# Introduction to Data Science

## with Python

Alexis Bogroff

May 15, 2022

Alexis Bogroff
Lecturer and Mentor in Data Science
at Paris 1 Panthéon-Sorbonne, ESILV,
Openclassrooms, EM-Lyon

# Presenter



Alexis Bogroff
Lecturer and Mentor in Data Science
at Paris 1 Panthéon-Sorbonne, ESILV,
Openclassrooms, EM-Lyon

- 4 years Teaching Assistant and lecturer in VBA, Python for finance, SQL, Data Analysis and Data Science
- 9 months Researcher Assistant at Paris 1 Panthéon-Sorbonne within H2020 European Project
- 1 year Data Scientist at Pléiade Asset Management

# Why Python?

- Versatile
- Simple
- Open Source
- Most used for Data Science

- Interactive console



```
(eda) alexisbogroff@Alexiss-MacBook-Pro src % python
Python 3.9.2 | packaged by conda-forge | (default, Feb 21 2021, 05:00:30)
[Clang 11.0.1 ] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> print("hello")
hello
>>> x = 1
>>> x += 2
>>> x
3
>>> def say_hello():
...     print("hello")
...
>>> say_hello()
hello
>>>
```

# Programming Environment



- Interactive console
- Scripts

# Programming Environment



- Interactive console
- Scripts
- Jupyter Notebooks

# Programming Environment



- Interactive console
- Scripts
- Jupyter Notebooks
- Code editors
  (VS Code)

- Interactive console
- Scripts
- Jupyter Notebooks
- Code editors
  (VS Code)
- Packages managers
  (Conda, Pip)

- Interactive console
- Scripts
- Jupyter Notebooks
- Code editors
  (VS Code)
- Packages managers
  (Conda, Pip)
- Packages / Libraries
  (Pandas, Matplotlib)



Python 3.7

Pandas    25%

NumPy    25-40%

learn    50-75%

matpl&tlib    20-25%

Seaborn Plots

# Programming Environment

- Interactive console
- Scripts
- Jupyter Notebooks
- Code editors
  (VS Code)
- Packages managers
  (Conda, Pip)
- Packages / Libraries
  (Pandas, Matplotlib)
- Virtual Environments

- Virtualenv
- Pipenv
- Venv
- Poetry

- Interactive console
- Scripts
- Jupyter Notebooks
- Code editors
  (VS Code)
- Packages managers
  (Conda, Pip)
- Packages / Libraries
  (Pandas, Matplotlib)
- Virtual Environments
- Version Control Systems
  (Git, Github, Gitlab)

# Essentials in the Python language

- Data types and structures
  - Numbers
  - Text (strings)
  - Iterables
  - Other

```python
# Text and numbers
12          # int (integer)
1.5         # float
'hola'      # str (string)
"hola"
"""hola"""

# Iterables
[42, 58, 209, 42]  # list
(42, 58, 209, 42)  # tuple
{42, 58, 209}      # set
{'name': ['akiko', 'julie'], 'age': [12, 43]}  # dict (dictionary)
```

- Operators
  - Greater than, lower than

Class methods:

- Greater than \_\_gt\_\_
- Lower than \_\_lt\_\_

```python
print(1 > 2)
print(1 < 2)
print(1 < 1)
print('a' < 'b')
print('a' > 'b')
print([1] < [2, 3])
print([1] > [2, 3])
```
✓ 0.2s

```
False
True
False
True
False
False
True
```

- Operators
    - Greater than, lower than
    - Greater or equal than, lower or equal than

Class methods:

- Greater than or equal to __ge__
- Lower than or equal to __le__

```
print(1 >= 2)
print(1 <= 2)
print(1 <= 1)
```
✓ 0.2s

```
False
True
True
```

# Essentials in the Python language

- Operators
  - Greater than, lower than
  - Greater or equal than, lower or equal than
  - Equals to, different from

Class methods:

- Equal to __eq__
- Different from __ne__

```
print(1 == 2)
print(1 != 2)
print(1 == 1)
print(1 != 1)
✓ 0.3s
```

```
False
True
True
False
```

- Operators
    - Greater than, lower than
    - Greater or equal than, lower or equal than
    - Equals to, different from
    - in

Class method:

