

AGENDA

CHOPR

Created by: Matthew Dye, Emily Foster, Brendan Graham,
Joe Mirizio, Ashley Oliver, Jake Riley, Emily Schriver, Siying Wang,
Paul Wildenhain

Date

6/7/2018 – 6/8/2018

Day 1: 6/7/2018

		<u>Rooms</u>
9:00 AM – 9:30 AM	Standard/ Set-up R standards checklist R project overview	8110/7105
9:30 AM – 11:00 AM	dplyr	8110/7105
11:00 AM – 12:00 PM	lubridate	8110/7105
12:00 PM – 1:00 PM	Lunch	8110/7105
1:00 PM – 3:00 PM	ggplot2 / rocqi	8110/8105
3:00 PM – 3:30 PM	Break	8110/8105
3:30 PM – 4:30 PM	tidyr	8110/8105

Day 2: 6/8/2018

9:00 AM – 9:30 AM	Chopr overview Kernel orientation Release to R Studio Connect	4110/4105
9:30 AM – 12:00 PM	Build Session 1	4110/4105
12:00 PM – 1:00 PM	Lunch	4110/4105
1:00 PM – 3:30 PM	Build Session 2	4110/4105
3:30 PM – 4:30 PM	Presentation	4110/4105

Additional Instructions:

- 1) Bring your data example!
- 2) Please complete the feedback survey found at <https://goo.gl/forms/ISTjXLdYPZrUVzD93>

R Adoption: Master Exercise Sheet

Set-up Session:

1. Create a directory and R project on your H drive that you will use throughout the chopr sessions
2. Create a folder in that directory called 'data'
3. Copy and paste the blood culture dataset from the Q drive to your 'data' folder

Dplyr Session:

1. Filter blood culture data to only look at blood cultures ordered in the NICU and store as a new dataframe
2. Limit the blood culture data set to only patient information, blood culture type, blood culture department, and blood culture date and store as a new dataframe

Create new variables for each of #3-#6 and save in a new dataframe:

3. Convert blood culture date and date of birth into date format
4. Calculate patient age when the blood culture was done
5. Flag patients who were under 90 days when the blood culture was done
6. Label the order that blood cultures were drawn for each patient (hint: you might need to use "arrange" to specify the order that the blood cultures are in)
7. In the blood culture data, determine the count for each line type (central vs non-central line) and store as a new dataframe
8. BONUS: combine the above 7 exercises in ONE statement using pipes (%>%)

Lubridate Session:

1. Convert BLD_CX_DT from a factor to a date variable type
2. Convert DOB from a factor to a date variable type
3. Extract the month from the BLD_CX_DT date field you created in exercise #1
4. Create a data frame that counts the number of unique blood cultures each calendar month (regardless of year)
5. Create a variable that is the first day of the month shown in BLD_CX_DT
6. Create a dataframe that counts the number of blood cultures each month-year

Create a variable that calculates the days between birth and blood culture

7. Determine the mean and median days between birth and blood culture
8. Use dplyr and format to create an indicator that describes whether or not the blood culture was taken after 12 PM

Ggplot Session:

- 1) Create a barplot of number of central vs. non central line blood cultures in FY17
- 2) Create a stacked barplot by patients under 90 days and patients over 90 days for the graph you just created (central vs non-central lines)
- 3) Create a side-by-side barplot of graph #2
- 4) Create a line graph of counts of blood cultures over time (months)

Rocqi Session:

1. Connect to CDW UAT
2. Using dbGetQuery(), run the sql in blood_culture.sql and assign to a dataframe called "blood(cx)_cohort"
3. Create a p chart that shows the percent of blood cultures drawn before noon per week
 1. Hint: use dplyr to create some necessary additional variables and flags
4. Make a run chart with spc() that shows the percent of blood cultures drawn before noon by month, stratified by culture source.
 1. Hint: the facet argument may be useful for stratification
5. Add the OCQI theme and colors to the graph
 1. Note: You can add layers with "+" after spc() just like with other ggplot graphs.

