

Super Study Guide: Data Science Tools

Afshine AMIDI and Shervine AMIDI

August 21, 2020

Contents

1	Data retrieval with SQL	2
1.1	General concepts	2
1.2	Aggregations	2
1.3	Window functions	3
1.4	Advanced functions	4
1.5	Table manipulation	5
2	Working with data with R	6
2.1	Data manipulation	6
2.1.1	Main concepts	6
2.1.2	Data preprocessing	6
2.1.3	Data frame transformation	7
2.1.4	Aggregations	8
2.1.5	Window functions	9
2.2	Data visualization	9
2.2.1	General structure	9
2.2.2	Advanced features	10
2.2.3	Last touch	11
3	Working with data with Python	13
3.1	Data manipulation	13
3.1.1	Main concepts	13
3.1.2	Data preprocessing	13
3.1.3	Data frame transformation	14
3.1.4	Aggregations	15
3.1.5	Window functions	16
3.2	Data visualization	16
3.2.1	General structure	16
3.2.2	Advanced features	17
3.2.3	Last touch	17

4	Engineering productivity tips with Git, Bash and Vim	18
4.1	Working in groups with Git	18
4.1.1	Overview	18
4.1.2	Main commands	18
4.1.3	Project structure	19
4.2	Working with Bash	20
4.3	Automating tasks	21
4.4	Mastering editors	21

Appendix A	Conversion between R and Python: data manipulation	22
A.1	Main concepts	22
A.2	Data preprocessing	22
A.3	Data frame transformation	22

Appendix B	Conversion between R and Python: data visualization	23
B.1	General structure	23
B.2	Advanced features	23

SECTION 1

Data retrieval with SQL

1.1 General concepts

Structured Query Language – Structured Query Language, abbreviated as SQL, is a language that is largely used in the industry to query data from databases.

Query structure – Queries are usually structured as follows:

SQL

```
-- Select fields                                mandatory
SELECT
  col_1,
  col_2,
  ... ,
  col_n

-- Source of data                                mandatory
FROM table t

-- Gather info from other sources                optional
JOIN other_table ot
  ON (t.key = ot.key)

-- Conditions                                    optional
WHERE some_condition(s)

-- Aggregating                                  optional
GROUP BY column_group_list

-- Sorting values                              optional
ORDER BY column_order_list

-- Restricting aggregated values                optional
HAVING some_condition(s)

-- Limiting number of rows                      optional
LIMIT some_value
```

Remark: the `SELECT DISTINCT` command can be used to ensure not having duplicate rows.

Condition – A condition is of the following format:

SQL

```
some_col some_operator some_col_or_value
```

where `some_operator` can be among the following common operations:

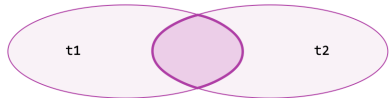

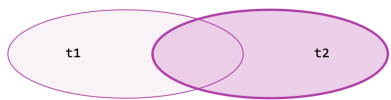
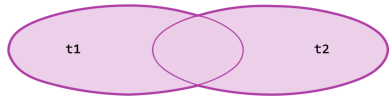
Category	Operator	Command
General	Equality / non-equality	<code>= / !=, <></code>
	Inequalities	<code>>=, >, <, <=</code>
	Belonging	<code>IN (val_1, ..., val_n)</code>
	And / or	<code>AND / OR</code>
	Check for missing value	<code>IS NULL</code>
	Between bounds	<code>BETWEEN val_1 AND val_2</code>
Strings	Pattern matching	<code>LIKE '%val%'</code>

Joins – Two tables `table_1` and `table_2` can be joined in the following way:

SQL

```
...
FROM table_1 t1
type_of_join table_2 t2
  ON (t2.key = t1.key)
...
```

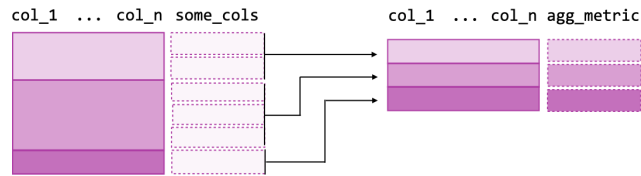
where the different `type_of_join` commands are summarized in the table below:

Type of join	Illustration
INNER JOIN	
LEFT JOIN	
RIGHT JOIN	
FULL JOIN	

Remark: joining every row of table 1 with every row of table 2 can be done with the `CROSS JOIN` command, and is commonly known as the cartesian product.

1.2 Aggregations

Grouping data – Aggregate metrics are computed on grouped data in the following way:



The SQL command is as follows:

SQL

```
SELECT
  col_1,
  agg_function(col_2)
FROM table
GROUP BY col_1
```

□ **Grouping sets** – The `GROUPING SETS` command is useful when there is a need to compute aggregations across different dimensions at a time. Below is an example of how all aggregations across two dimensions are computed:

SQL

```
SELECT
  col_1,
  col_2,
  agg_function(col_3)
FROM table
GROUP BY (
  GROUPING SETS
    (col_1),
    (col_2),
    (col_1, col_2)
)
```

□ **Aggregation functions** – The table below summarizes the main aggregate functions that can be used in an aggregation query:

Category	Operation	Command
Values	Mean	<code>AVG(col)</code>
	Percentile	<code>PERCENTILE_APPROX(col, p)</code>
	Sum / # of instances	<code>SUM(col) / COUNT(col)</code>
	Max / min	<code>MAX(col) / MIN(col)</code>
	Variance / standard deviation	<code>VAR(col) / STDEV(col)</code>
Arrays	Concatenate into array	<code>collect_list(col)</code>

Remark: the median can be computed using the `PERCENTILE_APPROX` function with `p` equal to 0.5.

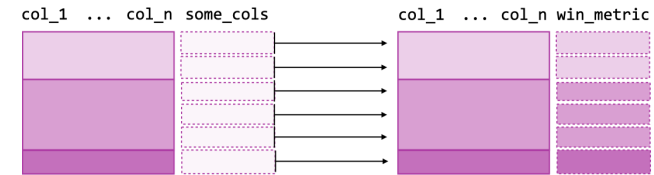
□ **Filtering** – The table below highlights the differences between the `WHERE` and `HAVING` commands:

WHERE	HAVING
- Filter condition applies to individual rows	- Filter condition applies to aggregates
- Statement placed right after <code>FROM</code>	- Statement placed right after <code>GROUP BY</code>

Remark: if `WHERE` and `HAVING` are both in the same query, `WHERE` will be executed first.

1.3 Window functions

□ **Definition** – A window function computes a metric over groups and has the following structure:



The SQL command is as follows:

SQL

```
some_window_function() OVER(PARTITION BY some_col ORDER BY another_col)
```

Remark: window functions are only allowed in the `SELECT` clause.

□ **Row numbering** – The table below summarizes the main commands that rank each row across specified groups, ordered by a specific column:

Command	Description	Example
<code>ROW_NUMBER()</code>	Ties are given different ranks	1, 2, 3, 4
<code>RANK()</code>	Ties are given same rank and skip numbers	1, 2, 2, 4
<code>DENSE_RANK()</code>	Ties are given same rank and don't skip numbers	1, 2, 2, 3

□ **Values** – The following window functions allow to keep track of specific types of values with respect to the partition:

Command	Description
<code>FIRST_VALUE(col)</code>	Takes the first value of the column
<code>LAST_VALUE(col)</code>	Takes the last value of the column
<code>LAG(col, n)</code>	Takes the n^{th} previous value of the column
<code>LEAD(col, n)</code>	Takes the n^{th} following value of the column
<code>NTH_VALUE(col, n)</code>	Takes the n^{th} value of the column

1.4 Advanced functions

□ **SQL tips** – In order to keep the query in a clear and concise format, the following tricks are often done:

Operation	Command	Description
Renaming columns	<code>SELECT operation_on_column AS col_name</code>	New column names shown in query results
Abbreviating tables	<code>FROM table_1 t1</code>	Abbreviation used within query for simplicity in notations
Simplifying group by	<code>GROUP BY col_number_list</code>	Specify column position in <code>SELECT</code> clause instead of whole column names
Limiting results	<code>LIMIT n</code>	Display only n rows

□ **Sorting values** – The query results can be sorted along a given set of columns using the following command:

```
SQL
... [query] ...
ORDER BY col_list
```

Remark: by default, the command sorts in ascending order. If we want to sort it in descending order, the `DESC` command needs to be used after the column.

