

✓ Lab 1: Basic Python Programming


CPE232 Data Models

Owner : 66070501043 - Phoorin Chinphuad

✓ [1] Variable

✓ 1.1 Number Variable


```
1 num = 100 #integer variable
2 num2 = 12.5 #float variable
3 print(num)
4 print(num2)
5
6 print(num + num2)    #addition
7 print(num - num2)    #subtraction
8 print(num * num2)    #multiplication
9 print( num / num2)   #division
```



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

✓ 1.2 String Variable


```
1 #string variable
2 string = "Data Models"
3 print(string) #print complete string
4
5 print("Hello " + string)    #print concatenated string
6 print(string[0])           #print first character of the string
7 print(string[:4])          #print first to 4th character of the string
8 print(string[5:])          #print 6th to last character of the string
9 print(string[1:4])         #print 2nd to 4th character of the string
10 print(string * 2)         #print string 2 time
11
```



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

✓ 1.3 Boolean Variable

```
1 #boolean variable
2 boolean = True
3 boolean2 = False
4
5 print(boolean)            #print boolean variable
6 print(not boolean)        #print opposite of boolean variable
7 print(boolean and boolean2) #print boolean and boolean2
8 print(boolean or boolean2) #print boolean or boolean2
```



```
True
False
False
True
```

✓ 1.4 List Variable

```

1 #list variable
2 list = ["Data",20,123.23,40,50]
3 another_list = ["Models",60]
4
5 print(list)           #print complete list
6 print(list[0])        #print first element of the list
7 print(list[1:3])      #print 2nd to 3rd element of the list
8 print(list[2:])       #print 3rd to last element of the list
9 print(another_list)   #print complete another_list
10 print(another_list * 2) #print another_list two times
11 print(list + another_list) #print concatenated list
12
13 list[0] = "CPE232"    #change first element of the list
14 print(list)          #print complete list
15

```

```

➦ ['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]

```

✓ 1.5 Tuple Variable

```

1 #tuple variable
2 tuple = ("Data",20,123.23,40,50)
3 another_tuple = ("Models",60)
4
5 print(tuple)           #print complete tuple
6 print(tuple[0])        #print first element of the tuple
7 print(tuple[1:3])      #print 2nd to 3rd element of the tuple
8 print(tuple[2:])       #print 3rd to last element of the tuple
9 print(tuple * 2)       #print tuple two times
10 print(tuple + another_tuple) #print concatenated tuple
11

```

```

➦ ('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)

```

```

1 tuple[0] = "CPE232"    #trying to change first element of the tuple but it cannot be changed so it gives error

```

✓ 1.6 Dictionary Variable

```

1 #dictionary variable
2 dictionary = {"name":"Alice","age":21}
3 another_dictionary = {}
4 another_dictionary["name"] = "Bob"
5 another_dictionary["age"] = 21
6
7 print(dictionary)      #print complete dictionary
8 print(dictionary["name"]) #print value for specific key
9 print(dictionary.keys()) #print all the keys
10 print(dictionary.values()) #print all the values
11 print(dictionary.items()) #print all the items
12 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}

```

✓ [2] Control Flow

✓ 2.1 IF ... ELIF ... ELSE

```
1 number = 123
2 number2 = 34
3
4 if number > number2:
5     print("number is greater thanu number2")
6 elif number < number2:
7     print("number is less than number2")
8 else:
9     print("number is equal to number2")
10
```

➞ number is greater thanu number2

✓ [3] Loop

✓ 3.1 For Loop

```
1 #for loops
2 for num in range(0,10):
3     print(num)
```

➞ 0
1
2
3
4
5
6
7
8
9

```
1 #for loop with list
2
3 list = ["Alice","Bob","Charlie","Daisy"]
4
5 for name in list:
6     print(name)
```

➞ Alice
Bob
Charlie
Daisy

```
1 #continue in for loop
2
3 list = [1,23,7,"hello",True,1123,43,23,12]
4
5 for element in list:
6     if type(element) != int:
7         continue
8     print(element)
```

➞ 1
23
7
1123
43
23
12

```
1 #break in for loop
2
3 list = [1,23,7,"hello",True,1123,43,23,12]
4
5 for element in list:
6     if type(element) != int:
7         break
8     print(element)
```

```
1
23
7
```

