# Data analytics in Python Programming
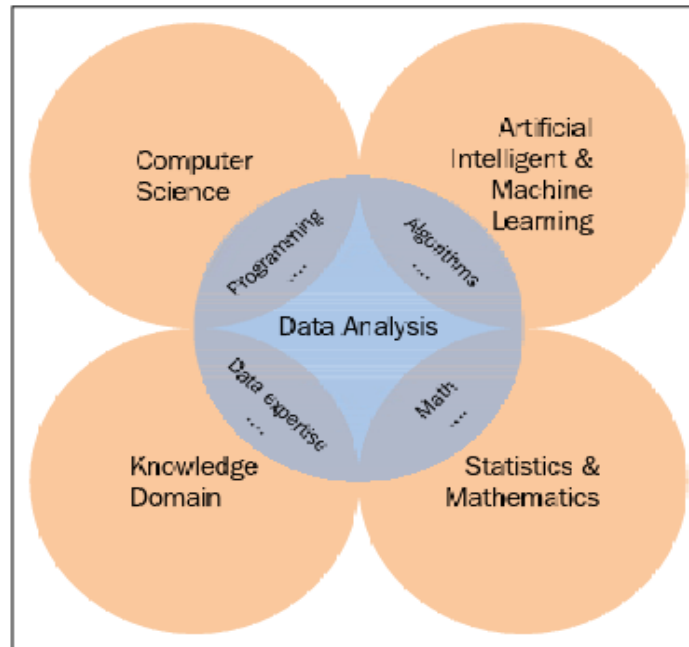
# Week 10:Data Analytics

## Learning objectives

- ❑ Basics of data analytics
- ❑ Data analytics from the programmer's perspectives
- ❑ To know programming styles
- ❑ Data analytics and Python programming

At the conclusion of this lecture, students will be able to understand the fundamentals of data analytics and writing simple data analytics based programs in Python

Ref: **Python : data analytics and visualization : understand, evaluate, visualize data**
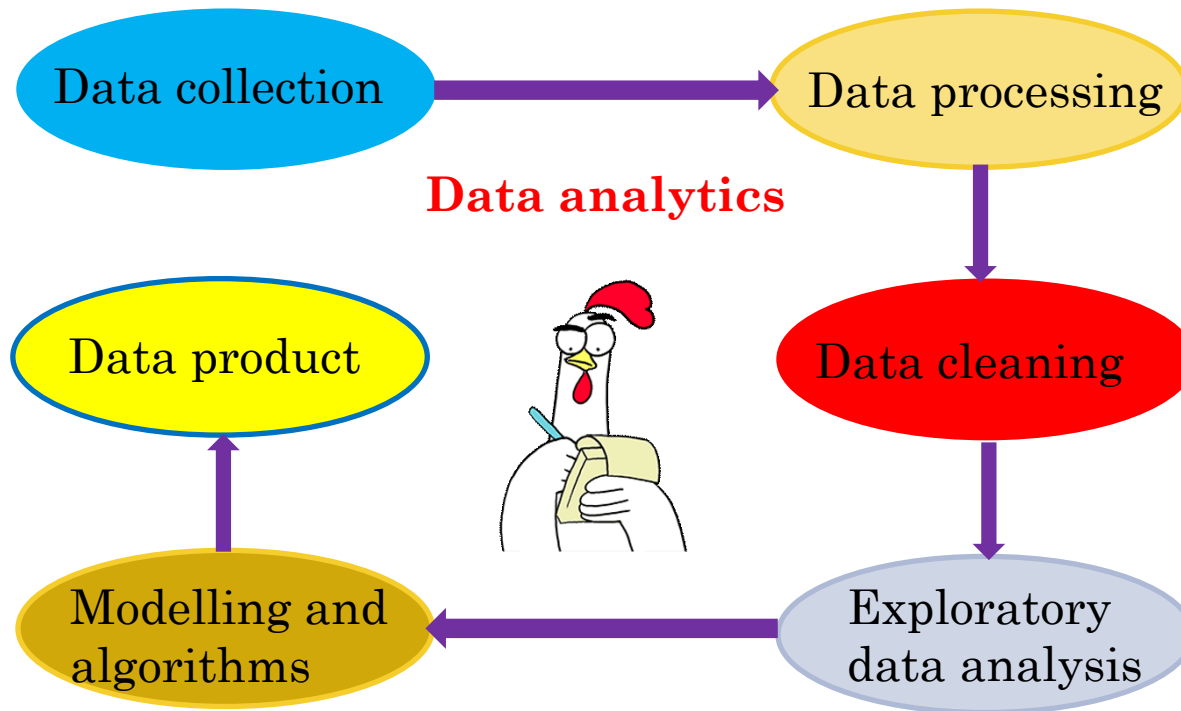Vo. T. H, Phuong, author.; Vo. T. H, Phuong, author.

