

Swiss Institute of  
Bioinformatics

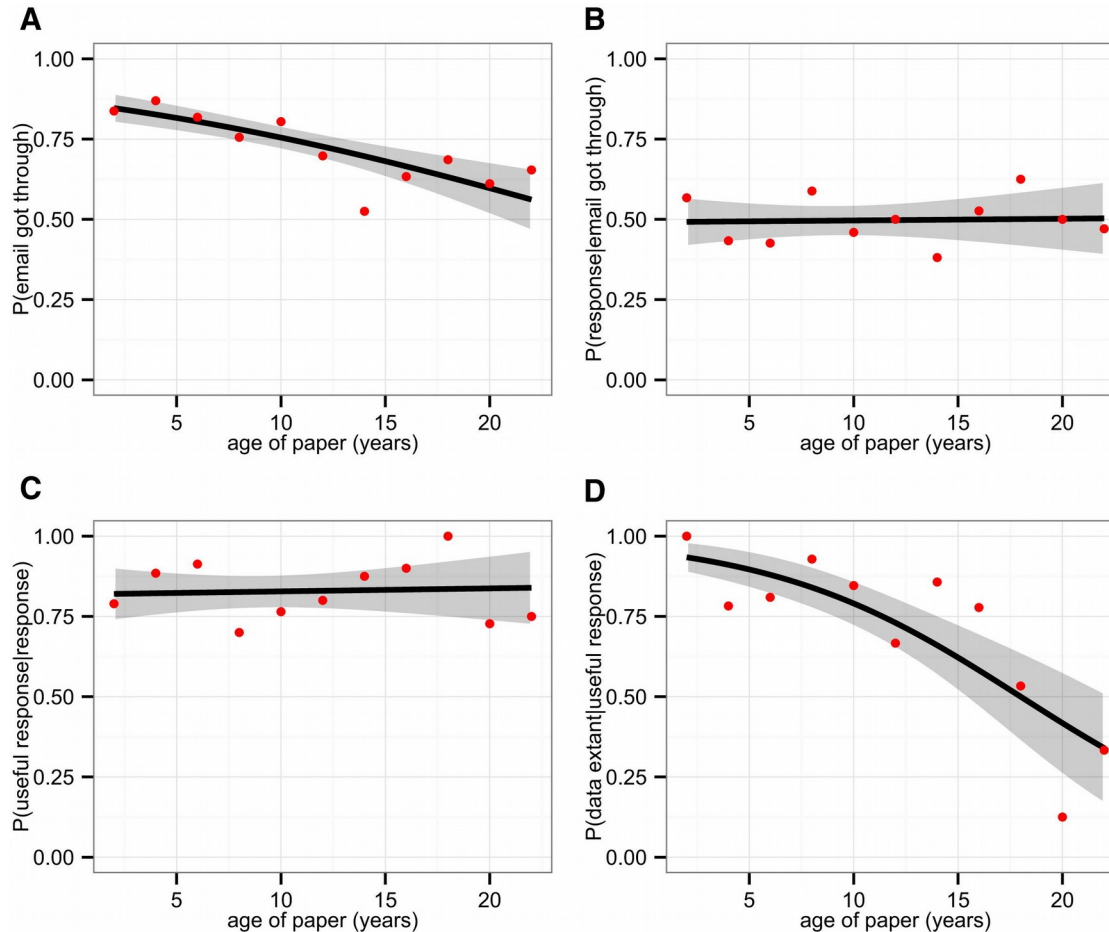
# First steps with Python in life sciences

Wandrille Duchemin

Robin Engler

Orlin Topalov

# The reproducibility crisis



Survey of 516 studies:

- - 17% data availability per year.
- Only 19% retrieval rate after 10 year...

# The **FAIR** guiding principles

---



**Findable:** Metadata and data should be easy to find for both humans and computers (unique global identifier, rich description, machine readable and searchable).



