

67070501021 นายธนบุญ ธิแก้ว

## Lab 1: Basic Python Programming

CPE232 Data Models

---

### [1] Variable

#### 1.1 Number Variable

```
In [1]: num = 100 #integer variable
num2 = 12.5 #float variable
print(num)
print(num2)

print(num + num2) #addition
print(num - num2) #subtraction
print(num * num2) #multiplication
print( num / num2) #division
```

```
100
12.5
112.5
87.5
1250.0
8.0
```

#### 1.2 String Variable

```
In [2]: string = "Data Models"
print(string) #print complete string

print("Hello " + string) #print concatenated string

# slicing string
print(string[0]) #print first character of the string
print(string[:4]) #print first to 4th character of the string
print(string[5:]) #print 6th to last character of the string
print(string[1:4]) #print 2nd to 4th character of the string
print(string * 2) #print string 2 times
```

```
Data Models
Hello Data Models
D
Data
Models
ata
Data ModelsData Models
```

```
In [3]: # format string
code = "CPE232"

print(f"{code}: Data Models")
```

CPE232: Data Models

## 1.3 Boolean Variable

```
In [4]: #boolean variable
boolean = True
boolean2 = False

print(boolean)           #print boolean variable
print(not boolean)       #print opposite of boolean variable
print(boolean and boolean2) #print boolean and boolean2
print(boolean or boolean2) #print boolean or boolean2
```

```
True
False
False
True
```

## 1.4 List Variable

```
In [5]: #list variable
list = ["Data",20,123.23,40,50]
another_list = ["Models",60]

print(list)           #print complete list
print(list[0])        #print first element of the list
print(list[1:3])      #print 2nd to 3rd element of the list
print(list[2:])       #print 3rd to last element of the list
print(another_list)   #print complete another_list
print(another_list * 2) #print another_list two times
print(list + another_list) #print concatenated list

list[0] = "CPE232"    #change first element of the list
print(list)           #print complete list

# add element at the end
list.append(True)
print(list)

# pop element by index
list.pop(1)
print(list)
```

```

['Data', 20, 123.23, 40, 50]
Data
[20, 123.23]
[123.23, 40, 50]
['Models', 60]
['Models', 60, 'Models', 60]
['Data', 20, 123.23, 40, 50, 'Models', 60]
['CPE232', 20, 123.23, 40, 50]
['CPE232', 20, 123.23, 40, 50, True]
['CPE232', 123.23, 40, 50, True]

```

## 1.5 Tuple Variable

```

In [6]: #tuple variable
tuple = ("Data",20,123.23,40,50)
another_tuple = ("Models",60)

print(tuple)                #print complete tuple
print(tuple[0])             #print first element of the tuple
print(tuple[1:3])          #print 2nd to 3rd element of the tuple
print(tuple[2:])            #print 3rd to last element of the tuple
print(tuple * 2)            #print tuple two times
print(tuple + another_tuple) #print concatenated tuple

```

```

('Data', 20, 123.23, 40, 50)
Data
(20, 123.23)
(123.23, 40, 50)
('Data', 20, 123.23, 40, 50, 'Data', 20, 123.23, 40, 50)
('Data', 20, 123.23, 40, 50, 'Models', 60)

```

```

In [7]: tuple[0] = "CPE232"          # tuple is immutable dtype, we can't change its value

```

```

-----
TypeError                                Traceback (most recent call last)
Cell In[7], line 1
----> 1 tuple[0] = "CPE232"          # tuple is immutable dtype, we can't change its
value after declared

TypeError: 'tuple' object does not support item assignment

```

## 1.6 Dictionary Variable

```

In [8]: #dictionary variable
dictionary = {"name":"Alice","age":21}
another_dictionary = {}
another_dictionary["name"] = "Bob"
another_dictionary["age"] = 21

print(dictionary)           #print complete dictionary
print(dictionary["name"])   #print value for specific key
print(dictionary.keys())    #print all the keys
print(dictionary.values())  #print all the values

```

```
print(dictionary.items())      #print all the items
print(another_dictionary)     #print complete another_dictionary
```

```
{'name': 'Alice', 'age': 21}
Alice
dict_keys(['name', 'age'])
dict_values(['Alice', 21])
dict_items([('name', 'Alice'), ('age', 21)])
{'name': 'Bob', 'age': 21}
```

## 1.7 Set Variable

```
In [9]: #set variable
my_set = {"Data", 20, 123.23, 60, "Data"}
another_set = {"Models", 60}

print(my_set)                    #print complete set (duplicates are removed)

#set operations
print(my_set.union(another_set)) # print all unique elements from both sets
print(my_set.intersection(another_set)) # print common elements from both sets

# add element
my_set.add("New Item")
my_set.add(20)                   # try adding a duplicate again
print(my_set)

# remove element
my_set.remove(123.23)
print(my_set)
```

```
{123.23, 'Data', 20, 60}
{'Models', 'Data', 20, 123.23, 60}
{60}
{'New Item', 'Data', 20, 123.23, 60}
{'New Item', 'Data', 20, 60}
```

## [2] Control Flow

### 2.1 IF ... ELIF ... ELSE

```
In [10]: number = 123
number2 = 34

if number > number2:
    print("number is greater than number2")
elif number < number2:
    print("number is less than number2")
else:
    print("number is equal to number2")
```

number is greater than number2

## 2.2 try ... except

```
In [11]: number1 = "number"
         number2 = 2

         try:
             number1 + number2
         except:                     #error handling
             print("What are you doing?")
```

What are you doing?

