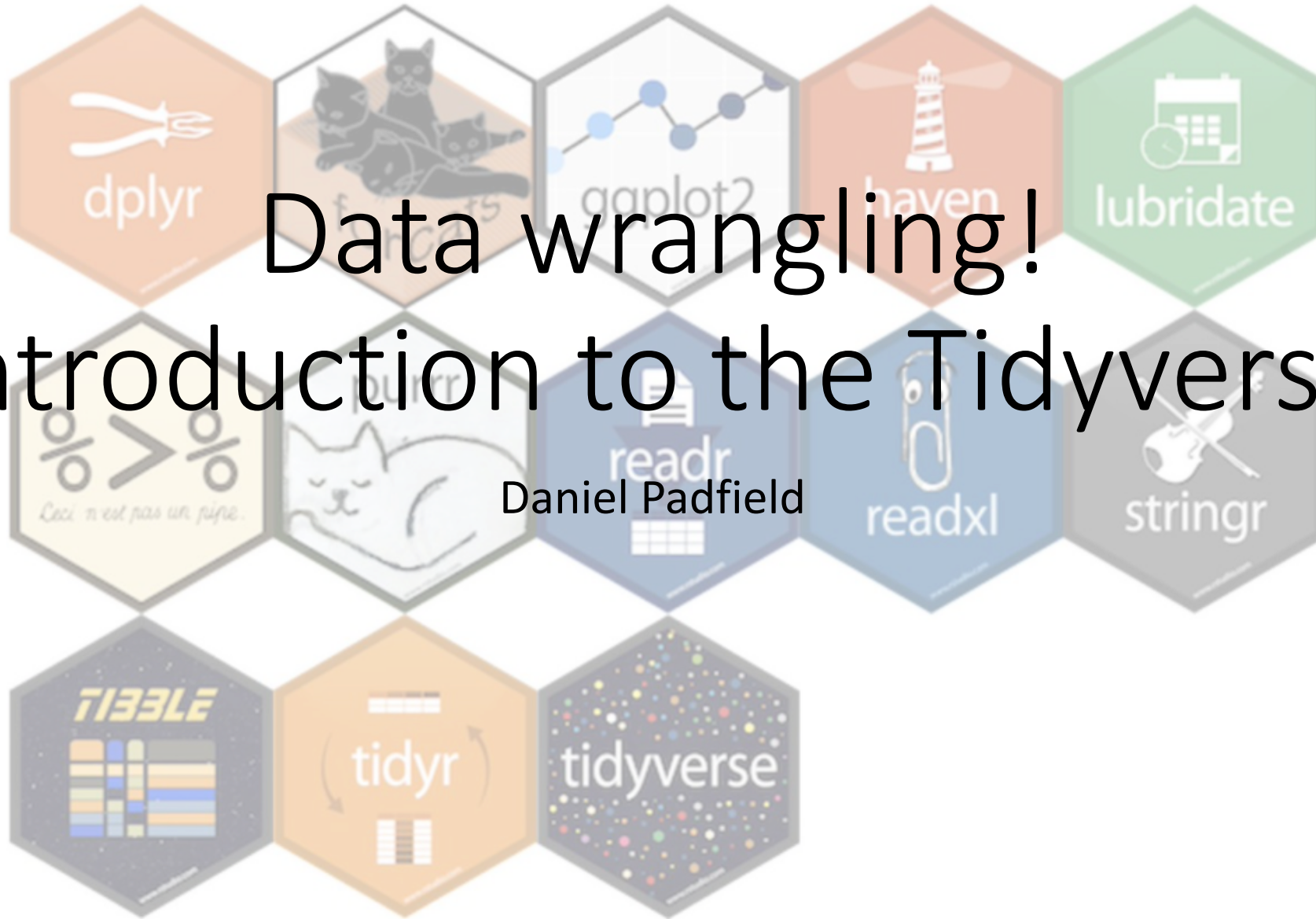


Components

Data wrangling!

Introduction to the Tidyverse

Daniel Padfield



What is the tidyverse?