BONUS: Perform the following operations:

6. Run the data_summary() function on the blood(cx)_cohort data frame
7. Using the network_drive() and paste() functions, read the blood_culture.csv into R as a data frame
8. Make a run chart that shows the number of cultures per week. Then use the ocqi_colors() function to make that line blue

Tidyr Session:

1. Find the total number of central line uses by culture source by department
2. From this, create a wide-form dataset having 1 row per department and a column for each culture source to show the totals, use a zero for any NAs
3. Pivot the data back into a long-form dataset that has the same # of rows and columns as step 1

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 

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 

Debug Mode

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

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

Highlighted line shows where execution has paused

Run commands in environment where execution has paused

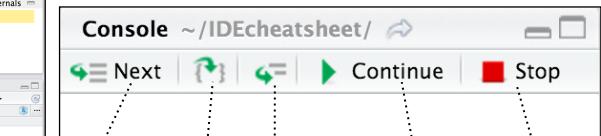
Examine variables in executing environment

Select function in traceback to debug

Launch debugger mode from origin of error

Open traceback to examine the functions that R called before the error occurred

 Error in get_digit(num, x): ...
Error!

 Step through code one line at a time

Version Control with Git or SVN



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

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

A Added  D Deleted  M Modified  R Renamed  ? Untracked  Environment  History  Git  master  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**

 Environment  History  Build  Reload  Check  More  Load All  Clean and Rebuild  Test Package  Check Package  Build Source Package  Build Binary Package Document Configure Build Tools...



Pro Features

Share Project Active shared with Collaborators... 

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

PROJECT SYSTEM File > New 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

Files  Plots  Packages  Help  Viewer  Navigate  Open in recent plots  Export  Delete plot  Delete all plots

GUI Package manager lists every installed package

Files  Plots  Packages  Help  Viewer  Install  Packages  Update  Create reproducible package library for your project  scales 0.3.0 shiny 0.12.2 shinydashboard 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

Files  Plots  Packages  Help  Viewer  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

Files  Plots  Packages  Help  Viewer  Stop Shiny app  Publish to shinyapps.io, rpubs, RSConnect, ... Refresh 

View(<data>) opens spreadsheet like view of data set  Filter  Sepal.Length  Sepal.Width  Petal.Length  Petal.Width  Species  All  1 5.1 3.5 1.4 0.2 setosa 2 3 4 Filter rows by value or value range Sort by values Search for value

1 LAYOUT

Move focus to Source Editor
Move focus to Console
Move focus to Help
Show History
Show Files
Show Plots
Show Packages
Show Environment
Show Git/SVN
Show Build

Windows/Linux Mac

Ctrl+1
Ctrl+2
Ctrl+3
Ctrl+4
Ctrl+5
Ctrl+6
Ctrl+7
Ctrl+8
Ctrl+9
Ctrl+0

2 RUN CODE

Search command history

Navigate command history
Move cursor to start of line
Move cursor to end of line
Change working directory

Interrupt current command

Clear console

Quit Session (desktop only)

Restart R Session

Run current (retain cursor)
Run from current to end
Run the current function
Source a file

Source the current file

Source with echo

Windows/Linux Mac

Ctrl+↑
↑/↓
Home
End
Ctrl+Shift+H

Esc**Ctrl+L**

Ctrl+Q

Ctrl+Shift+F10**Ctrl+Enter**

Alt+Enter
Ctrl+Alt+E
Ctrl+Alt+F
Ctrl+Alt+G

Ctrl+Shift+S

Ctrl+Shift+Enter

3 NAVIGATE CODE

Goto File/Function

Fold Selected
Unfold Selected
Fold All
Unfold All
Go to line
Jump to
Switch to tab
Previous tab
Next tab
First tab
Last tab
Navigate back
Navigate forward
Jump to Brace
Select within Braces
Use Selection for Find
Find in Files
Find Next
Find Previous
Jump to Word
Jump to Start/End
Toggle Outline

Windows /Linux

Mac

Ctrl+.
Alt+L
Shift+Alt+L
Alt+O
Shift+Alt+O
Shift+Alt+G
Shift+Alt+J
Ctrl+Shift+.
Ctrl+F11
Ctrl+F12
Ctrl+Shift+F11
Ctrl+Shift+F12
Ctrl+F9
Ctrl+F10
Ctrl+P
Ctrl+Shift+Alt+E
Ctrl+F3
Ctrl+Shift+F
Win: F3, Linux: Ctrl+G
Cmd+G
Cmd+Shift+G
Option+←/→
Ctrl+↑/↓
Ctrl+Shift+O

4 WRITE CODE

Attempt completion
Navigate candidates
Accept candidate
Dismiss candidates
Undo
Redo
Cut
Copy
Paste
Select All
Delete Line

Select

Select Word
Select to Line Start
Select to Line End
Select Page Up/Down
Select to Start/End
Delete Word Left
Delete Word Right
Delete to Line End
Delete to Line Start

Indent

Outdent
Yank line up to cursor
Yank line after cursor

Insert <->

Insert yanked text
Insert <->

Insert %>%

Show help for function
Show source code

New document
New document (Chrome)

