

PYTHON FOR DATA SCIENCE – **NOTES**

BY SWOPNIM GHIMIRE

Comprehensive handwritten notes from
my completed Python for Data Science course.



By Swopnima Ghimire

Python for Data Science

Gurteg & Gurjeet

Date _____
Page _____

Focus: Problem Solving

Why python

Python (Zero to Detailed)

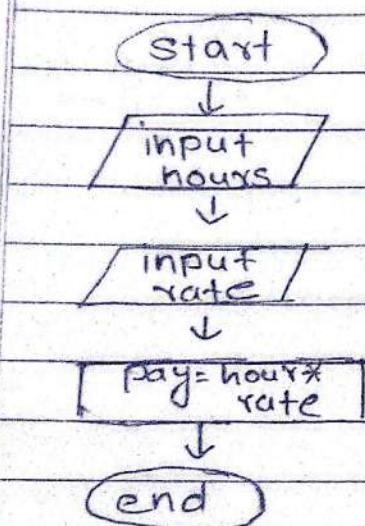
Data science.

Problem solving: finding Generic solution to complex problem.

Algorithm is step by step solution to a problem.

flowcharts

Pseudo Code



BEGIN

input hour

input rate

pay = hour * rate

print pay

END



Pseudocode for making tea

PROGRAM make tea

put teabag in a cup

→ WHILE (Sugar needed)

WHILE (water not boiling)

add sugar

boil water

ENDWHILE

ENDWHILE

Pour water in cup →

Serve and finish

- ⑧ Pseudo code for finding minimum number from a list of numbers.

Eg : list $L = [23, -4, 0, 73, -10, 13]$

- Search Minimum

```
searchMinfromlist(L, n)
```

$\text{minValue} = L[1] \rightarrow \text{assignment}$

$\text{counter} = 2$

while ($\text{counter} \leq n$)

$V = L[\text{counter}]$

if ($V < \text{minValue}$)

$\text{minValue} = V$ & increment
counter

else

pass

endif

increment counter

endWhile

return minValue

endSearchMinfromlist

- ✳ Pseudo code for finding minimum number and the index of that number.

first visualize list $L = [1, 2, 3, 4, 5]$

• Search Minimum and its index.

searchMinfromlist(L,n)

minValue = L[1]

Idx = 1

counter = 2

while (counter <= n)

v = L[counter]

if (v < minValue)

minValue = v

Idx = counter

else

pass

endif

 counter = counter + 1

endwhile

return minValue, idx

end searchMinfromlist

Q) Pseudo code for sorting algorithm, in ascending order.

- Selection sort

```
sortlist (L, n)
```

$L = [] \rightarrow$ This is supposed to be the newlist

counter = 1

while (counter $\leq n$)

 min^{x-value} idx = searchMinfromList (L, n)

 insertInL2 = min

 delete L [idx]

 n = n - 1

endwhile

return L

endSortlist

Here we are
using the
algorithm
which we
executed
previously.

searchMinfromList (L, n)

minValue = L[1]

idx = 1

counter = 2

while (counter $\leq n$)

 v = L[counter]

 if (v < minValue)

 minValue = v

 idx = counter

 else

 pass

 endif if

 counter = counter + 1

endwhile

return minValue, idx

HD Copy endSearchMinfromList .

(*) Converting Pseudo code into python code

- Selection sort

```
def sortlist(L,n):
    L2 = []
    used for defining
    a variable.    counter = 0
    while (counter <= n):
        minValue, idx = searchMinfromList(L,n)
        L2.append(minValue)
        del L[idx]
        n = n - 1
    return L2
```

```
def searchMinfromList(L,n):
    minValue = L[0]
    idx = 0
    counter = 1
    while (counter < minValue):
        v = L[counter]
        if (v < minValue):
            minValue = v
            idx = counter  → add counter too
        else:
            pass
    return minValue, idx
```

⑧ Why chose Python for data Science.

Ans: It's Beginner friendly, versatile & flexible,
Most powerful in Machine learning world
and Most mature libraries around.

IDE (Integrated Development Environment)

Jupyter is the most popular IDE for Data Science.

Download anaconda IDE and enter jupyter notebook to set up working environment.

You can just use ipython in anaconda to set up working environment too.

VARIABLE & OPERATOR

Variable is of different type like int, float and the variable is assigned a value.

They should be written in a descriptive way.

Multiple assignment $a, b = 4, 5.0$

\downarrow \downarrow
int float

Φ Declaring variable

eg :-

$x=3$ (%whos) tells abt all the
print(type(x)) variable in memory.
result: <class 'int'>

Multiple assignment

a, b, c = 3, 4.5, 5

?whois (result: Variable type Data/Info)

a	int	3
b	float	4.5
c	int	5

C = 2 + 4j is complex datatype.

del c is used to delete a variable here it's c.

Arithmetic operators

+ → Addition - = subtraction

/ → division $x \% y = \text{mod}$

* → Multiplication // floor division ie $3.33 // 2 = 1$

** → to the power of $2 ** 3 = 8$

These can be applied to different datatype.

Variable cannot be declared using number and special symbols in the initial position but -a is a proper variable.

Type bool and comparisons.

== True, if it is equal

!= for True, if not equal.

< less than

> greater than

<= less than or equal to

>= greater than or equal to

Bool is a valid data type and it can be either true or false.

If we use and & or gate in the same line without bracket then and will execute first.

e.g.: True or false and False
result: True

find result of print((not(2!=3)&True)or
(False and True))

result: False

Some useful function in python.