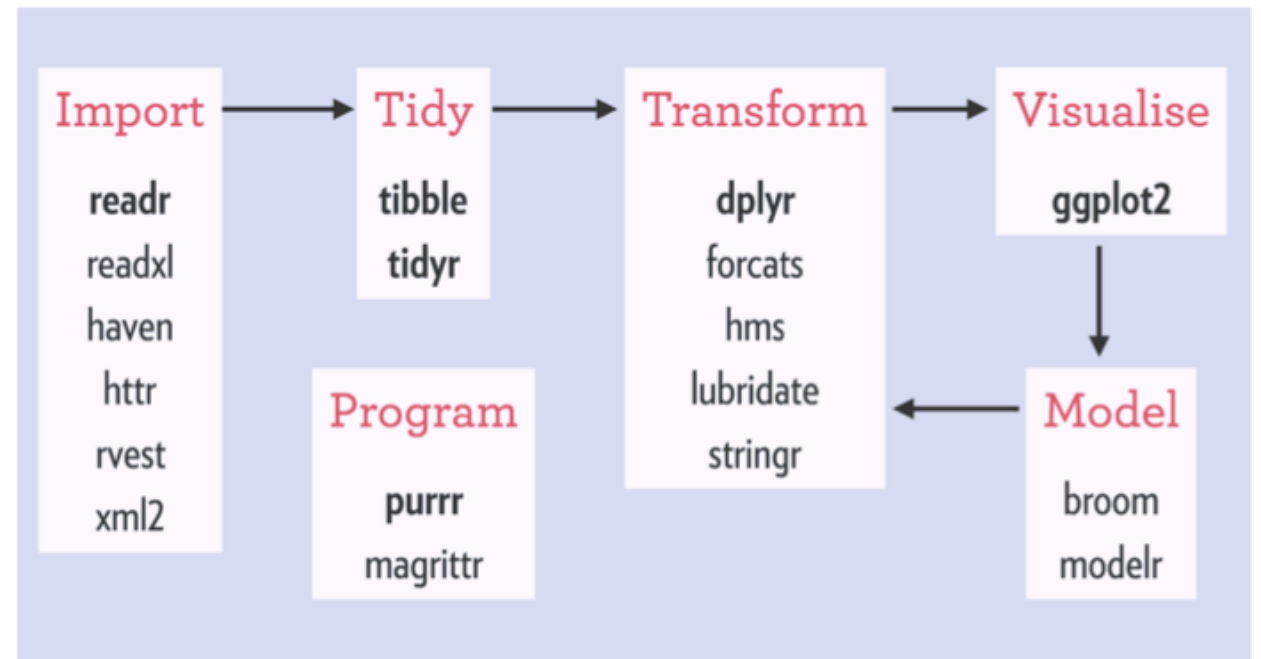
- Set of packages designed by Hadley Wickham and his team at Rstudio
- Try to create a set of tools for full-stack analysis that use the same syntax and set of rules
- Install the complete tidyverse using:
 - `install.packages("tidyverse")`

Components



What is the tidyverse?

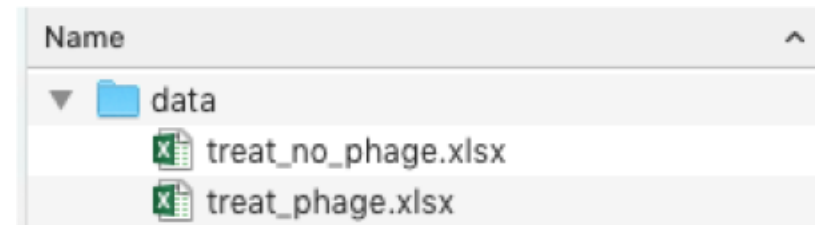
- Can do almost all my programming with the same set of packages
- Makes my code more readable so future me can understand
- In general, drops the need for \$ and “ ”
- Great documentation and support online



Read in files with purrr and readxl



- `read_excel()` handles excel files
- `map_df()` takes a vector a vector of file names and returns a dataframe



treat_phage.xlsx

Treatment	CLONE	10 degrees	15 degrees	20 degrees
phage	1	75	85	88
phage	2	64	64	74
phage	3	68	68	70
phage	4	89	94	99
phage	5	110	115	116
phage	6	48	52	57
phage	7	69	69	69
phage	8	140	141	151
phage	9	100	110	117
phage	10	99	107	107

Read in files with purrr and readxl



```
# list files
```

```
files <- list.files('data', full.names = TRUE)
```

```
# read in many files using map() and read_excel()
```

```
d <- map_df(files, read_excel)
```

	Treatment	CLONE	10 degrees	15 degrees	20 degrees
1	no phage	1	56	73	109
2	no phage	2	131	145	202
3	no phage	3	127	141	188
4	no phage	4	91	111	158
5	no phage	5	114	128	162
6	no phage	6	165	180	223
7	no phage	7	77	92	136
8	no phage	8	168	186	246
9	no phage	9	98	116	155
10	no phage	10	45	56	113
11	phage	1	75	78	82
12	phage	2	64	70	73
13	phage	3	68	71	71
14	phage	4	89	95	103
15	phage	5	110	111	115
16	phage	6	48	53	53
17	phage	7	69	78	82
18	phage	8	140	146	147
19	phage	9	100	100	105
20	phage	10	99	100	102

Tidy column names with janitor



- Writing code, odd column names can often be confusing
- No CAPITALS, numbers as first letter or spaces
- `janitor::clean_names()` automatically renames columns to be all lower case and replaces spaces with “_”, puts an X in front of the first number
- Can do one at a time using:

```
dplyr::rename(d, new_column_name = old_column name)
```

Tidy column names with janitor



```
# change column names with janitor
d <- clean_names(d)

# or a single column with rename
d <- rename(d, clone = CLONE)
```

• From this...