**Accessible:** The data can be retrieved using a standard communication protocol (e.g. https or sftp). Where needed, authentication and authorization procedures are available and documented.



**Interoperable:** the (meta)data should be based on standardized vocabulary and ontologies (categories and their relations), so that it integrates with existing applications and workflows.



**Reusable:** Metadata and data should be well described so that data can be replicated and/or combined in different research settings (rich metadata, clear license term, origin of data, data meets domain-relevant community standards).

# FAIR applied to code

```
1 def f(x, y):  
2     if x==[]: return y  
3     if y==[]: return x  
4     return [x[0]] + f(x[1:], y) if x[0]<y[0] else [y[0]] + f(x, y[1:])  
5  
6 def S(a, n):  
7     m = n//2  
8     return a if n<=1 else f(S(a[:m], m), S(a[m:], n-m))  
9  
10
```

← Don't do this

# FAIR applied to code

```
1 def f(x, y):
2     if x==[]: return y
3     if y==[]: return x
4     return [x[0]] + f(x[1:], y) if x[0]<y[0] else [y[0]] + f(x, y[1:])
5
6 def S(a, n):
7     m = n//2
8     return a if n<=1 else f(S(a[:m], m), S(a[m:], m))
9
10
```

← Don't do this

Do that →

```
2 def merge(x, y):
3     """
4     *RECURSIVE FUNCTION*
5     Takes 2 sorted lists and
6     returns the sorted list merging them.
7     """
8     ## if one of the 2 lists are empty, just return the other one
9     if x==[]:
10         return y
11     if y==[]:
12         return x
13
14     # x and y are sorted
15     if x[0]<y[0]:
16         # lowest element of x goes first
17         # --> merged list is first element of x + merging of y and the remainder of x
18         return [x[0]] + merge(x[1:], y)
19     else:
20         # lowest element of y goes first
21         # same logic as above
22         return [y[0]] + merge(x, y[1:])
23
24
25 def sort(a):
26     ''' *RECURSIVE FUNCTION*
27     Takes a list <a>, and returns a sorted version.
28     Uses the merge sort algorithm.
29     Implementation freely inspired from: http://www.rosettacode.org/wiki/Sorting\_algorithms/Merge\_sort#Python
30     '''
31     n = len(a)
32     m = n//2
33
34     if n<=1: # list empty or with a single element --> return as is
35         return a
36     else: ## split the list in two. Sort separately and then merge the 2 sorted halves
37         return merge(sort(a[:m]), sort(a[m:]))
```

