Before we start...

Access the RStudio Cloud Project (Link in https://datafest.psu.edu/go/)

Workshop - Introduction to R & Data Wrangling

OR

- 1. Download R
- 2. Download RStudio
- 3. Download the `titanic` dataset

https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/problem12.html

Introduction to R and Data Wrangling (with tidyverse)



Penn State, University Libraries, Research Informatics and Publishing

Workshop Housekeeping



Questions?

Use the Q&A, chat, or raise hand feature.



Feedback Survey

After the session, please fill out the Qualtrics survey (link in chat).

Credits

Many images are sourced from the teaching team at Harvard Chan Bioinformatics Core (HBC).

Content was similarly inspired by HBC.

Slides with `(HBC Source)` in the bottom corner indicate the image/table source.

Original source: https://hbctraining.github.io/Training-modules/IntroR/

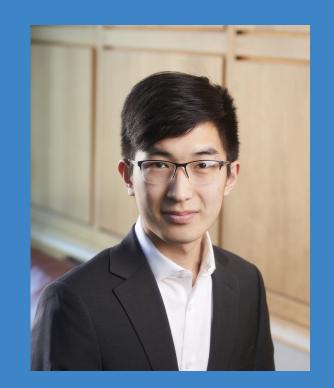
About Me

- Master's in Applied Statistics
- Bachelor's in Computational Statistics

R Experience:

- Self-taught for research (2017-Present)
- Statistics Courses (STAT 184/380)
- Other Projects

Research Consultant, University Libraries



David Chen dzc89@psu.edu

Research Consultant github.com/TheDavidChen

Introduce Yourself!

Name, Major, Year, Favorite Letter

Agenda

Introduction to R

- Why R?
- The RStudio Interface
- Running Code
- Variable Assignment
- Data Types
- Data Structures
- Importing Datasets

Agenda

Introduction to R

- Why R?
- The RStudio Interface
- Running Code
- Variable Assignment
- Data Types
- Data Structures
- Importing Datasets

Data Wrangling with Tidyverse

- Why Tidyverse?
- The Fundamental Functions
 - Subset data
 - Create new variables
 - Grouped summaries
 - Combining datasets
- Exporting Datasets

Why R?

- Free/Open Source
- Platform Agnostic
- Designed for Data Analysis
- Customizable Data
 Visualizations
- Reproducibility
- Big Data!



Logo: The R Foundation

Motivation - Max Number of Rows

Excel: 1,048,576

R: 2⁴⁸ = 281,474,976,710,656

- Free/OpenSource
- Platform agnostic
- Interface for R

Reasons to use Base R instead:



Pace vs Content tradeoff

General Comments

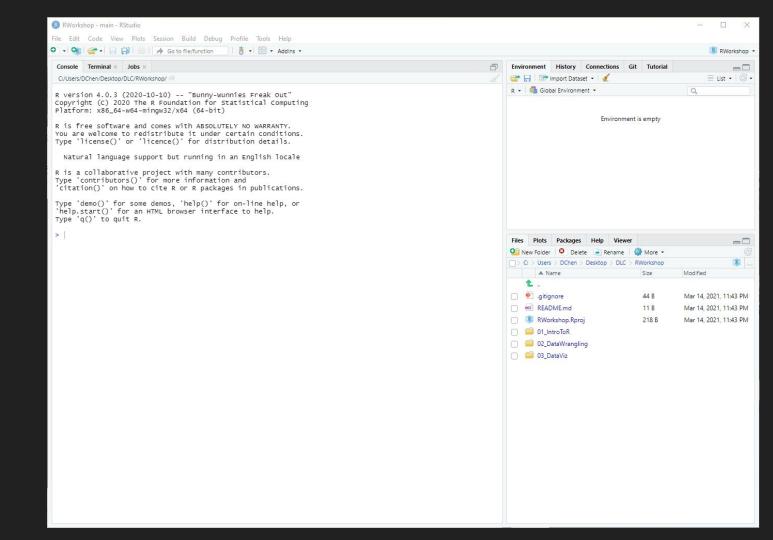
Expect a learning curve. Expect to struggle.

Errors are the best for learning. If they happen:

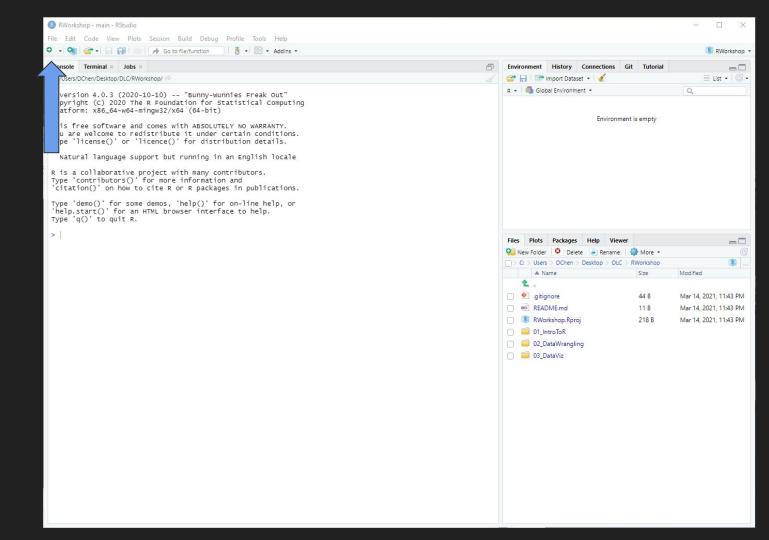
- Read the error and try to solve it
- Share the code that gave you an error
- Share the error message



Open RStudio!



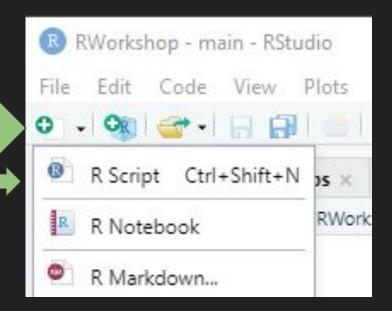
Open RStudio!

