

NUMPY , PANDAS AND MATPLOTLIB

ADVANTAGES OF PANDAS

- FAST AND EFFICIENT DATAFRAME OBJECT WITH DEFAULT AND CUSTOMIZED INDEXING.
- TOOLS FOR LOADING DATA INTO IN-MEMORY DATA OBJECTS FROM DIFFERENT FILE FORMATS.
- DATA ALIGNMENT AND INTEGRATED HANDLING OF MISSING DATA.
- RESHAPING AND PIVOTING OF DATE SETS.
- LABEL-BASED SLICING, INDEXING AND SUBSETTING OF LARGE DATA SETS.

ADVANTAGES CONTINUED

- COLUMNS FROM A DATA STRUCTURE CAN BE DELETED OR INSERTED.
- GROUP BY DATA FOR AGGREGATION AND TRANSFORMATIONS.
- HIGH PERFORMANCE MERGING AND JOINING OF DATA.
- TIME SERIES FUNCTIONALITY.

DATA TYPES

- PANDAS DEALS WITH THE FOLLOWING THREE DATA STRUCTURES:
 1. SERIES
 2. DATAFRAME
 3. PANEL
- THESE DATA STRUCTURES ARE BUILT ON TOP OF NUMPY ARRAY, WHICH MEANS THEY ARE FAST.

DIMENSION AND DESCRIPTION

- THE BEST WAY TO THINK OF THESE DATA STRUCTURES IS THAT THE HIGHER DIMENSIONAL DATA STRUCTURE IS A CONTAINER OF ITS LOWER DIMENSIONAL DATA STRUCTURE. FOR EXAMPLE, DATAFRAME IS A CONTAINER OF SERIES, PANEL IS A CONTAINER OF DATAFRAME.

Data Structure	Dimensions	Description
Series	1	1D labeled homogeneous array, size immutable.
Data Frames	2	General 2D labeled, size-mutable tabular structure with potentially heterogeneously typed columns.
Panel	3	General 3D labeled, size-mutable array.

NOTES

- ALL PANDAS DATA STRUCTURES ARE VALUE MUTABLE (CAN BE CHANGED) AND EXCEPT SERIES ALL ARE SIZE MUTABLE. SERIES IS SIZE IMMUTABLE.
- DATAFRAME IS WIDELY USED AND ONE OF THE MOST IMPORTANT DATA STRUCTURES. PANEL IS USED MUCH LESS.

SERIES

- SERIES IS A ONE-DIMENSIONAL ARRAY LIKE STRUCTURE WITH HOMOGENEOUS DATA. FOR EXAMPLE, THE FOLLOWING SERIES IS A COLLECTION OF INTEGERS 10, 23, 56, ...

10	23	56	17	52	61	73	90	26	72
----	----	----	----	----	----	----	----	----	----

- HOMOGENEOUS DATA
- SIZE IMMUTABLE
- VALUES OF DATA MUTABLE

DATAFRAME

- DATAFRAME IS A TWO-DIMENSIONAL ARRAY WITH HETEROGENEOUS DATA.

Name	Age	Gender	Rating
Steve	32	Male	3.45
Lia	28	Female	4.6
Vin	45	Male	3.9
Katie	38	Female	2.78

PANEL

- PANEL IS A THREE-DIMENSIONAL DATA STRUCTURE WITH HETEROGENEOUS DATA. IT IS HARD TO REPRESENT THE PANEL IN GRAPHICAL REPRESENTATION. BUT A PANEL CAN BE ILLUSTRATED AS A CONTAINER OF DATAFRAME.
- HETEROGENEOUS DATA
- SIZE MUTABLE
- DATA MUTABLE

CREATING A SERIES

- PANDAS.**SERIES**(DATA, INDEX, DTYPE, COPY)
- A SERIES CAN BE CREATED USING VARIOUS INPUTS LIKE:
 1. ARRAY
 2. DICT
 3. SCALAR VALUE OR CONSTANT

EXAMPLE

```
IMPORT PANDAS AS PD  
S = PD.SERIES()  
  
IMPORT PANDAS AS PD  
IMPORT NUMPY AS NP  
DATA = NP.ARRAY(['A','B','C','D'])  
#S = PD.SERIES(DATA, INDEX=[100,101,102,103])  
S = PD.SERIES(DATA)  
PRINT (S)
```

CREATING A SERIES FROM A DICT

```
DATA = {'A' : 0., 'B' : 1., 'C' : 2.}  
S = PD.SERIES(DATA)  
  
DATA = {'A' : 0., 'B' : 1., 'C' : 2.}  
S = PD.SERIES(DATA, INDEX=['B', 'C', 'D', 'A'])
```