- **What is data analytics?**

- It refers techniques that applied to analyse and to build data products that receive data input and generate output according to the problem requirements.

- The function of data analytics includes data collection, data processing, data cleaning, exploratory data analysis, modelling and algorithms, and data product.
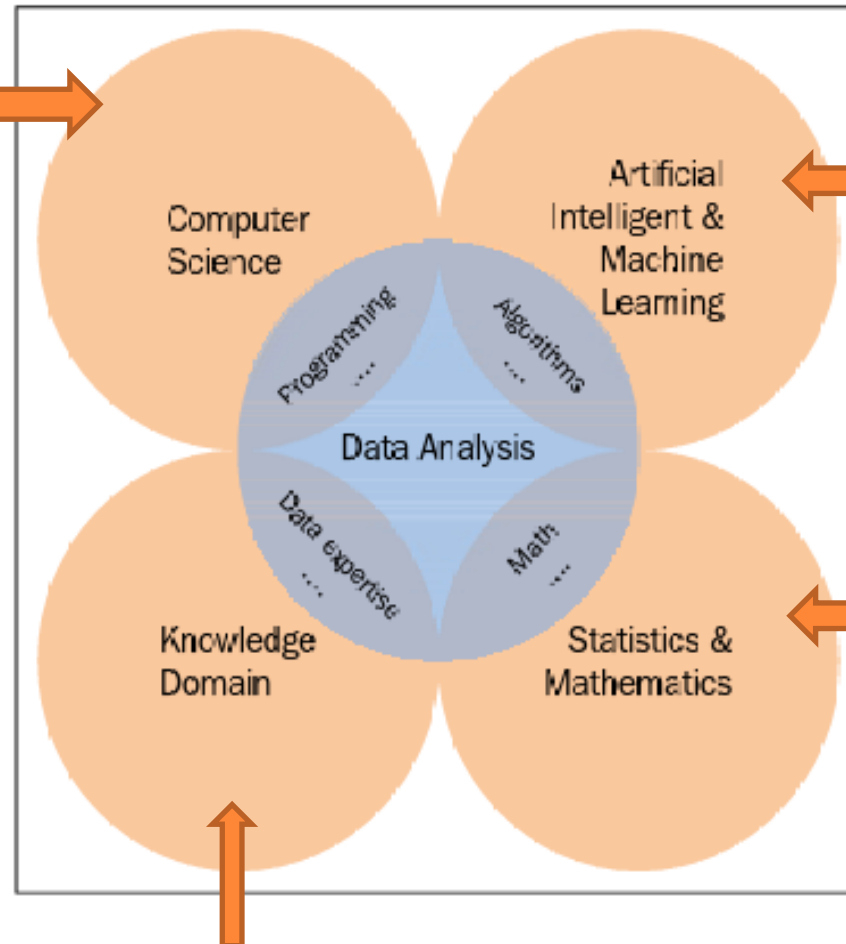
**Data analytics**

```
Data collection ──────► Data processing
                                │
                                ▼
Data product              Data cleaning
     ▲                          │
     │                          ▼
Modelling and    ◄────── Exploratory
algorithms                data analysis
```

# Domain knowledge required for data analysis

**Computer science** knowledge required to develop data analysis tools, in order to provide abstractions for efficient data processing.