# FAIR applied to code

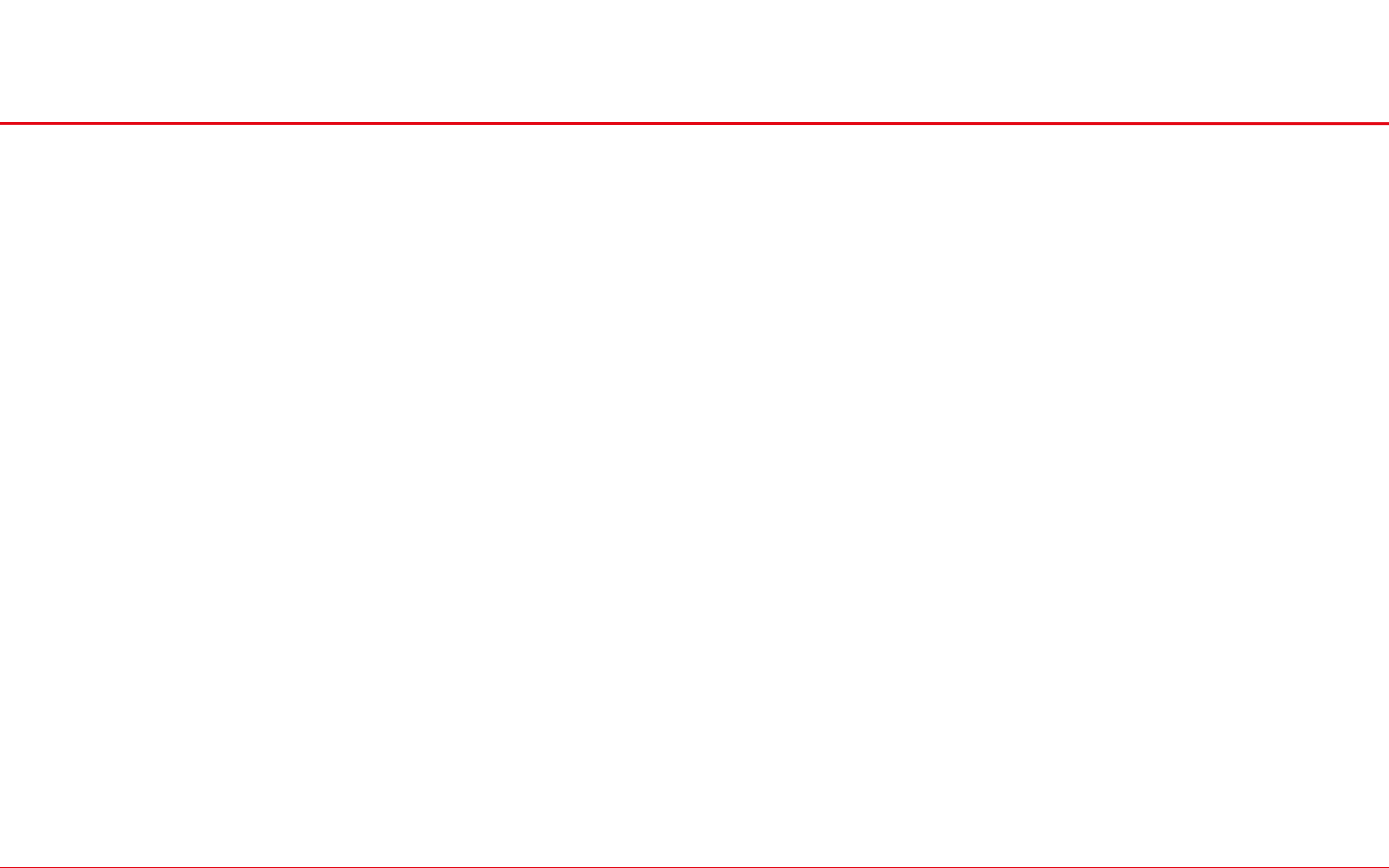
---

- Write code that also acts as documentation, and clearly communicates the analysis.
- Apply the standard you would expect of a 'wet-lab' protocol.
- Will a reasonably competent colleague understand your code ?
- Will you understand your code in 6 month ??

## **To achieve this, you should:**

- Comment as much as needed.
- Use explicit names when naming things (variables, functions, classes).
- Possibly use a support that allows to easily mix code and text, e.g. Jupyter notebook or Jupyter-lab.
- Also, you can look at:

Schwen LO, Rueschenbaum S (2018) Ten quick tips for getting the most scientific value out of numerical data. PLoS Comput Biol 14(10): e1006141.  
<https://doi.org/10.1371/journal.pcbi.1006141>.



# How to execute python code

---

Three main modes of interaction:

- ❑ Interactive console.
- ❑ Python code file (.py files).
- ❑ Jupyter Notebook (.ipynb files).

# How to execute python code

---

Three main modes of interaction:

- Interactive console.
  - Nice for quick debugging, testing syntax,...
- Python code file (.py files).
- Jupyter Notebook (.ipynb files).

# How to execute python code

---

Three main modes of interaction:

- Interactive console.
- Python code file (.py files).
  - Scripts, program, ...
  - most used form
- Jupyter Notebook (.ipynb files).

