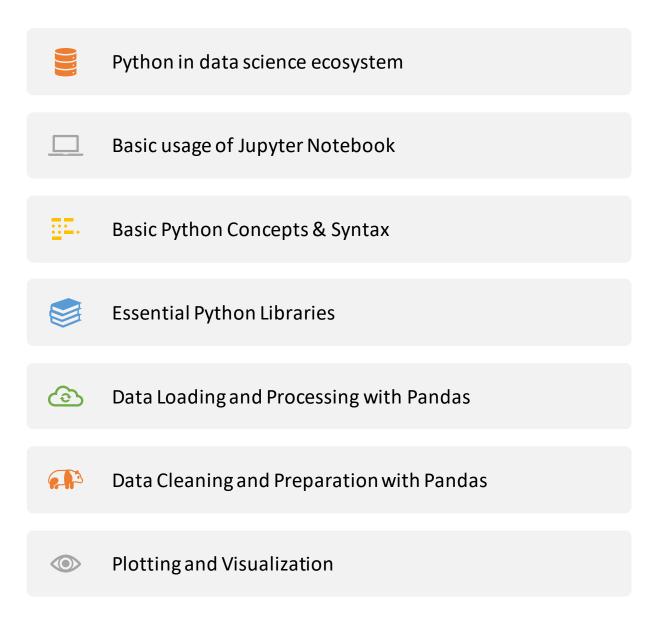


Python Workshop for Data Processing and Visualization
Nov 2020

# Workshop Topics



Python in data science ecosystem

#### **Platform**





#### **Tools**









#### **Frameworks**







#### Language









#### **Format**







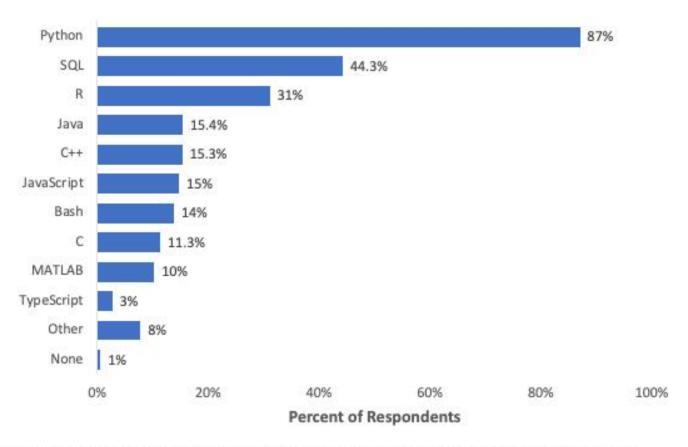
#### SaaS







#### What programming languages do you use on a regular basis?

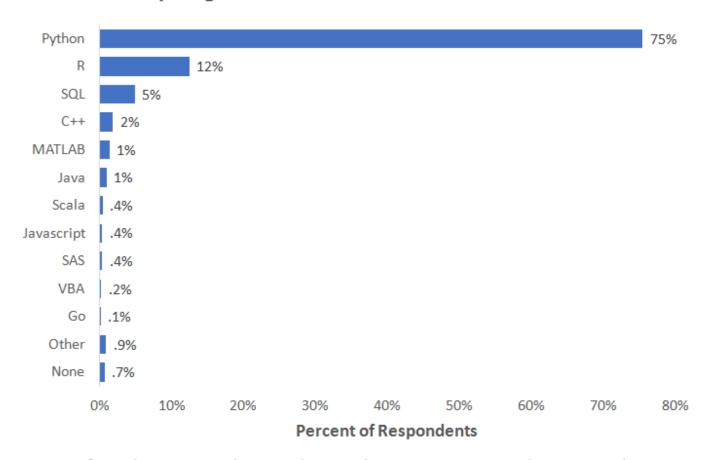


Note: Data are from the 2019 Kaggle ML and Data Science Survey. You can learn more about the study here: https://www.kaggle.com/c/kaggle-survey-2019.

A total of 19717 respondents completed the survey; the percentages in the graph are based on a total of 14762 respondents who provided an answer to this question.



### What programming language would you recommend an aspiring data scientist to learn first?



Note: Data are from the 2018 Kaggle ML and Data Science Survey. You can learn more about the study here: http://www.kaggle.com/kaggle/kaggle-survey-2018.

A total of 23859 respondents completed the survey; the percentages in the graph are based on a total of 18788 respondents who provided an answer to this question.