3.2 While loop

```
1 #while loop
2
3 list = ["Alice","Bob","Charlie","Daisy"]
4 count = 0
5
6 while count < len(list):
7     print(list[count])
8     count += 1
```

```
Alice
Bob
Charlie
Daisy
```

```
1 #continue in while loop
2
3 list = [1,23,7,"hello",True,1123,43,23,12]
4 count = 0
5
6 while count < len(list):
7     if type(list[count]) != int:
8         count += 1
9         continue
10    print(list[count])
11    count += 1
```

```
1
23
7
1123
43
23
12
```

```
1 #break in while loop
2
3 list = [1,23,7,"hello",True,1123,43,23,12]
4 count = 0
5
6 while count < len(list):
7     if type(list[count]) != int:
8         break
9     print(list[count])
10    count += 1
```

```
1
23
7
```

[4] Function

```
1 #define function
2 def function_name (arg1, arg2):
3     return arg1 + arg2
4
5 #calling function
6 function_name(1,2)
```

```
3
```

```
1 #define function with default argument
2 def function_with_default_arg(arg1, arg2 = 10, arg3 = 20 , arg4 = 30):
3     return arg1 + arg2 + arg3 + arg4
4
5 result_1 = function_with_default_arg(1)
6 result_2 = function_with_default_arg(1,2,5)
```

```

7 result_3 = function_with_default_arg(1,2,5,10)
8
9 print(result_1)
10 print(result_2)
11 print(result_3)

```

```

↔ 61
   38
   18

```

```

1 #multiple agument
2 def function_with_multiple_arg(*args):
3     print(args)
4     print(type(args))
5     sum = 0
6     for num in args:
7         sum += num
8
9     return sum
10
11 function_with_multiple_arg(1,2,3,4,5)

```

```

↔ (1, 2, 3, 4, 5)
   <class 'tuple'>
   15

```

```

1 #lambda function
2 lambda_function = lambda arg1, arg2: arg1 + arg2
3
4 print(lambda_function(1,2))

```

```

↔ 3

```

✓ [5] File Handling

✓ 5.1 Text File

```

1 with open("test.txt","w") as file:
2     file.write("Hello World")

```

```

1 with open("test.txt","r") as file:
2     print(file.read())

```

```

↔ Hello World

```

✓ 5.2 CSV File

```

1 import csv
2
3 with open("test.csv","w",newline='') as file:
4     writer = csv.writer(file)
5     writer.writerow(["Name","Surname"])
6     writer.writerow(["Alice","Johnson"])
7     writer.writerow(["Bob","Smith"])

```

```

1 import csv
2
3 with open("test.csv","r") as file:
4     reader = csv.reader(file)
5     for row in reader:
6         print(row)

```

```

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

```

✓ [4] Libraries

✓ 4.1 Numpy

✓ import numpy library

```
1 import numpy as np
```

✓ ndarray initialization

Construct using python list

```
1 # 1d ndarray from 1d python list
2 list_a1=[1,2,3.5]
3 arr_a1=np.array(list_a1)
4 arr_a1
```

```
↔ array([1. , 2. , 3.5])
```

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

```
↔ array([[1, 2],
        [3, 4],
        [5, 6]])
```

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

```
↔ array([[[1, 2],
         [2, 3]],
        [[3, 4],
         [4, 5]]])
```

or construct using some numpy classes and functions

```
1 np.zeros(5)
```

```
↔ array([0., 0., 0., 0., 0.])
```

```
1 np.ones((3,4),dtype=float)
```

```
↔ array([[1., 1., 1., 1.],
        [1., 1., 1., 1.],
        [1., 1., 1., 1.]])
```

```
1 np.full((4,),999)
```

```
↔ array([999, 999, 999, 999])
```

```
1 np.arange(3,10,2)
```

```
↔ array([3, 5, 7, 9])
```

```
1 np.linspace(10,15,11)
```

```
↔ array([10. , 10.5, 11. , 11.5, 12. , 12.5, 13. , 13.5, 14. , 14.5, 15. ])
```

```
1 np.random.choice(['a','b'],9)
```

```
↔ array(['a', 'a', 'a', 'b', 'a', 'b', 'a', 'a', 'b'], dtype='<U1')
```

```
1 np.random.randn(10)
```

```
↔ array([-0.82222123,  0.7728091 ,  0.07231017,  0.58226338,  0.25735609,
         0.96683725,  1.56468755, -0.36298275, -0.69294433, -1.30699885])
```

