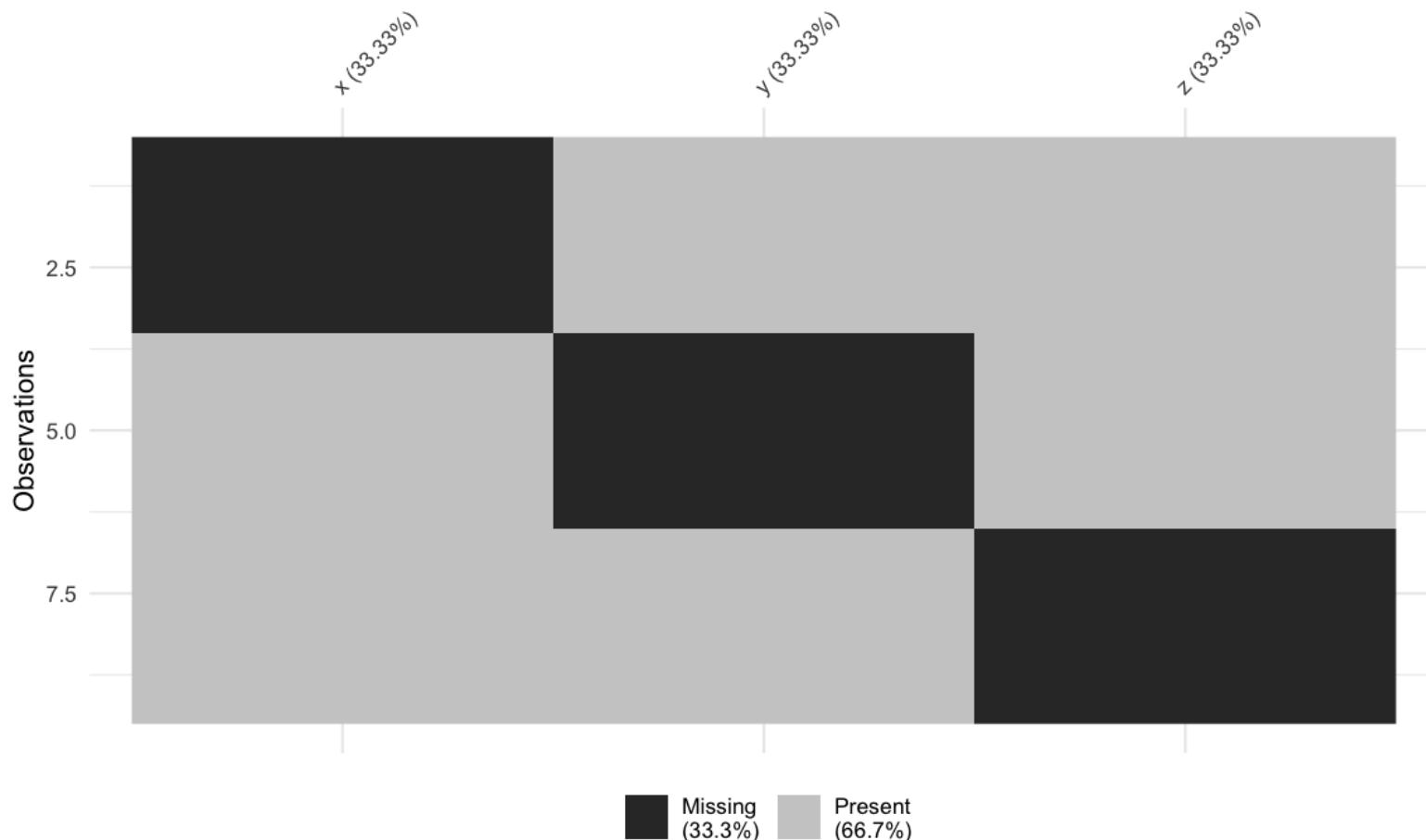
You can remove most of or all of your data:

| x | y | z |
|----|----|----|
| NA | 1 | 1 |
| NA | 2 | 2 |
| NA | 3 | 3 |
| 4 | NA | 4 |
| 5 | NA | 5 |
| 6 | NA | 6 |
| 7 | 7 | NA |
| 8 | 8 | NA |
| 9 | 9 | NA |

