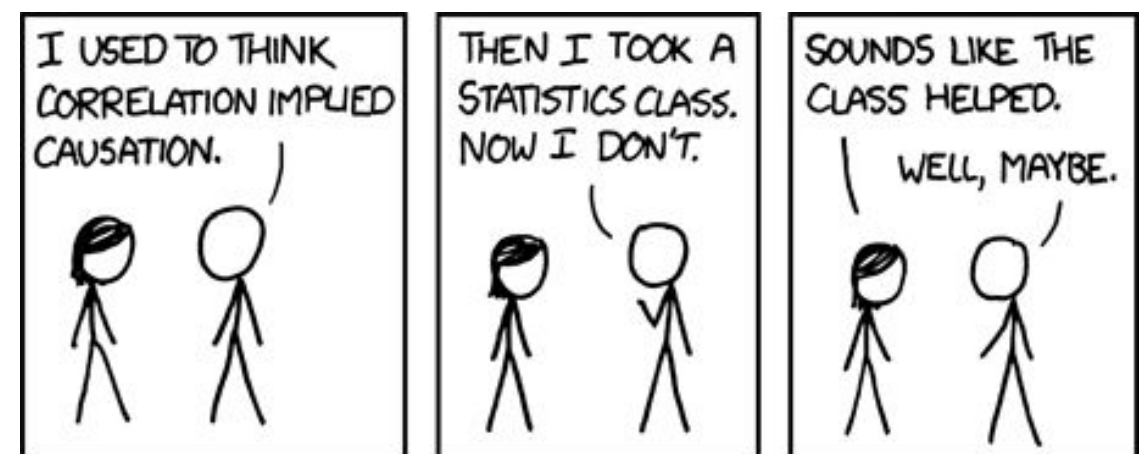


# Python for Data Science

## Python Core

# Hi !

- I am **Koby**
- **Data Scientist** @ Equancy (Marketing Consultancy)
  - Previously: Data Scientist at GDF SUEZ (Engie)
- Organizer of the **Kaggle** Meetup in Paris
- Studied Robotics in Tel Aviv University (2007)
  - Machine Learning in UPMC (2012)
  - Master of Business Administration (2013)



