# APPLICATIONS



OF DATA SCIENCE

## Tidy Data Wrangling - Part A

**Applications of Data Science - Class 2** 

Giora Simchoni

gsimchoni@gmail.com and add #dsapps in subject

Stat. and OR Department, TAU

2021-01-17



# dplyr: Basic Data Verbs



## **Basic Data Verbs**

- filter() rows based on one or more conditions
- mutate() one or more columns, usually based on existing columns
- select() the column(s) you want
- arrange () rows by one or more columns order
- summarize() or summarise() that single quantity off a column
- pull () a column as a vector, don't want it as a column no more

And the much beloved group\_by(): do whatever by groups of one or more variables.



## Read in the data

```
library(tidyverse)
okcupid <- read csv("~/okcupid.csv.zip")</pre>
```

#### Reminder:

```
dim(okcupid)
## [1] 59946
                 31
colnames(okcupid)
```

```
##
    [1] "age"
                       "body type"
                                      "diet"
                                                     "drinks"
                                                                    "drugs"
    [6] "education"
                       "essav0"
                                      "essav1"
                                                     "essav2"
                                                                    "essav3"
##
   [11] "essay4"
                       "essav5"
                                      "essav6"
                                                     "essav7"
                                                                    "essav8"
                                      "height"
                                                     "income"
                                                                    "dob"
## [16] "essay9"
                       "ethnicity"
## [21] "last online" "location"
                                      "offspring"
                                                     "orientation" "pets"
   [26] "religion"
                                      "sian"
                                                     "smokes"
                                                                    "speaks"
                       "sex"
## [31] "status"
```



## mutate()

Add a column height cm, the height in centimeters:

```
okcupid <- okcupid %>%
  mutate(height_cm = 2.54 * height)
```

♀ if you also load the magrittr package you could do:

okcupid %<>% mutate(height\_cm = 2.54 \* height)



## filter() and select()

Filter only women, select only age and height:

```
okcupid %>%
  filter(sex == "f") %>%
  select(age, height)
## # A tibble: 24,117 x 2
##
      age height
##
     <dbl> <dbl>
##
       32
             65
##
     31 65
##
   3 24 67
##
   4 30 66
##
   5 29
             62
##
   6 39
             65
##
  7 26 64
   8 27 67
##
  9 22
             67
## 10
     27
             64
## # ... with 24,107 more rows
```





### Same but income over 100K, and select all essay questions:

```
okcupid %>%
  filter(sex == "f", income > 100000) %>%
  select(starts with("essay"))
## # A tibble: 208 x 10
##
     essay0 essay1 essay2 essay3 essay4 essay5 essay6 essay7 essay
     <chr> <chr> <chr> <chr> <chr> <chr> <chr>
##
                                                                    <chr>
   1 "i lov~ "being~ "scraw~ "my bi~ "music~ "veget~ "makin~ "kicki~ "wow,
##
   2 "i'm s~ "curre~ "eatin~ "my po~ "pride~ "nothi~ "my ne~ "eatin~ "i'm
##
   3 "welco~ "piano~ "singi~ "my he~ "books~ "touch~ "diffe~ <NA>
##
                                                                     <NA>
   4 "pureb~ "by da~ "being~ "my ha~ "to st~ "- wat~ "my ne~ "i try~ "ummm
##
   5 "i was~ "chick~ "using~ "lips ~ "armag~ "lust,~ "enter~ "makin~ <NA>
##
##
   6 "hello~ "i tal~ "anyth~ "my as~ "book:~ "my qu~ "every~ "i wor~ <NA>
   7 "life'~ "i'm j~ "getti~ "its b~ "otis ~ "1. so~ "the w~ "oh ma~ "i do
##
##
   8 "every~ "livin~ "being~ "my ey~ "dubst~ "dirty~ "how t~ "recov~ "i lo
   9 "love ~ "daily~ "i am ~ "my sm~ "love ~ "masca~ "if i ~ <NA>
##
                                                                    "i am
## 10 "<b>ph~ "i am ~ "pissi~ "my sm~ "book:~ "my do~ "who p~ "total~ "my d
```



