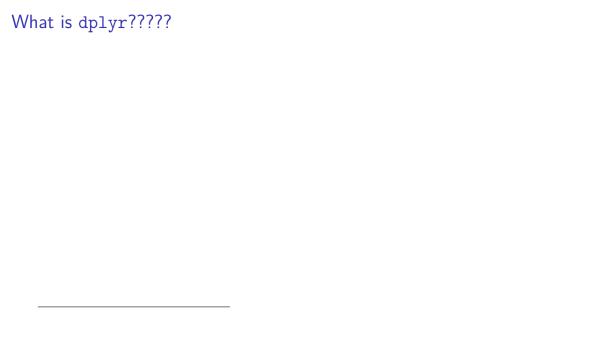
R Workshop Featuring dplyr Purdue Chapter of ASA

Jorge Loría

October 20th, 2020

Where can we find this presentation?

https://github.com/Jels95/Dplyr-Workshop



What is	dplyr?????
---------	------------

Dplyr is a package that permits a nice interface in R to work with data.frames.

What is dplyr?????

Dplyr is a package that permits a *nice* interface in R to work with data.frames.

Ok, nice, but what is a data.frame?

¹categorical

²https://vita.had.co.nz/papers/tidy-data.pdf

What is dplyr?????

Dplyr is a package that permits a *nice* interface in R to work with data.frames.

Ok, nice, but what is a data.frame?

A data frame is like a matrix of dimensions $n \times p$, where we have several different types of data. Each column corresponds to a single variable, and each variable has a specific type (numeric, string, logical, factor¹). Each row should correspond to a single observation².

¹categorical

²https://vita.had.co.nz/papers/tidy-data.pdf

What dataset are we going to use?

We are going to use a dataset from tidytuesday, about Himalayan Climbers

What dataset are we going to use?

We are going to use a dataset from tidytuesday, about Himalayan Climbers



Figure 1: Himalaya? It must be easy to survive there

How to access it?

```
install.packages('tidytuesdayR')
library(tidytuesdayR)
himalaya <- tidytuesdayR::tt_load('2020-09-22')
members <- himalaya$members</pre>
```

Let's start with questions

What questions can we ask this dataset that you think would be interesting to know? (Please write them in the chat!)

Let's start with questions

What questions can we ask this dataset that you think would be interesting to know? (Please write them in the chat!)

The idea is to be able to solve most, if not all of those questions by the end of the workshop!

Let's start with questions

What questions can we ask this dataset that you think would be interesting to know? (Please write them in the chat!)

The idea is to be able to solve most, if not all of those questions by the end of the workshop!

We will study some tools, and you should be able to answer them by the end of the workshop. If I can see that's not possible, I'll tell you what libraries are good to answer them. Ok?

Let's take a look:

#

#

```
members %>%
  head()
```

```
## # A tibble: 6 \times 21
##
     expedition id member id peak id peak name year season sex
                                                                       age ci
##
     <chr>
                    <chr>
                              <chr>
                                       <chr>
                                                 <dbl> <chr> <chr> <dbl> <chr> <dbl> <c
## 1 AMAD78301
                   AMAD7830~ AMAD
                                       Ama Dabl~
                                                  1978 Autumn M
## 2 AMAD78301
                   AMAD7830~ AMAD
                                       Ama Dabl~ 1978 Autumn M
```

40 Fr 41 Fr. ## 3 AMAD78301 AMAD7830~ AMAD Ama Dabl~ 1978 Autumn M 27 Fr Ama Dabl~ 1978 Autumn M ## 4 AMAD78301 AMAD7830~ AMAD 40 Fr AMAD7830~ AMAD ## 5 AMAD78301 Ama Dabl~ 1978 Autumn M 34 Fr ## 6 AMAD78301 AMAD7830~ AMAD Ama Dabl~ 1978 Autumn M

25 Fr ## # ... with 12 more variables: expedition_role <chr>, hired <lgl>, highpoint metres <dbl>, success <lgl>, solo <lgl>, oxygen used <lgl ## #

died <lgl>, death_cause <chr>, death_height_metres <dbl>, injured <</pre>

injury type <chr>, injury height metres <dbl>

%>%

This is a function that doesn't do much, but does a lot. It allows to compose functions (as the math people do) but in a way that permits an easy reading of the functions, and what is happening.

%>%

This is a function that doesn't do much, but does a lot. It allows to compose functions (as the math people do) but in a way that permits an easy reading of the functions, and what is happening. So, instead of writing: f(g(h(i(x))))

%>%

This is a function that doesn't do much, but does a lot. It allows to compose functions (as the math people do) but in a way that permits an easy reading of the functions, and what is happening. So, instead of writing: f(g(h(i(x)))), we write:

```
x %>%
i() %>%
h() %>%
g() %>%
f()
```

Which is arguably easier to read than the previous expression. Specially if some of those functions had extra arguments.