▼ ndarray properties

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

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

```
1 arr_a.ndim
```

```
↵ 2
```

```
1 arr_a.shape
```

```
↵ (3, 4)
```

```
1 arr_a.dtype
```

```
↵ dtype('int64')
```

```
1 arr_a.size
```

```
↵ 12
```

▼ Reshaping & Modification

from this original ndarray

```
1 arr_a
```

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

try to convert into 3D array

```
1 arr_a.reshape((2,2,3))
```

```
↵ 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

```
1 arr_a.reshape((-1,6))
```

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

Would you like to try this?

```
1 arr_a.reshape((-1,5))
```

```
↵ -----
ValueError                                Traceback (most recent call last)
<ipython-input-127-286d5aa6424c> in <cell line: 0>()
----> 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: **To reshape the array into a shape with 5 columns, the total number of elements must be divisible by 5. However, 12 is not divisible by 5.**

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

```
1 np.append(arr_a,13)
```

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

```
1 np.append(arr_a,arr_a[0])
```

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

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

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

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

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

```
1 np.concatenate([arr_a,arr_a])
```

```
↔ 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]])
```

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

```
↔ 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

```
1 arr_a
```

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

try to access all element at the first row

```
1 arr_a[1]
```

```
↔ array([5, 6, 7, 8])
```

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

```
1 arr_a[1][2]
```

```
↔ 7
```

```
1 arr_a[1,2]
```

```
↔ 7
```

Next, try to access all element start from 1th in the first row

```
1 arr_a[1,1:]
```

```
↔ array([6, 7, 8])
```



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

```
↵ array([[2, 3, 4],  
        [6, 7, 8]])
```

sometimes you may specify some row number using list within indexing

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

```
↵ array([[ 6,  7,  8],  
        [10, 11, 12],  
        [ 6,  7,  8]])
```

✓ Boolean slicing

based on this original array

```
1 arr_a
```

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

try to filter all elements which more than 5

```
1 arr_a>5
```

```
↵ 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

```
1 (arr_a>5)&(arr_a<10)
```

```
↵ array([[False, False, False, False],  
        [False,  True,  True,  True],  
        [ True, False, False, False]])
```

Run the cell below and answer a question.

```
1 arr_a[(arr_a>5)&(arr_a<10)]
```

```
↵ array([6, 7, 8, 9])
```

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

Ans: **The output consists of only the elements of arr_a that correspond to True in the combined boolean mask (arr_a > 5) & (arr_a < 10).**

Try running the cell below.

```
1 arr_a[(arr_a>5) and (arr_a<10)]
```

```
↵ -----  
ValueError                                Traceback (most recent call last)  
<ipython-input-149-78eb1746bbfd> in <cell line: 0>()  
----> 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: **The and operator in Python is a logical operator that works with single boolean values. While arr_a > 5 and arr_a < 10 are both boolean arrays, not single boolean values.**

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

Ans: **To perform element-wise logical AND on NumPy arrays, try using the & operator instead of and >> arr_a[(arr_a>5) & (arr_a<10)]**

▼ Basic operations

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

```
↕ array([[1, 2, 3, 4],
         [1, 2, 3, 4],
         [1, 2, 3, 4]])
```

This is some operations for only 1 array

```
1 np.sqrt(arr_b)
```

```
↕ 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

```
1 arr_a-arr_b
```

```
↕ array([[0, 0, 0, 0],
         [4, 4, 4, 4],
         [8, 8, 8, 8]])
```

```
1 np.add(arr_a,arr_b)
```

```
↕ array([[ 2,  4,  6,  8],
         [ 6,  8, 10, 12],
         [10, 12, 14, 16]])
```

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

```
1 arr_a*3
```

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

```
1 1+arr_a**2
```

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

Try to play with 2 arrays with different shape

```
1 arr_c=np.array([1,2,3])
2 arr_d=np.array([[3],[5],[8]])
3 print(arr_c)
4 print(arr_d)
```

```
↕ [1 2 3]
   [[3]
    [5]
    [8]]
```

```
1 arr_c-arr_d
```

