# 04 - Data Transformation

ST 597 | Spring 2017
University of Alabama

04-transform.pdf

Data Transformation

`dyplr` Package

Select rows with `filter()`

Arranging (ordering) rows with `arrange()`

Select columns with `select()`

Add or modify variables with `mutate()`

Other `dplyr` functions

# Required Packages and Data

```
library(tidyverse)
library(nycflights13)
```

Remember, if you are getting the error:
```
> Error in library(nycflights13) : there is no
package called 'nycflights13'
```

then you have not installed the `nycflights13` on your computer. You can do so by:

► typing `install.packages("nycflights13")` in console or
► `Tools -> Install Packages...` from RStudio.

# Practice

You need to **practice** to become proficient with the tools we are covering. The best way to do this is start analyzing data that is interesting to you. Here are some places:

- Many R packages have interesting data: `lahman`, `gapminder`, `acs`
- https://www.springboard.com/blog/free-public-data-sets-data-science-project/
- https://www.dataquest.io/blog/free-datasets-for-projects/

Look on-line and find something interests you. I can help you get the data into R if necessary, just ask.

# Data Transformation

# Working with data

When working with data you must:

1. Figure out what you want to do.
2. Precisely describe what you want to do in such a way that the compute can understand it (i.e. program it).
3. Execute the program.

The `dplyr` package makes some of these steps fast and easy:

- By constraining your options, it simplifies how you can think about common data manipulation tasks.
- It provides simple "verbs", functions that correspond to the most common data manipulation tasks, to help you translate those thoughts into code.
- It uses efficient data storage backends, so you spend less time waiting for the computer.

In this section you'll learn the key verbs of `dplyr` in the context of a new dataset on flights departing New York City in 2013.

# nycflights13

To explore the basic data manipulation verbs of `dplyr`, we'll use the `flights` data frame from the `nycflights13` package. This data frame contains all 336,776 flights that departed from New York City in 2013. The data comes from the US Bureau of Transportation Statistics, and is documented in `?nycflights13`.

# nycflights13

```r
#- Load the flights data from nycflights13 package
library(nycflights13)
flights
#> # A tibble: 336,776 × 19
#>    year month   day dep_time sched_dep_time dep_delay arr_time
#>   <int> <int> <int>    <int>          <int>     <dbl>    <int>
#> 1  2013     1     1      517            515         2      830
#> 2  2013     1     1      533            529         4      850
#> 3  2013     1     1      542            540         2      923
#> 4  2013     1     1      544            545        -1     1004
#> 5  2013     1     1      554            600        -6      812
#> 6  2013     1     1      554            558        -4      740
#> # ... with 3.368e+05 more rows, and 12 more variables:
#> #   sched_arr_time <int>, arr_delay <dbl>, carrier <chr>, flight <i
#> #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
#> #   distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

A `tibble` is a special data frame. See Chapter 10 of RDS for more details on the differences between `tibble` and `data.frame`.

# dyplr Package

# dyplr help

- Data Transformation Cheatsheet
- Introduction Vignette from package
  http://cran.r-project.org/web/packages/
  dplyr/vignettes/introduction.html

# dplyr single table verbs

1. `filter()`: find/keep certain rows
   - alternative to `base::subset()`
   - `slice()` to keep by row number
   - `between()`: numeric values in a range

2. `arrange()`: reorder rows
   - alternative to `base::order()`
   - `desc()` to use descending order

3. `select()`: find/keep certain columns
   - helper functions: `starts_with()`, `ends_with()`, `matches()`, `contains()`, `?select`

4. `mutate()`: add/create new variables
   - alternative to `base::transform()`
   - `transmute()`: only return new variables

# dplyr single table verbs

All verbs work similarly:

1. The first argument is a data frame.
2. The subsequent arguments describe what to do with the data frame. You can refer to columns in the data frame directly without using `$`.
3. The result is a new data frame.

Together these properties make it easy to chain together multiple simple steps to achieve a complex result.

> Again, the Data Transformation Cheatsheet is a handy reference.

Select rows with `filter()`

# Select rows by position with `slice()`

To select rows by position, use `slice()`:

```
slice(flights, 5:8)          # selects the 5th - 8th row
#> # A tibble: 4 × 19
#>    year month   day dep_time sched_dep_time dep_delay arr_time
#>   <int> <int> <int>    <int>          <int>     <dbl>    <int>
#> 1  2013     1     1      554            600        -6      812
#> 2  2013     1     1      554            558        -4      740
#> 3  2013     1     1      555            600        -5      913
#> 4  2013     1     1      557            600        -3      709
#> # ... with 12 more variables: sched_arr_time <int>, arr_delay <dbl>,
#> #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest
#> #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>,
#> #   time_hour <dttm>
```

# Select rows by values with `filter()`

`filter()` allows you to subset observations according to specific criteria.

► The first argument is the name of the data frame.
► The second and subsequent arguments are the expressions that filter the data frame (think **and**).
► For example, we can select all flights on January 1st with:

```
filter(flights, month == 1, day == 1)
```

> This is equivalent to the base `subset()` function:
>
> ```
> subset(flights, month == 1 & day == 1)
> ```
>
> `filter()` works similarly to `subset()` except that you can give it any number of filtering conditions, which are joined together with `&`.

## Relational Operators for Numeric Vectors

R provides the standard suite of *numeric* comparison operators:
$>$, $>=$, $<$, $<=$, $!=$ (not equal), and $==$ (equal).

When you're starting out with R, the easiest mistake to make is to use $=$ instead of $==$ when testing for equality. When this happens you'll get a somewhat uninformative error:

```
filter(flights, month = 1)
#> Error: filter() takes unnamed arguments. Do you need `==`?
```

Whenever you see this message, check for $=$ instead of $==$.

# Relational Operators for Character Vectors (and Factors)

- ► == equal to
- ► != not equal to
- ► %in% element of set (use: x %in% set)

```
x = c("aa", "bb", "aa", "bb", "aa", "cc", "dd")
x == "aa"
#> [1]  TRUE FALSE  TRUE FALSE  TRUE FALSE FALSE
x != "aa"
#> [1] FALSE  TRUE FALSE  TRUE FALSE  TRUE  TRUE
x %in% c("aa","bb")
#> [1]  TRUE  TRUE  TRUE  TRUE  TRUE FALSE FALSE
!(x %in% c("aa","bb"))  # x not in set
#> [1] FALSE FALSE FALSE FALSE FALSE  TRUE  TRUE
```

# Logical Operators

Multiple arguments to `filter()` are combined with "and".

```
#- select flights with dest of BHM *and* December
filter(flights, dest=="BHM", month == 12)
```

To get more complicated expressions, you can use Boolean operators. The `|` is read as "or"

```
#- select flights with Nov *or* Dec
filter(flights, month == 11 | month == 12)
```

```
#- dest of BHM *and* (Nov *or* Dec)
filter(flights, dest=="BHM", month == 11 | month == 12)
```

# Logical Dangers

Beware of a common mistake:

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

Note the order isn't like English. This expression doesn't find on months that equal 11 or 12. Instead it finds all months that equal 11 | 12, which is TRUE:

```
11 | 12
#> [1] TRUE
```

In a numeric context (like here), TRUE becomes one, so this finds all flights in January, not November or December.

# Values in a set

Instead you can use the helpful `%in%` shortcut:

```
filter(flights, month %in% c(11, 12))
```
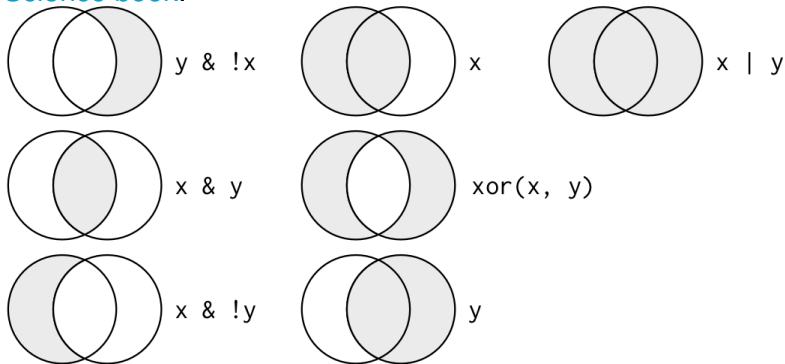
Or `between()`