- a) round() function rounds the input value to a specified number of places or to the nearest integer.

print(round(4.556)) result: 5

print(round(4.335)) result: 4

print(round(4.556789)) result: 4.56

print(round(4.556789, 3)) result: 4.557

- b) divmod(x,y) output the quotient and the remainder in a tuple
(we will see what a tuple is)

divmod(17,5)

Quotient=5 remainder=2 (5,2) → tuple

print(divmod(34,9)) result (3,7)

`G = divmod(24, 5)`

`print(G) → result (4, 4)`

`type(G) → B tuple` Quotient ↳ Remainder

`G[0] is 24 // 5`

`G[1] is 24 % 5`

- c) `isinstance()` returns True, if the first argument is an instance of that class. Multiple classes can also be checked at once.

eg:

`print(isinstance(1, int))` res: True

`print(isinstance(1.0, int))` res: False

`print(isinstance(1.5, float))` res: True

`print(isinstance(1.67, (int, float)))` res: True

`print(isinstance(2+3j, complex))` res: True

- d) `pow(x,y,z)` x raised to the power y and remainder by z
ie $x^y \% z$

for just power use `x**y`, `pow(x,y)`

eg:

`pow(2, 3, 4)` res: $2^3 \% 4 = 0$

- e) `Input()` `a = input("Enter Something")`
default value of a becomes string.
to convert a to int;

`a = int(a)`

`type(a)` result is integer.
or

`a = int(input("Enter a number"))`

but if user enters something else than no it shows error.

No Copy

- ⑥ If you know the function but aren't aware of the exact syntax then you can use `functionname?` or `help(functionname)`
- eg: `pow?` `help(pow)`
`input?` `help(input)`

CONTROL FLOW (IF-ELIF-ELSE)

`a = int(input())`

`b = int(input())`

Task :-

print only the greater number.

instead of task use if $b > a$

eg:-

`a = int(input())`

`b = int(input())`

`if a > b :`

`print(a)`

`print("a is the greater number")`

`elif b > a :`

`print(b)`

`print("b is the greater no")`

→ use, can use two if or else

`else:`

`print(b)`

`print("b is greater than a")`

eg:-

a = 10

b = 15

if (a == b):

print ("Equal")

elif (a > b):

print ("a is greater")

else

print ("b is greater")

④ A grading system.

a = int (input ("Enter your Marks"))

if (a >= 85):

print ("A Grade")

elif (a < 85) and (a >= 80):

print ("A- Grade")

elif (a < 80) and (a >= 75):

print ("B grade")

else (a < 75) and (a >= 70):

print ("B- grade")

else

print ("Below average")

instead of else

let a = 10

if (a > 9):

print ("A is greater than 9")

elif not (a > 9):

print ("Else part")

(Nested If)

eg:-

```

x = int(input("Enter a number"))
if x > 10:
    print("Your no. is greater than 10")
    if x > 20:
        print("and also above 20")
    else:
        print("but not above 20")
    print("outside the loop")

```

In python Indentation matters very much.

#comment → single line comment
 """ """ → Multi line comment.

User will enter a no 238.915. Your task is to find out the integer portion before the point (in this case 238) and check if it's even no or odd no

```

a = 238.915
a = int(a)
if (a % 2 == 0):
    print("a is even number")
else:
    print("a is odd number")

```

The question specifies just to take the integer portion, so if we use floor then the value will get rounded off so we can't use it.

Control Flow (Loops)

```

int
n = input('Max iterations: ')
i = 1
while (i <= n):
    print(i)
    i = i + 1 (i += 1)
print('done')
    
```

} printing
 natural no
 upto n

- ④ Print all even natural number upto n.

```

n = int(input("Max iteration: "))
i = 1
while (i <= n):
    if (i % 2 == 0):
        print(i)
        print("is even number")
    else:
        pass → used to say program to go on.
        i += 1 (i = i + 1)
    print("done")
    
```

Statements:

break inside a loop exists loop immediately & brings you outside the loop.

continue brings the next iteration loaded regardless of the next statement of the loop.

Eg :

i = 1

while True :

if i % 17 == 0 :

print('break')
break

else :

i += 1

continue

print('I am inside the loop')

print('done')

Current flow (for loop)

eg :- l = [] i value starts from 0 to 9

for i in range(10):
 print(i + 1)
 l.append(i ** 2)

result :-

$l = [0, 1, 4, 9, 16, 25, 36, 49, 64, 81]$
 if we print(l)

range(0, 10, 2) end point is not included.
 start end Jump

eg - for i in range(1, 20, 3):
 print(i + 1)
 l.append(i ** 2)

result

1, 4, 7, 10, 13, 16, 19
 $l = [1, 16, 49, 100, 169, 256, 361]$

Control flow (else in for loops)

eg:

```
S = {"apple", "y.g", "cherry"}
```

```
for x in S:
```

```
    print(x)
```

```
else:
```

```
    print('loop completes its iterations')
```

```
print("outside the loop")
```

result:

apple

y.g

cherry

loop completes with iteration

outside the loop

Case where this else doesn't execute.

i = 1

```
for x in s:
```

```
    print(x)
```

i = i + 1

```
    if (i == 3):
```

break

```
    else:
```

pass

```
else:
```

```
    print("for loop executed successfully")
```

```
print("outside the loop")
```

result:

apple

y.g

outside the loop.