There are 6 main verbs in dplyr that we will study:

function	action
filter arrange	keeps rows that satisfy a condition sorts the rows following the order
select mutate	keeps/eliminates the columns by name creates new variables from existing variables
summarise	summarises the data
groub_by*	groups under specific conditions

All the verbs work very similarly (in dplyr, in English verbs are more confusing):

▶ first argument is a dataframe,

- first argument is a dataframe,
- arguments describe what to do with the dataframe, using the existing variables (without quotations).

- first argument is a dataframe,
- arguments describe what to do with the dataframe, using the existing variables (without quotations).
- output is a dataframe

- first argument is a dataframe,
- arguments describe what to do with the dataframe, using the existing variables (without quotations).
- output is a dataframe

All the verbs work very similarly (in dplyr, in English verbs are more confusing):

- ▶ first argument is a dataframe,
- arguments describe what to do with the dataframe, using the existing variables (without quotations).
- output is a dataframe

This structure allows us to concatenate (%>%) simple operations to obtain complex results.



Figure 2: Another type of filters

This functions removes rows that don't satisfy a (or several) condition that we specify. The arguments it receives are logical, and will use it to do that removal:

```
library(dplyr) ## load the library
members %>%
  filter(oxygen_used)
```

```
# A tibble: 18.233 x 21
##
      expedition id member id peak id
                                       peak name
                                                   vear season sex
                                                                        age
      <chr>
                     <chr>
                                                  <dbl> <chr> <chr> <dbl>
##
                               <chr>>
                                        <chr>
##
    1 ANN170101
                     ANN17010~ ANN1
                                        Annapurn~
                                                   1970 Spring M
                                                                         35
##
    2 ANN170102
                     ANN17010~ ANN1
                                        Annapurn~
                                                   1970 Spring M
                                                                         28
    3 ANN170102
                                                                         32
##
                     ANN17010~ ANN1
                                        Annapurn~
                                                   1970 Spring M
##
    4 ANN177301
                     ANN17730~ ANN1
                                        Annapurn~
                                                   1977 Autumn M
                                                                         37
##
    5 ANN177301
                     ANN17730~ ANN1
                                        Annapurn~
                                                   1977 Autumn M
                                                                         29
##
    6 ANN177301
                     ANN17730~ ANN1
                                        Annapurn~
                                                   1977 Autumn M
                                                                         28
##
    7 ANN178301
                     ANN17830~ ANN1
                                        Annapurn~
                                                   1978 Autumn F
                                                                         35
```

We can use several columns to filter, and can even modify them. Let's see what people older than 75 years **didn't** need to use oxygen

```
members %>%
filter(!oxygen_used,age > 75)
```

```
# A tibble: 27 \times 21
##
      expedition id member id peak id
                                        peak name
                                                    year season sex
                                                                         age
##
      <chr>>
                     <chr>>
                                <chr>>
                                        <chr>>
                                                   <dbl> <chr> <chr> <dbl> <dbl>
##
    1 GTMM93301
                     GTMM9330~ GTMM
                                        Gimmigel~
                                                    1993 Autumn M
                                                                          80
    2 GTME93301
                     GIME9330~ GIME
                                        Gimmigel~
                                                    1993 Autumn M
                                                                          80
##
    3 PUM005104
                                                                          76
##
                     PUMO0510~ PUMO
                                        Pumori
                                                    2005 Spring M
                                        Ama Dabl~
                                                                          76
##
    4 AMAD06347
                     AMADO634~ AMAD
                                                    2006 Autumn M
##
    5 AMAD07314
                     AMADO731~ AMAD
                                        Ama Dabl~
                                                    2007 Autumn M
                                                                          77
                                                                          77
##
    6 YAKA08201
                     YAKA0820~ YAKA
                                        Yakawa K~
                                                    2008 Summer F
                                                                          78
##
    7 AMAD08331
                     AMADO833~ AMAD
                                        Ama Dabl~
                                                    2008 Autumn M
##
    8 TASH10302
                     TASH1030~ TASH
                                        Tashi Ka~
                                                    2010 Autumn M
                                                                           76
```

##

#

A tibble: 1 x 21

We can also use several columns at once to do a filter. Let's see what climber(s?) died a little bit after getting injured:

```
members %>%
filter(death_height_metres > injury_height_metres)
```

injury type <chr>, injury height metres <dbl>

expedition id member id peak id peak name vear season sex age ci

▶ From those who where not injured, how many died in 1979?

- From those who where not injured, how many died in 1979?
- ► How many of those that succeeded, were doing it solo and got injured above 1000 meters?

- From those who where not injured, how many died in 1979?
- ► How many of those that succeeded, were doing it solo and got injured above 1000 meters?
- ▶ How many people died on their climb below 1000 meters?

- From those who where not injured, how many died in 1979?
- ► How many of those that succeeded, were doing it solo and got injured above 1000 meters?
- ▶ How many people died on their climb below 1000 meters?