Dangers of removing missing values

You can remove most of or all of your data:

```
vis_miss(dat_df)
```



Dangers of removing missing values

You can remove most of or all of your data:

```
na.omit(dat_df)
## [1] x y z
## <0 rows> (or 0-length row.names)
```

wat?

na.omit / na.rm = listwise deletion

| x | y | z |
|----|----|----|
| NA | 1 | 1 |
| NA | 2 | 2 |
| NA | 3 | 3 |
| 4 | NA | 4 |
| 5 | NA | 5 |
| 6 | NA | 6 |
| 7 | 7 | NA |
| 8 | 8 | NA |
| 9 | 9 | NA |

na.omit / na.rm = listwise deletion

| x | y | z |
|----|----|----|
| NA | 1 | 1 |
| NA | 2 | 2 |
| NA | 3 | 3 |
| 4 | NA | 4 |
| 5 | NA | 5 |
| 6 | NA | 6 |
| 7 | 7 | NA |
| 8 | 8 | NA |
| 9 | 9 | NA |

na.omit / na.rm = listwise deletion

| x | y | z |
|----|----|----|
| NA | 1 | 1 |
| NA | 2 | 2 |
| NA | 3 | 3 |
| 4 | NA | 4 |
| 5 | NA | 5 |
| 6 | NA | 6 |
| 7 | 7 | NA |
| 8 | 8 | NA |
| 9 | 9 | NA |

na.omit / na.rm = listwise deletion

| x | y | z |
|----|----|----|
| NA | 1 | 1 |
| NA | 2 | 2 |
| NA | 3 | 3 |
| 4 | NA | 4 |
| 5 | NA | 5 |
| 6 | NA | 6 |
| 7 | 7 | NA |
| 8 | 8 | NA |
| 9 | 9 | NA |

Takehome:

- na.rm or na.omit can remove entire rows containing missings
- This is bad because you can lose data - sometimes all your data!
This might not be what you anticipate!
- It can also mean that you are removing / censoring observations.

Dangers of removing missing values

You can introduce bias - what happens when you remove the NAs?

| temp | location |
|------|----------|
| 27 | inside |
| 26 | inside |
| NA | outside |
| 29 | inside |
| NA | outside |
| 20 | outside |
| 21 | outside |
| 24 | inside |

Your turn:

- Open rstudio.
- go to `exercise-5a-intro-missing.Rmd`
- type the following:

```
# install.packages("usethis")
library(usethis)
use_course("https://mida.numbat.space/exercises/5a/mida-exercise-5a.zip")
```

Introduction to missingness summaries

Basic summaries of missingness:

- n_miss
- n_complete

Dataframe summaries of missingness:

- miss_var_summary
- miss_case_summary

These functions work with group_by

Missing data summaries: Variables

```
miss_var_summary(dat_sf_clean)
## # A tibble: 6 x 3
##   variable      n_miss pct_miss
##   <chr>        <int>    <dbl>
## 1 temp_min      70     17.3
## 2 temp_max      70     17.3
## 3 temp_avg      70     17.3
## 4 wind_speed_max 23     5.68
## 5 date          0      0
## 6 month         0      0
```

Missing data summaries: Cases

```
miss_case_summary(dat_sf_clean)
## # A tibble: 405 x 3
##       case n_miss pct_miss
##   <int>   <int>     <dbl>
## 1     89      4    66.7
## 2    182      4    66.7
## 3    188      4    66.7
## 4    271      4    66.7
## 5      6      3    50
## 6      7      3    50
## 7     10      3    50
## 8     29      3    50
## 9     37      3    50
## 10    39      3    50
## # ... with 395 more rows
```

Missing data tabulations: variables

```
miss_var_table(dat_sf_clean)
## # A tibble: 3 x 3
##   n_miss_in_var n_vars pct_vars
##       <int>    <int>     <dbl>
## 1          0      2     33.3
## 2         23      1     16.7
## 3         70      3     50
```

Missing data tabulations: cases

```
miss_case_table(dat_sf_clean)
## # A tibble: 4 x 3
##   n_miss_in_case n_cases pct_cases
##       <int>     <int>      <dbl>
## 1             0     316    78.0
## 2             1      19     4.69
## 3             3      66    16.3
## 4             4       4    0.988
```