Open document
Save document

Close document
Close document (Chrome)

Close all documents

Extract function
Extract variable

Reindent lines

(Un)Comment lines

Reflow Comment
Reformat Selection

Select within braces

Show Diagnostics

Transpose Letters

Move Lines Up/Down

Copy Lines Up/Down

Add New Cursor Above

Add New Cursor Below

Move Active Cursor Up

Move Active Cursor Down

Find and Replace

Use Selection for Find

Replace and Find

Windows /Linux

Tab or Ctrl+Space

↑/↓
Enter, Tab, or →
Esc
Ctrl+Z
Ctrl+Shift+Z
Ctrl+X
Ctrl+C
Ctrl+V
Ctrl+A
Ctrl+D

Shift+[Arrow]

Option+Shift+←/→

Cmd+Shift+←

Cmd+Shift+→

Shift+PageUp/Down

Shift+Alt+↑/↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Option+-

Cmd+Shift+M

F1

F2

Cmd+Shift+N

Ctrl+Alt+Shift+N

Ctrl+O

Ctrl+S

Cmd+W

Cmd+Option+W

Cmd+Shift+W

Ctrl+Alt+X

Ctrl+Option+V

Ctrl+I

Ctrl+Shift+C

Ctrl+Shift+/

Ctrl+Shift+A

Ctrl+Shift+E

Ctrl+Shift+E

Cmd+Shift+Opt+P

Ctrl+T

Alt+↑/↓

Shift+Alt+↑/↓

Cmd+Option+↑/↓

Ctrl+Option+Up

Ctrl+Option+Down

Ctrl+Option+Shift+Up

Ctrl+Opt+Shift+Down

Ctrl+Opt+Shift+J

Mac

Tab or Cmd+Space

↑/↓
Enter, Tab, or →
Esc
Cmd+Z
Cmd+Shift+Z
Cmd+X
Cmd+C
Cmd+V
Cmd+A
Cmd+D

Shift+[Arrow]

Option+Shift+←/→

Cmd+Shift+←

Cmd+Shift+→

Shift+PageUp/Down

Shift+Alt+↑/↓

Ctrl+Opt+Backspace

Option+Delete

Ctrl+K

Option+Backspace

Tab (at start of line)

Shift+Tab

Ctrl+U

Ctrl+K

Ctrl+Y

Alt+-

Option+-

Cmd+Shift+M

F1

F2

Cmd+Shift+N

Cmd+Shift+Opt+N

Cmd+O

Ctrl+S

Cmd+W

Cmd+Option+W

Cmd+Shift+W

Ctrl+Alt+X

Ctrl+Option+V

Ctrl+I

Cmd+Shift+C

Cmd+Shift+/

Cmd+Shift+A

Cmd+Shift+E

Cmd+Shift+E

Cmd+Shift+Opt+P

Ctrl+T

Option+↑/↓

Shift+Alt+↑/↓

Cmd+Option+↑/↓

Ctrl+Option+Up

Ctrl+Option+Down

Ctrl+Option+Shift+Up

Ctrl+Opt+Shift+Down

Ctrl+Opt+Shift+J

WHY RSTUDIO SERVER PRO?

RSP extends the open source server with a commercial license, support, and more:

- open and run multiple R sessions at once
 - tune your resources to improve performance
 - edit the same project at the same time as others
 - see what you and others are doing on your server
 - switch easily from one version of R to a different version
 - integrate with your authentication, authorization, and audit practices
- Download a free 45 day evaluation at www.rstudio.com/products/rstudio-server-pro/



5 DEBUG CODE

Toggle Breakpoint
Execute Next Line
Step Into Function
Finish Function/Loop
Continue
Stop Debugging

Windows/Linux Mac

Shift+F9
F10
Shift+F4
Shift+F6
Shift+F5
Shift+F8

6 VERSION CONTROL

Show diff
Commit changes
Scroll diff view
Stage/Unstage (Git)
Stage/Unstage and move to next

Windows/Linux Mac

Ctrl+Alt+D
Ctrl+Alt+M
Ctrl+↑/↓
Spacebar
Enter

7 MAKE PACKAGES

Build and Reload
Load All (devtools)
Test Package (Desktop)
Test Package (Web)
Check Package
Document Package

Windows/Linux Mac

Ctrl+Shift+B
Cmd+Shift+L
Ctrl+Shift+T
Cmd+Shift+T
Ctrl+Alt+F7
Ctrl+Shift+E
Ctrl+Shift+D

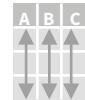
8 DOCUMENTS AND APPS

Preview HTML (Markdown, etc.)
Knit Document (knitr)
Compile Notebook
Compile PDF (TeX and Sweave)
Insert chunk (Sweave and Knitr)
Insert code section
Re-run previous region
Run current document
Run from start to current line
Run the current code section
Run previous Sweave/Rmd code
Run the current chunk
Run the next chunk
Sync Editor & PDF Preview
Previous plot
Next plot
Show Keyboard Shortcuts</

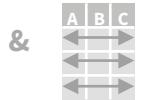
Data Transformation with dplyr :: CHEAT SHEET



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**



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).