Treatment	CLONE	10 degrees	15 degrees	20 degrees
165				
no phage	1	56	73	109
no phage	2	131	145	202
no phage	3	127	141	188
no phage	4	91	111	158
no phage	5	114	128	162
no phage	6	165	180	223
no phage	7	77	92	136
no phage	8	168	186	246
no phage	9	98	116	155
no phage	10	45	56	113

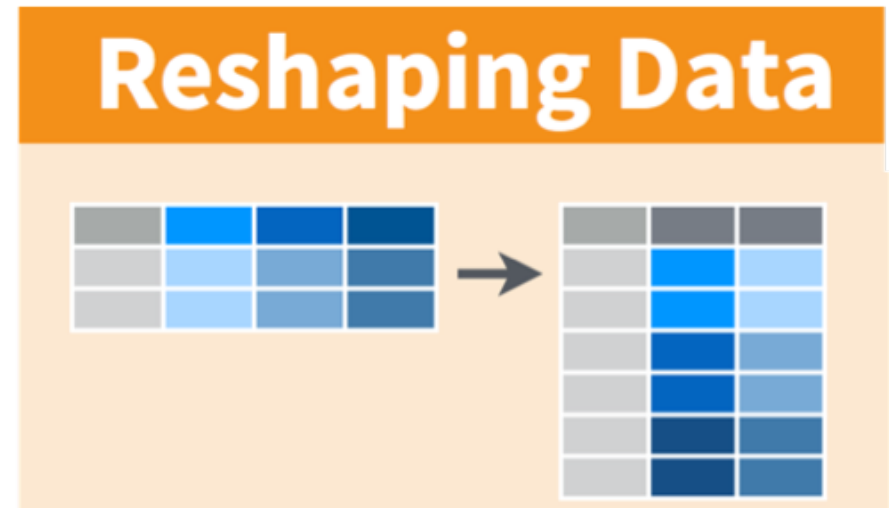
• To this...

treatment	clone	x10_degrees	x15_degrees	x20_degrees
no phage	1	56	73	109
no phage	2	131	145	202
no phage	3	127	141	188
no phage	4	91	111	158
no phage	5	114	128	162
no phage	6	165	180	223
no phage	7	77	92	136
no phage	8	168	186	246
no phage	9	98	116	155
no phage	10	45	56	113

Stack data (i.e. long format) using gather()



- Lots of packages (ggplot2, lme4 etc) need data to be in long format
- `tidyr::gather()` combines multiple columns into key value pairs
- Similar to `reshape2::melt()`



Stack data (i.e. long format) using gather()



```
# convert data to long format using  
gather
```

```
d <- gather(d, key = 'temp',  
            value = 'count',  
            contains('degree'))
```

• From this...

• To this...

treatment	clone	x10_degrees	x15_degrees	x20_degrees
no phage	1	56	73	109
no phage	2	131	145	202
no phage	3	127	141	188
no phage	4	91	111	158
no phage	5	114	128	162
no phage	6	165	180	223
no phage	7	77	92	136
no phage	8	168	186	246
no phage	9	98	116	155
no phage	10	45	56	113
phage	1	75	78	82
phage	2	64	70	73
phage	3	68	71	71
phage	4	89	95	103
phage	5	110	111	115
phage	6	48	53	53
phage	7	69	78	82
phage	8	140	146	147
phage	9	100	100	105
phage	10	99	100	102