- From those who where not injured, how many died in 1979?
- ► How many of those that succeeded, were doing it solo and got injured above 1000 meters?
- ▶ How many people died on their climb below 1000 meters?

filter saving

Let's now save the dataframe of the members that went on a expedition in 1905:

```
members1905 <- members %>%
filter(year == 1905)
```



arrange



Figure 3: Another type of arrangement



arrange

This verb sorts the data frame with the column(s) that we tell it to use

```
members %>%
arrange(year)
```

```
# A tibble: 76,519 x 21
##
      expedition id member id peak id
                                        peak name
                                                   vear season sex
                                                                        age
      <chr>>
                     <chr>
                               <chr>>
                                                  <dbl> <chr> <chr> <dbl>
##
                                        <chr>
##
    1 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         29
##
    2 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         NA
    3 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         NΑ
##
    4 KANG05201
                                                                         NΑ
##
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
##
    5 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         NA
##
    6 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         NΑ
##
    7 KANG05201
                     KANG0520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         36
##
    8 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         NΑ
##
    9 KANG05201
                     KANGO520~ KANG
                                        Kangchen~
                                                   1905 Summer M
                                                                         NA
```

arrange on characters

5 KANG05201

7 KANG05201

6 KANG05201

8 KANG05201

Let's go back to the 1905 dataset, and check how it orders when we use a string, instead of a number:

```
members1905 %>%
  arrange(member id)
```

## #	A tibble: 9 x	21			
##	expedition_id	member_id	peak_id peak_name	e year season sex	age

## #	A tibble: 9 x	21						
##	expedition_id	member_id	peak_id	peak_name	year	season	sex	age
##	<chr></chr>	<chr>></chr>	(chr)	<chr>></chr>	<dh1></dh1>	<chr>></chr>	<chr>></chr>	<dh1></dh1>

	_							
## #	A tibble: 9 x	21						
##	expedition_id	member_id	peak_id	peak_name	year	season	sex	age
	. 1 .	. 1	. 1	. 1	4 11 7 5	. 1 .	. 1 .	. 11 7 5

ci 1 KANG05201 KANGO520~ KANG Kangchen~ 1905 Summer M

<dbl> <c ## 2 KANG05201 KANGO520~ KANG 1905 Summer M

29 UK Kangchen~

36 Sw ## 3 KANG05201 KANGO520~ KANG 1905 Summer M NA Sw

Kangchen~

4 KANG05201 KANG0520~ KANG Kangchen~ NA Sw 1905 Summer M

Kangchen~

Kangchen~

Kangchen~

Kangchen~

1905 Summer M

1905 Summer M

1905 Summer M

1905 Summer M

NA It

NA Ne

NA Ne

NA Ne

KANGO520~ KANG

KANGO520~ KANG

KANGO520~ KANG

KANG0520~ KANG

arrange uses the lexicographic order, to sort when it encounters letters/characters, so: $a < A < b < B \cdots < z < Z$. If you're not sure if a number or a symbol is smaller than another one, you can try it out on the 'Console' in R for example type '1'<'a'. Is this true? Or false?

arrange uses the lexicographic order, to sort when it encounters letters/characters, so: $a < A < b < B \cdots < z < Z$. If you're not sure if a number or a symbol is smaller than another one, you can try it out on the 'Console' in R for example type '1'<'a'. Is this true? Or false? TRUE!

arrange uses the lexicographic order, to sort when it encounters letters/characters, so: $a < A < b < B \cdots < z < Z$. If you're not sure if a number or a symbol is smaller than another one, you can try it out on the 'Console' in R for example type '1'<'a'. Is this true? Or false? TRUE!How about '&' > '2'?

arrange uses the lexicographic order, to sort when it encounters letters/characters, so: $a < A < b < B \cdots < z < Z$. If you're not sure if a number or a symbol is smaller than another one, you can try it out on the 'Console' in R for example type '1'<'a'. Is this true? Or false? TRUE!How about '&' > '2'?FALSE!

arrange uses the lexicographic order, to sort when it encounters letters/characters, so: $a < A < b < B \cdots < z < Z$. If you're not sure if a number or a symbol is smaller than another one, you can try it out on the 'Console' in R for example type '1'<'a'. Is this true? Or false? TRUE!How about '&' > '2'?FALSE! Is ' ' < '1' (space less than one)?

arrange uses the lexicographic order, to sort when it encounters letters/characters, so: $a < A < b < B \cdots < z < Z$. If you're not sure if a number or a symbol is smaller than another one, you can try it out on the 'Console' in R for example type '1'<'a'. Is this true? Or false? TRUE!How about '&' > '2'?FALSE! Is ' ' < '1' (space less than one)? TRUE!

