

### Why Python

- Python is an open source, high-level, object-oriented, interpreted, and general-purpose dynamic programming language.
- It has an easy and clean syntax with simple semantics. Similar to C, Algol but it is also an OOP language; OOP is not mandatory (unlike CPP). No extra symbol to start or end a statement (no @, \$ symbol like in perl);
- Python finds its applications across various areas such as website development, mobile apps development, scientific and numeric
- computing, desktop GUI, and complex software development, machine learning and data science
- Development environments such as Ipython Notebook and Spyder that allow for a quick introspection of the data and enable developing of machine learning models interactively.
- No special character at the end (like; in C).
- The main reasons being:
- ✓ no declaration required. It is a high-level language
- ✓ many built-in data structures are already available: dictionary, lists...
- ✓ no need for memory handling: there is a memory garbage collector
- Python can be run, debugged and tested interactively in the python interpreter.
- Ideal for web scripting and data handling.
- There is a rich set of external libraries especially for science (matplotlib, pandas, numpy, Scikit, SciPy etc.)

#### **Drawbacks**

- Indentation probably intimidates beginners but after a couple of hours of coding it becomes natural.
- Slower than compiled languages.

# Python-Setup Installing Python

The most up-to-date and current source code, binaries, documentation, news, etc., is available on the official website of Python: <a href="http://www.python.org">http://www.python.org</a>

#### Windows Installation

- Here are the steps to install Python on Windows machine.
- Open a Web browser and go to http://www.python.org/download/
- Follow the link for the Windows installer *python-XYZ.msi* file where XYZ is the version you need to install.
- To use this installer *python-XYZ.msi*, the Windows system must support Microsoft Installer 2.0. Save the installer file to your local machine and then run it to find out if your machine supports MSI.
- Run the downloaded file. This brings up the Python install wizard, which is really easy to use. Just accept the default settings, wait until the install is finished, and you are done.

- Unix and Linux Installation
- Here are the simple steps to install Python on Unix/Linux machine.
- Open a Web browser and go to http://www.python.org/download/.
- Follow the link to download zipped source code available for Unix/Linux.
- Download and extract files.
- Editing the *Modules/Setup* file if you want to customize some options.
- run ./configure script
- make
- make install
- This installs Python at standard location /usr/local/bin and its libraries at /usr/local/lib/pythonXX where XX is the version of Python.

# Setting Path

To add the Python directory to the path for a particular session in Windows:

**At the command prompt:** type path %path%;C:\Python and press Enter.

**Note:** C:\Python is the path of the Python directory

To add the Python directory to the path for a particular session in Unix:

- In the csh shell: type setenv PATH "\$PATH:/usr/local/bin/python" and press Enter.
- In the bash shell (Linux): type export ATH="\$PATH:/usr/local/bin/python" and press Enter.
- In the sh or ksh shell: type PATH="\$PATH:/usr/local/bin/python" and press Enter.
- Note: /usr/local/bin/python is the path of the Python directory

# Running Python

• Three different ways to start Python:

#### 1. Interactive Interpreter:

Start Python from Unix, DOS, or any other system that provides you a command-line interpreter or shell window.

Enter python the command line.

Start coding right away in the interactive interpreter.

#### 2. Script from the Command-line:

A Python script can be executed at command line by invoking the interpreter on your application, as in the following:

\$python script.py # Unix/Linux or

python% script.py # Unix/Linux or

C:>python script.py # Windows/DOS

**Note:** Be sure the file permission mode allows execution.

3. Integrated Development Environment:

Run Python from a Graphical User Interface (GUI) environment as well, if you have a GUI application on your system that supports Python.

- Unix: IDLE is the very first Unix IDE for Python.
- Windows: PythonWin is the first Windows interface for Python and is an IDE with a GUI.
- Macintosh: The Macintosh version of Python along with the IDLE IDE is available from the main website, downloadable as either MacBinary or BinHex'd files.

# Using Python for Machine Learning

- Python is one of the most popular programming languages for data science and therefore enjoys a large number of useful add-on libraries developed by its great community.
- Being an interpreted high-level programming language, it may seem that Python was designed specifically for the process of trying out different things.
- Although the performance of interpreted languages, such as Python, for computation-intensive tasks is inferior to lower-level programming languages, extension libraries such as *NumPy* and *SciPy* have been developed that build upon lower layer Fortran and C implementations for fast and vectorized operations on multidimensional arrays.
- *scikit-learn* library, is one of the most popular and accessible open source machine learning libraries available in Python.

#### **Installing Python**

Python is available for all three major operating systems—Microsoft Windows, Mac OS X, and Linux—and the installer, as well as the documentation, can be downloaded from the official Python website: <a href="https://www.python.org">https://www.python.org</a>.

#### **Python Variants**

Python is available in 2.x and 3.x variants

A good summary about the differences between Python 3.4 and 2.7 can be found at <a href="https://wiki.python.org/moin/Python2orPython3">https://wiki.python.org/moin/Python2orPython3</a>.

- The Anaconda installer can be downloaded at <a href="http://continuum.io/downloads#py34">http://continuum.io/downloads#py34</a>, or <a href="https://www.anaconda.com/download/">https://www.anaconda.com/download/</a>
- Anaconda quick start-guide is available at <a href="https://store.continuum.io/static/img/Anaconda-Quickstart.pdf">https://store.continuum.io/static/img/Anaconda-Quickstart.pdf</a>.