- Find element in object __contains__

Available in iterables, not in numbers.

```
print(1 in [1, 4, 2])
print(1 in [4, 2])
print([1] in [1, 4, 2])
print([1] in [[1], 4, 2])
print('a' in 'oisj')
print('a' in 'oiasj')
```

✓ 0.2s

```
True
False
False
True
False
True
```

- Operators
    - Greater than, lower than
    - Greater or equal than, lower or equal than
    - Equals to, different from
    - in
    - not

Negation operator:

- not
- ~

```python
print(1 == 1)
print(not 1 == 1)
print(~ 1 == 1)
print(1 != 1)
print(not 1 != 1)
print(~ 1 != 1)
```

✓ 0.2s

```
True
False
False
False
True
True
```

- Operators
    - Greater than, lower than
    - Greater or equal than, lower or equal than
    - Equals to, different from
    - in
    - not
    - is

Not a class method:

- Check if element is (True, False, None, np.nan)

```python
print(True is True)
print(True is False)

print([1] is None)
x = 12
print(x is None)
x = None
print(x is None)

import numpy as np
print(x is np.nan)
```
✓ 0.2s

```
True
False
False
False
True
False
```

- Conditional structures
  - if else statement



Simple "if, else" condition

```
x = 3

if x == 3:
    print("Yes, x is equal to 3")
else:
    print("No, x is not equal to 3")
```
✓ 0.2s

```
Yes, x is equal to 3
```

- Control structures
  - for loop

```
for i in [2, 54, 39]:
    print(i)
✓ 0.1s

2
54
39
```

- Functions
  - A function can:
    - take no, to many arguments
    - arguments can be "Positional" or "Keyword" arguments
    - return nothing (None) or anything (to many things)
    - synonyms: arguments / parameters / inputs
  - One must:
    - Define the function
    - Call the function

# Functions

**Simplest form:**
- No argument required
- No return

```python
# Define the function
def say_something():
    print("Something")

# Call the function
say_something()
```
✓ 0.3s

Something

**Function with:**
- an argument
- no return

```python
def say_my_name(name):
    print(name)

say_my_name("Alexis")
```
✓ 0.2s

Alexis

**Function with:**
- an argument
- a return

```python
def square(x):
    return x**2

result = square(4)
result
```
✓ 0.3s

16

**Function with:**
- no argument
- multiple returns

```python
def return_many_things():
    return 'alexis', 'bogroff', 'data'

return_many_things()
```
✓ 0.2s

('alexis', 'bogroff', 'data')

**Function with:**
- multiple arguments
- a return

```python
def add(a, b, c):
    return a + b + c

result = add(4, 2, 9)
result
```
✓ 0.2s

15

**Function with:**
- no argument
- multiple returns

```python
def return_many_things():
    return 'alexis', 'bogroff', 'data'

return_many_things()
```
✓ 0.2s

('alexis', 'bogroff', 'data')

# Functions

Function with:

- a keyword argument
  - is thus optional
  - must be positioned after the positional arguments
- no return

```python
def say_what_you_doing(name, course='data'):
    print(f"{name} doing {course}")

say_what_you_doing("Alexis")
say_what_you_doing("Alexis", "writing the course")
```
✓ 0.2s

```
Alexis doing data
Alexis doing writing the course
```

Function with:

- a (positional) argument awaiting a function
- no return

```python
def complex_fct(func):
    print("This function will say")
    func()

complex_fct(say_something)
```
✓ 0.2s

```
This function will say
Something
```

```python
class Truc:
    # Define instanciator (init)
    def __init__(self):
        self.age = 10
        self.name = 'truc'
✓ 0.3s


truc_1 = Truc()
✓ 0.2s


truc_1
✓ 0.7s
<__main__.Truc at 0x1077e4340>


print(truc_1.name)
print(truc_1.age)

✓ 0.3s
truc
10
```

- Define init
- Define properties / attributes (internal variables)
- Access through **self**
- Instanciate
- Object reference
- Access properties

# Objects - method

```python
class Truc:

    # Define instanciator (init)
    def __init__(self):
        self.age = 10
        self.name = 'truc'

    def present(self):
        print(f"My name is: {self.name}, "
              f"I'm {self.age} years old")
```
✓ 0.3s

```python
truc_1 = Truc()
```
✓ 0.2s

```python
truc_1
```
✓ 0.2s

```
<__main__.Truc at 0x33bffdb80>
```

```python
truc_1.name
```
✓ 0.4s

```
'truc'
```

```python
truc_1.present()
```
✓ 0.4s

```
My name is: truc, I'm 10 years old
```

