

# **Week 5: Data Wrangling**

PM 566: Introduction to Health Data Science

Kelly Street



# Today's goals

We will learn about how to manipulate data, and in particular,

- Selecting variables.
- Filtering data.
- Creating variables.
- Summarize data.

Throughout the session we will see examples using:

- `data.table` in R,
- `dtplyr` in R, and
- `pydatatable`

All with the `MET` dataset.

# Disclaimer

There's a lot of extraneous information in these slides! While the `data.table` package and Python both have a lot of useful functionality, we strongly recommend sticking to the base R and `tidyverse` tools presented here. Slides covering material outside this scope will be marked with an asterisk (\*); you should be extremely cautious about using code from those slides!

# Data wrangling in R

Overall, you will find the following approaches:

- **base R**: Use only base R functions.
- **dplyr**: Using “verbs”.
- **data.table**: High-performing (ideal for large data)
- **dplyr + data.table = dtplyr**: High-performing + dplyr verbs.

Other methods involve, for example, using external tools such as [Spark](#), [sparkly](#).

We will be focusing on data.table because of [this](#)

Take a look at this very neat cheat sheet by [Erik Petrovski here](#).

# Selecting variables: Load the packages

```
1 library(data.table)
2 library(dtplyr)
3 library(dplyr)
4 library(ggplot2)
```

The `dtplyr` R package translates `dplyr` (`tidyverse`) syntax to `data.table`, so that we can still use **the dplyr verbs** while at the same time leveraging the performance of `data.table`.

# Loading the data

The data that we will be using is an already processed version of the MET dataset. We can download (and load) the data directly in our session using the following commands:

```
1 # Where are we getting the data from
2 met_url <- "https://github.com/USCbiostats/data-science-data/raw/master/02_met/met_all.gz"
3
4 # Downloading the data to a tempfile (so it is destroyed afterwards)
5 # you can replace this with, for example, your own data:
6 tmp <- tempfile(pattern = "met", fileext = ".gz")
7 # tmp <- "met.gz"
8
9 # We should be downloading this, ONLY IF this was not downloaded already.
10 # otherwise is just a waste of time.
11 if (!file.exists(tmp)) {
12   download.file(
13     url      = met_url,
14     destfile = tmp,
15     # method  = "libcurl", timeout = 1000 (you may need this option)
16   )
17 }
```

Now we can load the data using the `read.csv()` or `fread()` functions.

# Read the Data

## In R (base)

```
1 # Reading the data
2 dat <- read.csv(tmp)
3 head(dat)
```

## In R (`data.table`)

```
1 # Reading the data
2 dat <- fread(tmp)
3 head(dat)
4 dat <- as.data.frame(dat)
```

## In Python

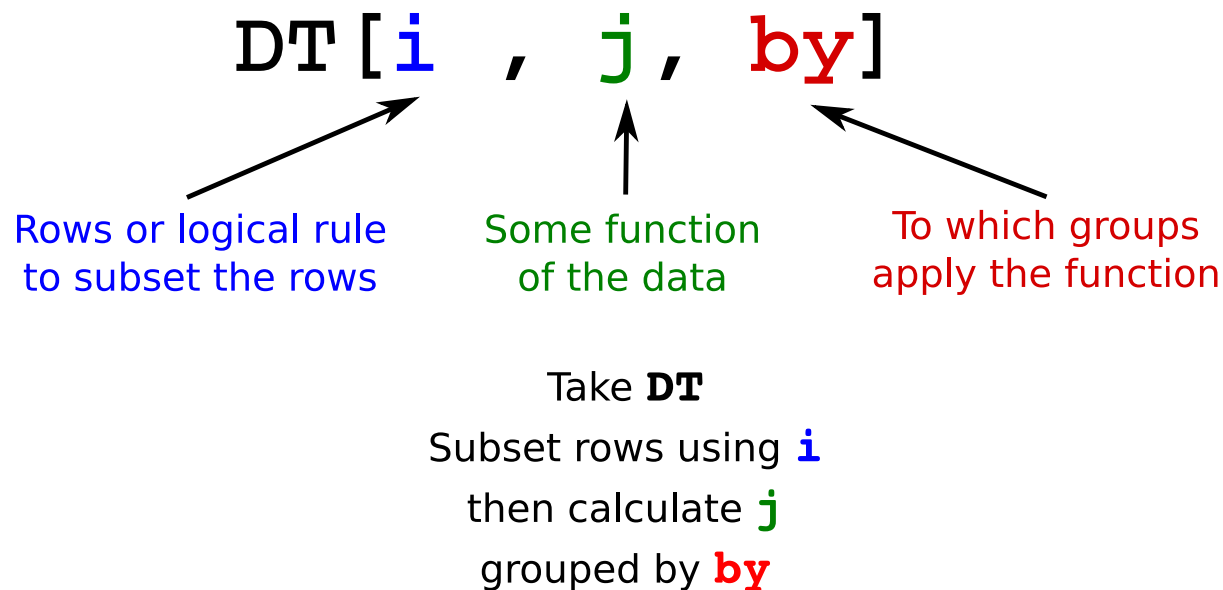
```
1 import datatable as dt
2 dat = dt.fread("met.gz")
3 dat.head(5)
```

Before we continue, let's learn a bit more on `data.table` and `dtplyr`



# \* `data.table` and `dtplyr`: Data Table's Syntax

- As you have seen in previous lectures, in `data.table` all happens within the square brackets. Here is common way to imagine DT:



- Any time that you see `:=` in `j` that is "Assignment by reference." Using `=` within `j` only works in some specific cases.

