

# Session 2 - Dataframes

Jongbin Jung

January 9-10, 2016

# Dependencies

- ▶ Latest version ( $\geq 3.1.2$ ) of R  
(*free* from <https://www.r-project.org/>)
- ▶ Latest version of Rstudio (also *free* from <https://www.rstudio.com/>)
- ▶ A bunch of *free* packages

```
install.packages('nycflights13') # sample data frame  
install.packages('dplyr')  
install.packages('tidyr')
```

# Data Frames: Introduction

- ▶ Data frames are the primary representation of data in R
- ▶ You can think of a data frame as a two-dimensional *table* of data
- ▶ It helps your sanity to always think of data frames as a table where

Each column represents a variable/feature

Each row represents an observation/instance

- ▶ Conceptually, a data frame is also a collection of vectors, i.e., each column is a vector that belongs to the (parent) data frame
- ▶ The fastest path to achieving R-ninja status is to get familiar with data frames

# Data Frames: First Impression

- ▶ Let's load an existing data frame to take a look at

```
# install data package (only need to do once)
install.packages('nycflights13')
```

```
# load data package to workspace
library('nycflights13')
```

- ▶ The `nycflights13` package contains a single data frame named `flights`
- ▶ Contains data (16 variables) on all 336,776 flights that departed NYC (i.e. JFK, LGA, or EWR) in 2013
- ▶ See documentation for details on what the 16 variables are

```
?flights
```

## Data Frames: First Impression (cont'd)

```
head(flights) # take a peek at the data frame
```

```
## Source: local data frame [6 x 16]
##
##   year month   day dep_time dep_delay arr_time
##   (int) (int) (int)   (int)      (dbl)   (int)
## 1  2013     1     1     517         2     830
## 2  2013     1     1     533         4     850
## 3  2013     1     1     542         2     923
## 4  2013     1     1     544        -1    1004
## 5  2013     1     1     554        -6     812
## 6  2013     1     1     554        -4     740
## Variables not shown: arr_delay (dbl), carrier
##   (chr), tailnum (chr), flight (int), origin
##   (chr), dest (chr), air_time (dbl), distance
##   (dbl), hour (dbl), minute (dbl)
```

# Some Question

- ▶ What questions could you ask (and answer) with this data?
  - ▶ how many flights were there each day?
  - ▶ what was the mean departure delay for flights every month/day?
  - ▶ what is the proportion of annual departures from each of the three airports?
  - ▶ what else?
- ▶ By the end of this session, we'll have the tools to answer most (if not all) of the questions you can come up with!

# Data Frame Basics

## Simple Example

- ▶ Use `data.frame()` function to create a data frame
- ▶ Arguments of `data.frame()` are vectors (of equal length) that constitute each column (variable)
- ▶ For example, let's create a data frame of the following table:

Age	Personality	Income
24	Good	2000
22	Bad	5800
23	Good	4200
25	Bad	1500
22	Good	6000



## Simple Example (cont'd)

- ▶ We'll save the data frame to an object (I'll call mine data)

```
data <- data.frame(  # start the data.frame()
  age = c(24, 22, 23, 25, 22),
  personality = c('g', 'b', 'g', 'b', 'g'),
  income = c(2000, 5800, 4200, 1500, 6000)
)  # finish the data.frame() function
```

- ▶ Note that the new lines are just a matter of coding style, i.e., it makes the code easier to read
- ▶ The same data frame can be created in a single line:

```
data <- data.frame(age = c(24, 22, 23, 25, 22),
  personality = c('g', 'b', 'g', 'b', 'g'), income
= c(2000, 5800, 4200, 1500, 6000))
```

## Simple Example (cont'd)

- ▶ Let's take a look at our new data frame

```
data
```

```
##   age personality income
## 1  24           g   2000
## 2  22           b   5800
## 3  23           g   4200
## 4  25           b   1500
## 5  22           g   6000
```

## Indexing: The \$ Operator

- ▶ The \$ operator lets you reference elements of an object (e.g., column vectors of a data frame) in R

```
data$age
```

```
## [1] 24 22 23 25 22
```

```
data$personality
```

```
## [1] g b g b g
```

```
## Levels: b g
```

- ▶ Similar to a . operation in other programming languages (but note that . has no special meaning in R!)

