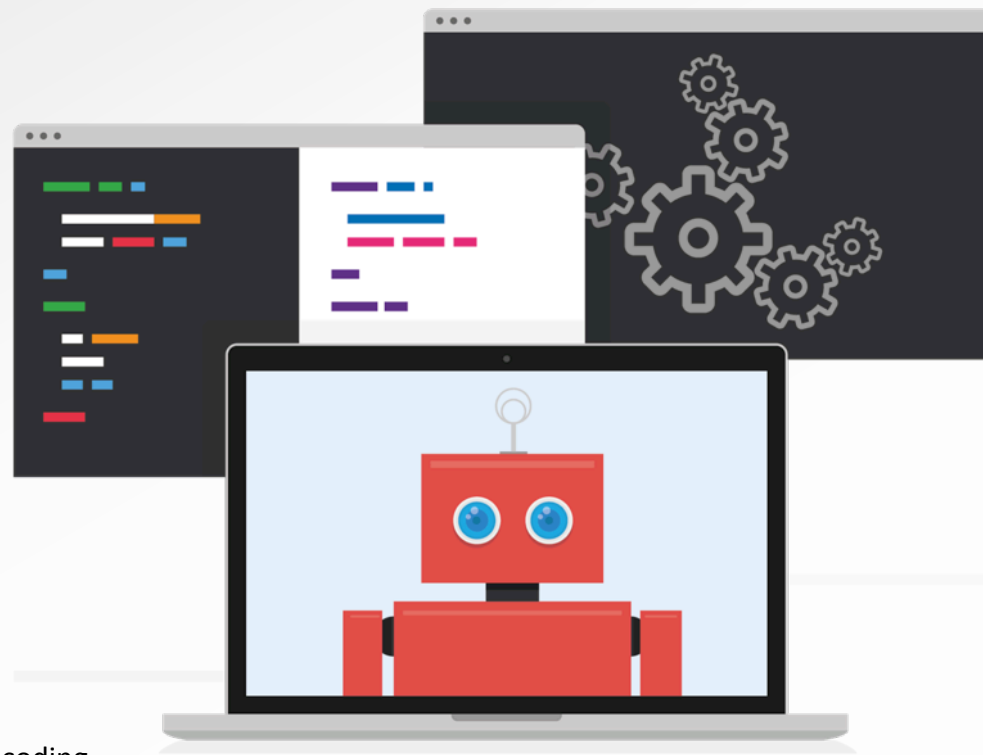


Cambridge Coding Academy

Exploratory data analysis & interactive figures with Plotly



📍 cambridgecoding.com

🐦 [@cambridgecoding](https://twitter.com/cambridgecoding)

📘 facebook.com/cambridgecoding



Agenda

- EDA - When to use it and why?
- The dataset and our tool belt
- Basics in plotting and data summarizing
- Hands-on coding session

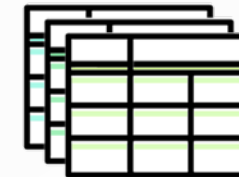
**UNDERSTAND THE
MACHINE LEARNING PROBLEM**



**COLLECT
REAL-WORLD DATA**



**PRE-PROCESS,
CLEAN & MINE**

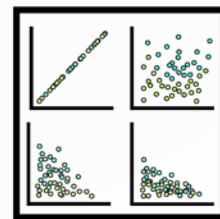


EDA

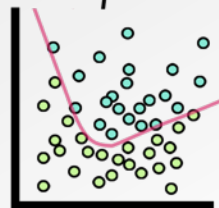


**VISUALISE &
ANALYSE**

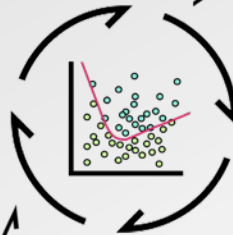
**BUILD
THE MODEL**



**TRAIN &
VALIDATE**



**EVALUATE &
OPTIMISE**



Online Platform

<http://online.cambridgecoding.com/>

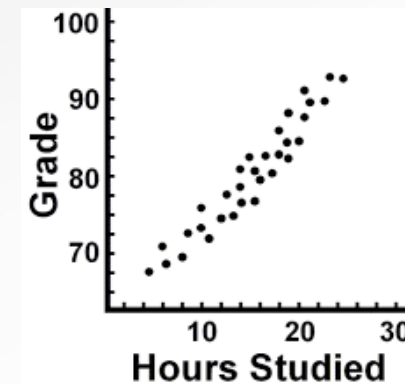
Interested in a two-day data science
bootcamp using Python?