# \* `data.table` and `dtplyr`: Data Table's Syntax

Operations applied in `j` are evaluated *within* the data, meaning that names work as symbols, e.g.,

```
1 data("USArrests")
2 USArrests_dt <- data.table(USArrests)
3 # This returns an error
4 USArrests[, Murder]
5 # This works fine
6 USArrests_dt[, Murder]
```

Furthermore, we can do things like this:

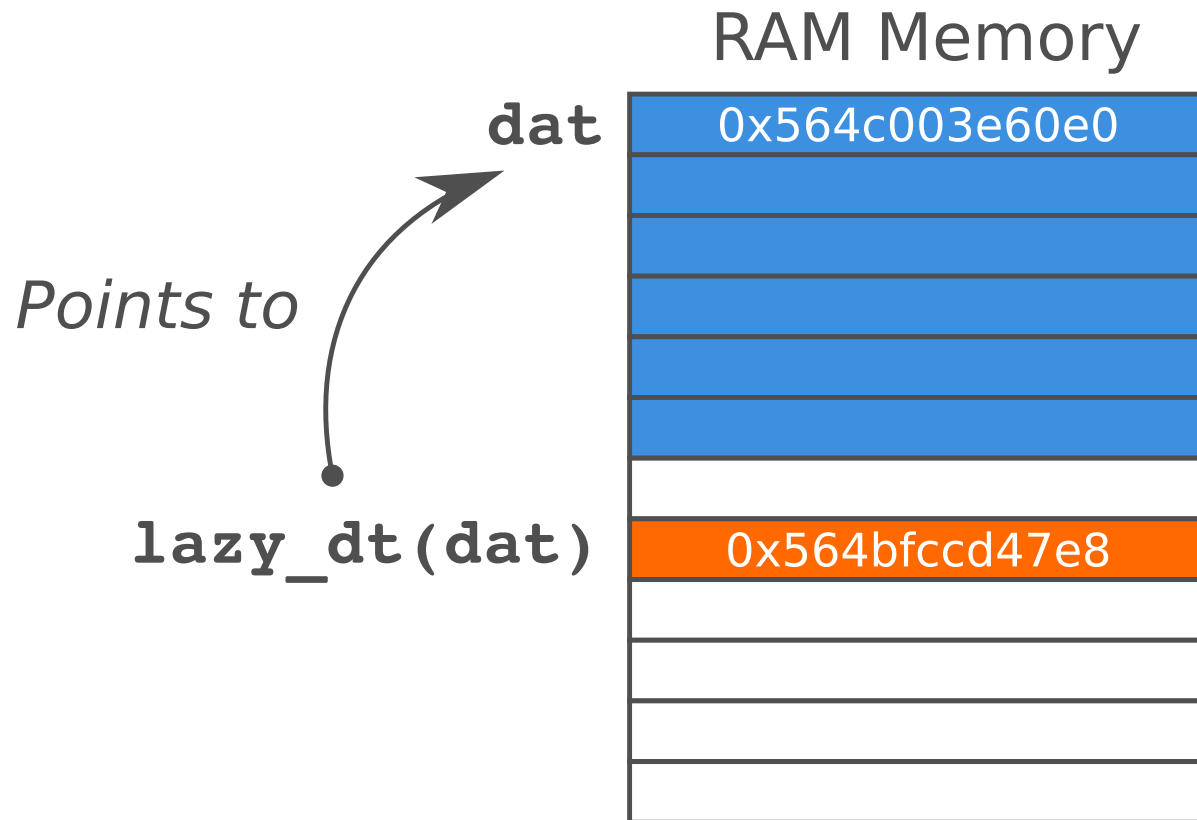
```
1 USArrests_dt[, plot(Murder, UrbanPop)]
```

NULL

# \* `data.table` and `dtplyr`: Lazy table

- The `dtplyr` package provides a way to translate `dplyr` verbs to `data.table` syntax.
- The key lies on the function `lazy_dt` from `dtplyr` (see ?  
`dtplyr::lazy_dt`).
- This function creates a wrapper that “points” to a `data.table` object

# \* `data.table` and `dtplyr`: Lazy table (cont.)



The lazy table only points to the actual data, avoiding duplicating memory.

Lazy tables have their own address, but the underlying object has the same address as the original `data.table`

Question: What is the `immutable = FALSE` option used for?

# \* Selecting columns

How can we select the columns `USAFID`, `lat`, and `lon`, using `data.table`:

```
1 dat[, list(USAFID, lat, lon)]  
2 # dat[, .(USAFID, lat, lon)]      # Alternative 1  
3 # dat[, c("USAFID", "lat", "lon")] # Alternative 2
```

What happens if instead of `list()` you used `c()`?