# 3-days plan

1

Pure Python - 1/2d

2

Pandas - 1/2d

3

More Pandas & Matplotlib - 1/2d

4

Scikit-Learn & Machine Learning - 1d

5

Spark - 1/2d

# Outline

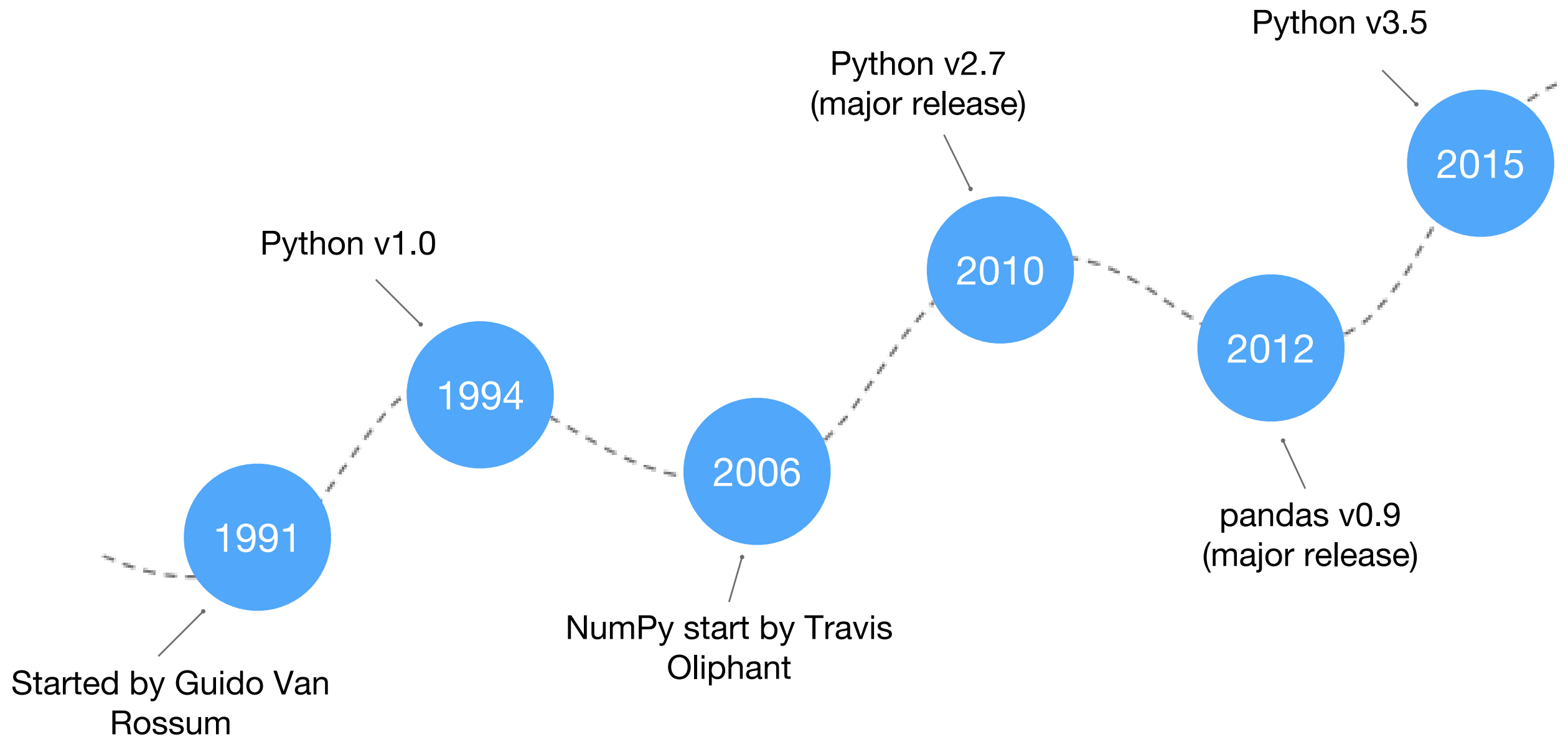
1

Overview

2

Basics & Exercises

# 1 Quick history



# 1 What is Python

- Python is a Programming Language
- Three major version:
  - 1.X (1994) - Historical version
  - 2.X (2000) - Most common version
  - 3.X (2008) - The version we should all be using
- Python is an interpreted version and not compiled:
  - Runs on-the-fly
  - Easy debugging
  - No optimizations be compiler (slower...)

# 1 How to run Python?

## 1 - Python Console

- Enter one line of Python code
- Execute
- Display Output

This is called REPL - Read-Eval-Print Loop

```
karp@K:~$ python
Python 2.7.6 (default, Jun 22 2015, 17:58:13)
[GCC 4.8.2] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> print('Hello')
Hello
>>> exit()
```

# 1 How to run Python?

## 2 - Launch scripts

- Python files have the extension “py”
- The file script.py contains a single line:

script.py

```
print('this is the first line of script.py')
```

- To launch it we simply type “python <script name>”:

```
karp@K:~$ python script.py  
this is the first line of script.py
```

- It's good to use an IDE (integrated development environment) such as PyCharm to make coding more efficient, fast and fun