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

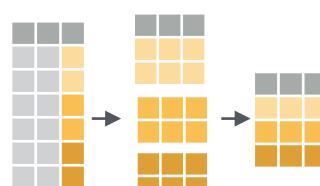
summarise_all() - Apply funs to every column.

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)`



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)`

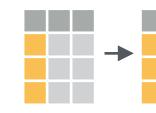
Manipulate Variables

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

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).



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)`



Vector Functions

TO USE WITH MUTATE ()

mutate() and **transmute()** apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function

OFFSETS

dplyr::lag() - Offset elements by 1
dplyr::lead() - Offset elements by -1

CUMULATIVE AGGREGATES

dplyr::cumall() - Cumulative all()
dplyr::cumany() - Cumulative any()
 cummax() - Cumulative max()
dplyr::cummean() - Cumulative mean()
 cummin() - Cumulative min()
 cumprod() - Cumulative prod()
 cumsum() - Cumulative sum()

RANKINGS

dplyr::cume_dist() - Proportion of all values <=
dplyr::dense_rank() - rank with ties = min, no gaps
dplyr::min_rank() - rank with ties = min
dplyr::ntile() - bins into n bins
dplyr::percent_rank() - min_rank scaled to [0,1]
dplyr::row_number() - rank with ties = "first"

MATH

+, -, *, /, ^, %/%, %% - arithmetic ops
log(), **log2()**, **log10()** - logs
<, <=, >, >=, !=, == - logical comparisons
dplyr::between() - x >= left & x <= right
dplyr::near() - safe == for floating point numbers

MISC

dplyr::case_when() - multi-case if_else()
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
 pmax() - element-wise max()
 pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

Summary Functions

TO USE WITH SUMMARISE ()

summarise() applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

COUNTS

dplyr::n() - number of values/rows
dplyr::n_distinct() - # of uniques
 sum(!is.na()) - # of non-NA's

LOCATION

mean() - mean, also **mean(!is.na())**
median() - median

LOGICALS

mean() - Proportion of TRUE's
sum() - # of TRUE's

POSITION/ORDER

dplyr::first() - first value
dplyr::last() - last value
dplyr::nth() - value in nth location of vector

RANK

quantile() - nth quantile
min() - minimum value
max() - maximum value

SPREAD

IQR() - Inter-Quartile Range
mad() - median absolute deviation
sd() - standard deviation
var() - variance

Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

A	B
1	a
2	b
3	c

rownames_to_column()

Move row names into col.
a <- rownames_to_column(iris, var = "C")

A	B	C
1	a	t
2	b	u
3	c	v

column_to_rownames()

Move col in row names.
column_to_rownames(a, var = "C")

Also **has_rownames()**, **remove_rownames()**

Combine Tables

COMBINE VARIABLES

X	A B C a t 1 b u 2 c v 3	+	y	A B D a t 3 b u 2 d w 1	=	A B C A B D a t 1 a t 3 b u 2 b u 2 c v 3 d w 1
---	----------------------------------	---	---	----------------------------------	---	--

Use **bind_cols()** to paste tables beside each other as they are.

bind_cols(...) Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

A B C D a t 1 3 b u 2 2 c v 3 NA	left_join(x, y, by = NULL, copy=FALSE, suffix=c("x","y"),...) Join matching values from y to x.
---	--

A B C D a t 1 3 b u 2 2 d w NA 1	right_join(x, y, by = NULL, copy = FALSE, suffix=c("x","y"),...) Join matching values from x to y.
---	---

A B C D a t 1 3 b u 2 2	inner_join(x, y, by = NULL, copy = FALSE, suffix=c("x","y"),...) Join data. Retain only rows with matches.
-------------------------------	---

A B C D a t 1 3 b u 2 2 d w NA 1	full_join(x, y, by = NULL, copy=FALSE, suffix=c("x","y"),...) Join data. Retain all values, all rows.
---	--

Use **by = c("col1", "col2", ...)** to specify one or more common columns to match on.
left_join(x, y, by = "A")