# Selecting columns (cont. 1)

Using base R:

```
1 dat[, c('USAFID', 'lat', 'lon')]
```

|    | USAFID | lat    | lon      |
|----|--------|--------|----------|
| 1  | 690150 | 34.300 | -116.166 |
| 2  | 690150 | 34.300 | -116.166 |
| 3  | 690150 | 34.300 | -116.166 |
| 4  | 690150 | 34.300 | -116.166 |
| 5  | 690150 | 34.300 | -116.166 |
| 6  | 690150 | 34.300 | -116.166 |
| 7  | 690150 | 34.300 | -116.166 |
| 8  | 690150 | 34.300 | -116.166 |
| 9  | 690150 | 34.300 | -116.166 |
| 10 | 690150 | 34.300 | -116.166 |
| 11 | 690150 | 34.296 | -116.162 |
| 12 | 690150 | 34.296 | -116.162 |
| 13 | 690150 | 34.300 | -116.166 |

# Selecting columns (cont. 2)

Using the **dplyr::select** verb:

```
1 dat |> select(USAFID, lat, lon)
```

|    | USAFID | lat    | lon      |
|----|--------|--------|----------|
| 1  | 690150 | 34.300 | -116.166 |
| 2  | 690150 | 34.300 | -116.166 |
| 3  | 690150 | 34.300 | -116.166 |
| 4  | 690150 | 34.300 | -116.166 |
| 5  | 690150 | 34.300 | -116.166 |
| 6  | 690150 | 34.300 | -116.166 |
| 7  | 690150 | 34.300 | -116.166 |
| 8  | 690150 | 34.300 | -116.166 |
| 9  | 690150 | 34.300 | -116.166 |
| 10 | 690150 | 34.300 | -116.166 |
| 11 | 690150 | 34.296 | -116.162 |
| 12 | 690150 | 34.296 | -116.162 |
| 13 | 690150 | 34.300 | -116.166 |

# \* Selecting columns (cont. 3)

In the case of `pydatatable`

```
1 dat[:, ["USAFID", "lat", "lon"]]
```

What happens if instead of `["USAFID", "lat", "lon"]` you used `{"USAFID", "lat", "lon"}` (vector vs set).



For the rest of the session we will be using these variables:  
USAFID, WBAN, year, month, day, hour, min, lat, lon, elev,  
wind.sp, temp, and atm.press.

```
1 # select only the relevant variables
2 dat <- dat |>
3   select(USAFID, WBAN, year, month, day,
4         hour, min, lat, lon, elev,
5         wind.sp, temp, atm.press)
```

# Data filtering: Logical conditions

- Based on logical operations, e.g. `condition 1 [and|or condition2 [and|or ...]]`
- Need to be aware of ordering and grouping of `and` and `or` operators.
- Fundamental **logical** operators:

| x     | y     | Negate<br>!x | And<br>x & y | Or<br>x   y | Xor<br>xor(x, y) |
|-------|-------|--------------|--------------|-------------|------------------|
| true  | true  | false        | true         | true        | false            |
| false | true  | true         | false        | true        | true             |
| true  | false | false        | false        | true        | true             |
| false | false | true         | false        | false       | false            |

- Fundamental **relational** operators, in R: `<`, `>`, `<=`, `>=`, `==`, `!=`.

# Questions 1: How many ways can you write an XOR operator?

Write a function that takes two arguments ( $x, y$ ) and applies the XOR operator element wise. Here you have a template:

```
1 myxor <- function(x, y) {  
2   res <- logical(length(x))  
3   for (i in 1:length(x)) {  
4     res[i] <- # do something with x[i] and y[i]  
5   }  
6   return(res)  
7 }
```

Or if vectorized (which would be better)

```
1 myxor <- function(x, y) {  
2   # INSERT YOUR CODE HERE  
3 }
```

Hint 1: Remember that negating ( $x \& y$ ) equals  $(!x \mid !y)$ .

Hint 2: Logical operators are distributive, meaning  $a * (b + c) = (a * b) + (a * c)$ , where  $*$  and  $+$  are  $\&$  or  $\mid$ .

# In R

```
1 myxor1 <- function(x,y) {(x & !y) | (!x & y)}
2 myxor2 <- function(x,y) {!((!x | y) & (x | !y))}
3 myxor3 <- function(x,y) {(x | y) & (!x | !y)}
4 myxor4 <- function(x,y) {!((!x & !y) | (x & y))}
5 cbind(
6   ifelse(xor(test[,1], test[,2]), "true", "false"),
7   ifelse(myxor1(test[,1], test[,2]), "true", "false"),
8   ifelse(myxor2(test[,1], test[,2]), "true", "false"),
9   ifelse(myxor3(test[,1], test[,2]), "true", "false"),
10  ifelse(myxor4(test[,1], test[,2]), "true", "false")
11 )
```

|      | [,1]    | [,2]    | [,3]    | [,4]    | [,5]    |
|------|---------|---------|---------|---------|---------|
| [1,] | "false" | "false" | "false" | "false" | "false" |
| [2,] | "true"  | "true"  | "true"  | "true"  | "true"  |
| [3,] | "true"  | "true"  | "true"  | "true"  | "true"  |
| [4,] | "false" | "false" | "false" | "false" | "false" |

# \* Or in python

```
1 # Loading the libraries
2 import numpy as np
3 import pandas as pa
4
5 # Defining the data
6 x = [True, True, False, False]
7 y = [False, True, True, False]
8 ans = {
9     'x'    : x,
10    'y'    : y,
11    'and'   : np.logical_and(x, y),
12    'or'    : np.logical_or(x, y),
13    'xor'   : np.logical_xor(x, y)
14 }
15 pa.DataFrame(ans)
```

## \* Or in python (bis)

```
1 def myxor(x,y):  
2     return np.logical_or(  
3         np.logical_and(x, np.logical_not(y)),  
4         np.logical_and(np.logical_not(x), y)  
5     )  
6  
7 ans['myxor'] = myxor(x,y)  
8 pa.DataFrame(ans)
```

We will now see applications using the [met](#) dataset

# Filtering (subsetting) the data

Say we need to select records according to some criteria. For example:

- First day of the month, and
- Above latitude 40, and
- Elevation outside the range 500 and 1,000.

The logical expressions would be

- `(day == 1)`
- `(lat > 40)`
- `((elev < 500) | (elev > 1000))`

Respectively.