# 1 How to run Python?

## IPython / Jupyter - interactive computing

Best alternative for analysis and debugging

➤ Command shell: *ipython*

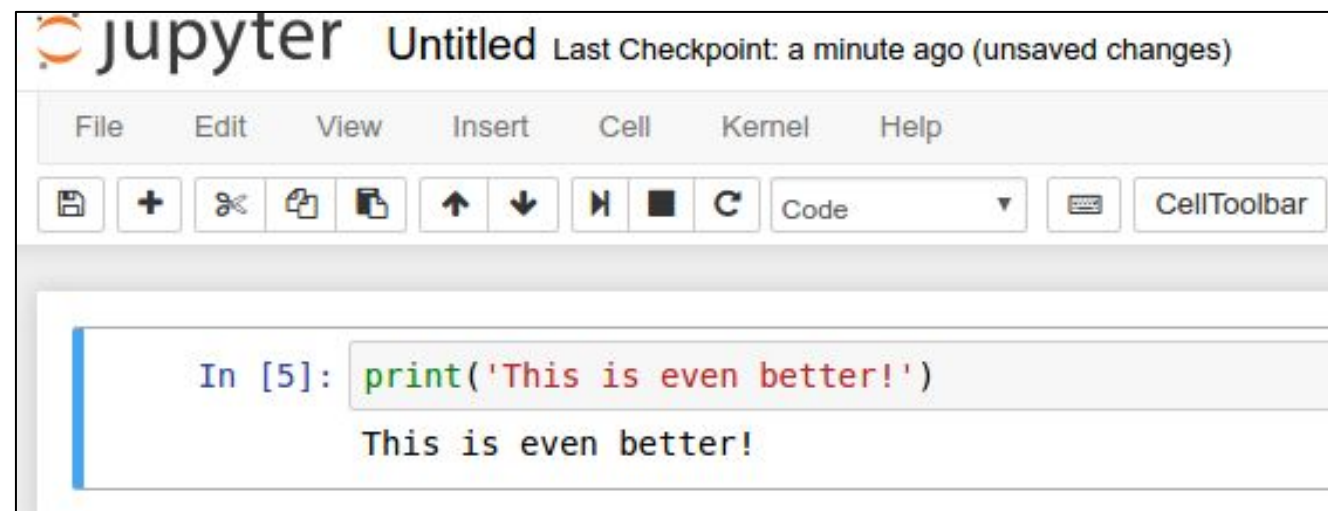
```
(kaggle)karp@K:~$ ipython
Python 3.4.3 (default, Oct 14 2015, 20:28:29)
Type "copyright", "credits" or "license" for more information.

IPython 4.1.2 -- An enhanced Interactive Python.
?                -> Introduction and overview of IPython's features.
%quickref        -> Quick reference.
help             -> Python's own help system.
object?         -> Details about 'object', use 'object??' for extra details.

In [1]: print('This is so much better than the python shell !')
This is so much better than the python shell !

In [2]:
```

➤ Browser-based: *ipython notebook*



# 1 Quick overcap: pros/cons

- Pros:
  - Lots of packages on PyPY for everything
  - Big community
  - Simple and fast learning: productivity gain
  - Easy to script, to read
  - Interpreted

# 1 Quick history: pros/cons

- Cons:
  - Low performance compared to C or Java
    - This is (also) due to dynamical typing

# 1 Let's code!

# 2 Basics: Variables

```
1 # We assign a with the 4 value
2 a = 4
3 # We can reassign a with a new value
4 a = 2
5
6 # So now a has 2 value inside
7 # We can assign 40 to b
8 b = 40
9
10 # Make operations and assign to another variable
11 c = a + b
12
13 # Which value has c
14 # You can output the value with 'print'
15 print(c)
16
17 # Multi assignation
18 d, e, f = 1, 2, 4
19
```

- variable name can't start with a number
- By convention:
  - constants are in CAPS
  - names are > 3 letters
  - spaces around '='
  - use clear names

# 2 Basics: Types

```
1 c = 42
2 # We can get the type of the variables with 'type'
3 type(c)
4
5
6 d = 2.4
7 type(d)
8
9
10 my_string = "Hello world!"
11 type(my_string)
12
13
14 my_list = [1, 1, 2, 3, 5, 8, 13, 21]
15 type(my_list)
16
17
18 my_dictionary = {"barack": 54, "françois": 61, "angela": 61}
19 type(my_dictionary)
20
21
22 my_range = range(10)
23 type(my_range)
24
25
```

- In python types are dynamic: types are set during assignment, we don't need to define them
- Null element is **None**

## 2 Basics: Types operations (int + float)

```
1 a = 4
2 b = -4
3 # Multiply 'a' with 'b' and print
4 print(a * b)
5
6
7
8 # You can get the rest of the division with %
9 23 % 2
10
11
12
13 c = 1.5
14 d = 2.5
15 # Multiply 'c' with 'd' and print
16 print(c * d)
17
18
19 # You can combine 'float' and 'int' to get a new float
20 print(a * c)
21
```

- Here we detailed base operations
- But you can combine variables to build a more complex expression

# 2 Basics: List

```
1 # We can create an empty list, two methods
2 a = []
3 a = list()
4
5 # or not
6 b = [1, 2, 4, 6, 8]
7 c = ["hello", "how", "are", "you?"]
8
9 # use different types in a list
10 d = [1, 1.5, "hello"]
11
12 # The lists are indexed to get items
13 B[0] # first element
14 B[-1] # last element
15 B[::-1] # reverse order
16 B[1:3] # subset starting from 1 (included) to 3 (not-included)
17
18 # some methods are available to manipulate lists
19 a.append(10) # append 10 to the list
20 d.pop() # pop the last item of the list
21
22 # or to operate on the list without modifying the list
23 len(a) # returns the length of the list
24
25 # concatenate strings
26 " ".join(c)
```