8/36

## # ... with 198 more rows

### Same but using a range of columns:

```
okcupid %>%
  filter(sex == "f", income > 100000) %>%
  select(essay0:essay9)
## # A tibble: 208 x 10
##
     essay0 essay1 essay2 essay3 essay4 essay5 essay6 essay7 essay
     <chr> <chr> <chr> <chr> <chr> <chr>
##
                                                                    <chr>
   1 "i lov~ "being~ "scraw~ "my bi~ "music~ "veget~ "makin~ "kicki~ "wow,
##
   2 "i'm s~ "curre~ "eatin~ "my po~ "pride~ "nothi~ "my ne~ "eatin~ "i'm
##
##
   3 "welco~ "piano~ "singi~ "my he~ "books~ "touch~ "diffe~ <NA>
                                                                     <NA>
   4 "pureb~ "by da~ "being~ "my ha~ "to st~ "- wat~ "my ne~ "i try~ "ummm
##
   5 "i was~ "chick~ "using~ "lips ~ "armag~ "lust,~ "enter~ "makin~ <NA>
##
##
   6 "hello~ "i tal~ "anyth~ "my as~ "book:~ "my qu~ "every~ "i wor~ <NA>
  7 "life'~ "i'm j~ "getti~ "its b~ "otis ~ "1. so~ "the w~ "oh ma~ "i do
##
##
   8 "every~ "livin~ "being~ "my ey~ "dubst~ "dirty~ "how t~ "recov~ "i lo
   9 "love ~ "daily~ "i am ~ "my sm~ "love ~ "masca~ "if i ~ <NA>
##
                                                                    "i am
## 10 "<b>ph~ "i am ~ "pissi~ "my sm~ "book:~ "my do~ "who p~ "total~ "my d
## # ... with 198 more rows
```

Many, many such gifts, see <a href="tidyselect">tidyselect</a>



## summarize()

### Find the average height of women

```
okcupid %>%
  filter(sex == "f") %>%
  summarize(avg_height = mean(height_cm, na.rm = TRUE))

## # A tibble: 1 x 1
## avg_height
## <dbl>
## 1 165.
```

Notice we got a tibble. We could either pull this single number:

```
okcupid %>%
  filter(sex == "f") %>%
  summarize(avg_height = mean(height_cm, na.rm = TRUE)) %>%
  pull()
```

## [1] 165.3638



### Or pull the vector of heights first, then calculate their mean:

```
okcupid %>%
  filter(sex == "f") %>%
  pull(height_cm) %>%
  mean(na.rm = TRUE)

## [1] 165.3638
```

### Amazingly, this would also work:

```
mean(pull(filter(okcupid, sex == "f"), height_cm), na.rm = TRUE)
## [1] 165.3638
```



## group\_by()

### But why settle for women only?

### And you might want to consider rename () ing sex!

```
okcupid %>%
  group_by(sex) %>%
  summarize(avg_height = mean(height_cm, na.rm = TRUE)) %>%
  rename(gender = sex)
```



Applications of Data Science 12 / 36

### Group by multiple variables, get more summaries, arrange by descending average height:

```
okcupid %>%
  group by(sex, status) %>%
  summarize(avg height = mean(height cm, na.rm = TRUE),
            med height = median(height cm, na.rm = TRUE),
            n = n()) %>%
  arrange (-med height)
```

```
## # A tibble: 10 \times 5
## # Groups: sex [2]
##
                     avg height med height
    sex status
##
  <chr> <chr>
                         <dbl> <dbl> <int>
                                  180. 1209
##
  1 m available
                          179.
## 2 m married
                          179.
                                   180. 175
##
                          179. 178. 1061
  3 m seeing someone
         single
##
                          179. 178. 33378
                                   177.
##
       unknown
                          177.
                                  166. 656
## 6 f available
                          166.
## 7 f married
                          166. 165. 135
## 8 f seeing someone 165. 165. 1003
## 9 f
         single
                         165. 165. 22319
## 10 f
         unknown
                          161.
                                   159.
```



