# Data Transformation

```
library(nycflights13)
```

```
## Warning: package 'nycflights13' was built under R version 3.5.3
```

```
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 3.5.3
```

```
## -- Attaching packages ------------------------------------------------
------------------------------------------------- tidyverse 1.3.0 --
```

```
## v ggplot2 3.3.2      v purrr   0.3.4
## v tibble  3.0.3      v dplyr   1.0.0
## v tidyr   1.1.0      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0
```

```
## Warning: package 'readr' was built under R version 3.5.3
```

```
## Warning: package 'purrr' was built under R version 3.5.3
```

```
## Warning: package 'stringr' was built under R version 3.5.3
```

```
## Warning: package 'forcats' was built under R version 3.5.3
```

```
## -- Conflicts ----------------------------------------------------------
----------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

The data set "flights" will be used for the following examples. To learn more about the data set, please use `?flights` . The data set is ni **tibble** format. The common variable types include -

- `int` stands for integers.
- `dbl` stands for doubles, or real numbers.
- `chr` stands for character vectors, or strings.
- `dttm` stands for date-times (a date + a time).
- `lgl` stands for logical, vectors that contain only TRUE or FALSE.

- `fctr` stands for factors, which R uses to represent categorical variables with fixed possible values.
- `date` stands for dates.

```
head(flights)
```

```
## # A tibble: 6 x 19
##    year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##   <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1  2013     1     1      517            515         2      830            819
## 2  2013     1     1      533            529         4      850            830
## 3  2013     1     1      542            540         2      923            850
## 4  2013     1     1      544            545        -1     1004           1022
## 5  2013     1     1      554            600        -6      812            837
## 6  2013     1     1      554            558        -4      740            728
## # ... with 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>
```

**major dplyr functions**

- Pick observations by their values `filter()`.
- Reorder the rows `arrange()`.
- Pick variables by their names `select()`.
- Create new variables with functions of existing variables `mutate()`.
- Collapse many values down to a single summary `summarise()`.
- change the scope `group_by()`

# filter()

print out the result and also assign the result to a new variable dec25.

```
(dec25 <- filter(flights, month == 12, day == 25))
```

```
## # A tibble: 719 x 19
##      year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##     <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013    12    25      456            500        -4      649            651
## 2   2013    12    25      524            515         9      805            814
## 3   2013    12    25      542            540         2      832            850
## 4   2013    12    25      546            550        -4     1022           1027
## 5   2013    12    25      556            600        -4      730            745
## 6   2013    12    25      557            600        -3      743            752
## 7   2013    12    25      557            600        -3      818            831
## 8   2013    12    25      559            600        -1      855            856
## 9   2013    12    25      559            600        -1      849            855
## 10  2013    12    25      600            600         0      850            846
## # ... with 709 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```
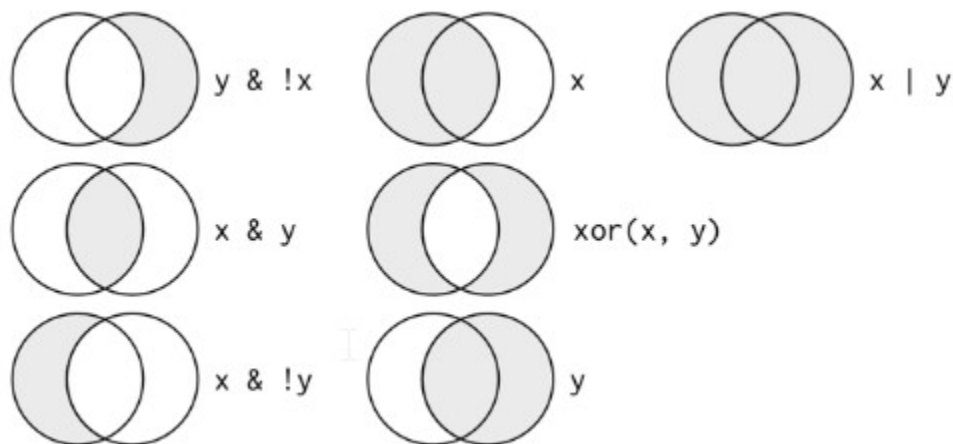
**Notes** `==` brings problem when dealing with floating point numbers. Please use `near(` instead

```
1/49 * 49 == 1
```

```
## [1] FALSE
```

```
near(1/49 * 49, 1)
```

```
## [1] TRUE
```