Python's frequently used libraries

- 1.NumPy----- Used to store and maniupluate data
- 2.Pandas------ Built on top of NumPy and provides higher level data manipulation tools that make working with tabular data more convenient.
- 3. Matplotlib----- To augment the learning experience and visualize quantitative data
- 4.SciPy-----Lots of numerical analysis tools
- 5.networkx----- Playing with graph
- 6.Igraph ----- Playing with graph

Other major libraries include:

Scikit learm, SciPy

#### Note:

Make sure that the version numbers of your installed packages are equal to, or greater than, those version numbers to ensure the code examples run correctly

- NumPy 1.9.1
- SciPy 0.14.0
- scikit-learn 0.15.2
- matplotlib 1.4.0
- pandas 0.15.2

# Basics of Python

#### Identifiers

- Names used to identify a variable, function, class, module or other object.
- Identifiers start with a letter A to Z or an underscore(\_) followed by zero or more letters, underscores and digits(0 to 9)
- Punctuation characters like @,\$ and % cannot be used with identifiers.
- Class names start with an uppercase letter. All other identifiers start with a lowercase letter.
- Identifier with single leading underscore indicates that it is a private identifier.
- Identifier with two leading underscores indicates a strongly private identifier.

# Keywords

<ul><li>And</li></ul>	exec	Not
<ul><li>Assert</li></ul>	finally	or
<ul><li>Break</li></ul>	for	pass
<ul><li>Class</li></ul>	from	print

- Continue global raise
- def if return
- del import try
- elif in while
- else is with
- except lambda yield

• Blocks of code are denoted by line indentation and Python provides no braces to indicate the same.

Ex: if True:

print "True"

else:

print "False"

Note: Lines indented with same number of spaces would form a block

• Statements in Python typically end with a new line. Python does, however, allow the use of the line continuation character (\) to denote that the line should continue. For example:

```
total = item_one + \
    item_two + \
    item_three
```

Statements contained within the [], {}, or () brackets do not need to use the line continuation character.

- Python accepts single ('), double (") and triple ("' or """) quotes to denote string literals, as long as the same type of quote starts and ends the string.
- The triple quotes are used to span the string across multiple lines.

```
word = 'word'
sentence = "This is a sentence."

paragraph = """This is a paragraph. It is
made up of multiple lines and sentences."""
```

- A hash sign(#) indicates a comment.
- '\n' creates a new line
- ';' allows multiple statements on a single line.

# Suites in Python

- A group of individual statements, which make a single code block are called **suites** in Python. Compound or complex statements, such as if, while, def, and class require a header line and a suite.
- Header lines begin the statement (with the keyword) and terminate with a colon (:) and are followed by one or more lines which make up the suite.

```
if expression:
    suite
elif expression:
    suite
else:
    suite
```

# Variable Types

- Python uses '=' sign to assign values to variables
- Variable declaration happens automatically when a value is assigned to it.

```
Ex: a = 10 # integer assignment

b = 5.3 # floating point

c = "Hello Python" # A string

a = b = c = 100 # Multiple Assignment
```

## Working with Python

- Python provides a shell, however, we strongly recommend to install <u>ipython</u> instead, which provides nice features such as tab completion.
- In interactive mode, the prompt is represented by the >>> signs
- Python codes are preceded by the >>> signs and outputs are shown without.
- Important points to remember while working with Python:
- 1. Python is a zero indexing language
- 2. No ';' character at the end of a line, just a return carriage
- 3. Indentation is important and used to separate code blocks
- 4. Everything is an object

#### Data Structures in Python

The built in data structures in Python include:

- Lists,
- Tuples,
- Dictionaries,
- Strings,
- Sets and
- Frozensets.
- ✓ Lists, strings and tuples are ordered sequences of objects.
- ✓ Unlike strings that contain only characters, list and tuples can contain any type of objects.
- ✓ Lists and tuples are like arrays.
- ✓ Tuples like strings are immutables.
- ✓ Lists are mutables so they can be extended or reduced at will.
- ✓ Sets are mutable unordered sequence of unique elements whereas frozensets are immutable sets.

- Most versatile data type available in Python which can be written as a list of comma separated values (items) between square brackets.
- Items in a list need not be of the same type:

Ex: >>>11 =[1,2,3]

To access list elements, use list name or list name followed by index

>>>11 >>>11[0]

>>>12=[1,2,"a"]

>>>12 >>>12[0]

To append elements to a list

>>>11.append(4)

>>>11

[1,2,3,4]

*Update existing value of list* 

Difference between append() and extend()

>>>list1 = ['a','b']

>>>list1.append('c') >>>list1

['a','b','c']

To extend 'list1' to contain more elements, 'extend' method can be used:

>>>list1.extend(['d', 'e', 'f']) >>>list1

['a', 'b', 'c', 'd', 'e', 'f']

Search for an element in the list using index() function

>>>my list = ['a', 'b', 'c', 'd', 'e'] >>>my list.index('b')

>>>my list1 =['a', 'b', 'c', 'd', 'b', 'c']

>>>my list1.index('b')