Control flow (Exploring a Dictionary)

D = {"apple": 44, "cherry": "game"}
for x in D:

print (x, D[x])

Result:

apple 44
cherry game

another eg:

D = {"A": 10, "B": -19, "C": "abc"}
for x in D:

print (x, D[x])

Result:

A 10
B -19
C abc

Q Given a list of numbers ie $l = [1, 2, 4, -5, 7, 9, 3, 2]$, make another list that contains all the items in sorted order from min to maximum. ie your result will be $l2 = [-5, 1, 2, 2, 3, 4, 7, 9]$

$$l = [1, 2, 4, -5, 7, 9, 3, 2]$$

for j in range($\text{len}(l)$)

$$m = l[j]$$

$$\text{id}x = j$$

$$c = j$$

for i in range($j, \text{len}(l)$):

if $l[i] < m$:

$$m = l[i]$$

$$\text{id}x = i$$

$$c = i$$

$$\text{temp} = l[i]$$

$$l[i] = m$$

$$l[\text{id}x] = \text{temp}$$

$\text{print}(l)$

Q finding min value. Q for finding index

$$l = [1, 2, 4, -5, 7, 9, 3, 2]$$

$$l = [1, 2, 3]$$

$$m = l[0]$$

$$m = l[0]$$

for i in l :

$$\text{id}x = 0$$

if $i < m$:

$$c = 0$$

$$m = i$$

for i in l :

$\text{print}(m)$

if $i < m$:

$$m = i$$

$$c += 1 - \text{id}x = c$$

$\text{print}(m), \text{id}x)$

Functions

```
def printSuccess():
```

```
    print("The task was successful")
```

```
    print("Moving to the next text")
```

```
    print("Send me the next text")
```

This is the initial phase of defining
a function then we need to call it
to execute the function.

```
, printSuccess()
```

function (Doc string) is one way of
describing the function functionality.

```
def printSuccess2():
```

"""This function is just used to
print the message hellow."""

```
print("hellow")
```

PrintSuccess2?

You'll get details in Docstring

functions (input arguments)

```
def printMessage(msg):
```

```
    print(msg)
```

function is dynamic it can be used to
print anything.

Eg :-

```
def printMsg(msg):
    if isinstance(msg, str):
        print(msg)
    else:
        print("Your msg isn't string")
        print("Your msg datatype is",
              type(msg))
```

Eg :-

```
printMsg("Hello world")
```

Result :-

Hello world

PrintMsg (23)

Result :-

Your msg isn't string

Your msg data type is : <class 'int'>

function (input arguments)

def mypow(a,b):

""" This function computes power just like
builltin power function """

c = a ** b

print(c)

Input :-

my pow(3,4) result :- 81

```

g) def checkArgs(a,b,c):
    if isinstance(a,(int,float)) and
       isinstance(b,(int,float)) and
       isinstance(c,(int,float))
        print((a+b+c)**2)
    else:
        print("Error! Argument
              passed is invalid")

```

input checkArgs(3,4,5) result: 144

input checkArgs(3,4,"g") result: Error

functions (Order of input arguments)

```

def f(c,c1,c3):
    print((c,c1,c3))

```

f(2,3,"game") result: 3,2,"game"

```

def f(a,b,c):
    print("A is: ", a)
    print("B is: ", b)
    print("C is: ", c)

```

Input f(3,2,"bird") result:

A is : 3
 B is : 2
 C is : bird

functions (variables inside)

```
def add(x,y):
```

```
c=x+y # I need value of c?
```

Input # The value of c is available just within the function so we use return

```
def add(x,y):
```

```
return x+y
```

so, if we use d = add(2,3)

```
print(d)
```

it can print the sum but c=x+y cannot print the sum as the value of the variable is just available within a function.

If you return some value out of a function then it returns the value but if you don't pass any value then calling it prints none and its datatype becomes 'NoneType'.

Return can also be used as a break function code below doesn't run:

```
def x():
```

```
a=5 b=7 d="something"
```

```
return a,b,d
```

```
x,y,z=x()
```

```
print(x,y,z)
```

result: 5 7 something

functions (Variable number of input arguments)

```
def add(*args):
    sum = 0
    for i in range(0, len(args)):
        sum += args[i]
    return sum
```

add(3, 4, 5, 6, 7) result: 25

Dictionary

```
def f(**c):
    for n in c:
        print(n, c[n])
```

f(c1="A", c2="B", c3="C")

Eg:

```
def printAllVariableNamesAndValues(*args):
    for n in args:
        print("Variable Name is:", n,
              "And Value is:", args[n])
```

printAllVariableNamesAndValues(a=3, b="B",
 c="CCC",
 y=6.7)

result:

Variable Name is : a And value is: 3

|| || : b || || || : B

|| || : c || || || : CCC

|| || : y || || || : 6.7

function (Default values)

```
def f(sum=0):
    print(sum)
```

f(3) returns 3 f() returns 0

If $l=[1, 2, 3]$ and $l_2=l$ Then in python they don't share different memory location, it's just another variable name for the same memory location.

def f(l=['']) so at the time of every compilation default value is assigned and the default value does not changes.

Eg: def gg(s=4):
 print(s)

input gg() res=4
input gg(56) res=56

Eg: l=[1, 2, 3]
 $l_2=l$

$l_2[0]=-9$
print(l)
result:
[-9, 2, 3]

Eg: def ff(l=[1, 2]):
 for i in l:
 print(i)

$l_2=[12, 3, 4]$
ff() res=1 2
ff(l2) res=12 3 4

Note: Even after multiple manipulation of data the default value remains same.

Modules

Modules is a python file that contains function that you want to use in several coding projects or the functions that you don't want to write.