# \* data.table

In R with `data.table`:

```
1 dat[(day == 1) & (lat > 40) & ((elev < 500) | (elev > 1000))] |>  
2   nrow()
```



## In base R:

```
1 dat[dat$day == 1 &
2     dat$lat > 40 &
3     (dat$elev < 500) | (dat$elev > 1000), ]
```

|      | USAFID | WBAN  | year | month | day | hour | min | lat    | lon     | elev | wind.sp | temp |
|------|--------|-------|------|-------|-----|------|-----|--------|---------|------|---------|------|
| 3009 | 720113 | 54829 | 2019 | 8     | 1   | 0    | 15  | 42.543 | -83.178 | 222  | 1.5     | 25.0 |
| 3010 | 720113 | 54829 | 2019 | 8     | 1   | 0    | 39  | 42.543 | -83.178 | 222  | 2.6     | 24.4 |
| 3011 | 720113 | 54829 | 2019 | 8     | 1   | 0    | 57  | 42.543 | -83.178 | 222  | 4.1     | 24.0 |
| 3012 | 720113 | 54829 | 2019 | 8     | 1   | 1    | 15  | 42.543 | -83.178 | 222  | 4.1     | 23.5 |
| 3013 | 720113 | 54829 | 2019 | 8     | 1   | 1    | 16  | 42.543 | -83.178 | 222  | 3.6     | 23.5 |
| 3014 | 720113 | 54829 | 2019 | 8     | 1   | 1    | 35  | 42.543 | -83.178 | 222  | 3.1     | 23.0 |
| 3015 | 720113 | 54829 | 2019 | 8     | 1   | 1    | 36  | 42.543 | -83.178 | 222  | 3.6     | 23.0 |
| 3016 | 720113 | 54829 | 2019 | 8     | 1   | 1    | 55  | 42.543 | -83.178 | 222  | 5.1     | 22.4 |
| 3017 | 720113 | 54829 | 2019 | 8     | 1   | 2    | 18  | 42.543 | -83.178 | 222  | 4.1     | 21.9 |
| 3018 | 720113 | 54829 | 2019 | 8     | 1   | 2    | 36  | 42.543 | -83.178 | 222  | 4.6     | 21.4 |
| 3019 | 720113 | 54829 | 2019 | 8     | 1   | 2    | 55  | 42.543 | -83.178 | 222  | 5.1     | 21.0 |
| 3020 | 720113 | 54829 | 2019 | 8     | 1   | 2    | 56  | 42.543 | -83.178 | 222  | 6.2     | 21.0 |
| 3021 | 720113 | 54829 | 2019 | 8     | 1   | 3    | 15  | 42.543 | -83.178 | 222  | 3.1     | 20.5 |

## In R with **dplyr::filter()**:

```
1 dat |>  
2   filter(day == 1, lat > 40, (elev < 500) | (elev > 1000)) |>  
3   collect() |> # Notice this line!  
4   nrow()
```

```
[1] 27623
```

# In Python

```
1 import datatable as dt
2 dat = dt.fread("met.gz")

1 dat[(dt.f.day == 1) & (dt.f.lat > 40) & ((dt.f.elev < 500) | (dt.f.elev > 1000))]
2 # dat[dt.f.day == 1,:][dt.f.lat > 40,:][(dt.f.elev < 500) | (dt.f.elev > 1000)]
```

In the case of pydatatable we use `dt.f.` to refer to a column. `df.` is what we use to refer to datatable's `namespace`.

The `f.` is a `symbol` that allows accessing column names in a datatable's `Frame`.

# Questions 2

1. How many records have a temperature within 18 and 25?
2. Some records have missing values. Count how many records have `temp` as `NA`.
3. Following the previous question, plot a sample of 1,000 pairs of (`lat`, `lon`) coordinates for (a) the stations with `temp` as `NA` and (b) those with data.

# Solutions

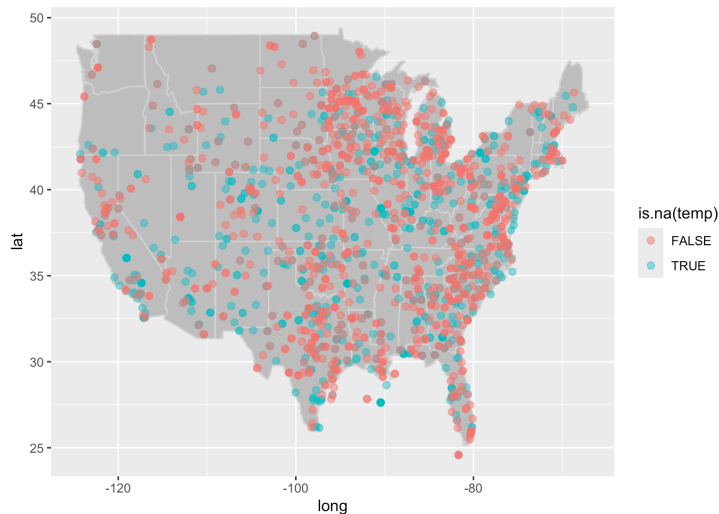
```
1 # Question 1
2 nrow(dat[dat$temp < 25 & dat$temp > 18, ])
3 # dat[temp %between% c(18, 25), .N]
4 # dat |> filter(between(temp, 18, 25)) |> collect() |> nrow()
5
6 # Question 2
7 nrow(dat[is.na(dat$temp), ])
8 # more succinct: sum(is.na(dat$temp))
```

```
[1] 968136
```

```
[1] 60089
```