□ **Column types** – In order to ensure that a column or value is of one specific data type, the following command is used:

```
SQL
CAST(some_col_or_value AS data_type)
```

where `data_type` is one of the following:

Data type	Description	Example
INT	Integer	2
DOUBLE	Numerical value	2.0
STRING	String	'teddy bear'
VARCHAR		
DATE	Date	'2020-01-01'
TIMESTAMP	Timestamp	'2020-01-01 00:00:00.000'

Remark: if the column contains data of different types, the `TRY_CAST()` command will convert unknown types to `NULL` instead of throwing an error.

□ **Column manipulation** – The main functions used to manipulate columns are described in the table below:

Category	Operation	Command
General	Take first non-NULL value	<code>COALESCE(col_1, col_2, ..., col_n)</code>
	Create a new column combining existing ones	<code>CONCAT(col_1, ..., col_n)</code>
Value	Round value to n decimals	<code>ROUND(col, n)</code>
String	Converts string column to lower / upper case	<code>LOWER(col) / UPPER(col)</code>
	Replace occurrences of old in col to new	<code>REPLACE(col, old, new)</code>
	Take the substring of col, with a given start and length	<code>SUBSTR(col, start, length)</code>
	Remove spaces from the left / right / both sides	<code>LTRIM(col) / RTRIM(col) / TRIM(col)</code>
	Length of the string	<code>LENGTH(col)</code>
Date	Truncate at a given granularity (year, month, week)	<code>DATE_TRUNC(time_dimension, col_date)</code>
	Transform date	<code>DATE_ADD(col_date, number_of_days)</code>

□ **Conditional column** – A column can take different values with respect to a particular set of conditions with the `CASE WHEN` command as follows:

```
SQL
CASE WHEN some_condition THEN some_value
      ...
      WHEN some_other_condition THEN some_other_value
      ELSE some_other_value_n END
```

□ **Combining results** – The table below summarizes the main ways to combine results in queries:

Category	Command	Remarks
Union	<code>UNION</code>	Guarantees distinct rows
	<code>UNION ALL</code>	Potential newly-formed duplicates are kept
Intersection	<code>INTERSECT</code>	Keeps observations that are in all selected queries

□ **Common table expression** – A common way of handling complex queries is to have temporary result sets coming from intermediary queries, which are called common table expressions (abbreviated CTE), that increase the readability of the overall query. It is done thanks to the `WITH ... AS ...` command as follows:

```
SQL
WITH cte_1 AS (
  SELECT ...
),
```

```
...
cte_n AS (
SELECT ...
)

SELECT ...
FROM ...
```

1.5 Table manipulation

❑ **Table creation** – The creation of a table is done as follows:

SQL

```
CREATE [table_type] TABLE [creation_type] table_name(
  col_1 data_type_1,
  ...
  col_n data_type_n
)
[options];
```

where [table_type], [creation_type] and [options] are one of the following:

Category	Command	Description
Table type	Blank	Default table
	EXTERNAL TABLE	External table
Creation type	Blank	Creates table and overwrites current one if it exists
	IF NOT EXISTS	Only creates table if it does not exist
Options	location 'path_to_hdfs_folder'	Populate table with data from hdfs folder
	stored as data_format	Stores the table in a specific data format, e.g. parquet, orc or avro

❑ **Data insertion** – New data can either append or overwrite already existing data in a given table as follows:

SQL

```
WITH ... -- optional
INSERT [insert_type] table_name -- mandatory
SELECT ...; -- mandatory
```

where [insert_type] is among the following:

Command	Description
OVERWRITE	Overwrites existing data
INTO	Appends to existing data

❑ **Dropping table** – Tables are dropped in the following way:

SQL

```
DROP TABLE table_name;
```

❑ **View** – Instead of using a complicated query, the latter can be saved as a view which can then be used to get the data. A view is created with the following command:

SQL

```
CREATE VIEW view_name AS complicated_query;
```

Remark: a view does not create any physical table and is instead seen as a shortcut.

SECTION 2

Working with data with R

2.1 Data manipulation

2.1.1 Main concepts

□ **File management** – The table below summarizes the useful commands to make sure the working directory is correctly set:

Category	Action	Command
Paths	Change directory to another path	<code>setwd(path)</code>
	Get current working directory	<code>getwd()</code>
	Join paths	<code>file.path(path_1, ..., path_n)</code>
Files	List files and folders in a given directory	<code>list.files(path, include.dirs = TRUE)</code>
	Check if path is a file / folder	<code>file.test('-f', path)</code>
		<code>file.test('-d', path)</code>
	Read / write csv file	<code>read.csv(path_to_csv_file)</code>
		<code>write.csv(df, path_to_csv_file)</code>

□ **Chaining** – The symbol `%>`, also called "pipe", enables to have chained operations and provides better legibility. Here are its different interpretations:

- `f(arg_1, arg_2, ..., arg_n)` is equivalent to `arg_1 %> f(arg_2, arg_3, ..., arg_n)`, and also to:
 - `arg_1 %> f(., arg_2, ..., arg_n)`
 - `arg_2 %> f(arg_1, ., arg_3, ..., arg_n)`
 - `arg_n %> f(arg_1, ..., arg_n-1, .)`
- A common use of pipe is when a dataframe `df` gets first modified by `some_operation_1`, then `some_operation_2`, until `some_operation_n` in a sequential way. It is done as follows:

R

```
# df gets some_operation_1, then some_operation_2, ...,
# then some_operation_n
df %>%
  some_operation_1 %>%
  some_operation_2 %>%
  ... %>%
  some_operation_n
```

□ **Exploring the data** – The table below summarizes the main functions used to get a complete overview of the data:

Category	Action	Command
Look at data	Select columns of interest	<code>df %>% select(col_list)</code>
	Remove unwanted columns	<code>df %>% select(-col_list)</code>
	Look at n first rows / last rows	<code>df %>% head(n) / df %>% tail(n)</code>
	Summary statistics of columns	<code>df %>% summary()</code>
Data types	Data types of columns	<code>df %>% str()</code>
	Number of rows / columns	<code>df %>% NROW() / df %>% NCOL()</code>

□ **Data types** – The table below sums up the main data types that can be contained in columns:

Data type	Description	Example
character	String-related data	<code>'teddy bear'</code>
factor	String-related data that can be put in bucket, or ordered	<code>'high'</code>
numeric	Numerical data	<code>24.0</code>
int	Numeric data that are integer	<code>24</code>
Date	Dates	<code>'2020-01-01'</code>
POSIXct	Timestamps	<code>'2020-01-01 00:01:00'</code>

2.1.2 Data preprocessing

□ **Filtering** – We can filter rows according to some conditions as follows:

R

```
df %>%
  filter(some_col some_operation some_value_or_list_or_col)
```

where `some_operation` is one of the following:

Category	Operation	Command
Basic	Equality / non-equality	<code>== / !=</code>
	Inequalities	<code><, <=, >=, ></code>
	And / or	<code>& / </code>
Advanced	Check for missing value	<code>is.na()</code>
	Belonging	<code>%in% (val_1, ..., val_n)</code>
	Pattern matching	<code>%like% 'val'</code>

Remark: we can filter columns with the `select_if` command.