>>>my list1.index('b',2)

4

11[2]=5

- To insert element wherever into a list, use insert()
- >>>my\_list.insert(2, 'a')
- >>>my\_list

insert() method inserts an object before the index provided

 To remove first occurrence of an element, use remove()

Ex:

If my\_list has elements as a,b,c,d,e,a,b then

- >>>my\_list.remove('a') will result in
- >>>my\_list
  ['b', 'c', 'd', 'e', 'a', 'b']

Ex:

- To remove the last element of the list use pop()
- >>>my\_list.pop()

Removes the last element

- >>>my\_list = ['a', 'b', 'c', 'd']
- >>>my\_list.pop()

To count the number of elements of a kind use count()

Ex:

- >>>my\_list = ['a', 'b', 'c', 'd', 'a', 'b']
- >>>my\_list. count('b')
- To perform in-place sorting use sort()
- >>>my\_list = ['c', 'a', 'd', 'b', 'e']
- >>>my\_list.sort()
- >>>my\_list
- ['a', 'b', 'c', 'd', 'e']
- To sort in reverse order use 'reverse=true'
- >>>my\_list.sort(reverse=True)
  >>>my\_list
- To reverse the list
- >>>my\_list.reverse = ['a', 'b', 'c']
- >>>my\_list.reverse()
- >>>my list

## • Slicing the list

Slicing uses the symbol: to access to part of a list:

- >>>list[first index: last index: step]
- >>> a = [0,1,2,3,4,5]
- >>>a
- [0,1,2,3,4,5]

>>>list[:]

- >>>a[2:] [2,3,4,5]
- >>>a[:2]
- >>>a[:2]
  [0,1]
- >>>a[2:-1]
- [2,3,4]
- >>>a[:]
  [0,1,2,3,4,5]

• List Comprehension

>>>evens = []

>>> for i in range(10):

if i % 2 == 0:

evens.append(i)

- >>>evens
- [0,2,4,6,8]

Rather than the above code

- >>>[i for I in range(10) if i % 2 == 0]
- [0,2,4,6,8]

- In Python, **tuples** are part of the standard language.
- Very similar to the **list** data structure.
- Tuples are faster and consume less memory
- The main difference being that tuple manipulation are faster than list because tuples are immutable
- Tuples are enclosed in parentheses

#### Ex:

$$>>t = (1,2,"a")$$

- >>>l[0]
- Tuples can also be created without parentheses, by using commas:

- >>>1
- (1,2)
- To create a tuple with single element:
- >>>singleton =(1, )
- Tuple can be repeated by multiplying by a number
- >>>(1,) \* 5

(1,1,1,1,1)

Concatenate tuples

$$>>>s1=(1.0)$$

$$>>> s1 += (1,)$$

• Tuples can be indexed and elements in tuple can be counted

• Signing multiple values

• Tuples can be unpacked

$$>>> x, y, z = data$$

- Tuples can be swapped
- >>(x,y)=(1,2)>>(x,y) = (y,x)
- To find length of the tuple:
- >>t=(1,2)
- >>>len(t)
- 2
- Slicing a tuple
- >>t=(1,2,3,4,5)
- >>>t[2:]
- (3,4,5)
- Copying a tuple
- >>t=(1,2,3,4,5)
- >>>newt = t
- >>>newt
- (1,2,3,4,5)

- Mathematical operations
- >>t=(1,2,3)
- >>>max(t)
- 3

- A dictionary is a sequence of items. Each item is a pair made of a key and a value.
- Dictionaries are not sorted. You can access to the list of keys or values independently.
- Items in a dictionary are not sorted
- Dictionaries are built with curly brackets: Ex:

>>>d = {"a":1, "b":2}

>>>d.keys()

['a', 'b']

>>>d.values()
[1,2]

• To query information in dictionary, use

item() method

Ex:

>>>d = {"a":1, "b":2}

>>>d.items() [('a',1),('b',2)] • To check the existence of a specific key use 'has\_key'

>>>d.has\_key('a')
True

To get value corresponding to a key>>d.get('a') #pop can

>>>d.get('a')
also be used

Pop removes the corresponding item from dictionary

To remove an item(i.e. key-value pair) use popitem() method

>>>d1={"a":1,"b":2,"c":3} >>>d1.popitem()

To clear elements in a dictionary use clear() method

>>>d1.clear()

T . 1.1.4. ...1

<ul> <li>To create a dictionary from an existing or empty dictionary:</li> <li>&gt;&gt;&gt;d2.fromkeys(['b', 'c'])</li> <li>Or</li> <li>&gt;&gt;&gt;{}.fromkeys(['a', 'b'])</li> <li>Note:</li> <li>Default values for the keys would be 'None'</li> </ul>	<ul> <li>Strings</li> <li>Strings are immutable sequence of characters.</li> <li>To define strings we can use"</li> <li>1.Single quotes</li> <li>2.Double quotes</li> <li>3.Triple quotes</li> <li>Latest string convention is triple quotes</li> </ul>
• Combining dictionaries  >>> d1 = {'a':1}  >>>d2= {'a':1,'b':2}  >>>d1.update(d2)  >>>d1  >>>d1  >>>d1['a']  >>>d2['a']	Ex:  >>>text ="""Welcome to Python""  >>>text  'Welcome to Python'  Triple quotes allow to specify multi-line strings while double and single quotes allow to insert single quote character.