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

▼ Basic aggregations

```
1 arr_a
```

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

```
1 arr_a.sum()
```

```
↕ 78
```

```
1 arr_a.mean()
```

```
↕ 6.5
```

```
1 arr_a.min()
```

```
↕ 1
```

```
1 arr_a.max()
```

```
↕ 12
```

```
1 arr_a.std() #standard division
```

```
↕ 3.452052529534663
```

▼ ndarray axis

```
1 arr_a
```

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

```
1 arr_a.sum(axis=0)
```

```
↕ array([15, 18, 21, 24])
```

```
1 arr_a.sum(axis=1)
```

```
↕ 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: **axis=0 : is performed column-wise. and axis=1 : is performed row-wise.**

▼ 4.2 Pandas

▼ Series

```
1 import pandas as pd
```

```
2 import numpy as np
```


```
1 pd.Series(np.random.randn(6))
```

```
↕
```

	0
0	0.637273
1	1.327869
2	-0.620961
3	-0.203078
4	0.378510
5	0.036953

dtype: float64

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



	0
a	-0.087000
b	-1.044680
c	-0.597356
d	0.285617
e	0.990374
f	0.968748

dtype: float64

Constructing Dataframe

Constructing DataFrame from a dictionary

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

```
1 df = pd.DataFrame(data=d)
2 df
```



	col1	col2
0	1	3
1	2	4

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

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



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

Constructing DataFrame from a List

```
1 marks_list = [85.10, 77.80, 91.54, 88.78, 60.55]
```

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




	Marks
0	85.10
1	77.80
2	91.54
3	88.78
4	60.55

Creating DataFrame from file

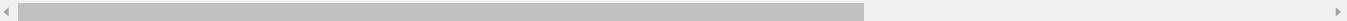
```
1 # Read csv file from path and store to df for create dataframe
2 df = pd.read_csv('nss15.csv')
```

```
1 df
```



	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
0	150733174	7/11/2015	15.7762	V	5.0	Male	NaN	57.0	33.0	1.0	9.0	1267.0
1	150734723	7/6/2015	83.2157	S	36.0	Male	White	57.0	34.0	1.0	1.0	1439.0
2	150817487	8/2/2015	74.8813	L	20.0	Female	NaN	71.0	94.0	1.0	0.0	3274.0
3	150717776	6/26/2015	15.7762	V	61.0	Male	NaN	71.0	35.0	1.0	0.0	611.0
4	150721694	7/4/2015	74.8813	L	88.0	Female	Other	62.0	75.0	1.0	0.0	1893.0
...
234085	150968928	9/22/2015	15.7762	V	23.0	Male	Black	64.0	92.0	1.0	1.0	1141.0
234086	150965850	9/24/2015	83.2157	S	37.0	Female	NaN	64.0	31.0	1.0	1.0	4014.0
234087	150971407	9/26/2015	5.6748	C	13.0	Male	White	71.0	79.0	1.0	9.0	1211.0
234088	151026924	10/6/2015	5.6748	C	1.0	Male	White	53.0	75.0	1.0	1.0	4057.0
234089	15100638	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

234090 rows × 12 columns



Viewing DataFrame information


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

Check simple information


- 1 # Check dimension by .shape
- 2 df.shape

 (234090, 12)


- 1 # Display the first 5 rows by default
- 2 df.head()




	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
0	150733174	7/11/2015	15.7762	V	5.0	Male	NaN	57.0	33.0	1.0	9.0	1267.0
1	150734723	7/6/2015	83.2157	S	36.0	Male	White	57.0	34.0	1.0	1.0	1439.0
2	150817487	8/2/2015	74.8813	L	20.0	Female	NaN	71.0	94.0	1.0	0.0	3274.0
3	150717776	6/26/2015	15.7762	V	61.0	Male	NaN	71.0	35.0	1.0	0.0	611.0
4	150721694	7/4/2015	74.8813	L	88.0	Female	Other	62.0	75.0	1.0	0.0	1893.0




- 1 # Display the first 3 rows
- 2 df.head(3)




	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
0	150733174	7/11/2015	15.7762	V	5.0	Male	NaN	57.0	33.0	1.0	9.0	1267.0
1	150734723	7/6/2015	83.2157	S	36.0	Male	White	57.0	34.0	1.0	1.0	1439.0
2	150817487	8/2/2015	74.8813	L	20.0	Female	NaN	71.0	94.0	1.0	0.0	3274.0



- 1 # Display the last 5 rows by default
- 2 df.tail()



	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
234085	150968928	9/22/2015	15.7762	V	23.0	Male	Black	64.0	92.0	1.0	1.0	1141.0
234086	150965850	9/24/2015	83.2157	S	37.0	Female	NaN	64.0	31.0	1.0	1.0	4014.0
234087	150971407	9/26/2015	5.6748	C	13.0	Male	White	71.0	79.0	1.0	9.0	1211.0
234088	151026924	10/6/2015	5.6748	C	1.0	Male	White	53.0	75.0	1.0	1.0	4057.0
234089	15100638	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN



- 1 # Overview information of dataframe
- 2 df.info()

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 234090 entries, 0 to 234089
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   caseNumber             234090 non-null int64
1   treatmentDate          234089 non-null object
2   statWeight             234089 non-null float64
3   stratum                234089 non-null object
4   age                    234089 non-null float64
5   sex                    234087 non-null object
6   race                   143318 non-null object
7   diagnosis              234089 non-null float64
8   bodyPart               234089 non-null float64
9   disposition            234089 non-null float64
10  location               234089 non-null float64
11  product                234089 non-null float64
dtypes: float64(7), int64(1), object(4)
memory usage: 21.4+ MB