Use a named vector, **by = c("col1" = "col2")**, to match on columns that have different names in each table.
left_join(x, y, by = c("C" = "D"))

Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

COMBINE CASES

X	A B C a t 1 b u 2 c v 3	+	y	A B C C v 3 d w 4
---	----------------------------------	---	---	-------------------------

Use **bind_rows()** to paste tables below each other as they are.

bind_rows(..., .id = NULL)
Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured)

intersect(x, y, ...)
Rows that appear in both x and y.

setdiff(x, y, ...)
Rows that appear in x but not y.

union(x, y, ...)
Rows that appear in x or y.
(Duplicates removed). **union_all()** retains duplicates.

Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

EXTRACT ROWS

X	A B C a t 1 b u 2 c v 3	+	y	A B D a t 3 b u 2 d w 1	=
---	----------------------------------	---	---	----------------------------------	---

Use a "**Filtering Join**" to filter one table against the rows of another.

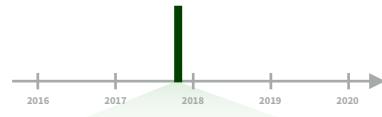
semi_join(x, y, by = NULL, ...)
Return rows of x that have a match in y.
USEFUL TO SEE WHAT WILL BE JOINED.

anti_join(x, y, by = NULL, ...)
Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.

Dates and times with lubridate :: CHEAT SHEET



Date-times



2017-11-28 12:00:00

2017-11-28 12:00:00

A **date-time** is a point on the timeline, stored as the number of seconds since 1970-01-01 00:00:00 UTC

```
dt <- as_datetime(1511870400)
## "2017-11-28 12:00:00 UTC"
```

PARSE DATE-TIMES (Convert strings or numbers to date-times)

1. Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data
2. Use the function below whose name replicates the order. Each accepts a wide variety of input formats.

2017-11-28T14:02:00

ymd_hms(), ymd_hm(), ymd_h().
ymd_hms("2017-11-28T14:02:00")

2017-22-12 10:00:00

ydm_hms(), ydm_hm(), ydm_h().
ydm_hms("2017-22-12 10:00:00")

11/28/2017 1:02:03

mdy_hms(), mdy_hm(), mdy_h().
mdy_hms("11/28/2017 1:02:03")

1 Jan 2017 23:59:59

dmy_hms(), dmy_hm(), dmy_h().
dmy_hms("1 Jan 2017 23:59:59")

20170131

ymd(), ydm(). ymd(20170131)

July 4th, 2000

mdy(), myd(). mdy("July 4th, 2000")

4th of July '99

dmy(), dym(). dmy("4th of July '99")

2001: Q3

yq() Q for quarter. yq("2001: Q3")

2:01

hms::hms() Also lubridate::hms(), hm() and ms(), which return periods.* hms::hms(sec = 0, min = 1, hours = 2)

2017.5

date_decimal(decimal, tz = "UTC") Q for quarter. date_decimal(2017.5)



now(tzone = "") Current time in tz (defaults to system tz). now()

today(tzone = "") Current date in a tz (defaults to system tz). today()

fast.strptime() Faster strftime.
fast.strptime('9/1/01', '%y/%m/%d')

parse_date_time() Easier strftime.
parse_date_time("9/1/01", "ymd")

2017-11-28

A **date** is a day stored as the number of days since 1970-01-01

```
d <- as_date(17498)
## "2017-11-28"
```

12:00:00

An hms is a **time** stored as the number of seconds since 00:00:00

```
t <- hms::as.hms(85)
## 00:01:25
```

GET AND SET COMPONENTS

Use an accessor function to get a component. Assign into an accessor function to change a component in place.

```
d ## "2017-11-28"
day(d) ## 28
day(d) <- 1
d ## "2017-11-01"
```

2018-01-31 11:59:59

date(x) Date component. date(dt)

2018-01-31 11:59:59

year(x) Year. year(dt)
isoyear(x) The ISO 8601 year.
epiyear(x) Epidemiological year.

2018-01-31 11:59:59

month(x, label, abbr) Month.
month(dt)

2018-01-31 11:59:59

day(x) Day of month. day(dt)
wday(x, label, abbr) Day of week.
qday(x) Day of quarter.

2018-01-31 11:59:59

hour(x) Hour. hour(dt)

2018-01-31 11:59:59

minute(x) Minutes. minute(dt)

2018-01-31 11:59:59

second(x) Seconds. second(dt)

2018-01-31 11:59:59

week(x) Week of the year. week(dt)
isoweek() ISO 8601 week.
epiweek() Epidemiological week.