```
filter(flights, between(month, 11, 12))
```

> The function `between(x, left, right)` is a shortcut for `x >= left &`
> `x<= right` (inclusive).

# More Logical and Relational Operators

- ▶ I have compiled a list of some common logical and relational operators
- ▶ Complete set of Boolean operations from the R for Data Science book:

## Your Turn #1 : filter()

Find all the flights that:

1. Departed in July
2. That flew to Houston (`IAH` or `HOU`)
3. Departed in July and flew to Houston
4. Flew to `Hou` or Originated from 'JFK'
5. That were delayed by more than two hours
6. That arrived more than two hours late, but didn't leave late
7. Had an arrival time earlier than departure time

Understand how each variable is coded (e.g. the integer 1 = January, the integer 517 = 5:17am, etc.).

# Solutions

Arranging (ordering) rows with `arrange()`

# Arrange rows with `arrange()`

- `arrange()` works similarly to `filter()` except that instead of filtering or selecting rows, it reorders them.
- It takes a data frame, and a set of column names (or more complicated expressions) to order by.
- If you provide more than one column name, each additional column will be used to break ties in the values of preceding columns.
- Order by `year`, then `month`, then `day`:

```
arrange(flights, year, month, day)
#> # A tibble: 336,776 × 19
#>     year month   day dep_time sched_dep_time dep_delay arr_time
#>    <int> <int> <int>    <int>          <int>     <dbl>    <int>
#> 1   2013     1     1      517            515         2      830
#> 2   2013     1     1      533            529         4      850
#> 3   2013     1     1      542            540         2      923
#> 4   2013     1     1      544            545        -1     1004
#> 5   2013     1     1      554            600        -6      812
#> 6   2013     1     1      554            558        -4      740
#> # ... with 3.368e+05 more rows, and 12 more variables:
#> #   sched_arr_time <int>, arr_delay <dbl>, carrier <chr>, flight <i
#>     tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
```

# Descending Order

- ▸ By default, `arrange()` orders from smallest to largest
- ▸ Use `desc()` to order a column in descending order:

```
arrange(flights, desc(dep_time))
#> # A tibble: 336,776 × 19
#>    year month   day dep_time sched_dep_time dep_delay arr_time
#>   <int> <int> <int>    <int>          <int>     <dbl>    <int>
#> 1  2013    10    30     2400           2359         1      327
#> 2  2013    11    27     2400           2359         1      515
#> 3  2013    12     5     2400           2359         1      427
#> 4  2013    12     9     2400           2359         1      432
#> 5  2013    12     9     2400           2250        70       59
#> 6  2013    12    13     2400           2359         1      432
#> # ... with 3.368e+05 more rows, and 12 more variables:
#> #   sched_arr_time <int>, arr_delay <dbl>, carrier <chr>, flight <i
#> #   tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>,
#> #   distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

- ▸ This works on categorical data too (alphabetical order)
- ▸ This works on factors too (ordered by `levels`)

# Base R function `order()`

The `dplyr::arrange()` function is a replacement for `order()`

You can accomplish the same thing in base R by the more verbose:

```
flights[order(flights$year, flights$month, flights$day),
        drop = FALSE]
```

**Your Turn #2 : arrange()**

1. Sort flights to find the most delayed flights
2. Sort flights to find the least delayed flights
3. Sort flights by destination and break ties by arrival delay
4. Sort flights to find highest average flight speed
   (distance/air_time)

