

INTRODUCTION TO DATA SCIENCE

Lecture 6

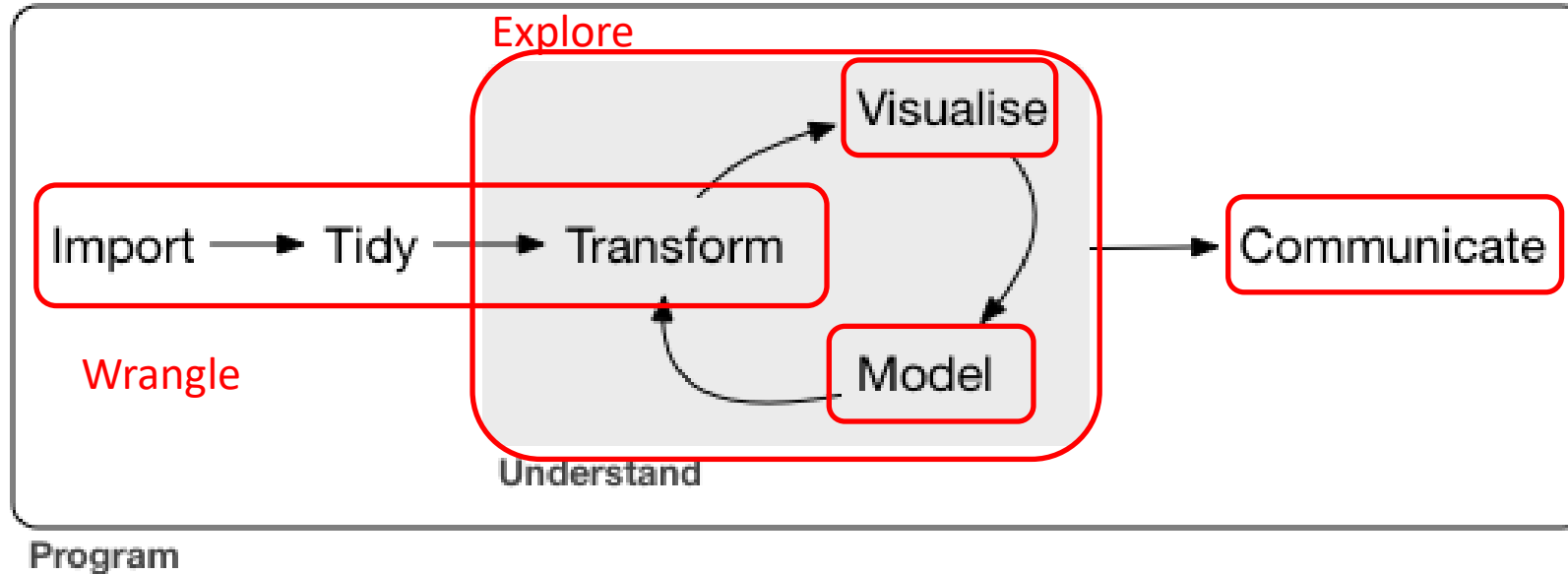
Dr. Ibrahim Radwan

OUTLINE

- Data Science, a Practical View
- Data Wrangling
- Data Manipulation

DATA SCIENCE; A PRACTICAL VIEW

Program Steps



- 1- Reading Data
- 2- Data Wrangling
- 3- Data Exploratory
- 4- Modelling
- 5- Result Communication

R for Data Science, by Garrett Golemund and Hadley Wickham

- Data wrangling is the process of transforming the data into a suitable format to conduct the modelling or analysis processes.
- Data wrangling is important: without it, the data may not be ready to be analysed.
- Data wrangling involves the following five steps:

import

1. Gather your data from different sources

Transform

2. Clean up the duplicates, blanks, and other logical errors

3. Join all your data into a single table (e.g. data frame)

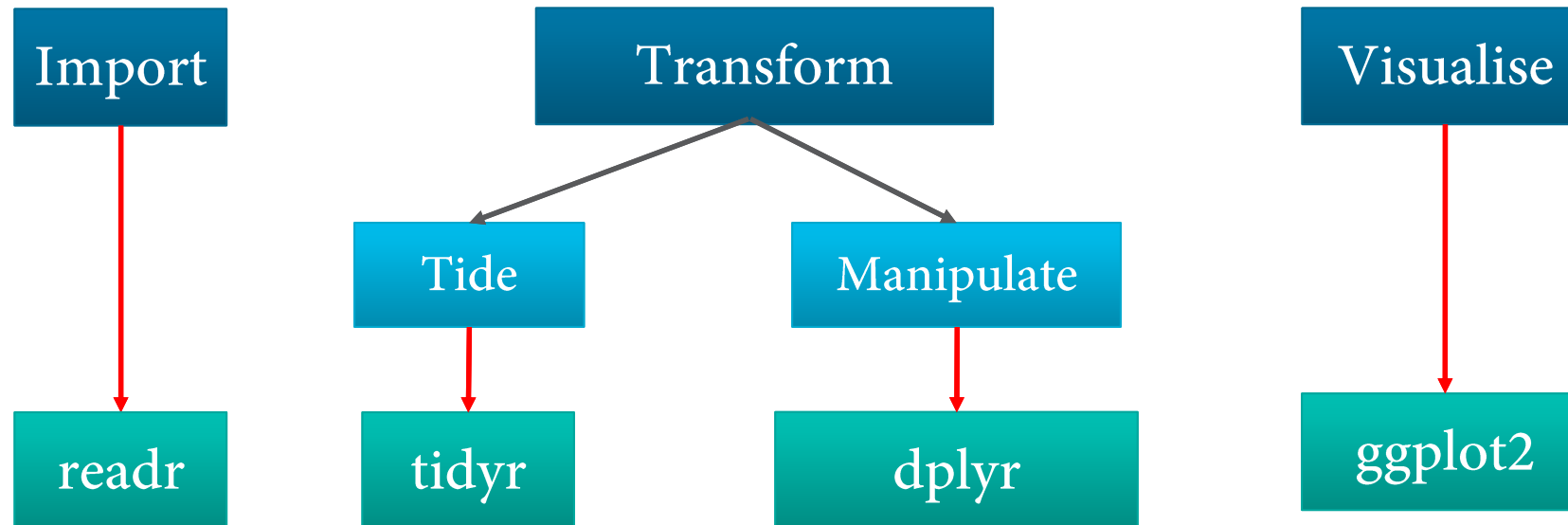
4. Adding new variables/records by calculating new fields and re-ordering

Visualise

5. View or visualize the data to remove outliers and illogical results

DATA WRANGLING (2)

- Practically, we have three main processes to wrangle the data



- These four packages provide the *grammar* of the data importing, tidying, manipulation and data visualisation.

DATA IMPORT

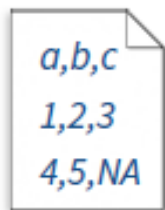
- ▶ Reading flat files such as CSV or text files
- ▶ Suppose we have the following CSV file:
 - ▶ We will use the `'readr'` library to read the csv files
 - ▶ Then, we can use the `read_csv("filename.csv")` function to read the contents.
 - ▶ The returned object is a tibble.
- ▶ It reads the data with 10x faster than the base R functions.

ID	Name	Age
23424	Ana	45
11234	Charles	23
77654	Susanne	76

```
> read_csv("data_sample.csv")
Parsed with column specification:
cols(
  ID = col_double(),
  Name = col_character(),
  Age = col_double()
)
# A tibble: 3 x 3
  ID Name      Age
  <dbl> <chr>    <dbl>
1 23424 Ana      45
2 11234 Charles  23
3 77654 Susanne  76
> |
```

DATA IMPORT (2)

```
read_*(file, col_names = TRUE, col_types = NULL, locale = default_locale(), na = c("", "NA"),  
quoted_na = TRUE, comment = "", trim_ws = TRUE, skip = 0, n_max = Inf, guess_max = min(1000,  
n_max), progress = interactive())
```

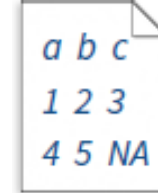


a,b,c
1,2,3
4,5,NA



A	B	C
1	2	3
4	5	NA

Comma Delimited Files
`read_csv("file.csv")`

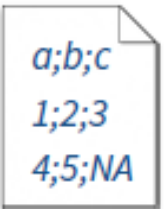


a b c
1 2 3
4 5 NA



A	B	C
1	2	3
4	5	NA

Tab Delimited Files
`read_tsv("file.tsv")`
Also `read_table()`

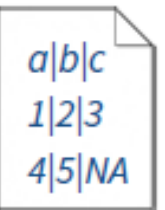


a;b;c
1;2;3
4;5;NA



A	B	C
1	2	3
4	5	NA

Semi-colon Delimited Files
`read_csv2("file2.csv")`



a|b|c
1|2|3
4|5|NA



A	B	C
1	2	3
4	5	NA

Files with Any Delimiter
`read_delim("file.txt", delim = "|")`

To save data into csv or txt file

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)`

DATA IMPORT (3)

```
1 # Import data into R
2
3 # It is part of the `tidyverse` package
4 if (!("tidyverse" %in% rownames(installed.packages()))){
5   install.packages("tidyverse")
6 }
7
8 library("tidyverse")
9
10 heights <- read_csv("data/heights.csv")
11
12 # supply inline csv file
13 read_csv("a,b,c
14 1,2,3
15 4,5,6")
16
17 # skip first two lines, i.e. meta data the beginning of a file
18 read_csv("The first line of metadata
19 The second line of metadata
20 x,y,z
21 1,2,3", skip = 2)
22 # or skip when comment
23 read_csv("# A comment I want to skip
24 x,y,z
25 1,2,3", comment = "#")
26
27 # by default the first row is for column, we can skip this rule
28 read_csv("1,2,3\n4,5,6", col_names = FALSE) # \n is convenient shortcut for new line
29 # pass column names
30 read_csv("1,2,3\n4,5,6", col_names = c("x", "y", "z"))
31
32 # deal with na values
33 read_csv("a,b,c\n1,2,.", na = ".")
34
```


DATA IMPORT – COL SPECIFICATIONS

- The functions of `readr` package guess the types of each column and convert them as appropriate (but will NOT convert strings to factors automatically). A message shows the type of each column will appear in the result.