□ **Changing columns** – The table below summarizes the main column operations:

Action	Command
Add new columns on top of old ones	<code>df %>% mutate(new_col = operation(other_cols))</code>
Add new columns and discard old ones	<code>df %>% transmute(new_col = operation(other_cols))</code>
Modify several columns in-place	<code>df %>% mutate_at(vars, funs)</code>
Modify all columns in-place	<code>df %>% mutate_all(funs)</code>
Modify columns fitting a specific condition	<code>df %>% mutate_if(condition, funs)</code>
Unite columns	<code>df %>% unite(new_merged_col, old_cols_list)</code>
Separate columns	<code>df %>% separate(col_to_separate, new_cols_list)</code>

□ **Conditional column** – A column can take different values with respect to a particular set of conditions with the `case_when()` command as follows:

R

```
case_when(condition_1 ~ value_1, # If condition_1 then value_1
          condition_2 ~ value_2, # If condition_2 then value_2
          ...
          TRUE ~ value_n)       # Otherwise, value_n
```

Remark: the `ifelse(condition_if_true, value_true, value_other)` can be used and is easier to manipulate if there is only one condition.

□ **Mathematical operations** – The table below sums up the main mathematical operations that can be performed on columns:

Operation	Command
\sqrt{x}	<code>sqr(x)</code>
$\lfloor x \rfloor$	<code>floor(x)</code>
$\lceil x \rceil$	<code>ceiling(x)</code>

□ **Datetime conversion** – Fields containing datetime values can be stored in two different POSIXt data types:

Action	Command
Converts to datetime with seconds since origin	<code>as.POSIXct(col, format)</code>
Converts to datetime with attributes (e.g. time zone)	<code>as.POSIXlt(col, format)</code>

where `format` is a string describing the structure of the field and using the commands summarized in the table below:

Category	Command	Description	Example
Year	<code>'%Y' / '%y'</code>	With / without century	2020 / 20
Month	<code>'%B' / '%b' / '%m'</code>	Full / abbreviated / numerical	August / Aug / 8
Weekday	<code>'%A' / '%a'</code>	Full / abbreviated	Sunday / Sun
	<code>'%u' / '%w'</code>	Number (1-7) / Number (0-6)	7 / 0
Day	<code>'%d' / '%j'</code>	Of the month / of the year	09 / 222
Time	<code>'%H' / '%M'</code>	Hour / minute	09 / 40
Timezone	<code>'%Z' / '%z'</code>	String / Number of hours from UTC	EST / -0400

Remark: data frames only accept datetime in POSIXct format.

□ **Date properties** – In order to extract a date-related property from a datetime object, the following command is used:

R

```
format(datetime_object, format)
```

where `format` follows the same convention as in the table above.

2.1.3 Data frame transformation

□ **Merging data frames** – We can merge two data frames by a given field as follows:

R

```
merge(df_1, df_2, join_field, join_type)
```

where `join_field` indicates fields where the join needs to happen:

Case	Fields are equal	Different field names
Command	<code>by = 'field'</code>	<code>by.x = 'field_1', by.y = 'field_2'</code>

and where `join_type` indicates the join type, and is one of the following:

Join type	Option	Illustration
Inner join	default	
Left join	all.x = TRUE	
Right join	all.y = TRUE	
Full join	all = TRUE	

Remark: if the by parameter is not specified, the merge will be a cross join.

□ **Concatenation** – The table below summarizes the different ways data frames can be concatenated:

Type	Command	Illustration
Rows	<code>rbind(df_1, ..., df_n)</code>	
Columns	<code>cbind(df_1, ..., df_n)</code>	

□ **Common transformations** – The common data frame transformations are summarized in the table below:

Type	Command	Illustration	
		Before	After
Long to wide	<pre>spread(df, key = 'key', value = 'value')</pre>		
Wide to long	<pre>gather(df, key = 'key', value = 'value', c(key_1, ..., key_n))</pre>		

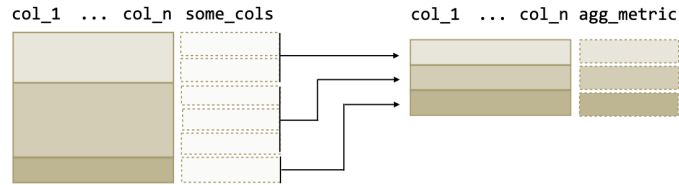
□ **Row operations** – The following actions are used to make operations on rows of the data frame:

Action	Command	Illustration	
		Before	After
Sort with respect to columns	<pre>df %>% arrange(col_1, ..., col_n)</pre>		
Dropping duplicates	<pre>df %>% unique()</pre>		
Drop rows with at least a null value	<pre>df %>% na.omit()</pre>		

Remark: by default, the arrange command sorts in ascending order. If we want to sort it in descending order, the - command needs to be used before a column.

2.1.4 Aggregations

□ **Grouping data** – Aggregate metrics are computed across groups as follows:



The R command is as follows:

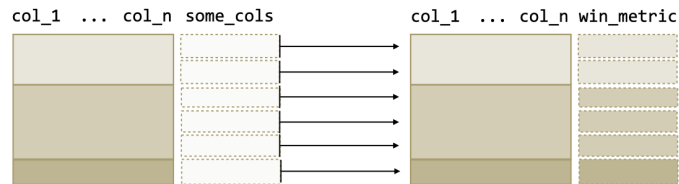
```
R
df %>%
  group_by(col_1, ..., col_n) %>%
  summarize(agg_metric = some_aggregation(some_cols))
# Ungrouped data frame
# Group by some columns
# Aggregation step
```

□ **Aggregate functions** – The table below summarizes the main aggregate functions that can be used in an aggregation query:

Category	Action	Command
Properties	Count of observations	<code>n()</code>
Values	Sum of values of observations	<code>sum()</code>
	Max / min of values of observations	<code>max()</code> / <code>min()</code>
	Mean / median of values of observations	<code>mean()</code> / <code>median()</code>
	Standard deviation / variance across observations	<code>sd()</code> / <code>var()</code>

2.1.5 Window functions

□ **Definition** – A window function computes a metric over groups and has the following structure:



The R command is as follows:

```
R
df %>%
  group_by(col_1, ..., col_n) %>%
  mutate(win_metric = window_function(col))
# Ungrouped data frame
# Group by some columns
# Window function
```

Remark: applying a window function will not change the initial number of rows of the data frame.

