



Data Science using Python

Dr.Vani Vasudevan

Outline



Introduction to Data Science



Data Processing



Data Visualization



Introduction to Machine Learning



Foundation of Neural Networks and Deep Learning

Expected Learning Outcomes

1. Acquire a comprehensive understanding in the fundamental concepts of Data Science Using Python
2. Gain proficiency in utilizing the concepts of Data Processing.
3. Demonstrate and Analyse Data Visualization using Python.
4. Utilize the basic concepts in Supervised and Unsupervised Machine Learning methods.
5. Realize the basics and develop regression & classification applications in Deep Learning.



Unit I

Refer:

1. Chapter 2 , Joel Grus, Data Science from Scratch: First Principles with Python,
Second Edition.
2. Chapter 2,3,5 , Jake VanderPlas ,Python Data Science Handbook: Essential Tools for
Working with Data, Second Edition



Introduction to Data Science

Data	Data Science	Python for Data Science
<ul style="list-style-type: none">• What?	<ul style="list-style-type: none">• What?• Define• Why Python for Data Science	<ul style="list-style-type: none">• Data types• Variables• Basic operations I• Control structures (if, else, loops)• Functions,• Python Libraries for Data Analysis• NumPy• Pandas• Scikit-learn• Introduction to NumPy for numerical computing,• working with arrays and basic operations in NumPy,• Introduction to Pandas for data manipulation and analysis,• Data manipulation and exploration with Pandas,• Introduction to Python Scikit-learn



What is Data?

- Data is a collection of facts
- According to the Oxford “Data is distinct pieces of information, usually formatted in a special way”.

What is Data?

- Data can be measured, collected, reported, and analyzed, whereupon it is often visualized using graphs, images, or other analysis tools.
- Raw data (“unprocessed data”) may be a collection of numbers or characters before it’s been “cleaned” and corrected by researchers.
- Data must be corrected so that we can remove outliers, instruments, or data entry errors.

What is Data?

- Data processing commonly occurs in stages, and therefore the “processed data” from one stage could also be considered the “raw data” of subsequent stages.
- Field data is data that’s collected in an uncontrolled environment.
- Experimental data is the data that is generated within the observation of scientific investigations.

What is Data?

- Data can be generated by:
 - Humans
 - Machines
 - Human-Machine combines.
- Data can often be generated anywhere where any information is generated and stored in structured or unstructured formats.

What is Data Science?

- Data science is an interconnected field that involves the use of statistical and computational methods **to extract insightful information and knowledge from data.**

Define Data Science

- The term “data science” combines two key elements: “data” and “science.”
- **Data:** It refers to the raw information that is collected, stored, and processed
- **Science:** It refers to the systematic study and investigation of phenomena using scientific methods and principles. Science involves forming hypotheses, conducting experiments, analyzing data, and drawing conclusions based on evidence.



Define Data Science

- **Data Science** refers to the scientific study of data.
 - Data Science involves applying scientific methods, statistical techniques, computational tools, and domain expertise to explore, analyze, and extract insights from data. The term emphasizes the rigorous and systematic approach taken to understand and derive value from vast and complex datasets.
 - Data science is about using scientific methods to unlock the potential of data, uncover patterns, make predictions, and drive informed decision-making across various domains and industries.



Why Python for Data Science?

- Python is a popular and versatile programming language, presently has become a popular choice among data scientists for its ease of use, extensive libraries, and flexibility.
- Python programming language provide an efficient and streamlined approach to handing complex data structure and extracts insights.

Python Basics

```
    --> or object to mirror
    mirror_mod.mirror_object = True
    if operation == "MIRROR_X":
        mirror_mod.use_x = True
        mirror_mod.use_y = False
        mirror_mod.use_z = False
    elif operation == "MIRROR_Y":
        mirror_mod.use_x = False
        mirror_mod.use_y = True
        mirror_mod.use_z = False
    elif operation == "MIRROR_Z":
        mirror_mod.use_x = False
        mirror_mod.use_y = False
        mirror_mod.use_z = True

    #selection at the end -add
    mirror_ob.select= 1
    mirrorer_ob.select=1
    context.scene.objects.active = mirrorer_ob
    ("Selected" + str(modifier))
    mirror_ob.select = 0
    bpy.context.selected_objects = []
    data.objects[one.name].select = 1
    int("please select exactly one object")
    print("operator classes")
    --> OPERATOR CLASSES ----
```

```
types.Operator):
    X mirror to the selected
    object.mirror_mirrror_x"
    mirror X"
```

Demos...