## Indexing: Numeric Row/Column

- ▶ Since a data frame is a table of data, you can treat it like a matrix, and index its entries by [row #, col #] notation

```
data[2, 3]  # item in row 2 column 3
```

```
## [1] 5800
```

```
data[, 2]  # entire column 2
```

```
## [1] g b g b g
```

```
## Levels: b g
```

```
data[4, ]  # entire row 4
```

```
##   age personality income
```

```
## 4   25             b   1500
```

# Indexing: Named Variables

- ▶ Since the columns represent variables with names, you can index columns by a string representing variable names

```
data[, 'age'] # entire 'age' column
```

```
## [1] 24 22 23 25 22
```

```
# entries 3~5 of 'personality' column  
data[3:5, 'personality']
```

```
## [1] g b g
```

```
## Levels: b g
```

## Indexing: Vectors

- ▶ As with vectors/matrices, you can index a data frame with vectors (either numeric or string)

```
data[1:3, c('age', 'income')]
```

```
##   age income
## 1  24   2000
## 2  22   5800
## 3  23   4200
```

```
data[c(1,4), 2:3]
```

```
##   personality income
## 1             g   2000
## 4             b   1500
```

# Conditional Indexing

- Pick-out entries that match specific criteria by first creating a binary vector for indexing

```
# find the 22-year-olds  
ind <- data$age == 22  
data[ind, ] # index rows by binary vector ind
```

```
##   age personality income  
## 2  22             b   5800  
## 5  22             g   6000
```

# Chained Indexing

- ▶ Note that
  - ▶ when you index rows of a single column, the result is a vector
  - ▶ when you index multiple columns, the result is a new data frame
- ▶ You can chain indices to pin-point elements of a data frame
- ▶ For example, all of the following operations are equivalent

```
# Equivalent operations to get the age of  
# third observation (row 3)  
data[3, 1] # if you know that 'age' is column 1  
data[3, 'age']  
data[3,]$age # get 'age' of row 3  
data$age[3] # get third observation of 'age' variable
```



# Write Data Frames to Files

- ▶ Use `write.table()` to write data frames to (text) files
- ▶ The syntax is

```
write.table(x, file = "", append = FALSE,  
quote = TRUE, sep = " ",  
row.names = TRUE, col.names = TRUE)
```

- ▶ For example, to save our sample data to a file named `data.tsv` with the entries of each row separated by a tab character, write

```
write.table(data, file='data.tsv', sep='\t',  
row.names=FALSE) # row names are rarely needed
```

- ▶ Recall, the default directory is the current working directory, specified with `setwd()`, and retrieved with `getwd()`
- ▶ For more options, see documentation

```
?write.table
```

# Read Data Frames from Files

- ▶ To read data frames that exist as text files, use the general `read.table()` function
- ▶ Note that specific options for `read.table()` will depend on the structure of the text file you wish to read (e.g., comma-separated or tab-separated)
- ▶ For example, to read the file we just saved,

```
data <- read.table('data.tsv', header=TRUE, sep='\t')
```

- ▶ Some shortcuts for pre-defined (commonly used) formats

```
read.csv(file) # comma-separated values (.csv)  
read.delim(file) # tab-separated values (.tsv)
```

- ▶ See the documentation for more details

```
?read.table
```

## Read Data from Online Database

- ▶ `read.table()` can also load data frames from an online database
- ▶ While loading data directly from the web is not recommended, this can be useful when making a local copy of an online database
- ▶ For example, to make a local copy of the dataset saved in <http://goo.gl/6fV7UT>

```
address <- 'http://goo.gl/6fV7UT'  
data <- read.table(address, header=TRUE)  
write.table(data, file='data.tsv', sep='\t')
```

- ▶ Note that you can read data in one format (e.g., comma-separated) and save the local copy in another (e.g., tab-separated)