R Markdown image

- use `|` and `&`, don't use `||` and `&&`.
- use `%in%` instead of `|` to simplify
- `!(x & y)` is the same as `!x | !y`
- `!(x | y)` is the same as `!x & !y`

```
nov_dec <- filter(flights, month %in% c(11, 12))

# filter(flights, month == 11 | month == 12)
```

**missing values** is contagious.

```
NA > 5
```

```
## [1] NA
```

```
#> [1] NA
10 == NA
```

```
## [1] NA
```

```
#> [1] NA
NA + 10
```

```
## [1] NA
```

```
#> [1] NA
NA / 2
```

```
## [1] NA
```

```
#> [1] NA
NA == NA
```

```
## [1] NA
```

```
#> [1] NA
```

to check if the value is NA, `is.na()`

`filter()` only includes rows where the condition is TRUE; it excludes both FALSE and NA values. If you want to preserve missing values, ask for them explicitly:

```
df <- tibble(x = c(1, NA, 3))
filter(df, x > 1)
```

```
## # A tibble: 1 x 1
##       x
##   <dbl>
## 1     3
```

```
filter(df, is.na(x) | x > 1)
```

```
## # A tibble: 2 x 1
##       x
##   <dbl>
## 1    NA
## 2     3
```

# arrange()

When sorting columns, NA are sorted at the end.

```
arrange(flights, year, month, day)
```

```
## # A tibble: 336,776 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
##  1  2013     1     1      517            515         2      830            819
##  2  2013     1     1      533            529         4      850            830
##  3  2013     1     1      542            540         2      923            850
##  4  2013     1     1      544            545        -1     1004           1022
##  5  2013     1     1      554            600        -6      812            837
##  6  2013     1     1      554            558        -4      740            728
##  7  2013     1     1      555            600        -5      913            854
##  8  2013     1     1      557            600        -3      709            723
##  9  2013     1     1      557            600        -3      838            846
## 10  2013     1     1      558            600        -2      753            745
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

```
arrange(flights, desc(dep_delay))
```

```
## # A tibble: 336,776 x 19
##      year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##     <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013     1     9      641            900      1301     1242           1530
## 2   2013     6    15     1432           1935      1137     1607           2120
## 3   2013     1    10     1121           1635      1126     1239           1810
## 4   2013     9    20     1139           1845      1014     1457           2210
## 5   2013     7    22      845           1600      1005     1044           1815
## 6   2013     4    10     1100           1900       960     1342           2211
## 7   2013     3    17     2321            810       911      135           1020
## 8   2013     6    27      959           1900       899     1236           2226
## 9   2013     7    22     2257            759       898      121           1026
## 10  2013    12     5      756           1700       896     1058           2020
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

# select()

```
select(flights, year, month, day)
```

```
## # A tibble: 336,776 x 3
##      year month   day
##     <int> <int> <int>
## 1   2013     1     1
## 2   2013     1     1
## 3   2013     1     1
## 4   2013     1     1
## 5   2013     1     1
## 6   2013     1     1
## 7   2013     1     1
## 8   2013     1     1
## 9   2013     1     1
## 10  2013     1     1
## # ... with 336,766 more rows
```

```
# select(flights, year:day)

select(flights, -(year:day))
```

```
## # A tibble: 336,776 x 16
##    dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier
##       <int>          <int>     <dbl>    <int>          <int>     <dbl> <chr>
##  1      517            515         2      830            819        11 UA
##  2      533            529         4      850            830        20 UA
##  3      542            540         2      923            850        33 AA
##  4      544            545        -1     1004           1022       -18 B6
##  5      554            600        -6      812            837       -25 DL
##  6      554            558        -4      740            728        12 UA
##  7      555            600        -5      913            854        19 B6
##  8      557            600        -3      709            723       -14 EV
##  9      557            600        -3      838            846        -8 B6
## 10      558            600        -2      753            745         8 AA
## # ... with 336,766 more rows, and 9 more variables: flight <int>,
## #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>,
## #   hour <dbl>, minute <dbl>, time_hour <dttm>
```