□ **Row numbering** – The table below summarizes the main commands that rank each row across specified groups, ordered by a specific field:

Join type	Command	Example
<code>row_number(x)</code>	Ties are given different ranks	1, 2, 3, 4
<code>rank(x)</code>	Ties are given same rank and skip numbers	1, 2.5, 2.5, 4
<code>dense_rank(x)</code>	Ties are given same rank and do not skip numbers	1, 2, 2, 3

□ **Values** – The following window functions allow to keep track of specific types of values with respect to the group:

Command	Description
<code>first(x)</code>	Takes the first value of the column
<code>last(x)</code>	Takes the last value of the column
<code>lag(x, n)</code>	Takes the n^{th} previous value of the column
<code>lead(x, n)</code>	Takes the n^{th} following value of the column
<code>nth(x, n)</code>	Takes the n^{th} value of the column

2.2 Data visualization

2.2.1 General structure

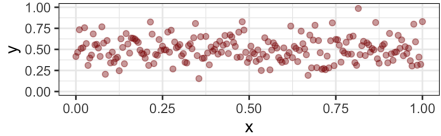
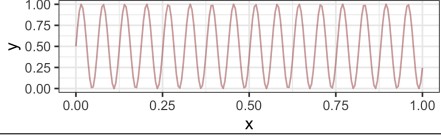
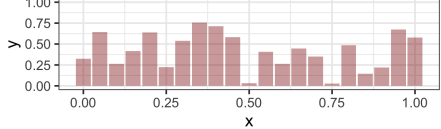
□ **Overview** – The general structure of the code that is used to plot figures is as follows:

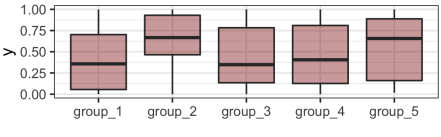
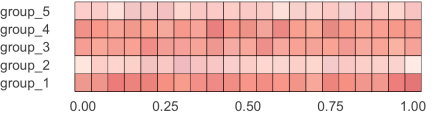
```
R
ggplot(...) +
  geom_function(...) + # Initialization
  facet_function(...) + # Main plot(s)
  labs(...) + # Facets (optional)
  scale_function(...) + # Legend (optional)
  theme_function(...) # Scales (optional)
# Theme (optional)
```

We note the following points:

- The `ggplot()` layer is mandatory.
- When the data argument is specified inside the `ggplot()` function, it is used as default in the following layers that compose the plot command, unless otherwise specified.
- In order for features of a data frame to be used in a plot, they need to be specified inside the `aes()` function.

□ **Basic plots** – The main basic plots are summarized in the table below:

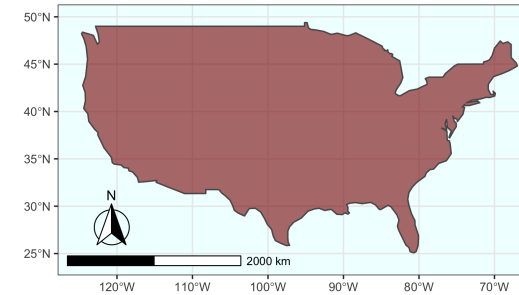
Type	Command	Illustration
Scatter plot	<code>geom_point(x, y, params)</code>	
Line plot	<code>geom_line(x, y, params)</code>	
Bar chart	<code>geom_bar(x, y, params)</code>	

Type	Command	Illustration
Box plot	<code>geom_boxplot(x, y, params)</code>	
Heatmap	<code>geom_tile(x, y, params)</code>	

where the possible parameters are summarized in the table below:

Command	Description	Use case
color	Color of a line / point / border	'red'
fill	Color of an area	'red'
size	Size of a line / point	4
shape	Shape of a point	4
linetype	Shape of a line	'dashed'
alpha	Transparency, between 0 and 1	0.3

□ **Maps** – It is possible to plot maps based on geometrical shapes as follows:



The following table summarizes the main commands used to plot maps:

Category	Action	Command
Map	Draw polygon shapes from the geometry column	<code>geom_sf(data)</code>
Additional elements	Add and customize geographical directions	<code>annotation_north_arrow(1)</code>
	Add and customize distance scale	<code>annotation_scale(1)</code>
Range	Customize range of coordinates	<code>coord_sf(xlim, ylim)</code>

□ **Animations** – Plotting animations can be made using the `gganimate` library. The following command gives the general structure of the code:

```
R
# Main plot
ggplot() +
  ... +
  transition_states(field, states_length)

# Generate and save animation
animate(plot, duration, fps, width, height, units, res, renderer)
anim_save(filename)
```

2.2.2 Advanced features

□ **Facets** – It is possible to represent the data through multiple dimensions with facets using the following commands:

Type	Command	Illustration
Grid (1 or 2D)	<code>facet_grid(row_var ~ column_var)</code>	
Wrapped	<code>facet_wrap(vars(x1, ..., xn), nrow, ncol)</code>	

□ **Text annotation** – Plots can have text annotations with the following commands:

Command	Illustration
<code>geom_text(x, y, label, hjust, vjust)</code>	
<code>geom_label_repel(x, y, label, nudge_x, nudge_y)</code>	

□ **Additional elements** – We can add objects on the plot with the following commands:

Type	Command	Illustration
Line	<code>geom_vline(xintercept, linetype)</code> <code>geom_hline(yintercept, linetype)</code>	
Curve	<code>geom_curve(x, y, xend, yend)</code>	
Rectangle	<code>geom_rect(xmin, xmax, ymin, ymax)</code>	

2.2.3 Last touch

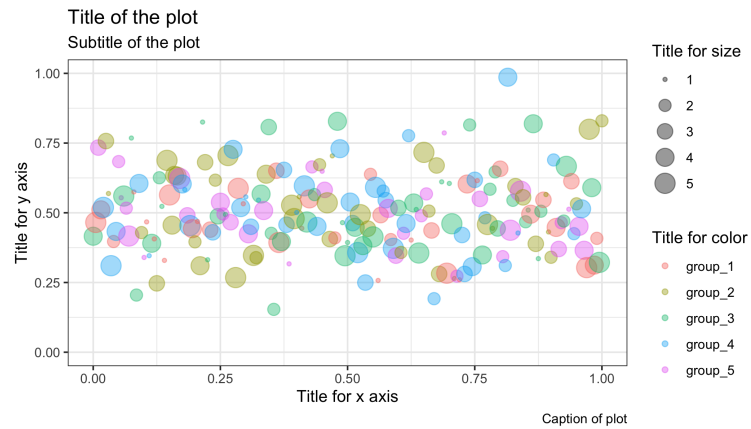
□ **Legend** – The title of legends can be customized to the plot with the following command:

R
plot + labs(params)

where the params are summarized below:

Element	Command
Title / subtitle of the plot	<code>title = 'text' / subtitle = 'text'</code>
Title of the x / y axis	<code>x = 'text' / y = 'text'</code>
Title of the size / color	<code>size = 'text' / color = 'text'</code>
Caption of the plot	<code>caption = 'text'</code>

This results in the following plot:



□ **Plot appearance** – The appearance of a given plot can be set by adding the following command:

Type	Command	Illustration
Black and white	<code>theme_bw()</code>	
Classic	<code>theme_classic()</code>	
Minimal	<code>theme_minimal()</code>	
None	<code>theme_void()</code>	

In addition, `theme()` is able to adjust positions/fonts of elements of the legend.

Remark: in order to fix the same appearance parameters for all plots, the `theme_set()` function can be used.

□ **Scales and axes** – Scales and axes can be changed with the following commands:

Category	Action	Command
Range	Specify range of x / y axis	<code>xlim(xmin, xmax)</code> <code>ylim(ymin, ymax)</code>
Nature	Display ticks in a customized manner	<code>scale_x_continuous()</code> <code>scale_x_discrete()</code> <code>scale_x_date()</code>
Magnitude	Transform axes	<code>scale_x_log10()</code> <code>scale_x_reverse()</code> <code>scale_x_sqrt()</code>

Remark: the `scale_x()` functions are for the x axis. The same adjustments are available for the y axis with `scale_y()` functions.

