# 1. Introduction to the Tidyverse
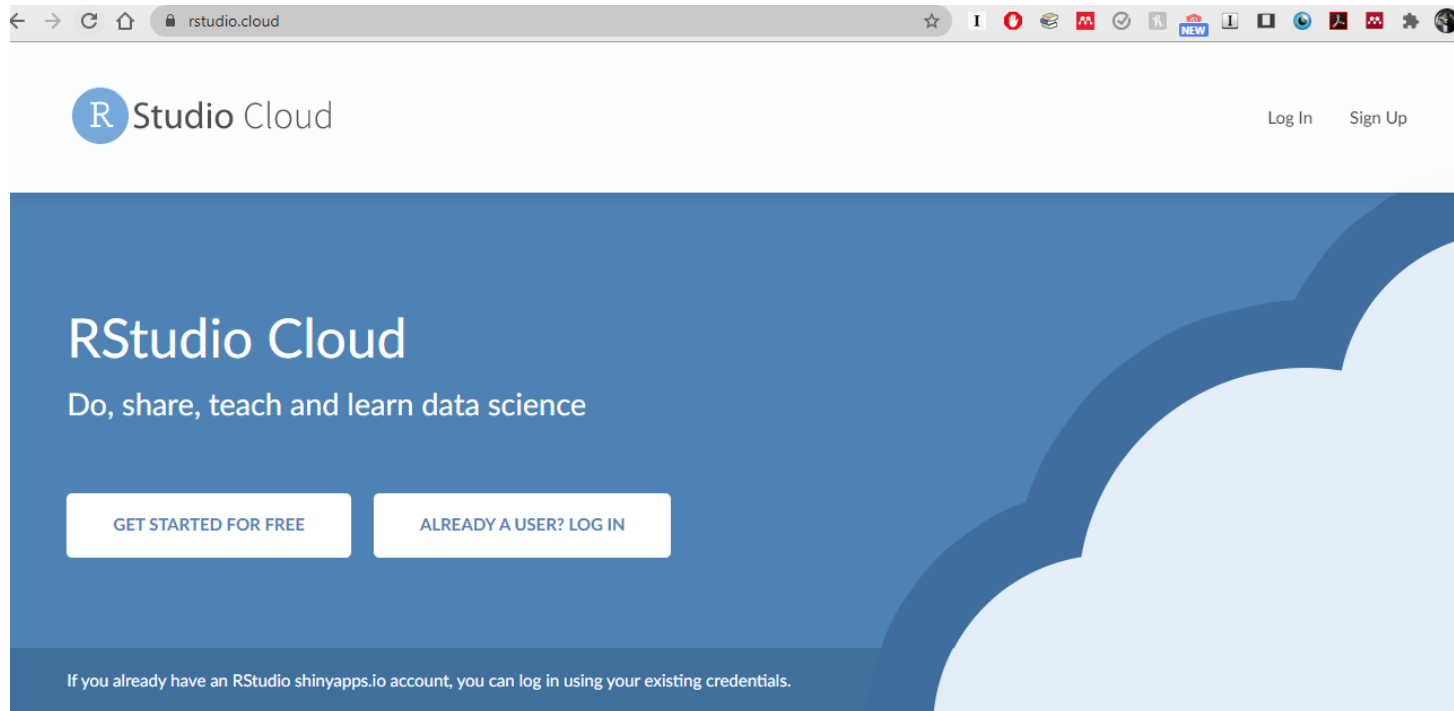
Jonathan Hersh, PhD (Chapman Argyros School of Business)

12/7/20

# Section 1: Outline

- Loading data
- Glimpse to view
- Pipe operator
- slice() to select rows
- arrange() to order data frame
- select() to choose variables
- rename() to rename variables
- filter() to select rows matching characteristics
- mutate() to create new variables
- group_by() and summarize()  to create group

# R Studio Cloud



- Go to rstudio.cloud if your version of R is ever not working

# R Studio Cloud



- R Studio Cloud is a full featured version of R in your browser!

# R Studio Projects



Be sure to work in a project file in R Studio

# Creating an R Studio Project



- Be sure to work in a project file in R Studio

- Click file -> new project -> Existing Directory

# Creating an R Studio Project



- Be sure to work in a project file in R Studio

- Click file -> new project -> Existing Directory

- Then navigate to the folder where you are storing all of today's files

# Reference Guides for Data Transformation, Import, Plotting and RStudio

# Data Import :: **CHEAT SHEET**

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.

The front side of this sheet shows how to read text files into R with **readr**.