## Protip: count()

When all you want is, well, count, no need to group by:

```
okcupid %>% count(body type, sort = TRUE)
## # A tibble: 13 \times 2
     body type
                        n
     <chr>
##
                  <int>
##
  1 average
                    14652
##
  2 fit
                    12711
##
                    11819
  3 athletic
## 4 <NA>
                    5296
## 5 thin
                    4711
## 6 curvy
                    3924
## 7 a little extra 2629
## 8 skinny
                    1777
## 9 full figured 1009
## 10 overweight
                    444
## 11 jacked
                     421
## 12 used up
                     355
## 13 rather not say
                      198
```



## Pro tip: add\_count()

Add count without first creating an initial table, joining etc.:

```
okcupid %>%
  mutate(id = row_number()) %>%
  select(id, body_type, sex) %>%
  add_count(body_type, name = "n_bt") %>%
  filter(n_bt > 10000) %>%
  head(5)
```



# **Beyond Basics**



## A simple answer to the religion question?

okcupid %>% count(religion)

```
## # A tibble: 46 \times 2
##
      religion
                                                         n
##
      \langle chr \rangle
                                                    \langle int \rangle
                                                     2724
    1 agnosticism
    2 agnosticism and laughing about it
                                                     2496
##
    3 agnosticism and somewhat serious about it 642
##
    4 agnosticism and very serious about it
                                                      314
##
    5 agnosticism but not too serious about it
                                                     2636
##
   6 atheism
                                                     2175
## 7 atheism and laughing about it
                                                     2074
##
   8 atheism and somewhat serious about it
                                                      848
                                                      570
    9 atheism and very serious about it
## 10 atheism but not too serious about it
                                                     1318
## # ... with 36 more rows
```



## Recoding with case\_when()

```
okcupid <- okcupid %>% mutate(religion2 = case_when(
   str_detect(religion, "agnosticism") | str_detect(religion, "athe
   str_detect(religion, "buddhism") ~ "buddhist",
   str_detect(religion, "christianity") | str_detect(religion, "cat
   str_detect(religion, "judaism") ~ "jewish",
   str_detect(religion, "hinduism") ~ "hindu",
   str_detect(religion, "islam") ~ "muslim",
   TRUE ~ "NA"))

okcupid %>% count(religion2, sort = TRUE)
```



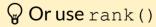
# Getting extreme observations with top\_n() and top frac()

```
okcupid %>%
  select(sex, age) %>%
  group by (sex) %>%
  top n(3, age)
## # A tibble: 33 x 2
## # Groups: sex [2]
##
     sex
         age
  <chr> <dbl>
## 1 f
       110
## 2 m
            69
##
           69
       69
       69
            69
## 7 m
            69
##
             69
             69
## 10 m
        69
## # ... with 23 more rows
```



# Unfortunately top\_n does not deal with ties, could use slice for that:

```
okcupid %>%
  select(sex, age) %>%
  group by(sex) %>%
  top \overline{n(3, age)} \%
  slice(1:3)
## # A tibble: 6 x 2
## # Groups: sex [2]
## sex age
## <chr> <dbl>
## 1 f 110
## 2 f 69
## 3 f 69
## 4 m 69
          69
## 5 m
## 6 m
         69
```





### Remove duplicates with distinct()

```
okcupid %>%
  filter(diet == "kosher") %>%
  distinct (body type, drugs)
## # A tibble: 7 x 2
## body type drugs
##
  <chr> <chr>
## 1 fit
               <NA>
## 2 <NA>
           never
## 3 used up <NA>
               never
## 4 fit
## 5 skinnv
           never
## 6 a little extra never
## 7 jacked
           never
```

♥ distinct() is much more powerful than unique(), see ?distinct.

To count number of distinct obs look at n\_distinct()



## The \_at(), \_if() and \_all() families

Many of the verbs we've seen come with these suffixes:

```
okcupid %>%
  select if(is.numeric)
## # A tibble: 59,946 x 4
##
     age height income height cm
                     <dbl>
##
    <dbl> <dbl> <dbl>
      22
           75
                      190.
##
                -1
##
  2
    35 70 80000 178.
  3 38 68
##
                -1 173.
##
  4 23
           71 20000 180.
  5 29
                   168.
##
           66
                -1
##
  6 29 67
               -1
                     170.
##
  7 32 65 -1 165.
##
  8 31 65 -1 165.
##
  9 24 67 -1 170.
## 10 37
                      165.
        65 -1
## # ... with 59,936 more rows
```

Do you see something strange?