- To access any character in a string To produce formatted output use '%' use slicing symbol >>>print("%s" % "Welcome") >>>text[0] Will return the first character from the Welcome string at index 0 >>>text[-1] Operators on Strings Will return the first character from the Mathematical operators '+' and '\*' can end of the string or simply returns the be used to create new strings: last character in the string.
- >>>text[0:]
  Will return the complete string starting at index 0
- Characters in given string cannot be changed

>>>text[0]='a'
Will return an error

Ex:

>>>t1 = """Welcome" >>>t2="""to""" >>>t3="""Pyhton""" >>>t4=t1+t2+t3

To get spaces between strings

>>t4=t1+""+t2+""+t3

>>>t4

Methods on Strings Count the occurrence of a character in a string >>>mystr =" This is a string" Few methods are available Now to count the occurrences of the character 'i' to check the type of alpha >>>mystr.count('i') numeric characters present in a string: To find length of a string >>>len(mystr) >>>"123".isdigit() 16 returns True >>>mystr.title() >>>"abc".isalpha()---->>>mystr.capitalize() returns True >>>mystr.upper() >>>"123".isalpha() ----->>>mystr.lower() returns False >>> >>>mystr.center(40) >>>"123abc".isalnum()->>>mystr.ljust(30) returns True >>>mystr.rjust(30,'-') >>>"aa".islower() ---->>>mystr = "This is a string" returns True >>>mystr.find('is') >>>mystr.index('is') >>>"AA".isupper() ---->>>mystr.split() returns True >>>"Aa".istitle()---returns True

hecause string starts with

#### Sets Frozen sets Sets are constructed from a sequence Sets are mutable, and may therefore (or some other iterable object). not be used, for example, as keys in dictionaries. Sets cannot have duplicated, there are usually used to build sequence of Another problem is that unique items (e.g., set of identifiers). themselves only may immutable (hashable) values, and thus >>a= set([1,2,3,4,5]) may not contain other sets. >>b = set([3,4,5,6])Because sets of sets often occur in >>>a | b #Union practice, there is the frozenset type, >>>a & b #Intersection which represents immutable (and, therefore, hashable) sets. >>>a < b #Subset >>>a=frozenset([1,2,3,4]) >>>a – b #Difference >>b=frozenset([2,3,4,5])>>> a ^ b #Symmetric Difference >>>a.union(b) >>>c = a.intersection(b)frozenset([1,2,3,4,5])>>>c >>> a=set([1,2,3])>>>c.issubset(a) >>b=set([2,3,4])>>>c.issuperset(a) >>>a.add(b)

Will return error

>>>a add(frozenset(h))

sets

contain

• Using frozenset as key in a dictionary

>>>fa = {frozenset([1,2]):1}

>>fa[ frozenset([1,2]) ]

- To clear screen
- >>>import os >>>os.system('cls')
- To exit
- >>>exit()
- Or
  Use Ctrl +z and press enter
- >>>^z

# Programming Fundamentals

# **Operators**

```
Arithmetic:
    +, -,*,/, %, **(exponent), //(Floor Division)
Comparison:
    ==,!=,<>,>,<,>=,<=
Assignment:
   = , +=, -=, *=, /=, %=, **=, //=
Logical:
    and, or, not
Membership:
     in, not in
Identity:
    is, is not
```

# Conditional Constructs: Decision Making

- Anticipating conditions occurring during program execution and specifying actions to be taken based on conditions.
- Actions to be executed are determined by evaluating expressions to TRUE or FALSE.
- Python programming language assumes any **non-zero** and **non-null** values as TRUE, and if it is either **zero** or **null**, then it is assumed as FALSE value.

#### If Statement

Contains a logical expression using which data is compared and a decision is made based on the result of the comparison.

#### Syntax:

if expression:

statement(s)

If the expression evaluates to TRUE, then the block of statement(s) inside the if statement is executed.

If expression evaluates to FALSE, then the first set of code after the end of the if statement(s) is executed.

```
Ex: >>>a = input("Enter a number") #reads input from console
>>>a = int(a) # input function returns user's response as a string
which is to be converted to int.
```

```
>>>if (a > 5):

print("a is greater than 5")

print(a)
```

#### If... else statement

- Extension of if statement.
- An **else** statement contains the block of code that executes if the conditional expression in the if statement resolves to 0 or a FALSE value.
- The *else* statement is an optional statement and there could be at most only one **else** statement following **if**.

```
Syntax:
     if expression:
                    statement(s)
     else:
                     statement(s)
Ex: >>> n = input("Enter a number")
        Enter a number 10
>>if n % 2 == 0:
         print (" even number")
     else:
         print("odd number")
```

#### elif statment

- Used to check multiple expressions for TRUE and executes a block of code.
- elif statement is optional.

statement(s)

• Any arbitrary number of elif statements can follow a if statement.

#### Iteration

• The for statement in Python has the ability to iterate over the items of any sequence, such as a list or a string.