2018-01-31 11:59:59

quarter(x, with_year = FALSE)
Quarter. quarter(dt)

2018-01-31 11:59:59

semester(x, with_year = FALSE)
Semester. semester(dt)

2018-01-31 11:59:59

am(x) Is it in the am? am(dt)
pm(x) Is it in the pm? pm(dt)

2018-01-31 11:59:59

dst(x) Is it daylight savings? dst(dt)

2018-01-31 11:59:59

leap_year(x) Is it a leap year?
leap_year(dt)

2018-01-31 11:59:59

update(object, ..., simple = FALSE)
update(dt, mday = 2, hour = 1)

Round Date-times



floor_date(x, unit = "second") Round down to nearest unit.
floor_date(dt, unit = "month")



round_date(x, unit = "second") Round to nearest unit.
round_date(dt, unit = "month")



ceiling_date(x, unit = "second", change_on_boundary = NULL) Round up to nearest unit.
ceiling_date(dt, unit = "month")

rollback(dates, roll_to_first = FALSE, preserve_hms = TRUE) Roll back to last day of previous month. **rollback**(dt)

Stamp Date-times

stamp() Derive a template from an example string and return a new function that will apply the template to date-times. Also **stamp_date()** and **stamp_time()**.

1. Derive a template, create a function
`sf <- stamp("Created Sunday, Jan 17, 1999 3:34")`

Tip: use a date with day > 12

2. Apply the template to dates
`sf(ymd("2010-04-05"))`
`## [1] "Created Monday, Apr 05, 2010 00:00"`

Time Zones

R recognizes ~600 time zones. Each encodes the time zone, Daylight Savings Time, and historical calendar variations for an area. R assigns one time zone per vector.

Use the **UTC** time zone to avoid Daylight Savings.

OlsonNames() Returns a list of valid time zone names. **OlsonNames()**

5:00 Mountain 6:00 Central
4:00 Pacific 7:00 Eastern

PT MT CT ET

7:00 Pacific 7:00 Mountain
7:00 Central

with_tz(time, tzone = "") Get the same date-time in a new time zone (a new clock time).
with_tz(dt, "US/Pacific")

7:00 Pacific 7:00 Mountain
7:00 Central

force_tz(time, tzone = "") Get the same clock time in a new time zone (a new date-time).
force_tz(dt, "US/Pacific")



Math with Date-times

— Lubridate provides three classes of timespans to facilitate math with dates and date-times

Math with date-times relies on the **timeline**, which behaves inconsistently. Consider how the timeline behaves during:

A normal day
`nor <- ymd_hms("2018-01-01 01:30:00", tz="US/Eastern")`

The start of daylight savings (spring forward)
`gap <- ymd_hms("2018-03-11 01:30:00", tz="US/Eastern")`

The end of daylight savings (fall back)
`lap <- ymd_hms("2018-11-04 00:30:00", tz="US/Eastern")`

Leap years and leap seconds
`leap <- ymd("2019-03-01")`

Periods track changes in clock times, which ignore time line irregularities.

`normal + minutes(90)`

`gap + minutes(90)`

`lap + minutes(90)`

`leap + years(1)`

Durations track the passage of physical time, which deviates from clock time when irregularities occur.

`normal + dminutes(90)`

`gap + dminutes(90)`

`lap + dminutes(90)`

`leap + dyears(1)`

Intervals represent specific intervals of the timeline, bounded by start and end date-times.

`interval(normal, normal + minutes(90))`

`interval(gap, gap + minutes(90))`

`interval(lap, lap + minutes(90))`

`interval(leap, leap + years(1))`

Not all years are 365 days due to **leap days**.

Not all minutes are 60 seconds due to **leap seconds**.

It is possible to create an imaginary date by adding **months**, e.g. February 31st

`jan31 <- ymd(20180131)`
`jan31 + months(1)`
`## NA`

`%m+%` and `%m-%` will roll imaginary dates to the last day of the previous month.

`jan31 %m+% months(1)`
`## "2018-02-28"`

`add_with_rollback(e1, e2, roll_to_first = TRUE)` will roll imaginary dates to the first day of the new month.

`add_with_rollback(jan31, months(1), roll_to_first = TRUE)`
`## "2018-03-01"`

PERIODS

Add or subtract periods to model events that happen at specific clock times, like the NYSE opening bell.

Make a period with the name of a time unit **pluralized**, e.g.

```
p <- months(3) + days(12)
p
## "3m 12d 0H 0M 0S"
```

Number of months Number of days etc.

