

purrr for biostatisticians



with examples

Daniel D. Sjoberg

Memorial Sloan Kettering Cancer Center
Department of Epidemiology and Biostatistics

February, 19 2019

purrr



www.rstudio.com

purrr's map functions

let's get started

purrr package

purrr enhances R's functional programming toolkit (read: the apply family of functions) by providing a complete and consistent set of tools for working with functions and vectors

best place to start is the family of `map` functions which allow you to replace many for loops* with code that is more succinct and easier to read

`map` functions transform their input by applying a function to each element and returning a vector the same length as the input

[*] And much much more

purrr::map vs base::apply

base apply is to _____ as purrr map is to _____ ?

purrr::map vs base::apply

base apply is to

as purrr map is to



www.impawards.com



base::apply

- first argument to `apply` is the data; the first argument to `MARGIN` is the function
- no consistent way to pass additional arguments; most use `DATA`, `MARGIN` uses `DATA`, and some require you to create a new anonymous function
- output from `apply` is not consistent

base::apply

- first argument to `apply` is the data; the first argument to `lapply` is the function
- no consistent way to pass additional arguments; most use `...args`, `lapply` uses `parlargs`, and some require you to create a new anonymous function
- output from `apply` is not consistent

- `lapply`, `lapply`, and `lapply` uses `lapply` to suppress names in output; `lapply` does not have this argument



purrr::map

- the `map` family has greater consistency among functions
- `map`, `map2`, and `map2_chr` inputs are the same and allow for flexible input
- consistent methods for passing additional arguments

purrr::map

- the `map` family has greater consistency among functions
- `map`, `map2`, and `map3` inputs are the same and allow for flexible input
- consistent methods for passing additional arguments

- the output from the map family of functions is predictable and easily modifiable



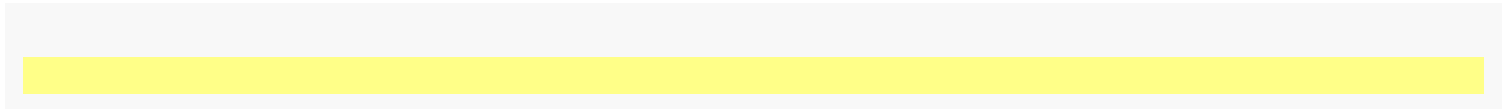
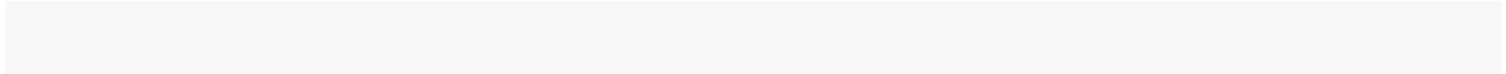
use cases

1. subgroup analyses
2. sensitivity analyses
3. reading all files in a folder
4. bootstrap analyses
5. other purrr package functions