The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

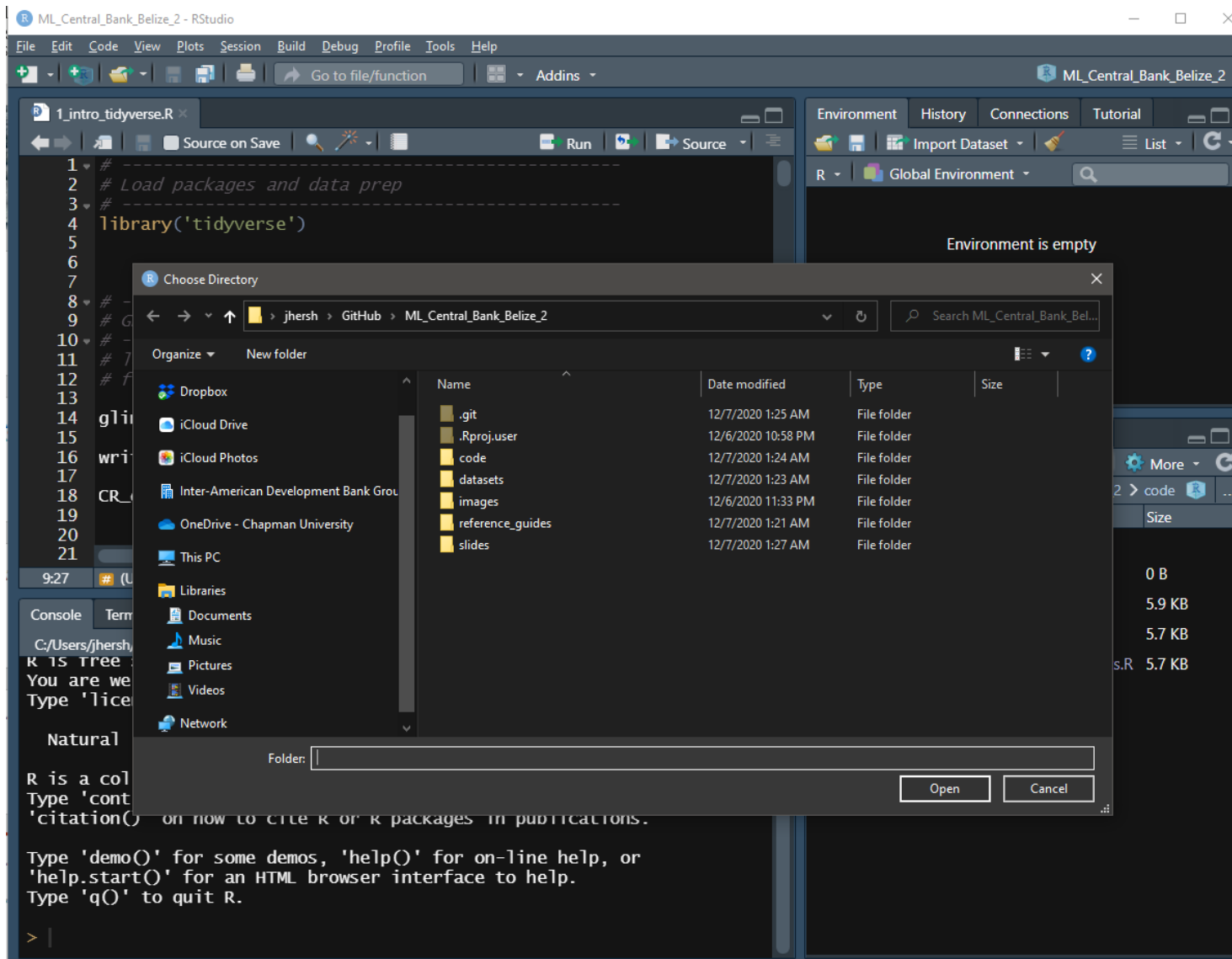**OTHER TYPES OF DATA**
Try one of the following packages to import other types of files
- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (.xls and .xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

## Save Data

Save **x**, an R object, to **path**, a file path, as:

**Comma delimited file**
  **write_csv**(x, path, na = "NA", append = FALSE, col_names = !append)

**File with arbitrary delimiter**
  **write_delim**(x, path, delim = " ", na = "NA", append = FALSE, col_names = !append)