# Life Cycle of a data analysis project based on CRISP- DM methodology



### **Business Issue Understanding**

- Business Objectives
- Information needed
- Type of analysis
- Scope of work
- Deliverables



#### Data Understanding

- Initial Data collection
- Data requirements
- Data availability
- Data exploration and characteristics



#### Data Preparation

- Gather data from multiple sources
- Cleanse
- Format
- Blend
- Sample



#### Exploratory Analysis/ Modeling

- Develop methodology
- Determine Important variables
- Build model
- Assess model



#### Validation

- Evaluate results
- Review process
- Determine next steps:
- ➢ OK
- ➢ NOT OK



### Visualization & Presentation

- Communicate results
- Determine best method/ graph to present insights based on analysis and audience.
- Craft a compelling story
- Make recommendations

Reference: Problem Solving with advanced analyti

Basic usage of Jupyter Notebook

# Python Editors/IDEs

General Editors and IDEs with Python Support

- Eclipse + PyDev
- Atom
- Vi / Vim
- Visual Studio Code

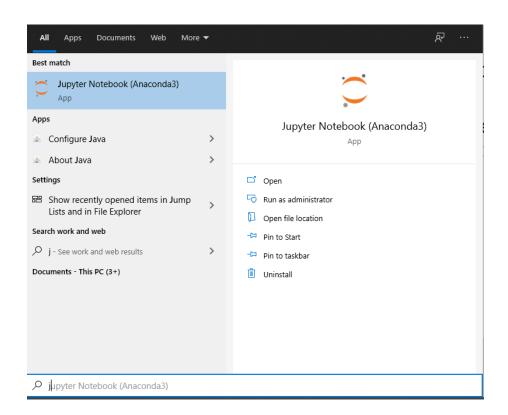
Python-Specific Editors and IDEs

- PyCharm
- Spyder
- IPython/Jupyter

### Starting Jupyter

- .ipynb file is notebook file.
- A kernel is a "computational engine" that executes the code contained in a notebook document.
- A cell is a container for text to be displayed in the notebook or code to be executed by the notebook's kernel.

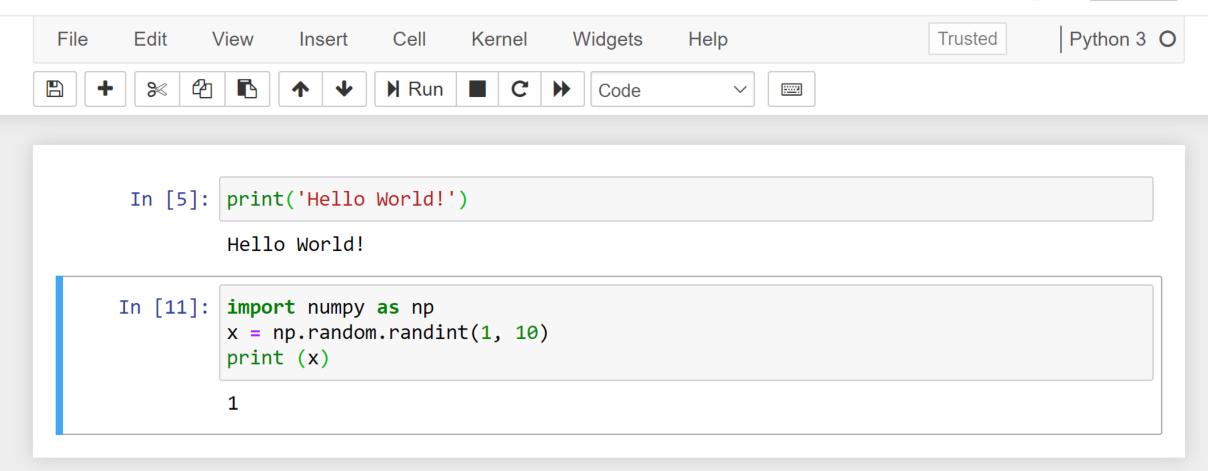




	Upload	New <b>▼</b>	2
Noteboo	k:		2
Python	1 3		d
Other:			10
Text Fi	Text File		10
Folder			
Termin	als Unava	ilable	