usage



usage



usage

pass a function name to

additional function arguments can be passed as well

usage

pass a function name to

additional function arguments can be passed as well

create a new function with

usage

pass a function name to

additional function arguments can be passed as well

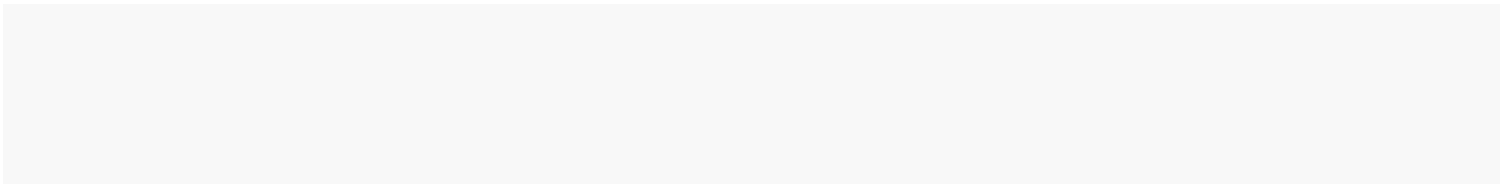
create a new function with

use the shortcut to create a function

usage



usage

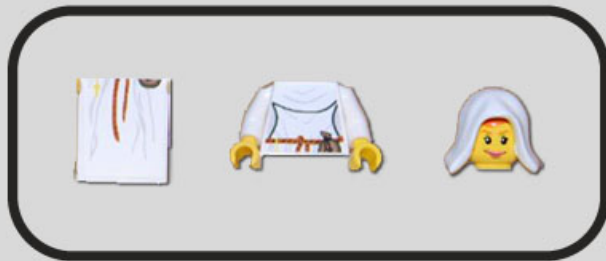


usage

pmap (.l, embody)



pmap (.l, embody)



trial dataset

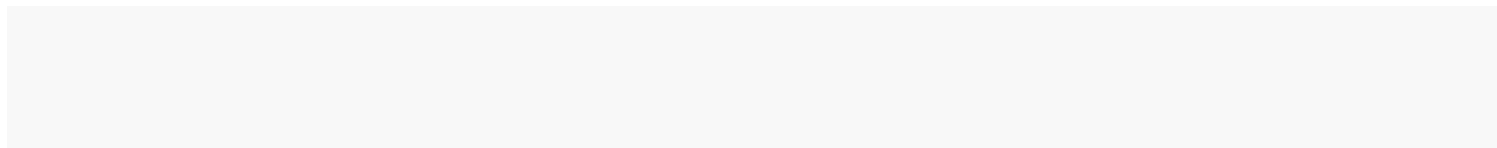
		N = 107	N = 93
Age, yrs	192	47 (39, 58)	46 (36, 54)
Marker Level, ng/mL	192	0.61 (0.22, 1.20)	0.72 (0.22, 1.63)
T Stage	200		
T1		25 (23%)	26 (28%)
T2		26 (24%)	23 (25%)
T3		29 (27%)	13 (14%)
T4		27 (25%)	31 (33%)
Grade	200		
I		38 (36%)	29 (31%)
II		34 (32%)	24 (26%)
III		35 (33%)	40 (43%)
Tumor Response	191	52 (51%)	30 (33%)

use cases

- 1.
2. sensitivity analyses
3. read all files in a folder
4. bootstrap analyses
5. other purrr package functions

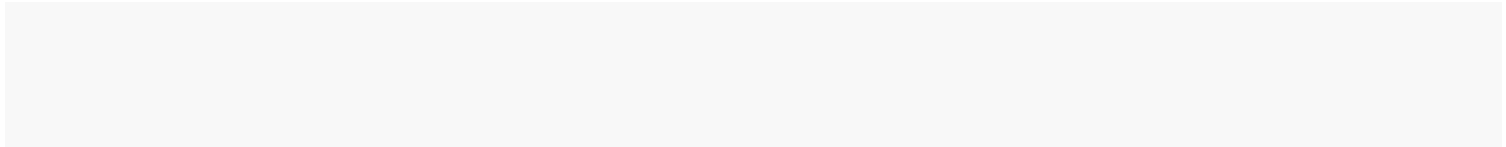
subgroup analysis

`tidyr::nest + purrr::map`



subgroup analysis

`tidyr::nest + purrr::map`



tibbles share the same structure as data frames

tibbles are a list of vectors, and it is possible to have a list column

very useful because a list can contain any other object: this means you can put any object in a tibble!

allows you to keep related objects together in a row, no matter how complex the individual objects are

subgroup analysis

`tidyr::nest + dplyr::mutate + purrr::map`

subgroup analysis

`tidyr::nest + dplyr::mutate + purrr::map`

subgroup analysis

`tidyr::nest + dplyr::mutate + purrr::map_dbl`

output types

the default output of `map()` is a list

we can coerce the output type with

<code>map()</code>	list
<code>map_dbl()</code>	double
<code>map_int()</code>	integer
<code>map_lgl()</code>	logical
<code>map_dfr()</code>	tibble (bind_rows)
<code>map_df()</code>	tibble (bind_cols)

when using the `map_*()` functions, `map_dfr()` runs as it typically would with the added step of coercing the output at the end

tip: make sure your code works with `map_dfr()` before adding `map_dfr()`.