**CSV for excel**
  **write_excel_csv**(x, path, na = "NA", append = FALSE, col_names = !append)

**String to file**
  **write_file**(x, path, append = FALSE)

**String vector to file, one element per line**
  **write_lines**(x,path, na = "NA", append = FALSE)

**Object to RDS file**
  **write_rds**(x, path, compress = c("none", "gz", "bz2", "xz"), ...)

**Tab delimited files**
  **write_tsv**(x, path, na = "NA", append = FALSE, col_names = !append)

## Read Tabular Data - These functions share the common arguments:

read_*(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"), quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000, n_max), progress = interactive())

**Comma Delimited Files**
  **read_csv**("file.csv")
    To make file.csv run:
    write_file(x = "a,b,c\n1,2,3\n4,5,NA", path = "file.csv")

**Semi-colon Delimited Files**
  **read_csv2**("file2.csv")
    write_file(x = "a;b;c\n1;2;3\n4;5;NA", path = "file2.csv")

**Files with Any Delimiter**
  **read_delim**("file.txt", delim = "|")
    write_file(x = "a|b|c\n1|2|3\n4|5|NA", path = "file.txt")

**Fixed Width Files**
  **read_fwf**("file.fwf", col_positions = c(1, 3, 5))
    write_file(x = "a b c\n1 2 3\n4 5 NA", path = "file.fwf")

**Tab Delimited Files**
  **read_tsv**("file.tsv") Also **read_table()**.
    write_file(x = "a\tb\tc\n1\t2\t3\n4\t5\tNA", path = "file.tsv")

### USEFUL ARGUMENTS

**Example file**
  write_file("a,b,c\n1,2,3\n4,5,NA","file.csv")
  f <- "file.csv"

**No header**
  read_csv(f, **col_names = FALSE**)

**Provide header**
  read_csv(f, **col_names = c("x", "y", "z")**)

**Skip lines**
  read_csv(f, **skip = 1**)

**Read in a subset**
  read_csv(f, **n_max = 1**)

**Missing Values**
  read_csv(f, **na = c("1", ".")**)

## Read Non-Tabular Data

**Read a file into a single string**
  **read_file**(file, locale = default_locale())

**Read each line into its own string**
  **read_lines**(file, skip = 0, n_max = -1L, na = character(), locale = default_locale(), progress = interactive())

**Read Apache style log files**
  **read_log**(file, col_names = FALSE, col_types = NULL, skip = 0, n_max = -1, progress = interactive())

**Read a file into a raw vector**
  **read_file_raw**(file)

**Read each line into a raw vector**
  **read_lines_raw**(file, skip = 0, n_max = -1L, progress = interactive())

## Data types

readr functions guess the types of each column and convert types when appropriate (but will NOT convert strings to factors automatically).

A message shows the type of each column in the result.

```
## Parsed with column specification:
## cols(
##   age = col_integer(),
##   sex = col_character(),
##   earn = col_double()
## )
```

age is an integer

sex is a character

earn is a double (numeric)

1. Use **problems()** to diagnose problems.
   *x <- read_csv("file.csv"); problems(x)*

2. Use a col_ function to guide parsing.
   - **col_guess()** - the default
   - **col_character()**
   - **col_double()**, **col_euro_double()**
   - **col_datetime**(format = "") Also **col_date**(format = ""), **col_time**(format = "")
   - **col_factor**(levels, ordered = FALSE)
   - **col_integer()**
   - **col_logical()**
   - **col_number()**, **col_numeric()**
   - **col_skip()**

   *x <- read_csv("file.csv", col_types = cols(*
     *A = col_double(),*
     *B = col_logical(),*
     *C = col_factor()))*

3. Else, read in as character vectors then parse with a parse_ function.
   - **parse_guess()**
   - **parse_character()**
   - **parse_datetime()** Also **parse_date()** and **parse_time()**
   - **parse_double()**
   - **parse_factor()**
   - **parse_integer()**
   - **parse_logical()**
   - **parse_number()**

   *x$A <- parse_number(x$A)*

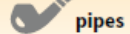# Data Transformation with dplyr :: CHEAT SHEET

**dplyr** functions work with pipes and expect **tidy data**. In tidy data:

Each **variable** is in its own **column**

Each **observation**, or **case**, is in its own **row**

**pipes**
x %>% f(y) becomes f(x, y)

## Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

**summary function**

**summarise**(.data, …)
Compute table of summaries.
*summarise(mtcars, avg = mean(mpg))*

**count**(x, …, wt = NULL, sort = FALSE)
Count number of rows in each group defined by the variables in … Also **tally**().
*count(iris, Species)*

### VARIATIONS

**summarise_all()** - Apply funs to every column.
**summarise_at()** - Apply funs to specific columns.
**summarise_if()** - Apply funs to all cols of one type.