In addition, you can use **selection helpers**. Some helpers select specific columns:

- `starts_with()` : Starts with a prefix.

- `ends_with()` : Ends with a suffix.

- `contains()` : Contains a literal string.

- `matches()` : Matches a regular expression.

- `num_range()` : Matches a numerical range like x01, x02, x03.

```
select(flights, matches('time'))
```

```
## # A tibble: 336,776 x 6
##    dep_time sched_dep_time arr_time sched_arr_time air_time time_hour
##       <int>          <int>    <int>          <int>    <dbl> <dttm>
##  1      517            515      830            819      227 2013-01-01 05:00:00
##  2      533            529      850            830      227 2013-01-01 05:00:00
##  3      542            540      923            850      160 2013-01-01 05:00:00
##  4      544            545     1004           1022      183 2013-01-01 05:00:00
##  5      554            600      812            837      116 2013-01-01 06:00:00
##  6      554            558      740            728      150 2013-01-01 05:00:00
##  7      555            600      913            854      158 2013-01-01 06:00:00
##  8      557            600      709            723       53 2013-01-01 06:00:00
##  9      557            600      838            846      140 2013-01-01 06:00:00
## 10      558            600      753            745      138 2013-01-01 06:00:00
## # ... with 336,766 more rows
```

**rename* columns

```
rename(flights, tail_num = tailnum)
```

```
## # A tibble: 336,776 x 19
##     year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
##    <int> <int> <int>    <int>          <int>     <dbl>    <int>          <int>
## 1   2013     1     1      517            515         2      830            819
## 2   2013     1     1      533            529         4      850            830
## 3   2013     1     1      542            540         2      923            850
## 4   2013     1     1      544            545        -1     1004           1022
## 5   2013     1     1      554            600        -6      812            837
## 6   2013     1     1      554            558        -4      740            728
## 7   2013     1     1      555            600        -5      913            854
## 8   2013     1     1      557            600        -3      709            723
## 9   2013     1     1      557            600        -3      838            846
## 10  2013     1     1      558            600        -2      753            745
## # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
## #   carrier <chr>, flight <int>, tail_num <chr>, origin <chr>, dest <chr>,
## #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

**everything()** move columns to the beginning

```
select(flights, time_hour, air_time, everything())
```

```
## # A tibble: 336,776 x 19
##    time_hour           air_time  year month   day dep_time sched_dep_time
##    <dttm>                 <dbl> <int> <int> <int>    <int>          <int>
## 1  2013-01-01 05:00:00      227  2013     1     1      517            515
## 2  2013-01-01 05:00:00      227  2013     1     1      533            529
## 3  2013-01-01 05:00:00      160  2013     1     1      542            540
## 4  2013-01-01 05:00:00      183  2013     1     1      544            545
## 5  2013-01-01 06:00:00      116  2013     1     1      554            600
## 6  2013-01-01 05:00:00      150  2013     1     1      554            558
## 7  2013-01-01 06:00:00      158  2013     1     1      555            600
## 8  2013-01-01 06:00:00       53  2013     1     1      557            600
## 9  2013-01-01 06:00:00      140  2013     1     1      557            600
## 10 2013-01-01 06:00:00      138  2013     1     1      558            600
## # ... with 336,766 more rows, and 12 more variables: dep_delay <dbl>,
## #   arr_time <int>, sched_arr_time <int>, arr_delay <dbl>, carrier <chr>,
## #   flight <int>, tailnum <chr>, origin <chr>, dest <chr>, distance <dbl>,
## #   hour <dbl>, minute <dbl>
```

# mutate()