- Lists have lot of built-in methods
- We can **iterate** over a list (see after)
- We can **sort** a list
- Some methods:
  - Change (or not) the list
  - Return (or not) a result
- All types can be addable in a list (like a list for instance)
- A list is **always ordered**
- ~~xx~~ Indexing starts at 0



# 2 Basics: Strings

```
1 # An empty string
2 a = ""
3
4 # Or not
5 b = "Koby"
6
7 # Strings support indexing (careful indexing starts at 0)
8 b[0] # returns ?
9 len(b) # returns ?
10
11 # We can format a string, two syntax:
12 c = "Hello my name is %s" % b # Old way
13 d = "Hello my name is {0}".format(b) # New way
14
15 # Of split a string
16 columns = "Age;Name;FirstName".split(";")
17
18 # Make upper
19 b_up = b.upper()
```

- A str has a behaviour similar to a list
  - Indexed (can get item and has a length)
  - Can be iterated
  - Can be sorted

## 2 Basics: More Strings

```
1 # strings are created with " or '  
2 "This is a string."  
3 'This is also a string.'  
4  
5 # joining strings  
6 mylist = ['spam', 'ham', 'eggs']  
7 print(', '.join(mylist))  
8  
9 # replace  
10 print('I like to eat popcorn'.replace('popcorn', 'chocolate'))  
11  
12 # contains  
13 animals = 'cats and dogs'  
14 if 'cats' in animals:  
15     print('we have cats')  
16  
17 # can put special characters in string  
18 print('if we want to have an apostrophe we just need to do \' <- this')  
19 print('this is a line\nand this is another one\t that\'s a tab')
```

# 2 Basics: Dictionaries

```
1 # An empty dict
2 my_dict = {}
3 my_dict = dict()
4
5 # Or with values
6 my_dictionary = {"barack": 54, "françois": 61, "angela": 61}
7
8 # returns an iterator
9 my_dictionary.values()
10 my_dictionary.keys()
11
12 # To combines keys and values
13 my_dictionary.items()
14
15 # To get an element
16 my_dictionary.get("angela") # will throw exception if key doesn't exist
17 my_dictionary["angela"] # will return None if key doesn't exist
18
19 # To update the dict
20 my_dictionary.update({"barack": 28})
21 my_dictionary["barack"] = 28
22
23 # To get the length
24 len(my_dictionary)
```

- You can do almost everything with Python base structures (list + dict)
- A dict **is not ordered**

# 2 Basics: For / While Loops

```
1 # You can iterate over a simple list
2 a = [1, 2, 3, 4, 5]
3 for element in a:
4     print(element * 2)
5
6 # Over an iterator
7 for element in range(10):
8     print(element)
9
10 # Over a string
11 for element in "Christophe":
12     print(element.upper())
13
14 # Over a dict
15 for key, item in my_dictionary.items():
16     print "%s is %s years old" % (key, item)
17
18 # And you can use while (but be careful!)
19 value = 10
20 while value > 0:
21     print(value / 2)
22     value = value - 1
```

- You can loop over lists or iterators
- element is a variable usable only in the 'for'
- In python you must **indent** your code. The indentation must be consistent (e.g. 4 spaces or tabs)

# 2 Basics: If statements

```
1 # Simple if with a simple condition
2 a = 10
3 if a > 20:
4     print("Yes!")
5
6 # An if with an else, if 'if' is false then else is executed
7 b = 34
8 if a < 20 and a < b:
9     pass
10 else:
11     print('Bouh :(!)')
12
13 # Same as previous but we test two conditions
14 if a == 1:
15     a = 4
16 elif a is not None:
17     print('Cool')
18 else:
19     print('Sad!')
20
21 # If can be used for evaluation
22 b = 'Great!' if a == 1 else 'Not Great!'
```

- If statements are based on **boolean** value
- It's only logic and you can combine everything to make **logical** expression:
  - or, and, in
  - not, is
  - >, <, <=, >=, ==
- Python syntax is verbose and simple
- You can use parenthesis to factorise expressions
- **True** and **False** are capitalised

## 2 Exercise 1

If we list all the natural numbers below 10 that are multiples of 3 or 5, we get 3, 5, 6 and 9. The sum of these multiples is 23.