CREATING A SERIES FROM A SCALAR

```
s = pd.Series(5, index=[0, 1, 2, 3])
```

OUTPUT:

```
0    5
```

```
1    5
```

```
2    5
```

```
3    5
```

DTYPE: INT64

ACCESSING DATA FROM SERIES WITH POSITION

```
s = pd.Series([1,2,3,4,5], index = [ 'A' , 'B' , 'C' , 'D' , 'E' ])
```

```
#RETRIEVE THE FIRST ELEMENT
```

```
print (s[0])
```

```
s = pd.Series([1,2,3,4,5], index = [ 'A' , 'B' , 'C' , 'D' , 'E' ])
```

```
#RETRIEVE THE FIRST THREE ELEMENT
```

```
print s[:3]
```

RETRIEVE DATA USING LABEL (INDEX)

```
s = pd.Series([1,2,3,4,5], index = ['A','B','C','D','E'])

#RETRIEVE A SINGLE ELEMENT

print (s['A'])

s = pd.Series([1,2,3,4,5], index = ['A','B','C','D','E'])

#RETRIEVE MULTIPLE ELEMENTS

print s[['A','C','D']]
```

FEATURES OF DATAFRAME

- POTENTIALLY COLUMNS ARE OF DIFFERENT TYPES
- SIZE – MUTABLE
- LABELED AXES (ROWS AND COLUMNS)
- CAN PERFORM ARITHMETIC OPERATIONS ON ROWS AND COLUMNS

CREATING A DATAFRAME

- `PANDAS.DataFrame(DATA, INDEX, COLUMNS, DTYPES, COPY)`
- A PANDAS DataFrame CAN BE CREATED USING VARIOUS INPUTS LIKE:
 1. LISTS
 2. DICT
 3. SERIES
 4. NUMPY NDARRAYS
 5. ANOTHER DataFrame

EXAMPLES

```
DF = PD.DataFrame()
```

```
DATA = [1,2,3,4,5]
```

```
DF = PD.DataFrame(DATA)
```

```
DATA = [['ALEX',10],['BOB',12],['CLARKE',13]]
```

```
DF = PD.DataFrame(DATA,COLUMNS=['NAME','AGE'])
```

EXMAPLES CONTINUED

```
DATA = [[ 'ALEX' ,10],[ 'Bob' ,12],[ 'CLARKE' ,13]]  
DF = PD.DATAFRAME(DATA,COLUMNS=[ 'NAME' , 'AGE' ],DTYPE=FLOAT)
```

CREATE A DATAFRAME FROM DICT OF NDARRAYS AND LISTS

```
DATA = {'NAME': ['Tom', 'Jack', 'Steve', 'Ricky'], 'AGE':[28,34,29,42]}

DF = PD.DATAFRAME(DATA)
```

```
DATA = {'NAME': ['Tom', 'Jack', 'Steve', 'Ricky'], 'AGE':[28,34,29,42]}

DF = PD.DATAFRAME(DATA, INDEX=['RANK1','RANK2','RANK3','RANK4'])
```

CREATE A DATAFRAME FROM LIST OF DICTS

```
DATA = [ {'A': 1, 'B': 2},{'A': 5, 'B': 10, 'C': 20}]  
DF = PD.DATAFRAME(DATA)
```

```
DATA = [ {'A': 1, 'B': 2},{'A': 5, 'B': 10, 'C': 20}]  
DF = PD.DATAFRAME(DATA, INDEX=['FIRST', 'SECOND'])
```

CONTINUED

```
DATA = [{ 'A': 1, 'B': 2}, { 'A': 5, 'B': 10, 'C': 20}]

#WITH TWO COLUMN INDICES, VALUES SAME AS DICTIONARY KEYS
DF1 = PD.DATAFRAME(DATA, INDEX=[ 'FIRST', 'SECOND'], COLUMNS=[ 'A', 'B'])

#WITH TWO COLUMN INDICES WITH ONE INDEX WITH OTHER NAME
DF2 = PD.DATAFRAME(DATA, INDEX=[ 'FIRST', 'SECOND'], COLUMNS=[ 'A', 'B1'])
```

CREATE A DATAFRAME FROM DICT OF SERIES

```
D = { 'ONE' : PD.SERIES([1, 2, 3], INDEX=['A', 'B', 'C']),  
      'TWO' : PD.SERIES([1, 2, 3, 4], INDEX=['A', 'B', 'C', 'D'])}  
  
DF = PD.DATAFRAME(D)
```

CONTINUED

```
D = {'ONE' : pd.Series([1, 2, 3], index=['A', 'B', 'C']),  
      'TWO' : pd.Series([1, 2, 3, 4], index=['A', 'B', 'C', 'D'])}  
  
DF = pd.DataFrame(D)
```