D:/mymodules/my funcs.py

```
def printMe (msg='No message  
was supplied'):  
    print (msg)
```

```
def printlist (l=[]):  
    for n in l:  
        print (n)
```

```
import sys  
sys.path.append ('D:/mymodules')  
import myfuncs as f  
f.printMe ('hellow')
```

Create a file in jupyter named anything
eg - name be my_universal_functions.

```
def checkType def checkInfoNumeric(**args):  
    output = True  
    for n in args:  
        if isinstance  
            (n,int, float)
```

Download this file as .py file.

```
def checkIfNotNumeric(*args):
    for n in args:
        if not (isinstance(n, (float, int))):
            return False
    return True
```

```
def addAllNumerics(*args):
    s = 0
    for n in args:
        s += n
    return s
```

myName = "Python Course"

After downloading the file create a folder in a drive and paste the downloaded file there e.g. :

creating a folder in win C named ABC and then pasting the file there.

NOW,

in jupyter;

```
import sys
sys.path.append('C:/ABC/')
```

```
import my_universal_functions as myfs
myfs.addAllNumerics(2, 3, 4, 6)
.myfs.checkIfNotNumeric(1, 3, A)
myfs.checkIfNotNumeric(1, 2, 3)
```

result for myfs.myfs.myName

15 false True 'Python Course'

If you have to import a part of the function and give it a new name then we use something like

```
import sys
sys.path.append('C:/ABC/')
```

From my_universal function import addAllNumerics as sum.

```
Sum(3,4,5,6) result: 18
```

Practice function

Q """" Give a list of no ie [1,2,4,-5,7] make another list that contains all the items in sorted order from min to max. ie your result will be another list ie [-5,1,2,4,7] """"

```
def findMin(L):
```

```
m=L[0]
```

```
idx=0
```

```
counter=0
```

```
f for x in L:
```

```
    if n < m if counter < m
```

```
        m=counter
```

```
        idx=
```

```
def findMin():
```

```
m = l[0]
```

```
idx = 0
```

```
i = 0
```

```
for x in l:
```

```
    if x < m:
```

```
        m = x
```

```
        idx = i
```

```
    else:
```

```
        pass
```

```
i += 1
```

```
return m, idx
```

This function defines minimum for a single data ie find the lowest value and return it with its index.
So, everytime same value is returned.

eg:-

```
a, b = findMin([2, 3, 4, 0, 9])
```

```
print(a, b) result 0, 3
```

```
def swapvalues(l, idx1, idx2):
```

```
tmp = l[idx1]
```

```
l[idx1] = l[idx2]
```

```
l[idx2] = tmp
```

```
return l
```

```
l = [2, 3, 6, 7]
```

```
l2 = swapvalues(l, 1, 3)
```

```
print(l2)
```

result :-

```
[2, 7, 6, 3]
```

function to find different min in diff cases.

```
def findMin(l, startIndx):
    m = l[startIndx]
    idk = startIndx
    for i in range(startIndx, len(l)):
        n = l[i]
        if n < m:
            m = n
            idk = i
        else:
            pass
    return m, idk
```

```
def checkIfNotNumeric(l):
    for n in l:
        if not (isinstance(n, (int, float))):
            return False
    return True
```

```
def sortList(l):
    if not (checkIfNotNumeric(l)):
        print("Error: No numeric val")
        return
    else:
        c = 0
        for n in l:
            m, idk = findMin(l, c)
            c += 1
            l[idk] = m
        return l
```

```

Input: l2 = sorted([2, 3, 6, 5, 4, -1])
      print(l2)
      result:
      [-1, 2, 3, 4, 5, 6]
  
```

DATA STRUCTURE IN PYTHON

STRING

"python" & 'python' both are valid string.

```

a = "Hello"
b = "World"
c = a + " " + b
d = a + " messed up" + b
print(c)   res: Hello World
print(d)   res: Hello messed up world.
  
```

This is called concatenation of string.

You can convert a integer to string like this :

```

a = 15
a = str(a)
print(type(a)) res: str
  
```

You can not concat an integer with a string
 You need to convert convert the int into string to concat.

String (Multiline String)

`multilineString = """ This is a very
good course
That I am learning
from """`

`print(multilineString).`

`result: This is a very good course that I am
learning from.`

`result: This is a very
good course`

`That I am learning
from.`

`print("Hey mate !`

`Hi buddy,`

`Hope you're okay")`

`result:`

`Hey mate !`

`Hi buddy`

`Hope you're okay`

(Indexing & Slicing)

`a = "Game of programming"`

`print(a[3:8])` → access from start

`print(a[-8:-3])` → access from end

`result:`

`end of.`

`SLSL`

Starting point
(inclusive) Ending point
(exclusive)

a = "Abra ka Dabra"

print(a[0, 5, 2]) ie 0, 2, 4

result: AYLI

If you don't mention the start [: 5] it starts from initial point and if you don't mention end it starts from there goes till the end.

If you want to totally reverse the string do [:: -1]

functions of string

a = "A lot of spaces at start"

b = a.strip()

print(b.lower())

print(b.upper())

print(a.replace(" ", ", "))

String to list

l = "game, and, no".split(",")

returns list

result = [game, and, no]

[0] = game [1] = and [2] = no

To find if a substring is at a string or not

print ("at" in "andhraatlas")

res: True

print ("lol" in "Lotesh")

res: False

print ("we are learning" "string" "here")

res: we are learning "string" here

print ("We are\n now on newline")

res: We are

now on new line

print ("Hello\n world")

res: Hello world

To print C:\name\drive

print (r "C:\name\drive")

Data Structure

List	ordered	changable	Duplicates
Tuple	ordered	unchangable	No Dup
Set	unordered	addable/removable	No dup
Dictionary	unordered	changable	No dup

Eg :

L = [1, 3, 4, 9, "name", 3]

T = (1, 3, 4, 9, "name", 3)

S = {1, 3, 4, 9, "name", 3}

D = {2: "two", 3: "three", 'B': 4, 'C': 'CD'}

Accessing

for list use index L[0], L[1] etc

for tuple use index T[0], T[1] etc

for set use x in S ie 3 in S, 4 in S

↳ It returns either True or False

for Dictionary Use D["Val"] eg: D[2] D['B']

The properties of string like slicing and indexing applies to list and tuple.