treatment	clone	temp	count
no phage	1	x10_degrees	56
no phage	2	x10_degrees	131
no phage	3	x10_degrees	127
no phage	4	x10_degrees	91
no phage	5	x10_degrees	114
no phage	6	x10_degrees	165
no phage	7	x10_degrees	77
no phage	8	x10_degrees	168
no phage	9	x10_degrees	98
no phage	10	x10_degrees	45
phage	1	x10_degrees	75
phage	2	x10_degrees	64
phage	3	x10_degrees	68
phage	4	x10_degrees	89
phage	5	x10_degrees	110
phage	6	x10_degrees	48
phage	7	x10_degrees	69
phage	8	x10_degrees	140
phage	9	x10_degrees	100
phage	10	x10_degrees	99
no phage	1	x15_degrees	73
no phage	2	x15_degrees	145
no phage	3	x15_degrees	141
no phage	4	x15_degrees	111
no phage	5	x15_degrees	128
no phage	6	x15_degrees	180
no phage	7	x15_degrees	92
no phage	8	x15_degrees	186
no phage	9	x15_degrees	116
no phage	10	x15_degrees	56
phage	1	x15_degrees	78
phage	2	x15_degrees	70
phage	3	x15_degrees	71
phage	4	x15_degrees	95
phage	5	x15_degrees	111
phage	6	x15_degrees	53
phage	7	x15_degrees	78
phage	8	x15_degrees	146
phage	9	x15_degrees	100
phage	10	x15_degrees	100



Select columns with select()

- Can select or deselect columns
- Via position
`select(d, c(1,3,4))`
- Via name
`select(d, c(clone, treatment, temp, count))`
- Drop columns
`select(d, -c(clone, treatment))`
- Helper functions
`select(d, c(clone, everything()))`
`select(d, starts_with('t'))`
`?select_helpers` for all functions



Add columns with mutate()

- Create new columns with mutate()

change and add columns

```
d <- mutate(d, temp = readr::parse_number(temp),  
            temp = as.numeric(temp),  
            log_count = log10(count))
```

treatment	clone	temp	count	log_count
no phage	1	10	56	1.748188
no phage	2	10	131	2.117271
no phage	3	10	127	2.103804
no phage	4	10	91	1.959041
no phage	5	10	114	2.056905
no phage	6	10	165	2.217484
no phage	7	10	77	1.886491
no phage	8	10	168	2.225309
no phage	9	10	98	1.991226



Do the same function on many columns using mutate_at()



- We can do the same function on multiple columns using mutate_at()
- Change character or numeric columns to a factor

change multiple columns using mutate_at()

```
d <- mutate_at(d, c('treatment', 'clone'), as.factor)
```

 d	60 obs. of 5 variables	
treatment: chr "no phage" "no phage" "no phage" "no phage" ...		
clone : num 1 2 3 4 5 6 7 8 9 10 ...		

Do the same function on many columns using mutate_at()



- Do the same function on multiple columns using mutate_at()
- Change character or numeric columns to a factor

change multiple columns using mutate_at()

```
d <- mutate_at(d, c('treatment', 'clone'), as.factor)
```

d	60 obs. of 5 variables
treatment:	Factor w/ 2 levels "no phage","phage": 1 1 1 1 1..
clone :	Factor w/ 10 levels "1","2","3","4",...: 1 2 3 4 5 6..

No more nested if_else() with case_when()



- Different dilution factors
 - $10^{\circ}\text{C} = 10^{-5}$, $15^{\circ}\text{C} = 10^{-4}$, $20^{\circ}\text{C} = 10^{-6}$
- Different bacteria
 - Clones 1-5 = LacZ, Clones 6-10 = wild type
- Common to want to add columns based on other conditions

No more nested if_else() with case_when()



- Different dilution factors
 - $10^{\circ}\text{C} = 10^{-5}$, $15^{\circ}\text{C} = 10^{-4}$, $20^{\circ}\text{C} = 10^{-6}$
- Different bacteria
 - Clones 1-5 = LacZ, Clones 6-10 = wild type

```
# using ifelse() and mutate()
d <- mutate(d, type = ifelse(clone <= 5, 'lacz', 'wt'),
              df = ifelse(temp == 10, 10^-5,
                          ifelse(temp == 15, 10^-4,
                                10^-6)))
```

No more nested if_else() with case_when()



- Different dilution factors
 - $10^{\circ}\text{C} = 10^{-4}$, $15^{\circ}\text{C} = 10^{-5}$, $20^{\circ}\text{C} = 10^{-6}$
- Different bacteria
 - Clones 1-5 = LacZ, Clones 6-10 = wild type

```
# using ifelse() and mutate()
d <- mutate(d, type = ifelse(clone <= 5, 'lacz', 'wt'),
              df = ifelse(temp == 10, 10^-4,
                          ifelse(temp == 15, 10^-5,
                                10^-6)))
```

```
# using case_when() and mutate()
d <- mutate(d, type = ifelse(clone <= 5, 'lacz', 'wt'),
              df = case_when(temp == 10 ~ 10^-4,
                             temp == 15 ~ 10^-5,
                             temp == 20 ~ 10^-6))
```




Subset rows with filter()

- Similar to subset()
- removes rows of a dataframe based on a condition

```
# examples of using filter
```

```
# keep just 15 and 20 degrees
```

```
filter(d, temp > 10)
```

```
# keep just 10 degrees LacZ
```

```
filter(d, temp == 10 & type == 'lacz')
```

Piping with %>%



- Read a pipe as “and then”
- Allows your code to be read like a recipe
- Removes the need to constantly re-assign objects
- The “.” operator acts as an indicator for whatever is coming from the left hand side of the pipe

Put it all together! Pipe example!