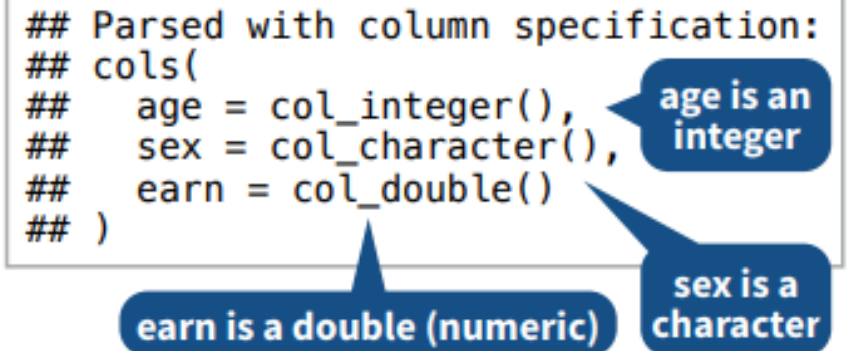
For each `col_*` function, there is a `parse_*` function

`col_guess` function uses the first 1000 rows from each column to guess the column type.

What about if:

- *The first 1000 rows are just special case and the rest are different?*
- *The first 1000 rows are `na` ?*

```
## Parsed with column specification:
## cols(
##   age = col_integer(),
##   sex = col_character(),
##   earn = col_double()
## )
```



The solution is to use
`problems()` function

DATA IMPORT – COL SPECIFICATIONS (2)

```
1 # col specifications
2 # logicals
3 lgl <- parse_logical(c("TRUE", "FALSE", "NA"))
4 str(lgl)
5 # integers
6 intgr <- parse_integer(c("1", "2", "3"))
7 str(intgr)
8 # date
9 dt <- parse_date(c("2010-01-01", "1979-10-14"))
10 str(dt)
11 # which strings should be specified as missing
12 intgr_missing <- parse_integer(c("1", "231", ".", "456"), na = ".")
13 str(intgr_missing)
14 # sometimes the parsing fails, NAs will be used when failure is occurring
15 x <- parse_integer(c("123", "345", "abc", "123.45", "221"))
16 x
17 # here you can use problems
18 problems(x)
19
20 # parse numbers with . or , or $ or %
21 dbl <- parse_double("1.24")
22 str(dbl)
23 num <- parse_number("$123,456,789")
24 num
25 num <- parse_number("$123,456,789")
26 num
27 num <- parse_number("$123.456.789", locale=locale(grouping_mark = '.'))
28 num
29 num <- parse_number("123'456'789", locale = locale(grouping_mark = ''))
30 num
```

DATA IMPORT – COL SPECIFICATIONS (3)

- Using the *problems()* function with the `read_*` functions from the `'readr'` package to check the problems in parsing and guessing the column types.

