

Przetwarzanie i wizualizacja danych

prof. UAM dr hab. Tomasz Górecki

tomasz.gorecki@amu.edu.pl

Zakład Statystyki Matematycznej i Analizy Danych
Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza w Poznaniu



Podstawowe funkcje to: `Sys.time` (data wraz z godziną), `Sys.Date` (data bez godziny). Do wprowadzania danych jako dat służy funkcja `as.Date`, której argumentem jest data. Domyślny format daty, to cztery cyfry na rok, dwie na miesiąc i dwie na dzień, oddzielone kreską lub ukośnikiem. Jeśli chcemy użyć niestandardowego formatu, należy go wyspecyfikować jako wartość parametru `format` według oznaczeń zawartych w poniższej tabeli.

Symbol	Meaning	Example
%a	Abbreviated weekday name	Tue
%A	Full weekday name	Tuesday
%b	Abbreviated month name	Apr
%B	Full month name	April
%C	Century: the integer part of the year divided by 100	20
%d	Day of the month	09
%H	Hours as decimal number (00–23)	13
%I	Hours as decimal number (01–12)	1
%m	Month as number (01–12)	04
%M	Minute as number (00–59)	12
%p	AM/PM indicator for 12-hour time (%I)	PM
%S	Second as integer (00–61)	12
%u	Weekday as a decimal number (1–7, Monday is 1)	2
%w	Weekday as decimal number (0–6, Sunday is 0)	2
%y	2-Digit Year (00–99)	19
%Y	4-Digit Year	2019

Data przechowywana jest jako liczba dni, jaka upłynęła od 1 stycznia 1970 roku. Można również podać datę jako liczbę dni, która upłynęła od pewnej daty początkowej. Jeśli chcemy się dowiedzieć, jakim dniem, miesiącem lub kwartałem jest dana data możemy użyć funkcji **weekdays**, **months** oraz **quarters**. Często możemy być zainteresowani jaka była różnica pomiędzy dwoma datami. W R różnicę tę możemy wyrazić w sekundach, minutach, godzinach, dniach i miesiącach używając funkcji **difftime** i określając parametr **units** na *secs*, *mins*, *hours*, *days*, *weeks* odpowiednio. Przy konstrukcji szeregów czasowych potrzebne nam są sekwencje dat. Można je z łatwością utworzyć korzystając z poznanej wcześniej funkcji **seq** z wykorzystaniem jej parametru **by**, który może przyjmować wartości będące jednostkami czasowymi.

Pakiet **lubridate** służy do pracy z datami i stanowi część ekosystemu **tidyverse**.



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.

```
d <- as_datetime(2511870400)
## "2017-11-28 12:00:00 UTC"
```

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

- Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data.
- 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|_ymd_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(), mdyd(), 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")
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 strptime.
fast_strptime("3/1/01", "%m/%d/%Y")
```

```
parse_date_time() Easier strptime.
parse_date_time("3/1/01", "ymd")
```

2017-11-28

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

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

12:00:00

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

```
t <- hms::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.

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(x) ISO 8601 week.
epiweek(x) 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)
```

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



2017-11-28 12:00:00



2017-11-28 12:00:00



2017-11-28 12:00:00

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

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

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

```
change_on_boundary = NULL)
Roll back to last day of previous month. rollback(dt)
```

Stamp Date-times

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

- Derive a template, create a function `if = stamp("Created Sunday, Jan 17, 1999 3:34")`
- Apply the template to dates `stamp("2010-04-05")`

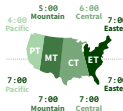
Tip: use a date with day > 15

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()**



```
with_tz(time, tzone = "") Get the same date-time in a new time zone (a new clock-time).
```

```
with_tz(dt, "US/Pacific")
```

```
force_tz(time, tzone = "") Get the same clock time in a new time zone (a new date-time).
```

```
force_tz(dt, "US/Pacific")
```

R Studio

RStudio® is a trademark of RStudio, Inc. • CC BY SA RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more at lubridate.tidyverse.org • lubridate 1.6.0 • Updated: 2017-12



Math with Date-times

lubridate provides three classes of timesteps 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:

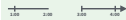
A normal day

