

Data Wrangling in R - Data Wrangling - 4

One should look for what is and not what he thinks should be. (Albert Einstein)

Module completion checklist

Objective	Complete
Select specific variables using the select command	
Derive new variables from the existing variables using the mutate and transmute commands	

Select

select() helps you select specific
 columns within your dataframe

- We often use this function with pipes(%>%) which we will cover later in the course
- The selection criteria can be written in multiple ways, as shown in the next couple of slides

Usage

```
select(.data, ...)
```

Arguments

- data A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See *Methods*, below, for more details.
- <tidy-select</p>
 One or more unquoted expressions separated by commas. Variable names can be used as if they were positions in the data frame, so expressions like x:y can be used to select a range of variables.

Select a subset

Simply specify the column name(s)

```
# Select columns from `flights` data frame.
select(flights, #<- specify the data frame
    year, #<- specify the 1st column
    month, #<- specify the 2nd column
    day) #<- specify the 3rd column</pre>
```

 You can also specify a range of columns with the range operator (i.e.:)

```
# Select columns from `flights` data frame
select(flights, #<- specify the data frame
    year:day) #<- specify the range of columns</pre>
```

Select by excluding

• Finally, you can select by excluding certain columns using the exclusion operator (i.e. -)

```
# A tibble: 336,776 x 16
  dep_time sched_dep_time dep_delay arr_time sched_arr_time arr_delay carrier
                                                                <dbl> <chr>
     <int>
                    <int>
                              <dbl>
                                       <int>
                                                      <int>
       517
                                                        819
                      515
                                         830
                                                                   11 UA
       533
                      529
                                                        830
                                                                   20 UA
                                         850
       542
                      540
                                        923
                                                        850
                                                                  33 AA
       544
                      545
                                        1004
                                                       1022
                                                                  -18 B6
       554
                                                        837
                                                                  -25 DL
                      600
                                       812
                                     740
                                                        728
       554
                      558
                                                                  12 UA
       555
                      600
                                                        854
                                 -5 913
                                                                 19 B6
       557
                                 -3 709
                                                        723
                                                                  -14 EV
                      600
       557
                      600
                                 -3 838
                                                        846
                                                                   -8 B6
10
       558
                      600
                                         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>
```

Select - helper functions

- Helpers are multiple functions you can use to select variables based on their names
- They act like regular expressions, but in a more simplified manner
- Here are some of the more commonly used helper functions:

Helper Function	Use Case
starts_with("abc")	matches names that begin with "abc"
ends_with("xyz")	matches names that end with "xyz"
contains("ijk")	matches names that contain "ijk"
num_range("x", 1:3)	matches "x1", "x2" and "x3"

Select: Helper functions example

To select columns whose names start with 'arr':

```
select(flights, starts_with("arr"))
# A tibble: 336,776 x 2
   arr_time arr_delay
      <int>
                <dbl>
        830
                   11
        850
                   33
        923
                  -18
       1004
                  -25
       812
        740
                  12
       913
        709
                  -14
        838
        753
```

... with 336,766 more rows

Module completion checklist

Objective	Complete
Select specific variables, sometimes using specific rules, using the select command	
Derive new variables from the existing variables using the mutate and transmute commands	

Mutate

- mutate() is an essential function of dplyr
- It allows us to create new variables using the current data and append these variables to the existing dataframe

```
?dplyr::mutate
mutate(df, # <- dataframe</pre>
       new_col1, # <- rule(s) for the new column</pre>
```

Note: mutate() always adds columns to the end of the dataset, so make sure you are able to see the last columns

```
mutate {dplyr}
                                                            R Documentation
```

Create, modify, and delete columns

Description

mutate() adds new variables and preserves existing ones; transmute() adds new variables and drops existing ones. New variables overwrite existing variables of the same name. Variables can be removed by setting their value to NULL.

Usage

```
mutate(.data, ...)
## S3 method for class 'data.frame'
mutate(
  .data,
  .keep = c("all", "used", "unused", "none"),
  .before = NULL,
  .after = NULL
transmute(.data, ...)
```

Arguments

.data

A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See Methods, below, for more details.

First - create a new dataset

Create the dataset using select ()

```
# A tibble: 336,776 x 7
   year month day dep_delay arr_delay distance air_time
  <int> <int> <int>
                                 <dbl>
                                         <dbl>
                   <dbl>
                                                  <dbl>
   2013
                                          1400
                                                    227
                                   11
                                      1416
                                                227
   2013
                                 20
                                33 1089 160
   2013
                             -18 1576 183
   2013
                             2013
                                                116
   2013
                                12 719 150

      -5
      19
      1065
      158

      -3
      -14
      229
      53

   2013
   2013
                          -3
   2013
                                 -8 	 944 	 140
   2013
                                           733
                                                    138
# ... with 336,766 more rows
```