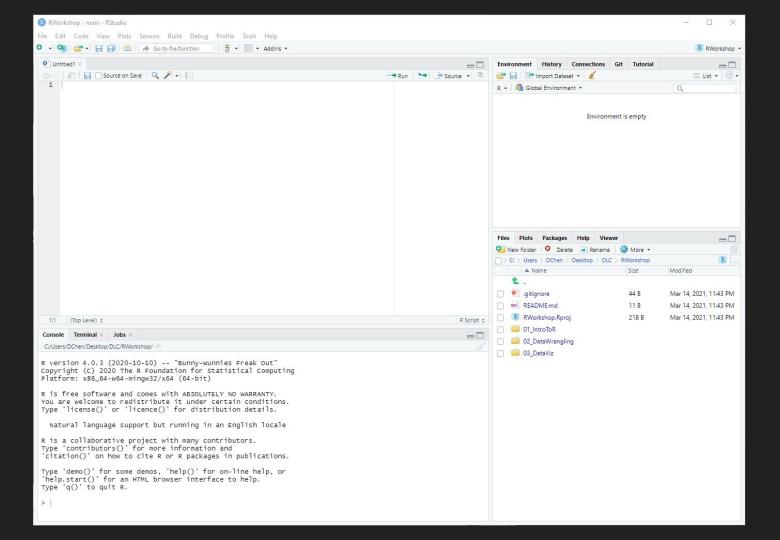


Create an R Script



- Click the top left button
- Click `File -> New File -> R Script`
- Press `Ctrl`+`Shift`+`N`





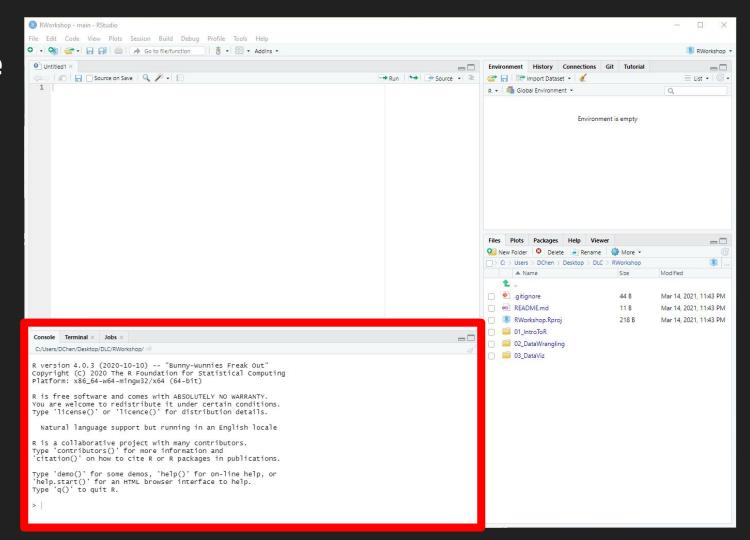
The Console

Run Code

Show Output

Show Errors

Code Unsaved

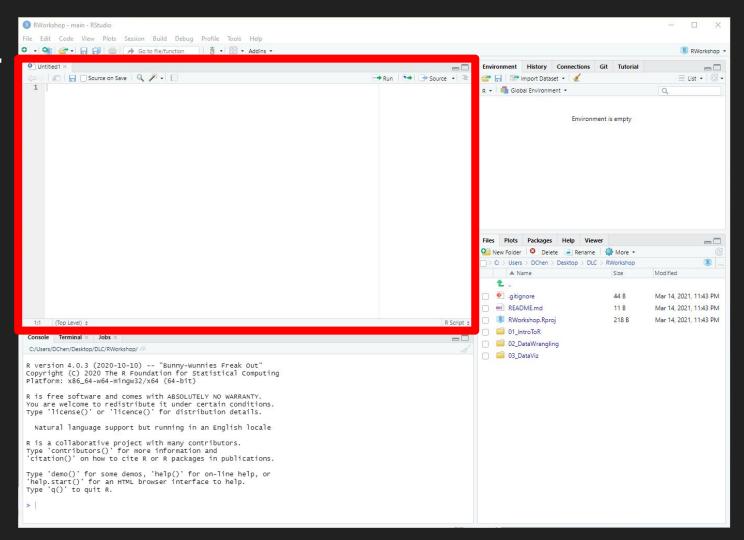


Script Editor

Type code

Run code (output in console)

Saves all code

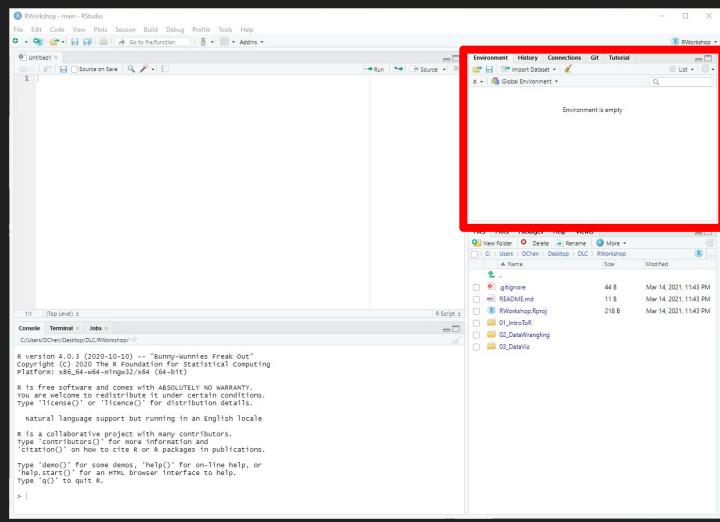


Environment

Import Datasets

Shows all datasets and defined variables

View imported datasets

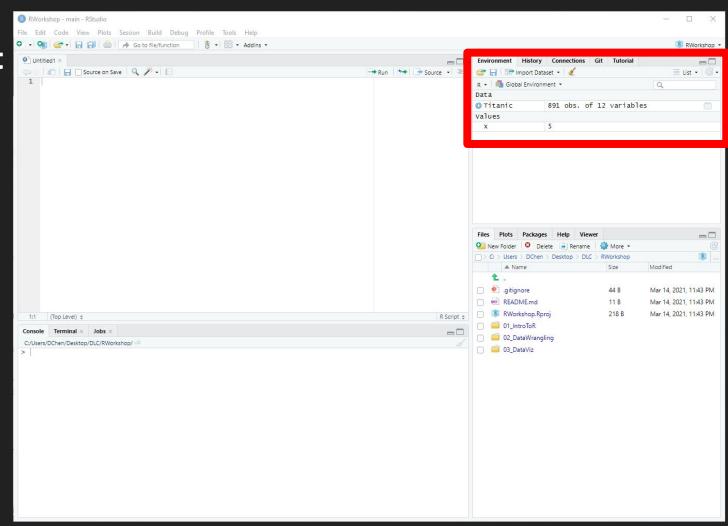


Environment

Import Datasets

Shows all datasets and defined variables

View imported datasets

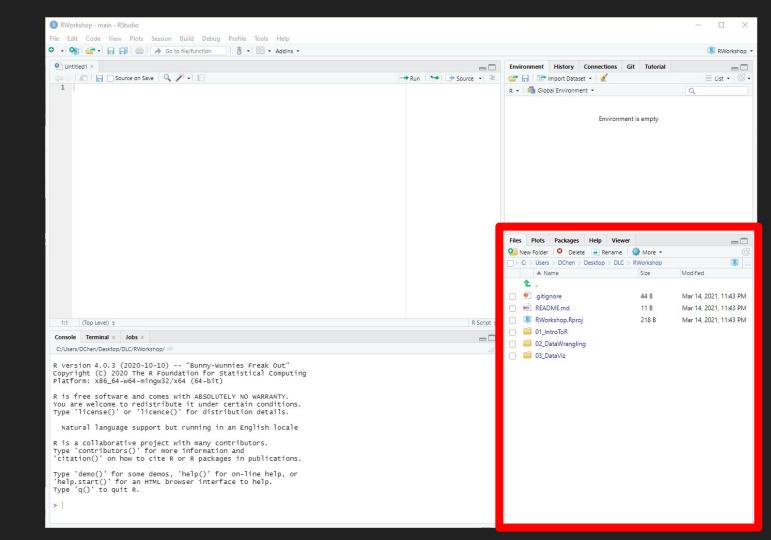


Misc.