□ **Double axes** – A plot can have more than one axis with the `sec.axis` option within a given scale function `scale_function()`. It is done as follows:

R

```
scale_function(sec.axis = sec_axis(~ .))
```

□ **Saving figure** – It is possible to save figures with predefined parameters regarding the scale, width and height of the output image with the following command:

R

```
ggsave(plot, filename, scale, width, height)
```

SECTION 3

Working with data with Python

3.1 Data manipulation

3.1.1 Main concepts

□ **File management** – The table below summarizes the useful commands to make sure the working directory is correctly set:

Category	Action	Command
Paths	Change directory to another path	<code>os.chdir(path)</code>
	Get current working directory	<code>os.getcwd()</code>
	Join paths	<code>os.path.join(path_1, ..., path_n)</code>
Files	List files and folders in a directory	<code>os.listdir(path)</code>
	Check if path is a file / folder	<code>os.path.isfile(path)</code>
		<code>os.path.isdir(path)</code>
	Read / write csv file	<code>pd.read_csv(path_to_csv_file)</code> <code>df.to_csv(path_to_csv_file)</code>

□ **Chaining** – It is common to have successive methods applied to a data frame to improve readability and make the processing steps more concise. The method chaining is done as follows:

Python

```
# df gets some_operation_1, then some_operation_2, ..., then some_operation_n
(df
 .some_operation_1(params_1)
 .some_operation_2(params_2)
 ...
 .some_operation_n(params_n))
```

□ **Exploring the data** – The table below summarizes the main functions used to get a complete overview of the data:

Category	Action	Command
Look at data	Select columns of interest	<code>df[col_list]</code>
	Remove unwanted columns	<code>df.drop(col_list, axis=1)</code>
	Look at n first rows / last rows	<code>df.head(n)</code> / <code>df.tail(n)</code>
	Summary statistics of columns	<code>df.describe()</code>
Paths	Data types of columns	<code>df.dtypes</code> / <code>df.info()</code>
	Number of (rows, columns)	<code>df.shape</code>

□ **Data types** – The table below sums up the main data types that can be contained in columns:

Data type	Description	Example
object	String-related data	'teddy bear'
float64	Numerical data	24.0
int64	Numeric data that are integer	24
datetime64	Timestamps	'2020-01-01 00:01:00'

3.1.2 Data preprocessing

□ **Filtering** – We can filter rows according to some conditions as follows:

Python

```
df[df['some_col'] some_operation some_value_or_list_or_col]
```

where some_operation is one of the following:

Category	Operation	Command
Basic	Equality / non-equality	<code>==</code> / <code>!=</code>
	Inequalities	<code><</code> , <code><=</code> , <code>>=</code> , <code>></code>
	And / or	<code>&</code> / <code> </code>
Advanced	Check for missing value	<code>pd.isnull()</code>
	Belonging	<code>.isin([val_1, ..., val_n])</code>
	Pattern matching	<code>.str.contains('val')</code>

□ **Changing columns** – The table below summarizes the main column operations:

Operation	Command
Add new columns on top of old ones	<code>df.assign(new_col=lambda x: some_operation(x))</code>
Rename columns	<code>df.rename(columns={ 'current_col': 'new_col_name' })</code>
Unite columns	<code>df['new_merged_col'] = (df[old_cols_list].agg('-', join, axis=1))</code>

□ **Conditional column** – A column can take different values with respect to a particular set of conditions with the `np.select()` command as follows:

Python

```
np.select(
    [condition_1, ..., condition_n], # If condition_1, ..., condition_n
    [value_1, ..., value_n],        # Then value_1, ..., value_n respectively
    default=default_value           # Otherwise, default_value
)
```

Remark: the `np.where(condition_if_true, value_true, value_other)` command can be used and is easier to manipulate if there is only one condition.

□ **Mathematical operations** – The table below sums up the main mathematical operations that can be performed on columns:

Operation	Command
\sqrt{x}	<code>np.sqrt(x)</code>
$\lfloor x \rfloor$	<code>np.floor(x)</code>
$\lceil x \rceil$	<code>np.ceil(x)</code>

□ **Datetime conversion** – Fields containing datetime values are converted from string to date-time as follows:

Python

```
pd.to_datetime(col, format)
```

where `format` is a string describing the structure of the field and using the commands summarized in the table below:

Category	Command	Description	Example
Year	<code>'%Y' / '%y'</code>	With / without century	2020 / 20
Month	<code>'%B' / '%b' / '%m'</code>	Full / abbreviated / numerical	August / Aug / 8
Weekday	<code>'%A' / '%a'</code>	Full / abbreviated	Sunday / Sun
	<code>'%u' / '%w'</code>	Number (1-7) / Number (0-6)	7 / 0
Day	<code>'%d' / '%j'</code>	Of the month / of the year	09 / 222
Time	<code>'%H' / '%M'</code>	Hour / minute	09 / 40
Timezone	<code>'%Z' / '%z'</code>	String / Number of hours from UTC	EST / -0400

□ **Date properties** – In order to extract a date-related property from a datetime object, the following command is used:

Python

```
datetime_object.strftime(format)
```

where `format` follows the same convention as in the table above.

3.1.3 Data frame transformation

□ **Merging data frames** – We can merge two data frames by a given field as follows:

Python

```
df1.merge(df2, join_field, join_type)
```

where `join_field` indicates fields where the join needs to happen:


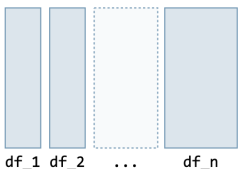
Case	Fields are equal	Fields are different
Command	<code>on='field'</code>	<code>left_on='field_1', right_on='field_2'</code>

and where `join_type` indicates the join type, and is one of the following:

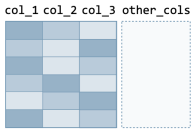
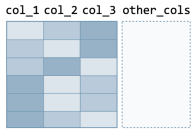
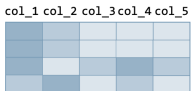
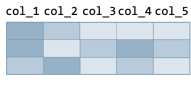
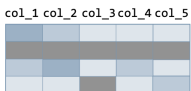
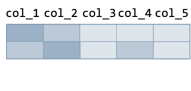
Join type	Option	Illustration
Inner join	<code>how='inner'</code>	
Left join	<code>how='left'</code>	
Right join	<code>how='right'</code>	
Full join	<code>how='outer'</code>	

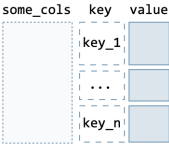

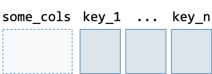
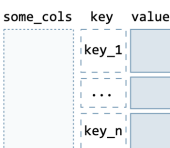
Remark: a cross join can be done by joining on an undifferentiated column, typically done by creating a temporary column equal to 1.

□ **Concatenation** – The table below summarizes the different ways data frames can be concatenated:

Type	Command	Illustration
Rows	<code>pd.concat([df_1, ..., df_n], axis=0)</code>	
Columns	<code>pd.concat([df_1, ..., df_n], axis=1)</code>	

□ **Common transformations** – The common data frame transformations are summarized in the table below:

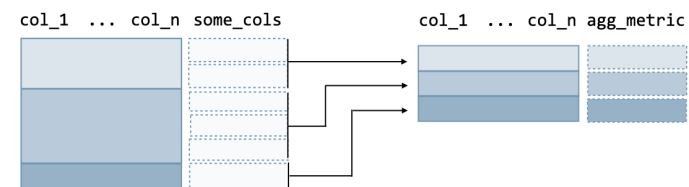
Action	Command	Illustration	
		Before	After
Sort with respect to columns	<code>df.sort_values(by=['col_1', ..., 'col_n'], ascending=True)</code>		
Dropping duplicates	<code>df.drop_duplicates()</code>		
Drop rows with at least a null value	<code>df.dropna()</code>		