```
flights_sml <- select(flights,
  year:day,
  ends_with("delay"),
  distance,
  air_time
)
mutate(flights_sml,
  gain = dep_delay - arr_delay,
  speed = distance / air_time * 60
)
```

```
## # A tibble: 336,776 x 9
##     year month   day dep_delay arr_delay distance air_time  gain speed
##    <int> <int> <int>     <dbl>     <dbl>    <dbl>    <dbl> <dbl> <dbl>
##  1  2013     1     1         2        11     1400      227    -9  370.
##  2  2013     1     1         4        20     1416      227   -16  374.
##  3  2013     1     1         2        33     1089      160   -31  408.
##  4  2013     1     1        -1       -18     1576      183    17  517.
##  5  2013     1     1        -6       -25      762      116    19  394.
##  6  2013     1     1        -4        12      719      150   -16  288.
##  7  2013     1     1        -5        19     1065      158   -24  404.
##  8  2013     1     1        -3       -14      229       53    11  259.
##  9  2013     1     1        -3        -8      944      140     5  405.
## 10  2013     1     1        -2         8      733      138   -10  319.
## # ... with 336,766 more rows
```

```
# to keep the new variables only
transmute(flights,
  gain = dep_delay - arr_delay,
  hours = air_time / 60,
  gain_per_hour = gain / hours
)
```

```
## # A tibble: 336,776 x 3
##      gain hours gain_per_hour
##     <dbl> <dbl>         <dbl>
## 1      -9 3.78          -2.38
## 2     -16 3.78          -4.23
## 3     -31 2.67         -11.6
## 4      17 3.05           5.57
## 5      19 1.93           9.83
## 6     -16 2.5           -6.4
## 7     -24 2.63          -9.11
## 8      11 0.883         12.5
## 9       5 2.33           2.14
## 10    -10 2.3           -4.35
## # ... with 336,766 more rows
```

available operators include

- `+, -, ^`

- `sum(), mean()`

- `%/%` (integer division), `%%` (remainder)

- `log(), log2(), log10()`

```
transmute(flights,
  dep_time,
  hour = dep_time %/% 100,
  minute = dep_time %% 100
)
```

```
## # A tibble: 336,776 x 3
##    dep_time  hour minute
##       <int> <dbl>  <dbl>
## 1       517     5     17
## 2       533     5     33
## 3       542     5     42
## 4       544     5     44
## 5       554     5     54
## 6       554     5     54
## 7       555     5     55
## 8       557     5     57
## 9       557     5     57
## 10      558     5     58
## # ... with 336,766 more rows
```

- `lead(), lag()`

```
(x <- 1:10)
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10
```

```
#> [1]  1  2  3  4  5  6  7  8  9 10
lag(x)
```

```
##  [1] NA  1  2  3  4  5  6  7  8  9
```

```
#> [1] NA  1  2  3  4  5  6  7  8  9
lead(x)
```

```
##  [1]  2  3  4  5  6  7  8  9 10 NA
```

```
#> [1]  2  3  4  5  6  7  8  9 10 NA
```

- `cumsum()`, `cumprod()`, `cummin()`, `cummax()`; and dplyr provides `cummean()` for cumulative means.