```

Select column, multiple column, with condition

1 df.columns

```

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

```

1 #select single column

2 df['age']

```

      age
0      5.0
1     36.0
2     20.0
3     61.0
4     88.0
...     ...
234085  23.0
234086  37.0
234087  13.0
234088   1.0
234089  NaN
234090 rows x 1 columns

```

dtype: float64

1 df.age


```

      age
0      5.0
1     36.0
2     20.0
3     61.0
4     88.0
...     ...
234085  23.0
234086  37.0
234087  13.0
234088   1.0
234089  NaN
234090 rows x 1 columns

```

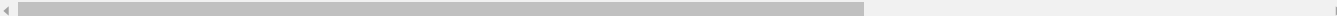
dtype: float64

```
1 #select multiple column
2 df[['treatmentDate', 'statWeight', 'age', 'sex']]
```




	treatmentDate	statWeight	age	sex
0	7/11/2015	15.7762	5.0	Male
1	7/6/2015	83.2157	36.0	Male
2	8/2/2015	74.8813	20.0	Female
3	6/26/2015	15.7762	61.0	Male
4	7/4/2015	74.8813	88.0	Female
...
234085	9/22/2015	15.7762	23.0	Male
234086	9/24/2015	83.2157	37.0	Female
234087	9/26/2015	5.6748	13.0	Male
234088	10/6/2015	5.6748	1.0	Male
234089	NaN	NaN	NaN	NaN

234090 rows × 4 columns



Viewing the unique value

```
1 df.race.unique()
```

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

Describe

```
1 df['age'].describe()
```



	age
count	234089.000000
mean	31.323274
std	26.077750
min	0.000000
25%	10.000000
50%	23.000000
75%	51.000000
max	106.000000

dtype: float64



Select row with condition

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

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
0	150733174	7/11/2015	15.7762	V	5.0	Male	NaN	57.0	33.0	1.0	9.0	1267.0
1	150734723	7/6/2015	83.2157	S	36.0	Male	White	57.0	34.0	1.0	1.0	1439.0
3	150717776	6/26/2015	15.7762	V	61.0	Male	NaN	71.0	35.0	1.0	0.0	611.0
6	150713483	6/8/2015	15.7762	V	25.0	Male	Black	51.0	33.0	4.0	9.0	1138.0
7	150704114	6/14/2015	83.2157	S	53.0	Male	White	57.0	30.0	1.0	0.0	5040.0
...
234081	150953450	9/19/2015	5.6748	C	7.0	Male	White	57.0	80.0	4.0	0.0	3286.0
234084	150906288	8/27/2015	5.6748	C	14.0	Male	White	57.0	33.0	1.0	9.0	1329.0
234085	150968928	9/22/2015	15.7762	V	23.0	Male	Black	64.0	92.0	1.0	1.0	1141.0
234087	150971407	9/26/2015	5.6748	C	13.0	Male	White	71.0	79.0	1.0	9.0	1211.0
234088	151026924	10/6/2015	5.6748	C	1.0	Male	White	53.0	75.0	1.0	1.0	4057.0