#### Syntax:

for iterating\_var in sequence: statements(s)

- If a sequence contains an expression list, it is evaluated first.
- Then, the first item in the sequence is assigned to the iterating variable *iterating var*.
- Next, the statements block is executed.
- Each item in the list is assigned to *iterating\_var*, and the statement(s) block is executed until the entire sequence is exhausted.
- range() function can be used to iterate over a sequence of numbers.
- range() generates an iterator to progress integers starting with 0

```
Ex:1
                               Ex2:
>>>for i in list(range(5)): >>> n=int(input("Enter
Range")
       print(i)
                               Enter Range 5
                               >>> for i in list(range(n)):
Output
                                         if i \% 2 == 0:
                                            print(i)
                                else:
                                          i = i + 1
```

## Loop Control Statements

• The Loop control statements change the execution from its normal sequence. When the execution leaves a scope, all automatic objects that were created in that scope are

destroyed.

Control Statement	Description
break statement	Terminates the loop statement and transfers execution to the statement immediately following the loop.
continue statement	Causes the loop to skip the remainder of its body and immediately retest its condition prior to reiterating.
pass statement	The pass statement in Python is used when a statement is required syntactically but you do not want any command or code to execute.

# Examples

#### **Break statement Continue statement** >>> n = 10>>> n = 10if i<=5: if i<=5: print(i) print(i) break i = i + 1Output block") 0 continue Output 0

#### pass statement

```
>>> n = 10
>>> for i in range(n)
if i == 3:
pass
print("pass
```

# While Loop

• Executes a target statement as long as a given condition is true.

```
Syntax:
       while expression:
         statement(s)
                               Ex:>>> count = 0
Ex: >> count = 0
>>>while (count <10):
                              >>> while count<5:
      print(count)
                                  print(count, "is less than 5")
       count = count + 1
                                  count = count + 1
                               else:
                                   print(count, "is not less
than 5")
```

# Data Analysis using Python

## thon Libraries

Following libraries of Python are useful for scientific computations

	-
Library	Purpose
NumPy	Numerical Python. Contains basic linear algebra functions, Fourier transforms, advanced random number capabilities and tools for integration with other languages.

Scientific Python. It is built on NumPy .Used to perform discrete Fourier

transform, Linear Algebra optimization and Sparse matrices.

Matplotlib Used for plotting variety of graphs like histograms, line plots etc. **Pandas** Used for structured data operations and manipulations and used for data

munging and preparation. Scikit Learn Used for machine learning. Built on NumPy, SciPy and Matplotlib. Contains tools for Machine learning and statistical modeling including

Used for creating interactive plots, dashboards on modern web-browsers

Used to access data from multitude of sources likeMongoDB,

classification, regression, clustering and dimensionality reduction

Used for statistical data visualization

COI Alabamary frame layers about a of data

and large streaming datasets.

Statsmodels Used for statistical modeling to explore data, estimate statistical models

and perform statistical tests.

Seaborn

Bokeh

Blaze

SciPy

Library	Purpose
Scrapy	Used for web crawling for getting specific patterns of data
SymPy	Used for symbolic computation including basic arithmetic to calculus, algebra, discrete mathematics etc.
os	Used for operating system and file operations
networkx and igraph	Used for graph based data manipulations
regular expressions	Used for finding patterns in text data

# Key activities

• Data Exploration:

Finding out more about data.

• Data Munging:

Cleaning the data to make it better suitable for statistical modeling

• Predictive Modeling:

Running the algorithms

## Introduction to Jupyter Notebook

- Jupyter is and IDE used for working with Python.
- Note book environment such as Jupyter are very useful for a variety of purpose.
- Interactive and Scripting support
- Browser based.
- Supports all imports and exports
- Supports variety of data types within the same window such as text, code, graphs, videos, pictures
- Great for Visualizations
- Enables parallel computing

## Importing dataset into Jupyter

- Before we import a dataset into Jupyter notebook:
- 1. Start Jupyter notebook using Windows Start button.
- 2. Click Anaconda.
- 3. Click Jupyter Notebook. Now Jupyter notebook opens.
- 4. Click New button and Click Pyhton2/Python3. A new notebook appears.
- 5. To see the current working directory use the command 'pwd' (print working directory) and use Shift + Enter keys. Now the current working directory is displayed.
- 6. To change the working directory, use jupyter magic command as below:

*Method1*: %cd "C:\Users\pc\Desktop\Data Sets"