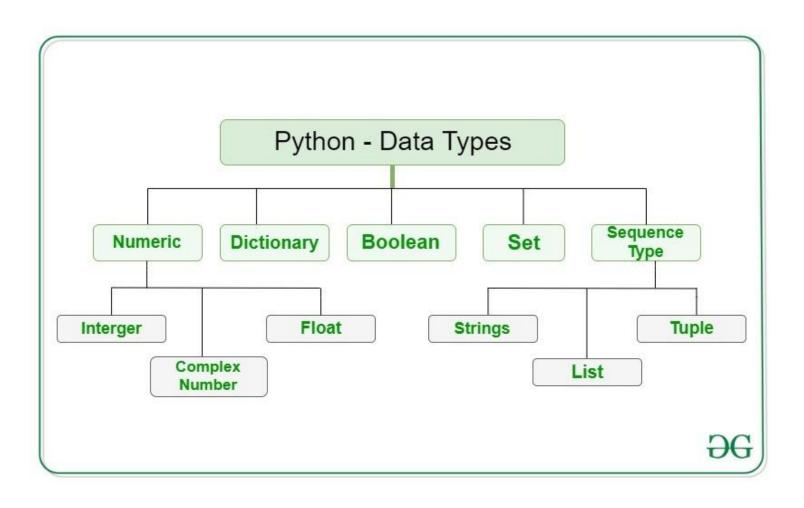


Logout



Basic Python Concepts & Syntax

# Basic Data types



### Int, float, Boolean

```
In []: 1 a = 5 #int
    print(a, "is of type", type(a))
3
4 a = 2.0 #float
    print(a, "is of type", type(a))
6
7 a= True
8 print(a, "is of type", type(a))
9
```

# String

```
In [ ]:
          1 # Strings immutable
          2 s= "Hello World"
            print (type(s))
            print("s[4] = ", s[4])
            print("s[6:11] = ", s[6:11])
            s2 = s.replace("World", "universe")
            print(s)
            print(s2)
         12
         13 # Strings are immutable in Python
            s[5] = 'd'
         14
         15
```

# List, Tuple, Set, Dict

#### List

General purpose

Most widely used data structure
Grow and shrink size as needed
Sequence type
Sortable

#### Tuple

Immutable (can't add/change)
Useful for fixed data
Faster than Lists
Sequence type

#### Set

Store non-duplicate items
Very fast access vs Lists
Math Set ops (union, intersect)
Unordered

#### Dict

Key/Value pairs
Associative array, like Java HashMap
Unordered

# List, Tuple, Set, Dict

Data Structure	Ordered	Mutable	Constructor	Example
List	Yes	Yes	[ ] or list()	[5.7, 4, 'yes', 5.7]
Tuple	Yes	No	( ) or tuple()	(5.7, 4, 'yes', 5.7)
Set	No	Yes	{}* or set()	{5.7, 4, 'yes'}
Dictionary	No	Yes**	{ } or dict()	{'Jun': 75, 'Jul': 89}

### List

```
1 #list (mutable)
In [ ]:
            a = [5,10,15,20,25,30,35,40]
            print("a[2] = ", a[2])
            print("a[0:3] = ", a[0:3])
            print("a[5:] = ", a[5:])
            a[2]= 1500
            print(a)
         11
         12
         13
            a.append(50)
         14 print(a)
         15
         16 a.pop(3)
         17
            print(a)
         18
         19 len(a)
```

# Tuple

### Set

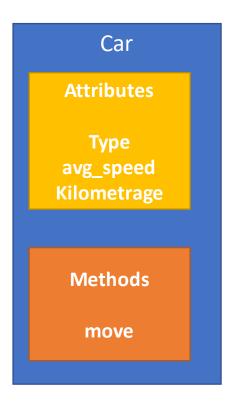
```
1 #Set unordered & unique
2 a = {5,2,3,1,4}
3
4 print("a = ", a)
5
6 print(type(a))
7
8 a.add(4)
9 a.add(5)
10
11 for x in a:
    print(x)
13
14 print (a[1])
```