# How to execute python code

---

Three main modes of interaction:

- Interactive console.
- Python code file (.py files).
- Jupyter Notebook (.ipynb files).
  - Great for data analysis and teaching

# Python – using the console

```
(base) wandrille@wandrille-Latitude-7400:~$ python
Python 3.7.4 (default, Aug 13 2019, 20:35:49)
[GCC 7.3.0] :: Anaconda, Inc. on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> a=3
>>> b=4
>>> a*b + 7
19
>>> print("this is an interactive console")
this is an interactive console
>>> █
```

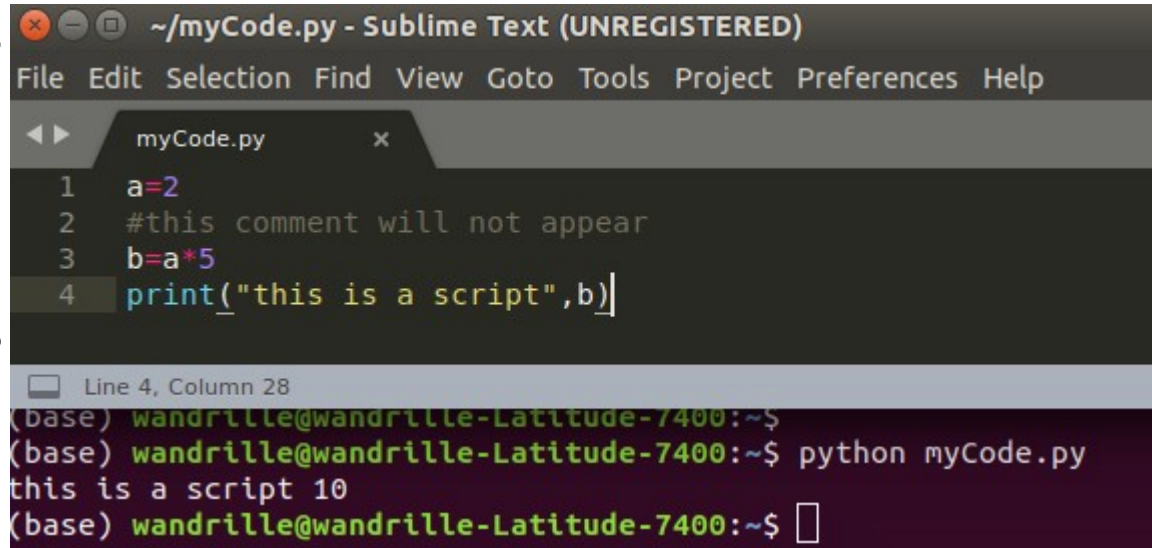
- ❑ **Interactive** : the code is executed as you press 'Enter'
- ❑ Great for quickly **testing** things out.
- ❑ But, you keep **no trace** of your workflow/environment.



Only use it to test little bits of code

# Python – writing a code in a .py file

- ❑ Write a script in a .py text file, then execute it.
- ❑ Main way python code is shared.
- ❑ Ideal for standalone programs and modules.
- ❑ Code and results are kept separate (may be a good or a bad thing).



The screenshot shows a Sublime Text editor window titled `~/myCode.py - Sublime Text (UNREGISTERED)`. The editor displays a Python script in `myCode.py` with the following code:

```
1 a=2
2 #this comment will not appear
3 b=a*5
4 print("this is a script",b)
```

Below the editor, a terminal window shows the execution of the script:

```
(base) wandrille@wandrille-Latitude-7400:~$
(base) wandrille@wandrille-Latitude-7400:~$ python myCode.py
this is a script 10
(base) wandrille@wandrille-Latitude-7400:~$
```