```
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")
```



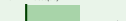
Leap years and leap seconds

```
leap ~ ymd("2019-03-01")
```



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

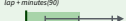
```
nor + minutes(90)
```



```
gap + minutes(90)
```



```
lap + minutes(90)
```



```
leap + years(1)
```



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

```
nor ~ 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(nor, nor + 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(2018|01|31)
jan31 + months(1)
## NA
```

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

```
jan31 %m-%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"
```



- years(x = 1)** x years.
- months(x = 1)** 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 nanoseconds.
- picoseconds(x = 1)** x picoseconds.

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

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

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. **Durations** 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
"142086000s (~2 weeks)"
```



- days(x = 1)** 31536000x seconds.
- weeks(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(S, unit = "years")

as.duration(x, ...) Coerce a timespan to a duration. Also **is.duration()**, **is.duration()**, **as.duration()**

make_diffime(x) Make diffime with the specified number of units.
make_diffime(S9999)

INTERVALS

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

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

```
l <- interval(ymd("2017-01-01"), d)
l
#> 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 timespan to an interval with the start date-time. Also **is.interval()**, **as.interval(i, start = now())**



Pakiet **stringr** służy do pracy z napisami (tekstem) i stanowi część ekosystemu **tidyverse**. Bazuje on na pakiecie `stringi`, którego autorem jest prof. Marek Gągolewski.





Work with strings with stringr : : CHEAT SHEET

The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.

Detect Matches



str_detect(string, pattern) Detect the presence of a pattern match in a string. `str_detect("fruit", "o")`



str_which(string, pattern) Find the indexes of strings that contain a pattern match. `str_which("fruit", "o")`



str_count(string, pattern) Count the number of matches in a string. `str_count("fruit", "o")`



str_locate(string, pattern) Locate the positions of pattern matches in a string. Also **str_locate_all**(string, pattern)

Subset Strings



str_sub(string, start = 1L, end = -1L) Extract substrings from a character vector. `str_sub("fruit", 1, 3); str_sub("fruit", -2)`



str_subset(string, pattern) Return only the strings that contain a pattern match. `str_subset("fruit", "o")`



str_extract(string, pattern) Return the first pattern match found in each string, as a vector. Also **str_extract_all** to return every pattern match. `str_extract("fruit", "[aeoiu]*")`



str_match(string, pattern) Return the first pattern match found in each string, as a matrix with a column for each () group in pattern. Also **str_match_all**. `str_match("sentences", "[a](the) ([^]+)"`)

Manage Lengths



str_length(string) The width of strings (i.e. number of code points, which generally equals the number of characters). `str_length("fruit")`



str_pad(string, width, side = c("left", "right", "both"), pad = " ") Pad strings to constant width. `str_pad("fruit", 17)`