arrange on several columns

When we include several arguments, it first sorts the first one, then the second one *within* the first order, then the third one *within* the second one, and so on...

```
members %>%
arrange(desc(year),citizenship,hired,peak_name)
```

```
# A tibble: 76.519 x 21
      expedition id member id peak id
##
                                        peak name
                                                    vear season sex
                                                                         age
##
      <chr>>
                     <chr>>
                               <chr>>
                                        <chr>>
                                                   <dbl> <chr> <chr> <dbl> <dbl>
##
    1 AMAD19107
                     AMAD1910~ AMAD
                                        Ama Dabl~
                                                    2019 Spring M
                                                                          55
    2 EVER19116
                     EVER1911~ EVER
                                                    2019 Spring M
                                                                          53
##
                                        Everest
    3 EVER19117
                                                    2019 Spring M
                                                                          38
##
                     EVER1911~ EVER
                                        Everest
##
    4 EVER19117
                     EVER1911~ EVER
                                        Everest
                                                    2019 Spring M
                                                                          40
##
    5 EVER19148
                     EVER1914~ EVER
                                        Everest
                                                    2019 Spring M
                                                                          51
##
    6 KANG19102
                     KANG1910~ KANG
                                        Kangchen~
                                                    2019 Spring M
                                                                          43
##
    7 EVER19183
                     EVER1918~ EVER
                                                    2019 Spring M
                                                                          48
                                        Everest
##
    8 AMAD19106
                     AMAD1910~ AMAD
                                        Ama Dabl~
                                                    2019 Spring F
                                                                          33
```

arrange on transformations

A tibble: 1.713 x 21

5 HIML13301

6 KANG00106

##

##

First, let's only take those that got injured, and then see who got furthest away with respect to their injury. Is there anything weird going on? Is my code correct? Is the data correct?

```
members %>%
  filter(injured) %>%
  arrange(highpoint_metres - injury_height_metres)
```

```
expedition id member id peak id
##
                                       peak name
                                                  vear season sex
                                                                       age
##
      <chr>
                    <chr>>
                               <chr>>
                                       <chr>
                                                 <dbl> <chr> <chr> <dbl>
    1 CHOY02327
                    CHOY0232~ CHOY
                                       Cho Oyu
##
                                                  2002 Autumn M
                                                                        29
    2 MAKA15109
                                                                        50
##
                    MAKA1510~ MAKA
                                       Makalu
                                                  2015 Spring M
##
    3 MAKA15109
                    MAKA1510~ MAKA
                                       Makalu
                                                  2015 Spring M
                                                                        46
                                                  2003 Autumn M
##
    4 CHOY03301
                    CHOY0330~ CHOY
                                       Cho Oyu
                                                                        53
```

Himlung ~

Kangchen~

2013 Autumn M

2000 Spring M

44

35

HTML1330~ HTML

KANGOO10~ KANG



Exercises for arrange:

► From those that died, sort by those who got the highest before dying w.r.t. where they died.

Exercises for arrange:

- ► From those that died, sort by those who got the highest before dying w.r.t. where they died.
- ➤ Sort by year and then season. Does this make sense? If it doesn't, how could we fix it?

Exercises for arrange:

- ► From those that died, sort by those who got the highest before dying w.r.t. where they died.
- ➤ Sort by year and then season. Does this make sense? If it doesn't, how could we fix it?
- ▶ If you arrange by died, can you tell how arrange interprets the logicals?



select

This verbs selects the columns that we want to keep. Sometimes, we only need a couple of variables, and it's good to get rid of the rest:

```
members %>%
  select(year,sex,citizenship)

## # A tibble: 76,519 x 3
## year sex citizenship
```

```
##
     <dbl> <chr> <chr>
  1 1978 M
##
                 France
## 2 1978 M
                 France
   3 1978 M
##
                 France
##
   4 1978 M
                 France
##
   5 1978 M
                 France
##
   6 1978 M
                 France
## 7 1978 M
                 France
##
      1978 M
                 France
```

select ranges

Sometimes we want to select all columns between two other columns, so we can use the colon (:, not the organ) to do this:

select ranges

Sometimes we want to select all columns between two other columns, so we can use the colon (:, not the organ) to do this:

```
members %>%
select(age:success)
```

```
## # A tibble: 76.519 x 6
##
        age citizenship expedition role hired highpoint metres success
##
      <dbl> <chr>
                         <chr>>
                                          <lgl>
                                                            <dbl> <lgl>
         40 France
                         Leader
                                          FALSE
                                                               NA FALSE
##
         41 France
                                          FALSE
                                                             6000 FALSE
##
    2
                         Deputy Leader
    3
                                          FALSE
                                                               NA FALSE
##
         27 France
                         Climber
##
    4
         40 France
                         Exp Doctor
                                          FALSE
                                                             6000 FALSE
##
    5
         34 France
                         Climber
                                          FALSE
                                                               NA FALSE
##
    6
         25 France
                         Climber
                                          FALSE
                                                             6000 FALSE
##
         41 France
                         Climber
                                          FALSE
                                                             6000 FALSE
##
    8
         29 France
                         Climber
                                          FALSE.
                                                             6000 FALSE
```

Eliminating with select

##

8 Ama Dabl~

Other times we want to get rid of specific variables. For this, we can use the - (minus) symbol.

```
members %>%
  select(-expedition_id,-member_id,-peak_id)
```

```
## # A tibble: 76 519 \times 18
```

		0,010							
##	$peak_name$	year	season	sex	age	${\tt citizenship}$	expedition	role	hire
##	<chr></chr>	<dbl></dbl>	<chr></chr>	<chr></chr>	<dbl></dbl>	<chr></chr>	<chr></chr>		<lgl< th=""></lgl<>