Type	Command	Illustration	
		Before	After
Long to wide	<code>pd.pivot_table(df, values='value', index=some_cols, columns='key', aggfunc=np.sum)</code>		
Wide to long	<code>pd.melt(df, var_name='key', value_name='value', value_vars=['key_1', ..., 'key_n'], id_vars=some_cols)</code>		

□ **Row operations** – The following actions are used to make operations on rows of the data frame:

3.1.4 Aggregations

□ **Grouping data** – A data frame can be aggregated with respect to given columns as follows:



The Python command is as follows:

Python

```
(df
.groupby(['col_1', ..., 'col_n'])
.agg({'col': builtin_agg})
```

where builtin_agg is among the following:

Category	Action	Command
Properties	Count of observations	'count'
Values	Sum of values of observations	'sum'
	Max / min of values of observations	'max' / 'min'
	Mean / median of values of observations	'mean' / 'median'
	Standard deviation / variance across observations	'std' / 'var'

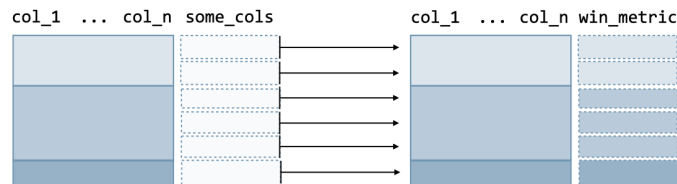
❑ **Custom aggregations** – It is possible to perform customized aggregations by using lambda functions as follows:

Python

```
df_agg = (
    df
    .groupby(['col_1', ..., 'col_n'])
    .apply(lambda x: pd.Series({
        'agg_metric': some_aggregation(x)
    }))
)
```

3.1.5 Window functions

❑ **Definition** – A window function computes a metric over groups and has the following structure:



The Python command is as follows:

Python

```
(df
 .assign(win_metric = lambda x:
        x.groupby(['col_1', ..., 'col_n'])['col_1'].window_function(params))
```

Remark: applying a window function will not change the initial number of rows of the data frame.

❑ **Row numbering** – The table below summarizes the main commands that rank each row across specified groups, ordered by a specific field:

Join type	Command	Example
<code>x.rank(method='first')</code>	Ties are given different ranks	1, 2, 3, 4
<code>x.rank(method='min')</code>	Ties are given same rank and skip numbers	1, 2.5, 2.5, 4
<code>x.rank(method='dense')</code>	Ties are given same rank and do not skip numbers	1, 2, 2, 3

❑ **Values** – The following window functions allow to keep track of specific types of values with respect to the group:

Command	Description
<code>x.shift(n)</code>	Takes the n^{th} previous value of the column
<code>x.shift(-n)</code>	Takes the n^{th} following value of the column

3.2 Data visualization

3.2.1 General structure

❑ **Overview** – The general structure of the code that is used to plot figures is as follows:

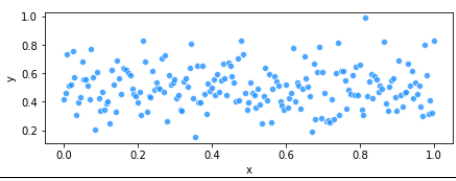
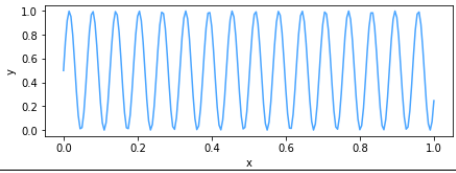
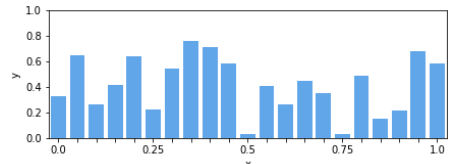
Python

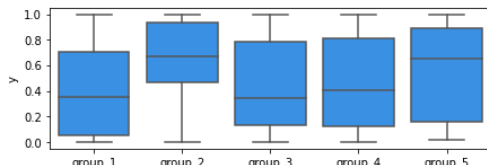
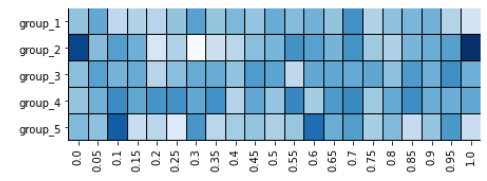
```
# Plot
f, ax = plt.subplots(...)
ax = sns...

# Legend
plt.title()
plt.xlabel()
plt.ylabel()
```

We note that the `plt.subplots()` command enables to specify the figure size.

❑ **Basic plots** – The main basic plots are summarized in the table below:

Type	Command	Illustration
Scatter plot	<code>sns.scatterplot(x, y, params)</code>	
Line plot	<code>sns.lineplot(x, y, params)</code>	
Bar chart	<code>sns.barplot(x, y, params)</code>	

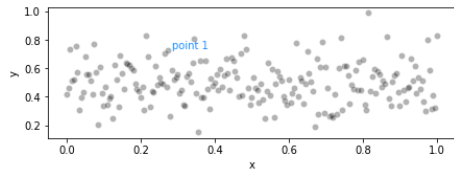
Type	Command	Illustration
Box plot	<code>sns.boxplot(x, y, params)</code>	
Heatmap	<code>sns.heatmap(data, params)</code>	

where the meaning of parameters are summarized in the table below:

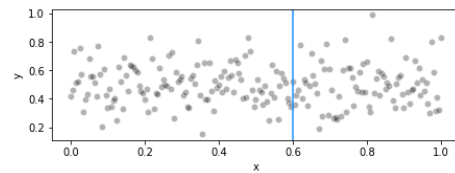
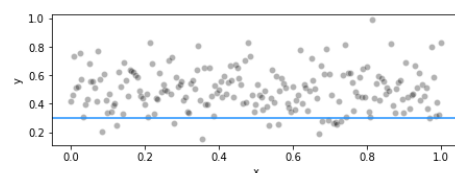
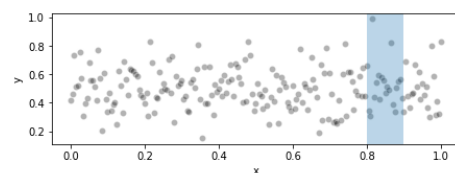
Command	Description	Use case
hue	Color of a line / point / border	'red'
fill	Color of an area	'red'
size	Size of a line / point	4
linetype	Shape of a line	'dashed'
alpha	Transparency, between 0 and 1	0.3

3.2.2 Advanced features

□ **Text annotation** – Plots can have text annotations with the following commands:

Type	Command	Illustration
Text	<code>ax.text(x, y, s, color)</code>	

□ **Additional elements** – We can add objects on the plot with the following commands:

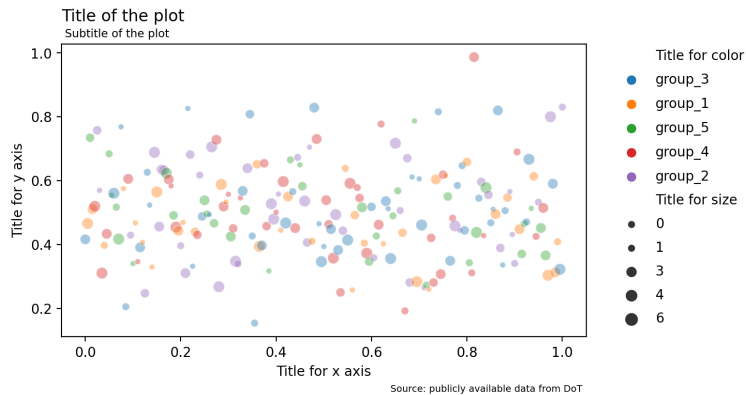
Type	Command	Illustration
Line	<code>ax.axvline(x, ymin, ymax, color, linewidth, linestyle)</code>	
	<code>ax.axhline(y, xmin, xmax, color, linewidth, linestyle)</code>	
Rectangle	<code>ax.axvspan(xmin, xmax, ymin, ymax, color, fill, alpha)</code>	