## Group Cases

Use **group_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.

mtcars %>%

group_by(cyl) %>%
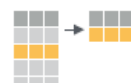
summarise(avg = mean(mpg))

**group_by**(.data, …, add = FALSE)
Returns copy of table grouped by …
*g_iris <- group_by(iris, Species)*

**ungroup**(x, …)
Returns ungrouped copy of table.
*ungroup(g_iris)*
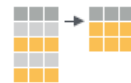
## Manipulate Cases

### EXTRACT CASES

Row functions return a subset of rows as a new table.

**filter**(.data, …) Extract rows that meet logical criteria. *filter(iris, Sepal.Length > 7)*
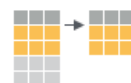
**distinct**(.data, …, .keep_all = FALSE) Remove rows with duplicate values. *distinct(iris, Species)*

**sample_frac**(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select fraction of rows. *sample_frac(iris, 0.5, replace = TRUE)*

**sample_n**(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select size rows. *sample_n(iris, 10, replace = TRUE)*

**slice**(.data, …) Select rows by position. *slice(iris, 10:15)*

**top_n**(x, n, wt) Select and order top n entries (by group if grouped data). *top_n(iris, 5, Sepal.Width)*
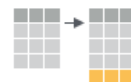
**Logical and boolean operators to use with filter()**

| | | | |
|---|---|---|---|
| < | is.na() | %in% | \| | xor() |
| > | >= | !is.na() | ! | & |

See **?base::Logic** and **?Comparison** for help.

### ARRANGE CASES

**arrange**(.data, …) Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low.
arrange(mtcars, mpg)
arrange(mtcars, desc(mpg))

### ADD CASES

**add_row**(.data, …, .before = NULL, .after = NULL) Add one or more rows to a table.
*add_row(faithful, eruptions = 1, waiting = 1)*

## Manipulate Variables

### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

**pull**(.data, var = -1) Extract column values as a vector. Choose by name or index.
*pull(iris, Sepal.Length)*

**select**(.data, …)
Extract columns as a table. Also **select_if()**.
*select(iris, Sepal.Length, Species)*

**Use these helpers with select (),**
*e.g. select(iris, starts_with("Sepal"))*

**contains**(match)
**ends_with**(match)
**matches**(match)
**num_range**(prefix, range)
**one_of**(…)
**starts_with**(match)
:, e.g. mpg:cyl
-, e.g, -Species

### MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).

**vectorized function**

**mutate**(.data, …)
Compute new column(s).
*mutate(mtcars, gpm = 1/mpg)*

**transmute**(.data, …)
Compute new column(s), drop others.
*transmute(mtcars, gpm = 1/mpg)*

**mutate_all**(.tbl, .funs, …) Apply funs to every column. Use with **funs()**. Also **mutate_if()**.
*mutate_all(faithful, funs(log(.), log2(.)))*
*mutate_if(iris, is.numeric, funs(log(.)))*

**mutate_at**(.tbl, .cols, .funs, …) Apply funs to specific columns. Use with **funs()**, **vars()** and the helper functions for select().
*mutate_at(iris, vars( -Species), funs(log(.)))*

**add_column**(.data, …, .before = NULL, .after = NULL) Add new column(s). Also **add_count()**, **add_tally()**. *add_column(mtcars, new = 1:32)*

**rename**(.data, …) Rename columns.
*rename(iris, Length = Sepal.Length)*

RStudio

# RStudio IDE :: **CHEAT SHEET**

## Documents and Apps

Open Shiny, R Markdown, knitr, Sweave, LaTeX, .Rd files and more in Source Pane

Check spelling
Render output
Choose output format
Choose output location
Insert code chunk

Jump to previous chunk
Jump to next chunk
Run selected lines
Publish to server
Show file outline

Access markdown guide at
**Help > Markdown Quick Reference**

Jump to chunk
Set knitr chunk options
Run this and all previous code chunks
Run this code chunk

```
17 - ```{r pressure, echo=FALSE}
18   plot(pressure)
19 - ```
```

RStudio recognizes that files named **app.R**, **server.R**, **ui.R**, and **global.R** belong to a shiny app

Run app
Choose location to view app
Publish to shinyapps.io or server
Manage publish accounts

## Write Code

Navigate tabs
Open in new window
Save
Find and replace
Compile as notebook
Run selected code
**Import data** with wizard
History of past commands to run/copy
Display .RPres slideshows
**File > New File > R Presentation**

```
1  # Good Start...
2  
3  Cursors of shared users
4  
5  
6  "P0030001"
7  "P0030002"
8  "P0030003"
9  "P0030004"
10 
11 
12 get_digit <- function() {
13   ("num" %% (10 ^ n))
14   %/% (10 ^ (n - 1))
15 }}
16 
17 fo
18   for        {snippet}
19   foo    {.GlobalEnv}
20   force      {base}
21 Jump to function in file
22 
```

Re-run previous code
Source with or without Echo
Show file outline

Multiple cursors/column selection with **Alt + mouse drag**.