- Data types
- Variables
- Basic operations in Python
- Control structures (if, else, loops)
- Functions

<https://www.w3schools.com/python/default.asp>

<https://docs.python.org/3.11/tutorial/index.html>

<https://www.geeksforgeeks.org/python-cheat-sheet/?ref=lbp>



Python For Data Science Cheat Sheet

Python Basics

Learn More Python for Data Science Interactively at www.datacamp.com



Variables and Data Types

Variable Assignment

```
>>> x=5
>>> x
5
```

Calculations With Variables

<code>>>> x+2 7</code>	Sum of two variables
<code>>>> x-2 3</code>	Subtraction of two variables
<code>>>> x*2 10</code>	Multiplication of two variables
<code>>>> x**2 25</code>	Exponentiation of a variable
<code>>>> x%2 1</code>	Remainder of a variable
<code>>>> x/float(2) 2.5</code>	Division of a variable

Types and Type Conversion

<code>str()</code>	'5', '3.45', 'True'	Variables to strings
<code>int()</code>	5, 3, 1	Variables to integers
<code>float()</code>	5.0, 1.0	Variables to floats
<code>bool()</code>	True, True, True	Variables to booleans

Asking For Help

```
>>> help(str)
```

Strings

```
>>> my_string = 'thisStringIsAwesome'
>>> my_string
'thisStringIsAwesome'
```

String Operations

```
>>> my_string * 2
'thisStringIsAwesomethisStringIsAwesome'
>>> my_string + 'Innit'
'thisStringIsAwesomeInnit'
>>> 'm' in my_string
True
```

Lists

Also see NumPy Arrays

```
>>> a = 'is'
>>> b = 'nice'
>>> my_list = ['my', 'list', a, b]
>>> my_list2 = [[4,5,6,7], [3,4,5,6]]
```

Selecting List Elements

Index starts at 0

Subset	Select item at index 1 Select 3rd last item
<code>>>> my_list[1]</code>	<code>Select item at index 1</code>
<code>>>> my_list[-3]</code>	<code>Select 3rd last item</code>
Slice	Select items at index 1 and 2 Select items after index 0
<code>>>> my_list[1:3]</code>	<code>Select items at index 1 and 2</code>
<code>>>> my_list[1:]</code>	<code>Select items after index 0</code>
<code>>>> my_list[:-3]</code>	<code>Select items before index 3</code>
<code>>>> my_list[:]</code>	<code>Copy my_list</code>
Subset Lists of Lists	<code>my_list[list][itemOfList]</code>
<code>>>> my_list2[1][0]</code>	
<code>>>> my_list2[1][1:2]</code>	

List Operations

```
>>> my_list + my_list
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list * 2
['my', 'list', 'is', 'nice', 'my', 'list', 'is', 'nice']
>>> my_list2 > 4
True
```

List Methods

	Get the index of an item Count an item
<code>>>> my_list.index(a)</code>	<code>Count an item</code>
<code>>>> my_list.count(a)</code>	<code>Append an item at a time</code>
<code>>>> my_list.append('!')</code>	<code>Remove an item</code>
<code>>>> my_list.remove('!')</code>	<code>Reverse the list</code>
<code>>>> del(my_list[0:1])</code>	<code>Append an item</code>
<code>>>> my_list.reverse()</code>	<code>Remove an item</code>
<code>>>> my_list.extend('!')</code>	<code>Insert an item</code>
<code>>>> my_list.pop(-1)</code>	<code>Sort the list</code>
<code>>>> my_list.insert(0, '!')</code>	
<code>>>> my_list.sort()</code>	

String Operations

Index starts at 0

```
>>> my_string[3]
>>> my_string[4:9]
```

String Methods

<code>>>> my_string.upper()</code>	String to uppercase
<code>>>> my_string.lower()</code>	String to lowercase
<code>>>> my_string.count('w')</code>	Count String elements
<code>>>> my_string.replace('e', 'i')</code>	Replace String elements
<code>>>> my_string.strip()</code>	Strip whitespaces

Libraries

Import libraries

```
>>> import numpy
>>> import numpy as np
Selective import
>>> from math import pi
```

pandas

Data analysis

lemon

Machine learning

NumPy

Scientific computing

matplotlib

2D plotting

Install Python



Leading open data science platform
powered by Python



Free IDE that is included
with Anaconda



Create and share
documents with live code,
visualizations, text, ...