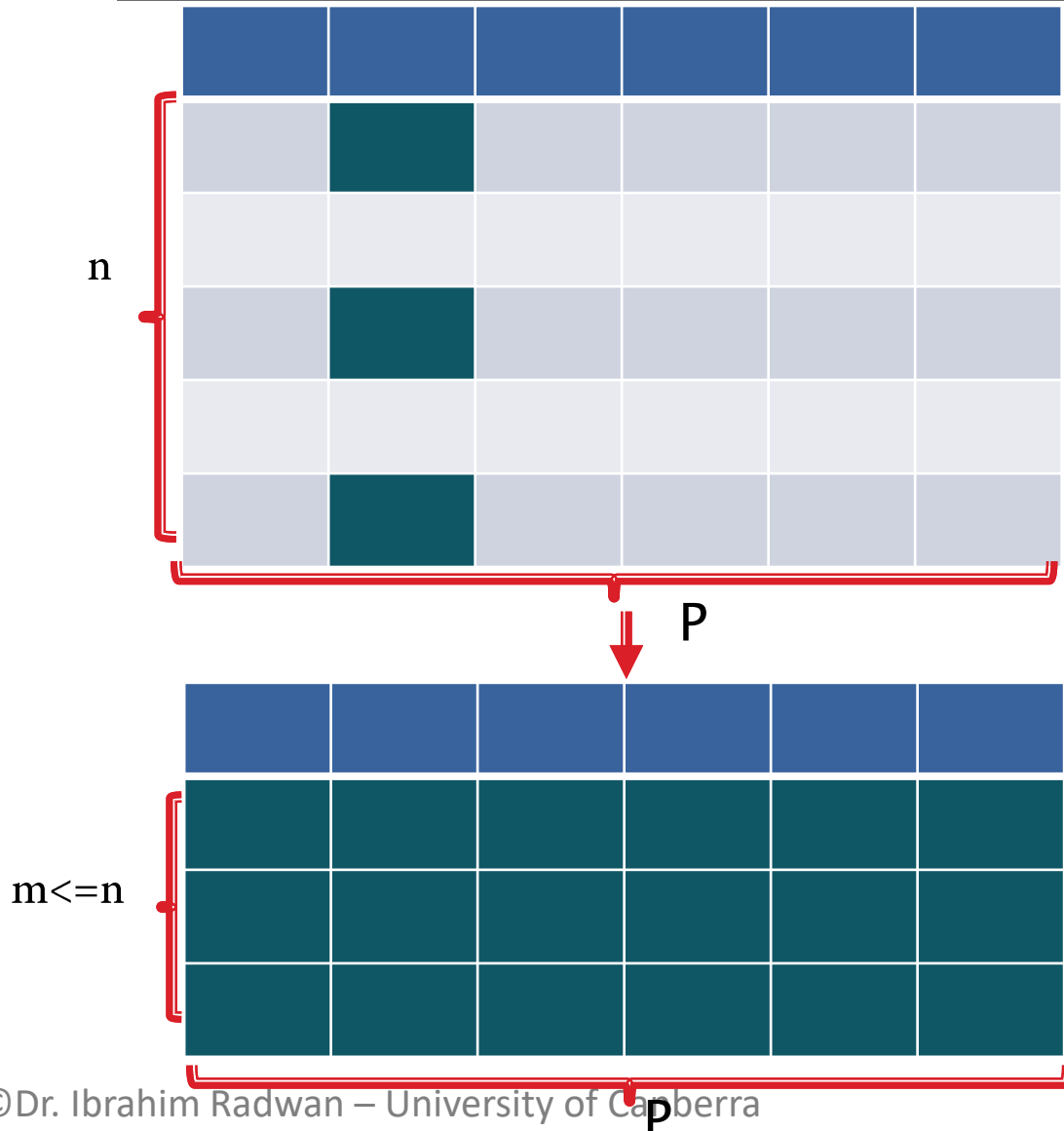
```
problems(tibble_object)
```

DATA IMPORT – COL SPECIFICATIONS (4)

```
1 # base function to read data.frame
2 df <- read.csv(readr_example("challenge.csv"))
3 df
4 view(df)
5 df[1001,]
6 str(df)
7 problems(df) # no problems as everything are chars
8 #####
9 # readr functions to read csv
10 tbl <- read_csv(readr_example("challenge.csv"))
11 tbl
12 view(tbl)
13 tbl[1001,]
14 str(tbl)
15 #####
16 # use problems to check the problems in the tbl parsing
17 problems(tbl)
18
19 # how to fix the problems
20 # first let us fix the problem by hard-specifying the column types
21 tbl2 <- read_csv(
22   readr_example("challenge.csv"),
23   col_types = cols(
24     x = col_double(),
25     y = col_character()
26   )
27 )
28 tbl2
29 tail(tbl2) # there are dates stored with this column
30 # change the character to date
31 tbl3 <- read_csv(
32   readr_example("challenge.csv"),
33   col_types = cols(
34     x = col_double(),
35     y = col_date()
36   )
37 )
```

- The `dplyr` package in `tidyverse` library presents five *verbs* for manipulating the data in data frames:
 1. `filter()` extracts a subset of the rows (i.e., observations) based on some criteria
 2. `select()` extracts a subset of the columns (i.e., features, variables) based on some criteria
 3. `mutate()` adds or modifies existing columns
 4. `arrange()` sorts the rows
 5. `summarise()` aggregates the data across rows (e.g., group them according to some criteria)
- Each of these functions takes a data frame as its first argument and returns a data frame.