L[1:3] res: [3, 4, 9]

L[::-1] res: [3, "name", 4, 9, 1, 0]

T[:3] res: (1, 3, 4, 9)

Tuple is immutable you cannot change it.

Adding element.

for list $l = l + ["\text{game"}]$ or $l[1] = "orange"$
 for tuple $T_3 = T_1 + T_2$,
 for set $S.add("newItem")$ individual
 $S.update(["hi", "she"])$ multiple
 for dictionary $DE["val"] = newValue$
 $DE["newValue"] = "newVal"$

Deletion of element

for list $del l[1]$ or $del l$
 for tuple $DEmuttable del T$
 for set $S.remove("banana")/del S$
 for dictionary $del DE["val"]/del D$

$l = [1, 2, 3]$

$l_2 = l + [4, 5, 6]$

$\text{print}(l_2) \quad \text{res: } [1, 2, 3, 4, 5, 6]$

$T_1 = (1, 2, 3)$

$T_2 = (4, 5, 6)$

$T_3 = T_1 + T_2$

$\text{print}(T_3) \quad \text{res: } (1, 2, 3, 4, 5, 6)$

$l_2.append(6, 8)$

$\text{print}(l_2) \quad \text{res: } [1, 2, 3, 4, 5, 6, 6, 8]$

$S = \{1, 2, 3, 4, "name"\}$

$S.add(5, 6) \quad \text{res: } \{1, 2, 3, 4, "name", 5, 6\}$

$S.update(\{2, 3, "he"\}) \quad \text{res: } \{1, 2, 3, 4, "name", 5, 6, 2, 3, "he"\}$

`D = {2: "twothree", "B": 43, 3}`

It's possible to concat two dictionary but not using plus operator rather you use update.

`D2 = { "B": "DB", "C": "CC" }`

`D.update(D2)`
`print(D)`

`res =`

`{2: 'twothree', "B": 43, "B": "BB",
 "C": "CC"}`

Copy function

In python when you assign a variable to another variable it doesn't copy all the content rather just simply points to the same memory location.

If you change any content of l2 then l1 gets affected too.

`l1 = [1, 2, 3]`

`l2 = l1`

`l2`

to fix it you should call copy function

`l2 = l1.copy()`

`l2[2] = "hey" # change just happens
 in l2, l1 remains same.`

Same concept goes for the Set & Dictionary.

Since the content of the tuple can't be changed so copy function isn't available for it.

$L = [1, 2, 3, 4, 5]$

$L3 = L[1:3]$ yes $\therefore L3 = [2, 3]$

Now L3 is a new living List.

$L3[0] \in "Three"$

print(L3) yes $\therefore L3 = [0, 3]$

$L2 = L.copy()$

$S2 = S.copy()$

$D2 = D.copy()$

for Query

help(L.append)

L.clear?

L.pop?

L.reverse() / [P::-1]

In $D2 = \{ 'A': L, 'B': T, 'C': S, 'D': D3 \}$

print $D2['A']$

result: [68, 'you', 6, 'are', 'how',
3, 'four point nine', 3, 1]

print [$'A'$][3]

result

print $D2['A'][3]$ result: are

$K = D2['D']$

print (K)

for n in k:

print (X, K[N])

result:

23 - twothree

B 43

newkey newValue

y yy

z 10

L3 = [L, T, D, D3, "game"]

type (L3[2])

result: diet

L3 = [n**2 for n in range (10)]

L3

[0, 1, 4, 9, 16, 25, 36, 49, 64, 81]

S3 = {n**2 for n in range (2, 2013)}

S3

{4, 25, 64, 121, 196, 289, ...}

Problem of Data structures :-

lets say you are a teacher and you have different student records containing id of a student and the list in each subject where diff students have cap taken different number of subjects. All these records are in handwry. You want to enter all the data in computer and want to compute the average marks of each student & display.

```
def getDatafromUser():
    D = {}
    while True:
        StudentID = input("Enter
        Student ID: ")
        marksList = input("Enter the
        marks by comma separate")
        moreStudents = input(
            "Enter yes/No for adding
            more students:")
        if StudentID in D:
            print(StudentID,
                  "is already inserted")
        else:
            D[StudentID] = []
            marklist = split(",")
            if moreStudents.lower() == "no":
                return D
```

```
def getAvgMarks(D):
    avgMarksDict = {}
    for n in D:
        L = D[n]
        S = 0
        for marks in L:
            L = D[n]
            S = 0
            for marks in L:
                S += int(marks)
            avgMarks[n] = S / len(L)
    return avgMarks
```

avgM = getAvgMarks(Student Data)

for n in avgM:

```
print("Students:", n, "got avg Marks  
as: ", avgM[n])
```

NumPy (python library)

It's very fast. You can import it as:

```
import numpy as np
```

```
a = np.array([1, 2, 3, 4, 5])
```

```
b = np.array([1, 6, 9, 12])
```

```
print(a)
```

```
print(b)
```

You can use both [] () to define an array. Here type of a and b is numpy.ndarray

You can define the type of the array

```
a = np.array([1, 2, 3], dtype='i')
```

```
b = np.array([4, 5, 6], dtype='f')
```

```
a.dtype res = dtype('int32')
```

```
b.dtype res = dtype('float32')
```

Numpy (Dimensions)

`import numpy as np.`

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

`a.ndim` # this gives the dimensions

result = 2

`a[0, 2]` result = 30
 ↑

This defines which array to choose

This defines which component of array is chosen.