# Exploring Data Frames

## Example Data

- ▶ We'll use a sample dataset from <http://goo.gl/6fV7UT>
- ▶ First, load the data into your workspace

```
address <- 'http://goo.gl/6fV7UT'  
autompg <- read.table(address, header=TRUE)
```

- ▶ The data contains fuel consumption data of 398 vehicles
- ▶ Originally from the UCI Machine Learning Repository
- ▶ See documentation here
  - ▶ <http://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.names>

## Display Structure with str()

- ▶ The str() function is useful for exploring the overall structure of a data frame

```
str(autompg)
```

```
## 'data.frame':      398 obs. of  10 variables:
##  $ mpg          : num  18 15 18 16 17 15 14 14 14..
##  $ cylinders     : int   8 8 8 8 8 8 8 8 8 8 ...
##  $ displacement: num   307 350 318 304 302 429 45..
##  $ horsepower   : Factor w/ 94 levels "?","100.0"..
##  $ weight        : int   3504 3693 3436 3433 3449 4..
##  $ accel         : num   12 11.5 11 12 10.5 10 9 8...
##  $ year          : int   70 70 70 70 70 70 70 70 70..
##  $ origin        : int    1 1 1 1 1 1 1 1 1 1 ...
##  $ model         : Factor w/ 305 levels "amc amba"..
##  $ make          : Factor w/ 37 levels "amc","aud"..
```

# Factors

- ▶ Note that some variables are factors
- ▶ A factor is a data frame representation of categorical variables
- ▶ The entries of a factor variable is defined by `levels`

```
levels(autompg$make)
```

- ▶ Use `unique()` to list the unique values of any variable

```
unique(autompg$year)
```

```
## [1] 70 71 72 73 74 75 76 77 78 79 80 81 82
```

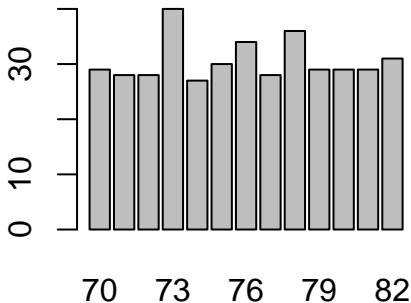
- ▶ Use `factor()` to make a factor variable from non-factor variables

```
autompg$year <- factor(autompg$year)
```

## Basic plots

- ▶ Use `plot()` to generate quick and dirty (but often helpful) plots
- ▶ By default, `plot()` will generate histograms of categorical variables (factors) and scatter plots (with respect to row index) of continuous variables

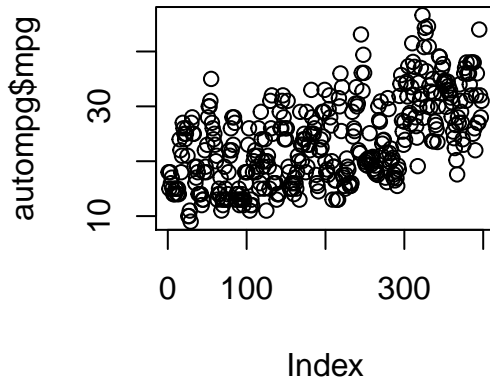
```
plot(autompg$year)
```





## Basic plots (cont'd)

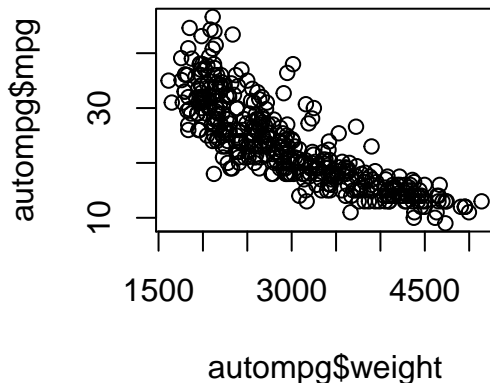
```
plot(autompg$mpg)
```



## Basic plots (cont'd)

- Use syntax `plot(x, y)` to plot two variables

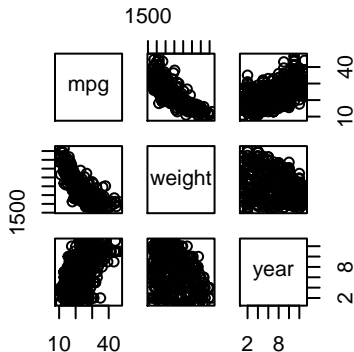
```
plot(autompg$weight, autompg$mpg)
```