Numpy Arrays

Also see Lists

```
>>> my_list = [1, 2, 3, 4]
>>> my_array = np.array(my_list)
>>> my_2darray = np.array([[1,2,3],[4,5,6]])
```

Selecting Numpy Array Elements

Index starts at 0

Subset	Select item at index 1
<code>>>> my_array[1]</code>	
<code>>>> my_array[1:2]</code>	

Slice

```
>>> my_array[0:2]
```

```
array([1, 2])
```

Subset 2D Numpy arrays

```
>>> my_2darray[:,0]
```

```
array([1, 4])
```

my_2darray[rows, columns]

NumPy Array Operations

```
>>> my_array > 3
array([False, False, False, True], dtype=bool)
>>> my_array * 2
array([2, 4, 6, 8])
>>> my_array + np.array([5, 6, 7, 8])
array([6, 8, 10, 12])
```

NumPy Array Functions

<code>>>> my_array.shape</code>	Get the dimensions of the array
<code>>>> np.append(other_array)</code>	Append items to an array
<code>>>> np.insert(my_array, 1, 5)</code>	Insert items in an array
<code>>>> np.delete(my_array, [1])</code>	Delete items in an array
<code>>>> np.mean(my_array)</code>	Mean of the array
<code>>>> np.median(my_array)</code>	Median of the array
<code>>>> my_array.correlcoef()</code>	Correlation coefficient
<code>>>> np.std(my_array)</code>	Standard deviation

DataCamp

Learn Python for Data Science Interactively



NITTE
EDUCATION TRUST

Python Libraries for Data Analysis

- Numpy
- Pandas
- Scikit-Learn

Numpy

- For numerical computing
 - Working with arrays
 - Basic operations

https://numpy.org/doc/stable/user/absolute_beginners.html

<https://media.geeksforgeeks.org/wp-content/uploads/20240104182515/NumPy-Cheat-Sheet.pdf>



Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science interactively at www.DataCamp.com



NumPy

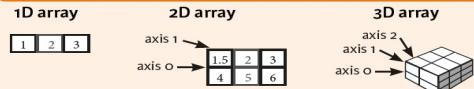
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```



NumPy Arrays



Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1,5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]),
      dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3,4))          Create an array of zeros
>>> np.ones((2,3),dtype=np.int16) Create an array of ones
>>> d = np.arange(10,25,5)
>>> np.linspace(0,2,9)
>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savetxt('array.npy', a, b)
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter="")
```

Data Types

```
>>> np.int64           Signed 64-bit integer types
>>> np.float32          Standard double-precision floating point
>>> np.complex          Complex numbers represented by 128 floats
>>> np.bool              Boolean type storing TRUE and FALSE values
>>> np.object            Python object type
>>> np.string_          Fixed-length string type
>>> np.unicode_          Fixed-length unicode type
```

Inspecting Your Array

```
>>> a.shape           Array dimensions
>>> len(a)            Length of array
>>> b.ndim             Number of array dimensions
>>> e.size             Number of array elements
>>> b.dtype            Data type of array elements
>>> b.dtype.name       Name of data type
>>> b.astype(int)      Convert an array to a different type
```

Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

Array Mathematics

Arithmetic Operations

```
>>> q = a - b
>>> array([[-0.5, 0., 0.],
         [-3., -3., -3., 1.]])
>>> np.subtract(a,b)
>>> b + a
>>> array([[2.5, 4., 6.],
         [5., 7., 9.]]])
>>> np.add(b,a)
>>> a / b
>>> array([[0.66666667, 1.,
         1., 1.],
         [0.25, 0.4, 0.5, 1.]])
>>> np.divide(a,b)
>>> a * b
>>> array([[ 1.5, 4., 9.],
         [ 4., 10., 18.]])
>>> np.multiply(a,b)
>>> np.sqrt(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> e.dot(f)
>>> array([[ 7., 7.],
         [ 7., 7.]])
```

Subtraction

Addition

Division

Multiplication

Exponentiation

Square root

Print sines of an array

Element-wise cosine

Element-wise natural logarithm

Dot product

Comparison