`B = np.array`

when you're defining one dimensional array the number of elements should remain constant.

Defining three dimensional array.

`B = np.array([[[1, 2, 3], [4, 5, 6], [7, 8, 9]], [[-1, -2, -3], [-5, -6, -7], [-8, -9, -10]]])`

`B.ndim` (result = 3)

`B[1, 0, 2]`

Defines whether it's first or second group of array

chose a particular array within the group choose index.

$B = [[[1, 2, 3], [4, 5, 6], [7, 8]], [[-1, -2, -3], [-4, -5, -6], [-7, -8, -9]]]$

$B.shape$ res: $(2, 3, 3)$

There are 2
2-d array

There are 3
1-d array

There are 3
elements in 1-d
array.

$C.shape[0]$ res: 3

$C.shape[1]$ res: 3

$A = np.array([2])$

$A.ndim$ res: 1

$B = np.array(3)$

$B.ndim$ res: 0

$B.size$ res: 18

$B.nbytes$ res: 72

Numpy (np·arange
 np·random·permutation
 np·reshape)

A = np·arange (100)
 [0, 1, 2, 3, ..., 97, 98, 99]

A = np·arange (20, 100)
 [20, 21, 22, ..., 97, 98, 99]

A = np·arange (20, 100, 3)
 [20, 23, 26, 27, ..., 92, 95, 98]

Iteration : process of executing a block of code repeatedly.

```
print (range(10))
range (0,10)
print (list (range(10)))
[0,1,2,3,4,5,6,7,8,9]
```

np·random has a lot of package one of them is permutation.

A = np·random·permutation (np·arange(10))
 print (A)
 [3, 9, 5, 7, 6, 4, 2, 1, 0, 8]

A = np·random·randint (20, 30)
 res = 27 # returns random no within range.
 type (A) is integer.

`np.reshape`

`D = np.arange(100).reshape(4,25)`

`D.shape`

$(4, 25)$

D will be 4 by 25 matrix with
4 rows and 25 columns. This will be
a two dimensional array.

`D = np.arange(100).reshape(4,5,5)`

`D.shape`

$(4, 5, 5)$

`np.zeros?`

`np.ones?`

`A = np.random.rand(1000)`

`print(A)`

prints 1000 random numbers
between the range 0-1

(we'll import plotting package)

`import matplotlib.pyplot as plt`
`plt.hist(A)`

(you get a uniform graph)
`plt.hist(A, bins=100)`

↓

100 individual data.

`B = np.random.randn(1000)`

`plt.hist(B, bins=200)`

creates a \sim this natured curved.

rand → uniform distribution

randn → standard normal distribution.

You can also create random matrix array using.

```
n = np.random.rand(2,3)
```

```
print(n)
```

```
res = [[0.44, 0.53, 0.26], [0.32, 0.42, 0.5]]
```

```
n.ndim res // 2
```

and,

```
c = np.random.rand(2, 3, 4, 2)
```

```
c.ndim res // 4
```

NUMPY SLICING

a[1:5] # index 1 till 5 but not 5

a[:5] # index 0 till 5 but not 5

a[2:] # index 2 till end including last element.

a[::-1] # from end till start (reverses the array)

a[::-2] # from start till end every other element.

a[:, :, ?]

```
A = np.arange(100) # imp
```

```
b = A[3:10]
```

```
print b [3 4 5 6 7 8 9]
```

```
b[0] = -1200
```

```
P [-1200 4 5 6 7 8 9]
```

Changing the element of b also changes the element of a as despite of having different variable name they share same memory location.

If you don't want that to happen then

$b = a[3:10].copy()$

Then if you change any element in b then the original value of a is not altered.

$b = a[3:10].copy()$

$A[::5]$

array ([0 5 10 15 ...
... 85 90 95])

$A[:::-5]$

array ([99 94 89 84 ...
... 14 9 4])

A

array ([0, 1, 2, 3, 4, ..., 97, 98, 99])

$idx = np.argwhere(A == 3)[0][0]$
 $idx \approx 3$
 $A[idx] = 400$

Here we found the index of array.

`A = np.random(10* np.random.rand(5,4))`

```
array([[3., 6., 10., 4.],
       [4., 5., 3., 7.],
       [1., 3., 4., 5.],
       [4., 9., 8., 7.],
       [2., 3., 1., 4.]])
```

`A[1,2]` i.e 3

`A[1,:]` # accessing the whole 2nd row.

```
array([4., 5., 3., 7.])
```

`A[:,1]` # accessing the whole 2nd col

```
array([6., 5., 3., 9., 3.])
```

`A[1:3,2:4]` row 1 to row 3

array([[3, 7], [4, 5]]) column 2 to 4

(end not included)

Endpoint is exclusive.

`A.T` gives the transpose of the matrix

`import numpy.linalg as lg`

`(A).inv(np.random.rand(3,3))`

(returns the inverse of the matrix)

A

```
array [[ 3., 6., 10., 4.],
      [ 6., 4., 8., 8.],
      [ 5., 2., 7., 8.],
      [ 5., 10., 8., 10.],
      [ 4., 1., 5., 7.]]
```

for sorting each column individually.

A. sort (axis=0)

```
A
array [[ 3., 1., 5., 0.],
      [ 4., 2., 7., 4.],
      [ 5., 4., 8., 7.],
      [ 5., 6., 8., 8.],
      [ 6., 10., 10., 10.]]
```

for sorting each row individually

A. sort (axis=1)

```
A
array [[ 0., 1., 13., 5.],
      [ 2., 4., 4., 7.],
      [ 4., 5., 7., 8.],
      [ 5., 6., 8., 8.],
      [ 6., 8., 10., 10.]]
```

Numpy (More Indexing)

A [index_array]

A[[1, 4, 6]] # index 1, 4 and 6 elements

a [[True, false, True, True, False]]

Assuming array has 5 elements, the above returns all the elements corresponding to True Index.