Using summaries with group_by

```
dat_sf_clean %>%  
  group_by(month) %>%  
  miss_var_summary()  
## # A tibble: 60 x 4  
## # Groups: month [12]  
##   month variable     n_miss pct_miss  
##   <dbl> <chr>       <int>    <dbl>  
## 1 1     temp_min      7     11.5  
## 2 1     temp_max      7     11.5  
## 3 1     temp_avg      7     11.5  
## 4 1     wind_speed_max 4     6.56  
## 5 1     date          0      0  
## 6 2     temp_min      5     12.8  
## 7 2     temp_max      5     12.8  
## 8 2     temp_avg      5     12.8  
## 9 2     wind_speed_max 4     10.3  
## 10 2    date          0      0  
## # ... with 50 more rows
```

Your Turn

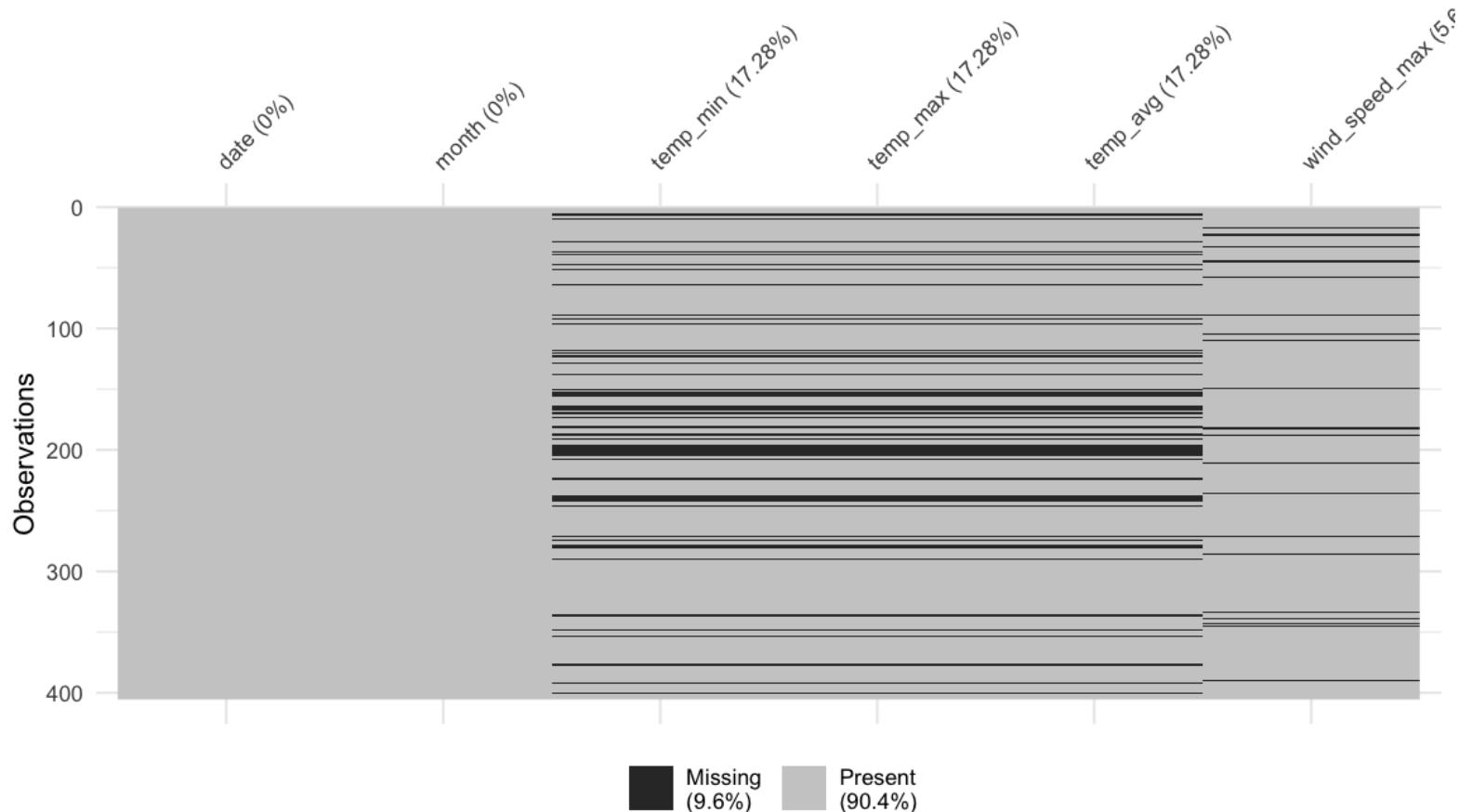
Open exercise-5a-summarise-missings.Rmd

Introduction to missing data visualisations in naniar

- Visualisation can quickly capture an idea or thought.
- naniar provides a friendly family of missing data visualization functions.
- Each visualization corresponds to a data summary.
- Visualisations help you operate closer to the speed of thought.

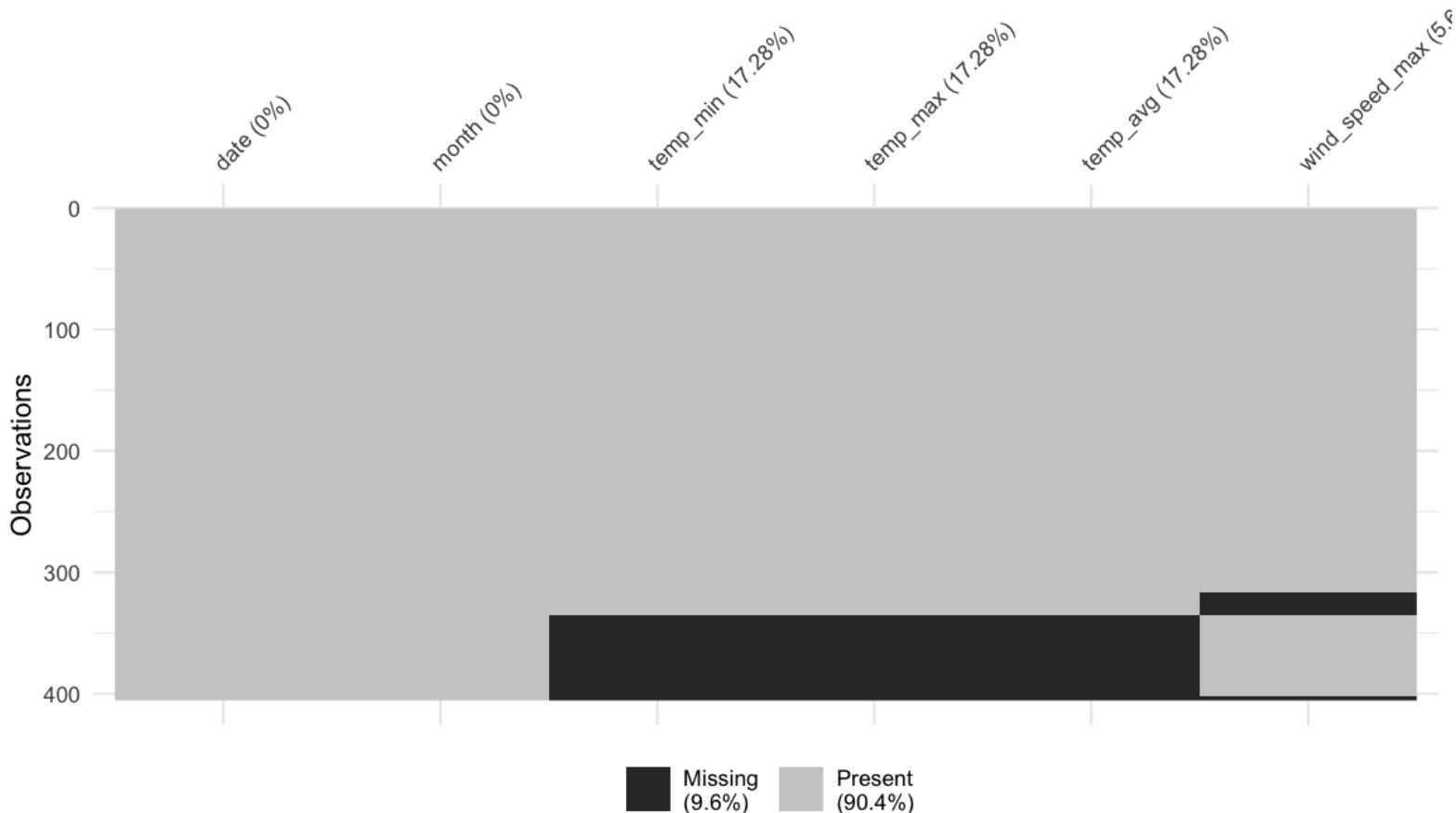
Get a bird's eye view of the missing data

```
vis_miss(dat_sf_clean)
```