# Take care of those missing observations for me without breaking the pipe:

```
okcupid %>%
  na if (-1) %>%
  select if(is.numeric)
## # A tibble: 59,946 x 4
##
       age height income height cm
     <dbl> <dbl> <dbl>
                          <dbl>
##
##
       22
              75
                           190.
                    NA
##
     35
              70 80000
                           178.
##
     38
             68
                          173.
                    NA
##
     23
             71 20000
                         180.
##
       29
                          168.
             66
                    NA
##
       29
              67
                           170.
                    NA
##
     32
              65
                          165.
                    NA
##
   8 31
              65
                    NA
                           165.
   9 24
##
              67
                           170.
                    NA
## 10 37
              65
                           165.
                    NA
```



## # ... with 59,936 more rows

### Transform all my numeric columns with log:

```
okcupid %>%
  na if (-1) %>%
  select if(is.numeric) %>%
  mutate all(log)
## # A tibble: 59,946 x 4
      age height income height cm
##
##
  <dbl> <dbl> <dbl>
                        <dbl>
##
  1 3.09 4.32 NA 5.25
  2 3.56 4.25 11.3 5.18
##
  3 3.64 4.22 NA
##
                       5.15
  4 3.14 4.26 9.90
##
                      5.19
##
  5 3.37 4.19 NA
                      5.12
## 6 3.37 4.20 NA 5.14
## 7 3.47 4.17 NA
                      5.11
## 8 3.43 4.17 NA
                       5.11
## 9 3.18 4.20 NA
                       5.14
## 10 3.61 4.17 NA
                         5.11
## # ... with 59,936 more rows
```

```
\bigcirc Also see mutate_if()
```



### Same but add sqrt and keep original columns:

```
okcupid %>%
  na if (-1) %>%
  select if(is.numeric) %>%
  mutate all(list(logged = log, sgrted = sgrt))
## # A tibble: 59,946 x 12
##
       age height income height cm age logged height logged income logged
                           <dbl>
##
     <dbl>
           <dbl>
                  <dbl>
                                     <dbl>
                                                  <dbl>
                                                              <dbl>
        22
##
              75
                    NA
                            190.
                                     3.09
                                                  4.32
                                                              NA
     35
                                     3.56
##
              70 80000
                            178.
                                                  4.25
                                                              11.3
                           173.
##
       38
              68
                                     3.64
                                                  4.22
                    NA
                                                              NA
##
       23
              71 20000
                         180.
                                     3.14
                                                  4.26
                                                               9.90
##
       29
                                                  4.19
              66
                    NA
                        168.
                                   3.37
                                                              NA
##
   6 29
              67
                    NA
                        170. 3.37
                                                  4.20
                                                              NA
##
  7 32
              65
                    NA
                           165.
                                   3.47
                                                  4.17
                                                              NA
  8 31 65
##
                           165.
                                     3.43
                                                  4.17
                    NA
                                                              NA
##
       24
              67
                           170.
                                     3.18
                                                  4.20
                    NA
                                                              NA
## 10
        37
              65
                    NA
                            165.
                                      3.61
                                                  4.17
                                                              NA
## # ... with 59,936 more rows, and 5 more variables: height cm logged <dbl
####
      age sgrted <dbl>, height sgrted <dbl>, income sgrted <dbl>,
      height cm sqrted <dbl>
####
```