COLUMN ADDITION

```
IMPORT PANDAS AS PD

D = { 'ONE' : PD.SERIES([1, 2, 3], INDEX=['A', 'B', 'C']),
      'TWO' : PD.SERIES([1, 2, 3, 4], INDEX=['A', 'B', 'C', 'D'])}

DF = PD.DATAFRAME(D)

# ADDING A NEW COLUMN TO AN EXISTING DATAFRAME OBJECT WITH COLUMN LABEL BY
PASSING NEW SERIES

PRINT ("ADDING A NEW COLUMN BY PASSING AS SERIES:")
DF['THREE']=PD.SERIES([10,20,30],INDEX=['A','B','C'])

PRINT DF

PRINT ("ADDING A NEW COLUMN USING THE EXISTING COLUMNS IN DATAFRAME:")
DF['FOUR']=DF['ONE']+DF['THREE']
```

COLUMN DELETION

```
IMPORT PANDAS AS PD

D = {'ONE' : PD.SERIES([1, 2, 3], INDEX=['A', 'B', 'C']),
      'TWO' : PD.SERIES([1, 2, 3, 4], INDEX=['A', 'B', 'C', 'D']),
      'THREE' : PD.SERIES([10,20,30], INDEX=['A','B','C'])}

DF = PD.DATAFRAME(D)
PRINT ("OUR DATAFRAME IS:")
PRINT DF

# USING DEL FUNCTION
PRINT ("DELETING THE FIRST COLUMN USING DEL FUNCTION:")
DEL DF['ONE']
PRINT DF

# USING POP FUNCTION
PRINT ("DELETING ANOTHER COLUMN USING POP FUNCTION:")
DF.pop('TWO')
```

ROW SELECTION, ADDITION, AND DELETION

- ROWS CAN BE SELECTED BY PASSING ROW LABEL TO A **LOC** FUNCTION.

```
D = { 'ONE' : pd.Series([1, 2, 3], index=['A', 'B', 'C']),  
      'TWO' : pd.Series([1, 2, 3, 4], index=['A', 'B', 'C', 'D'])}
```

```
df = pd.DataFrame(D)
```

```
print (df.loc['B'])
```

- THE RESULT IS A SERIES WITH LABELS AS COLUMN NAMES OF THE DATAFRAME. AND, THE NAME OF THE SERIES IS THE LABEL WITH WHICH IT IS RETRIEVED.

CONTINUED

- ROWS CAN BE SELECTED BY PASSING INTEGER LOCATION TO AN **ILOC** FUNCTION.

```
D = {'ONE' : pd.Series([1, 2, 3], index=['A', 'B', 'C']),  
      'TWO' : pd.Series([1, 2, 3, 4], index=['A', 'B', 'C', 'D'])}
```

```
df = pd.DataFrame(D)  
print(df.iloc[2])
```

SLICE ROWS

- MULTIPLE ROWS CAN BE SELECTED USING ‘ : ’ OPERATOR

```
D = { 'ONE' : pd.Series([1, 2, 3], index=['A', 'B', 'C']),  
      'TWO' : pd.Series([1, 2, 3, 4], index=['A', 'B', 'C', 'D'])}
```

```
df = pd.DataFrame(D)
```

```
print(df[2:4])
```

ADDITION OF ROWS

```
DF = PD.DataFrame([[1, 2], [3, 4]], columns = ['A', 'B'])
```

```
DF2 = PD.DataFrame([[5, 6], [7, 8]], columns = ['A', 'B'])
```

```
DF = DF.append(DF2)
```

DELETION OF ROWS

```
IMPORT PANDAS AS PD

DF = PD.DataFrame([[1, 2], [3, 4]], columns = ['A','B'])
DF2 = PD.DataFrame([[5, 6], [7, 8]], columns = ['A','B'])

DF = DF.append(DF2)

# DROP ROWS WITH LABEL 0
DF = DF.drop(0)
```

SERIES FUNCTIONS

```
s = pd.Series(np.random.randn(4))
```

```
s = pd.Series(np.random.randn(4))
```

```
print ("THE AXES ARE:")
```

```
print (s.axes)
```

- RETURNS THE LIST OF THE LABELS OF THE SERIES

CONTINUED

```
S = PD.SERIES(NP.RANDOM.RANDN(4))  
PRINT ("IS THE OBJECT EMPTY?")  
PRINT S.EMPTY  


- RETURNS A BOOLEAN IF THE SERIES IS EMPTY OR NOT.