## Plotting pairs

- To plot more than two variables against each other, use `pairs()`

```
pairs(autompg[, c('mpg', 'weight', 'year')])
```



- Note that you can plot the entire data frame with `pairs(autompg)`

# Data Frame Basics: Exercise

- ▶ From the `autompg` data
  - ▶ create a new data frame with all the buick vehicles (i.e., `make=="buick"`)
  - ▶ generate a `summary()` of the buick vehicles' `mpg`
  - ▶ make the `cylinders` variable of the buick data frame into a factor
  - ▶ plot a histogram of the buick's `cylinders`
- ▶ These are just (very) basic operations
- ▶ For more complicated operations, we'll use `dplyr` and `tidyr` (covered next)
- ▶ For more sophisticated plots, we'll use `ggplot2` (covered in the next session)

## Exercise Solution

# WARNING

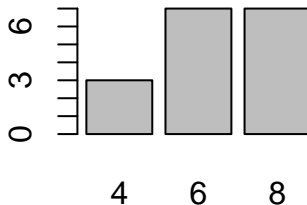
- ▶ Solutions to the exercise are presented in the next slide
- ▶ Try the exercise before proceeding!

## Solution

```
buick_index <- autmpg$make == 'buick'  
buick <- autmpg[buick_index, ]  
summary(buick$mpg)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##    12.00   14.00   17.70   19.18   22.40   30.00
```

```
buick$cylinders <- factor(buick$cylinders)  
plot(buick$cylinders)
```



## Munging Data with dplyr



# Introduction to dplyr

- ▶ dplyr is a package that provides a convenient framework (along with a handful of useful functions) for wrangling data (frames)
- ▶ Install and load the dplyr package like you would any other R package

```
# Install, if you haven't already.  
# Only need to do this once on a single machine.  
install.packages('dplyr')  
# load package into workspace  
library('dplyr')
```

- ▶ We'll primarily use the `flights` data frame from the `nycflights13` package in this part

# Verbs

- ▶ A *verb* in the world of `dplyr` is a function that takes a data frame as its first argument, and returns another data frame as a result
- ▶ For example, the `head()` function can be considered a verb

```
head(flights, n = 10)
```

- ▶ Note that the result of the `head()` function is another data frame (in this case, with 3 rows)
- ▶ The **core idea of `dplyr`** is that most of your data manipulation needs can be satisfied with 5 basic verbs (or 4, depending on how you categorize them)

## Five basic verbs

- ▶ The five basic verbs of `dplyr` and associated actions are presented below

verb	action
<code>filter()</code>	select a subset of <i>rows</i> by specified conditions
<code>select()</code>	select a subset of <i>columns</i>
<code>mutate()</code>	create a <i>new column</i> (usually by operations of existing columns)
<code>arrange()</code>	reorder (sort) <i>rows</i> by values of specified <i>column(s)</i>
<code>summarize()</code>	aggregate values and reduce to single value

- ▶ Some verbs have additional options or convenient wrappers

## Selecting Rows: `filter()`

- ▶ Select a subset of *rows*
- ▶ Multiple conditions can be used
- ▶ Use `&` to specify AND conditions
- ▶ Use `|` to specify OR conditions
- ▶ AND(`&`)/OR(`|`) operations can be used together (where default behavior for multiple conditions is AND)

```
filter(flights, tailnum == 'N14228' & arr_delay > 10)
filter(flights,
       tailnum == 'N14228' | tailnum == 'N24211')
filter(flights,
       tailnum == 'N14228' | tailnum == 'N24211',
       arr_delay > 10)
```