## [3] Loop

### 3.1 For Loop

```
In [12]: #for Loops
         for num in range(0,10):
             print(num)
```

0  
1  
2  
3  
4  
5  
6  
7  
8  
9

```
In [13]: #for Loops with step
         for num in range(5,15,2):
             print(num)
```

5  
7  
9  
11  
13

```
In [14]: #for Loop with enumerate

         list = ["Alice", "Bob", "Charlie", "Daisy"]

         for index, name in enumerate(list):
             print(f"{index}: {name}")
```

0: Alice  
1: Bob  
2: Charlie  
3: Daisy

In [15]: *#for loop with list*

```
for name in list:  
    print(name)
```

Alice  
Bob  
Charlie  
Daisy

In [16]: *#continue in for loop*

```
list = [1,23,7,"hello",True,1123,43,23,12]  
  
for element in list:  
    if type(element) != int:  
        continue  
    print(element)
```

1  
23  
7  
1123  
43  
23  
12

In [17]: *#break in for loop*

```
list = [1,23,7,"hello",True,1123,43,23,12]  
  
for element in list:  
    if type(element) != int:  
        break  
    print(element)
```

1  
23  
7

## 3.2 While loop

In [18]: *#while loop*

```
list = ["Alice","Bob","Charlie","Daisy"]  
count = 0  
  
while count < len(list):  
    print(list[count])  
    count += 1
```

Alice  
Bob  
Charlie  
Daisy

In [19]: *#continue in while loop*

```
list = [1,23,7,"hello",True,1123,43,23,12]
count = 0

while count < len(list):
    if type(list[count]) != int:
        count += 1
        continue
    print(list[count])
    count += 1
```

1  
23  
7  
1123  
43  
23  
12

In [20]: *#break in while loop*

```
list = [1,23,7,"hello",True,1123,43,23,12]
count = 0

while count < len(list):
    if type(list[count]) != int:
        break
    print(list[count])
    count += 1
```

1  
23  
7

## [4] Function

In [21]: *#define function*

```
def function_name (arg1, arg2):
    return arg1 + arg2

#calling function
function_name(1,2)
```

Out[21]: 3

In [22]: *#define function with default argument*

```
def function_with_default_arg(arg1, arg2 = 10, arg3 = 20 , arg4 = 30):
    return arg1 + arg2 + arg3 + arg4

result_1 = function_with_default_arg(1)
result_2 = function_with_default_arg(1,2,5)
result_3 = function_with_default_arg(1,2,5,10)
```

```
print(result_1)
print(result_2)
print(result_3)
```

61  
38  
18

```
In [23]: #multiple arguments
def function_with_multiple_arg(*args):
    print(args)
    print(type(args))
    sum = 0
    for num in args:
        sum += num

    return sum

function_with_multiple_arg(1,2,3,4,5)
```

(1, 2, 3, 4, 5)  
<class 'tuple'>

Out[23]: 15

```
In [24]: #lambda function
lambda_function = lambda arg1, arg2: arg1 + arg2

print(lambda_function(1,2))
```

3

## [5] OOP



```
In [25]: class Car:
        """A simple attempt to represent a car."""

        def __init__(self, make, model, year):
            """Initialize attributes to describe a car."""
            self.make = make
            self.model = model
            self.year = year
            self.odometer_reading = 0

        def get_description_name(self):
            """Return a neatly formatted descriptive name."""
            long_name = f"{self.year} {self.make} {self.model}"
            return long_name.title()

        def read_odometer(self):
            """Print a statement showing the car's mileage."""
            print(f"This car has {self.odometer_reading} miles on it.")

        def update_odometer(self, mileage):
            """Set the odometer reading to the given value."""
            self.odometer_reading = mileage
```

```
In [26]: my_car = Car("Honda", "Civic", 2016)
        my_car.read_odometer()

        # update odometer
        my_car.update_odometer(50)
        my_car.read_odometer()
```

This car has 0 miles on it.  
This car has 50 miles on it.

## [6] File Handling

### 6.1 Text File

```
In [27]: # write file
        with open("test.txt", "w") as file:
            file.write("Hello World")

        # read file
        with open("test.txt", "r") as file:
            print(file.read())
```

Hello World

### 6.2 CSV File

```
In [28]: import csv

        with open("test.csv", "w", newline='') as file:
            writer = csv.writer(file)
```

```
writer.writerow(["Name", "Surname"])
writer.writerow(["Alice", "Johnson"])
writer.writerow(["Bob", "Smith"])
```

```
In [29]: import csv

with open("test.csv", "r") as file:
    reader = csv.reader(file)
    for row in reader:
        print(row)
```

```
['Name', 'Surname']
['Alice', 'Johnson']
['Bob', 'Smith']
```

## 6.3 JSON file

```
In [30]: import json

data = {
    "username": "admin",
    "password": "admin",
    "role": "sleeping",
}

with open('test.json', 'w') as f:
    json.dump(data, f)
```

```
In [31]: with open("test.json", "r") as f:
    data = json.load(f)

data
```

```
Out[31]: {'username': 'admin', 'password': 'admin', 'role': 'sleeping'}
```

## [7] Libraries

### 7.1 Numpy

import numpy library

```
In [32]: import numpy as np
```

#### ndarray initialization

Construct using python list

```
In [33]: # 1d ndarray from 1d python list
list_a1=[1,2,3.5]
```

```
arr_a1=np.array(list_a1)
arr_a1
```

Out[33]: array([1. , 2. , 3.5])

```
In [34]: # 2d ndarray from 2d python list (list of list)
list_a2=[[1,2],[3,4],[5,6]]
arr_a2=np.array(list_a2)
arr_a2
```

Out[34]: array([[1, 2],  
[3, 4],  
[5, 6]])

```
In [35]: list_a3=[ [[1,2],[2,3]], [[3,4],[4,5]] ]
arr_a3=np.array(list_a3)
arr_a3
```

Out[35]: array([[ [1, 2],  
[2, 3]],  
  
[[3, 4],  
[4, 5]]])

or construct using some numpy classes and functions

```
In [36]: np.zeros(5)
```