Code diagnostics that appear in the margin. Hover over diagnostic symbols for details.

Syntax highlighting based on your file's extension

Tab completion to finish function names, file paths, arguments, and more.

Multi-language code snippets to quickly use common blocks of code.

Change file type

```
Console   Compile PDF   R Markdown
~/IDEcheatsheet/
> foo(1)
[1]_2
> foo <- function(x) x + 1
> foo(2)
  foo(2)
> foo(1)
```

Working Directory
Maximize, minimize panes
Press ↑ to see command history
Drag pane boundaries

## R Support

Load workspace
Save workspace
Delete all saved objects
Search inside environment

Choose environment to display from list of parent environments
Display objects as list or grid

```
Data
 iris           150 obs. of 5 variables
Values
 a              1
Functions
 foo            function (x)
```

Displays saved objects by type with short description
View in data viewer
View function source code

Files   Plots   Packages   Help   Viewer
New Folder   Upload   Delete   Rename   More
Home   IDEcheatsheet
Name
Copy...
Move...
Export...
Set As Working Directory
Go To Working Directory

Create folder
Upload file
Delete file
Rename file
Change directory

Path to displayed directory

hello.R          19 B          Apr 13, 2016, 11:17 AM

A File browser keyed to your working directory. Click on file or directory name to open.

## Pro Features

**Share Project** with Collaborators
Active shared collaborators

Start **new R Session** in current project
Close R Session in project
**Select R Version**

New Project...
Open Project...
Close Project
Share Project...
IDEcheatsheet
RStudio-Essentials
Essentials
shiny-examples
Clear Project List
Project Options...

R version 3.2.2
R version 3.1.3
R version 3.0.3
R version 2.15.3

### PROJECT SYSTEM
**File > New Project**

Name of current project

RStudio saves the call history, workspace, and working directory associated with a project. It reloads each when you re-open a project.

RStudio opens plots in a dedicated Plots pane

Navigate recent plots
Open in window
**Export plot**
Delete plot
Delete all plots

GUI Package manager lists every installed package

Install Packages
Update Packages
Create reproducible package library for your project

scales          Scale Functions for Visualization          0.3.0
shiny           Web Application Framework for R            0.12.2
shinydashboard  Create Dashboards with 'Shiny'            0.5.1

Click to load package with **library()**. Unclick to detach package with **detach()**
Package version installed
Delete from library

RStudio opens documentation in a dedicated Help pane

Home page of helpful links
Search within help file
Search for help file

Viewer Pane displays HTML content, such as Shiny apps, RMarkdown reports, and interactive visualizations

Stop Shiny app
Publish to shinyapps.io, rpubs, RSConnect, …
Refresh

**View(<data>)** opens spreadsheet like view of data set

| | Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|---|---|---|---|---|---|
| | All | All | All | All | All |
| 1 | 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 2 | | | | | |
| 3 | Filter rows by value or value range | | Sort by values | Search for value | |

## Debug Mode

Open with **debug()**, **browser()**, or a breakpoint. RStudio will open the debugger mode when it encounters a breakpoint while executing code.

Launch debugger mode from origin of error
Open traceback to examine the functions that R called before the error occurred

Click next to line number to add/remove a breakpoint.

Highlighted line shows where execution has paused

```
Console   ~/IDEcheatsheet/
> foo()
Error in get_digit(num, x) :
Error!
```
Show Traceback
Rerun with Debug

```
Console   ~/IDEcheatsheet/
  Next    Continue   Stop
```

Run commands in environment where execution has paused
Examine variables in executing environment
Select function in traceback to debug
Step through code one line at a time
Step into and out of functions to run
Resume execution
Quit debug mode

## Version Control with Git or SVN

Turn on at **Tools > Project Options > Git/SVN**

Stage files:
Show file diff
Commit staged files
Push/Pull to remote
View History

A Added
D Deleted
M Modified
R Renamed
? Untracked

Open shell to type commands
current branch

## Package Writing

**File > New Project > New Directory > R Package**

Turn project into package,
Enable roxygen documentation with
**Tools > Project Options > Build Tools**

Roxygen guide at
**Help > Roxygen Quick Reference**