```
x
```

```
##  [1]  1  2  3  4  5  6  7  8  9 10
```

```
#> [1]  1  2  3  4  5  6  7  8  9 10
cumsum(x)
```

```
##  [1]  1  3  6 10 15 21 28 36 45 55
```

```
#> [1]  1  3  6 10 15 21 28 36 45 55
cummean(x)
```

```
##  [1] 1.000000 1.000000 1.333333 1.750000 2.200000 2.666667 3.142857 3.625000
## [9] 4.111111 4.600000
```

- `min_rank()`, `row_number()`, `dense_rank()`, `percent_rank()`, `cume_dist()`, `ntile()`

```
y <- c(1, 2, 2, NA, 3, 4)
min_rank(y)
```

```
## [1]  1  2  2 NA  4  5
```

```
#> [1]  1  2  2 NA  4  5
min_rank(desc(y))
```

```
## [1]  5  3  3 NA  2  1
```

```
#> [1]  5  3  3 NA  2  1
row_number(y)
```

```
## [1]  1  2  3 NA  4  5
```

```
#> [1]  1  2  3 NA  4  5
dense_rank(y)
```

```
## [1]  1  2  2 NA  3  4
```

```
#> [1]  1  2  2 NA  3  4
percent_rank(y)
```

```
## [1] 0.00 0.25 0.25   NA 0.75 1.00
```

```
#> [1] 0.00 0.25 0.25   NA 0.75 1.00
cume_dist(y)
```

```
## [1] 0.2 0.6 0.6  NA 0.8 1.0
```

```
#> [1] 0.2 0.6 0.6  NA 0.8 1.0
```

## summarise()

```
summarise(flights, delay = mean(dep_delay, na.rm = TRUE))
```

```
## # A tibble: 1 x 1
##   delay
##   <dbl>
## 1  12.6
```

```
by_day <- group_by(flights, year, month, day)
summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))
```

```
## `summarise()` regrouping output by 'year', 'month' (override with `.groups` argumen
t)
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##     year month   day delay
##    <int> <int> <int> <dbl>
##  1  2013     1     1 11.5
##  2  2013     1     2 13.9
##  3  2013     1     3 11.0
##  4  2013     1     4  8.95
##  5  2013     1     5  5.73
##  6  2013     1     6  7.15
##  7  2013     1     7  5.42
##  8  2013     1     8  2.55
##  9  2013     1     9  2.28
## 10  2013     1    10  2.84
## # ... with 355 more rows
```
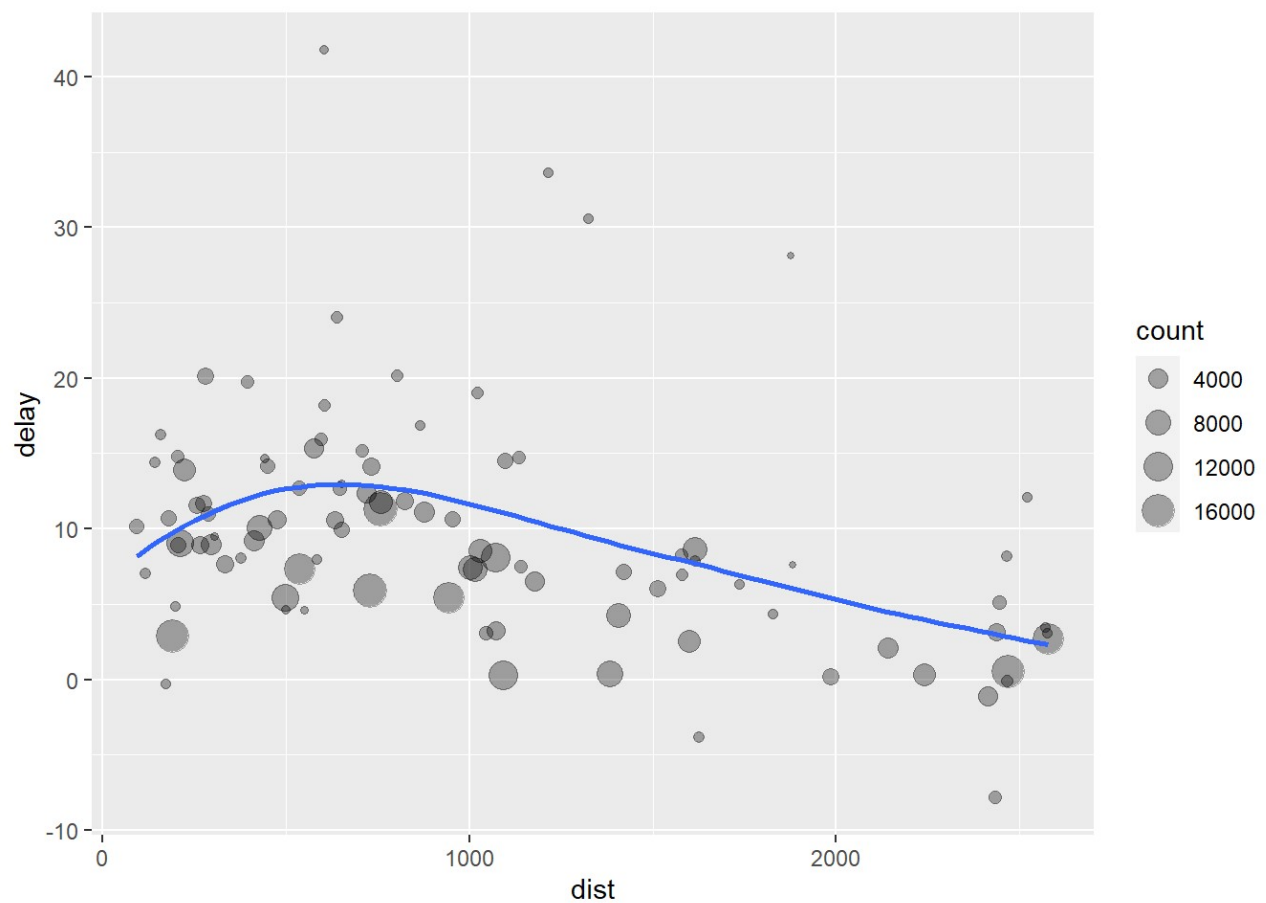
**pipe** %>%