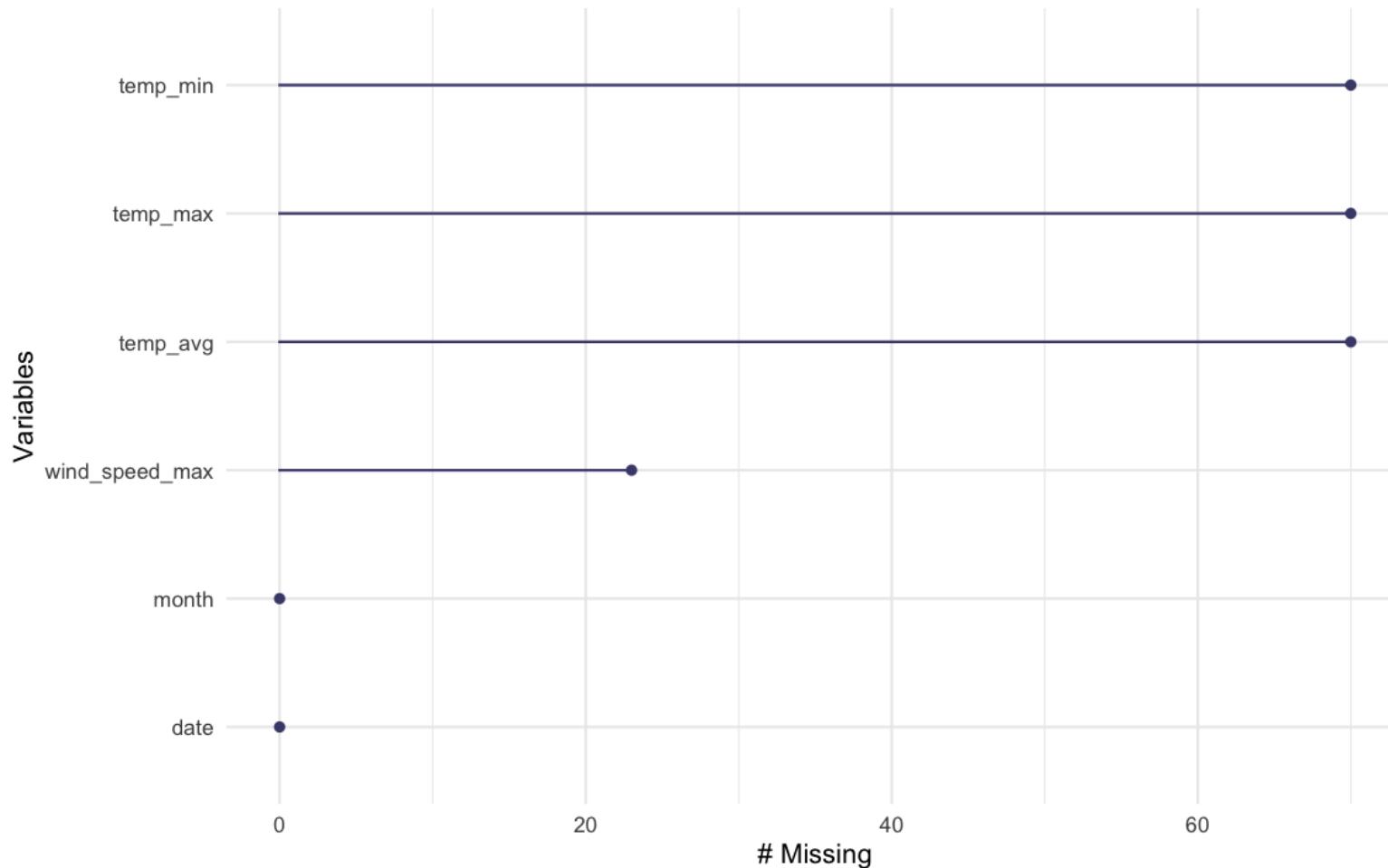
Get a bird's eye view of the missing data

```
vis_miss(dat_sf_clean, cluster = TRUE)
```