Out[36]: array([0., 0., 0., 0., 0.])

```
In [37]: np.ones((3,4),dtype=float)
```

Out[37]: array([[1., 1., 1., 1.],  
[1., 1., 1., 1.],  
[1., 1., 1., 1.]])

```
In [38]: np.eye(3, 3)
```

Out[38]: array([[1., 0., 0.],  
[0., 1., 0.],  
[0., 0., 1.]])

```
In [39]: np.full((4,),999)
```

Out[39]: array([999, 999, 999, 999])

```
In [40]: np.arange(3,10,2)
```

Out[40]: array([3, 5, 7, 9])

```
In [41]: np.linspace(0, 1, 5)
```

Out[41]: array([0. , 0.25, 0.5 , 0.75, 1. ])

```
In [42]: np.random.choice(['a','b'], 9)
```

```
Out[42]: array(['a', 'b', 'b', 'a', 'b', 'a', 'b', 'b', 'a'], dtype='<U1')
```

```
In [43]: np.random.randn(10)
```

```
Out[43]: array([ 1.35591984, -0.91629096,  1.58048896,  1.18457964, -1.28496854,
                 0.30538524,  1.74836835,  0.41450679, -0.35972875,  0.2995627 ])
```

## ndarray properties

```
In [44]: list_a=[[1,2,3,4],[5,6,7,8],[9,10,11,12]]
arr_a=np.array(list_a)
arr_a
```

```
Out[44]: array([[ 1,  2,  3,  4],
                [ 5,  6,  7,  8],
                [ 9, 10, 11, 12]])
```

```
In [45]: print("dimension:", arr_a.ndim)
print("shape:", arr_a.shape)
print("data type:", arr_a.dtype)
print("array's size (elements)", arr_a.size)
```

```
dimension: 2
shape: (3, 4)
data type: int64
array's size (elements) 12
```

## Reshaping & Modification

from this original ndarray

```
In [46]: arr_a
```

```
Out[46]: array([[ 1,  2,  3,  4],
                [ 5,  6,  7,  8],
                [ 9, 10, 11, 12]])
```

try to convert into 3D array

```
In [47]: arr_a.reshape((2,2,3))
```

```
Out[47]: array([[[ 1,  2,  3],
                 [ 4,  5,  6]],
                [[ 7,  8,  9],
                 [10, 11, 12]])
```

sometimes you may resize for same dimension where only known some dimension, insert -1 for unknown len

```
In [48]: arr_a.reshape((-1,6))
```

```
Out[48]: array([[ 1,  2,  3,  4,  5,  6],
               [ 7,  8,  9, 10, 11, 12]])
```

Would you like to try this?

```
In [49]: arr_a.reshape((-1,5))
```

```
-----
ValueError                                Traceback (most recent call last)
Cell In[49], line 1
----> 1 arr_a.reshape((-1,5))

ValueError: cannot reshape array of size 12 into shape (5)
```

[Q1] From the above cell, explain in your own words why it worked or did not work.

Ans: It did not work because array in shape 5 can be formed by Array with elements which is multiplicant of 5 (5,10,15,...) but arr\_a have 12 elements

Next, try to append any value(s) into existing 2darray

```
In [50]: np.append(arr_a, 13)
```

```
Out[50]: array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12, 13])
```

```
In [51]: np.append(arr_a, arr_a[0])
```

```
Out[51]: array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 11, 12,  1,  2,  3,  4])
```

```
In [52]: np.append(arr_a, arr_a[0].reshape((1,-1)),axis=0)
```

```
Out[52]: array([[ 1,  2,  3,  4],
               [ 5,  6,  7,  8],
               [ 9, 10, 11, 12],
               [ 1,  2,  3,  4]])
```

```
In [53]: np.append(arr_a, arr_a[:,0].reshape((-1,1)),axis=1)
```

```
Out[53]: array([[ 1,  2,  3,  4,  1],
               [ 5,  6,  7,  8,  5],
               [ 9, 10, 11, 12,  9]])
```

```
In [54]: np.concatenate([arr_a, arr_a])
```

```
Out[54]: array([[ 1,  2,  3,  4],
               [ 5,  6,  7,  8],
               [ 9, 10, 11, 12],
               [ 1,  2,  3,  4],
               [ 5,  6,  7,  8],
               [ 9, 10, 11, 12]])
```

```
In [55]: np.concatenate([arr_a, arr_a],axis=1)
```

```
Out[55]: array([[ 1,  2,  3,  4,  1,  2,  3,  4],
                [ 5,  6,  7,  8,  5,  6,  7,  8],
                [ 9, 10, 11, 12,  9, 10, 11, 12]])
```

## indexing & slicing

from this original array again

```
In [56]: arr_a
```

```
Out[56]: array([[ 1,  2,  3,  4],
                [ 5,  6,  7,  8],
                [ 9, 10, 11, 12]])
```

try to access all element at the first row

```
In [57]: arr_a[1]
```

```
Out[57]: array([5, 6, 7, 8])
```

then you would like to access the second element from the first row

```
In [58]: arr_a[1][2]
arr_a[1,2]
```

```
Out[58]: np.int64(7)
```

Next, try to access all element start from the 2nd element in the first row

```
In [59]: arr_a[1,1:]
```

```
Out[59]: array([6, 7, 8])
```

```
In [60]: arr_a[:2, 1:]
```

```
Out[60]: array([[2, 3, 4],
                [6, 7, 8]])
```

sometimes you may specify some row number using list within indexing

```
In [61]: arr_a[[1,2,1], 1:]
```

```
Out[61]: array([[ 6,  7,  8],
                [10, 11, 12],
                [ 6,  7,  8]])
```

## Boolean slicing

based on this original array

```
In [62]: arr_a
```

```
Out[62]: array([[ 1,  2,  3,  4],
               [ 5,  6,  7,  8],
               [ 9, 10, 11, 12]])
```

try to filter all elements which more than 5

```
In [63]: arr_a>5
```

```
Out[63]: array([[False, False, False, False],
               [False,  True,  True,  True],
               [ True,  True,  True,  True]])
```

Next, try to filter all elements which more than 5 and less than 10

```
In [64]: (arr_a>5)&(arr_a<10)
```

```
Out[64]: array([[False, False, False, False],
               [False,  True,  True,  True],
               [ True, False, False, False]])
```

Run the cell below and answer a question.

```
In [65]: arr_a[(arr_a>5)&(arr_a<10)]
```

```
Out[65]: array([6, 7, 8, 9])
```

[Q2] From the above cell, explain in your own words how the output came about?

