

Practical Tutorial on Data Manipulation with Numpy and Pandas in Python:-

The pandas library has emerged into the open house of data manipulation used in python since it was developed in 2008. With its intuitive syntax and flexible data structure, it's easy to learn and enables faster data computation. The development of numpy and pandas libraries has extended python's multi-purpose nature to solve machine learning problems as well. The acceptance of python language in machine learning has been phenomenal since then.

6 Important things about Numpy and Pandas:-

1. The data manipulation capabilities of pandas are built on top of the numpy library. In a way, numpy is a dependency of the pandas library.
2. Pandas is best at handling tabular data sets comprising different variable types (integer, float, double, etc.). In addition, the pandas library can also be used to perform even the most naive of tasks such as loading data or doing feature engineering on time series data.
3. Numpy is most suitable for performing basic numerical computations such as mean, median, range, etc. *Alongside*, it also supports the creation of multi-dimensional arrays.
4. Numpy library can also be used to prepare C/C++ and Fortran code.
5. Remember, python is a zero indexing language unlike R where indexing starts at one.
6. The best part of learning pandas and numpy is the strong active community support you'll get from around the world.

a) STARTING WITH NUMPY:-

```
In [5]: import numpy as np

In [6]: L = 1141*(range(10))

In [7]: #converting integers to string - this style of handling lists is known as list comprehension.
#list comprehension offers a very useful way to handle list manipulation tasks easily, we'll learn about them in future tutorials. Here's an example.
[str(i) for c in L]
['0', '1', '2', '3', '4', '5', '6', '7', '8', '9']
[types] For item in L

Out[7]: [int, int, int, int, int, int, int, int, int, int]

a) Creation of array:-

In [20]: #creating arrays
np.zeros(10, dtype='int')
np.array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0])

#creating a 3 row x 5 column matrix
np.ones((3,5), dtype=float)
np.array([[1., 1., 1., 1., 1.],
         [1., 1., 1., 1., 1.],
         [1., 1., 1., 1., 1.]])

#creating a matrix with a predefined value
np.full((2,5),1.23)
np.array([[1.23, 1.23, 1.23, 1.23, 1.23],
         [1.23, 1.23, 1.23, 1.23, 1.23]])

#create an array with a set sequence
np.arange(0, 20, 2)
np.array([0, 2, 4, 6, 8,10,12,14,16,18])

#create an array of even space between the given range of values
np.linspace(0, 1, 5)
np.array([ 1.,0.25, 0.5, 0.75, 1.])

#create a 3D array with mean 0 and standard deviation 1 in a given dimension
np.random.normal(0, 1, (2,3))
np.array([[ 0.48430229, 1.4096656, -1.0467500],
        [-0.42939094, -1.0208444, -1.50730593]])

#create an identity matrix
np.eye(3)
np.array([[1., 0., 0.],
        [0., 1., 0.],
        [0., 0., 1.]])

#set a random seed
np.random.seed(0)

x1 = np.random.randint(10, size=6) #one dimension
x2 = np.random.randint(10, size=(3,4)) #two dimension
x3 = np.random.randint(10, size=(3,4,5)) #three dimension

print("x3 ndim:", x3.ndim)
print("x3 shape:", x3.shape)
print("x3 size:", x3.size)

x3.ndim: 3
x3.shape: (3, 4, 5)
x3.size: 60
b) Array indexing:-

In [35]: x1 = np.array([4, 3, 4, 4, 8, 4])
x1
np.array([4, 3, 4, 4, 8, 4])
#access value to index zero
x1[0]
4

#access first value
x1[4]
8

#get the last value
x1[-1]
4

#get the second last value
x1[-2]
8

#in a multidimensional array, we need to specify row and column index
x2
np.array([[3, 7, 5, 8],
        [0, 5, 5, 0],
        [3, 0, 5, 0]])

#first row and 3rd column value
x2[0,3]
8

#3rd row and last value from the 3rd column
x2[2,-1]
0

#evaluate value at 0,0 index
x2[0,0] = 12
x2
Out[35]: array([[12, 7, 5, 8],
        [ 0, 5, 5, 0],
        [ 3, 0, 5, 0]])

c) Array Slicing:-

In [37]: x1 = np.array([4, 3, 4, 4, 8, 4])
x1
np.array([4, 3, 4, 4, 8, 4])
#access value to index zero
x1[0]
4

#access first value
x1[4]
8

#get the last value
x1[-1]
4

#get the second last value
x1[-2]
8

#in a multidimensional array, we need to specify row and column index
x2= np.array([[3, 7, 5, 8],
             [0, 5, 5, 0],
             [3, 0, 5, 0]])

#first row and 3rd column value
x2[0,3]
8

#3rd row and last value from the 3rd column
x2[2,-1]
0

#evaluate value at 0,0 index
x2[0,0] = 12
np.array([[12, 7, 5, 8],
        [ 0, 5, 5, 0],
        [ 3, 0, 5, 0]])

#array slicing
#now, we'll learn to access multiple or a range of elements from an array.

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

#from start to 4th position
x[5]
np.array([6, 1, 2, 3, 4])

#from 4th position to end
x[4:]
np.array([4, 5, 6, 7, 8, 9])

#from 0th to 4th position
x[:4]
np.array([4, 5, 6, 9])

#return elements at even place
x[::2]
np.array([0, 2, 4, 6, 8])

#return elements from first position step by two
x[1::2]
np.array([1, 3, 5, 7, 9])

#reverse the array
x[::-1]
np.array([9, 8, 7, 6, 5, 4, 3, 2, 1, 0])

Out[37]: array([9, 8, 7, 6, 5, 4, 3, 2, 1, 0])

d) Array Concatenation:-

In [38]: #You can concatenate two or more arrays at once.
x = np.array([1, 2, 3])
y = np.array([2, 2, 1])
z = [2,2,2,2]
np.concatenate(x,y,z)
np.array([1, 2, 3, 2, 2, 1, 2, 2, 2, 2, 2, 2, 2])

#You can also use this function to create 2-dimensional arrays.
grid = np.array([[1,2,3],[4,5,6]])
np.concatenate([grid,grid])
np.array([[1, 2, 3],
        [4, 5, 6],
        [1, 2, 3],
        [4, 5, 6]])

#Using its axis parameter, you can define row-wise or column-wise matrix
np.concatenate([grid*grid,grid*grid])

Out[38]: array([[1, 2, 3, 1, 2, 3],
        [4, 5, 6, 4, 5, 6]])
```