use cases

1. subgroup analyses
- 2.
3. read all files in a folder
4. bootstrap analyses
5. other purrr package functions

sensitivity analyses

run your analysis among

- all patients ()
- excluding low grade patients ()

TRUE

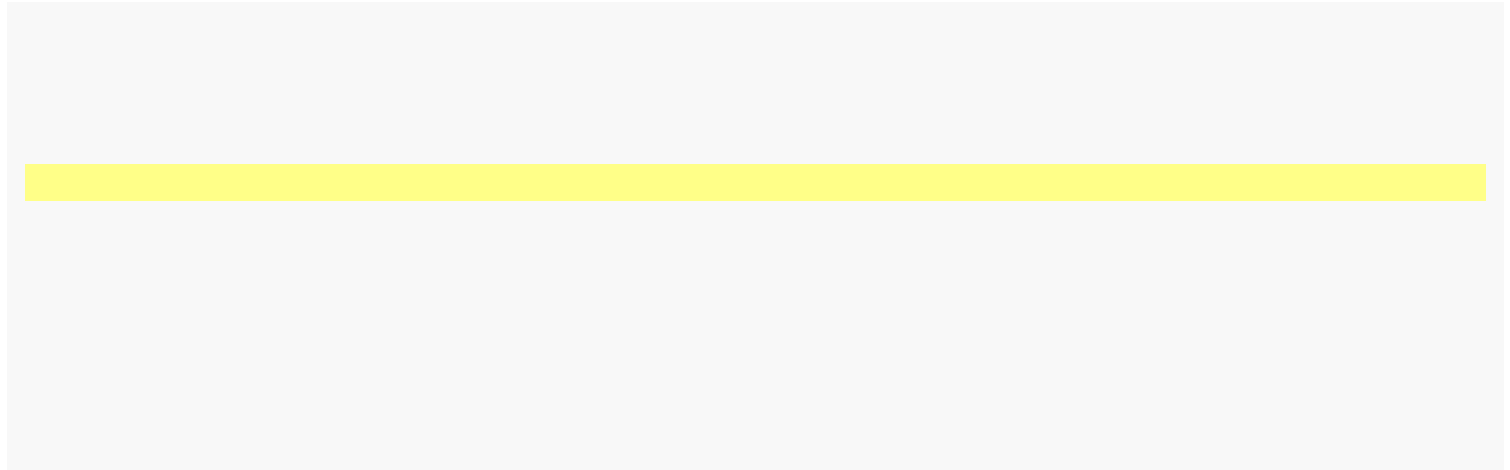
grade != 'I'

sensitivity analyses

run your analysis among

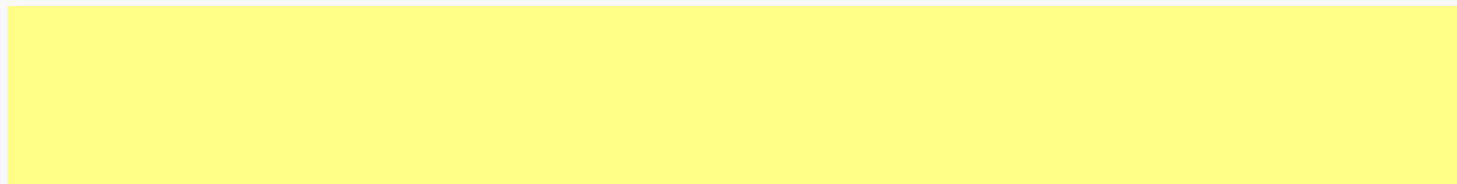
- all patients ()
- excluding low grade patients ()

sensitivity analyses



sensitivity analyses

you can also save figures in a tibble



use cases

1. subgroup analyses
2. sensitivity analyses
- 3.
4. bootstrap analyses
5. other purrr package functions

read files

store vector of the files you want to import



read files

store vector of the files you want to import

use `read_csv()` to read the files.