Access Files

View created plots

Read help documentation



```
Se[b]()})}var c=function(b){this.element=a(b)};c.VERSION="3.3.7",c.TRANSITION_DURATION=150,c.prot
орdоwn-menu)"),d=b.data("target");if(d||(d=b.attr("href"),d=d&&d.replace(/.*(?=#[^\s]*$)/,"")),!
st a"),f=a.Event("hide.bs.tab",{relatedTarget:b[0]}),g=a.Event("show.bs.tab",{relatedTarget:e[0]
faultPrevented()){var h=a(d);this.activate(b.closest("li"),c),this.activate(h,h.parent(),functio
rigger({type:"shown.bs.tab",relatedTarget:e[0]})})}}},c.prototype.activate=function(b,d,e){func
aria-expanded",!1), attr("aria-expanded",!1),
ia-expanded",!0),h?(b[0].offsetWidth,b.addClass("in")):b.removeClass("fade"),b.parent(".dropdou
().find('[data-toggle="tab"]'\ a++n/"ania avnandad" | 12\ a00-/\lambda.
                                                                        active"),h=e&&
de")||!!d.find("> .fade").
                          Running Code
                                                                        FransitionEnd
;var d=a.fn.tab;a.fn.tab=l
                                                                        {return a.fn.t
"show")};a(document).on("c
se strict";function b(b){return this.each(function(){var d=a(this),e=d.data("bs.affix"),f="ob
-typeof b&&e[b]()})}var c=function(b,d){this.options=a.extend({},c.DEFAULTS,d),this.$target=a
,a.proxy(this.checkPosition,this)).on("click.bs.affix.data-api",a.proxy(this.checkPositionWi
null,this.pinnedOffset=null,this.checkPosition()};c.VERSION="3.3.7",c.RESET="affix affix-top
State=function(a,b,c,d){var e=this.$target.scrollTop(),f=this.$element.offset(),g=this.$targ
"bottom"==this.affixed)return null!=c?!(e+this.unpin<=f.top)&&"bottom":!(e+g<=a-d)&&"bottom"
!!=c&&e<=c?"top":null!=d&&i+j>=a-d&&"bottom"},c.prototype.getPinnedOffset=function(){if(this
.RESET).addClass("affix");var a=this.$target.scrollTop().b=this.$element_offset():return
```

General Comments

If you see: > print("Hello World")

Run print("Hello World") in R Script. Do not include the `>`.

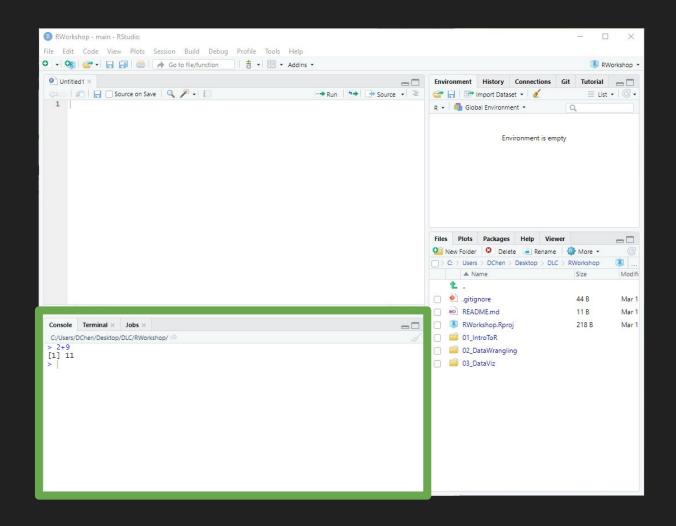
Use `#` to write comments - code after # is not run.

> # This is not run

Type the following calculations:

- > 2+9
- > 15*20
- > 20/5
- > 2^2

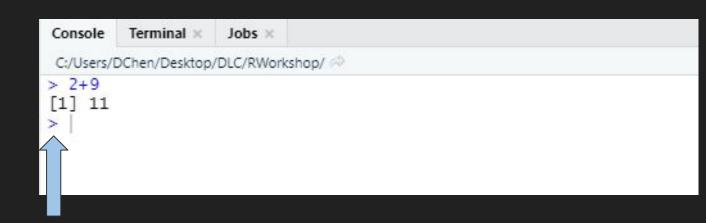
Press 'Enter' to run



Type the following calculations:

- > 2+9
- > 15*20
- > 20/5
- > 2^2

Press 'Enter' to run

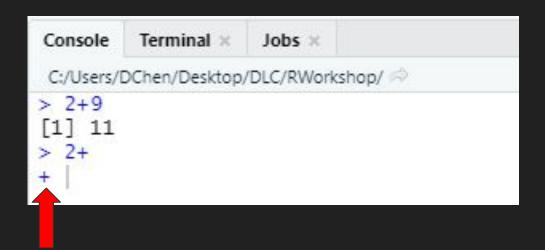


`>` means it's ready for new code

Type the following calculations:

- > 2+9
- > 15*20
- > 20/5
- > 2^2

Press 'Enter' to run



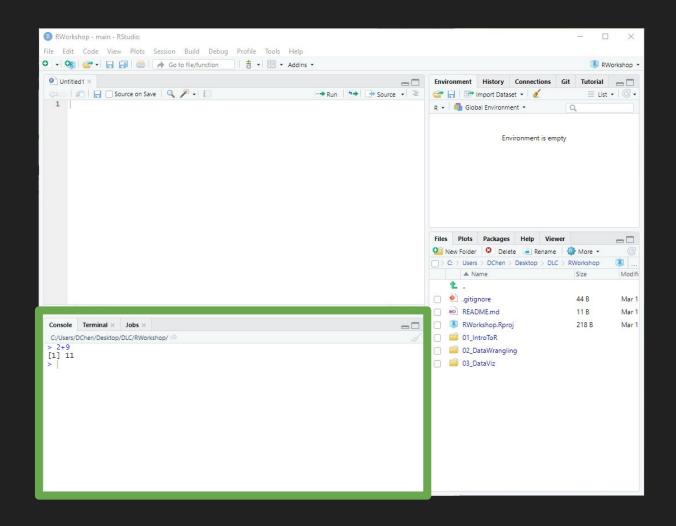
If you see + on the new line:

- Reset with the `Esc` key
- Continue from the previous code

Type the following calculations:

- > 2+9
- > 15*20
- > 20/5
- > 2^2

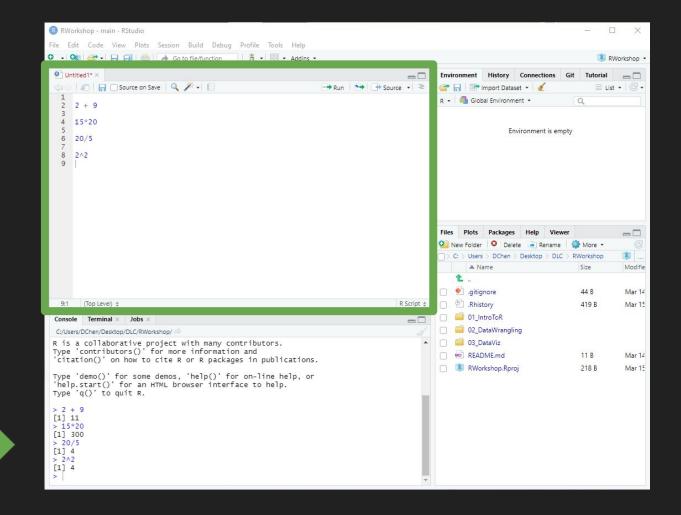
Press 'Enter' to run



Running Code in the Script

Code in the script can be saved and ran repeatedly

Output shows up in the console

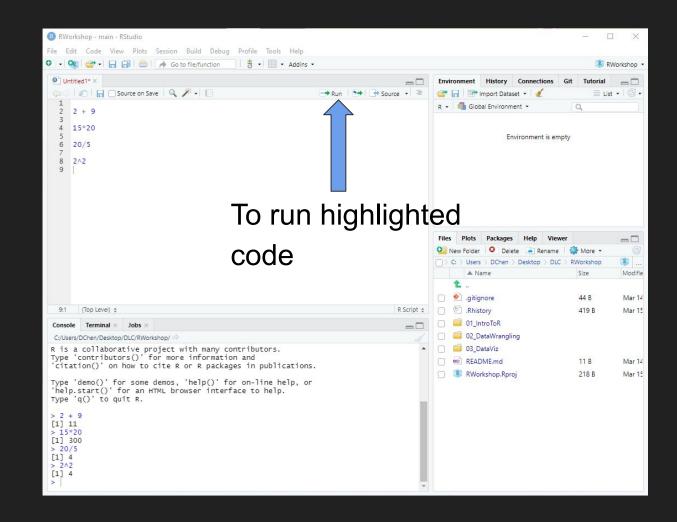


Running Code in the Script

Type code into the R Script

Click or highlight the line of code and

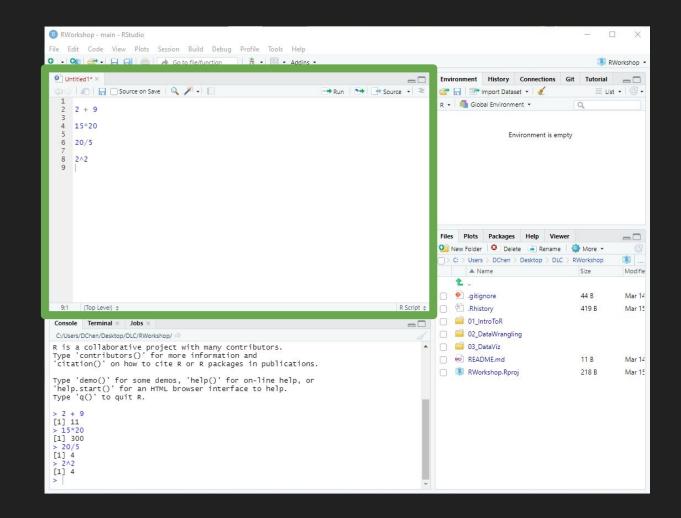
- press `Run`
- press `CTRL` +`ENTER`

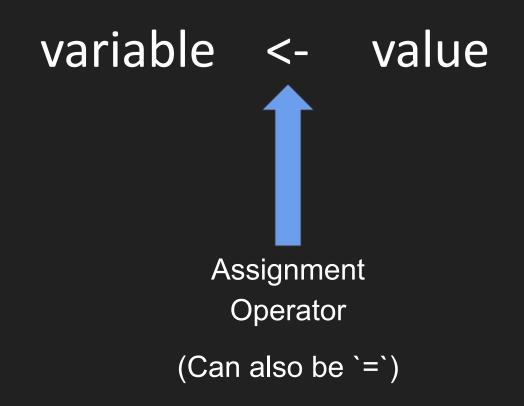


Running Code in the Script

Type the following calculations into the R Script:

- > 2+9
- > 15*20
- > 20/5
- > 2^2





variable <- value

Replaces values

$$> x + 10$$

variable <- value

Replaces values

$$> x + 10$$

Redefine over and over

variable <- value

Replaces values

$$> x + 10$$

Redefine over and over

Multiple variables

Variable Assignment - Comments

R is case sensitive

- > value <- 5
- > VALUE + 5

Variable Assignment - Comments

R is case sensitive

Variable names must be one line and start with a character

> value <- 5

> 1x <- 3

not valid

> VALUE + 5

> x y <- 3

not valid

Functions

function_name(input)

- > sqrt(4)
- > print("Hello World")

Learn more about functions

- Google "<function name> in R"
- Use `help()` or `?` in R to see documentation

- > help(sqrt)
- > ?sqrt

Data Types

Data Type	Examples
Numeric	-5, 1, 3.33, 100, pi

Data Type	Examples	
Numeric	-5, 1, 3.33, 100, pi	
Integer	-5L, 1L, 100L	

Data Type	Examples
Numeric	-5, 1, 3.33, 100, pi
Integer	-5L, 1L, 100L
Character	'words', "3.33", 'TRUE', "1L"

Data Type	Examples	
Numeric	-5, 1, 3.33, 100, pi	
Integer	-5L, 1L, 100L	
Character	'words', "3.33", 'TRUE', "1L"	
Boolean/Logical	TRUE, FALSE, T, F	

(HBC Source)

Data Type	Examples	
Numeric	-5, 1, 3.33, 100, pi	
Integer	-5L, 1L, 100L	
Character	'words', "3.33", 'TRUE', "1L"	
Boolean/Logical	TRUE, FALSE, T, F	

Common Data Types - Boolean

Operator	Description	
<	Less than	
<=	Less than or equal to	
>	Greater than	
>=	Greater than or equal to	
==	Equal to	
!=	Not equal to	

Common Data Types - Boolean

Examples

TRUE

- 5 > 3
- 5!=3

FALSE

- 5 <= 3
- 8 == 10

Description

<=

<

Less than or equal to

>

Greater than

Less than

>=

Greater than or equal to

==

Equal to

!=

Not equal to

Common Data Types - Boolean

Try	ıit ر	you	rself!
-	/	,	

Applies to variables!

Operator

<=

>

>=

__

!=

Description

Less than

Less than or equal to

Greater than

Greater than or equal to

Equal to

Not equal to

Check the type of a variable/object with `class()`:

- > class("4")
- > class(4)

Certain functions require specific data types:

> sqrt("4")

If a number is in quotations, think of it as the word instead of the number.

• "5" != 5

If a number is in quotations, think of it as the word instead of the number.

• "5" != 5

TRUE/FALSE are actually encoded as 1/0

TRUE == 1; FALSE == 0

If a number is in quotations, think of it as the word instead of the number.

• "5" != 5

TRUE/FALSE are actually encoded as 1/0

• TRUE == 1 ; FALSE == 0

Characters without quotation marks are variables!