Example of using pipes %>%

```
d <- list.files('data', full.names = TRUE) %>%  
  map_df(., read_excel) %>%  
  clean_names(.) %>%  
  gather(., key = 'temp', value = 'count',  
         contains('degree')) %>%  
  mutate(., temp = readr::parse_number(temp),  
         temp = as.numeric(temp),  
         log_count = log10(count)) %>%  
  mutate_at(., c('treatment', 'clone'),  
            as.factor) %>%  
  ...
```

The recipe version

d <- list all the files in the “data” directory (and then)
 read them into a dataframe (and then)
 change their names (and then)
 make the data long format (and then)
 add a bunch of columns (and then)
 do the same function to treatment and clone ...

Summary data using group_by() and summarise()



- Get summary stats based on groupings of your data
- summarise() - get means and standard deviations, or anything that returns a single value (i.e. max, min, number of rows, mode, median)

```
# get summary data
d_means <- group_by(d, type, treatment) %>%
  summarise(mean = mean(log_count),
             sd = sd(log_count)) %>%
  ungroup()
```

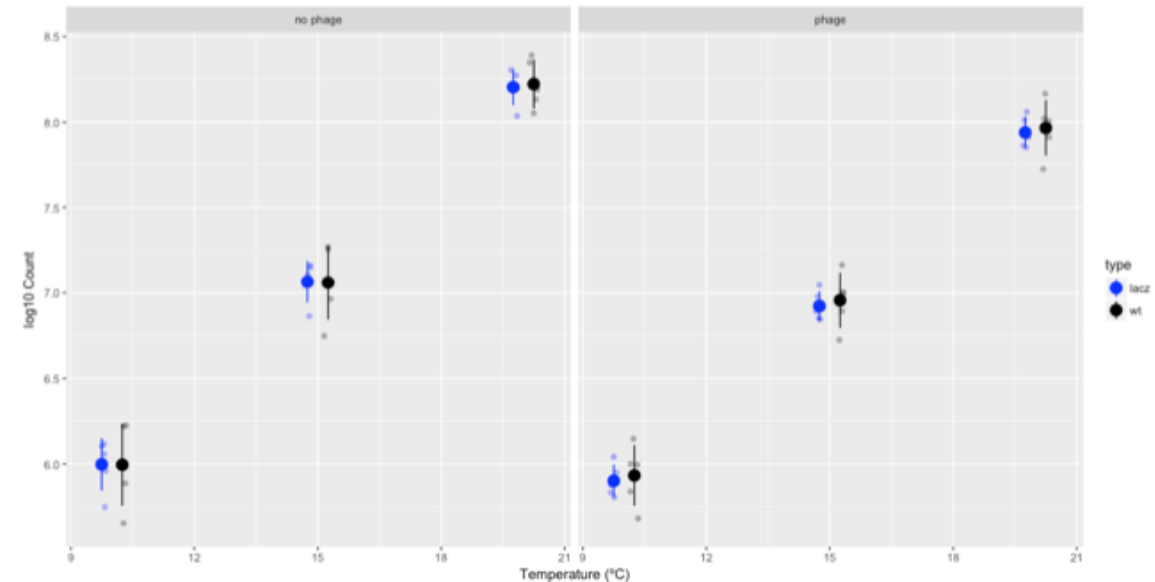
type	treatment	temp	mean	sd
lacz	no phage	10	5.997042	0.15232289
lacz	no phage	15	7.065289	0.12167762
lacz	no phage	20	8.205022	0.10367582
lacz	phage	10	5.900907	0.09540223
lacz	phage	15	6.922300	0.08678401
lacz	phage	20	7.940386	0.09261839
wt	no phage	10	5.994745	0.24038944
wt	no phage	15	7.060244	0.21710026
wt	no phage	20	8.223238	0.14304463
wt	phage	10	5.932371	0.17753916
wt	phage	15	6.956145	0.16206695
wt	phage	20	7.967039	0.16320068