##		peak_name	year	season	sex	age	Citizenship	expedition_role	nirea
##		<chr></chr>	<dbl></dbl>	<chr></chr>	<chr></chr>	<dbl></dbl>	<chr></chr>	<chr></chr>	<1g1>
##	1	Ama Dabl~	1978	${\tt Autumn}$	M	40	France	Leader	FALSE

		Pouri_mamo	your	Doabon	0011	ago	ототдопритр	onpour or on_roro	
##	<chr></chr>		<dbl></dbl>	<chr></chr>	<chr></chr>	<dbl></dbl>	<chr></chr>	<chr></chr>	<lg1< th=""></lg1<>
##	1	Ama Dabl~	1978	${\tt Autumn}$	M	40	France	Leader	FALS
44	0	Ama Dabl.	1070	A +	M	11	Emamaa	Domiter I and an	TATO

##	Т	Ama	Dabi~	19/0	Autumn	ľľ	40	France	Leader	FALSE
##	2	\mathtt{Ama}	Dabl~	1978	${\tt Autumn}$	M	41	France	Deputy Leader	FALSE
##	3	\mathtt{Ama}	Dabl~	1978	${\tt Autumn}$	M	27	France	Climber	FALSE
##	4	Ama	Dabl~	1978	${\tt Autumn}$	M	40	France	Exp Doctor	FALSE
##	5	Ama	Dabl~	1978	Autumn	M	34	France	Climber	FALSE

##	4	Ama	Dabl~	1978	Autumn	M	40	France	Exp Doctor	FALSE
##	5	Ama	Dabl~	1978	Autumn	M	34	France	Climber	FALSE
##	6	Ama	Dabl~	1978	Autumn	M	25	France	Climber	FALSE
##	7	Ama	Dabl~	1978	Autumn	M	41	France	Climber	FALSE

29 France

Climber

FALSE

1978 Autumn M

##	1 Ama	Dab1~	1978	Autumn	ľ	40	France	Leader	FAL
##	2 Ama	Dabl~	1978	${\tt Autumn}$	M	41	France	Deputy Leader	FALS
##	3 Ama	Dabl~	1978	${\tt Autumn}$	M	27	France	Climber	FALS
##	4 Ama	Dabl~	1978	${\tt Autumn}$	M	40	France	Exp Doctor	FALS
##	5 Amo	Dahl~	1079	11+11mn	М	3/	Franco	Climbor	EVI (

More of select.

##

##

##

6 Ama Dablam Autumn M

7 Ama Dablam Autumn M

8 Ama Dablam Autumn M

We can also use the number of the column to indicate which columns to select, and combine it with the names.

```
members %>%
  select(4,6:10, highpoint metres)
```

```
## # A tibble: 76.519 x 7
##
     peak name season sex
                               age citizenship expedition role highpoint
```

<chr> <chr>> 1 Ama Dablam Autumn M Leader ## 40 France ## 2 Ama Dablam Autumn M 41 France Deputy Leader

27 France ## 3 Ama Dablam Autumn M Climber ## 4 Ama Dablam Autumn M 40 France Exp Doctor ## 5 Ama Dablam Autumn M 34 France Climber

25 France

41 France

29 France

Climber

Climber

Climber

<chr> <chr> <dbl> <chr>

There's a family of functions designed to work with select, so we can work more easily. Among them we have:

starts_with('some') that looks for all columns that begin with 'some'

- starts_with('some') that looks for all columns that begin with 'some'
- ends_with('thing') that looks for all the columns that start with 'thing'

- starts_with('some') that looks for all columns that begin with 'some'
- ends_with('thing') that looks for all the columns that start with 'thing'
- contains('else') that looks for all the columns that contain 'else' in their name

- starts_with('some') that looks for all columns that begin with 'some'
- ends_with('thing') that looks for all the columns that start with 'thing'
- contains('else') that looks for all the columns that contain 'else' in their name
- matches(...) that looks for all the columns that match with a regular expression (see ?regexp) that we indicate.

- starts_with('some') that looks for all columns that begin with 'some'
- ends_with('thing') that looks for all the columns that start with 'thing'
- contains('else') that looks for all the columns that contain 'else' in their name
- matches(...) that looks for all the columns that match with a regular expression (see ?regexp) that we indicate.
- ▶ num_range('x',1:15) that looks for all columns that are like x1, x2,...,x15

- ▶ starts with('some') that looks for all columns that begin with 'some'
- ends_with('thing') that looks for all the columns that start with 'thing'
- contains('else') that looks for all the columns that contain 'else' in their name
- matches(...) that looks for all the columns that match with a regular expression (see ?regexp) that we indicate.
- ▶ num_range('x',1:15) that looks for all columns that are like x1, x2,...,x15

There's a family of functions designed to work with select, so we can work more easily. Among them we have:

- starts_with('some') that looks for all columns that begin with 'some'
- ends_with('thing') that looks for all the columns that start with 'thing'
- contains('else') that looks for all the columns that contain 'else' in their name
- ▶ matches(...) that looks for all the columns that match with a regular expression (see ?regexp) that we indicate.
- ▶ num_range('x',1:15) that looks for all columns that are like x1, x2,...,x15

Select all the columns that end with 'ed', and the column solo. Is there something that stands out?