b) STARTING WITH PANDAS:-

```
In [39]: import pandas as pd

In [42]: #create a data frame - dictionary is used here where keys get converted to column names and values to row values.
data = pd.DataFrame({'Country': ['Russia','Columbia','China','Equador','Nigeria'],
                    'data': [[121,40,100,130,11]]})

Out[42]:
   Country  Rank
0  Russia    121
1  Columbia  40
2    China   100
3  Equador   130
4  Nigeria    11

In [43]: #we can do a quick analysis of any data set using:
data.describe()

Out[43]:
   Rank
count  5.000000
mean    80.400000
std     52.300000
min     11.000000
25%    40.000000
50%    100.000000
75%    121.000000
max    130.000000

In [55]: #let's create another data frame.
pd.DataFrame({'group':['a','a','a','b','b','b','b','c','c','c'],'ounces':[4, 3, 12, 6, 7, 8, 3, 5, 6]})

Out[55]:
   group  ounces
0      a      4.0
1      a     3.0
2      a    12.0
3      b     6.0
4      b     7.5
5      b     8.0
6      c     3.0
7      c     5.0
8      c    12.0

In [57]: #let's sort the data frame by ounces - inplace = True will make changes to the data
data.sort_values(by=['ounces'],ascending=True,inplace=False)

Out[57]:
   group  ounces
1      a     3.0
6      c     3.0
0      a     4.0
3      b     6.0
2      b     6.0
8      c     6.0
4      b     7.5
5      b     8.0
2      a    12.0

In [58]: #data.sort_values(by=['group','ounces'],ascending=[True,False],inplace=False)

Out[58]:
   group  ounces
2      a    12.0
1      a     3.0
1      a     4.0
3      b     6.0
5      b     7.5
3      b     8.0
6      c     6.0
7      c     5.0
0      c     3.0

In [61]: #remove duplicates - is dar
data.drop_duplicates()

Out[61]:
   k1 k2
0  one  1
1  one  2
2  one  1
3  two  3
4  two  2
5  two  4

In [62]: data.drop_duplicates(subset='k1')

Out[62]:
   k1 k2
3  two  3

In [63]: data = pd.DataFrame({'food':['bacon','pulled pork','bacon','Pastrami','corned beef','bacon','pastrami','honey ham','nova lox'],
                             'ounces': [4, 3, 12, 6, 7, 8, 3, 5, 6]})

Out[63]:
   food  ounces
0  bacon     4.0
1 pulled pork  3.0
2  bacon    12.0
3  Pastrami    6.0
4 corned beef  7.5
5  Bacon      8.0
6 pastrami    3.0
7 honeyham    5.0
8 nova lox    6.0

In [64]: meat_to_animal = {
    'bacon': 'pig',
    'pulled pork': 'pig',
    'pastrami': 'cow',
    'corned beef': 'cow',
    'honey ham': 'hog',
    'nova lox': 'salmon'
}

def meat_2_animal(series):
    if series['food'] == 'bacon':
        return 'pig'
    elif series['food'] == 'pulled pork':
        return 'pig'
    elif series['food'] == 'pastrami':
        return 'cow'
    elif series['food'] == 'corned beef':
        return 'cow'
    elif series['food'] == 'honey ham':
        return 'hog'
    elif series['food'] == 'nova lox':
        return 'salmon'

#create a new variable
data['animal'] = data['food'].map(str.lower).map(meat_to_animal)

Out[64]:
   food  ounces animal
0  bacon     4.0    pig
1 pulled pork  3.0    pig
2  bacon    12.0    pig
3  Pastrami    6.0   cow
4 corned beef  7.5   cow
5  Bacon      8.0   cow
6 pastrami    3.0   cow
7 honeyham    5.0    pig
8 nova lox    6.0 salmon

In [65]: #another use of doing if is: convert the food values to the lower case and apply the function
lower = lambda x: x.lower()
data['food'] = data['food'].apply(lower)
data['animal2'] = data.apply(meat_2_animal, axis='columns')
data

Out[65]:
   food  ounces animal
0  bacon     4.0    pig
1 pulled pork  3.0    pig
2  bacon    12.0    pig
3  Pastrami    6.0   cow
4 corned beef  7.5   cow