str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...") Truncate the width of strings, replacing content with ellipsis. `str_trunc("fruit", 3)`



str_trim(string, side = c("both", "left", "right")) Trim whitespace from the start and/or end of a string. `str_trim("fruit")`

Mutate Strings



str_sub() <- value. Replace substrings by identifying the substrings with `str_sub()` and assigning into the results. `str_sub("fruit", 1, 3) <- "str"`



str_replace(string, pattern, replacement) Replace the first matched pattern in each string. `str_replace("fruit", "o", "y")`



str_replace_all(string, pattern, replacement) Replace all matched patterns in each string. `str_replace_all("fruit", "o", "y")`



str_to_lower(string, locale = "en") Convert strings to lower case. `str_to_lower("sentences")`



str_to_upper(string, locale = "en") Convert strings to upper case. `str_to_upper("sentences")`



str_to_title(string, locale = "en") Convert strings to title case. `str_to_title("sentences")`

Join and Split



str_c(), `sep = ""`, `collapse = NULL` Join substrings into a single string. `str_c("letters", LETTERS)`



str_c(), `sep = ""`, `collapse = NULL` Collapse a vector of strings into a single string. `str_c("letters", collapse = "")`



str_dup(string, times) Repeat strings times times. `str_dup("fruit", times = 2)`



str_split_fixed(string, pattern, n) Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also **str_split** to return a list of substrings. `str_split_fixed("fruit", "", n = 2)`



glue(..., `sep = ""`, `envir = parent.frame()`, `open = "{"`, `close = "}"`) Create a string from strings and (expressions) to evaluate. `glue("glue{Pi is {pi}}")`



glue_data(x, ..., `sep = ""`, `envir = parent.frame()`, `open = "{"`, `close = "}"`) Use a data frame, list, or environment to create a string from strings and (expressions) to evaluate. `glue_data(mtcars, "rownames{mtcars} has {hp} hp")`

Order Strings



str_order(x, decreasing = FALSE, `na_last = TRUE`, locale = "en", numeric = FALSE, ...) Return the vector of indexes that sorts a character vector. `x[str_order(x)]`



str_sort(x, decreasing = FALSE, `na_last = TRUE`, locale = "en", numeric = FALSE, ...) Sort a character vector. `str_sort(x)`

Helpers



str_conv(string, encoding) Override the encoding of a string. `str_conv("fruit", "ISO-8859-1")`



str_view(string, pattern, match = NA) View HTML rendering of first regex match in each string. `str_view("fruit", "[aeoiu]*")`



str_view_all(string, pattern, match = NA) View HTML rendering of all regex matches. `str_view_all("fruit", "[aeoiu]*")`



str_wrap(string, width = 80, indent = 0, `exact = 0`) Wrap strings into nicely formatted paragraphs. `str_wrap("sentences", 20)`

↑ See bit.ly/ISO639-1 for a complete list of locales.



Need to Know

Pattern arguments in `stringr` are interpreted as regular expressions after any special characters have been parsed.

In R, you write regular expressions as strings, sequences of characters surrounded by quotes ("") or single quotes ('').

Some characters cannot be represented directly in an R string. These must be represented as special characters, sequences of characters that have a specific meaning, e.g.

Special Character	Represents
\\	backslash
\"	double quote
\"'	single quote
\\n	new line

Run `?""` to see a complete list

Because of this, whenever a `\` appears in a regular expression, you must write it as `\\` in the string that represents the regular expression.

Use `writeLines()` to see how R views your string after all special characters have been parsed.