There's a family of functions designed to work with select, so we can work more easily. Among them we have:

- ▶ starts with('some') that looks for all columns that begin with 'some'
- ▶ ends_with('thing') that looks for all the columns that start with 'thing'
- contains('else') that looks for all the columns that contain 'else' in their name
- ▶ matches(...) that looks for all the columns that match with a regular expression (see ?regexp) that we indicate.
- ▶ num_range('x',1:15) that looks for all columns that are like x1, x2,...,x15

Select all the columns that end with 'ed', and the column solo. Is there something that stands out? Define a dataframe named members_n that only has members that died, and no columns end with 'ed' or 'id'.

Renaming with select

We can use select to rename the columns we are choosing:

```
members %>%
  select(gender = sex,used oxygen = oxygen used,role = expedition role)
## # A tibble: 76,519 \times 3
##
     gender used oxygen role
##
     <chr> <lgl>
                      <chr>>
## 1 M
           FALSE.
                      Leader
## 2 M FALSE
                      Deputy Leader
## 3 M FALSE
                      Climber
## 4 M
       FALSE
                      Exp Doctor
## 5 M
           FALSE
                      Climber
## 6 M
           FALSE
                      Climber
           FALSE
## 7 M
                      Climber
## 8 M
           FALSE
                      Climber
##
   9 M
           FALSE
                      Climber
```

rename

What did you notice from the previous slide?

rename

##

##

##

5 AMAD78301

6 AMAD78301

7 AMAD78301

What did you notice from the previous slide? We lost all of the columns we didn't mentioned. We can modify this by including everything() next to the last column. Or we can use the function rename.

```
members %>%
  rename(gender = sex,used_oxygen = oxygen_used,role = expedition_role)
```

```
## # A tibble: 76,519 x 21
##
      expedition id member id peak id
                                      peak name
                                                  year season gender
                                                                       age
      <chr>
                    <chr>
                                                 <dbl> <chr> <chr>
                                                                     <dbl>
##
                              <chr>
                                       <chr>
    1 AMAD78301
                    AMAD7830~ AMAD
                                       Ama Dabl~ 1978 Autumn M
##
                                                                        40
    2 AMAD78301
                    AMAD7830~ AMAD
                                                                        41
##
                                       Ama Dabl~ 1978 Autumn M
    3 AMAD78301
                                                                        27
##
                    AMAD7830~ AMAD
                                       Ama Dabl~ 1978 Autumn M
##
    4 AMAD78301
                    AMAD7830~ AMAD
                                       Ama Dabl~ 1978 Autumn M
                                                                        40
```

Ama Dabl~ 1978 Autumn M

Ama Dabl~ 1978 Autumn M

1978 Autumn M

Ama Dabl~

34

25

41

AMAD7830~ AMAD

AMAD7830~ AMAD

AMAD7830~ AMAD



mutate



Figure 4: Another type of mutate

mutate

Let's concatenate the role and the citizenship, using the function paste, which binds together two strings:

mutate

Let's concatenate the role and the citizenship, using the function paste, which binds together two strings:

```
members n %>%
 mutate(Role citizenship = paste(expedition role,citizenship))
## # A tibble: 1.106 x 15
                                   age citizenship expedition role
##
     peak name year season sex
##
     <chr>
               <dbl> <chr> <chr> <dbl> <chr>
                                                   <chr>>
   1 Ama Dabl~ 1979 Autumn M
                                    23 New Zealand Climber
##
   2 Ama Dabl~ 1983 Autumn M
                                    31 Switzerland Leader
##
   3 Ama Dabl~ 1983 Autumn F
##
                                    28 Switzerland Climber
##
   4 Ama Dabl~ 1985 Spring M
                                    32 Japan
                                                   Climber
##
   5 Ama Dabl~ 1988 Spring M
                                    33 Canada
                                                   Climber
##
   6 Ama Dabl~ 1992 Spring M
                                    36 Spain
                                                  Leader
##
   7 Annapurn~ 1970 Spring M
                                    32 UK
                                                   Climber
##
   8 Annapurn~
                1973 Spring M
                                    37 Japan
                                                   Climber
```

mutate

We can also do numeric operations, for example getting differences explicitly:

```
members n %>%
 mutate(Difference mts died = death height metres - highpoint metres)
## # A tibble: 1,106 x 15
                                  age citizenship expedition role
##
     peak name year season sex
     <chr>
              <dbl> <chr> <chr> <dbl> <chr>
##
                                                 <chr>
##
   1 Ama Dabl~ 1979 Autumn M
                                   23 New Zealand Climber
##
   2 Ama Dabl~ 1983 Autumn M
                                   31 Switzerland Leader
   3 Ama Dabl~ 1983 Autumn F
##
                                   28 Switzerland Climber
   4 Ama Dabl~ 1985 Spring M
##
                                   32 Japan
                                                 Climber
##
   5 Ama Dabl~ 1988 Spring M
                                   33 Canada
                                                 Climber
##
   6 Ama Dabl~ 1992 Spring M
                                   36 Spain
                                                Leader
##
   7 Annapurn~ 1970 Spring M
                                   32 UK
                                                 Climber
##
   8 Annapurn~ 1973 Spring M
                                   37 Japan
                                                 Climber
##
   9 Annapurn~
               1973 Spring M
                                   36 Japan
                                                 Climber
```