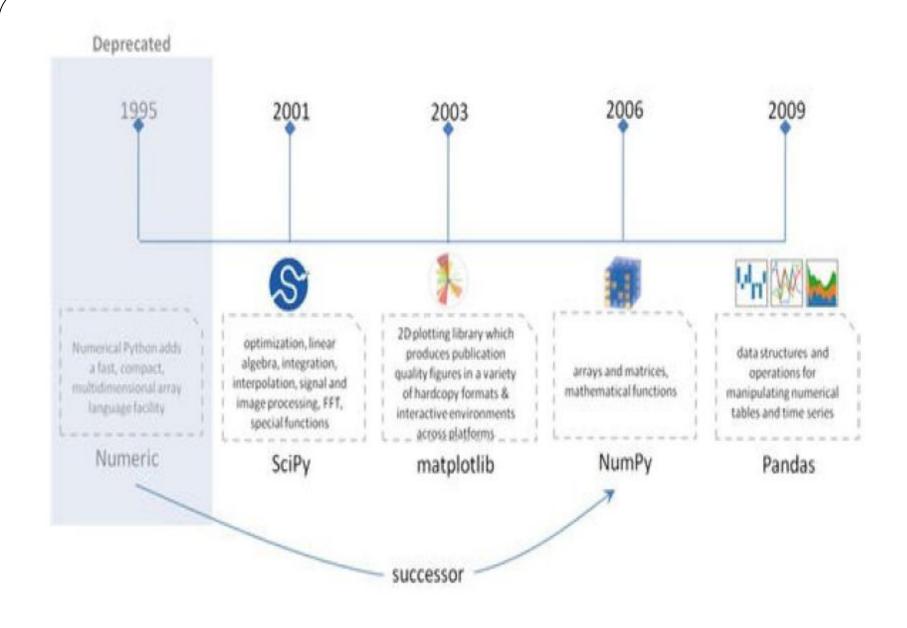
Method2: import os

os.getcwd()

os.chdir("C:/Users/pc/Desktop/Data Sets")

## Machine Learning Python Packages

- There is a rich number of open source libraries available to facilitate practical machine learning. These are mainly known as scientific Python
- At a high level we can divide these libraries into data analysis and core machine learning libraries based on their usage/purpose.
- 1. Data analysis packages: These are the sets of packages that provide us the mathematic and scientific functionalities that are essential to perform data preprocessing and transformation.
- 2. Core Machine learning packages: These are the set of packages that provide us with all the necessary machine learning algorithms and functionalities that can be applied on a given dataset to extract the patterns.
- There are four key packages that are most widely used for data analysis.
- 1. NumPy
- 2. SciPy
- 3. Matplotlib
- 4. Pandas
- Pandas, NumPy, and Matplotlib play a major role and have the scope of usage in almost all data analysis tasks.



## NumPy

- Stands for Numerical Python. It is the core library for scientific computing in Python.
- Numpy is most suitable for performing basic numerical computations such as mean, median, range, etc. Alongside, it also supports the creation of multi-dimensional arrays.
- Provides techniques for loading, storing and manipulating in-memory data in Python.
- Usually real world data is heterogeneous comprising numerical measurements, categorical data, collection of images and may even include audio clips. All varied data can inherently be treated as arrays of numbers.

#### Ex:

An image can be an array of numbers representing pixel brightness across the area. An audio clip can be an array of intensity versus time.

- Efficient storage and manipulation of numerical arrays is fundamental to the process of data science.
- NumPy arrays are like Python's built-in list type and provide much more efficient storage and data operations as the arrays grow larger in size.

#### To work with numpy, import it:

## import numpy as np

;where np is an alias of numpy

To display NumPy's built-in documentation, use:

np?

To create arrays from Python lists:

```
Test the time limits of creating normal list array in python as well as using numpy %%timeit

l1=list(range(5000000))

l2=list(range(5000000))
```

np\_array\_1=np.arange(5000000) np\_array\_2=np.arange(5000000)

To explicitly set the data type of the resulting array, use the 'dtype' keyword:

np.array([1, 2, 3, 4], dtype='float32')

NumPy arrays can explicitly be multi-dimensional which can be initialized using a list of lists.

np.array([range(i, i + 3) for i in [2, 4, 6]])

To index the first element of the array To index the last element To index the first 10 elements To index odd positions

## Finding mean, variance and covaraince

## Finding mean

from numpy import array from numpy import mean

v = array([1,2,3,4,5,6])print(v)

result = mean(v)
print(result)

Finding variance

from numpy import array from numpy import var

v = array([1,2,3,4,5,6])print(v)

result = var(v, ddof=1)
print(result)

from numpy import array

from numpy import cov

x = array([1,2,3,4,5,6,7,8,9])

print(x)
y =

array([9,8,7,6,5,4,3,2,1]) print(y)

Sigma = cov(x,y)[0,1]print(Sigma)

```
from numpy import array
from numpy import cov
x = array([1,2,3,4,5,6,7,8,9])
print(x)
y = array([9,8,7,6,5,4,3,2,1])
print(y)
Sigma = cov(x,y)[0,1]
print(Sigma)
```

## using Pandas

- Pandas are an open source Python package providing fast, flexible, and expressive data structures designed to make working with "relational" or "labeled" data both easy and intuitive
  Handles tabular data sets comprising different variable types (integer,
- float, double, etc.).

  Can also be used to perform leading data or daing feature engineering.
- Can also be used to perform loading data or doing feature engineering on time series data.
- A Python library of data analysis and data manipulation tool built on NumPy
- Fast, intuitive data structure
  - Before data is to be understood, one has to understand two key data structures in Pandas:
    - SeriesDataFrames
  - Series is a one dimensional indexed array. Elements in the series can be accessed using labels.(1-D labeled NumPy array)

## Focus on

- Descriptive statistics, function application
- GroupBy: aggregation and transforming datasets
- Reshaping, pivot tables
- Joining, merging, concatenating
- Time series functionality

# Example for Series and Dataframe

Creating a series by passing a list of values, and a custom index label.

and a custom index label. s = pd.Series([1,2,3,np.nan,5,6],

s = pd.Series([1,2,3,np.nan,5,6], index=

['A','B','C','D','E','F'])
print s

print s

Create a dataframe
data = {'Gender': ['F', 'M',
'M'], 'Emp\_ID': ['E01', 'E02', 'E03'],