Plot with ggplot2

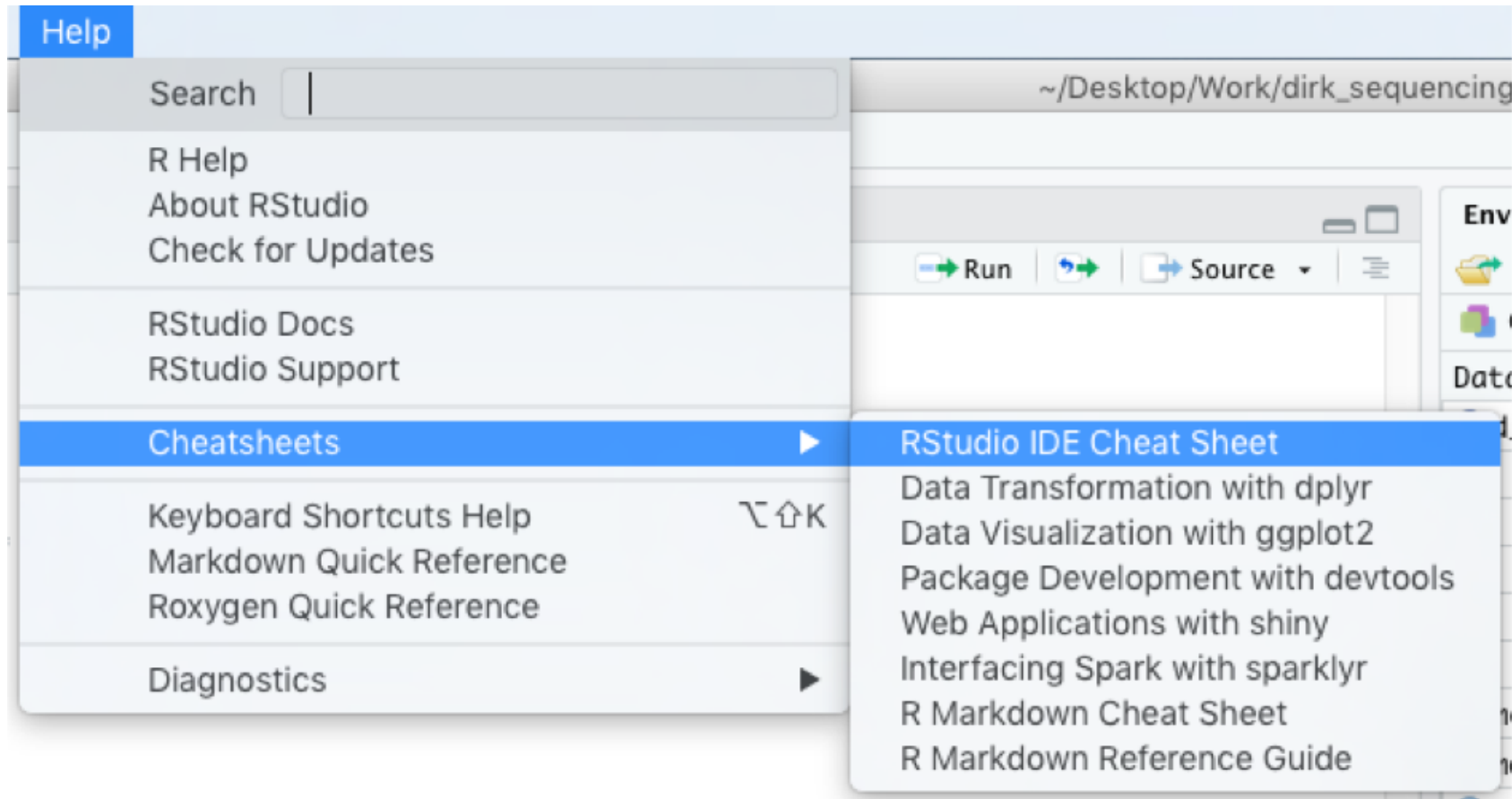


- Plot summary data with raw data

```
# plot with ggplot2
ggplot() +
  geom_linerange(aes(x = temp, ymin = mean - sd, ymax = mean + sd, col = type, group = type), d_means,
    position = position_dodge(width = 1)) +
  geom_point(aes(temp, mean, col = type), d_means,
    size = 4, position = position_dodge(width = 1)) +
  geom_point(aes(temp, log_count, col = type),
    position = position_jitterdodge(dodge.width = 1,
    jitter.width = 0.2), alpha = 0.3, d) +
  facet_wrap(~ treatment) +
  scale_color_manual(values = c('blue', 'black')) +
  ylab('log10 Count') +
  xlab('Temperature (°C)')
```



Cheatsheets for the tidyverse



Cheatsheets for the tidyverse



Data Import :: CHEAT SHEET

R 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 **tidyverse** and to layout tidy data with **tidyverse**.