```

CONTINUED

- `NDIM` , RETURNS DIMENSION.
- `SIZE` , RETURNS NUMBER OF ELEMENTS IN THE SERIES.
- `VALUES`, RETURNS ACTUAL ELEMENTS IN THE SERIES.
- `HEAD()` AND `TAIL()`, RETURNS VALUES FROM BOTH SIDES.

DATAFRAME BASIC FUNCTIONALITIES

- T (TRANSPOSE OF A DATAFRAME)

```
D = {'NAME':pd.Series(['TOM','JAMES','RICKY','VIN','STEVE','SMITH','JACK']),  
'AGE':pd.Series([25,26,25,23,30,29,23]),  
'RATING':pd.Series([4.23,3.24,3.98,2.56,3.20,4.6,3.8])}
```

```
# CREATE A DATAFRAME  
DF = pd.DataFrame(D)  
PRINT ("THE transpose OF THE DATA SERIES IS:")  
PRINT (DF.T)
```

CONTINUED

- AXES
- DTYPES
- EMPTY
- NDIM
- SHAPE
- SIZE
- VALUES
- HEAD
- TAILS

DESCRIPTIVE STATISTICS

- `SUM()` , YOU CAN SPECIFY THE AXIS
- `MEAN()`
- `STD()`
- `DESCRIBE()`

ROW OR COLUMN WISE FUNCTION APPLICATION

```
IMPORT PANDAS AS PD  
IMPORT NUMPY AS NP  
  
DF = PD.DATAFRAME(NP.RANDOM.RANDN(5,3),COLUMNS=[ 'COL1','COL2','COL3' ])  
DF.APPLY(NP.MEAN)  
PRINT DF.APPLY(NP.MEAN)
```

ITERATING A DATAFRAME

- ITERATING A DATAFRAME GIVES COLUMN NAMES

N=20

```
DF = PD.DATAFRAME({  
    'A': PD.DATE_RANGE(START='2016-01-01',PERIODS=N,FREQ='D'),  
    'X': NP.LINSPACE(0,STOP=N-1,NUM=N),  
    'Y': NP.RANDOM.RAND(N),  
    'C': NP.RANDOM.CHOICE(['Low','MEDIUM','HIGH'],N).TOLIST(),  
    'D': NP.RANDOM.NORMAL(100, 10, SIZE=(N)).TOLIST()  
})
```

```
FOR COL IN DF:  
    PRINT (COL)
```

CONTINUED

- TO ITERATE OVER THE ROWS OF THE DATAFRAME, WE CAN USE THE FOLLOWING FUNCTIONS:
 1. **ITERITEMS()** – TO ITERATE OVER THE (KEY,VALUE) PAIRS
 2. **ITERROWS()** – ITERATE OVER THE ROWS AS (INDEX,SERIES) PAIRS
 3. **ITERTUPLES()** – ITERATE OVER THE ROWS AS NAMEDTUPLES

ITERITEMS()

- ITERATES OVER EACH COLUMN AS KEY, VALUE PAIR WITH LABEL AS KEY AND COLUMN VALUE AS A SERIES OBJECT

```
IMPORT PANDAS AS PD
```

```
IMPORT NUMPY AS NP
```

```
DF = PD.DATAFRAME(NP.RANDOM.RANDN(4,3),COLUMNS=[ 'COL1','COL2','COL3' ])  
FOR KEY,VALUE IN DF.ITERITEMS():  
    PRINT (KEY,VALUE)
```

ITERROWS()

- ITERROWS() RETURNS THE ITERATOR YIELDING EACH INDEX VALUE ALONG WITH A SERIES CONTAINING THE DATA IN EACH ROW.

```
IMPORT PANDAS AS PD
```

```
IMPORT NUMPY AS NP
```

```
DF = PD.DATAFRAME(NP.RANDOM.RANDN(4,3),COLUMNS = [ 'COL1','COL2','COL3' ])  
FOR ROW_INDEX,ROW IN DF.ITERROWS():  
    PRINT (ROW_INDEX,ROW)
```

ITERTUPLES()

- ITERTUPLES() METHOD WILL RETURN AN ITERATOR YIELDING A NAMED TUPLE FOR EACH ROW IN THE DATAFRAME. THE FIRST ELEMENT OF THE TUPLE WILL BE THE ROW'S CORRESPONDING INDEX VALUE, WHILE THE REMAINING VALUES ARE THE ROW VALUES.

```
IMPORT PANDAS AS PD
```

```
IMPORT NUMPY AS NP
```

```
DF = PD.DATAFRAME(NP.RANDOM.RANDN(4,3),COLUMNS = [ 'COL1','COL2','COL3' ])
```

```
FOR ROW IN DF.ITERTUPLES():
```

```
    PRINT (ROW)
```