'Age': [25, 27, 25]}

df = pd.DataFrame(data,

columns=['Emp\_ID','Gender','Age'])

Reading and Writing Data

To read from a csv file:

df=pd.read\_csv('Data/mtcars.cs

df=pd.read\_csv('Data/mtcars.txt', sep='\t')

To read from Excel:

df=pd.read\_excel('Data/mtcars.xlsx','
Sheet1')

'Sheet2')

To read from a text file:

To read from multiple sheets of Excel into different dataframes

xlsx = pd.ExcelFile('file\_name.xls')
sheet1\_df = pd.read\_excel(xlsx,
'Sheet1')
sheet2\_df = pd.read\_excel(xlsx,

To write to a file df.to\_csv('Data/mtcars\_new.csv', index=False)

index=False)

df.to\_csv('Data/mtcars\_new.txt',
sep='\t',index=False)

## Example

Read the csv file "sales" into jupyter notebook.

#Set current working directory to the location where 'sales' data is available #import pandas library

# import pandas as pd

data1 = pd.read csv("sales.csv")

#View the imported data under 'data1' dataframe data1

#### Note:

1.If the csv file contains any headers(with out values) we can skip the header based on its location in the file

Ex: data1 = read.csv("sales.csv", skiprows=1)

2.If the csv file doesn't contain header then during import, mention the names of columns using the names attribute.

Ex:

data1 = read.csv("sales.csv", header=None, names=["SalesExecutive","Name","Gender","Age","Location","Sales (in Rs)"])

```
To read select rows(say 3 rows) from dataframe:

data1 = pd.read_csv("sales.csv", nrows=3)

data1

To read dataframe containing 'na' (not available) values, use:

data1 = pd.read_csv(sales.csv", na_values =[" na"])

data1  #All the 'na' values will be turned to NaN in the dataframe.
```

Pandas has some built-in functions that enable to get better understanding of data using basic statistical summary methods.

### describe()-

It returns the quick stats such as count, mean, std (standard deviation), min, first quartile, median, third quartile, max on each column of the dataframe

Ex:

#### *cov()* –

Covariance indicates how two variables are related. A positive covariance means the variables are positively related, while a negative covariance means the variables are inversely related.

Ex:

data1.cov()

#### corr()-

Correlation is another way to determine how two variables are related. It tells whether variables are positively or inversely related.

Correlation also tells the degree to which the variables tend to move together.

Correlation ranges between -1 and 1.

data1.corr()

head()-

It displays by default top 5 records. Also it takes user specified number of rows as input.

Ex:

data1.head()

```
tail()-
It displays the bottom 'n' records.
data1.tail()
To extract column names from the dataframe
data1.columns
To extract column data types
data1.dtypes
To get unique column values
data1[column name].unique()
single quotes Ex:data1['Age'].unique()
To view by column name
data1[column name]
To view by row number
data1[0:3]
To display only number of rows
data1.shape[0]
```

#column name must be in

To rename a specific column name df.rename(columns={'old\_columnname':'new\_columnname'},inplace=True)

To rename all column names of DataFrame

df.columns = ['col1\_new\_name', 'col2\_new\_name'....]

To create a new column from existing column

df['new\_column\_name'] = df['existing\_column\_name'] + 5

To create new column from elements of two columns

df['new\_column\_name'] = df['existing\_column1'] + '\_' +

df['existing\_column2']

To add a list or a new column to DataFrame

df['new\_column\_name'] =pd.Series(mylist)

To drop missing rows and columns having missing values

```
To replace missing values with last valid observation df.fillna(method='ffill', inplace=True, limit = 1)
```

Check missing value condition and return Boolean value of true or false for each cell pd.isnull(df)

To replace all missing values for a given column with its mean mean=df['column\_name'].mean(); df['column\_name'].fillna(mean)

To return mean, min, max, sum, count, cumulative sum for each column df.mean()

df.max()
df.min()

df.sum()

df.count()
df.cumsum()

To apply a function along any axis of the DataFrame df.apply(np.cumsum)

## Merge/Join DataFrames

Pandas provide various facilities for easily combining together Series, DataFrame, and Panel objects with various kinds of set logic for the indexes and relational algebra functionality in the case of join merge-type operations.

```
data2 = {\text{"emp id": ["1","2","3","4","5"],"emp name":}}
["Ajay","Vijay","Atul","Sneha","Ravi"]}
```

data2 df =pd.DataFrame(data2, columns =['emp id','emp name'])

```
data3 = {\text{"emp id": ["6","7","8","9","10"],"emp name":}}
["Vimal","Rani","Sirisha","Ashok","Pavan"]}
data3 df =pd.DataFrame(data3, columns = ['emp id', 'emp name'])
```

```
df=pd.concat([data2 df,data3 df])
df
print data2 df.append(data3 df)
pd.concat([data2 df, data2 df], axis=1)
```

dataframes along columns

Pandas offer SQL style merges as well. Left join produces a complete set of records from Table A. with the matching