Ans: The output came from elements in arr\_a which is intersection between arr\_a > 5 (6,7,8,9,10,11,12) and arr\_a < 10 (1,2,3,4,5,6,7,8,9) so the intersection is (6,7,8,9) or index 5 to index 8 and we filter arr\_a using the intersection index

Try running the cell below.

```
In [119... arr_a[(arr_a>5) and (arr_a<10)]
```

```
-----
ValueError                                Traceback (most recent call last)
Cell In[119], line 1
----> 1 arr_a[(arr_a>5) and (arr_a<10)]

ValueError: The truth value of an array with more than one element is ambiguous. Use
a.any() or a.all()
```

[Q3] Explain in your own words why the above cell gives an error.

Ans: "and" operator is using on single element in this case `arr_a > 5` have 12 elements which is `array([[False, False, False, False], [False, True, True, True], [ True, True, True, True]])` so "and" operator don't know what is Truth value of `arr_a > 5` and occurs the error

while "&" operator check element by element(bitwise) so it works on `arr_a > 5 & arr_a < 10` both have same 12 elements it intersect it one by one so it can give the right output

[Q4] And what should be written instead so that the code is error-free?

Ans: `arr_a[(arr_a>5).any() and (arr_a<10).any()]`

```
In [120]: arr_a[(arr_a>5).any() and (arr_a<10).any()]
```

```
Out[120]: array([[ 1,  2,  3,  4],
                [ 5,  6,  7,  8],
                [ 9, 10, 11, 12]])
```

## Basic operations

```
In [67]: list_b=[[1,2,3,4],[1,2,3,4],[1,2,3,4]]
arr_b=np.array(list_b)
arr_b
```

```
Out[67]: array([[1, 2, 3, 4],
                [1, 2, 3, 4],
                [1, 2, 3, 4]])
```

This is some operations for only 1 array

```
In [68]: np.sqrt(arr_b)
```

```
Out[68]: array([[1.         , 1.41421356, 1.73205081, 2.         ],
                [1.         , 1.41421356, 1.73205081, 2.         ],
                [1.         , 1.41421356, 1.73205081, 2.         ]])
```

This is some operations for 2 arrays with the same shape

```
In [69]: arr_a-arr_b
```

```
Out[69]: array([[0, 0, 0, 0],
                [4, 4, 4, 4],
                [8, 8, 8, 8]])
```

```
In [70]: np.add(arr_a,arr_b)
```

```
Out[70]: array([[ 2,  4,  6,  8],
                [ 6,  8, 10, 12],
                [10, 12, 14, 16]])
```

Next, try to operate with 1 array and one numeric variable



```
In [71]: arr_a*3
```

```
Out[71]: array([[ 3,  6,  9, 12],
               [15, 18, 21, 24],
               [27, 30, 33, 36]])
```

```
In [72]: 1+arr_a**2
```

```
Out[72]: array([[ 2,  5, 10, 17],
               [26, 37, 50, 65],
               [82, 101, 122, 145]])
```

## Broadcasting

```
In [73]: arr_c=np.array([1,2,3])
arr_d=np.array([[3],[5],[8]])
```

```
In [74]: arr_c-arr_d
```

```
Out[74]: array([[ -2,  -1,   0],
               [-4,  -3,  -2],
               [-7,  -6,  -5]])
```

## Basic aggregations

```
In [75]: arr_a

print(arr_a)
print("sum:", arr_a.sum())
print("mean:", arr_a.mean())
print("min:", arr_a.min())
print("max:", arr_a.max())
print("standard deviation:", arr_a.std())
```

```
[[ 1  2  3  4]
 [ 5  6  7  8]
 [ 9 10 11 12]]
sum: 78
mean: 6.5
min: 1
max: 12
standard deviation: 3.452052529534663
```

## ndarray axis

```
In [76]: # column axis
arr_a.sum(axis=0)
```

```
Out[76]: array([15, 18, 21, 24])
```

```
In [77]: # row axis
arr_a.sum(axis=1)
```

```
Out[77]: array([10, 26, 42])
```

[Q5] Summarize the value of the argument *axis*, what is the value for row-wise summation and column-wise summation, respectively?

Ans: row-wise summation axis = 1 and column-wise summation axis = 0

## 7.2 Pandas

### Series

```
In [84]: import pandas as pd  
import numpy as np
```

```
In [85]: pd.Series(np.random.randn(6))
```

```
Out[85]: 0    -0.890746  
1    -0.391122  
2    -0.844234  
3    -0.615343  
4     0.572527  
5     0.880289  
dtype: float64
```

```
In [86]: pd.Series(np.random.randn(6), index=['a', 'b', 'c', 'd', 'e', 'f'])
```

```
Out[86]: a    -0.593951  
b    -0.670971  
c     0.217029  
d    -0.504141  
e     0.010253  
f    -1.194180  
dtype: float64
```