- returns a list where each element is a tibble

read files

append each of the data sets with the

- after files have been imported,

function

will create one final tibble

read files

append each of the data sets with the

- after files have been imported,

function

will create one final tibble

include an identifier with a piped

use cases

1. subgroup analyses
2. sensitivity analyses
3. read all files in a folder
- 4.
5. other purrr package functions

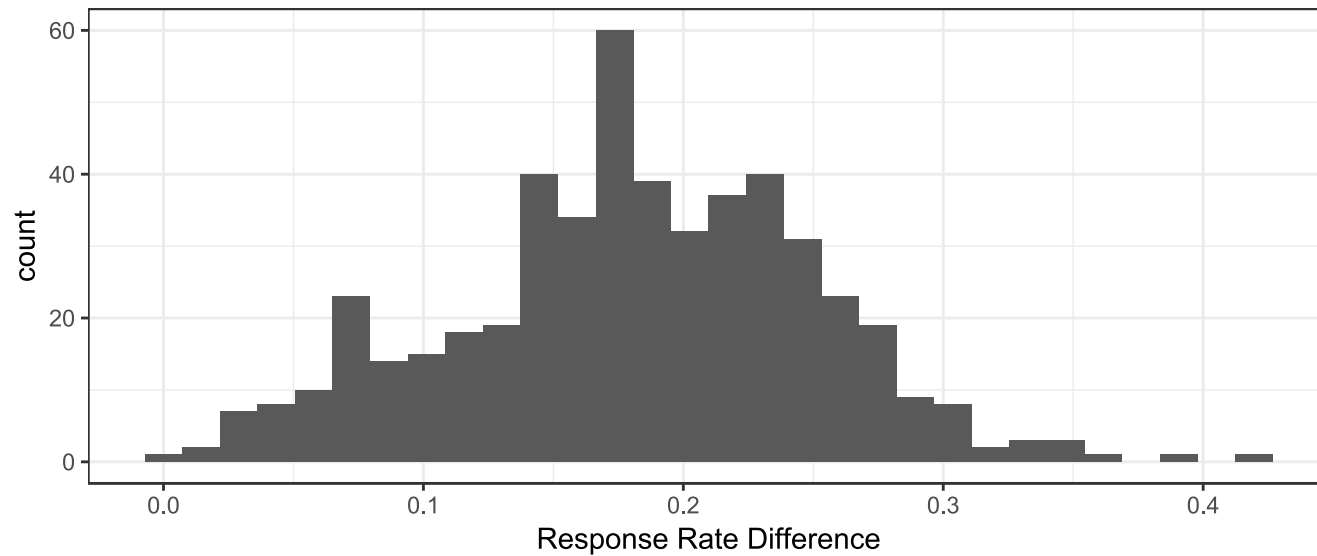
bootstrap analyses

use bootstrap re-sampling to estimate the difference in response rate by treatment

we'll use 500 re-sampled data sets to estimate the standard deviation of the response rate difference

assuming normality of the response rate difference, construct a 95% confidence interval for the difference

bootstrap analyses



bootstrap analyses

the result

- 18% (95% CI 4.6%, 32%)

bootstrap analyses

the result

- 18% (95% CI 4.6%, 32%)

how do the results compare to the the Wald CI?

bootstrap analyses

the result

- 18% (95% CI 4.6%, 32%)

how do the results compare to the the Wald CI?