127729 rows × 12 columns

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

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
8	150736558	7/16/2015	83.2157	S	98.0	Male	Black	59.0	76.0	1.0	1.0	1807.0
63	150418623	1/12/2015	15.0591	V	97.0	Male	Other	62.0	75.0	4.0	1.0	4076.0
97	150700375	6/28/2015	83.2157	S	85.0	Male	NaN	59.0	92.0	1.0	0.0	478.0
131	150940801	9/14/2015	15.7762	V	96.0	Male	NaN	62.0	75.0	1.0	5.0	1807.0
177	160110774	12/19/2015	85.7374	S	81.0	Male	White	59.0	82.0	1.0	1.0	3278.0
...
233609	150117332	1/2/2015	78.5926	S	82.0	Male	NaN	59.0	76.0	4.0	9.0	3223.0
233867	150822235	7/25/2015	15.7762	V	93.0	Male	NaN	57.0	79.0	4.0	5.0	1679.0
233932	151029869	6/17/2015	74.8813	L	85.0	Male	NaN	71.0	36.0	4.0	1.0	4076.0
233985	151015969	10/6/2015	83.2157	S	82.0	Male	NaN	62.0	75.0	1.0	5.0	1807.0
234052	150827615	7/30/2015	15.7762	V	85.0	Male	NaN	59.0	33.0	1.0	1.0	1884.0

4407 rows × 12 columns

Select row with .iloc

```
1 # select row by .iloc
2 df.iloc[10:15]
```

	caseNumber	treatmentDate	statWeight	stratum	age	sex	race	diagnosis	bodyPart	disposition	location	product
10	150734952	7/4/2015	15.7762	V	20.0	Male	Black	59.0	82.0	1.0	1.0	1894.0
11	150821622	7/20/2015	83.2157	S	20.0	Female	White	57.0	36.0	1.0	9.0	1267.0
12	150713631	7/4/2015	15.7762	V	11.0	Male	NaN	60.0	88.0	1.0	0.0	3274.0
13	150666343	6/27/2015	15.7762	V	26.0	Female	White	62.0	75.0	1.0	1.0	1807.0
14	150748843	7/16/2015	37.6645	L	33.0	Male	Asian	53.0	93.0	1.0	1.0	4057.0

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




	caseNumber	treatmentDate	statWeight	stratum	age
0	150733174	7/11/2015	15.7762	V	5.0
1	150734723	7/6/2015	83.2157	S	36.0
2	150817487	8/2/2015	74.8813	L	20.0
3	150717776	6/26/2015	15.7762	V	61.0
4	150721694	7/4/2015	74.8813	L	88.0

Select column and row with .loc

```
234085    150968928    9/22/2015    15.7762    V    23.0
```

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



	treatmentDate	statWeight	stratum	age	sex	race	diagnosis
0	7/11/2015	15.7762	V	5.0	Male	NaN	57.0
1	7/6/2015	83.2157	S	36.0	Male	White	57.0
2	8/2/2015	74.8813	L	20.0	Female	NaN	71.0
3	6/26/2015	15.7762	V	61.0	Male	NaN	71.0
4	7/4/2015	74.8813	L	88.0	Female	Other	62.0
5	7/2/2015	5.6748	C	1.0	Female	White	71.0
6	6/8/2015	15.7762	V	25.0	Male	Black	51.0

```
1 # select row by condition
2
3 df.loc[df['age'] > 80, ['treatmentDate', 'age']]
```



	treatmentDate	age
4	7/4/2015	88.0
8	7/16/2015	98.0
39	5/3/2015	88.0
46	4/15/2015	91.0
63	1/12/2015	97.0
...
233985	10/6/2015	82.0
234047	8/4/2015	83.0
234049	6/16/2015	85.0
234052	7/30/2015	85.0
234069	8/29/2015	83.0

14171 rows × 2 columns

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

Ans: .iloc is Integer-based indexing (position-based). While .loc is Label-based indexing (name-based).