**Find the sum of all the multiples of 3 or 5 below 1000.**

**Answer: 233168**

# 2 Basics: More Lists

```
1 ITEMS = [1, 2, 3, 4, 5]
2 ITEMS_BIS = [6, 7, 8, 9, 10]
3
4 # You can multiply a list if you need to repeat
5 print(ITEMS * 2)
6 NEW_ITEMS = ITEMS + ITEMS_BIS
7
8 # You can in one instruction iterate over the list
9 # It called list comprehension
10 NEW_ITEMS = [i * 2 for i in ITEMS]
11
12 # More complex list comprehension
13 NEW_ITEMS_EVEN = [i * 2 for i in ITEMS + ITEMS_BIS if i % 2 == 0]
14
15 # It's the same as:
16 NEW_ITEMS = []
17 for i in ITEMS:
18     NEW_ITEMS.append(i * 2)
19
20 # Some operations are useful when we use lists like map and filter, sorted
21 NEW_ITEMS = map(lambda x: x * 2, ITEMS)
22 FILTERED_ITEMS = filter(lambda x: x % 2 == 0, ITEMS)
23 SORTED_ITEMS = sorted(ITEMS, reverse=True)
```

- Lists are the most useful structure in Python for data analyses
- You need to master the list operations
- Index starts at 0 (reminder!)

# 2 Basics: Functions

```
1 # You first function
2 def greetings(name):
3     sentence = "Hello %s" % name
4     print(sentence)
5
6 greetings("Koby")
7 greetings("Jacques")
8 # /\ Variable greetings not accessible from this part of the code
9
10
11 # Functions can return a value
12 def add(a, b):
13     """ This method returns a sum between two variables a and b
14     :param a: first param to add
15     :param b: second param to add
16     :return: the sum between a and b
17     """
18     return a + b
19
20 c = add(3, 5)
21 print(c)
22
23 # We can name our function variables
24 def identity(name, age):
25     print(name, age)
26
27 identity(name="Roméo", age=31)
28 identity(age=34, name="Juliette")
29
```

- Functions can return a value (if not it returns **None**)
- We can say also a **method** or a **procedure**
- By convention we give explicit names to functions and we comment a lot the code (e.g. with docstrings)



## 2 Exercise 2

Write a function:

1. It takes a first name as variable
2. It will print the length of the name
3. It will print the first letter of the name only if name's length is even (2, 4, 6, ...)
4. It will return the name in reverse order

## 2 Exercise 3

Write a function

1. It receives 3 arguments: `func(digit1, digit2, operation)`
2. `digit1` and `digit2` are strings representing digits
3. Operation is a string:
  - a. `add`
  - b. `subtract`
  - c. `multiply`
  - d. `divide`
4. The function will return the answer for the operation

## 2 Basics: Main block

```
1 # You first function
2 def greetings(name):
3     sentence = "Hello %s" % name
4     print(sentence)
5
6 # Main block
7 if __name__ == "__main__":
8     greetings("Koby")
9     greetings("Jacques")
10    greetings("Emmanuelle")
11
```

- We use the “main” block to separate the code execution and the import call
- i.e. if we import a module all code except the main will be executed
- If we call the file himself the main will be executed

## 2 Basics: Open and read files

```
1 # Method 1
2 with open("file.csv", "r") as my_file:
3     data = my_file.read()
4     rows = data.split('\n')
5
6 print(rows)
7
8 # Method 2
9 f = open("file.csv", "r")
10 data = f.read()
11 rows = data.split('\n')
12 f.close() # don't forget to close it
13
14 print(rows)
```

- There are two methods to **open** files
- In data analyses you often use files, so this snippet on code is very important
- `open()` second parameter is the open mode: here we read the file so "r"

## 2 Basics: Open and write files

```
1 data = [  
2     ["Name", "Gender", "Age"],  
3     ["A", "Male", 5],  
4     ["B", "Male", 10],  
5     ["C", "Female", 20],  
6     ["D", "Female", 30],  
7 ]  
8  
9 with open("file.txt", "w") as f:  
10     for item in data:  
11         f.write("%s\n" % ";".join(item))  
12
```

- The file will be created in the path given
- Here we write the file so “w” for the open mode

# 2 Basics: Imports

```
1 # You can import a simple package
2 import datetime
3
4 # Import with an alias
5 import pandas as pd
6
7 # Import a specific method or module in a package
8 from sys import path
9 from os.path import splitext
10
11 # Or all but it's not advised because everything will be in you code
12 from sys import *
13
14 # So for usage
15 today = datetime.datetime.now()
16 dataframe = pd.DataFrame()
17 name, ext = splitext("path")
18
19
```