# Solutions (cont.)

```
1 # Question 3
2 # Drawing a sample
3 set.seed(123)
4 idx1 <- sample(which(is.na(dat$temp)), 1000)
5 idx2 <- sample(which(!is.na(dat$temp)), 1000)
6
7 # Visualizing the data
8 # make a map of the US, as we did last class
9 ggplot(map_data("state"), aes(x = long, y = lat)) +
10   geom_map(aes(map_id = region), map = map_data("state"), col = "lightgrey", fill = "gray") +
11   geom_jitter(
12     data      = dat[c(idx1, idx2), ],
13     mapping   = aes(x = lon, y = lat, col = is.na(temp)),
14     inherit.aes = FALSE, alpha = .5, cex = 2
15   )
```



# Creating variables: Data types

- **logical**: Bool true/false type, e.g. dead/alive, sick/healthy, good/bad, yes/no, etc.
- **strings**: string of characters (letters/symbols), e.g. names, text, etc.
- **integer**: Numeric variable with no decimal (discrete), e.g. age, days, counts, etc.
- **double**: Numeric variable with decimals (continuous), e.g. distance, expression level, time.

In C (and other languages), strings, integers, and doubles may be specified with size, e.g. in [python](#) integers can be of 9, 16, and 32 bits. This is relevant when managing large datasets, where saving space can be fundamental ([more info](#)).

# Creating variables: Special data types

Most programming languages have special types which are built using basic types. A few examples:

- **time**: Could be date, date + time, or a combination of both. Usually it has a reference number defined as date 0. In R, the `Date` class has as reference 1970-01-01, in other words, "days since January 1st, 1970".
- **categorical**: Commonly used to represent strata/levels of variables, e.g. a variable "country" could be represented as a factor, where the data is stored as numbers but has a label.
- **ordinal**: Similar to factor, but it has ordering, e.g. "satisfaction level: 5 very satisfied, ..., 1 very unsatisfied".

Other special data types could be ways to represent missings (usually described as `na` or `NA`), or special numeric types, e.g. `+Inf` and Undefined (`NaN`).

When storing/sharing datasets, it is a good practice to do it along a dictionary describing each column data type/format.



# Questions 3: What's the best way to represent the following

- 0, 1, 1, 0, 0, 1
- Diabetes type 1, Diabetes type 2, Diabetes type 1, Diabetes type 2
- on, off, off, on, on, on
- 5, 10, 1, 15, 0, 0, 1
- 1.0, 2.0, 10.0, 6.0
- high, low, medium, medium, high
- -1, 1, -1, -1, 1,
- .2, 1.5, .8,  $\pi$
- $\pi$ , exp 1,  $\pi$ ,  $\pi$

# Variable creation

If we wanted to create two variables, `elev^2` and the scaled version of `wind.sp` by its standard error, we could do the following

```
1 dat$elev2 <- dat$elev^2
2 dat$windsp_scaled <- dat$wind.sp / sd(dat$wind.sp, na.rm = TRUE)
```

# Variable creation (cont. 1)

With the verb **dplyr::mutate()**:

|    | USAFID | WBAN  | year | month | day | hour | min | lat    | lon      | elev | wind.sp | temp |
|----|--------|-------|------|-------|-----|------|-----|--------|----------|------|---------|------|
| 1  | 690150 | 93121 | 2019 | 8     | 1   | 0    | 56  | 34.300 | -116.166 | 696  | 5.7     | 37.2 |
| 2  | 690150 | 93121 | 2019 | 8     | 1   | 1    | 56  | 34.300 | -116.166 | 696  | 8.2     | 35.6 |
| 3  | 690150 | 93121 | 2019 | 8     | 1   | 2    | 56  | 34.300 | -116.166 | 696  | 6.7     | 34.4 |
| 4  | 690150 | 93121 | 2019 | 8     | 1   | 3    | 56  | 34.300 | -116.166 | 696  | 5.1     | 33.3 |
| 5  | 690150 | 93121 | 2019 | 8     | 1   | 4    | 56  | 34.300 | -116.166 | 696  | 2.1     | 32.8 |
| 6  | 690150 | 93121 | 2019 | 8     | 1   | 5    | 56  | 34.300 | -116.166 | 696  | 0.0     | 31.1 |
| 7  | 690150 | 93121 | 2019 | 8     | 1   | 6    | 56  | 34.300 | -116.166 | 696  | 1.5     | 29.4 |
| 8  | 690150 | 93121 | 2019 | 8     | 1   | 7    | 56  | 34.300 | -116.166 | 696  | 2.1     | 28.9 |
| 9  | 690150 | 93121 | 2019 | 8     | 1   | 8    | 56  | 34.300 | -116.166 | 696  | 2.6     | 27.2 |
| 10 | 690150 | 93121 | 2019 | 8     | 1   | 9    | 56  | 34.300 | -116.166 | 696  | 1.5     | 26.7 |
| 11 | 690150 | 93121 | 2019 | 8     | 1   | 10   | 56  | 34.296 | -116.162 | 625  | 1.5     | 26.7 |
| 12 | 690150 | 93121 | 2019 | 8     | 1   | 11   | 56  | 34.296 | -116.162 | 625  | 2.6     | 25.6 |
| 13 | 690150 | 93121 | 2019 | 8     | 1   | 12   | 56  | 34.300 | -116.166 | 696  | 0.0     | 25.6 |

# Variable creation (cont. 2)

Imagine that we needed to scale multiple variables by their SD and didn't want to copy-paste this code several times. Here's how we could do it automatically for a given list of variable names:

```
1 # Listing the names
2 names <- c("wind.sp", "temp", "atm.press")
3
4 for(var in names){
5   dat[,paste0(var, '_scaled')] <- dat[,var] / sd(dat[,var], na.rm = TRUE)
6 }
```

Why can't we use `dat$var` inside the loop?