### Constructing Dataframe

Constructing DataFrame from a dictionary

```
In [81]: d = {'col1': [1, 2], 'col2': [3, 4]}
```

```
In [87]: df = pd.DataFrame(data=d)  
df
```

```
Out[87]:
```

	col1	col2
0	1	3
1	2	4

```
In [88]: d2 = {'Name': ['Joe', 'Nat', 'Harry', 'Sam', 'Monica'],  
              'Age': [20, 21, 19, 20, 22]}
```

```
df2 = pd.DataFrame(data=d2)
df2
```

Out[88]:

	Name	Age
0	Joe	20
1	Nat	21
2	Harry	19
3	Sam	20
4	Monica	22

Constructing DataFrame from a List

```
In [89]: marks_list = [85.10, 77.80, 91.54, 88.78, 60.55]

df3 = pd.DataFrame(marks_list, columns=['Marks'])
df3
```

Out[89]:

	Marks
0	85.10
1	77.80
2	91.54
3	88.78
4	60.55

Creating DataFrame from file

```
In [90]: # Read csv file from path and store to df for create dataframe
df = pd.read_csv('nss15.csv')
df
```

Out[90]:

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	57
1	150734723	7/6/2015	83.2157	S	36	Male	White	57
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	71
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	71
4	150721694	7/4/2015	74.8813	L	88	Female	Other	62
...	...	...	...	...	...	...	...	...
334834	150739278	5/31/2015	15.0591	V	7	Male	NaN	59
334835	150733393	7/11/2015	5.6748	C	3	Female	Black	68
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN	71
334837	150823002	8/8/2015	97.9239	M	38	Female	White	59
334838	150723074	6/20/2015	49.2646	M	5	Female	White	57

334839 rows × 12 columns



## Viewing DataFrame information

(.shape, .head, .tail, .info, select column, .unique, .describe, select row with .loc and .iloc)

Check simple information

```
In [91]: # Check dimension by .shape
df.shape
```

Out[91]: (334839, 12)

```
In [92]: # Display the first 5 rows by default
df.head()
```

Out[92]:


	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	body
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	57	
1	150734723	7/6/2015	83.2157	S	36	Male	White	57	
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	71	
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	71	
4	150721694	7/4/2015	74.8813	L	88	Female	Other	62	



```
In [93]: # Display the first 3 rows
df.head(3)
```

```
Out[93]:
```


	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	body
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	57	
1	150734723	7/6/2015	83.2157	S	36	Male	White	57	
2	150817487	8/2/2015	74.8813	L	20	Female	NaN	71	



```
In [94]: # Display the last 5 rows by default
df.tail()
```

```
Out[94]:
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	
334834	150739278	5/31/2015	15.0591	V	7	Male	NaN	59	
334835	150733393	7/11/2015	5.6748	C	3	Female	Black	68	
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN	71	
334837	150823002	8/8/2015	97.9239	M	38	Female	White	59	
334838	150723074	6/20/2015	49.2646	M	5	Female	White	57	



```
In [95]: # Overview information of dataframe
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 334839 entries, 0 to 334838
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   caseNumber      334839 non-null  int64
1   treatmentDate   334839 non-null  object
2   statWeight      334839 non-null  float64
3   stratum         334839 non-null  object
4   age             334839 non-null  int64
5   sex             334837 non-null  object
6   race            205014 non-null  object
7   diagnosis       334839 non-null  int64
8   bodyPart        334839 non-null  int64
9   disposition     334839 non-null  int64
10  location         334839 non-null  int64
11  product         334839 non-null  int64
dtypes: float64(1), int64(7), object(4)
memory usage: 30.7+ MB
```

Select column, multiple column, with condition

```
In [96]: df.columns
```

```
Out[96]: Index(['caseNumber', 'treatmentDate', 'statWeight', 'stratum', 'age', 'sex',  
              'race', 'diagnosis', 'bodyPart', 'disposition', 'location', 'product'],  
              dtype='object')
```

```
In [97]: # column slicing  
df['age']
```

```
Out[97]: 0          5  
         1         36  
         2         20  
         3         61  
         4         88  
         ..  
334834    7  
334835    3  
334836   38  
334837   38  
334838    5  
Name: age, Length: 334839, dtype: int64
```

```
In [98]: df.age
```

```
Out[98]: 0          5  
         1         36  
         2         20  
         3         61  
         4         88  
         ..  
334834    7  
334835    3  
334836   38  
334837   38  
334838    5  
Name: age, Length: 334839, dtype: int64
```

Viewing the unique value

```
In [99]: df.race.unique()
```

```
Out[99]: array([nan, 'White', 'Other', 'Black', 'Asian', 'American Indian'],  
              dtype=object)
```

Null values

```
In [100]: df.isnull().sum()
```

```
Out[100...] caseNumber      0
            treatmentDate  0
            statWeight     0
            stratum        0
            age            0
            sex            2
            race           129825
            diagnosis      0
            bodyPart       0
            disposition    0
            location       0
            product        0
            dtype: int64
```

Describe

```
In [101...] df['age'].describe()
```

```
Out[101...] count      334839.000000
            mean        31.385451
            std         26.105098
            min         0.000000
            25%         10.000000
            50%         23.000000
            75%         51.000000
            max         107.000000
            Name: age, dtype: float64
```

## Slicing dataframe

```
In [102...] # select multiple column
            df[['treatmentDate', 'statWeight', 'age', 'sex']]
```

Out[102...

	treatmentDate	statWeight	age	sex
0	7/11/2015	15.7762	5	Male
1	7/6/2015	83.2157	36	Male
2	8/2/2015	74.8813	20	Female
3	6/26/2015	15.7762	61	Male
4	7/4/2015	74.8813	88	Female
...	...	...	...	...
334834	5/31/2015	15.0591	7	Male
334835	7/11/2015	5.6748	3	Female
334836	7/24/2015	15.7762	38	Male
334837	8/8/2015	97.9239	38	Female
334838	6/20/2015	49.2646	5	Female

334839 rows × 4 columns

In [103...