```
>>> a == b
>>> array([[False, True, True],
         [False, False, False]], dtype=bool)
>>> a < 2
>>> array([True, False, False], dtype=bool)
>>> np.array_equal(a, b)
```

Element-wise comparison

Element-wise comparison

Array-wise comparison

Aggregate Functions

```
>>> a.sum()           Array-wise sum
>>> a.min()           Array-wise minimum value
>>> b.max(axis=0)    Maximum value of an array row
>>> b.cumsum(axis=1) Cumulative sum of the elements
>>> a.mean()          Mean
>>> a.median()        Median
>>> a.corrcoef()      Correlation coefficient
>>> np.std(b)          Standard deviation
```

Copying Arrays

```
>>> h = a.view()
>>> np.copy(a)
>>> h = a.copy()
```

Create a view of the array with the same data

Create a copy of the array

Create a deep copy of the array

Sorting Arrays

```
>>> a.sort()
>>> c.sort(axis=0)
```

Sort an array

Sort the elements of an array's axis

Subsetting, Slicing, Indexing

Also see Lists

Subsetting

```
>>> a[2]
>>> 3
>>> b[1,2]
>>> 6,0
```

1	2	3
1	2	3
4	5	6

Select the element at the 2nd index
Select the element at row 1 column 2
(equivalent to b[1][2])

Slicing

```
>>> a[0:2]
>>> array([[1, 2, 3],
         [4, 5, 6]])
>>> b[0:2,1]
>>> array([[ 1.5, 2., 3.],
         [ 4., 5., 6.]])
>>> b[:,1]
>>> array([[1.5, 2., 3.],
         [ 4., 5., 6.]])
```

1	2	3
1.5	2	3
4	5	6

Select items at index 0 and 1
Select items at rows 0 and 1 in column 1

```
>>> b[:,1,:]
>>> array([[ 1.5, 2., 3.],
         [ 4., 5., 6.]])
```

1	2	3
1.5	2	3
4	5	6

Select all items at row 0
(equivalent to b[0,:,:])
Same as [1,:,:]

Boolean Indexing

```
>>> a[a<2]
>>> array([1])
```

1	2	3
1	2	3
4	5	6

Reversed array a
Select elements from a less than 2

Fancy Indexing

```
>>> b[[1, 0, 1, 0, 0], [0, 1, 2, 1, 0]]
>>> b[[1, 0, 1, 0, 0], [0, 1, 2, 1, 0]]
>>> array([[ 1.5, 2., 3.],
         [ 4., 5., 6.]])
```

1	2	3
1.5	2	3
4	5	6

Select elements (1,0),(0,1),(1,2) and (0,0)
Select a subset of the matrix's rows and columns

Array Manipulation

Transposing Array

```
>>> i = np.transpose(b)
>>> i.T
```

1	2	3
1.5	2	3
4	5	6

Permute array dimensions
Permute array dimensions

Changing Array Shape

```
>>> g.ravel()
```

1	2	3
1.5	2	3
4	5	6

Flatten the array
Reshape, but don't change data

Adding/Removing Elements

```
>>> h.resize((2,6))
>>> np.append(h,g)
>>> np.insert(a, 1, 5)
>>> np.delete(a, [1])
```

1	2	3	4	5	6
1.5	2	3	4	5	6

Return a new array with shape (2,6)
Append items to an array
Insert items in an array
Delete items from an array

Combining Arrays