#Joining two

```
print pd.merge(data2_df, data3_df, on='emp_id', how='inner')
#inner join

print pd.merge(data2_df, data3_df, on='emp_id', how='outer')
#outer join

data1.groupby(['Location','Sales (in Rs)']).max()
#Splitting data into groups using 'groupby'
```

#### Pivot Tables

Pandas provides a function 'pivot\_table' to create MS-Excel spreadsheet style pivot tables. It can take following arguments:

- data: DataFrame object,
- values: column to aggregate,
- index: row labels,
- columns: column labels,
- Aggfunc: aggregation function to be used on values, default is NumPy.mean

Ex:

## Matplotlib

- Matplotlib is a numerical mathematics extension NumPy and a great package to view or present data in a pictorial or graphical format.
- It enables analysts and decision makers to see analytics presented visually, so they can grasp difficult concepts or identify new patterns.
- The most common way for building a pictorial or graphical format is by using global functions to build and display a global figure using matplotlib
- 1. plt.bar creates a bar chart
- 2. plt.scatter makes a scatter plot
- 3. plt.boxplot makes a box and whisker plot
- 4. plt.hist makes a histogram
- 5. plt.plot creates a line plot

#### Ex:

To create a bar chart of test scores of five students import matplotlib.pyplot as plt x=np.arange(5)

y=(35,40,25,70,56)

plt.bar(x,y)

plt.show() commands will not use same figure.

#show()/close() implies that any follow up plot

```
plt.scatter(x,y)
plt.show()
```

- Histogram, line graph, and boxplot can be used directly on a dataframe.
- A simple histogram can be a great first step in understanding a dataset.

#### Ex:

```
%matplotlib inline
Import matplotlib.pyplot as plt
data1.histogram()
data1.plot()
data1.boxplot()
```

#### Customizing a histogram

The hist() function has many options to tune both the calculation and the display:

# Example

Read a csv file of sales data into python using pandas library import os os.getcwd()
 Move the sales.csv file into the current working directory (or)
 To change the current working directory

os.chdir('path of the required directory')
import pandas as pd
data1 =pd.read\_csv("sales.csv")
data1

# Start Exploration

• Start iPython interface in Inline Pylab mode

%pylab inline plot(arange(5))

• Importing libraries and the data set:

Following are the libraries we will use during this project:

- numpy
- matplotlib
  - pandas
- No need to import matplotlib and numpy because of Pylab environment.
- I have still kept them in the code, in case you use the code in a different environment.
- After importing the library, you read the dataset using function read csv().

#### **Understanding numerical variables:**

```
import pandas as pd
import numpy as np
import matplotlib as plt
data1 = pd.read csv("sales.csv")
data1
data1.head()
data1.tail()
data1.info()
                      #Get summary of numerical variables
data1.describe()
```

describe() function would provide count, mean, standard deviation (std), min, quartiles(25%,50%,75%) and max in its output

## Understanding non numerical variables

- The data set also has two non numerical variables i.e. Gender and Location
- Look at frequency distribution to understand whether they make sense or not.

data1.shape #Get no. of Rows and Columns data1.dtypes

data1['Gender'].value\_counts()
data1['Location'].value\_counts()

data1[data1['Gender'].str.contains('Male')] #Get data related to

Males

data1[data1['Gender'].str.contains('Female')] #Get data related to

Females

data1[data1['Gender'].str.contains('Female') & data1['Location'].str.contains('Hyderabad')] #Get data related

# Descriptive Statistics

```
data1['Sales ( in Rs)'].sum() #Get total sale value of company data1['Sales ( in Rs)'].mean() # Get average sale value of company data1['Sales ( in Rs)'].min() # Get minimum sale value data1['Sales ( in Rs)'].max() # Get maximum sale value data1['Sales ( in Rs)'].median() # Get median sale value data1.corr() #To understand correlation significance of numerical variables
```

# Distribution Analysis

- After understanding the basic nature of the data, next step is to understand distribution of variables.
- This can be achieved by plotting the histogram on numerical and non-numerical variables.

For numeric variables Age and Sales

```
data1['Age'].hist(bins = 10)
data1['Sales (in Rs)'].hist(bins=10)
data1.boxplot(column = 'Sales (in Rs)')
data1.boxplot(column = 'Age')
data1.boxplot(column = 'Sales (in Rs)', by = 'Age')
data1.boxplot(column = 'Sales (in Rs)', by = 'Location')
data1.boxplot(column = 'Sales (in Rs)', by = 'Gender')
```

# Interpretation

#### Interpretation:

- ✓ Analysis shows the maximum sale is near to 80(i.e.77) by an age 27.
- ✓ Minimum sale value is 11 at an age 20.
- ✓ Median sale is somewhere near to 40(around 34).
- ✓ Age group 24 is performing average sales for the company.
- ✓ Age group 27 is performing best sales for the company.
- ✓ Age has positive correlation over sales.
- ✓ Median sale value of Bangalore is higher than Mumbai.
- ✓ Least and highest sales happen in Bangalore.
- ✓ Region Hyderabad attributes to increased sale value making it a prospective sales horizon.
- ✓ Males have generated least and highest sales
- Median sale value of Males is higher then Females.
- ✓ Sales distribution appears uniform for all three regions with a sale value ranging from 25 to 75. Hence Bangalore region attributes to outliers in sales with sale values lying below and above the normal sale distribution.