- We you want to develop you will always have to use external packages **imports** are the key
- Hint: order your imports at the top of the file alphabetically

## 2 Basics: datetime

```
1 import datetime
2
3 # We can have the today datetime
4 today = datetime.datetime.now()
5 today_date = today.date()
6
7 # We can parse datetimes with a given format
8 date = datetime.datetime.strptime("2015-01-01", "%Y-%m-%d")
9
10 # And we can format datetime
11 date.strftime("%Y-%m")
12
13 # We can also subtract or add days to a given datetime
14 tomorrow = today + datetime.timedelta(days=1)
15 yesterday = today - datetime.timedelta(days=1)
16
```

- Date times formats are described in the official docs

## 2 Basics: Other basic types

```
1 # tuples
2 point1 = (1, 2)
3 point2 = (2, 5)
4 point1 + point2 # behaves as a list
5
6 point1[0] = 3
7
8
9 # sets
10 items = {1, 2, 'three'}
11
12 len(items)
13
14 set([1, 2, 3, 3])
15
```



# 2 Example

```
1 # My CSV:
2 # Name;Age;Gender
3 # Max;23;Male
4 # Lou;29;Female
5 # Paul;67;Male
6 # Marion;12;Female
7
8 # Open the file with 'with' syntax
9 with open('people.csv', 'r') as f:
10     data = f.read()
11     rows = data.split('\n')
12
13 full_data = []
14 count_row = 0
15 count_columns = 0
16 for row in rows:
17     count_row += 1
18     full_data.append(row.split(';'))
19
20 count_columns = len(full_data[0])
```

1. Open the CSV
2. Create a list of list of data
3. Count with a for #columns
4. Count with a for #rows

## 2 Exercise 4

Make a program that creates the file *today.txt* and writes inside today's date

# 2 Basics: Exceptions

ZeroDivisionError

TypeError

SyntaxError

IndexError

- An exception is an issue **raise** by the code
- You can choose to **catch** exceptions and so write limit cases of your code
- If the exception is not caught somewhere the code will fail and you will get a **traceback**

# 3 Basics: Exceptions

```
1 # A try except block is to catch exceptions
2 # In other languages we call him try/catch block sometimes
3 try:
4     data = [1, 2, 3]
5     last = data[4]
6 except IndexError as er:
7     print("Yes! We caught you: %s" % er)
8 except Exception as er:
9     print("Unexpected error happened")
10 finally:
11     # This block code is always executed at the end of try except
12     pass
13
```

- An exception is an issue **raise** by the code
- You can choose to **catch** exceptions and so write limit cases of your code
- If the exception is not caught somewhere the code will fail and you will get a **traceback**
- **Pass** is necessary if nothing appears in the close (applicable in loops and functions too)

# 2 Basics: PEP8

```
1 # A line can't contains more than 80 characters
2 # At the end of your file you must have a blank line
3 # You have to use all your imports
4 import os
5
6
7 # Here we have two spaces between import and function
8 def hello(name):
9     a = name # Spaces around '=' (note the two spaces before '#')
10    print("Hello %s" % a) # Here spaces around '%'
11
12    if name == "Koby":
13        print("Oh yeah! We have the same name.")
14
15 hello("Koby")
16 hello(name="Koby") # But here no space around '='
17
18
```

- PEP8 is a convention to write clean code and readable by anyone

## 2 Basics: Misc

- Python driven by the **indentation** in your code
- Don't hesitate to **comment** your code
- Code **slowly** and test at every step your code
- Exceptions and errors in Python are very clear, so **read them please** 🐼

## 2 Exercise 5

A palindromic number reads the same both ways. The largest palindrome made from the product of two 2-digit numbers is  $9009 = 91 \times 99$ .

**Find the largest palindrome made from the product of two 3-digit numbers.**

# A Resources

- Python website: <https://www.python.org/>
- Codecademy: <https://www.codecademy.com/tracks/python>
- Python 3 docs: <https://docs.python.org/3/contents.html>
- Project Euler: <https://projecteuler.net/>
- CodeWars: <https://www.codewars.com/>
- Coding Game: <https://www.codinggame.com/>
- <http://dataquest.io>



# Thanks a lot!

[kkarp@equancy.com](mailto:kkarp@equancy.com)