Mutate - arguments

- The first argument is the dataframe
- The following arguments are the columns we would like to add to the data frame

```
# A tibble: 336,776 x 9
                     day dep_delay arr_delay distance air_time
    year month
                                                                           gain speed
   <int> <int> <int>
                                                                   <dbl> <dbl> <dbl>
                               <dbl>
                                            <dbl>
                                                       <dbl>
                                                        1400
                                                                     227
                                                                                  370.
    2013
                                               11
                                                       1416
                                                                     227 16 374.
    2013
                                               20
                                                       1089
                                                                          31 408.
    2013
                                               33
                                                                    160
                                                                    183
    2013
                                              -18
                                                      1576
                                                                            -17 517.
                                   -6
                                                                            -19 394.
    2013
                                              -25
                                                       762
                                                                    116

      12
      719
      150
      16
      288.

      19
      1065
      158
      24
      404.

      -14
      229
      53
      -11
      259.

                                             12
                                                        719
                                                                    150
                                                                          16 288.
    2013
    2013
    2013
                                  -3
                                                      944
                                   -3
    2013
                                            -8
                                                                    140
                                                                          -5 405.
    2013
                                                                          10 319.
                                                                     138
                                                         733
# ... with 336,766 more rows
```

Transmute

 transmute() is a function that does the same thing as mutate() except it will only keep the new columns

```
transmute(df,  # <- dataframe
    new_col1, # <- rule(s) for new column
    ...)</pre>
```

- The first argument is the dataframe
- The following arguments are the columns that will be included in your new data frame

mutate {dplyr}

R Documentation

Create, modify, and delete columns

Description

mutate() adds new variables and preserves existing ones; transmute() adds new variables and drops existing ones. New variables overwrite existing variables of the same name. Variables can be removed by setting their value to NULL.

Usage

```
mutate(.data, ...)

## S3 method for class 'data.frame'
mutate(
    .data,
    ...,
    .keep = c("all", "used", "unused", "none"),
    .before = NULL,
    .after = NULL
)

transmute(.data, ...)
```

Arguments

.data

A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See *Methods*, below, for more details.

Transmute example

With the same arguments as in the mutate() example, we can see transmute()
 function only returns new columns

Mutate and transmute - useful functions

• When creating new variables with mutate/transmute, there are many helpful widgets and functions that can assist in creating interesting features:

Useful Functions	Explanation
+, -, *, /, ^	all mathematical operators can be used on variables
log, log2, log10	logarithmic functions for variable transformation can be used
%/% and %%	modulus and remainder are useful when converting time
lag(x) and lead(x)	lag and lead allow reference to leading or lagging values - useful for detecting changes in values.
cumsum(x),cummean(x),	cumulative, running functions, mins, max, prod, mean, etc.
cummax(x),cumprod(x)	

Mutate and transmute - useful functions (cont'd)

- Ranking functions are very helpful in data manipulation
- There are several within the dplyr package such as row_number(),
 ntile() and dense_rank()

ranking {dplyr}

R Documentation

Windowed rank functions.

Description

Six variations on ranking functions, mimicking the ranking functions described in SQL2003. They are currently implemented using the built in rank function, and are provided mainly as a convenience when converting between R and SQL. All ranking functions map smallest inputs to smallest outputs. Use desc() to reverse the direction.

Usage

```
row_number(x)
ntile(x = row_number(), n)
min_rank(x)
dense_rank(x)
percent_rank(x)
cume_dist(x)
```

Arguments

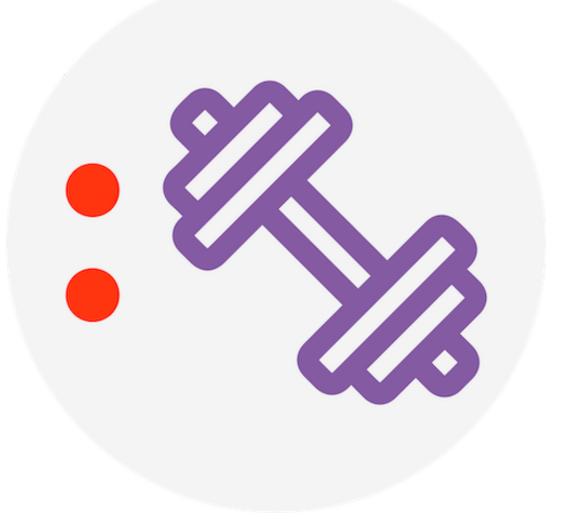
- x a vector of values to rank. Missing values are left as is. If you want to treat them as the smallest or largest values, replace with Inf or -Inf before ranking.
- n number of groups to split up into.

Knowledge check



Exercise

You are now ready to try Tasks 13-18 in the Exercise for this topic



Module completion checklist

Objective	Complete
Select specific variables, sometimes using specific rules, using the select command	
Derive new variables from the existing variables using the mutate and transmute commands	

Data Wrangling: Topic summary

In this part of the course, we have covered:

- Data wrangling basics
- Use dplyr for data wrangling

Congratulations on completing this module!