If you want to get all the element that are less than 8 in the index_array then we use :

A[a < 8]

A[a < 8 & a > 4] # difference between (and, &)?

A[a < 8 and a > 4]

*, A = np.arange(100)

B = A[[3, 5, 6]]

B[0] = -4

print(B) res: [-4, 5, 6]

(= A[A < 40])

print(() ie [0, 1, 2, ..., 37, 38, 39])

$$B = A[(A < 40) \& (A > 30)]$$

B

array ([31, 32, 33, 34, 35, 36, 37, 38,
39])

and : It's used when both of the sides is one objects and has one value either true or false.

& : When the left side and the right side can be arrays & each element of that array can be true or false.

&, and # !, or # ~, not

NUMPY (Broadcasting)

$$A = A + 5$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Normally you can't do this
but here it's converted
to $\begin{bmatrix} 5 & 5 \\ 5 & 5 \end{bmatrix}$

$$A = A + \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

↓

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

This is turned to $\begin{bmatrix} 1 & 1 \\ 2 & 2 \end{bmatrix}$

NUMPY

np.hstack

np.vstack

np.sort

A = np.round(10 * np.random.rand(2,3))

A
array([[1.1, 2.1, 3.1], [0.1, 0.1, 2.1]])

B = np.arange(12).reshape(2,6)

array([[1.1, 2.1, 3.1],
[4.1, 5.1, 6.1]])

B = np.round(10 * np.random.rand(2,2))

A
array([[2.1, 3.1],
[0.1, 1.1]])

B
array([[4.1, 2.1],
[5.1, 1.1]])

C = np.hstack((A, B))

C
array([[2.1, 3.1, 4.1, 2.1],
[0.1, 1.1, 5.1, 1.1]])

A = np.random.permutation(np.arange(10)
[0, 1, 4, 3, 7, 8, 6, 5, 2, 9])

A.sort()

[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]

np.sort(A)

Numpy Ueed: UFuncs

B = np.random.rand(1000000)

%timeit sum(B) → 300ms

%timeit np.sum(B) → 3ms

def mysum: %timeit
 s = 0
 for x in b:
 s += x
 return s

def mysum(c):
 s = 0
 for x in c:
 s += x
 return s

%timeit mysum(B)

Pandas

import pandas as pd

print(pd.__version__)

A = pd.Series([2, 3, 4, 5], index=['a', 'b', 'c', 'd'])

type(A.values)

res: numpy.ndarray

type(A)

res: pandas.core.series.Series

A.index

Index(['a', 'b', 'c', 'd'], dtype='object')

A['a'] res: 2

A['a':'c'] res: 2
b 3
c 4

grades_dict = {'A': 4, 'A-': 3.5, 'B': 3,
'B-': 2.5, 'C': 2}

grades = pd.Series(grades_dict)

marks_dict = {'A': 85, 'A-': 80, 'B': 75,
'B-': 70, 'C': 65}

marks = pd.Series(marks_dict)

print(grades.values)

array [4., 3.5, 3., 2.5]

print

```
marks_dict = {'A': 85, 'B': 75, 'C': 65,
              'D': 55}
```

```
marks = pd.Series(marks_dict)
```

A	85
B	75
C	65
D	55

```
dtype: int64
```

```
marks[A] res: 85
```

```
marks[0:2] res: A 85  
          B 75
```

```
datatype: int64
```

```
D = pd.DataFrame({'Marks': marks,
                  'Grades': grads})
```

```
print(D)
```

	Marks	Grades
A	85	4.0
B	75	3.5
C	65	3.0
D	55	2.5

```
print(D.T)
```

	A	B	C	D
Marks	85.0	75.0	65.0	55.0
Grades	4.0	3.5	3.0	2.5

D-values

```
array ([[ 85., 4.],
       [ 75., 3.5],
       [ 65., 3.],
       [ 55., 2.5]])
```

D-values[2,0] → 65.0

D.columns

```
Index(['Marks', 'Grades'],
      dtype='object')
```

D['scaledMarks'] = (D['Marks'] / 90) * 100

D

	Marks	Grades	scaledMarks
A	85	4.0	94.44
B	75	3.5	83.33
C	65	3.0	72.22
D	55	2.5	61.11

del D['scaledMarks']

D

	Marks	Grades	
A	85	4.0	
B	75	3.5	
C	65	3.0	
D	55	2.5	

G = D[D['Marks'] > 70]

G

Marks

Grades

A 85 4.0

B 75 3.5

Pandas (NaN)

`pd.DataFrame([{'a': 1, 'b': 2, 'c': 3}, {'b': 3, 'c': 4}])`

	a	b	c
0	1.0	2	NaN
1	NaN	3	4.0

`A = pd.DataFrame([{'a': 10, 'b': 15}, {'a': 13, 'c': 9}])`

	a	b	c
0	10	15	NaN
1	NaN	13	9

`A.fillna(6.8)`

	a	b	c
0	10	15	6.8
1	6.8	13	9

`A.dropna?` (Remove missing value)

Pandas (Indexing)

`data = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])`

`data[1]` # explicit index, use loc instead

`data[1:3]` # implicit index, use iloc instead.