Hi <- 5

If a number is in quotations, think of it as the word instead of the number.

"5" != 5

TRUE/FALSE are actually encoded as 1/0

• TRUE == 1 ; FALSE == 0

Characters without quotation marks are variables!

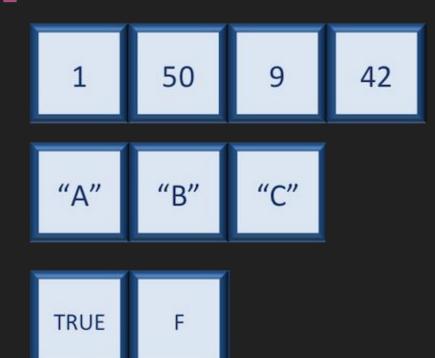
Hi <- 5

Be careful of single `=`; outside of functions, it is assignment!

- 5 = 5 --> error!
- x = 5 is equivalent to x < -5, but not recommended

Data Structures

Vectors



A column in excel

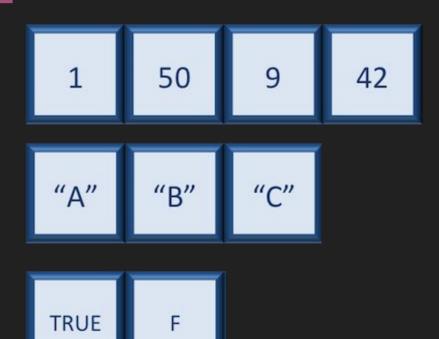
Any length from 1 onwards

All values have to be same type (all numeric, boolean, character, etc.)

Defined with function `c()`

(HBC Source)

Vectors



A column in excel

Any length from 1 onwards

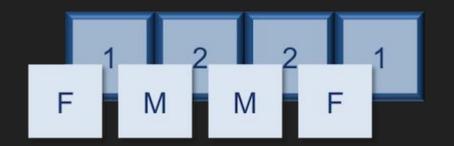
All values have to be same type (all numeric, boolean, character, etc.)

Defined with function `c()`

> lengths <- c(1, 50, 9, 42)

(HBC Source)

Factors



A special vector with values assigned to each factor level (category)

Data Frame



Most common data format; equivalent to an excel sheet.

Multiple vectors combined together

Columns must be same type

All columns must have equal number of rows

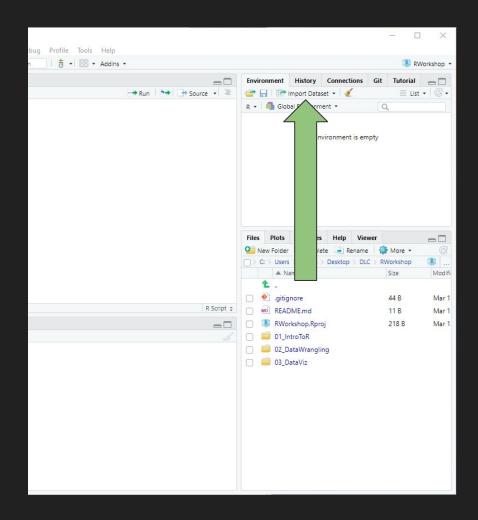
(HBC Source)

Reading in Data

Data type	Extension	Function	Package
Comma separated values	CSV	read.csv()	utils (default)
		read_csv()	readr (tidyverse)
Tab separated values	tsv	read_tsv()	readr
Other delimited formats	txt	read.table()	utils
Stata version 7-12	dta	read.dta()	foreign
SPSS	sav	read.spss()	foreign
SAS	sas7bdat	read.sas7bdat()	sas7bdat
Excel	xlsx, xls	read_excel()	readxl (tidyverse)

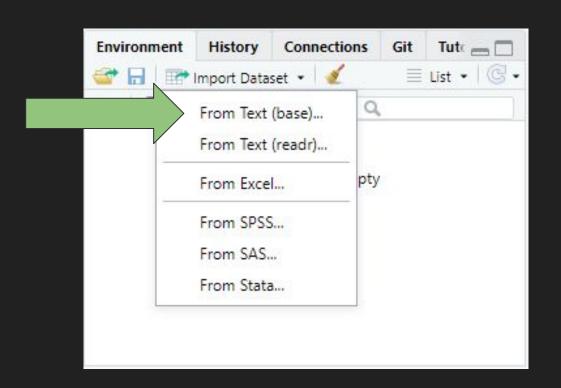
(HBC Source)

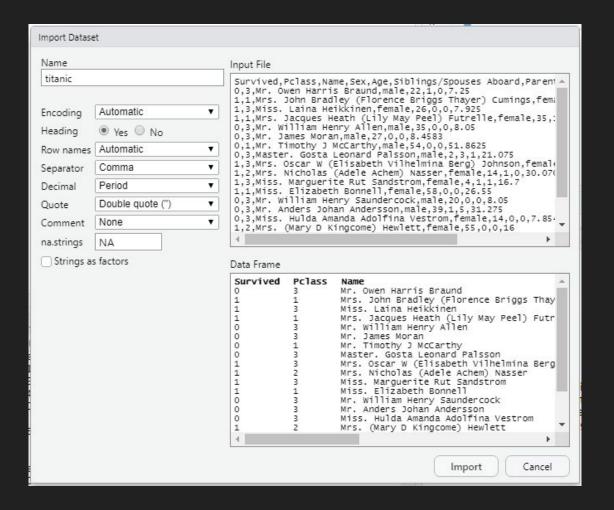
'Import Dataset' in the environment will automatically generate the code for you!



From Text works for almost all basic data files (.csv, .txt, etc.)

Find and select the dataset

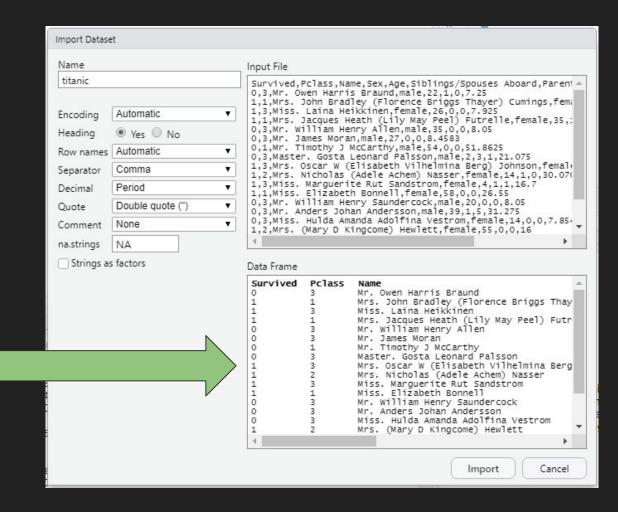




Make sure the data is being imported correctly!