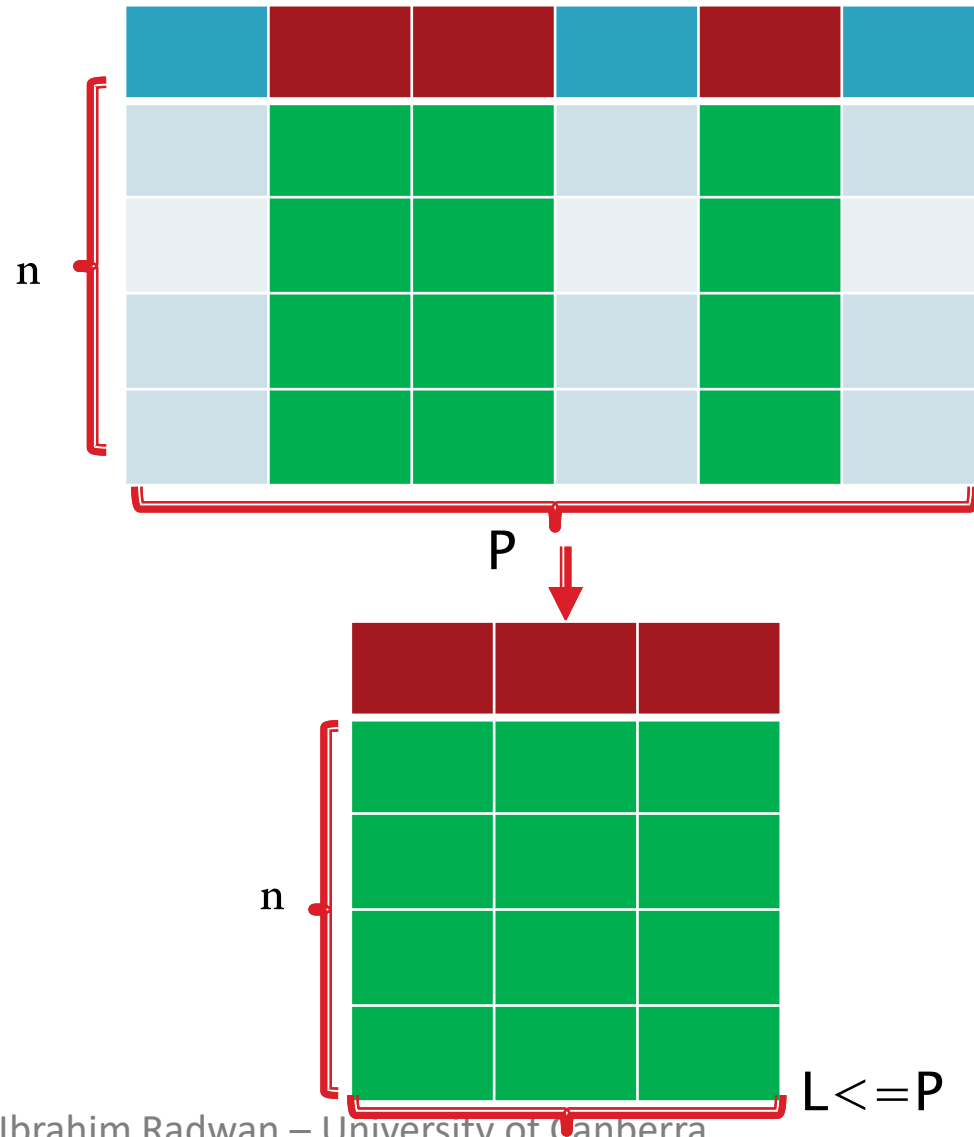
SUBSET BY ROW (FILTER)



```
filter(dataframe, criteria)
```

```
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# filter based on conditions
filter(df, month == 1, day == 1)
filter(df, month == 12, day == 25)
# you may combine conditions using logical operators
filter(df, month == 1 | month == 12)
# or by combining variables in vector
filter(df, month %in% c(11, 12))
# comma means and (&)
# for example, extract all records for flights,
# that were not delayed (arr and dep) more than 2 hrs
filter(flights, arr_delay <= 120, dep_delay <= 120)
```