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
- readr** - HTML (Web Scraping)

Save Data

Save 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 vector to file, one element per line

```
write_lines(x, path, na = "NA", append = FALSE)
```

Object to RDS file

```
saveRDS(x, path, append = FALSE)
```

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, locale2, na = c("","NA"),
skip_rows = 0, skip_cols = 0, guess_max = 100, guess_max2 = 100,
n_max_col, progress = interactive()
```

Comma Delimited Files

```
read_csv("file.csv")
```

Excel files

```
read_excel("file.xlsx")
```

SPSS files

```
read_spss("file.sav")
```

Stata files

```
read_stata("file.dta")
```

SAS files

```
read_sas("file.sas7bcat")
```

Apply Functions

Map functions apply a function iteratively to each element of a vector.

```
map(x, f, .args = list())
```

map2(x, y, f, .args = list())

map2_l(x, y, f, .args = list())

map2_df(x, y, f, .args = list())

map2_dfr(x, y, f, .args = list())

map2_dfc(x, y, f, .args = list())

map2_dtbl(x, y, f, .args = list())

map2_dttbl(x, y, f, .args = list())

map2_dttbl2(x, y, f, .args = list())

map2_dttbl3(x, y, f, .args = list())

map2_dttbl4(x, y, f, .args = list())

map2_dttbl5(x, y, f, .args = list())

map2_dttbl6(x, y, f, .args = list())

map2_dttbl7(x, y, f, .args = list())

map2_dttbl8(x, y, f, .args = list())

map2_dttbl9(x, y, f, .args = list())

map2_dttbl10(x, y, f, .args = list())

Useful R Functions

map(x, f, .args = list()) Apply function to each list element of a list

map2(x, y, f, .args = list()) Apply function to each element of a list and a vector

map2_l(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a list

map2_df(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data frame

map2_dfr(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data frame, rowwise

map2_dfc(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data frame, columnwise

map2_dtbl(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table

map2_dttbl(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise

map2_dttbl2(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 2 columns

map2_dttbl3(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 3 columns

map2_dttbl4(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 4 columns

map2_dttbl5(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 5 columns

map2_dttbl6(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 6 columns

map2_dttbl7(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 7 columns

map2_dttbl8(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 8 columns

map2_dttbl9(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 9 columns

map2_dttbl10(x, y, f, .args = list()) Apply function to each element of a list and a vector, returning a data table, rowwise, with 10 columns

Data Wrangling with dplyr and tidy

Cheat Sheet

Syntax - Helpful conventions for wrangling

dplyr::tbl_dfiris()

Converts data to tbl class. tbls are easier to examine than data frames. R displays only the data that fits on screen:

```
Source: local data frame [150 x 5]
  Sepal.Length Sepal.Width Petal.Length
1 5.1 3.5 1.4
2 4.9 3.0 1.4
3 5.4 4.7 1.5
4 4.3 3.0 1.1
5 5.7 4.4 1.5
...
Variables not shown: Petal.Width (dbl), Species (fctr)
```

dplyr::glimpseiris()

Information dense summary of tbl data.

utils::Viewiris()

View data set in spreadsheet-like display (note capital V).

dplyr::%>%

Passes object on left hand side as first argument (or argument) of function on righthand side.

```
x %>% f(y) is the same as f(x, y)
y %>% f(x, ..) is the same as f(x, y, z)
```

"Piping" with %>% makes code more readable, e.g.

```
iris %>%
  group_by(Species) %>%
  summarise(avg = mean(Sepal.Width)) %>%
  arrange(desc(avg))
```

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

Tidy Data - A foundation for wrangling in R

In a tidy data set:

- Each **variable** is saved in its own **column**
- Each **observation** is saved in its own **row**

Tidy data complements R's **vectorized operations**. R will automatically preserve observations as you manipulate variables. No other format works as intuitively with R.

Reshaping Data - Change the layout of a data set

gather() and **spread()**

```
tidyr::gather(cases, "year", "n", 2:4)
```

Gather columns into rows.

```
tidyr::spread(pollution, size, amount)
```

Spread rows into columns.

tidyr::separate() and **tidyr::unite()**

```
tidyr::separate(storms, date, c("y", "m", "d"))
```

Separate one column into several.

```
tidyr::unite(data, col, ..., sep)
```

Unite several columns into one.

tidyr::data_frame()

```
tidyr::data_frame(a = 1:3, b = 4:6)
```

Combine vectors into data frame (optimized).

tidyr::arrange()

```
tidyr::arrange(mtcars, mpg)
```

Order rows by values of a column (low to high).