## Selecting Rows: `slice()`

- ▶ To select rows by numerical index (position), use `slice()`
- ▶ For example, to select the first 10 rows

```
slice(flights, 1:10)
```

- ▶ or to select the last 10 rows

```
slice(flights, (n() - 9):n())
```

- ▶ Use `n()` inside a dplyr verb to indicate the *number of rows* of the data frame

## Selecting Columns: `select()`

- ▶ Select a subset of *columns*
- ▶ Either specify the columns that you want to select

```
select(flights, carrier, tailnum)
```

- ▶ Or specify the columns you wish to drop

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

## Selecting Columns: `select()` (cont'd)

- ▶ `dplyr` provides useful helper functions you can use to `select()` columns that match specific criteria such as
  - ▶ `starts_with(x)`: names that start with `x`
  - ▶ `ends_with(x)`: names that end with `x`
  - ▶ `contains(x)`: names that contain `x`
  - ▶ `matches(x)`: names that match the (regular expression) `x`
- ▶ See the documentation for more details

```
?dplyr::select
```

- ▶ While you can assign new column names with `select()` the convenience function `rename()` lets you rename columns while retaining the rest of the data frame

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

## Create New Columns: `mutate()`

- ▶ Create new columns, usually as a function of existing columns
- ▶ You can refer to new columns you just created, inside the same `mutate()` function

```
mutate(flights, gain = arr_delay - dep_delay,  
       speed = distance / air_time * 60,  
       # use the gain column we just created  
       # to create yet another gain_per_hour column  
       gain_per_hour = gain / (air_time / 60)  
       )
```

- ▶ Use `transmute()` to create a new data frame *just from* the new column(s)

```
transmute(flights, gain = arr_delay - dep_delay)
```



## Sorting Rows by Column Value: `arrange()`

- ▶ Reorder the rows of a data frame by the specified column's value
- ▶ Multiple conditions are arranged from left to right
- ▶ Use `desc()` to arrange in descending order

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

## Aggregate Data: summarize()

- ▶ Aggregate/collapse the data into a single row
- ▶ Think of as applying a function to columns

```
summarize(flights, delay = mean(dep_delay))  
# Note that the mean function need help  
# handling NA values  
summarize(flights,  
           delay = mean(dep_delay, na.rm = TRUE))
```

- ▶ More useful as a grouped operation (see next)

# Grouped Operations

- ▶ If a data frame is *grouped*, operations are applied to each group separately, and the results are combined back to a single data frame
- ▶ Use the `group_by()` verb to specify variables to use for generating groups

```
flights_by_day <- group_by(flights, day)
```

- ▶ Some verbs have specific behavior when applied to grouped data

verb	group specific action
<code>arrange()</code>	sort rows within each group
<code>slice()</code>	extract rows within each group
<code>summarize()</code>	aggregate values group-wise

## Grouped slice()

- Retrieve the first 2 departures (rows) of each day

```
slice(flights_by_day, 1:2)
```

```
## Source: local data frame [62 x 16]
```

```
## Groups: day [31]
```

```
##
```

```
##      year month   day dep_time dep_delay arr_time  
##      (int) (int) (int)   (int)      (dbl)    (int)  
## 1   2013     1     1     517         2      830  
## 2   2013     1     1     533         4      850  
## 3   2013     1     2      42        43      518  
## 4   2013     1     2    126       156      233  
## 5   2013     1     3      32        33      504  
## 6   2013     1     3      50       185      203  
## 7   2013     1     4      25        26      505  
## 8   2013     1     4    106       141      201  
## 9   2013     1     5      14        15      503
```

## Grouped summarize()

- Retrieve (1) number of departures (observations), (2) average distance, and (3) average arrival delay for each day (i.e., for flights grouped by day)

```
summarize(flights_by_day, count = n(),  
          dist = mean(distance, na.rm=TRUE),  
          delay = mean(arr_delay, na.rm=TRUE))
```