```
writeLines("i.")
# [1]
```

```
writeLines("\\") # is a backslash
# [1]
```

INTERPRETATION

Patterns in stringr are interpreted as regexes. To change this default, wrap the pattern in one of:

`regex(pattern, ignore_case = FALSE, multiline = FALSE, comments = FALSE, dotall = FALSE, ...)`
 Modifies a regex to ignore cases, match end of lines as well as end of strings, allow R comments within regex's, and/or to have . match everything including \n.
`str_detect("T", regex("T", TRUE))`

`fixed()` Matches raw bytes but will miss some characters that can be represented in multiple ways (fast).
`str_detect("u0130", coll("r", TRUE, locale = "tr"))`

`coll()` Matches raw bytes and will use locale specific collation rules to recognize characters that can be represented in multiple ways (slow).
`str_detect("u0130", coll("r", TRUE, locale = "tr"))`

`boundary()` Matches boundaries between characters, line_breaks, sentences, or words.
`str_split(sentences, boundary("word"))`

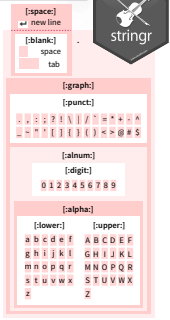
Regular Expressions -

Regular expressions, or regexps, are a concise language for describing patterns in strings.

MATCH CHARACTERS

string type	regex	matches	example
this	(to mean this)	(which matches this)	
	a (etc.)	a (etc.)	see("a")
\\.	.	any character	abc ABC 123 .?@0
\\[[any character in brackets	abc ABC 123 [?@0]
\\		any character in pipe	abc ABC 123 ?@0
\\{	{	any character in curly braces	abc ABC 123 {?@0}
\\((any character in parentheses	abc ABC 123 (?@0)
\\[[any character in square brackets	abc ABC 123 [?@0]
\\((any character in parentheses	abc ABC 123 (?@0)
\\n	\n	new line (return)	abc ABC 123 .?@0
\\t	\t	tab	abc ABC 123 .?@0
\\s	\s	any whitespace (\S for non-whitespaces)	abc ABC 123 .?@0
\\d	\d	any digit (\D for non-digits)	abc ABC 123 .?@0
\\w	\w	any word character (\W for non-word chars)	abc ABC 123 .?@0
\\b	\b	word boundaries	abc ABC 123 .?@0
	[digit]	digits	see("[digit]")
	[alpha]	letters	see("[alpha]")
	[lower]	lowercase letters	see("[lower]")
	[upper]	uppercase letters	see("[upper]")
	[alnum]	letters and numbers	see("[alnum]")
	[punct]	punctuation	see("[punct]")
	[graph]	letters, numbers, and punctuation	see("[graph]")
	[space]	space characters (i.e. \s)	see("[space]")
	[blank]	space and tab (but not new line)	see("[blank]")
	[^]	every character except a new line	see("[^]")

* Many base R functions require classes to be wrapped in a second set of [], e.g. `[[:digit:]]`



ALTERNATES

regex	matches	example
 	or	alt("a d")
[a b]	one of	alt("a b")
[a b c]	anything but	alt("a b c")
[a-c]	range	alt("a-c")

QUANTIFIERS

regex	matches	example
*	zero or more	quant("a")
+	one or more	quant("a+")
?	exactly n	quant("a{2}")
{n}	n or more	quant("a{2,}")
{n,m}	between n and m	quant("a{2,4}")

ANCHORS

regex	matches	example
^	start of string	anchor("a")
\$	end of string	anchor("a\$")

LOOK AROUNDS

regex	matches	example
(?=)	followed by	look("a(?=c)")
(?!)	not followed by	look("a(?!=c)")
(?=)	preceded by	look("(?=b)a")
(?!)	not preceded by	look("(?!b)a")

GROUPS

Use parentheses to set precedence (order of evaluation) and create groups

regex	matches	example
(a b)c	set precedence	alt("a b)c" abcdc

Use an escaped number to refer to and duplicate parentheses groups that occur earlier in a pattern. Refer to each group by its order of appearance

string (type this)	regex (to mean this)	matches (which matches this)	example (the result is the same as ref("abba"))
\\1	{n} (etc.)	first () group, etc.	ref("a(b)(\\1)\\1") abbaab