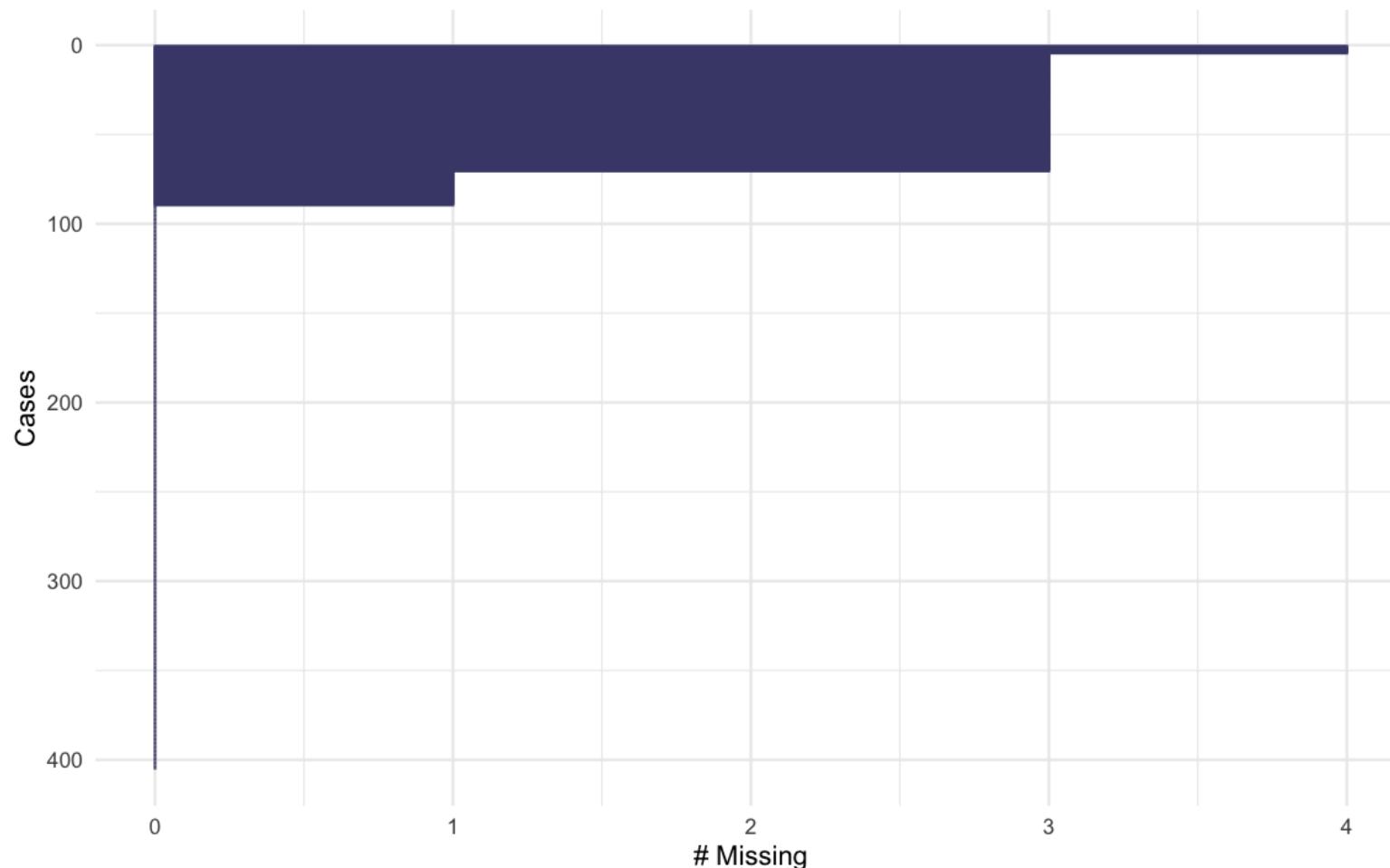
Look at missings in cases

```
gg_miss_var(dat_sf_clean)
```