```
#select by condition
df[df['sex'] == 'Male']
```

Out[103...

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	b
0	150733174	7/11/2015	15.7762	V	5	Male	NaN	57	
1	150734723	7/6/2015	83.2157	S	36	Male	White	57	
3	150717776	6/26/2015	15.7762	V	61	Male	NaN	71	
6	150713483	6/8/2015	15.7762	V	25	Male	Black	51	
7	150704114	6/14/2015	83.2157	S	53	Male	White	57	
...	...	...	...	...	...	...	...	...	...
334824	150607827	5/27/2015	5.6748	C	1	Male	White	71	
334825	150600190	5/28/2015	80.8381	S	5	Male	NaN	56	
334833	150747217	7/24/2015	83.2157	S	2	Male	NaN	62	
334834	150739278	5/31/2015	15.0591	V	7	Male	NaN	59	
334836	150819286	7/24/2015	15.7762	V	38	Male	NaN	71	

182501 rows × 12 columns



In [104...

```
# select by multiple condition
df[(df['sex'] == 'Male') & (df['age'] > 80)]
```



Out[104...

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	b
<b>8</b>	150736558	7/16/2015	83.2157	S	98	Male	Black	59	
<b>63</b>	150418623	1/12/2015	15.0591	V	97	Male	Other	62	
<b>97</b>	150700375	6/28/2015	83.2157	S	85	Male	NaN	59	
<b>131</b>	150940801	9/14/2015	15.7762	V	96	Male	NaN	62	
<b>177</b>	160110774	12/19/2015	85.7374	S	81	Male	White	59	
...	...	...	...	...	...	...	...	...	
<b>334616</b>	160104368	12/30/2015	74.8813	L	86	Male	Other	71	
<b>334677</b>	151115099	11/4/2015	16.5650	V	83	Male	NaN	63	
<b>334699</b>	150633387	5/29/2015	74.8813	L	84	Male	NaN	53	
<b>334701</b>	150515945	4/27/2015	97.9239	M	86	Male	NaN	57	
<b>334785</b>	150733286	7/11/2015	15.7762	V	86	Male	White	71	

6379 rows × 12 columns



Select row with .iloc

In [105...

```
# row slicing
df.iloc[10:15]
```

Out[105...

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bod
<b>10</b>	150734952	7/4/2015	15.7762	V	20	Male	Black	59	
<b>11</b>	150821622	7/20/2015	83.2157	S	20	Female	White	57	
<b>12</b>	150713631	7/4/2015	15.7762	V	11	Male	NaN	60	
<b>13</b>	150666343	6/27/2015	15.7762	V	26	Female	White	62	
<b>14</b>	150748843	7/16/2015	37.6645	L	33	Male	Asian	53	



In [106...

```
# column slicing
df.iloc[:, [0,1,2,3,4]]
```

Out[106...

	caseNumber	treatmentDate	statWeight	stratum	age
0	150733174	7/11/2015	15.7762	V	5
1	150734723	7/6/2015	83.2157	S	36
2	150817487	8/2/2015	74.8813	L	20
3	150717776	6/26/2015	15.7762	V	61
4	150721694	7/4/2015	74.8813	L	88
...	...	...	...	...	...
334834	150739278	5/31/2015	15.0591	V	7
334835	150733393	7/11/2015	5.6748	C	3
334836	150819286	7/24/2015	15.7762	V	38
334837	150823002	8/8/2015	97.9239	M	38
334838	150723074	6/20/2015	49.2646	M	5

334839 rows × 5 columns

Select column and row with .loc

In [107...

```
# select column and row by .loc
df.loc[:6, 'treatmentDate': 'diagnosis']
```

Out[107...

	treatmentDate	statWeight	stratum	age	sex	race	diagnosis
0	7/11/2015	15.7762	V	5	Male	NaN	57
1	7/6/2015	83.2157	S	36	Male	White	57
2	8/2/2015	74.8813	L	20	Female	NaN	71
3	6/26/2015	15.7762	V	61	Male	NaN	71
4	7/4/2015	74.8813	L	88	Female	Other	62
5	7/2/2015	5.6748	C	1	Female	White	71
6	6/8/2015	15.7762	V	25	Male	Black	51

In [108...

```
# select row by condition
df.loc[df['age'] > 80, ['treatmentDate', 'age']]
```

Out[108...

	treatmentDate	age
4	7/4/2015	88
8	7/16/2015	98
39	5/3/2015	88
46	4/15/2015	91
63	1/12/2015	97
...	...	...
334701	4/27/2015	86
334784	7/7/2015	82
334785	7/11/2015	86
334815	10/28/2015	85
334819	1/13/2015	85

20422 rows × 2 columns

[Q6] What is the difference between .iloc and .loc?

Ans: .iloc is index-based query while .loc is label-based query

for example df.loc[:6,'treatmentDate':'diagnosis'] equals to df.iloc[:7,1:8] loc using 'treatmentDate':'diagnosis' label to query but iloc using 1:8 to query

also loc can query using condition such as df.loc[df['age']>80, ['treatmentDate', 'age']]

In [121...

```
df.loc[:6,'treatmentDate':'diagnosis']
```

Out[121...

	treatmentDate	statWeight	stratum	age	sex	race	diagnosis
0	7/11/2015	15.7762	V	5	Male	NaN	57
1	7/6/2015	83.2157	S	36	Male	White	57
2	8/2/2015	74.8813	L	20	Female	NaN	71
3	6/26/2015	15.7762	V	61	Male	NaN	71
4	7/4/2015	74.8813	L	88	Female	Other	62
5	7/2/2015	5.6748	C	1	Female	White	71
6	6/8/2015	15.7762	V	25	Male	Black	51

In [125...