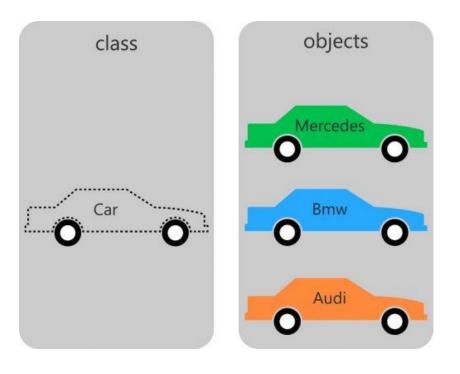
### Dict

```
1 # Dictionary key:value
3 # empty dictionary
4 my_dict = {}
6 # dictionary with integer keys
 7 my_dict = {1: 'apple', 2: 'ball'}
9 # get vs [] for retrieving elements
10 my_dict = {'name': 'Jack', 'age': 26}
11
12 # Output: Jack
13 print(my dict['name'])
14
15 # Output: 26
16 print(my_dict.get('age'))
17
18 # Trying to access keys which doesn't exist throws error
19 # Output None
20 print(my_dict.get('address'))
21
22 # update value
23 my_dict['age'] = 27
24
25 #Output: {'age': 27, 'name': 'Jack'}
26 print(my dict)
27
28 # add item
29 my_dict['address'] = 'Downtown'
30
31 # Output: {'address': 'Downtown', 'age': 27, 'name': 'Jack'}
32 print(my dict)
33
```

### Import package, comment, Indentation

# Class, Object, Function & Method





Function	s. Method
can have many parameters	the object is one of its parameters
exists on its own	belongs to a certain class
function()	object.method()

### Class, Object, Function & Method

```
1 #Class, Object, Function & Method
 2 class Car: # class
       def __init__ (self,car_type,avg_speed): # constructor
           self.type = car type
           self.avg speed = avg_speed
           self.kilometrage = 0
       def move(self, duration_h): # Method
            self.kilometrage += (self.avg speed*duration h)
 9
           return self.kilometrage
10
11
12
       def str (self):
            return self.type + " " + str(self.kilometrage)
13
       def repr (self):
14
           return self.type + " " + str(self.kilometrage)
15
16
17
18 def add(x,y): # Function
       return x+y
19
21 | ford car = Car('Ford',100) # creat object
   print (ford car.type)
24 | ford kilos = ford car.move(100)
25 print (ford kilos)
26
27 honda car = Car('Honda',200)
28 print (honda_car.type)
30 honda kilos = honda car.move(20)
31 print(honda kilos)
33 total_kilos = add(ford_kilos,honda_kilos)
                                                                       24 - 40
34 print(total kilos)
```



### Important Libraries

- Basic Libraries
  - Numpy (mathematical functions)
  - Pandas (data analysis)
  - Scipy (on top of Numpy)
  - Statsmodels (statistical models)
- Visualizations
  - Matplotlib
  - Seaborn (beautiful visualizations on top of Matplotlib)
- Machine Learning
  - Scikit-learn
  - Pytorch (by Facebook)
  - Tensorflow (by Google)
  - Keras (interface for Tensorflow)