Knowing **AI and machine learning algorithms** helps to proceed predictive analytics (Predictive models) in order to provide possible future outcomes
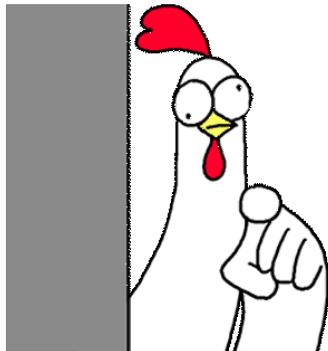


Knowledge in **mathematics and statistics** is required to convert raw data into meaningful information

**Domain knowledge** is required to understand the nature of data
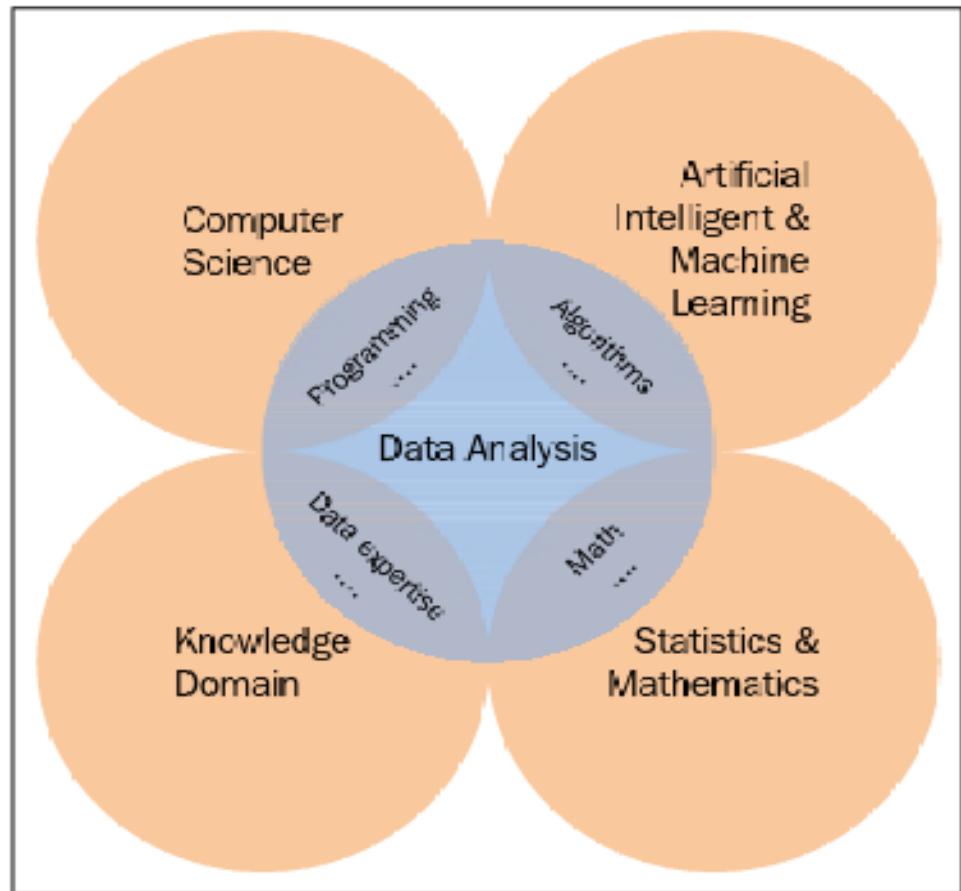Example: learning analytics, business analytics, health data analytics…..

# Data analytics-based jobs and responsibilities

Where do you fit here as a programmer/data analyst/data scientist/ data mining engineer/ data architect/ machine learning engineer?
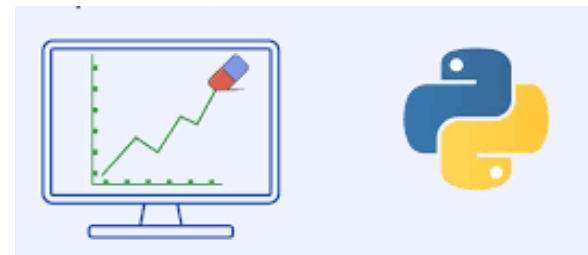
- **Data analytics and Python**
- In general Python and R programming languages are frequently used for data analytics.