```
>>> np.concatenate((a,d),axis=0)
>>> array([[ 1.,  2.,  3.,  4.,  5.,  6.],
         [ 7.,  8.,  9., 10., 11., 12.]])
>>> np.vstack((a,b))
>>> array([[ 1.,  2.,  3.],
         [ 4.,  5.,  6.],
         [ 1.5, 2., 3.],
         [ 4., 5., 6.]])
>>> np.hstack([e,f])
>>> array([[ 7.,  7.,  1.,  0.],
         [ 7.,  7.,  1.,  0.],
         [ 7.,  7.,  0.,  1.]])
>>> np.column_stack((a,d))
>>> array([[ 1.,  2.,  3.,  4.,  5.,  6.],
         [ 7.,  8.,  9., 10., 11., 12.]])
>>> np.c_[a,d]
```

1	2	3	4	5	6
1.5	2	3	4	5	6
4	5	6	7	8	9
7	8	9	10	11	12
1.5	2	3	4	5	6

Concatenate arrays
Stack arrays vertically (row-wise)
Stack arrays vertically (row-wise)
Stack arrays horizontally (column-wise)

Create stacked column-wise arrays
Create stacked column-wise arrays

Splitting Arrays

```
>>> np.hsplit(a,3)
>>> ([array([[1, 2, 3], [4, 5, 6]]), array([[ 1.,  2.,  3.], [ 4.,  5.,  6.]]), array([[ 1.5, 2., 3.], [ 4., 5., 6.]]])]
```

1	2	3
1.5	2	3
4	5	6

Split the array horizontally at the 3rd index
Split the array vertically at the 2nd index

DataCamp

Learn Python for Data Science interactively



NITTE
EDUCATION TRUST

Pandas

- For Data manipulation, exploration and Analysis

https://pandas.pydata.org/docs/getting_started/index.html
<https://www.geeksforgeeks.org/pandas-cheat-sheet/>



Learn Python for Data Science Interactively at www.DataCamp.com



Pandas

The Pandas library is built on NumPy and provides easy-to-use data structures and data analysis tools for the Python programming language.



Use the following import convention:

```
>>> import pandas as pd
```

Pandas Data Structures

Series

A one-dimensional labeled array capable of holding any data type

a	3
b	-5
c	7
d	4

```
>>> s = pd.Series([3, -5, 7, 4], index=['a', 'b', 'c', 'd'])
```

DataFrame

Columns

	Country	Capital	Population
0	Belgium	Brussels	11190846
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

Index

A two-dimensional labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   'Population': [11190846, 1303171035, 207847528]}

>>> df = pd.DataFrame(data,
   columns=['Country', 'Capital', 'Population'])
```

I/O

Read and Write to CSV

Selection

Getting

```
>>> s['b']
-5
```

Get one element

```
>>> df[1]
   Country    Capital  Population
1  India      New Delhi     1303171035
2  Brazil     Brasilia     207847528
```

Get subset of a DataFrame

Also see NumPy Arrays

```
>>> ar.drop('Country', axis=1) | Drop values from columns(axis=1)
```

Sort & Rank

```
>>> df.sort_index()
>>> df.sort_values(by='Country')
>>> df.rank()
```

Sort by labels along an axis
Sort by the values along an axis
Assign ranks to entries

Retrieving Series/DataFrame Information

Basic Information

<code>>>> df.shape</code>	(rows,columns)
<code>>>> df.index</code>	Describe index
<code>>>> df.columns</code>	Describe DataFrame columns
<code>>>> df.info()</code>	Info on DataFrame
<code>>>> df.count()</code>	Number of non-NA values

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iloc[[0], [0]]
'Belgium'
```

By Label

```
>>> df.loc[[0], ['Country']]
'Belgium'
```

By Label/Position

```
>>> df.ix[2]
   Country    Brazil
   Capital   Brasilia
   Population 207847528
```

Select single value by row & column

Select single value by row & column labels

Select single row of subset of rows

Select a single column of subset of columns

Select rows and columns

Boolean Indexing

```
>>> s[~(s > 1)]
>>> s[(s < -1) | (s > 2)]
>>> df[df['Population']>1200000000]
```

Setting

```
>>> s['a'] = 6
```

Series s where value is not >1
s where value is <-1 or >2
Use filter to adjust DataFrame

Set index a of Series s to 6

Applying Functions

```
>>> f = lambda x: x*x2
>>> df.apply(f)
>>> df.applymap(f)
```

Apply function
Apply function element-wise

Data Alignment

Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([-7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
   a    10.0
   b    NaN
   c    5.0
   d    7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

Scikit-learn

- For Machine Learning

<https://scikit-learn.org/stable/>