Check:

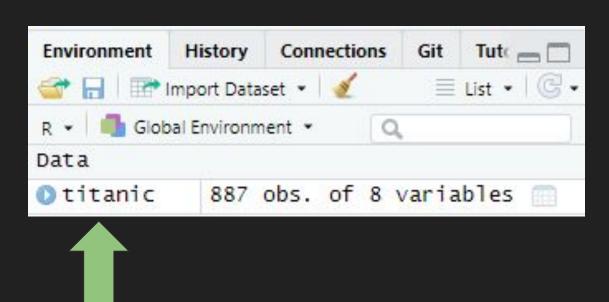
- Headers
- Columns
- Values



Viewing the Datasets

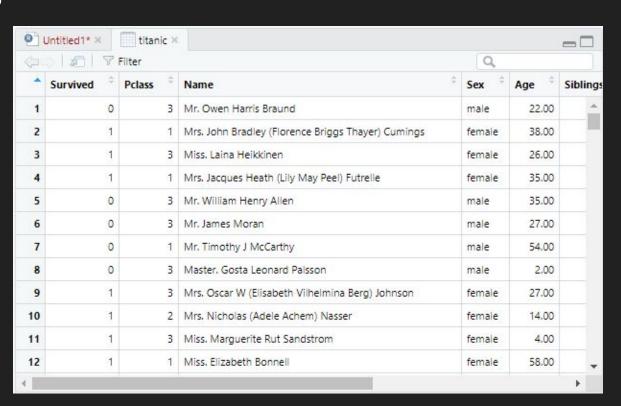
The data is now in the environment!

Click the name (titanic) to view the data!



Viewing the Datasets

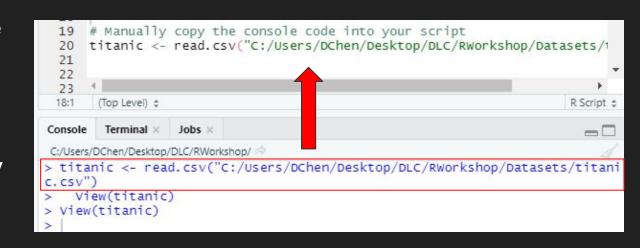
The data view after clicking the name in the environment



After Importing the Data

Copy the generated code into the R Script.

In the future, you will only need to run the code.



Do not include the `>` or the `View()` functions.

Examine the Data - Try These Functions!

Functions - Add `titanic` to ()	
--------------------------------	---	--

str() dim()

summary() nrow()

head() ncol()

tail() colnames()

> str(titanic)

> head(titanic)

> colnames(titanic)

Data Wrangling

Before we start...

Import the `titanic` dataset into R

https://web.stanford.edu/class/archive/cs/cs109/cs109.1166/problem12.html

Install the packages ONCE (to check, skip this step)

> install.packages("tidyverse")

After installing, you must load the packages every time you open R

> library(tidyverse)

What is the purpose of Data Wrangling?

What is the Tidyverse?

Over 87 packages combined together

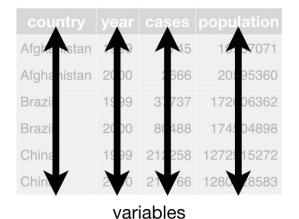
Packages align on data processing philosophy

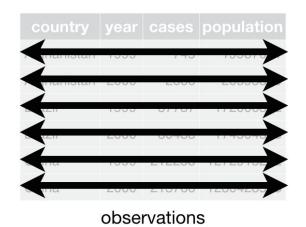
Encompasses all aspects of a project, including data processing, visualization, modeling

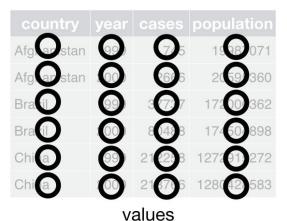


Tidyverse Philosophy

"Variables are in columns, observations are in rows, and values are in cells."







The Fundamental Functions

The Fundamental Functions

select Select columns

select - Select which variables to keep

select(data, `Column1`, `Column2`,...)

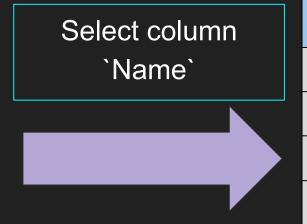
Columns to keep or remove (-)

Use column name, number, or colon

select - Keep columns

select(data, Name)

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87



Name		
Matt		
Mary		
Bill		
Brie		

select - Remove Columns



select(data, -Major)

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

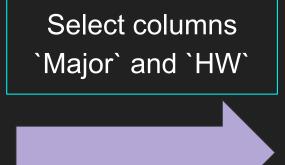
Select all columns except `Major`

Name	HW
Matt	60
Mary	90
Bill	85
Brie	87

select - Keep Multiple Columns

select(data, Major, HW)

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

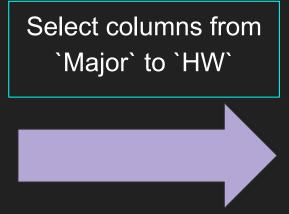


Major	HW
math	60
math	90
biology	85
biology	87

select - Keep Multiple Columns

select(data, Major:HW)

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87



Major	HW
math	60
math	90
biology	85
biology	87

select - Titanic Exercises

1. Select the variables `Survived`, `Sex`, and `Parents.Children.Aboard`

- > select(titanic, _____, ____, ____)
- 2. Remove the `Name` variable, and store the new results in `Titanic_nameless`
- > Titanic_nameless <- select(titanic, _____)
- 3. Store the first four columns of `titanic` in a new dataframe. How many different ways can you do this using `select()`?

select - Titanic Exercise Solutions

- 1. Select the variables 'Survived', 'Sex', and 'Parents.Children.Aboard'
- > select(titanic, Survived, Sex, Parents.Children.Aboard
- 2. Remove the `Name` variable, and store the new results in `Titanic_nameless`
- > Titanic_nameless <- select(titanic, -Name)
- 3. Store the first four columns of `titanic` in a new dataframe. How many different ways can you do this using `select()`?

select - Titanic Exercise Solutions

- 3. Store the first four columns of `titanic` in a new dataframe. How many different ways can you do this using `select()`?
- > t2 <- select(titanic, 1:4)
- > t2 <- select(titanic, Survived:Sex)
- > t2 <- select(titanic, Survived, Pclass, Name, Sex)</p>
- > t2 <- select(titanic, -Age, -Siblings.Aboard, -Parents.Children.Aboard, -Fare)
- > t2 <- select(titanic, -c(Age:Fare))

The Fundamental Functions

select Select columns

The Fundamental Functions

select Select columns

filter Filter rows

filter - filter the rows

filter(data, <Boolean condition>,...)

Condition on the variables

filter - filter the rows

filter(data, <Boolean condition>,...)