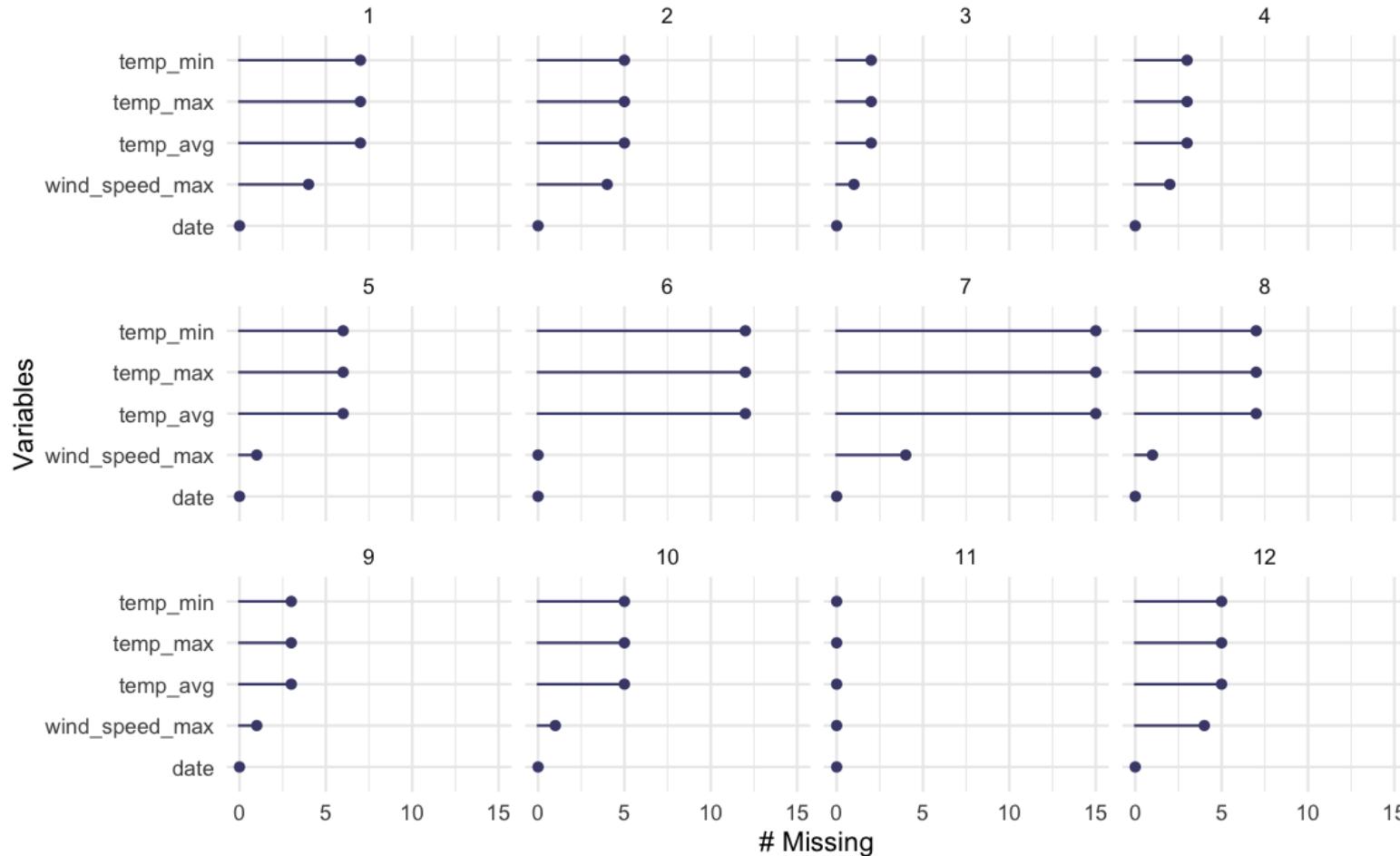
Look at missings in cases

```
gg_miss_case(dat_sf_clean)
```