- **Python libraries for data analysis**
- NumPy (Numerical Python)
- Pandas (Python data analysis)
- Matplotlib (Plotting library)
- And more…

## Numerical Python (NumPy)

- It contains powerful functions that handle 1-n-dimensional array objects

- In addition, it contain functions that can handle linear algebra operations, Fourier  transformations (decomposition of image), and random number capabilities.

```python
numpy1.py

1  import numpy as np
2  a = [1,2,3]
3  b = [3,4,5]
4  #print(a*b) # will throw an error as type <list>
5  #print([a * b for a, b in zip(a, b)])
6  #simple solution is using numpy array
7  c = np.array([1,2,3]) # single dimensional array
8  d = np.array([3,4,5])
9  print(c*d)
10 print(c+d)
11 print(c/d)
12 print(c-d)
13
```

| | 0 | 1 | 2 |
|---|---|---|---|
| | 1 | 2 | 3 |

```
Shell
Python 3.7.9 (bundled)
>>> %Run numpy1.py

[ 3  8 15]
[4 6 8]
[0.33333333 0.5         0.6        ]
[-2 -2 -2]
```

```python
1  import numpy as np
2  #two X three dimensional array
3  a = np.array([[1,2,3],[3,1,2]]) # one more square bracket[]
4  b = np.array([[3,4,5],[4,5,6]])
5  print(a)
6  print(b)
7  print(a*b)
8  print(a.shape) # display how many rows and columns
9  print(b.shape)
10 print(a.ndim) # display dimension (1D or 2D or..)
11 print(b.size) # number of elements in the array
12 print(a[1,2]) # fetching specific value/element
```

```
Shell
Python 3.7.9 (bundled)
>>> %Run numpy2.py

[[1 2 3]
 [3 1 2]]
[[3 4 5]
 [4 5 6]]
[[ 3  8 15]
 [12  5 12]]
(2, 3)
(2, 3)
2
6
2
```

| | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 1 | 2 | 3 |
| 1 | 3 | 1 | 2 |

# NymPy array floating point values

```
1  #numpy array float values
2  import numpy as np
3  x = np.array([129.45,130.0,173.23,89.23,0.45])
4  print(x)
5
6  new_x = np.delete(x, 4) #deleting a value from specified index
7
8  np.set_printoptions(formatter={'float': "{0:0.1f}".format}) #setting decimal points
9  print(new_x)
10 y = np.array([12.4567,13.230,13.000,891.234,0.456])
11 np.around(y,1)
12 print(y)
13
```

Shell ×

```
Python 3.7.9 (bundled)
>>> %Run numpyfloat.py

 [129.45 130.    173.23  89.23   0.45]
 [129.4 130.0 173.2 89.2]
 [12.5 13.2 13.0 891.2 0.5]

>>>
```

# Reading data from file and converting as NumPy array

```python
 1  import numpy as np
 2  a = np.loadtxt('height.txt',dtype=float)
 3  print(a)
 4  print(np.sort(a)) # sorting the 1D array
 5  print(np.sum(a)) # sum of all values of 1D array
 6  b = np.loadtxt('marks.txt', delimiter=",",dtype=int)
 7  print(b)
 8  print(np.sort(b)) # sorting each row
 9  print(np.sum(b)) # sum of all values
10  print(b)
11  print(np.sum(b,axis=1)) # 1 sum by row
12  print(np.sum(b,axis=0)) # 0 sum by column
```
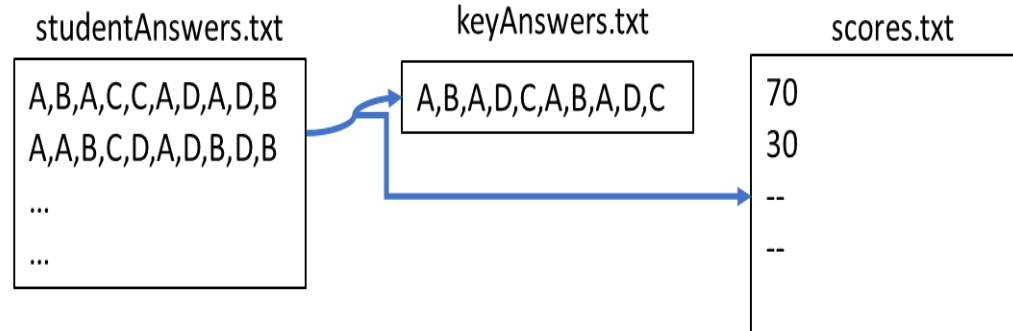
Shell ×

```
>>> %Run numpy4.py
 [166.5 170.   150.   162.2 180.   156.5 134.5 145.6]
 [134.5 145.6 150.   156.5 162.2 166.5 170.   180. ]
 1265.3000000000002
 [[ 45  85  90  20  44]
  [ 64  72 100  30  48]
  [ 89  35  90  13  19]
  [ 72  55  50  82  94]]
 [[ 20  44  45  85  90]
  [ 30  48  64  72 100]
  [ 13  19  35  89  90]
  [ 50  55  72  82  94]]
 1197
 [[ 45  85  90  20  44]
  [ 64  72 100  30  48]
  [ 89  35  90  13  19]
  [ 72  55  50  82  94]]
 [284 314 246 353]
 [270 247 330 145 205]
```