Condition on the variables

< Less than >= Greater than or equal to

<= Less than or equal to == Equal to</pre>

> Greater than != Not equal to

filter - Keep specific value

filter(data, Major == "math")

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

Keep rows where Major is "math"

Name	Major	HW
Matt	math	60
Mary	math	90

filter - Keep range of values

filter(data, HW <= 85)

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

Keep rows where `HW` is less than or equal to 85

Name	Major	HW
Matt	math	60
Bill	biology	85

filter - Multiple Conditions

filter(data, HW <= 85, Major == "math")

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

Keep rows where `HW` is less than or equal to 85, **AND** `Major` is math

Name	Major	HW
Matt	math	60

filter - Multiple Conditions



filter(data, HW <= 85 & Major == "math")

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

Keep rows where `HW` is less than or equal to 85, **AND** `Major` is math

Name	Major	HW
Matt	math	60

filter - Multiple Conditions



filter(data, HW <= 85 | Major == "biology")

Name	Major	HW
Matt	math	60
Mary	math	90
Bill	biology	85
Brie	biology	87

Keep rows where `HW` is less than or equal to 85, **OR** `Major` is biology

Name	Major	HW
Matt	math	60
Bill	biology	85
Brie	biology	87

filter - Titanic Exercises

1. Create a new dataset with only those who survived (`Survived` value of 1)

2. Only show rows with 1st class ('Pclass') females ('Sex')

3. How many rows are there if you keep only males between ages 20-30?

filter - Titanic Exercise Solutions

- 1. Create a new dataset with only those who survived (`Survived` value of 1)
- > survived <- filter(titanic, Survived == 1)
- > survived <- filter(titanic, Survived == "1")</pre>

filter - Titanic Exercise Solutions

- 1. Create a new dataset with only those who survived (`Survived` value of 1)
- > survived <- filter(titanic, Survived == 1)
- > survived <- filter(titanic, Survived == "1")
- 2. Only show rows with 1st class ('Pclass') females ('Sex')
- > filter(titanic, Pclass == 1 & Sex == "female")
- > filter(titanic, Pclass == 1, Sex == "female")

filter - Titanic Exercise Solutions

3. How many rows are there if you keep only males between ages 20-30?

- > tm <- filter(titanic, Sex == "male", Age >= 20, Age <= 30)
- > nrow(tm) # 223
- > t2 <- filter(titanic, Sex == "male")
- > t3 <- filter(t2, Age >= 20)
- > t4 <- filter(t3, Age <= 30)

The Fundamental Functions

select Select columns

filter Filter rows

Pipes

%>%



Pipes (%>%) insert the data into the next line







2 %>% sqrt()

Without pipes, multi-line processing is hard to read

df1 <- select(df, A, B, C)

df2 <- filter(df1, A == 1, B > 2)

df3 <- filter(df2, C == 1 | C == 5)

Creates unnecessary intermediate datasets

df1 <- select(df, A, B, C)

1, B > 2), C == 1 | C == 5)

Without pipes, multi-line processing is hard to read

```
df2 <- filter(df1, A == 1, B > 2)
df3 <- filter(df2, C == 1 | C == 5)

df1 <- filter(filter(select(df, A, B, C), A ==
```

Hard to read



With pipes, one line is one logical step

df1 <- df %>%

Start with dataset `df`

With pipes, one line is one logical step

df1 <- df %>%

select(A, B, C) %>%

Start with dataset `df`

Select variables A, B, C

With pipes, one line is one logical step

df1 <- df %>%

select(A, B, C) %>%

filter(A == 1, B > 2) %

Start with dataset `df`

Select variables A, B, C

Keep rows with A == 1 and B > 2

With pipes, one line is one logical step

df1 <- df %>%

select(A, B, C) %>%

filter(A == 1, B > 2) % > %

filter(C == 1 | C == 5)

Start with dataset `df`

Select variables A, B, C

Keep rows with A == 1 and B > 2

Keep rows with C == 1 or C == 5

(Store results in df1)



Pipe %>% - Titanic Exercise

1. Using pipes, store the first four columns and keep only 1st class (`Pclass`)

- > firstclass <- _____ %>%
- > select(1_4) ____
- > filter(Pclass ___ ___)

Pipe %>% - Titanic Exercise Solution

- 1. Using pipes, store the first four columns and keep only 1st class (`Pclass`)
- > firstclass <- titanic %>%
- > select(1:4) %>%
- > filter(Pclass == 1)

The Fundamental Functions

select Select columns

filter Filter rows

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

mutate - Create new variables/columns

Computation of new variable

Can be based on another variable, a function, a single value, etc.

mutate - Create var with constant value

data %>% mutate(Sect = 1)

Name	HW	Ex
Matt	60	75
Mary	90	100
Bill	85	60
Brie	87	85

Create new variable `Sect` where all values are 1

Name	HW	Ex	Sect
Matt	60	75	1
Mary	90	100	1
Bill	85	60	1
Brie	87	85	1

mutate - Create var based on another var

data %>% mutate(Ex2 = Ex/100)

Name	HW	Ex
Matt	60	75
Mary	90	100
Bill	85	60
Brie	87	85

Create new variable `Ex2` as `Ex` divided by 100

Name	HW	Ex	Ex2
Matt	60	75	0.75
Mary	90	100	1.00
Bill	85	60	0.60
Brie	87	85	0.85

mutate - Create var based on other variables

data %>% mutate(Final = (HW + Ex)/2)

Name	HW	Ex
Matt	60	75
Mary	90	100
Bill	85	60
Brie	87	85

Create new variable `Final` as average of `HW` and `Ex`

Name	HW	Ex	Final
Matt	60	75	67.5
Mary	90	100	95
Bill	85	60	72.5
Brie	87	85	86

mutate - Titanic Exercises

1. `Fare` is currently in pounds. Convert it to dollars (1 pound = 1.37 dollars)

2. How many family members total (combine `Siblings.Spouses.Aboard` and `Parents.Children.Aboard`) were on board for each passenger?

3. Using pipes: remove the name, keep only 1st class survivors, subtract 20 from `Age`. Store these results in a dataset named `AgeAdjusted`.

mutate - Titanic Exercise Solutions

- 1. `Fare` is currently in pounds. Convert it to dollars (1 pound = 1.37 dollars)
- > titanic %>% mutate(FareDollars = Fare*1.37)
- 2. How many family members total (combine `Siblings.Spouses.Aboard` and `Parents.Children.Aboard`) were on board for each passenger?
- > titanic %>% mutate(Family = Siblings.Spouses.Aboard
- + Parents.Children.Aboard)

mutate - Titanic Exercise Solutions

- 3. Using pipes: remove the name, keep only 1st class survivors, subtract 20 from `Age`. Store these results in a dataset named `AgeAdjusted`.
- > AgeAdjusted <- titanic %>%
- > select(-Name) %>%
- > filter(Pclass == 1, Survived == 1) %>%
- > mutate(Age = Age 20)

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

summarize Compute summary variable

summarize - Create new aggregated variables

summarize(., new_variable = function(_), ...)

Computation of new variable

Function that returns one value (e.g. mean(), median(), max())

summarize - Create new aggregated variables

summarize() computes a single value or column

mutate() adds a variable to existing dataset

summarize - Compute mean of dataset

data %>% summarize(HW_mean = mean(HW))

Name	HW
Matt	60
Mary	90
Bill	85
Brie	87

Calculate mean HW score



summarize - Count total number of observations

data %>% summarize(Count = n())

Name	HW
Matt	60
Mary	90
Bill	85
Brie	87

Calculate number of observations

Count

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

summarize Compute summary variable

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

summarize Compute summary variable

group_by Split data by groups

group_by(., `id_var`, ...)