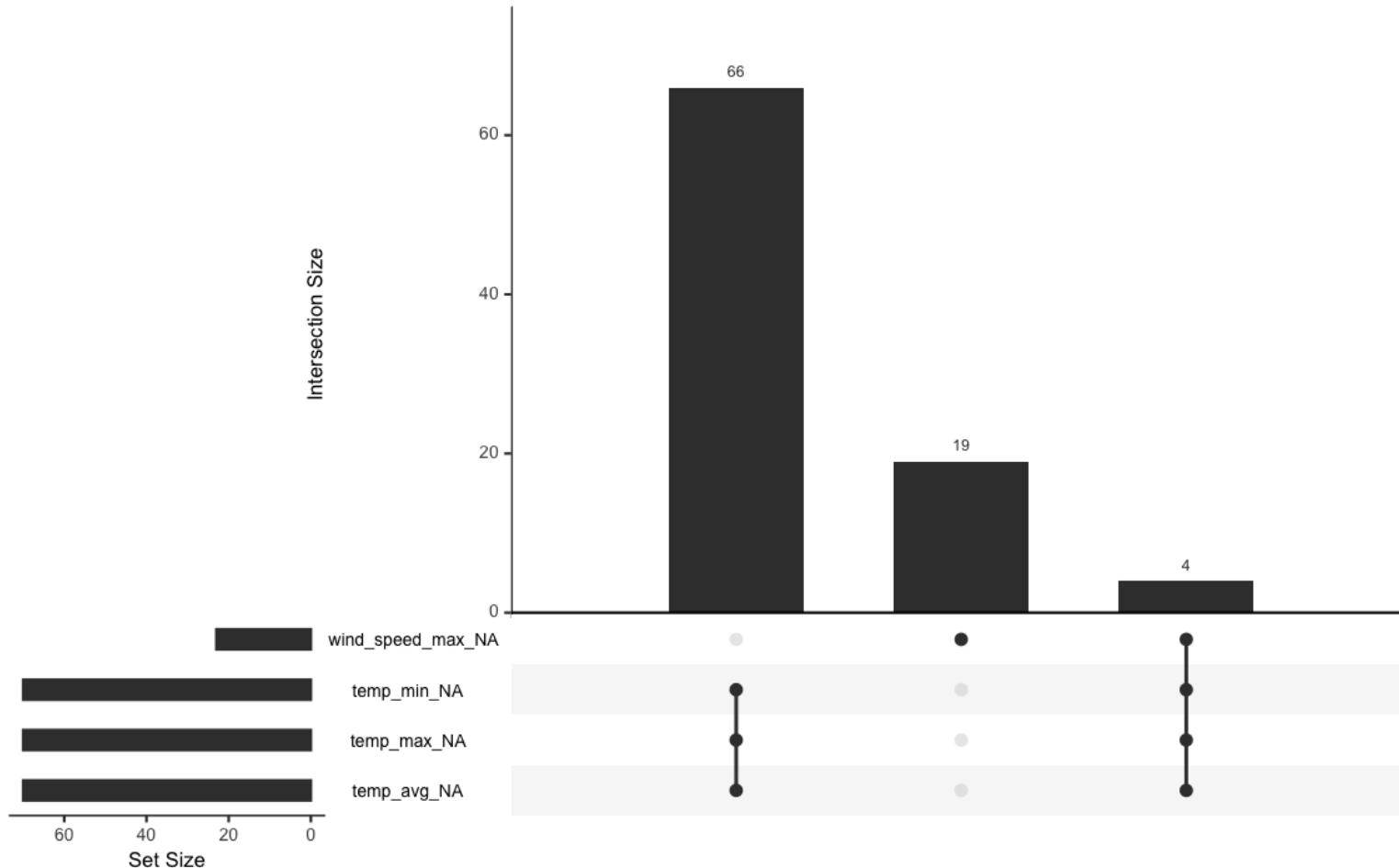
Look at missings in variables

```
gg_miss_var(dat_sf_clean, facet = month)
```



Visualizing missingness patterns

```
gg_miss_upset(dat_sf_clean)
```



Your Turn:

- lab quiz open (requires answering questions from Lab exercise)
- go to rstudio and finish final exercise

Resources

- [R-miss-Tastic](#)
- [naniar](#)
- [visdat](#)