`years(x = 1) x years.`
`months(x) x months.`
`weeks(x = 1) x weeks.`
`days(x = 1) x days.`
`hours(x = 1) x hours.`
`minutes(x = 1) x minutes.`
`seconds(x = 1) x seconds.`
`milliseconds(x = 1) x milliseconds.`
`microseconds(x = 1) x microseconds`
`nanoseconds(x = 1) x milliseconds.`
`picoseconds(x = 1) x picoseconds.`

`period(num = NULL, units = "second", ...)`
An automation friendly period constructor.
`period(5, unit = "years")`

`as.period(x, unit)` Coerce a timespan to a period, optionally in the specified units.
Also `is.period()`. `as.period(i)`

`period_to_seconds(x)` Convert a period to the "standard" number of seconds implied by the period. Also `seconds_to_period()`.
`period_to_seconds(p)`

DURATIONS

Add or subtract durations to model physical processes, like battery life. Durations are stored as seconds, the only time unit with a consistent length.

Diftimes are a class of durations found in base R.

Make a duration with the name of a period prefixed with a **d**, e.g.

```
dd <- ddays(14)
dd
## 1209600s (~2 weeks)"
```

Exact length in seconds Equivalent in common units

`dyears(x = 1) 31536000x seconds.`
`dweeks(x = 1) 604800x seconds.`
`ddays(x = 1) 86400x seconds.`
`dhours(x = 1) 3600x seconds.`
`dminutes(x = 1) 60x seconds.`
`dseconds(x = 1) x seconds.`
`dmilliseconds(x = 1) x × 10-3 seconds.`
`dmicroseconds(x = 1) x × 10-6 seconds.`
`dnanoseconds(x = 1) x × 10-9 seconds.`
`dpicoseconds(x = 1) x × 10-12 seconds.`

`duration(num = NULL, units = "second", ...)`
An automation friendly duration constructor. `duration(5, unit = "years")`

`as.duration(x, ...)` Coerce a timespan to a duration. Also `is.duration()`, `is.difftime()`. `as.duration(i)`

`make_difftime(x)` Make difftime with the specified number of units.
`make_difftime(99999)`

INTERVALS

Divide an interval by a duration to determine its physical length, divide and interval by a period to determine its implied length in clock time.

Make an interval with **interval()** or `%--%`, e.g.

```
i <- interval(ymd("2017-01-01"), d)
j <- d %--% ymd("2017-12-31")
## 2017-01-01 UTC--2017-11-28 UTC
## 2017-11-28 UTC--2017-12-31 UTC
```

Start Date End Date

a %within% b Does interval or date-time a fall within interval b? `now() %within% i`

`int_start(int)` Access/set the start date-time of an interval. Also `int_end()`. `int_start(i) <- now(); int_start(i)`

`int_aligns(int1, int2)` Do two intervals share a boundary? Also `int_overlaps()`. `int_aligns(i, j)`

`int_diff(times)` Make the intervals that occur between the date-times in a vector.
`v <- c(dt, dt + 100, dt + 1000)); int_diff(v)`

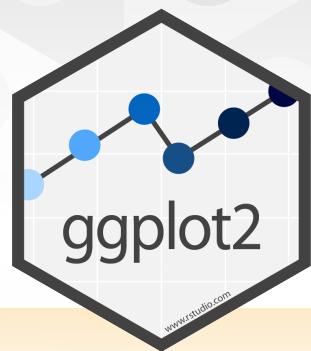
`int_flip(int)` Reverse the direction of an interval. Also `int_standardize()`. `int_flip(i)`

`int_length(int)` Length in seconds. `int_length(i)`

`int_shift(int, by)` Shifts an interval up or down the timeline by a timespan. `int_shift(i, days(-1))`

`as.interval(x, start, ...)` Coerce a timespans to an interval with the start date-time. Also `is.interval()`. `as.interval(days(1), start = now())`