- 18% (95% CI 4.1%, 31%)



use cases

1. subgroup analyses
2. sensitivity analyses
3. read all files in a folder
4. bootstrap analyses
- 5.

other purrr functions

-
-
- and
- and
-
- , and



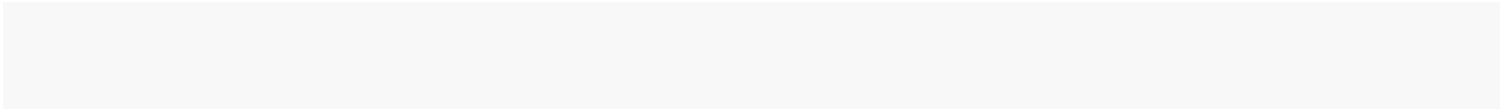
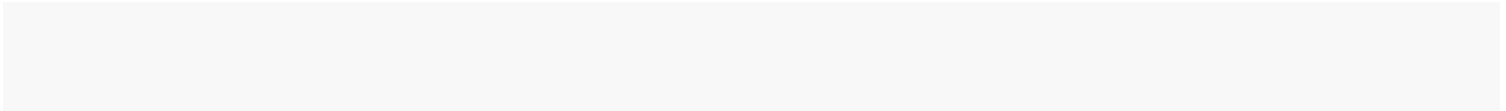
Yes, Lauren Hill was in !

unlike `map()` and its variants which always return a fixed object type (list for `map()`, integer vector for `map_int()`, etc), the `map_*()` family always returns the same type as the input object

names, or is short hand for if it does not if has

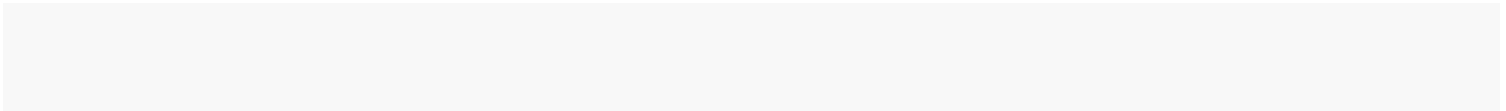
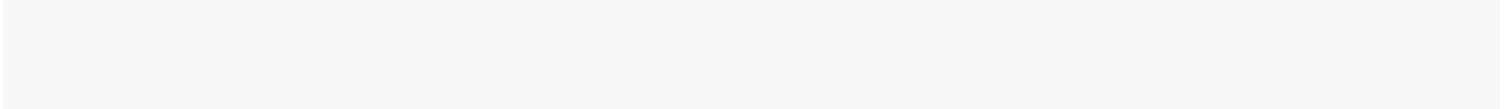


names, or is short hand for if it does not if has



and

keep or discard elements of a list or vector



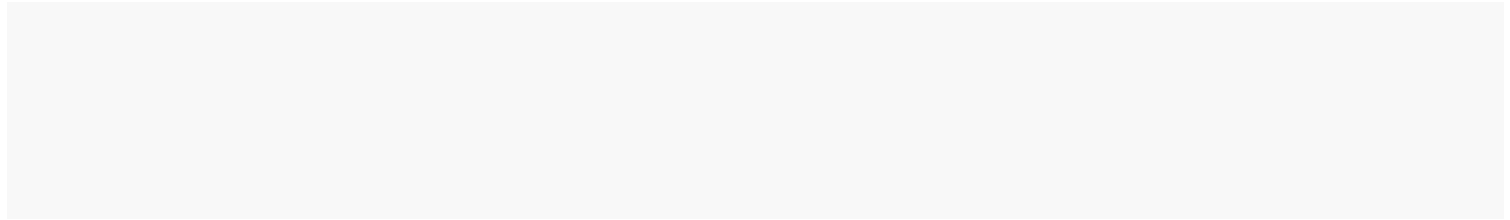
pluck is similar to `[[]]` and selects a single element from a list or vector

use position or name to select item

pluck is easier to read when used with the pipe (`%>%`)

like

without the factors...finally!



check out , ,

they are similar, but return lists rather than a tibble

, , and

these functions wrap functions so that instead of generating side effects through printed output, messages, warnings, and errors, they return enhanced output

similar to and

wrapped function returns a list with components and

wrapped function instead returns a list with components
, , and

wrapped function uses a default value (otherwise) whenever an error occurs

done!



questions?

▣ slides available at danieldsjoberg.com/purrr-for-biostatisticians

🔗 souce code available at github.com/ddsjoberg/purrr-for-biostatisticians