```
df.iloc[:7,1:8]
```

Out[125...

	treatmentDate	statWeight	stratum	age	sex	race	diagnosis
0	7/11/2015	15.7762	V	5	Male	NaN	57
1	7/6/2015	83.2157	S	36	Male	White	57
2	8/2/2015	74.8813	L	20	Female	NaN	71
3	6/26/2015	15.7762	V	61	Male	NaN	71
4	7/4/2015	74.8813	L	88	Female	Other	62
5	7/2/2015	5.6748	C	1	Female	White	71
6	6/8/2015	15.7762	V	25	Male	Black	51

### Basic aggregations

In [109...

```
# count elements in column via pandas Series method
df['sex'].value_counts()
```

Out[109...

```
sex
Male      182501
Female    152336
Name: count, dtype: int64
```

In [110...

```
# count elements in column via pandas function
pd.crosstab(df['sex'], columns='N')
```

Out[110...

col_0	N
sex	
Female	152336
Male	182501

In [111...

```
summary_by_sex = df.groupby('sex').agg({
    'age': 'mean',
    'statWeight': 'sum',
    'caseNumber': 'count',
    'stratum': lambda x: x.mode()[0]
}).rename(columns={'caseNumber': 'total_cases'})

print(summary_by_sex)
```

	age	statWeight	total_cases	stratum
sex				
Female	35.585699	6.208773e+06	152336	V
Male	27.879792	6.964776e+06	182501	V

In [112...

```
df['age'].apply('mean')
```

Out[112...

```
np.float64(31.38545091820248)
```

```
In [133... df['race'].fillna("").apply(len)
```

```
Out[133... 0      0
1      5
2      0
3      0
4      5
..
334834  0
334835  5
334836  0
334837  5
334838  5
Name: race, Length: 334839, dtype: int64
```

```
In [ ]: df['race'].apply(len) #data have NaN so len(Nan) is error because python treat data
```

-----  
TypeError

Traceback (most recent call last)

Cell In[132], line 1

```
----> 1 df[ ].apply(len)
```

File c:\Users\win25\Desktop\Desktop\work\CPE232 Data models\w1\.venv\Lib\site-packages\pandas\core\series.py:4943, in Series.apply(self, func, convert\_dtype, args, by\_row, \*\*kwargs)

```
    4808 def apply(
    4809     self,
    4810     func: AggFuncType,
    (...) 4815     **kwargs,
    4816 ) -> DataFrame | Series:
    4817     """
    4818     Invoke function on values of Series.
    4819
    (...) 4934     dtype: float64
    4935     """
    4936     return SeriesApply(
    4937         self,
    4938         func,
    4939         convert_dtype=convert_dtype,
    4940         by_row=by_row,
    4941         args=args,
    4942         kwargs=kwargs,
-> 4943     ).apply()
```

File c:\Users\win25\Desktop\Desktop\work\CPE232 Data models\w1\.venv\Lib\site-packages\pandas\core\apply.py:1422, in SeriesApply.apply(self)

```
    1419     return self.apply_compat()
    1421 # self.func is Callable
-> 1422 return self.apply_standard()
```

File c:\Users\win25\Desktop\Desktop\work\CPE232 Data models\w1\.venv\Lib\site-packages\pandas\core\apply.py:1502, in SeriesApply.apply\_standard(self)

```
    1496 # row-wise access
    1497 # apply doesn't have a `na_action` keyword and for backward compat reasons
    1498 # we need to give `na_action="ignore"` for categorical data.
    1499 # TODO: remove the `na_action="ignore"` when that default has been changed i
n
    1500 # Categorical (GH51645).
    1501 action = "ignore" if isinstance(obj.dtype, CategoricalDtype) else None
-> 1502 mapped = obj._map_values(
    1503     mapper=curried, na_action=action, convert=self.convert_dtype
    1504 )
    1506 if len(mapped) and isinstance(mapped[0], ABCSeries):
    1507     # GH#43986 Need to do list(mapped) in order to get treated as nested
    1508     # See also GH#25959 regarding EA support
    1509     return obj._constructor_expanddim(list(mapped), index=obj.index)
```

File c:\Users\win25\Desktop\Desktop\work\CPE232 Data models\w1\.venv\Lib\site-packages\pandas\core\base.py:925, in IndexOpsMixin.\_map\_values(self, mapper, na\_action, convert)

```
    922 if isinstance(arr, ExtensionArray):
    923     return arr.map(mapper, na_action=na_action)
--> 925 return algorithms.map_array(arr, mapper, na_action=na_action, convert=conver
```

t)

```
File c:\Users\win25\Desktop\Desktop\work\CPE232 Data models\w1\.venv\Lib\site-packages\pandas\core\algorithms.py:1743, in map_array(arr, mapper, na_action, convert)
    1741 values = arr.astype(object, copy=False)
    1742 if na_action is None:
-> 1743     return lib.map_infer(values, mapper, convert=convert)
    1744 else:
    1745     return lib.map_infer_mask(
    1746         values, mapper, mask=isna(values).view(np.uint8), convert=convert
    1747     )

File pandas/_libs/lib.pyx:2999, in pandas._libs.lib.map_infer()