Check out our upcoming bootcamp on 02-03 July in London

<http://cambridgecoding.com/datascience-bootcamp>

What is EDA?

Friend	Number of hours of studying per week	Grade Point Average (out of 5.0)
Allie	14	3.91
Samantha	42	4.98
Hayley	10	3.22
Jessica	32	4.81
Megan	5	2.0
Rachel	10	2.82
Briley	25	3.79
Lauren	18	3.48



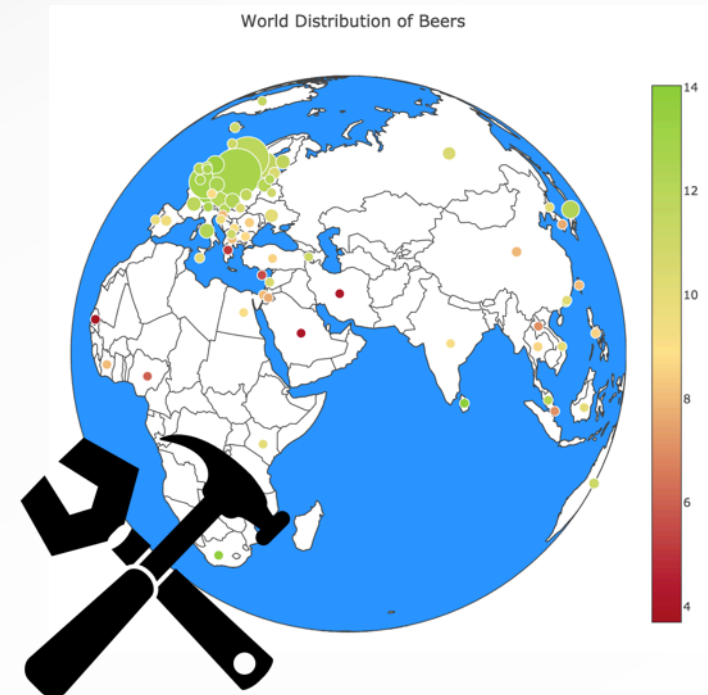
- Mapping data to visual object to make it relatable and understandable
- Goes hand in hand with statistics and machine learning
- Communicate your decision to somebody efficiently

The Dataset

What country in the world produces the most and the best beer?