```
delays <- flights %>%
  group_by(dest) %>%
  summarise(
    count = n(),
    dist = mean(distance, na.rm = TRUE),
    delay = mean(arr_delay, na.rm = TRUE)
  ) %>%
  filter(count > 20, dest != "HNL")
```

```
## `summarise()` ungrouping output (override with `.groups` argument)
```

```
ggplot(data = delays, mapping = aes(x = dist, y = delay)) +
  geom_point(aes(size = count), alpha = 1/3) +
  geom_smooth(se = FALSE)
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

```
flights %>%
  group_by(year, month, day) %>%
  summarise(mean = mean(dep_delay))
```

```
## `summarise()` regrouping output by 'year', 'month' (override with `.groups` argumen
t)
```

```
## # A tibble: 365 x 4
## # Groups:   year, month [12]
##     year month   day  mean
##    <int> <int> <int> <dbl>
##  1  2013     1     1    NA
##  2  2013     1     2    NA
##  3  2013     1     3    NA
##  4  2013     1     4    NA
##  5  2013     1     5    NA
##  6  2013     1     6    NA
##  7  2013     1     7    NA
##  8  2013     1     8    NA
##  9  2013     1     9    NA
## 10  2013     1    10    NA
## # ... with 355 more rows
```

without `na.rm`, we will get a lot of missing values. (If one of the input is NA, the calculation result will be NA.)

**Useful summary functions** * `mean(x)`, `median(x)`,

- `sd(x)` standard deviation, `IQR(x)` interquartile range, `mad(x)` median absolute deviation,

- `min(x)`, `quantile(x, 0.25)`, `max(x)`

- `first(x)`, `nth(x,2)`, `last(x)`

- `n()` count, `n_distinct(x)`, `sum(!is.na(x))`

- `sum(x < 10)`, `mean(y==0)`, converts TRUE to 1 and FALSE to 0,

remember to **ungrouping** `ungroup()`