# Solutions

Select columns with `select()`

# Select columns with `select()`

- ▶ It's not uncommon to get datasets with hundreds or even thousands of variables.
- ▶ In this case, the first challenge is often narrowing in on the variables you're actually interested in.
- ▶ `select()` allows you to rapidly zoom in on a useful subset using operations based on the names or positions of the variables.
- ▶ Select columns **by name**

```
select(flights, year, month, day) # keep year, month, and day columns
```

- ▶ Select columns **by position**

```
select(flights, 1:3) # keep first 3 columns
```

# Other ways to select columns

► *De*select or drop columns using the − symbol

```r
select(flights, -year, -month, -day)  # keep all except year, month, d
```

```r
select(flights, -(1:3))  # keep all except first 3 columns
```

► Select range of columns by name

```r
# Select all columns between year and day (inclusive)
select(flights, year:day)
# Select all columns except those from year to day (inclusive)
select(flights, -(year:day))
```

# Yet more ways to select columns

There are a number of helper functions you can use within `select()`:

- `starts_with("abc")`: matches names that begin with "abc".
- `ends_with("xyz")`: matches names that end with "xyz".
- `contains("ijk")`: matches name that contain "ijk".
- `matches("(.)\\1")`: selects variables that match a regular expression.
  This one matches any variables that contain repeated characters. You'll learn more about regular expressions later in the course
- `num_range("x", 1:3)` matches `x1`, `x2` and `x3`.
- `one_of(x)` selects any names in the vector `x`

See `?select` and Data Transformation Cheatsheet for more details.

# Related functionality

Use `rename()` function to rename a column

```
rename(flights, tail_number = tailnum)
#> # A tibble: 336,776 × 19
#>    year month   day dep_time sched_dep_time dep_delay arr_time
#>   <int> <int> <int>    <int>          <int>     <dbl>    <int>
#> 1  2013     1     1      517            515         2      830
#> 2  2013     1     1      533            529         4      850
#> 3  2013     1     1      542            540         2      923
#> 4  2013     1     1      544            545        -1     1004
#> 5  2013     1     1      554            600        -6      812
#> 6  2013     1     1      554            558        -4      740
#> # ... with 3.368e+05 more rows, and 12 more variables:
#> #   sched_arr_time <int>, arr_delay <dbl>, carrier <chr>, flight <i
#> #   tail_number <chr>, origin <chr>, dest <chr>, air_time <dbl>,
#> #   distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
```

- ▶ Note: this returns a full data frame. It does *not* modify the original.
- ▶ To apply the renaming, use `flights = rename(flights, tail_number = tailnum)`

# Re-arrange Columns

▶ The column order can be rearranged with `select()`. This is especially helpful for viewing on the screen/console

```
select(flights, distance, air_time, origin, dest, carrier)
#> # A tibble: 336,776 × 5
#>    distance air_time origin  dest carrier
#>       <dbl>    <dbl>  <chr> <chr>   <chr>
#> 1      1400      227    EWR   IAH      UA
#> 2      1416      227    LGA   IAH      UA
#> 3      1089      160    JFK   MIA      AA
#> 4      1576      183    JFK   BQN      B6
#> 5       762      116    LGA   ATL      DL
#> 6       719      150    EWR   ORD      UA
#> # ... with 3.368e+05 more rows
```

Add or modify variables with `mutate()`

# Add or modify variables with `mutate()`

- ▶ The job of `mutate()` is to add new (or modify) columns that are functions of existing columns.
- ▶ `mutate()` always adds the new columns at the end of the data frame in order created

```
flights_sml <- select(flights,              # reduce variables
  year:day,
  ends_with("delay"),
  distance,
  air_time
)

mutate(flights_sml,
  gain = arr_delay - dep_delay,            # add gain variable
  speed = distance / air_time * 60         # add speed variable
)
#> # A tibble: 336,776 × 9
#>     year month   day dep_delay arr_delay distance air_time  gain
#>    <int> <int> <int>     <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   -17 516
#>     2013     1     1        -6       -25      762      116   -19
```