Data Visualization with ggplot2 :: CHEAT SHEET



Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data** set, a **coordinate system**, and geoms—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot (data = <DATA>) +
<GEOM_FUNCTION>(mapping = aes(<MAPPINGS>),
stat = <STAT>, position = <POSITION>) +
<COORDINATE_FUNCTION> +
<FACET_FUNCTION> +
<SCALE_FUNCTION> +
<THEME_FUNCTION>
```

required

Not required, sensible defaults supplied

ggplot(data = mpg, **aes**(x = cty, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

aesthetic mappings **data** **geom**

qplot(x = cty, y = hwy, data = mpg, geom = "point") Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

last_plot() Returns the last plot

gsave("plot.png", **width** = 5, **height** = 5) Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES

- a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))
- a + geom_blank()**
(Useful for expanding limits)
- b + geom_curve(aes(yend = lat + 1, xend = long + 1, curvature = z))** - x, yend, alpha, angle, color, curvature, linetype, size
- a + geom_path(lineend = "butt", linejoin = "round", linemitre = 1)** - x, y, alpha, color, group, linetype, size
- a + geom_polygon(aes(group = group))** - x, y, alpha, color, fill, group, linetype, size
- b + geom_rect(aes(xmin = long, ymin = lat, xmax = long + 1, ymax = lat + 1))** - xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size
- a + geom_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

LINE SEGMENTS

- common aesthetics: x, y, alpha, color, linetype, size
- b + geom_abline(aes(intercept = 0, slope = 1))**
 - b + geom_hline(aes(yintercept = lat))**
 - b + geom_vline(aes(xintercept = long))**

- b + geom_segment(aes(yend = lat + 1, xend = long + 1))**
- b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

- c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
- c + geom_area(stat = "bin")** - x, y, alpha, color, fill, linetype, size
- c + geom_density(kernel = "gaussian")** - x, y, alpha, color, fill, group, linetype, size, weight
- c + geom_dotplot()** - x, y, alpha, color, fill
- c + geom_freqpoly()** - x, y, alpha, color, group, linetype, size
- c + geom_histogram(binwidth = 5)** - x, y, alpha, color, fill, linetype, size, weight
- c2 + geom_qq(aes(sample = hwy))** - x, y, alpha, color, fill, linetype, size, weight

discrete

- d <- ggplot(mpg, aes(f1))
- d + geom_bar()** - x, alpha, color, fill, linetype, size, weight

TWO VARIABLES

continuous x , continuous y

- e <- ggplot(mpg, aes(cty, hwy))
- e + geom_label(aes(label = cty), nudge_x = 1, nudge_y = 1, check_overlap = TRUE)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

- e + geom_jitter(height = 2, width = 2)** - x, y, alpha, color, fill, shape, size

- e + geom_point()** - x, y, alpha, color, fill, shape, size, stroke

- e + geom_quantile()** - x, y, alpha, color, group, linetype, size, weight

- e + geom_rug(sides = "bl")** - x, y, alpha, color, linetype, size

- e + geom_smooth(method = lm)** - x, y, alpha, color, fill, group, linetype, size, weight

- e + geom_text(aes(label = cty), nudge_x = 1, nudge_y = 1, check_overlap = TRUE)** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

discrete x , continuous y

- f <- ggplot(mpg, aes(class, hwy))

- f + geom_col()** - x, y, alpha, color, fill, group, linetype, size

- f + geom_boxplot()** - x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight

- f + geom_dotplot(binaxis = "y", stackdir = "center")** - x, y, alpha, color, fill, group

- f + geom_violin(scale = "area")** - x, y, alpha, color, fill, group, linetype, size, weight

discrete x , discrete y

- g <- ggplot(diamonds, aes(cut, color))

- g + geom_count()** - x, y, alpha, color, fill, shape, size, stroke

THREE VARIABLES

- seals\$z <- with(seals, sqrt(delta_long^2 + delta_lat^2))
l <- ggplot(seals, aes(long, lat))

- l + geom_contour(aes(z = z))** - x, y, z, alpha, colour, group, linetype, size, weight

continuous bivariate distribution

- h <- ggplot(diamonds, aes(carat, price))
- h + geom_bin2d(binwidth = c(0.25, 500))** - x, y, alpha, color, fill, linetype, size, weight

- h + geom_density2d()** - x, y, alpha, colour, group, linetype, size

- h + geom_hex()** - x, y, alpha, colour, fill, size

continuous function

- i <- ggplot(economics, aes(date, unemploy))

- i + geom_area()** - x, y, alpha, color, fill, linetype, size

- i + geom_line()** - x, y, alpha, color, group, linetype, size

- i + geom_step(direction = "hv")** - x, y, alpha, color, group, linetype, size

visualizing error

- df <- data.frame(grp = c("A", "B"), fit = 4.5, se = 1.2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))

- j + geom_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size

- j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width (also **geom_errorbarh()**)

- j + geom_linerange()** - x, ymin, ymax, alpha, color, group, linetype, size

- j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

maps

- data <- data.frame(murder = USArrests\$Murder, state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))

- k + geom_map(aes(map_id = state), map = map) + expand_limits(x = map\$long, y = map\$lat)** - map_id, alpha, color, fill, linetype, size