mutate to create brand new columns

<chr>

KANGO520~ KANG

##

<chr>>

2 KANG05201

3 KANG05201

4 KANG05201

5 KANG05201

6 KANG05201

7 KANG05201

1 KANG05201

We can also add and create our own columns, using our own values or from other places. But we have to be very careful they are in the appropriate order, otherwise we risk making a very dangerous mistake. How can it be dangerous?

```
members1905 %>% # this has 9 rows!
 mutate(my_row = (1:9)^2 + log(15)*(9-row_number()))
```

## #	A tibble: 9 x 22		
##	expedition id member id peak id peak name vear season sex	ge o	ci

<chr>

Kangchen~

Kangchen~

Kangchen~

Kangchen~

Kangchen~

Kangchen~

Kangchen~

<dbl> <chr> <chr> <dbl> <chr> <dbl> <c

1905 Summer M

29 UK

NA Sw

NA Ne

NA Ne

NA Ne

NA Ne

36 Sw

<chr>>

mutate + if else

6 KANG05201

R has a very neat function called if else that is just like an if, and checks wethere a condition is true, then do something, if it's not, then do something else:

```
members1905 %>% ## notice the difference between: ' and "
  mutate(Cheating = if else(expedition role == 'Leader',
                             'Cheated'.
                            "Didn't cheat"))
```

```
## # A tibble: 9 \times 22
     expedition id member id peak id peak name
##
                                                   year season sex
```

<chr> <chr> <chr>> <chr> ##

age ci <dbl> <chr> <chr> <dbl> <chr> <chr> <dbl> <c ## 1 KANG05201 KANGO520~ KANG Kangchen~ 1905 Summer M 29 UK

2 KANG05201 KANGO520~ KANG Kangchen~ 1905 Summer M NA Sw

3 KANG05201 KANGO520~ KANG

1905 Summer M Kangchen~

NA Ne ## 4 KANG05201 KANGO520~ KANG Kangchen~ 1905 Summer M NA Ne

KANGO520~ KANG

5 KANG05201 Kangchen~ 1905 Summer M NA Ne KANGO520~ KANG

Kangchen~

1905 Summer M

NA Ne

mutate + other functions

##

<lgl>

1 TRUE

2 TRUE

We can also combine mutate with other functions to obtain new columns that depend on all the values from specific columns of the dataframe:

```
members1905 %>%
  mutate(anyone died = any(died),
         max height = max(highpoint metres,na.rm = TRUE),
         last_citizenship = last(citizenship),
         youngest = min(age),
```

percent dead = mean(died),

6300 Switzerland

6300 Switzerland

<dbl> <chr>

```
## # A tibble: 9 x 6
##
     anyone died max height last citizenship youngest percent dead diff fr
```

<dbl>

NΑ

NΑ

dbl>

0.556

0.556

diff from average height death =death height metres- mean(death h select(-(1:21)) # to see the new variables



Exercises mutate

▶ What's the percentage of people that died out of those that got injured?

Exercises mutate

- ▶ What's the percentage of people that died out of those that got injured?
- Out of those that died, how many got injured before?

Exercises mutate

- ▶ What's the percentage of people that died out of those that got injured?
- Out of those that died, how many got injured before?
- Create a new column that is the concatenation of the member id and it's citizenship.

summarise

summarise

It's just like mutate, but the output has to be of only one row, and it eliminates (in the output) everything else that is not mentioned*:

summarise

A tibble: 1 x 5

It's just like mutate, but the output has to be of only one row, and it eliminates (in the output) everything else that is not mentioned*:

```
members1905 %>%
    summarise(anyone_died = any(died),
        max_height = max(highpoint_metres,na.rm = TRUE),
        last_citizenship = last(citizenship),
        youngest = min(age),
        percent_dead = mean(died))
```

We are only getting one row, and losing everything else. Using, filter, mutate and summarise, indicate how far the average of the highpoint metres for the women from Erange differs from the average from all the table (use no property).





Figure 5: Another type of group

```
members %>%
  group_by(died) %>%
  summarise(Percentage injured = mean(injured))
## 'summarise()' ungrouping output (override with '.groups' argument)
## # A tibble: 2 x 2
## died Percentage injured
    <lgl>
                    <dbl>
##
## 1 FALSE
                   0.0227
## 2 TRUE
                   0.00271
What do you think group by does?
```

On it's own, group_by doesn't do much, it really shines when we combine it with the other 5 verbs that we have studied.

On it's own, group_by doesn't do much, it really shines when we combine it with the other 5 verbs that we have studied.