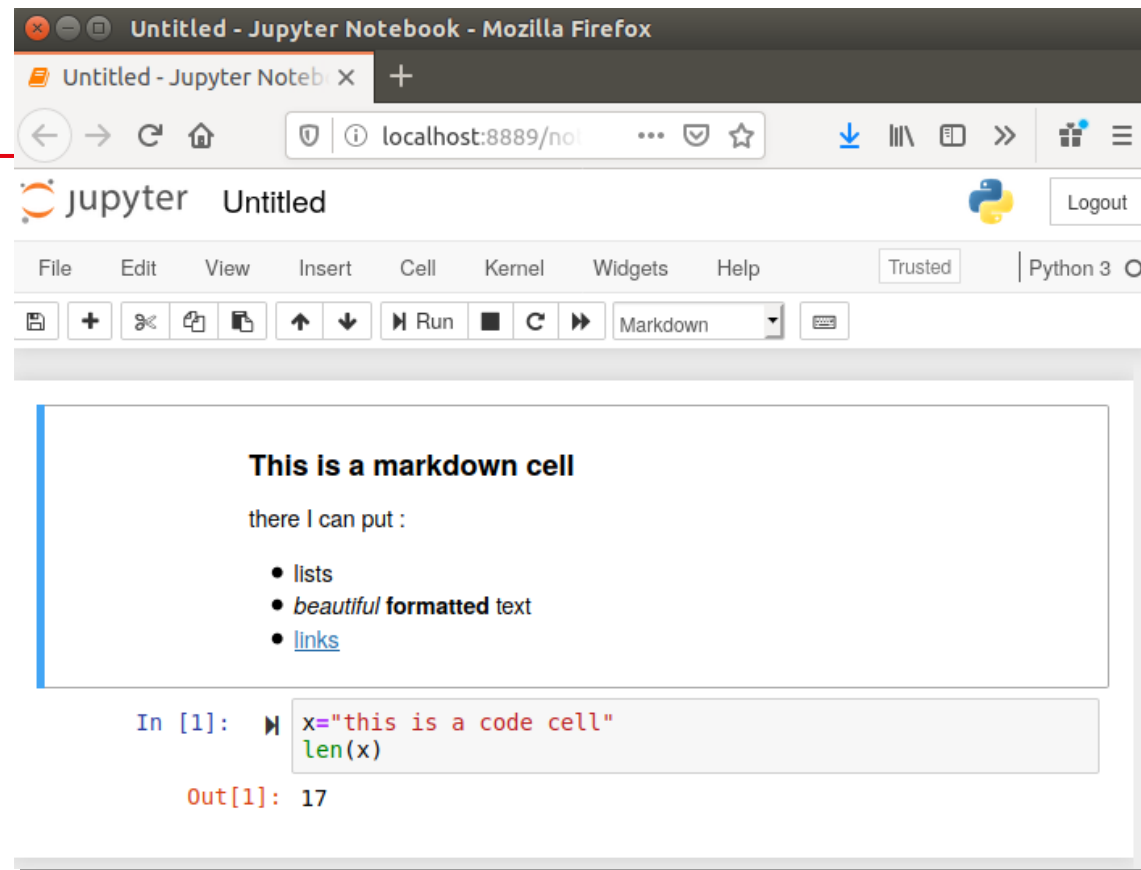
Good for general purpose coding, “operational script”.

# Jupyter notebook / Jupyter Lab

- Browser based interface (but runs locally on your machine).
- Interlace Markdown and code 'cells'.
- Execute code cell by cell.
- **Commentary, code, and results together in the same file.**
- Visually pleasant and fairly ergonomic.



The ability to have documentation, code and results in a single file makes it ideal for **data analysis**: helps reproducibility of results.



## Jupyter notebook / Jupyter Lab - magic commands

---

- Jupyter add a bunch of magic commands on top of the python syntax
- Start with `%` (applies to line) or `%%` (applies to whole cell)
- A lot of very useful recipes:
  - `%%time` : report execution time of a cell
  - `%%bash` : run this cell with bash instead of python
  - `%%writefile` : write cell content to a file  
(eg. to export some code to a script)

<https://ipython.readthedocs.io/en/stable/interactive/magics.html>