# Reading data and store it as NumPy array

The students' answers for multiple choice exam are stored in a file called *"studentAnswers.txt"*. The key answer for each question is stored in another file called *"keyAnswers.txt"*. <u>Do not modify the contents of these text files</u>. Write code that evaluates student answers by using that key answer text given in another file. Then write the final score computed for each student in a new file called *"scores.txt"*. Then print the final score by invoking the procedure called *"printScores()"* at main program. Check the example given below. Each row represents a student. In total 10 questions were given for the exam.

**studentAnswers.txt**

```
A,B,A,C,C,A,D,A,D,B
A,A,B,C,D,A,D,B,D,B
...

...
```

**keyAnswers.txt**

```
A,B,A,D,C,A,B,A,D,C
```

**scores.txt**

```
70
30
--

--
```

**numpy5.py**

```python
import numpy as np
stAns = np.loadtxt('studentAnswers.txt',delimiter=",",dtype=str)
print(stAns)
keyAns = np.loadtxt('keyAnswers.txt',delimiter=",",dtype=str)
for i in range(np.shape(stAns)[0]): # to move row by row
    marks = 0
    for j in range(np.shape(stAns)[1]): # to move column in each row
        if stAns[i,j] == keyAns[j]:
            marks =  marks+10
    print(marks)
```

**Shell**

```
Python 3.7.9 (bundled)
>>> %Run numpy5.py
 [['A' 'B' 'A' 'B' 'A' 'B' 'B' 'A' 'C' 'C']
  ['A' 'B' 'A' 'D' 'B' 'A' 'C' 'C' 'A' 'A']
  ['A' 'B' 'C' 'D' 'C' 'A' 'C' 'B' 'A' 'D']
  ['A' 'B' 'A' 'A' 'D' 'B' 'B' 'C' 'C' 'D']
  ['B' 'A' 'D' 'B' 'B' 'C' 'D' 'D' 'B' 'C']
  ['A' 'C' 'A' 'A' 'D' 'B' 'C' 'C' 'D' 'D']
  ['A' 'B' 'A' 'D' 'C' 'A' 'C' 'A' 'D' 'C']
  ['B' 'B' 'B' 'A' 'C' 'A' 'B' 'D' 'B' 'A']
  ['D' 'B' 'D' 'B' 'A' 'D' 'D' 'C' 'C' 'A']]
60
50
50
40
10
30
90
40
10
```

## Saving NumPy array as NumPy file(.npy)

```python
import numpy as np
a = np.array([],dtype=int)
x=1
while x!=0:
    x = int(input("Enter any value 0 to terminate:"))
    if x==0:
        break
    else:
        a = np.append(a,x)
print(a)
#writing into a file
np.save("numbers",a) #saving as numpy file in binary format
b = np.load("numbers.npy") #reading and storing as numpy array
print(b)
```

Shell ×

```
Python 3.7.9 (bundled)
>>> %Run numpy6.py
 Enter any value 0 to terminate:12
 Enter any value 0 to terminate:34
 Enter any value 0 to terminate:90
 Enter any value 0 to terminate:-23
 Enter any value 0 to terminate:78
 Enter any value 0 to terminate:0
 [ 12  34  90 -23  78]
 [ 12  34  90 -23  78]

>>> |
```

OK. **What does np.savetxt, np.savez does? Similarly, np.loadtxt means ?**

# LET'S THINK AS DATA ANALYST