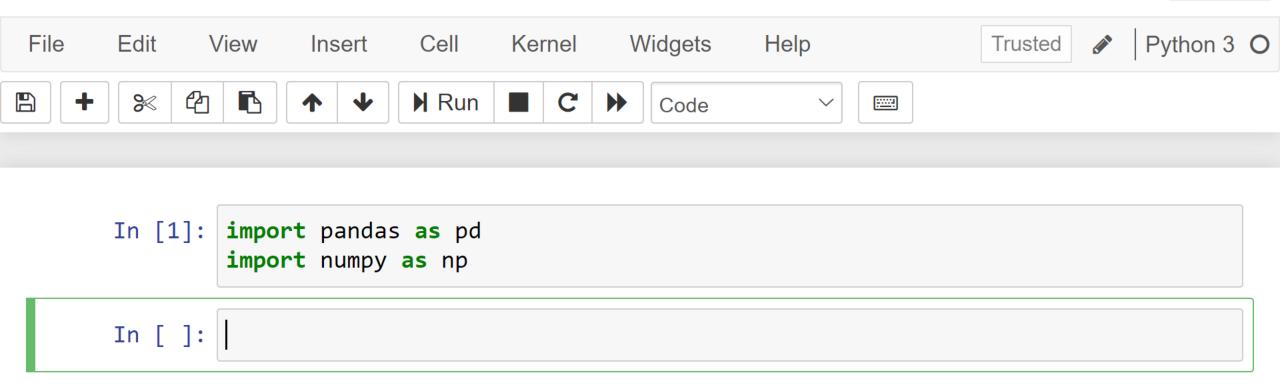
# Installing Packages

- conda install -c anaconda seaborn
- pip install seaborn





Logout



Data Loading, Storage, and File Formats

### Series & DataFrame

### Series

### **Series**

### **DataFrame**

	apples
0	3
1	2
2	0
3	1

	oranges
0	0
1	3
2	7
3	2

	apples	oranges
0	3	0
1	2	3
2	0	7
3	1	2

### Creating DataFrame

# Loading Files and basic DF info

```
In [36]: import pandas as pd
import numpy as np

url_data = 'https://raw.githubusercontent.com/owid/covid-19-data/master/public/data/owid-covid-data.csv'
df = pd.read_csv(url_data, low_memory=False)

print(df.shape)
print(df.columns)
print(df.describe())
print (df.info())
```

### Sub DataFrame based on columns

```
In [51]:

df_sub_col= df[['location','total_cases','total_deaths','date'] ]
    print(df_sub_col.shape)
    print (df_sub_col.head(10))
    print (df_sub_col.tail(10))
```

# DataFrame basic operations

```
In [ ]: print(df[['total_deaths_per_million','hospital_beds_per_thousand', 'human_development_index']].corr())
    print(df['location'].value_counts().head(10))
    print(df['new_cases'].sum()/2)
```



# **DataFrame slicing**

```
fd_eg = df_sub_col.loc[15180]
print (fd_eg)

fd_i = df_sub_col.iloc[15180]
print (fd_i)
fd_i = df_sub_col.iloc[15180:15185]
print (fd_i)
```

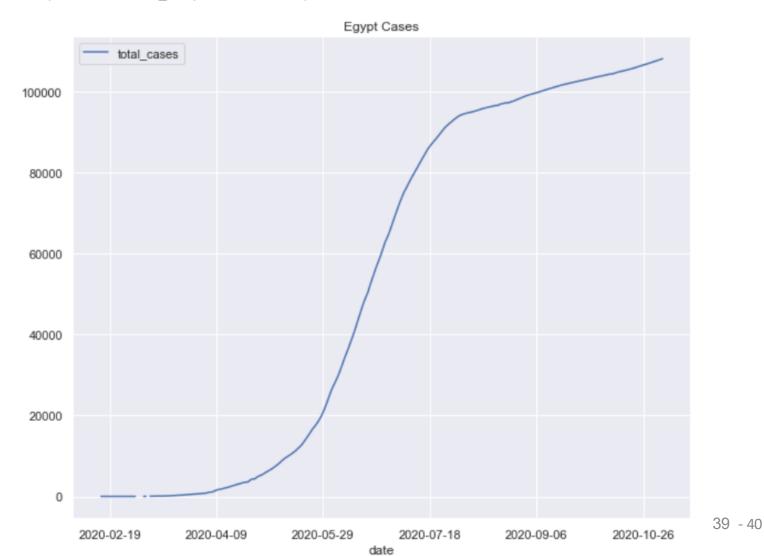
### **DataFrame Conditional selections**

```
condition = (df sub col['date'] == "2020-11-04")
print(condition)
df_last = df_sub_col[df_sub_col['date'] == "2020-11-04"]
print(df last.shape)
print(df_last.head(10))
df_last = df_sub_col[ (df_sub_col['date'] == "2020-11-04" ) & (df_sub_col['total_deaths']>10000)]
print(df last.shape)
print(df_last.head(10))
df_last = df_sub_col[ (df_sub_col['date'] == "2020-11-04" ) & (df_sub_col['location'].isin (['Egypt','Japan']))]
print(df last.shape)
print(df last.head(10))
```



```
In [72]: import matplotlib.pyplot as plt
plt.rcParams.update({'font.size': 20, 'figure.figsize': (10, 8)}) # set font and plot size to be larger
fd_eg = df_sub_col[df_sub_col['location'] == "Egypt"]
fd_eg.plot(kind='line',x='date',y='total_cases', title='Egypt Cases')
fd_eg.plot(kind='line',x='date',y='total_deaths', title='Egypt deaths')
```

Out[72]: <matplotlib.axes.\_subplots.AxesSubplot at 0x2618d8930d0>



```
In [75]: plt.plot( 'date', 'total_cases', data=fd_eg, marker='o', color='skyblue', linewidth=4)
  plt.plot( 'date', 'total_deaths', data=fd_eg, marker='', color='olive', linewidth=2)
  plt.legend()
```

Out[75]: <matplotlib.legend.Legend at 0x2618a6eb820>