## \* Or with data.table

```
1 in_names <- c("wind.sp", "temp", "atm.press")
2 out_names <- paste0(in_names, "_scaled")
3 dat[,
4     c(out_names) := lapply(.SD, function(x) x/sd(x, na.rm = TRUE)),
5     .SDcols = in_names
6 ]
7
8 # Looking at the first 6
9 head(dat[, .SD, .SDcols = out_names], n = 4)
```

Key things to notice here: **c(out\_names)**, **.SD**, and **.SDCols**.

# Variable creation (cont. 3)

In the case of dplyr, we could use the following

```
1 names <- c("wind.sp", "temp", "atm.press")
2 dat |>
3   mutate(
4     across(
5       all_of(names),
6       function(x) x/sd(x, na.rm = TRUE),
7       .names = "{col}_scaled2"
8     )
9   ) |>
10  # Just to print the last columns
11  select(ends_with("_scaled2")) |>
12  head(n = 4)
```

|   | wind.sp_scaled2 | temp_scaled2 | atm.press_scaled2 |
|---|-----------------|--------------|-------------------|
| 1 | 2.654379        | 6.139348     | 248.7889          |
| 2 | 3.818580        | 5.875290     | 248.8874          |
| 3 | 3.120059        | 5.677247     | 248.9613          |
| 4 | 2.374970        | 5.495707     | 249.2077          |

# Complex variable creation

Don't forget about loops! `for` loops and `sapply` may be slow on a dataset of this size, but they can be quite handy for creating variables that rely on complicated relationships between variables. Consider this a "brute force" approach. Vectorized methods will *always* be faster, but these can be easier to conceptualize and, in rare cases, may be the only option.

Consider the problem creating a weird variable: `wind.temp`. This will take on 4 possible values, based on the temperature and wind speed: cool & still, cool & windy, warm & still, or warm & windy. We will split each variable based on their median value. Note that this code is too slow to actually run on this large dataset.

# Complex variable creation (cont 1)

Here's how we would do that with the `sapply` function (and a custom, unnamed function):

```
1 # create the new variable one entry at a time
2 wind.temp <- sapply(1:nrow(dat), function(i){
3   if(is.na(dat$temp[i]) | is.na(dat$wind.sp[i])){
4     return(NA)
5   }
6   if(dat$temp[i] <= median(dat$temp, na.rm=TRUE)){
7     if(dat$wind.sp[i] <= median(dat$wind.sp, na.rm=TRUE)){
8       return('cool & still')
9     }else{
10      return('cool & windy')
11    }
12  }else{
13    if(dat$wind.sp[i] <= median(dat$wind.sp, na.rm=TRUE)){
14      return('warm & still')
15    }else{
16      return('warm & windy')
17    }
18  }
19 })
```

Check: what would we need to change to add this variable to our dataset?



# Complex variable creation (cont 2)

Here's the code for doing that with a `for` loop:

```
1 # initialize a variable of all missing values
2 wind.temp <- rep(NA, nrow(dat))
3 # fill in the values one at a time
4 for(i in 1:nrow(dat)){
5   if(is.na(dat$temp[i]) | is.na(dat$wind.sp[i])){
6     return(NA)
7   }else{
8     if(dat$temp[i] <= median(dat$temp, na.rm=TRUE)){
9       if(dat$wind.sp[i] <= median(dat$wind.sp, na.rm=TRUE)){
10        wind.temp[i] <- 'cool & still'
11      }else{
12        wind.temp[i] <- 'cool & windy'
13      }
14    }else{
15      if(dat$wind.sp[i] <= median(dat$wind.sp, na.rm=TRUE)){
16        wind.temp[i] <- 'warm & still'
17      }else{
18        wind.temp[i] <- 'warm & windy'
19      }
20    }
21  }
22 }
```

Check: why do we need to include `na.rm=TRUE` when calculating the medians?

# Complex variable creation (cont 3)

Here's a simple vectorized approach that will actually run on a large dataset. This works for our current case, but it's still a brute force approach, because we had to specifically assign every possible value of our new variable. You can imagine that as the number of possible values increases, this code will get increasingly cumbersome.

```
1 # initialize a variable of all missing values
2 wind.temp <- rep(NA, nrow(dat))
3 # assign every possible value by subsetting
4 wind.temp[dat$temp <= median(dat$temp, na.rm=TRUE) &
5           dat$wind.sp <= median(dat$wind.sp, na.rm=TRUE)] <- 'cool & still'
6 wind.temp[dat$temp <= median(dat$temp, na.rm=TRUE) &
7           dat$wind.sp > median(dat$wind.sp, na.rm=TRUE)] <- 'cool & windy'
8 wind.temp[dat$temp > median(dat$temp, na.rm=TRUE) &
9           dat$wind.sp <= median(dat$wind.sp, na.rm=TRUE)] <- 'warm & still'
10 wind.temp[dat$temp > median(dat$temp, na.rm=TRUE) &
11           dat$wind.sp > median(dat$wind.sp, na.rm=TRUE)] <- 'warm & windy'
12
13 head(wind.temp)
```

```
[1] "warm & windy" "warm & windy" "warm & windy" "warm & windy" "warm & still"
```

```
[6] "warm & still"
```

# Merging data

- While building the MET dataset, we dropped the State data.
- We can use the original Stations dataset and *merge* it to the MET dataset.
- But we cannot do it right away. We need to process the data somewhat first.