<http://www.cambridgebeerfestival.com/>



Our EDA setup

Python



Tool belt

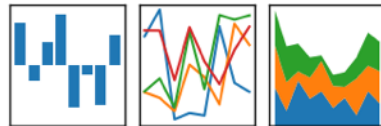
plotting



data processing

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Coding environment



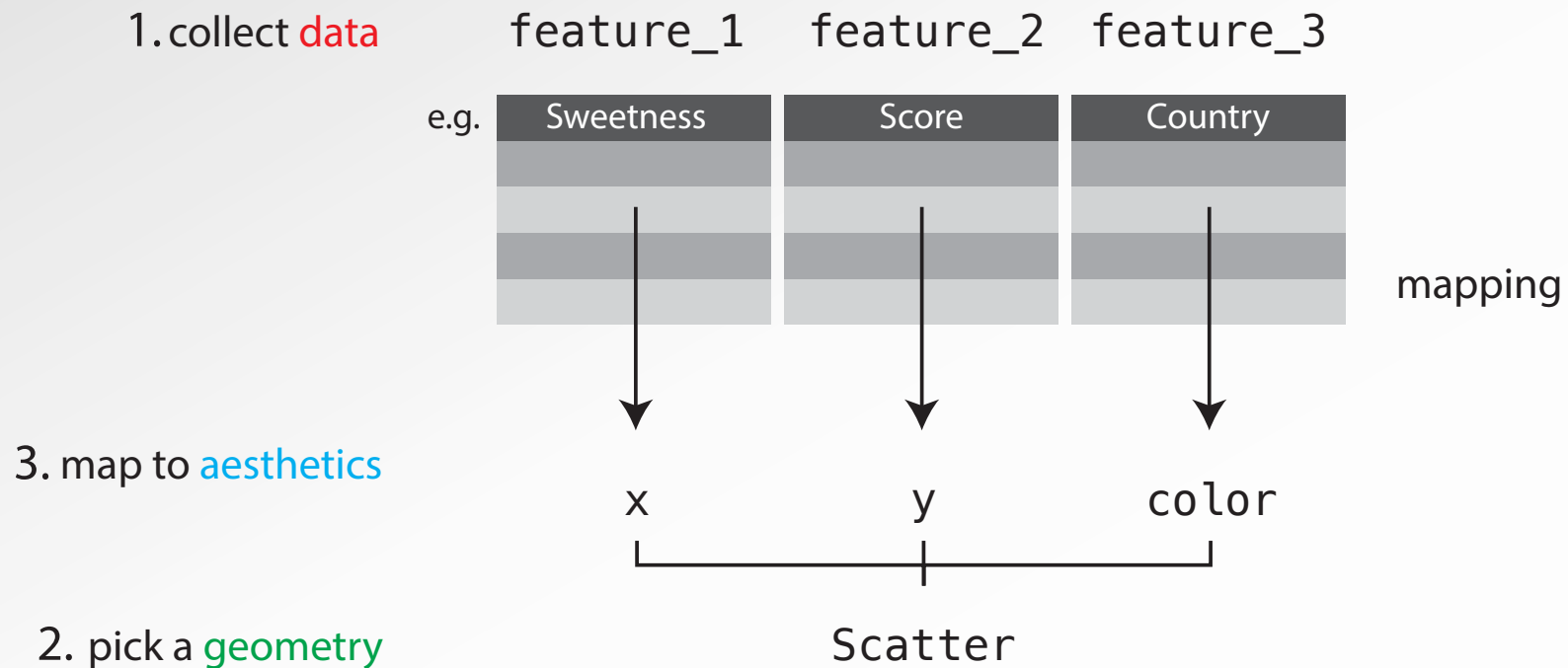
What is Plotly



<https://plot.ly/>

- Produce interactive graphics
- Share with plots with people (if registered)
- Integrates with other frameworks (R, Python...)
- Easy to pick up
- Many different plots types available

Plotting basics

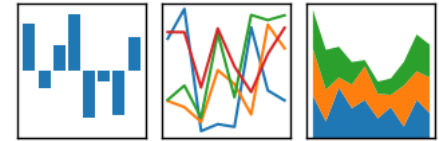


```
trace = go.Scatter(x = feature_1 , y = feature_2, color = feature_3)
```

What is Pandas

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$

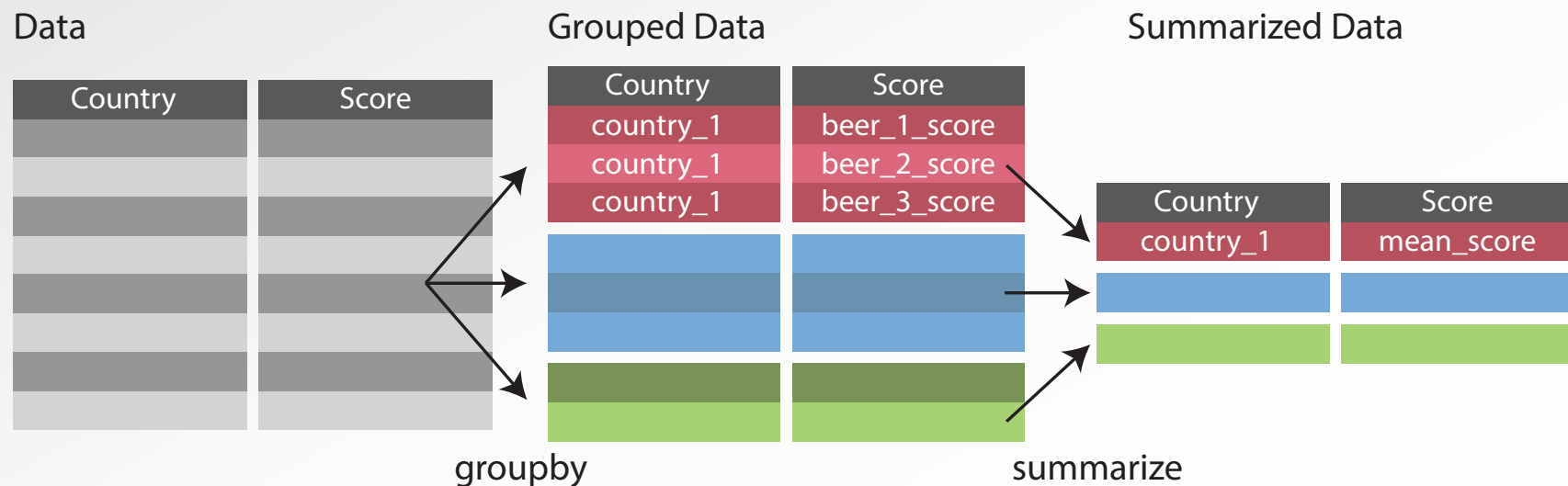


<http://pandas.pydata.org/>

- Convenient data analysis in python
- Easy to read in data
- Reshaping of and preprocessing data
- Filtering data
- Calculate summary statistics