`vs. iloc[2:3]`

eg:- A = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])

A[1] result :- (a)

A[1:3] result :-
3 b
5 c

dtype: object

A.loc[1:3]

result :-
1 a
3 b

dtype: object

A.iloc[1:3]

result :-
3 b
5 c

dtype: object

D result :-

Marks Grades

A	85	4.0
B	75	3.5
C	65	3.0
D	55	2.5

D.iloc[2, :]

Marks :- 65.0

Grades :- 3.0

Name :- C, dtype: float64

D.iloc[:: -1, :]

result :-

Marks Grades

D	55	2.5
C	65	3.0
B	75	3.5
A	85	4.0

Pandas (CSV file)

```

import pandas as pd
import numpy as np
from sklearn.impute import SimpleImputer
# df = pd.read_csv('covid_19_data.csv')
df = pd.read_csv('E:/covid/covid_19_data.csv')
df.head()
# first few records is shown
df.head(10)
# first 10 records
df.drop(['SNO', 'Last Update'],
        axis=1, inplace=True)
# deletes the SNO and last update
# from the CSV file and inplace=
# True assures that it's to be done
# in the original vari CSV file.
# Now if you call df.head() it
# won't show SNO and it will
# also not display the last update

```

```

df.rename(columns = {'Observation Date': 'Date',
                     'Province/State': 'Province',
                     'Country/Region': 'Country',
                     inplace=True)

```

`df['Date'] = pd.to_datetime(df['Date'])`

`df.head()` # to see the result

Initially

1/22/2020

New date format

2020-01-22

Use this to change any date format to pandas date format.

`df.describe()` # gives you details about the data

`df.info()` # more information about the file and all.

`df.fillna('NA')` # All missing values will be filled with NA.

`df.info()` # you'll see no null values.

Q Group by country and show country, Confirmed cases, death and recovered by their total number.

`df2 = df.groupby(['Country'])[['Confirmed', 'Deaths', 'Recovered']].sum().reset_index()`

`df3 = df2[df2['Confirmed'] > 100]`

Matplotlib

```
import matplotlib.pyplot as plt.
```

```
x=np.linspace(0,10,100)
```

```
plt.plot(x, np.sin(x));
```

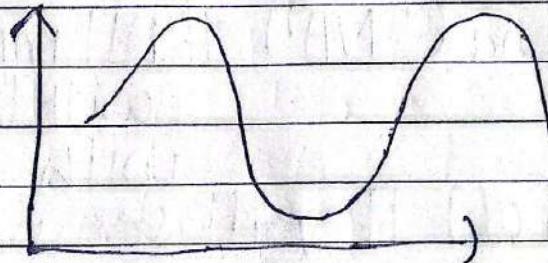
another

```
import matplotlib.pyplot as plt
```

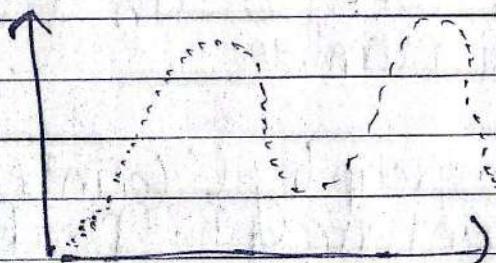
```
x=np.linspace(0,10,1000)
```

```
y=np.sin(x)
```

```
plt.plot(x,y)
```



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



more clear visualization of

```
plt.scatter(x[::10],y[::10])
```

You can also give a color = 'red'
just use comma and add it.

Some things to do:

```
plt.plot(x1,x1+0, 'g') # solid green
plt.plot(x1,x1+1, '--c') # dashed cyan
plt.plot(x1,x1+2, '-.k') # dashed black
plt.plot(x1,x1+3, ':r') # dotted red
```

Eg:-

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
from sklearn.impute import Simple
Imputer
```

```
df = pd.read_csv('E:/covid/covid-19_data.
csv')
df.head(50)
```

```
df.drop(['SNO', 'Last Update'], axis=1,
inplace=True)
df.rename(columns= {'ObservationDate': 'Date',
'Province/State': 'State', 'Country/
Region': 'Country'})
```

```
df['Date'] = pd.to_datetime(df['Date'])
imputer = SimpleImputer(strategy =
'constant')
```

```
df2 = pd.DataFrame(imputer.fit_transform(df),
columns=df.columns)
```

```
df3 = df2.groupby(['Country', 'Date'])
[[['Country', 'Date', ('Confirmed', 'Deaths', 'Recovered')].sum()]
.reset_index()
```

df3.head(70)

Countries = df3[['Country']].unstack(1)[len(Countries)].reset_index()

for idn in range(0, len(Countries)):

```
c = df3[df3['Country'] == Countries
[idn]].reset_index()
```

```
plt.scatter(np.arange(0, len()), c[['Confirmed']], color='blue', label
= 'Confirmed')
```

```
plt.scatter(np.arange(0, len()), c[['Recovered']], color='green',
label = 'Recovered')
```

```
plt.scatter(np.arange(0, len()), c[['Deaths']], color='red', label=
'Deaths'))
```

plt.title(Countries[idn])

plt.xlabel('Days since the first
suspect')

plt.ylabel('Number of cases')

plt.legend()

plt.show()

This is for each country individually

If we want overall trend, the overall growth and the overall death of the world then we can do:

```
df4 = df3.groupby(['Date'])[['Date',
    'Confirmed', 'Deaths', 'Recovered']].  
sum().reset_index()
```

C = df4

```
plt.scatter(np.arange(0, len()),  
C['Confirmed'], color='blue', label='confirmed')  
plt.scatter(np.arange(0, len()),  
C['Recovered'], color='green', label='recovered')
```

```
plt.scatter(np.arange(0, len()),  
C['Deaths'], color='red',  
label='Deaths')
```

plt.title('World')

plt.xlabel('Days since the first suspect')

plt.ylabel('Number of cases')

plt.legend()

plt.show()