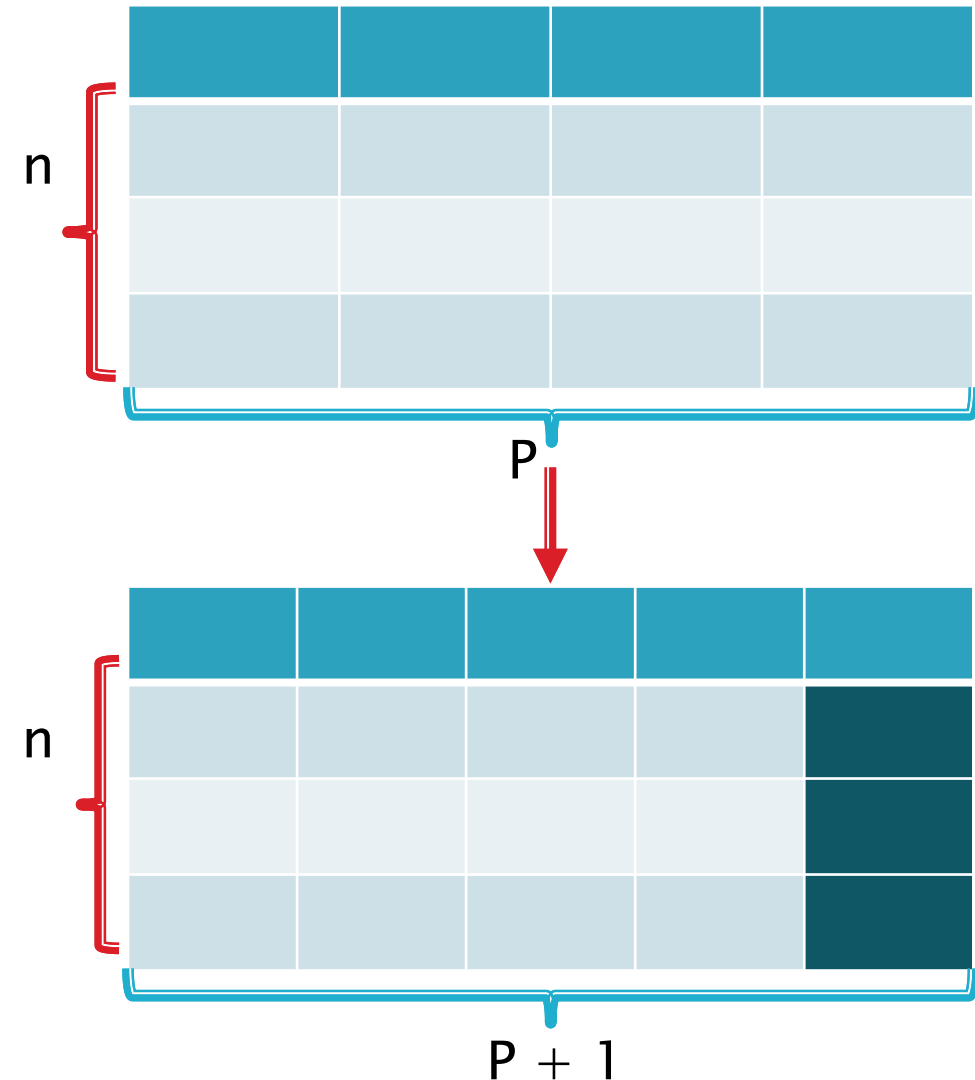
SUBSET BY COLUMN (SELECT)



```
select(dataframe, col_names or criteria)
```

```
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# Select columns by name
select(flights, year, month, day)
# Select all columns between year and day
(inclusive)
select(flights, year:day)
# Select all columns except those from year to day
(inclusive)
select(flights, -(year:day))
# rename() is a variant of select() that keeps all
the variables that aren't explicitly mentioned:
rename(flights, tail_num = tailnum)
# Move a variable to the start of the data frame.
select(flights, time_hour, air_time, everything())
```

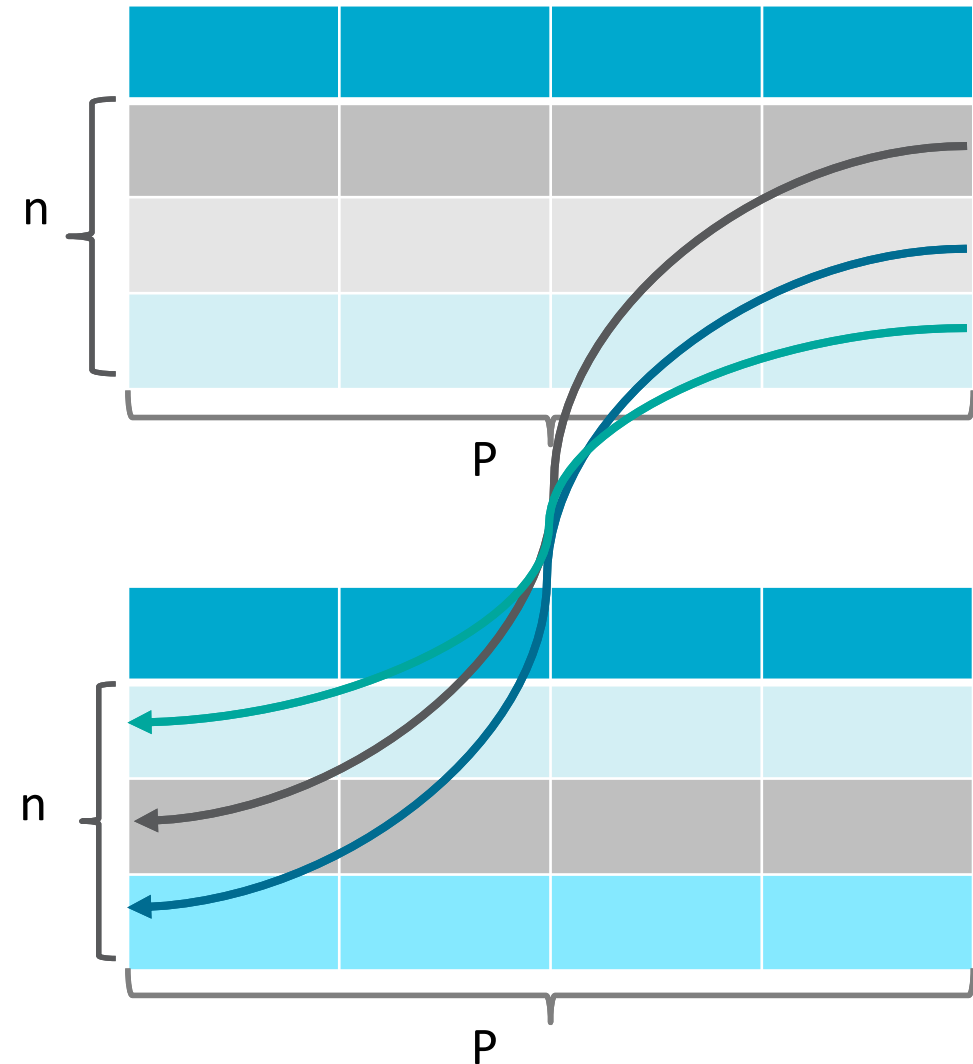
ADD - MODIFY COLUMN (MUTATE)



```
mutate(dataframe, new_cols or modified_cols)
```

```
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# Create a new dataset
flights_sml <- select(df, year:day, ends_with("delay"),
distance, air_time )
# add new columns to the data frame
mutate(flights_sml, gain = dep_delay - arr_delay, speed =
distance / air_time * 60 )
# Note that you can refer to columns that you've just
created:
mutate(flights_sml, gain = dep_delay - arr_delay, hours =
air_time / 60, gain_per_hour = gain / hours )
# If you only want to keep the new variables, use
transmute():
transmute(flights, gain = dep_delay - arr_delay, hours =
air_time / 60, gain_per_hour = gain / hours )
```