Basics of data processing

- Break process up into steps

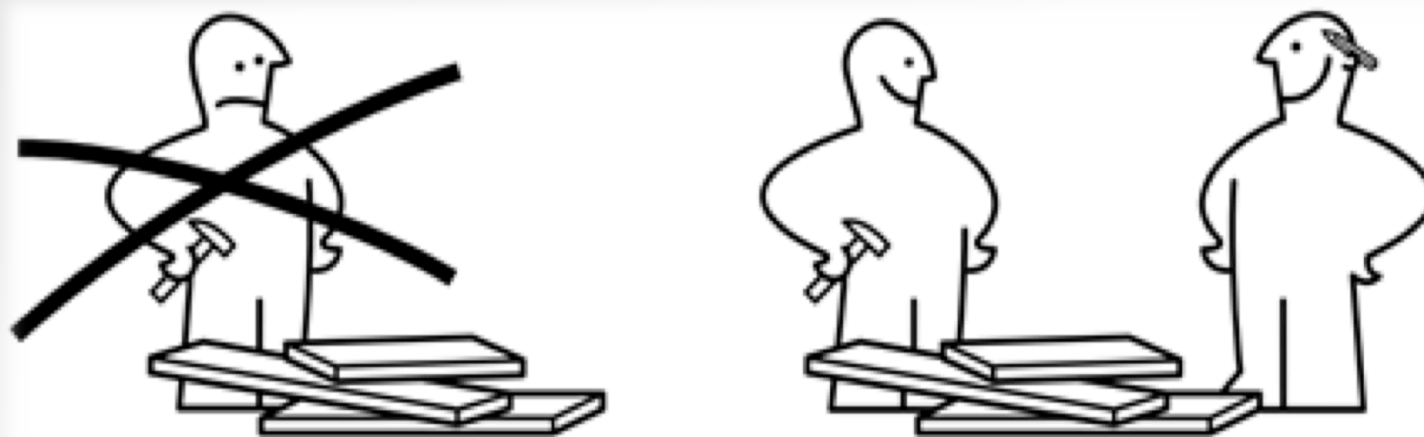


```
dataset['Score'].groupby(dataset['Country']).mean()
```



Hands-on session

Pair Programming



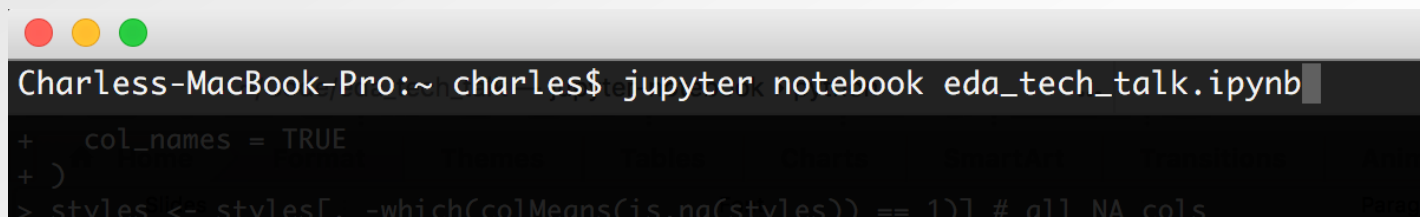
- Reduces risk, bug fixing and overall working time
- Shared knowledge, discussion and QA

Change code and see what happens.



- In the interest of time the code is all there for you
 - run it
 - change it
 - break it
 - learn from it
- Notebook is standalone - go at your own pace!
- Take the notebook home and run your own data

Using Jupyter

- Download the workshop from: https://github.com/cnjr2/eda_tech_talk
- Start the notebook:

A terminal window with a dark background and light text. The title bar shows three colored circles (red, yellow, green). The prompt is 'Charles-MacBook-Pro:~ charles\$'. The command entered is 'jupyter notebook eda_tech_talk.ipynb'. Below the command, there is a code snippet: '+ col_names = TRUE', '+)', and '> styles <- styles[-which(colMeans(is.na(styles)) == 1)] # all NA cols'.

```
Charles-MacBook-Pro:~ charles$ jupyter notebook eda_tech_talk.ipynb
+ col_names = TRUE
+ )
> styles <- styles[-which(colMeans(is.na(styles)) == 1)] # all NA cols
```

- Double-click a cell to edit its contents (its border will become green)
- Execute the code in a cell by pressing 'ctrl+ enter' or press  button.
- You can add cells with plus  button
- Put short chunks of code into each cell to go step-by-step
- Inactivate a line of code preceding it with a #; the re-run the chunk to see what changes!