Split dataset by `id_var` (or multiple variables)

Dataset split based on unique values of 'id_var'

data %>% group_by(., Major)

Major	HW
math	60
math	90
biology	85
biology	87

data %>% group_by(., Major)

Major	HW	
math	60	
math	90	
biology	85	
biology	87	

Major	HW
math	60
math	90

Major	HW
biology	85
biology	87

data %>% group_by(., Major)

Major	HW	
math	60	
math	90	
biology	85	
biology	87	

Major	HW	
math	60	
math	90	

Major	HW	
biology	85	
biology	87	



Major	HW	
math	60	
math	90	
biology	85	
biology	87	

group_by(., Major)

Major	HW
math	60
math	90
biology	85
biology	87

Major	HW
math	60
math	90

Major	HW	
biology	85	
biology	87	

group_by(., Major) %>% summarize(hw_m = mean(HW))

Major	HW
math	60
math	90
biology	85
biology	87

Major	HW	
math	60	
math	90	

Major	HW
biology	85
biology	87



Major	hw_m
math	75



Major	hw_	_m
biology		86

Major	hw_m
math	75
biology	86

group_by(., Major)

Major	HW
math	60
math	90
biology	85
biology	87

Major	HW
math	60
math	90

Major	HW
biology	85
biology	87

Take top HW value

group_by(., Major) %>% top_n(1, HW)

Major	HW
math	60
math	90
biology	85
biology	87

Major	HW
math	60
math	90

HW

85

87

Major

biology

biology

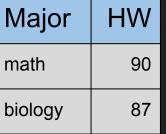


Major	HW
math	90



N
87

Major	HW
biology	87



summarize() and group_by() - Titanic Exercises

- 1. What was the average `Fare` cost?
- > titanic %>% summarize(FareMean = mean(_____))

- 2. What was the average 'Age' for each 'Pclass'? Store the results in 'ClassAge'.
- > ClassAge <- titanic %>%
- 3. What percent of each `Pclass` survived? (Hint: sum() and n() are helpful).

summarize() and group_by() - Titanic Exercise Solutions

- 1. What was the average `Fare` cost?
- > titanic %>% summarize(FareMean = mean(Fare))

summarize() and group_by() - Titanic Exercise Solutions

- 1. What was the average `Fare` cost?
- > titanic %>% summarize(FareMean = mean(Fare))

- 2. What was the average 'Age' for each 'Pclass'? Store the results in 'ClassAge'.
- > ClassAge <- titanic %>%
- > group_by(Pclass) %>%
- > summarize(AgeM = mean(Age))

summarize() and group_by() - Titanic Exercises

- 3. What percent of each 'Pclass' survived? (Hint: sum() and n() are helpful).
- > titanic %>%
- > group_by(Pclass) %>%
- > summarize(pct = sum(Survived)/n())

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

summarize Compute summary variable

group_by Split data by groups

The Fundamental Functions

select Select columns

filter Filter rows

mutate Create new variables

summarize Compute summary variable

group_by Split data by groups

left join Combine datasets

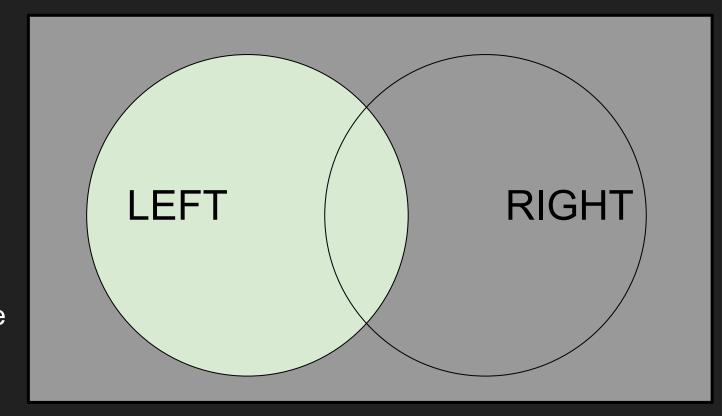
left_join(df_left, df_right, ...)

The two datasets to combine together

Join two tables, matching on a key column

All values in the left dataset are kept

Only matching values in the right dataset are kept



Student <- left_join(Majors, Grades)

Ma	Majors		Grad	Grades			Student		
Name	Major		Name	HW		Name	Major	¥	
Matt	math	_	Matt	60		Matt	math	60	
Mary	math	+	Mary	90		Mary	math	90	
Bill	biology		Bill	85	r	Bill	biology	85	
Brie	biology		Brie	87		Brie	biology	87	

Student <- left_join(Majors, Grades)

Majors Grades Student Major Major HW Name HW Name Name Matt 60 Matt math Matt 60 math Mary 90 Mary math Mary 90 math Bill biology Bill NA biology Brie biology biology Brie NA

Student <- left_join(Majors, Grades)

Majors

Name	Major
Matt	math
Mary	math

Grades

Name	HW
Matt	60
Mary	90
Bill	85
Brie	87

Student

Name	Major	HW
Matt	math	60
Mary	math	90

If matching columns have different names

Student <- left_join(Grades, Majors, by = c("NX" = "NY")

Majors		Grades			Student			
NX	Major		NY	HW		NX	Major	HW
Matt	math		Matt	60		Matt	math	60
Mary	math	+	Mary	90		Mary	math	90
Bill	biology		Bill	85	,	Bill	biology	85
Brie	biology		Brie	87		Brie	biology	87

Other joins

Function	Description
left_join	Keep all left
right_join	Keep all right
inner_join	Match left and right
full_join	Keep all values
nest_join	Create list with match
anti_join	Keep non-matching

left_join() - Titanic Exercise

- 1. Split the titanic dataset accordingly:
- > titanic1 <- titanic %>% select(1:3)
- > titanic2 <- titanic %>% select(3:5)
- 2. Combine the datasets back together based on 'Name'
- > titanic3 <- left_join(____, ____)

left_join() - Titanic Exercise

- 1. Split the titanic dataset accordingly:
- > titanic1 <- titanic %>% select(1:3)
- > titanic2 <- titanic %>% select(3:5)
- 2. Combine the datasets back together based on 'Name'
- > titanic3 <- left_join(titanic1, titanic2)</pre>

Exporting Data

Name in R

write_csv(dataset, file = "dataset.csv")

Output name and path to location

NOTE: MUST INCLUDE .CSV

Exporting Data - Titanic Exercises

- 1. Store the first four columns in 'TitanicReduced'
- > TitanicReduced <- titanic %>% select(____)

- 2. Export the previously created dataset
- > write_csv(_____, "RTitanic.csv")

Exporting Data - Titanic Exercise Answers

- 1. Store the first four columns in 'TitanicReduced'
- > TitanicReduced <- titanic %>% select(1:4)

- 2. Export the previously created dataset
- > write_csv(TitanicReduced, "RTitanic.csv")

Other Tips

Working with dates/times? Use the `lubridate` package!

Skim through the data wrangling cheat sheet to see what is possible.

Do not reinvent the wheel. Use google before you try to code something fancy.

Have fun!