##	day	count	dist	delay
## 1	1	11036	1039.478	7.3636956
## 2	2	10808	1046.753	6.7680540
## 3	3	11211	1041.299	4.4699187
## 4	4	11059	1037.793	-1.7827199
## 5	5	10858	1037.845	0.4925064
## 6	6	11059	1040.868	-1.7489044

# Multiple (Chained) Operations

- ▶ Consider the following task

*find days when the mean arrival delay OR departure delay was greater than 30*

- ▶ We can achieve the desired result with three operations
  1. `group_by` date (year, month, day)
  2. `summarize` mean arrival/departure delay
  3. `filter` summarized results (i.e., `mean arr_delay > 30 | mean dep_delay > 30`)
- ▶ Note that `dplyr` verbs do **not** modify the original data frame
  - ▶ This is generally a good thing, since it guarantees the integrity of your data
  - ▶ But it makes multiple operations on a data frame difficult
- ▶ There are two (acceptable) ways to apply multiple operations on a data frame, and one is definitely preferred to the other

## Multiple Operations: The OK Way

- ▶ One way to perform multiple operations is to save intermediate data frames as new data frames
- ▶ This method delivers desired results, but makes your workspace quite messy (i.e., you'll end up with a workspace full of intermediate results)

```
flights_by_date <- group_by(flights, year, month, day)
summary_by_date <- summarize(flights_by_date,
  arr = mean(arr_delay, na.rm=TRUE),
  dep = mean(dep_delay, na.rm=TRUE))
big_delay_dates <- filter(summary_by_date,
  arr > 30 | dep > 30)
```

- ▶ This method might be preferred if you need the intermediate results in the future
- ▶ If not, there is a better way to chain multiple operations in with dplyr

# The Pipe Operator %>%

- ▶ The pipe operator, aka the 'magic' operator, takes the output from the verb on its left-hand side, and uses it as the first argument (data frame) for the verb on the right-hand side

```
big_delay_dates <-  
  group_by(flights, year, month, day) %>%  
  summarize(arr = mean(arr_delay, na.rm=TRUE),  
            dep = mean(dep_delay, na.rm=TRUE)) %>%  
  filter(arr > 30 | dep > 30)
```

- ▶ No need to save intermediate results
- ▶ Easier to read (i.e., you can follow the operations step-by-step without too much mental accounting)



# dplyr: Exercise

- ▶ With the `flights` data
  1. find the average speed ( $\text{distance} / \text{air\_time} * 60$ ) by each carrier (ignore NA), and sort the data in descending order of average speed
  2. find the number of flights and average flight time of all flights greater than 10 hours by each carrier in April

## Exercise Solution

# WARNING

- ▶ Solutions to the exercise are presented in the next slide
- ▶ Try the exercise before proceeding!

# Solution 1

```
speed_by_carrier <-  
  group_by(flights, carrier) %>%  
  mutate(speed = distance / air_time * 60) %>%  
  summarize(avg_speed = mean(speed, na.rm=TRUE)) %>%  
  arrange(desc(avg_speed))  
speed_by_carrier
```

##	carrier	avg_speed
## 1	HA	480.3577
## 2	VX	446.1749
## 3	AS	443.6789
## 4	F9	425.1721
## 5	UA	420.8838
## 6	DL	418.4628
## 7	AA	417.4727
## 8	WN	400.5320

## Solution 2

```
april_long_flights <-  
  group_by(flights, month, carrier) %>%  
  filter(month == 4 & hour > 10) %>%  
  summarize(avg = mean(hour, na.rm=TRUE),  
            count = n())  
april_long_flights
```

##	month	carrier	avg	count
## 1	4	9E	16.72074	1085
## 2	4	AA	15.72994	1670
## 3	4	AS	18.36667	30
## 4	4	B6	16.98868	2916
## 5	4	DL	15.89183	2718
## 6	4	EV	16.38317	2876
## 7	4	F9	17.51613	31
## 8	4	FL	15.68398	231

## Reshape Data with tidyr

## Joins with merge

# Reference

- ▶ A great “cheat sheet” for wrangling data with dplyr and tidyr is available for free at <https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>