# `mutate()` function

- ▶ Note that you can refer to columns that you've just created:

```
mutate(flights_sml,
  gain = arr_delay - dep_delay,
  hours = air_time / 60,
  gain_per_hour = gain / hours        # used the newly created variable
)
#> # A tibble: 336,776 × 10
#>    year month   day dep_delay arr_delay distance air_time  gain
#>   <int> <int> <int>     <dbl>     <dbl>    <dbl>    <dbl> <dbl>
#> 1  2013     1     1         2        11     1400      227   9 3.7
#> 2  2013     1     1         4        20     1416      227  16 3.7
#> 3  2013     1     1         2        33     1089      160  31 2.6
#> 4  2013     1     1        -1       -18     1576      183 -17 3.0
#> 5  2013     1     1        -6       -25      762      116 -19 1.9
#> 6  2013     1     1        -4        12      719      150  16 2.5
#> # ... with 3.368e+05 more rows, and 1 more variables: gain_per_hour
```

> `mutate()` is also used to modify the columns (e.g. `recode()` or change
> data type). E.g., `mutate(flights, flight = as.character(flight)`
> will change `flight` column to a character.

# `transmute()` to only keep new columns

If you only want to keep the newly created columns, use `transmute()` instead of `mutate()` + `select()`

```
transmute(flights,
  gain = arr_delay - dep_delay,
  hours = air_time / 60,
  gain_per_hour = gain / hours
)
#> # A tibble: 336,776 × 3
#>    gain    hours gain_per_hour
#>   <dbl>    <dbl>         <dbl>
#> 1     9 3.783333      2.378855
#> 2    16 3.783333      4.229075
#> 3    31 2.666667     11.625000
#> 4   -17 3.050000     -5.573770
#> 5   -19 1.933333     -9.827586
#> 6    16 2.500000      6.400000
#> # ... with 3.368e+05 more rows
```

# Using aggregate functions in `mutate()`

- ► For statistical analysis, we often want to compare individual values to aggregates
- ► E.g., create the *Z* score for the `distance` column

```
transmute(flights,
          Zdist = (distance - mean(distance))/sd(distance))
#> # A tibble: 336,776 × 1
#>          Zdist
#>          <dbl>
#> 1  0.49109544
#> 2  0.51291660
#> 3  0.06694652
#> 4  0.73112827
#> 5 -0.37902357
#> 6 -0.43766796
#> # ... with 3.368e+05 more rows
```

For each element in the `distance` column, it subtracts the column mean and divides by the column standard deviation.

## Your Turn #3 : mutate()

1. Create a new data frame that contains only the flights that were less than 1000 miles (`distance`). Keep only the columns: `dep_delay`, `arr_delay`, `origin`, `dest`, `air_time`, and `distance`.
2. Add the *Z*-score for departure delays to the new data frame
3. Convert the departure and arrival delays into hours
4. Return only the average flight speed (in mph)
5. Calculate the mean speed

# Solutions

Other `dplyr` functions

# Honorable Mentions: Data frame functions

- ► `distinct()`: retain unique/distinct rows
- ► `sample_n()` and `sample_frac()`: randomly sample rows
- ► `top_n()`: selects and orders the top n rows according to `wt`
- ► `add_column()` add new column in particular position
- ► `add_row()` adds new row(s) to the table

# Honorable Mentions: Dealing with NA's (missing values)

Dealing with missing values (NA) is important, but tedious.
These can help

- `na_if(x, y)` converts the `y` valued elements in `x` to NA

```
x = c(1, 2, -99, 5, 5, -99)
na_if(x, -99)              # replace -99 with NA
#> [1]  1  2 NA  5  5 NA
```

- `coalesce(x, y)` replaces the NA in `x` with `y`

```
x = c(1, 2, NA, 5, 5, NA)
coalesce(x, 0)            # replace NA with 0
#> [1] 1 2 0 5 5 0
```

> These two functions can be used in `mutate()` to modify columns.