All it does is group by the variables we tell it to, and the following modifications that happen on the data.frame are done on each of the groups we defined, as if each group was a dataframe.

group_by + mutate

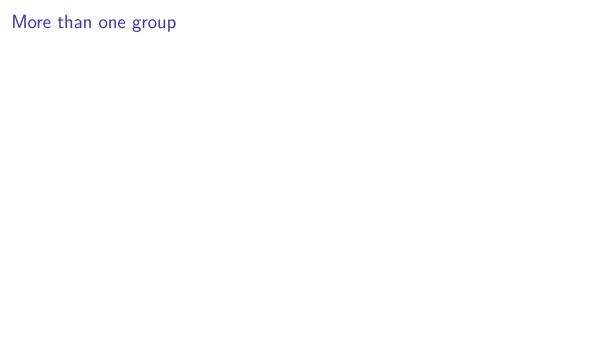
When we combine it with mutate, we get a new observation, that is only different for each of the grouped variables.

```
members_n %>% select(-(4:14)) %>%
  group_by(season) %>%
  mutate(Obs_per_season = n())
```

```
## # A tibble: 1,106 x 4
##
  # Groups: season [4]
##
     peak name vear season Obs per season
##
     <chr> <dbl> <chr>
                                      <int>
  1 Ama Dablam 1979 Autumn
##
                                       493
##
   2 Ama Dablam 1983 Autumn
                                       493
##
   3 Ama Dablam 1983 Autumn
                                       493
##
   4 Ama Dablam 1985 Spring
                                       555
##
   5 Ama Dablam
                 1988 Spring
                                       555
##
   6 Ama Dablam
                  1992 Spring
                                       555
```

group_by + mutate

Repeat what we did above, but grouping by year, instead of season.



More than one group

1 1905 Summer

2 1907 Autumn

3 1909 Autumn

4 1910 Spring

##

If you want to group with more than one variable, you can simply add the column in the argument:

```
members %>%
  group by (year, season) %>%
  summarise(Average height = mean(highpoint metres,
                                   na.rm = TRUE)
```

```
## 'summarise()' regrouping output by 'year' (override with '.groups' argu
## # A tibble: 239 \times 3
```

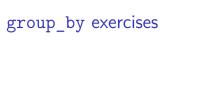
Groups: year [92] year season Average height ## <dbl> <chr> ## <dbl>

6300

7270

6965

6965



group_by exercises

▶ What year had the most dead people? (Hint: arrange)

group_by exercises

- ► What year had the most dead people? (Hint: arrange)
- Using only those with a role of climbers, compute the mean and standard deviation of the yearly number of climbers for each season. (Hint: use season as the first grouping variable)

group_by exercises

- ► What year had the most dead people? (Hint: arrange)
- Using only those with a role of climbers, compute the mean and standard deviation of the yearly number of climbers for each season. (Hint: use season as the first grouping variable)

group by + summarize Answer to the second question above:

A tibble: 5×3

season Mean Std ## <chr> <dbl> <dbl> ## 1 Autumn 321, 270, ## 2 Spring 256. 240.

##

```
members %>%
  filter(expedition_role =='Climber') %>%
  group by(season, year) %>%
  summarize(Total = n()) %>%
  summarise(Mean = mean(Total),
            Std = sd(Total)
## 'summarise()' regrouping output by 'season' (override with '.groups' ar
## 'summarise()' ungrouping output (override with '.groups' argument)
```

group_by + filter

##

##

5 AMAD79101

6 AMAD79101

If you want to eliminated the groups that don't have enough, or that have too many observations, you can do it by combining filter and group_by directly:

```
(surviving_members <- members %>%
  group_by(expedition_id) %>%
  filter(n() > 10))
```

```
## # A tibble: 38.762 x 21
## # Groups: expedition id [2,264]
##
     expedition id member id peak id peak name vear season sex
                                                                     age
##
      <chr>
                    <chr>
                              <chr>
                                      <chr>
                                                <dbl> <chr> <chr> <dbl>
    1 AMAD79101
##
                    AMAD7910~ AMAD
                                      Ama Dabl~ 1979 Spring M
                                                                      35
##
   2 AMAD79101
                    AMAD7910~ AMAD
                                      Ama Dabl~ 1979 Spring M
                                                                      37
                                      Ama Dabl~ 1979 Spring M
   3 AMAD79101
                    AMAD7910~ AMAD
                                                                      23
##
                                                                      44
##
   4 AMAD79101
                    AMAD7910~ AMAD
                                      Ama Dabl~ 1979 Spring M
```

Ama Dabl~ 1979 Spring M

Ama Dabl~

1979 Spring M

AMAD7910~ AMAD

AMAD7910~ AMAD

25

28

rowwise

This function, as it name tells us, is like doing a group_by, but operates on each row. This one is particular useful when you are creating your own functions and they have weird interactions with vectors, like using sum, mean, and such... But we would probably get different results on each row.

Thank you :)

Any questions?

Further references

- Check the help page and the vignettes of dplyr! (type ?dplyr, or: vignette('dplyr') on the console)
- R for Data Science, by Hadley Wickham
- Advanced R, by Hadley Wickham
- ► The R Inferno, by Patrick Burns