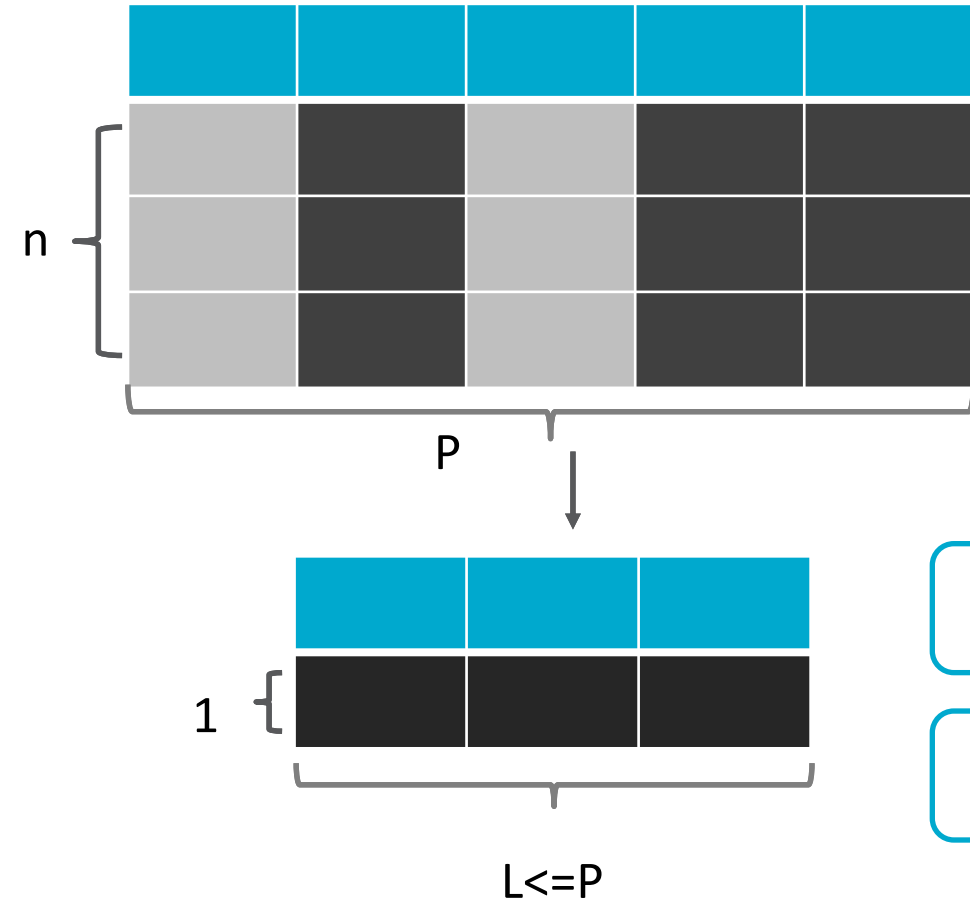

RE-ORDER ROWS (ARRANGE)



```
arrange(dataframe, col_name)
```

```
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# sort data by distance
arrange(df, distance)
# sort data by distance descendingly
arrange(df, desc(distance))
# Sort Data by Multiple Variables
arrange(df, dep_time, arr_time)
```