### Same but take care of zeros under log:

```
okcupid %>%
  na if (-1) %>%
  select if(is.numeric) %>%
  mutate all(list(logged = function(x) log(x + 1), sqrted = sqrt))
## # A tibble: 59,946 x 12
##
       age height income height cm age logged height logged income logged
##
                           <dbl>
     <dbl>
           <dbl>
                  <dbl>
                                      <dbl>
                                                   <dbl>
                                                                <dbl>
        22
                                      3.14
##
              75
                     NA
                            190.
                                                   4.33
                                                               NA
     35
##
              70 80000
                            178.
                                      3.58
                                                   4.26
                                                                11.3
                            173.
##
       38
              68
                                      3.66
                                                   4.23
                     NA
                                                               NA
##
        23
              71 20000
                          180.
                                      3.18
                                                   4.28
                                                                9.90
##
        29
                                                   4.20
              66
                     NA
                         168.
                                      3.40
                                                               NA
##
   6 29
              67
                     NA
                         170.
                                      3.40
                                                   4.22
                                                               NA
##
  7 32
              65
                    NA
                            165.
                                      3.50
                                                   4.19
                                                               NA
  8 31 65
##
                            165.
                                      3.47
                                                   4.19
                    NA
                                                               NA
##
        24
              67
                           170.
                                      3.22
                                                   4.22
                     NA
                                                               NA
## 10
        37
              65
                     NA
                            165.
                                      3.64
                                                   4.19
                                                               NA
## # ... with 59,936 more rows, and 5 more variables: height cm logged <dbl
####
      age sgrted <dbl>, height sgrted <dbl>, income sgrted <dbl>,
      height cm sqrted <dbl>
####
```



### Same but select only non-negative columns:

```
is non negative \leftarrow function(x) is.numeric(x) && (is.na(x) || x >=
okcupid %>%
  na if (-1) %>%
  select if (is non negative) %>%
  mutate all(list(logged = function(x) log(x + 1), sqrted = sqrt))
## # A tibble: 59,946 x 12
##
       age height income height cm age logged height logged income logged
##
     <dbl>
          <dbl>
                  <dbl>
                           <dbl>
                                     <dbl>
                                                 <dbl>
                                                              <dbl>
##
        22
              75
                           190.
                                     3.14
                                                  4.33
                    NA
                                                              NA
##
     35
              70 80000
                           178.
                                     3.58
                                                  4.26
                                                              11.3
##
     38
              68
                    NA
                           173.
                                     3.66
                                                  4.23
                                                              NA
##
       23
              71 20000
                          180.
                                     3.18
                                                  4.28
                                                              9.90
##
       29
              66
                        168.
                                   3.40
                                                  4.20
                    NA
                                                              NA
##
   6 29
              67
                        170.
                                   3.40
                                                  4.22
                    NA
                                                              NA
##
  7 32
              65
                    NA
                        165.
                                   3.50
                                                  4.19
                                                              NA
  8 31 65
##
                    NA
                           165.
                                     3.47
                                                  4.19
                                                              NA
##
     24
              67
                    NA
                          170.
                                     3.22
                                                  4.22
                                                              NA
## 10
        37
              65
                           165.
                                     3.64
                                                  4.19
                    NA
                                                              NA
## # ... with 59,936 more rows, and 5 more variables: height cm logged <dbl
####
      age sqrted <dbl>, height sqrted <dbl>, income sqrted <dbl>,
## #
      height cm sqrted <dbl>
```



# On second thought log would probably be appropriate just for income and height cm (not really, just for demo):

```
okcupid %>%
  na if (-1) %>%
  mutate at(c("income", "height cm"),
             list(logged = function(x) log(x + 1), sqrted = sqrt))
   select(ends with("logged"), ends with("sqrted"))
## # A tibble: 59,946 x 4
##
      income logged height cm logged income sqrted height cm sqrted
##
              <dbl>
                                <dbl>
                                               <dbl>
                                                                 <dbl>
##
   1
              NA
                                 5.25
                                                 NA
                                                                  13.8
##
              11.3
                                 5.19
                                                283.
                                                                  13.3
##
                                 5.16
                                                                  13.1
              NA
                                                 NA
##
              9.90
                                 5.20
                                                141.
                                                                  13.4
##
                                 5.13
                                                 NA
                                                                  12.9
              NA
##
                                 5.14
                                                                  13.0
              NA
                                                 NA
##
   7
                                 5.11
                                                                  12.8
              NA
                                                 NA
##
                                 5.11
                                                                  12.8
              NA
                                                 NA
##
                                 5.14
                                                                  13.0
              NA
                                                 NA
## 10
                                 5.11
                                                                  12.8
              NA
                                                 NA
## # ... with 59,936 more rows
```