Pakiet **purrr** zastępuje i ulepsza funkcje z rodziny `apply`, stanowi część ekosystemu **tidyverse**.

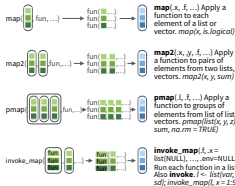


Apply functions with purrr : : CHEAT SHEET



Apply Functions

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



imap $[x, f, \dots]$ Apply function to each list-element of a list or vector.
imap $[x, f, \dots]$ Apply f to each element of a list or vector and its index.

OUTPUT

map(), **map2()**, **pmap()**, **imap** and **invoke_map** each return a list. Use a suffixed version to return the results as a specific type of flat vector: e.g. **map2_chr**, **map_dbl**, **map_lgl**, etc.

Use **walk**, **walk2**, and **walk3** to trigger side effects. Each return its input invisibly.

SHORTCUTS - within a purrr function:

"name" becomes **function(x) x["name"]**
 e.g. **map(x, "a")** extracts a from each element of x

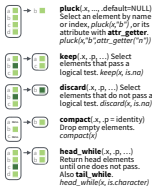
-x becomes **function(x) x**
 e.g. **map(1, -2 + x)** becomes **map(function(x) 2 + x)**

-x.y becomes **function(x, y) x.y**, e.g. **map2(l, p, -x.y)** becomes **map2(l, p, function(l, p) l + p)**

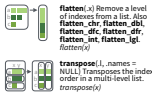
-1..2 etc becomes **function(x) x[1..2]** etc, e.g. **map(list(a, b, c), -3 + 1..2 etc)** becomes **pmap(list(a, b, c), function(a, b, c) c + a + b)**

Work with Lists

FILTER LISTS



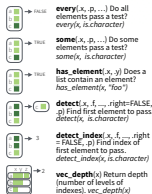
RESHAPE LISTS



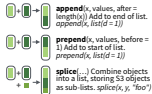
Reduce Lists



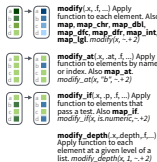
SUMMARISE LISTS



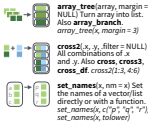
JOIN (TO) LISTS



TRANSFORM LISTS



WORK WITH LISTS



Modify function behavior

compose() Compose multiple functions.

lift() Change the type of input a function takes. Also **lift_dbl**, **lift_chr**, **lift_lgl**, **lift_int**, **lift_lgl**, **lift_int**, **lift_rn**. Return expression times.

negate() Negate a predicate function (a pipe friendly!)

partial() Create a version of a function that has some args preset to values.

safely() Modify fun to return list of results and errors.

quietly() Modify function to return list of results, output, messages, warnings.

possibly() Modify function to return default value whenever an error occurs (instead of error).





Nested Data

A nested data frame stores individual tables within the cells of a larger, organizing table.

nested data frame

Species	Site
setosa	setosa_001-01
versicolour	setosa_001-02
virginica	setosa_001-03

"cell" contents

Sepal.L	Sepal.W	Petal.L	Petal.W
5.1	3.5	1.4	0.2
4.9	3.0	1.4	0.2
4.7	3.2	1.3	0.2
4.6	3.1	1.5	0.2
5.0	3.6	1.4	0.2

n_iris\$data[[1]]

Species	Site	Sepal.L	Sepal.W	Petal.L	Petal.W
setosa	setosa_001-01	7.0	3.2	4.7	1.4
setosa	setosa_001-02	6.4	3.2	4.6	1.3
setosa	setosa_001-03	6.9	3.1	4.9	1.5
setosa	setosa_001-04	5.5	2.3	4.0	1.3
setosa	setosa_001-05	6.0	3.0	4.6	1.4

n_iris\$data[[2]]

Species	Site	Sepal.L	Sepal.W	Petal.L	Petal.W
setosa	setosa_002-01	7.1	3.0	5.0	1.7
setosa	setosa_002-02	6.3	2.9	5.0	1.6
setosa	setosa_002-03	6.4	3.0	5.6	1.9

n_iris\$data[[3]]

Species	Site	Sepal.L	Sepal.W	Petal.L	Petal.W
setosa	setosa_003-01	6.9	3.1	5.1	1.7
setosa	setosa_003-02	7.0	3.2	5.2	1.8
setosa	setosa_003-03	6.8	3.1	5.1	1.7

Use a nested data frame to:

- preserve relationships between observations and subsets of data
- manipulate many sub-tables at once with the **purrr** functions **map()**, **map2()**, or **pmap()**.

Use a two step process to create a nested data frame:

1. Group the data frame into groups with **dplyr::group_by()**
2. Use **nest()** to create a nested data frame with one row per group