3.2.3 Last touch

□ **Legend** – The title of legends can be customized to the plot with the commands summarized below:

Element	Command
Title / subtitle of the plot	<code>ax.set_title('text', loc, pad)</code> <code>plt.suptitle('text', x, y, size, ha)</code>
Title of the x / y axis	<code>ax.set_xlabel('text')</code> / <code>ax.set_ylabel('text')</code>
Title of the size / color	<code>ax.get_legend_handles_labels()</code>
Caption of the plot	<code>ax.text('text', x, y, fontsize)</code>

This results in the following plot:



□ **Double axes** – A plot can have more than one axis with the `plt.twinx()` command. It is done as follows:

Python

```
ax2 = plt.twinx()
```

□ **Figure saving** – There are two main steps to save a plot:

- Specifying the width and height of the plot when declaring the figure:

Python

```
f, ax = plt.subplots(1, figsize=(width, height))
```

- Saving the figure itself:

Python

```
f.savefig(fname)
```

SECTION 4

Engineering productivity tips with Git, Bash and Vim

4.1 Working in groups with Git

4.1.1 Overview

□ **Overview** – Git is a version control system (VCS) that tracks changes of different files in a given repository. In particular, it is useful for:

- keeping track of file versions
- working in parallel thanks to the concept of branches
- backing up files to a remote server

4.1.2 Main commands

□ **Getting started** – The table below summarizes the commands used to start a new project, depending on whether or not the repository already exists:

Case	Action	Command	Illustration
No existing repository	Initialize repository from local folder	<code>git init</code>	
Repository already exists	Copy repository from remote to local	<code>git clone git_address</code>	

□ **File check-in** – We can track modifications made in the repository, done by either modifying, adding or deleting a file, through the following steps:

Step	Command	Illustration
1. Add modified, new, or deleted file to staging area	<code>git add file</code>	
2. Save snapshot along with descriptive message	<code>git commit -m 'description'</code>	

Remark 1: `git add .` will have all modified files to the staging area.

Remark 2: files that we do not want to track can be listed in the `.gitignore` file.

❑ **Sync with remote** – The following commands enable changes to be synchronized between remote and local machines:

Action	Command	Illustration
Fetch most recent changes from remote branch	<code>git pull name_of_branch</code>	
Push latest local changes to remote branch	<code>git push name_of_branch</code>	

❑ **Parallel workstreams** – In order to make changes that do not interfere with the current branch, we can create another branch `name_of_branch` as follows:

Bash

```
git checkout -b name_of_new_branch # Create and checkout to that branch
```

Depending on whether we want to incorporate or discard the branch, we have the following commands:

Action	Command	Illustration
Merge with initial branch	<code>git merge initial_branch</code>	
Remove branch	<code>git branch -D name_of_branch</code>	

❑ **Tracking status** – We can check previous changes made to the repository with the following commands:

Action	Command	Illustration
Check status of modified file(s)	<code>git status</code>	
View last commits	<code>git log --oneline</code>	
Compare changes made between two commits	<code>git diff commit_1 commit_2</code>	
View list of local branches	<code>git branch</code>	

❑ **Canceling changes** – Canceling changes is done differently depending on the situation that we are in. The table below sums up the most common cases:

Case	Action	Command	Illustration
Unstaged	Revert file to last commit	<code>git checkout -- file</code>	
Staged	Remove file from staging area	<code>git reset HEAD file</code>	
Committed	Go back to a previous commit	<code>git reset --hard prev_commit</code>	

4.1.3 Project structure

❑ **Structure of folders** – It is important to keep a consistent and logical structure of the project. One example is as follows:

Terminal

```
my_project/
analysis/
graph/
notebook/
data/
```

```

query/
raw/
processed/
modeling/
method/
tests
README.md

```

4.2 Working with Bash

❑ **Basic terminal commands** – The table below sums up the most useful terminal commands:

Category	Action	Command
Exploration	Display list of files (including hidden ones)	<code>ls (-a)</code>
	Show current directory	<code>pwd</code>
	Show content of file	<code>cat path_to_file</code>
	Show statistics of file (lines/words/characters)	<code>wc path_to_file</code>
File management	Make new folder	<code>mkdir folder_name</code>
	Change directory to folder	<code>cd path_to_folder</code>
	Create new empty file	<code>touch filename</code>
	Copy-paste file (folder) from origin to destination	<code>scp (-R) origin destination</code>
	Move file/folder from origin to destination	<code>mv origin destination</code>
	Remove file (folder)	<code>rm (-R) path</code>
Compression	Compress folder into file	<code>tar -czvf comp_folder.tar.gz folder</code>
	Uncompress file	<code>tar -xzvf comp_folder.tar.gz</code>
Miscellaneous	Display message	<code>echo "message"</code>
	Overwrite / append file with output	<code>output > file.txt / output >> file.txt</code>
	Execute <code>command</code> with elevated privileges	<code>sudo command</code>
	Connect to a remote machine	<code>ssh remote_machine_address</code>

❑ **Chaining** – It is a concept that improves readability by chaining operations with the pipe `|` operator. The most common examples are summed up in the table below:

Action	Command
Count number of files in a folder	<code>ls path_to_folder wc -l</code>
Count number of lines in file	<code>cat path_to_file wc -l</code>
Show last n commands executed	<code>history tail -n</code>

❑ **Advanced search** – The `find` command allows the search of specific files and manipulate them if necessary. The general structure of the command is as follows:

Bash

```
find path_to_folder/. [conditions] [actions]
```

The possible conditions and actions are summarized in the table below:

Category	Action	Command
Conditions	Certain names, regex accepted	<code>-name 'certain_name'</code>
	Certain file types (d/f for directory/file)	<code>-type certain_type</code>
	Certain file sizes (c/k/M/G for B/kB/MB/GB)	<code>-size file_size</code>
	Opposite of a given condition	<code>-not [condition]</code>
Actions	Delete selected files	<code>-delete</code>
	Print selected files	<code>-print</code>

Remark: the flags above can be combined to make a multi-condition search.

❑ **Changing permissions** – The following command enables to change the permissions of a given file (or folder):

Bash

```
chmod (-R) three_digits file
```

with `three_digits` being a combination of three digits, where:

- the first digit is about the owner associated to the file
- the second digit is about the group associated to the file
- the third digit is anyone irrespective of their relation to the file

Each digit is one of (0, 4, 5, 6, 7), and has the following meaning:

Representation	Binary	Digit	Explanation
---	000	0	No permission
r--	100	4	Only read permission
r-x	101	5	Both read and execution permissions
rw-	110	6	Both read and write permissions
rwX	111	7	Read, write and execution permissions

For instance, giving read, write, execution permissions to everyone for a given_file is done by running the following command:

```
Bash
chmod 777 given_file
```

Remark: in order to change ownership of a file to a given user and group, we use the command `chown user:group file`.