- Create method
- Pass **self** argument
- Access attributes via **self.attribute**
- Re-instanciate object ⚠
- New reference
- Access method via **self.method**

```python
class Truc:

    # Define instanciator (init)
    def __init__(self):
        self.age = 10
        self.name = 'truc'

    def present(self):
        print(f"My name is: {self.name}, "
              f"I'm {self.age} years old")

    def dog_age(self):
        return self.age * 7
```
✓ 0.3s

```python
truc_1 = Truc()
```
✓ 0.3s

```python
dog_age = truc_1.dog_age()
```
✓ 0.3s

```python
dog_age
```
✓ 0.3s

```
70
```

- Create method with return
- Set variable using return value

# Objects - init with arguments

```python
class Truc:

    # Define instanciator (init)
    def __init__(self, name, age):
        self.age = age
        self.name = name

    def present(self):
        print(f"My name is: {self.name}, "
              f"I'm {self.age} years old")

    def dog_age(self):
        return self.age * 7
```
✓ 0.3s

```python
truc_1 = Truc('Yuko', 7)
truc_2 = Truc('Mila', 13)
```
✓ 0.5s

```python
print(f"{truc_1.name} is {truc_1.dog_age()} old")
print(f"{truc_2.name} is {truc_2.dog_age()} old")
```
✓ 0.2s

```
Yuko is 49 old
Mila is 91 old
```
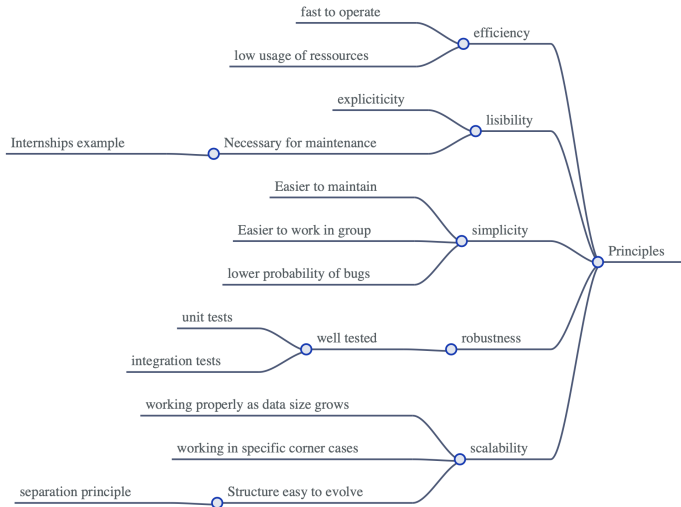
- Set specific init values
- Create different objects

# Coding conventions

- Names
  - variables: snake_style
  - constants: CAPITAL_SNAKE_STYLE
  - functions: snake_style
  - classes: First_letter_capital
- Spaces
- Max number characters by row: **79**
- Creation of iterables (lists, dicts)
- Comments
- File, functions, classes description

## Coding conventions

- Code order in a file:
  - Description
  - Imports
  - Constants
  - Functions and Classes alphabetically
  - Body (functions calls, loops, variables)
- Code organisation between files (script):
  - Main file
  - Functions and Classes file
  - Settings file
- Code organisation files (Jupyter Notebook):
  - Load and prepare data file
  - Analysis file
  - Predictions file

# Programming in general, good practices

# Programming in general, good practices

- Vectorization
- Don't use loops when its possible to vectorize
- Same in Python, Matlab
- *This code is explained in the next course*

```python
# /!\ Never do it this way (loops) /!\
for i, val in enumerate(df_temp['nums']):
    df_temp.loc[i, 'nums'] = val * 2
```

```python
# Do this way (vectorized)
df_temp['nums'] = df_temp['nums'] * 2
```

の

- Difference between programming for:
    - Analysis, statistics
    - Software development
        - Front-end
        - Back-end

# How to learn programming

- Trial and error
- Could be enough at first:
    - Python official documentation
    - Exercises Coding game
    - Google, Stackoverflow, Blogs
- Progress:
    - Choose project (company, Kaggle, personal)
    - Peers: open-source project, Data For Good
    - MOOC: advanced course