AGGREGATE (SUMMARISE)



```
summarise(dataframe, agg_func(col_name))
```

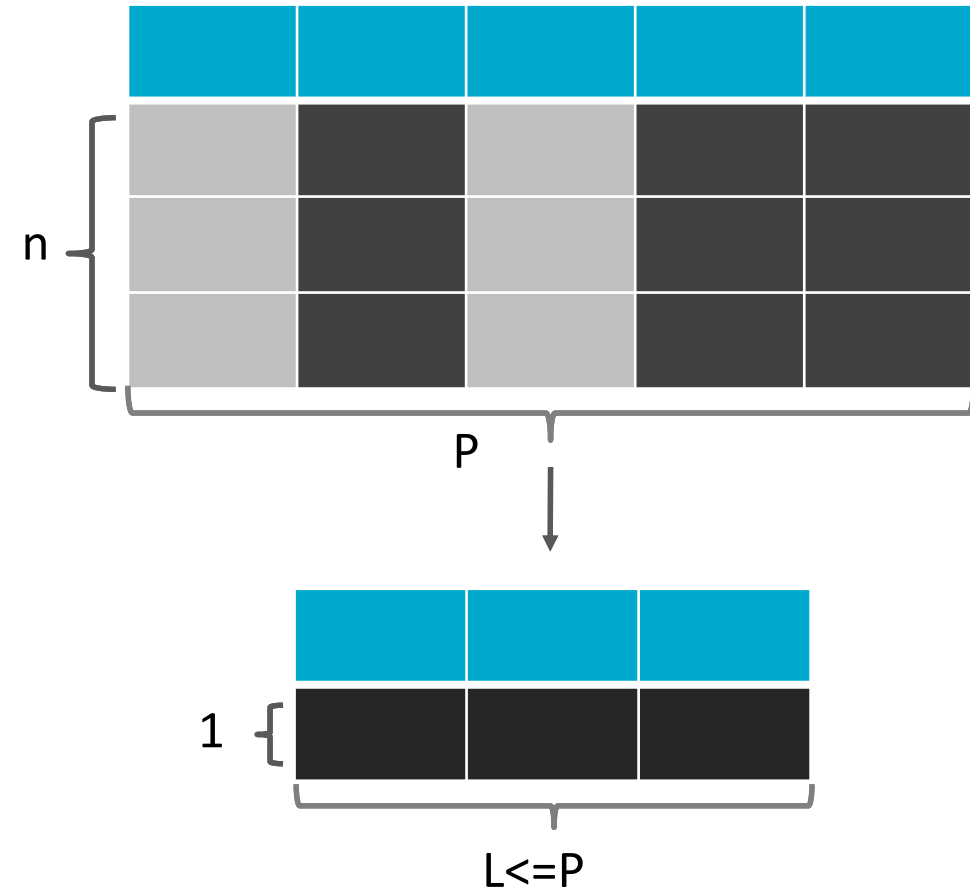
```
# call the required libraries
library(tidyverse)
library(nycflights13)
df <- flights
# extract a statistical metric from variable /
variables of the data
summarise(df, delay = mean(dep_delay, na.rm = TRUE))
```

Aggregation functions such as mean, sd, var, median, min and max

There is also summarise_each(dataframe, funs(aggregation_func))

summarise() is not terribly useful unless we pair it with
group_by()

AGGREGATE WITH GROUPING



```
group_by(dataframe, col_name)
```

```
# group the data of the flights by the date
by_day <- group_by(flights, year, month, day)
# get the average delay per date/day
summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))
# Imagine that we want to explore the relationship
between the distance and average delay for each
location.
```

```
by_dest <- group_by(flights, dest)
```

```
# extract the number of flights, average distance and
average delay for each destination
```

```
delay <- summarise(by_dest, count= n(), dist=
mean(distance, na.rm = TRUE), delay= mean(arr_delay,
na.rm = TRUE) )
```

```
# visualise to understand the relationship
```

```
ggplot(data= delay, mapping=aes(x= dist, y= delay)) +
geom_point(aes(size= count), alpha= 1/ 3) +
geom_smooth()
```

THE PIPE OPERATOR %>%

- In data wrangling, most likely, you need to perform series of operations (i.e. *verbs*) on the same data.
- This will need you to create intermediate tables temporarily to save the results to be processed with the next operations.
- R provides an elegant way to perform series of operations on the same data in one go via using the *pipe operator* %>%

original data → select → filter

```
F(x) is the same as  
x %>% F
```

```
16 %>% sqrt() %>% log2()  
[1] 2
```

THE PIPE OPERATOR %>% (2)

- In the previous ``nycflights`` example, there were three steps to extract the relationship between the distance and the delay of the flights per destination.
 1. Group flights by destination.
 2. Summarise to compute distance, average delay, and number of flights.
 3. Filter to remove noisy points

This can be achieved by using the pipe operator:

```
delay <- df %>%  
  group_by(dest) %>%  
  summarise(count= n(), dist= mean(distance, na.rm = TRUE),  
            delay= mean(arr_delay, na.rm = TRUE)) %>%  
  filter(count > 20, dest != 'HNL')
```

RECOMMENDED READING

- You are recommended to read chapters 5 & 12 from the “*R for Data Science*” book:
 - <https://r4ds.had.co.nz/transform.html>