□ **Terminal shortcuts** – The table below summarizes the main shortcuts when working with the terminal:

Action	Command
Search previous commands	Ctrl + r
Go to beginning / end of line	Ctrl + a / Ctrl + e
Remove everything after the cursor	Ctrl + k
Clear line	Ctrl + u
Clear terminal window	Ctrl + l

4.3 Automating tasks

□ **Create aliases** – Shortcuts can be added to the `~/.bash_profile` file by adding the following code:

```
Bash
shortcut="command"
```

□ **Bash scripts** – Bash scripts are files whose file name ends with `.sh` and where the file itself is structured as follows:

```
Bash
#!/bin/bash
... [bash script] ...
```

□ **Crontabs** – By letting the day of the month vary between 1-31 and the day of the week vary between 0-6 (Sunday-Saturday), a crontab is of the following format:

```
Terminal
*      *      *      *      *
minute hour  day  month  day
          of month of week
```

□ **tmux** – Terminal multiplexing, often known as `tmux`, is a way of running tasks in the background and in parallel. The table below summarizes the main commands:

Category	Action	Command
Session management	Open a new / last existing session	<code>tmux / tmux attach</code>
	Leave current session	<code>tmux detach</code>
	List all open sessions	<code>tmux ls</code>
	Remove session_name	<code>tmux kill-session -t session_name</code>
Window management	Open / close a window	<code>Cmd + b + c / Cmd + b + x</code>
	Move to n^{th} window	<code>Ctrl + b + n</code>

4.4 Mastering editors

□ **Vim** – Vim is a popular terminal editor enabling quick and easy file editing, which is particularly useful when connected to a server. The main commands to have in mind are summarized in the table below:

Category	Action	Command
File handling	Go to beginning / end of line	<code>0 / \$</code>
	Go to first / last line / i^{th} line	<code>gg / G / i G</code>
	Go to previous / next word	<code>b / w</code>
	Exit file with / without saving changes	<code>:wq / :q!</code>
Text editing	Copy line n line(s), where $n \in \mathbb{N}$	<code>nyy</code>
	Insert n line(s) previously copied	<code>p</code>
Searching	Search for expression containing name_of_pattern	<code>/name_of_pattern</code>
	Next / previous occurrence of name_of_pattern	<code>n / N</code>
Replacing	Replace old with new expressions with confirmation for each change	<code>:s/old/new/gc</code>

□ **Jupyter notebook** – Editing code in an interactive way is easily done through Jupyter notebooks. The main commands to have in mind are summarized in the table below:

Category	Action	Command
Cell transformation	Transform selected cell to text / code	Click on cell + <code>m / y</code>
	Delete selected cell	Click on cell + <code>dd</code>
	Add new cell below / above selected cell	Click on cell + <code>b / a</code>

SECTION A

Conversion between R and Python: data manipulation

A.1 Main concepts

□ **File management** – The table below summarizes the useful commands to make sure the working directory is correctly set:

Category	R Command	Python Command
Paths	<code>setwd(path)</code>	<code>os.chdir(path)</code>
	<code>getwd()</code>	<code>os.getcwd()</code>
	<code>file.path(path_1, ..., path_n)</code>	<code>os.path.join(path_1, ..., path_n)</code>
Files	<code>list.files(path, include.dirs = TRUE)</code>	<code>os.listdir(path)</code>
	<code>file_test('-f', path)</code>	<code>os.path.isfile(path)</code>
	<code>file_test('-d', path)</code>	<code>os.path.isdir(path)</code>
	<code>read.csv(path_to_csv_file)</code>	<code>pd.read_csv(path_to_csv_file)</code>
	<code>write.csv(df, path_to_csv_file)</code>	<code>df.to_csv(path_to_csv_file)</code>

□ **Exploring the data** – The table below summarizes the main functions used to get a complete overview of the data:

Category	R Command	Python Command
Look at data	<code>df %>% select(col_list)</code>	<code>df[col_list]</code>
	<code>df %>% head(n) / df %>% tail(n)</code>	<code>df.head(n) / df.tail(n)</code>
	<code>df %>% summary()</code>	<code>df.describe()</code>
Data types	<code>df %>% str()</code>	<code>df.dtypes / df.info()</code>
	<code>df %>% NROW() / df %>% NCOL()</code>	<code>df.shape</code>

□ **Data types** – The table below sums up the main data types that can be contained in columns:

R Data type	Python Data type	Description
character	object	String-related data
factor		String-related data that can be put in bucket, or ordered
numeric	float64	Numerical data
int	int64	Numeric data that are integer
POSIXct	datetime64	Timestamps

A.2 Data preprocessing

□ **Filtering** – We can filter rows according to some conditions as follows:

R

```
df %>%
  filter(some_col some_operation some_value_or_list_or_col)
```

where `some_operation` is one of the following:

Category	R Command	Python Command
Basic	<code>== / !=</code>	<code>== / !=</code>
	<code><, <=, >=, ></code>	<code><, <=, >=, ></code>
	<code>& / </code>	<code>& / </code>
Advanced	<code>is.na()</code>	<code>pd.isnull()</code>
	<code>%in% (val_1, ..., val_n)</code>	<code>.isin([val_1, ..., val_n])</code>
	<code>%like% 'val'</code>	<code>.str.contains('val')</code>

□ **Mathematical operations** – The table below sums up the main mathematical operations that can be performed on columns:

Operation	R Command	Python Command
\sqrt{x}	<code>sqrt(x)</code>	<code>np.sqrt(x)</code>
$\lfloor x \rfloor$	<code>floor(x)</code>	<code>np.floor(x)</code>
$\lceil x \rceil$	<code>ceiling(x)</code>	<code>np.ceil(x)</code>

A.3 Data frame transformation

□ **Common transformations** – The common data frame transformations are summarized in the table below:

Category	R Command	Python Command
Concatenation	<code>rbind(df_1, ..., df_n)</code>	<code>pd.concat([df_1, ..., df_n], axis=0)</code>
	<code>cbind(df_1, ..., df_n)</code>	<code>pd.concat([df_1, ..., df_n], axis=1)</code>
Dimension change	<code>spread(df, key, value)</code>	<code>pd.pivot_table(df, values='some_values', index='some_index', columns='some_column', aggfunc=np.sum)</code>
	<code>gather(df, key, value)</code>	<code>pd.melt(df, id_vars='variable', value_vars='other_variable')</code>

SECTION B

Conversion between R and Python: data visualization

B.1 General structure

□ **Basic plots** – The main basic plots are summarized in the table below:

Type	R Command	Python Command
Scatter plot	<code>geom_point(x, y, params)</code>	<code>sns.scatterplot(x, y, params)</code>
Line plot	<code>geom_line(x, y, params)</code>	<code>sns.lineplot(x, y, params)</code>
Bar chart	<code>geom_bar(x, y, params)</code>	<code>sns.barplot(x, y, params)</code>
Box plot	<code>geom_boxplot(x, y, params)</code>	<code>sns.boxplot(x, y, params)</code>
Heatmap	<code>geom_tile(x, y, params)</code>	<code>sns.heatmap(x, y, params)</code>

where the meaning of parameters are summarized in the table below:

Command	Description	Use case
color / hue	Color of a line / point / border	'red'
fill	Color of an area	'red'
size	Size of a line / point	4
linetype	Shape of a line	'dashed'
alpha	Transparency, between 0 and 1	0.3

B.2 Advanced features

□ **Additional elements** – We can add objects on the plot with the following commands:

Type	R Command	Python Command
Line	<code>geom_vline(xintercept, linetype)</code>	<code>ax.axvline(x, ymin, ymax, color, linewidth, linestyle)</code>
	<code>geom_hline(yintercept, linetype)</code>	<code>ax.axhline(y, xmin, xmax, color, linewidth, linestyle)</code>
Rectangle	<code>geom_rect(xmin, xmax, ymin, ymax)</code>	<code>ax.axvspan(xmin, xmax, ymin, ymax)</code>
Text	<code>geom_text(x, y, label, hjust, vjust)</code>	<code>ax.text(x, y, s, color)</code>