```
tidyr::arrange(mtcars, desc(mpg))
```

Order rows by values of a column (high to low).

```
tidyr::rename(tbl, y = year)
```

Rename the columns of a data frame.

Subset Observations (Rows)

dplyr::filter()

```
dplyr::filter(iris, Sepal.Length > 7)
```

Extract rows that meet logical criteria.

dplyr::distinct()

Remove duplicate rows.

dplyr::sample_frac()

```
dplyr::sample_frac(iris, 0.5, replace = TRUE)
```

Randomly select fraction of rows.

dplyr::sample_n()

```
dplyr::sample_n(iris, 10, replace = TRUE)
```

Randomly select n rows.

dplyr::slice()

```
dplyr::slice(iris, 10:15)
```

Select rows by position.

dplyr::top_n()

```
dplyr::top_n(storms, 2, date)
```

Select and order top n entries (by group if grouped data).

Subset Variables (Columns)

dplyr::select()

```
dplyr::select(iris, Sepal.Length, Petal.Length, Species)
```

Select columns by name or helper function.

Helper functions for select - select

- select_in()**: Select columns whose name contains a character string.
- select_at()**: Select columns with "Length" in the name.
- select_everything()**: Select columns whose name ends with a character string.
- select_everything()**: Select every column.
- select_matches()**: Select columns whose name matches a regular expression.
- select_num()**: Select columns named "num" or "y".
- select_named()**: Select columns named with a character string.
- select_one()**: Select one column.
- select_all()**: Select all columns.
- select_if()**: Select columns whose name starts with a character string.
- select_if_not()**: Select columns whose name does not start with a character string.
- select_if_else()**: Select columns whose name starts with a character string or does not.
- select_if_not_else()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else_if()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else_if_not()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else_if_not_else()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else_if_not_else_if()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else_if_not_else_if_not()**: Select columns whose name does not start with a character string or does not.
- select_if_not_else_if_not_else_if_not_else_if_not_else()**: Select columns whose name does

Moving from base R to the tidyverse



- [Blog post link!](#)

Base R command	Tidyverse Command	What it does and why you should use the tidyverse version	Comment
<code>read.csv()</code>	<code>read_csv()</code>	reads in a csv file, but its much faster, shows progress bar for large files, can automatically parse data types	also see <code>read_delim()</code> , <code>read_tsv()</code> and <code>readxl::read_xlsx()</code>
<code>sort()</code> , <code>order()</code>	<code>arrange()</code>	sort column(n) within a data frame	see also <code>order_by()</code>
<code>mtcars\$mpg = ...</code>	<code>mutate()</code>	modify a column	see also <code>transmute()</code> which drops existing variables
<code>mtcars[,c("mpg", "am")]</code> , <code>subset()</code>	<code>select()</code> , <code>rename()</code>	select or rename columns	see also <code>pull()</code>

MicrobioUoE

- Set of functions I commonly use, wrote into a package

Installation

Installation of packages from GitHub is relatively straightforward.

```
# firstly need to make sure devtools is installed
install.packages('devtools')

# install MicrobioUoE
devtools::install_github('padpadpadpad/MicrobioUoE')
```

Functions

- `bind_biolog_sheet()`
 - cleans and binds biolog data on a single sheet of an excel spreadsheet. Returns a dataframe where each substrate is a column and each row is a separate plate. The values in each column are the OD readings of that plate in that substrate.
- `bind_biolog_all()`
 - cleans and binds biolog data across multiple sheets of an excel spreadsheet. Returns a dataframe where each substrate is a column and each row is a separate plate. The values in each column are the OD readings of that plate in that substrate. Also includes columns for the sheet and plate of each row.

```
filename <- '~/Desktop/biolog_data.xlsx'
sheets <- paste0('S', 1:10, sep = '')

MicrobioUoE::bind_biolog_all(filename, sheets)
```

- `stock_sol_vol()`
 - calculates the amount of volume needed to dilute a stock solution to a required concentration and volume

```
MicrobioUoE::stock_sol_vol(stock_sol_conc = 0.75, new_sol_conc = 0.2, new_sol_vol = 1000)
#> [1] 266.6667
```

- `calc_norm_OD()`
 - takes an csv of OD readings and writes another csv of the required volumes to normalise the concentrations across samples

```
MicrobioUoE::calc_norm_OD('first_OD_readings.csv',
                           'df_volumes.csv',
                           new_sol_conc = 0.1,
                           new_sol_vol = 1000,
                           control = 0.035)
```

Cheers!

Components