TypeError: object of type 'float' has no len()
```

[Q7] In the above cell, Explain why pandas .agg()/apply() method accepts 'len' and 'lambda' without quotation marks.

Ans: Because len and lambda is function its mean we pass the function to make pandas calling it themself for each row

## 7.3 Matplotlib

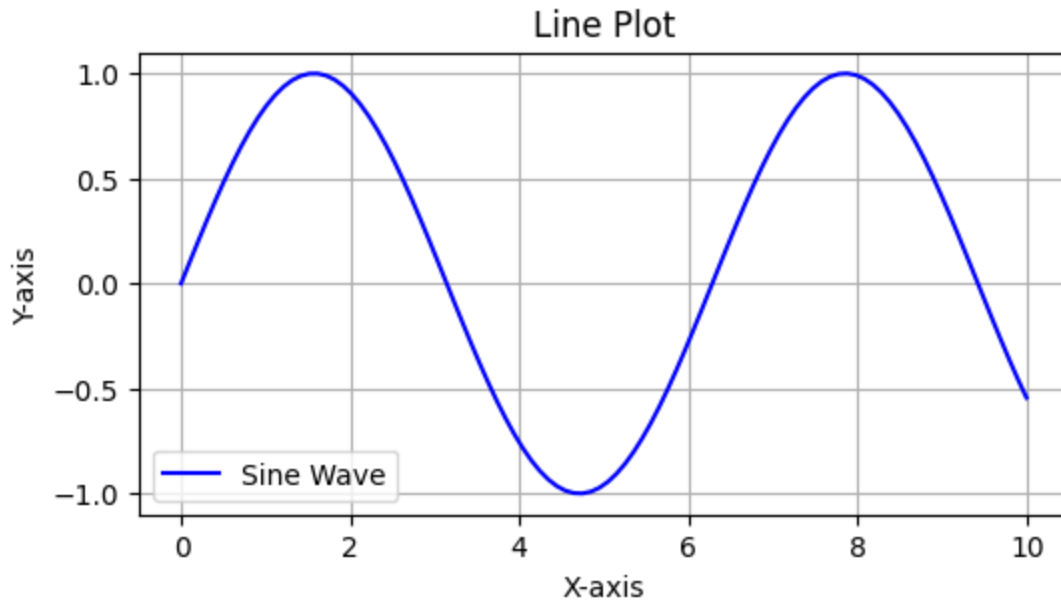
```
In [115... import matplotlib.pyplot as plt
import numpy as np
```

```
In [116... x = np.linspace(0, 10, 100)
y = np.sin(x)

# Line plot

plt.figure(figsize=(6, 3))
plt.plot(x, y, label='Sine Wave', color='blue')
plt.title('Line Plot')
plt.xlabel('X-axis')
plt.ylabel('Y-axis')
plt.legend()
plt.grid(True)

plt.savefig("sine.png")
plt.show()
```



In [117...

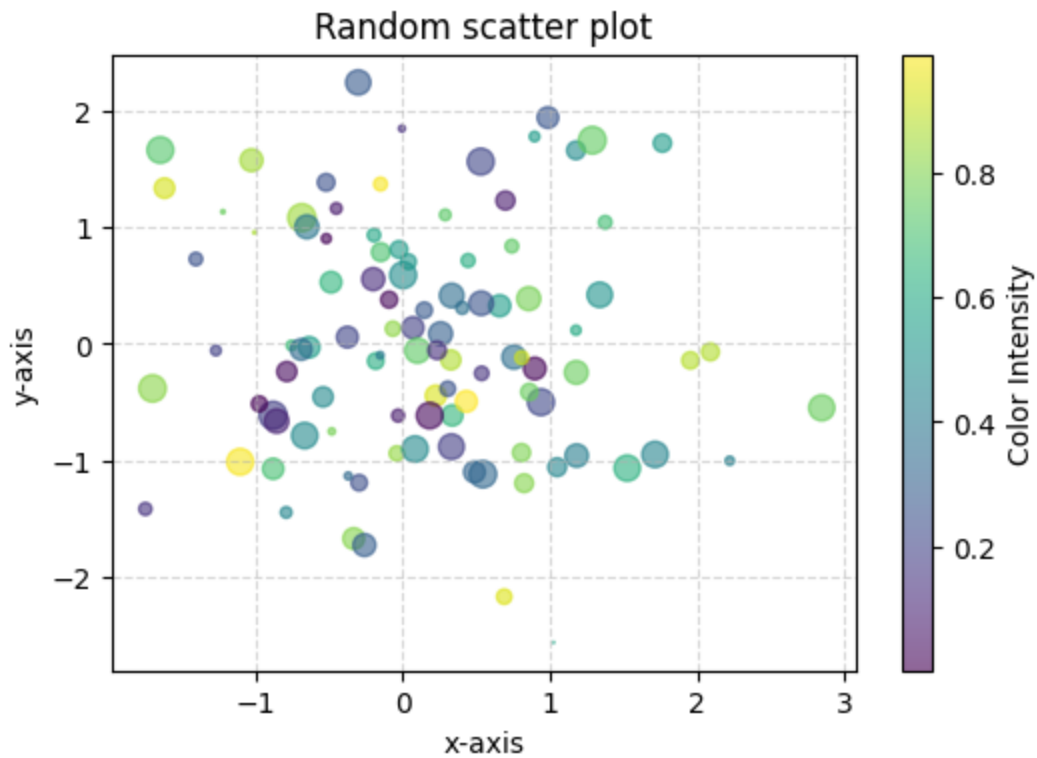
```
x = np.random.randn(100)
y = np.random.randn(100)
colors = np.random.rand(100) # Random colors
sizes = np.random.rand(100) * 100 # Random sizes

# Scatter plot
plt.figure(figsize=(6, 4))
scatter = plt.scatter(x, y, c=colors, s=sizes, alpha=0.6, cmap='viridis')

# Add styling and labels
plt.title('Random scatter plot')
plt.xlabel('x-axis')
plt.ylabel('y-axis')
plt.colorbar(scatter, label='Color Intensity')
plt.grid(True, linestyle='--', alpha=0.5)

# Show the plot
plt.show()
```





```
In [118... plt.figure(figsize=(5, 3))

# bar plot

df['sex'].value_counts().plot(kind='bar', color=['salmon', 'lightblue'])
plt.title('Gender Distribution')
plt.xlabel('Sex')
plt.ylabel('Count')
plt.xticks(rotation=0)
```

```
Out[118... (array([0, 1]), [Text(0, 0, 'Male'), Text(1, 0, 'Female')])
```