# Merging data (cont. 1)

```
1 stations <- fread("https://noaa-isd-pds.s3.amazonaws.com/isd-history.csv")
2 stations <- as.data.frame(stations)
3 stations$USAF <- as.integer(stations$USAF)
4
5 # Dealing with NAs and 999999
6 stations$USAF[stations$USAF == 999999] <- NA
7 stations$CTRY[stations$CTRY == ""] <- NA
8 stations$STATE[stations$STATE == ""] <- NA
9
10 # Selecting the three relevant columns, and keep unique records
11 stations <- unique(stations[, c('USAF', 'CTRY', 'STATE')])
12
13 # Dropping NAs
14 stations <- stations[!is.na(stations$USAF), ]
15
16 head(stations, n = 4)
```

|   | USAF | CTRY | STATE |
|---|------|------|-------|
| 1 | 7018 | <NA> | <NA>  |
| 2 | 7026 | AF   | <NA>  |
| 3 | 7070 | AF   | <NA>  |
| 4 | 8260 | <NA> | <NA>  |

# Merging data (cont. 2)

```
1 merge(  
2   # Data  
3   x     = dat,  
4   y     = stations,  
5   # List of variables to match  
6   by.x  = "USAFID",  
7   by.y  = "USAF",  
8   # Which obs to keep?  
9   all.x = TRUE,  
10  all.y = FALSE  
11 ) |> nrow()
```

```
[1] 2385443
```

This is more rows! The original dataset, `dat`, has 2377343. This means that the `stations` dataset has duplicated IDs. We can fix this:

```
1 stations <- stations[!duplicated(stations$USAF), ]
```

# Merging data (cont. 3)

We now can use the function `merge()` to add the extra data

```
1 dat <- merge(  
2   # Data  
3   x     = dat,  
4   y     = stations,  
5   # List of variables to match  
6   by.x  = "USAFID",  
7   by.y  = "USAF",  
8   # Which obs to keep?  
9   all.x = TRUE,  
10  all.y = FALSE  
11 )  
12  
13 head(dat[, c('USAFID', 'WBAN', 'STATE')], n = 4)
```

|   | USAFID | WBAN  | STATE |
|---|--------|-------|-------|
| 1 | 690150 | 93121 | CA    |
| 2 | 690150 | 93121 | CA    |
| 3 | 690150 | 93121 | CA    |
| 4 | 690150 | 93121 | CA    |

What happens when you change the options `all.x` and `all.y`?

# \* Aggregating data: Adding grouped variables

- Many times we need to either impute some data, or generate variables by strata.
- If we, for example, wanted to impute missing temperature with the daily state average, we could use **by** together with the **data.table::fcoalesce()** function:

```
1 dat[, temp_imp := fcoalesce(temp, mean(temp, na.rm = TRUE)),  
2   by = .(STATE, year, month, day)]
```

# Aggregating data: Adding grouped variables

- In the case of dplyr, we can do the following using `dplyr::group_by()` together with `dplyr::coalesce()`:

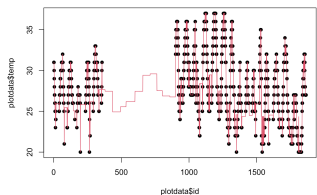
```
1 dat <- dat |>
2   group_by(STATE, year, month, day) |>
3   mutate(
4     temp_imp = coalesce(temp, mean(temp, na.rm = TRUE))
5   ) |> collect()
```



# Aggregating data: Adding grouped variables (cont.)

Let's see how it looks:

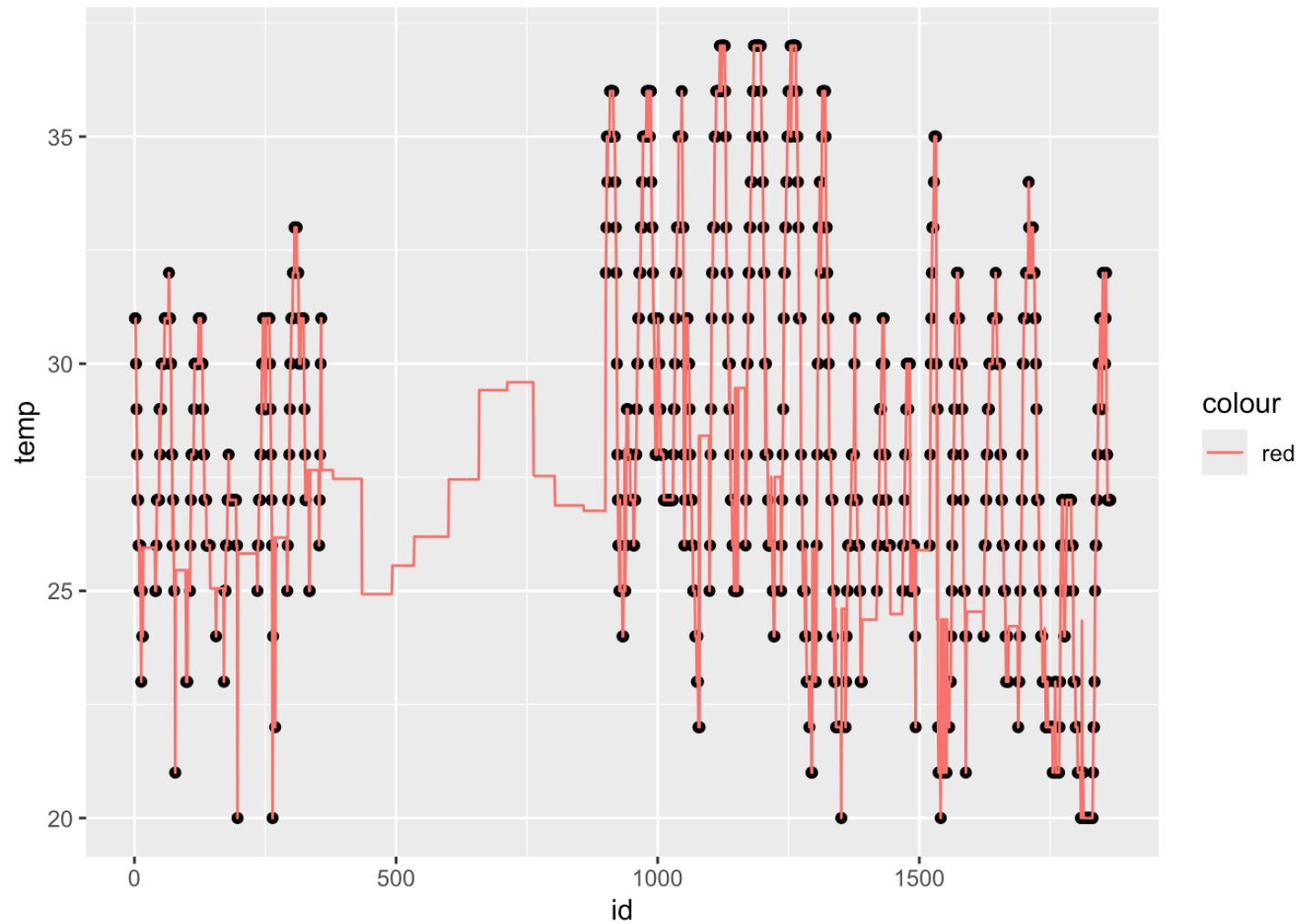
```
1 # Preparing for plotting
2 # select single station, sort by date
3 plotdata <- dat[dat$USAFID == 720172, ]
4 plotdata <- plotdata[order(plotdata$year, plotdata$month, plotdata$day, plo
5 # Generate an 'x' variable for time
6 plotdata$id <- 1:nrow(plotdata)
7
8
9 plot(plotdata$id, plotdata$temp, pch=16)
10 lines(plotdata$id, plotdata$temp_imp, col=2)
```



```

1 plotdata |>
2   ggplot(aes(x = id)) +
3   geom_point(aes(y = temp)) +
4   geom_line(aes(y = temp_imp, colour = 'red'))

```



# \* Aggregating data: Summary table

- Using `by` also allow us creating summaries of our data.
- For example, if we wanted to compute the average temperature, wind-speed, and atmospheric pressure by state, we could do the following

```
1 dat[, .(  
2   temp_avg      = mean(temp, na.rm=TRUE),  
3   wind.sp_avg   = mean(wind.sp, na.rm=TRUE),  
4   atm.press_avg = mean(atm.press, na.rm = TRUE)  
5   ),  
6   by = STATE  
7   ][order(STATE)] |> head(n = 4)
```

# \* Aggregating data: Summary table (cont. 1)

When dealing with too many variables, we can use the `.SD` special symbol in `data.table`:

```
1 # Listing the names
2 in_names  <- c("wind.sp", "temp", "atm.press")
3 out_names <- paste0(in_names, "_avg")
4
5 dat[,
6   setNames(lapply(.SD, mean, na.rm = TRUE), out_names),
7   .SDcols = in_names, keyby  = STATE
8 ] |> head(n = 4)
```

Notice the **keyby** option here: "Group by STATE and order by STATE".

# Aggregating data: Summary table (cont. 2)

- Using `dplyr` verbs

```
1 dat |>
2   group_by(STATE) |>
3   summarise(
4     temp_avg      = mean(temp, na.rm=TRUE),
5     wind.sp_avg   = mean(wind.sp, na.rm=TRUE),
6     atm.press_avg = mean(atm.press, na.rm = TRUE)
7   ) |>
8   arrange(STATE) |>
9   head(n = 4)
```

# A tibble: 4 × 4

|   | STATE | temp_avg | wind.sp_avg | atm.press_avg |
|---|-------|----------|-------------|---------------|
|   | <chr> | <dbl>    | <dbl>       | <dbl>         |
| 1 | AL    | 26.2     | 1.57        | 1016.         |
| 2 | AR    | 26.2     | 1.84        | 1015.         |
| 3 | AZ    | 28.8     | 2.98        | 1011.         |
| 4 | CA    | 22.4     | 2.61        | 1013.         |

Notice the `arrange()` function.

# \* Other data.table goodies

- `shift()` Fast lead/lag for vectors and lists.
- `ifelse()` Fast if-else, similar to base R's `ifelse()`.
- `fcoalesce()` Fast coalescing of missing values.
- `%between%` A short form of `(x < lb) & (x > up)`
- `%inrange%` A short form of `x %in% lb:up`
- `%chin%` Fast match of character vectors, equivalent to `x %in% X`, where both `x` and `X` are character vectors.
- `nafill()` Fill missing values using a constant, last observed value, or the next observed value.

# Benchmarks

- [H2O.ai's benchmark \(link\)](#): Designed by the lead developer of data.table [Matt Dowle](#)
- [RStudio's benchmark \(link\)](#): Designed as part of the benchmarks with the [vroom](#) package.