```
iris %>% group_by(Species) %>% nest()
```

Species	data
setosa	[[1]] [[2]] [[3]] [[4]] [[5]]
versicolour	[[1]] [[2]] [[3]] [[4]] [[5]]
virginica	[[1]] [[2]] [[3]] [[4]] [[5]]

```
n_iris <- iris %>% group_by(Species) %>% nest()
```

```
tidy::nest(data, ..., key = data)
```

For grouped data, moves groups into cells as data frames.

Unnest a nested data frame with **unnest()**:

```
n_iris %>% unnest()
```

```
tidy::unnest(data, ..., drop = NA, id=NULL, sep=NULL)
```

Unnests a nested data frame.



List Column Workflow

Nested data frames use a **list column**, a list that is stored as a column vector of a data frame. A typical **workflow** for list columns:

1 Make a list column

```
tibble(iris) %>% group_by(Species) %>% nest()
```

Species	iris
setosa	[[1]] [[2]] [[3]] [[4]] [[5]]
versicolour	[[1]] [[2]] [[3]] [[4]] [[5]]
virginica	[[1]] [[2]] [[3]] [[4]] [[5]]

```
n_iris <- iris %>% group_by(Species) %>% nest()
```

2 Work with list columns

```
mod_fun <- function(df) {
  lm(Sepal.Length ~ ., data = df)
}
m_iris <- n_iris %>% mutate(model = map(mod_fun))
```

Species	iris	model
setosa	[[1]] [[2]] [[3]] [[4]] [[5]]	[[1]] [[2]] [[3]] [[4]] [[5]]
versicolour	[[1]] [[2]] [[3]] [[4]] [[5]]	[[1]] [[2]] [[3]] [[4]] [[5]]
virginica	[[1]] [[2]] [[3]] [[4]] [[5]]	[[1]] [[2]] [[3]] [[4]] [[5]]

```
mod_fun <- function(df) {
  lm(Sepal.Length ~ ., data = df)
}
m_iris <- n_iris %>% mutate(model = map(mod_fun))
```

3 Simplify the list column

```
b_fun <- function(mod) {
  coefficients(mod)[1]
}
m_iris %>% transmute(Species, beta = map_dbl(model, b_fun))
```

Species	beta
setosa	[[1]] [[2]] [[3]] [[4]] [[5]]
versicolour	[[1]] [[2]] [[3]] [[4]] [[5]]
virginica	[[1]] [[2]] [[3]] [[4]] [[5]]

```
b_fun <- function(mod) {
  coefficients(mod)[1]
}
m_iris %>% transmute(Species, beta = map_dbl(model, b_fun))
```

1. MAKE A LIST COLUMN - You can create list columns with functions in the **tibble** and **dplyr** packages, as well as **tidyr**'s **nest()**

- tibble::tibble(...)** Makes list column when needed

```
tibble(max = c(3, 4, 5), seq = list(1:3, 1:4, 1:5))
```
- dplyr::mutate(...)** Also **transmute()** Returns list col when result returns list.

```
mtcars %>% mutate(seq = map(cyl, seq))
```
- dplyr::summarise(...)** Returns list col when result is wrapped with **list()**

```
summarise(q = list(quantile(mpg)))
```

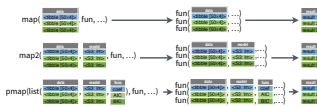
2. WORK WITH LIST COLUMNS - Use the purrr functions **map()**, **map2()**, and **pmap()** to apply a function that returns a single element-wise to the cells of a list column. **walk()**, **walk2()**, and **pwalk()** work the same way, but return a side effect.

- purrr::map(x, f, ...)** Apply f element-wise to x as f(x)

```
n_iris %>% mutate(n = map(data, dim))
```
- purrr::map2(x, y, f, ...)** Apply f element-wise to x and y as f(x, y)

```
m_iris %>% mutate(n = map2(data, model, list))
```
- purrr::pmap(l, f, ...)** Apply f element-wise to vectors saved in l

```
m_iris %>% mutate(n = pmap(list(data, model, data), list))
```



3. SIMPLIFY THE LIST COLUMN (into a regular column)

- purrr::map_lgl(x, f, ...)** Apply f element-wise to x, return a logical vector

```
n_iris %>% transmute(n = map_lgl(data, is.matrix))
```
- purrr::map_int(x, f, ...)** Apply f element-wise to x, return an integer vector

```
n_iris %>% transmute(n = map_int(data, nrow))
```
- purrr::map_dbl(x, f, ...)** Apply f element-wise to x, return a double vector

```
n_iris %>% transmute(n = map_dbl(data, nrow))
```
- purrr::map_chr(x, f, ...)** Apply f element-wise to x, return a character vector

```
n_iris %>% transmute(n = map_chr(data, nrow))
```