## Dealing with NAs

You've already seen na\_if(). We could simply, always, keep those NAs in income:

```
okcupid <- okcupid %>%
  mutate(income = ifelse(income == -1, NA, income))
```

#### Or:

```
okcupid <- okcupid %>%
  mutate(income = na_if(income, -1))
```

Dropping NAs with, well, drop na():

```
okcupid_no_nas <- okcupid %>% drop_na()
```



### Replacing NAs with, well, replace na():

```
okcupid back to minus1 <- okcupid %>% replace na(list(income = -1)
```

### Could be useful for imputing NAs, say the median:

```
okcupid na income imputed <- okcupid %>%
 replace na(list(income = median(.$income, na.rm = TRUE)))
```



# Sampling with sample\_n() and sample frac()

```
okcupid %>% select(drugs, age, income, sex) %>%
  group by (drugs) %>%
  sample n(3, replace = TRUE)
  \# A tibble: 12 x 4
## # Groups: drugs [4]
##
     drugs
                age income sex
##
     <chr> <dbl>
                    <dbl> <chr>
##
  1 never
                 22
                        NA m
##
                 36
                    NA f
  2 never
##
                52
  3 never
                       NA m
##
                 30
                    30000 f
  4 often
##
   5 often
            21
                        NA m
## 6 often
              26
                    NA m
##
  7 sometimes
              33
                       NA m
              32
                    50000 m
  8 sometimes
##
              28
   9 sometimes
                       NA f
               35
## 10 <NA>
                    NA f
## 11 <NA>
                 27
                    NA m
## 12 <NA>
                 29
                       NA m
```



# Put it in a function



# Compose a function which would accept an unquoted variable

```
count_var_for_gender <- function(var, gender) {
  okcupid %>%
    filter(sex == gender) %>%
    count({{var}}, sort = TRUE)
}

count_var_for_gender(body_type, "f") %>% head(9)
```

```
## # A tibble: 9 \times 2
## body type
                     n
## <chr>
                <int>
                 5620
## 1 average
## 2 fit
                 4431
               3811
## 3 curvy
                 2703
## 4 <NA>
## 5 thin
                 2469
## 6 athletic
                2309
## 7 full figured 870
## 8 a little extra 821
                  601
## 9 skinny
```



### Making a data. frame function pipeable

```
transform_all_my_numerics <- function(df, transformation) {
   df %>% mutate_if(is.numeric, transformation)
}

okcupid %>%
   transform_all_my_numerics(log) %>%
   select_if(is.numeric)
```

```
## # A tibble: 59,946 x 4
##
      age height income height cm
   <dbl> <dbl> <dbl>
##
                         <dbl>
  1 3.09 4.32 NA
##
                        5.25
  2 3.56 4.25 11.3
##
                       5.18
  3 3.64 4.22 NA
##
                        5.15
  4 3.14 4.26 9.90 5.19
##
##
  5 3.37 4.19 NA
                        5.12
## 6 3.37 4.20 NA
                         5.14
## 7 3.47 4.17 NA
                         5.11
## 8 3.43 4.17 NA
                       5.11
## 9 3.18 4.20 NA
                       5.14
## 10 3.61 4.17 NA
                         5.11
## # ... with 59,936 more rows
```



### invisible()

If your function does not return a data.frame make it!

```
print n rows <- function(df) {</pre>
  cat("number of rows: ", nrow(df), "\n")
  invisible(df)
okcupid %>%
  filter(sex == "m", body type %in% c("fit", "thin", "skinny")) %
  print n rows() %>%
  summarise(mean height = mean(height cm, trim = 0.025))
## number of rows: 11698
## # A tibble: 1 x 1
## mean height
##
          <dbl>
           179.
## 1
```



#### Or even better:

```
filter and print <- function(df, ...) {
  df filtered <- df %>% filter(...)
  cat("number of rows: ", nrow(df filtered), "\n")
  df filtered
okcupid %>%
  filter and print(sex == "m", body type %in% c("fit", "thin", "s)
  summarise(mean height = mean(height cm, trim = 0.025))
## number of rows: 11698
## # A tibble: 1 x 1
## mean height
##
         <dbl>
          179.
## 1

    for better living see glue::glue("number of rows: {nrow(df)}")
```