5  Bacon      8.0   cow
6 pastrami    3.0   cow
7 honeyham    5.0    pig
8 nova lox    6.0 salmon

In [66]: data.assign(new_variable = data['ounces']**18)

Out[66]:
   food  ounces animal new_variable
0  bacon     4.0    pig          68.0
1 pulled pork  3.0    pig          30.0
2  bacon    12.0    pig        120.0
3  Pastrami    6.0   cow         60.0
4 corned beef  7.5   cow        75.0
5  Bacon      8.0   cow         80.0
6 pastrami    3.0   cow         30.0
7 honeyham    5.0    pig         50.0
8 nova lox    6.0 salmon        60.0

In [67]: data.drop('animal2',axis='columns',inplace=True)
data

Out[67]:
   food  ounces animal
0  bacon     4.0    pig
1 pulled pork  3.0    pig
2  bacon    12.0    pig
3  Pastrami    6.0   cow
4 corned beef  7.5   cow
5  Bacon      8.0   cow
6 pastrami    3.0   cow
7 honeyham    5.0    pig
8 nova lox    6.0 salmon

In [68]: #Series function from pandas are used to create arrays
data = pd.Series([1, -999, 2, -999, -1000, 3,2])

Out[68]:
0      1.0
1    -999.0
2       2.0
3   -1000.0
4    -1000.0
dtype: float64

In [69]: #replace -999 with nan values
data.replace(-999, np.nan,inplace=True)
data

Out[69]:
0      1.0
1      NaN
2       2.0
3   -1000.0
4    -1000.0
dtype: float64

In [70]: #we can also replace multiple values at once.
data = pd.Series([1, -999, 2, -999, -1000, 3,2])
data.replace([-999,-1000],np.nan,inplace=True)
data

Out[70]:
0      1.0
1      NaN
2       2.0
3      NaN
4      NaN
dtype: float64

In [71]: data = pd.DataFrame(np.arange(12).reshape((3, 4)),index=['Ohio','Colorado','New York'],columns=['one','two','three','four'])
data

Out[71]:
   one  two  three  four
Ohio   0   1     2     3
Colorado  4   5     6     7
New York  8   9    10    11

In [72]: #Using reshape function
data.rename(index = {'Ohio':'Sanf'}, columns={'one':'one_p','two':'two_p'},inplace=True)
data

Out[72]:
   one_p two_p three  four
Sanf    0   1     2     3
Colorado  4   5     6     7
New York  8   9    10    11

In [73]: #We can also use string functions
data.rename(index = str.upper, columns=str.title,inplace=True)
data

Out[73]:
   One_P Two_P Three  Four
Sanf    0   1     2     3
Colorado  4   5     6     7
NEWYORK  8   9    10    11

In [77]: #get = [20, 22, 25, 27, 31, 33, 37, 31, 61, 40, 41, 32]

In [78]: #using the isin() - 'i' means the value is included in the list, 'I' means the value is excluded
bins = [10, 25, 35, 60, 100]
cats = pd.cut(Ages, bins)

Out[78]:
[(10, 25], (10, 25], (25, 35], (25, 35], ..., (10, 25], (60, 100], (60, 100], (35, 60], (25, 35)]
Categories (4, interval[int64, right]): [(10, 25] < (25, 35] < (35, 60] < (60, 100]]

In [79]: #to include the right bin value, we can do:
pd.cut(Ages,bins,include=True)

Out[79]:
[(10, 25], (10, 25], (25, 35], (25, 35], (10, 25], ..., (25, 35], (60, 100], (35, 60], (25, 35)]
length: 12
Categories (4, interval[int64, left]): [(10, 25] < (25, 35] < (35, 60] < (60, 100]]

In [82]: #let's check how many observations fall under each bin
pd.value_counts(cats)

Out[82]:
(10, 25]    5
(25, 35]    3
(60, 100]   3
(35, 60]    1
dtype: int64

In [83]: #Bins names = 'Youth', 'YoungAdult', 'MiddleAge', 'Senior'
new_cats = pd.cut(Ages, bins, labels=['a','b','c','d'])
pd.value_counts(new_cats)

Out[83]:
Youth      5
YoungAdult  3
MiddleAge   3
Senior      1
dtype: int64

In [84]: #we can also calculate their cumulative sum
pd.value_counts(new_cats.cumsum())

Out[84]:
Youth      5
YoungAdult  8
MiddleAge  11
Senior     12
dtype: int64

In [85]: #if = pd.DataFrame({'keys': ['a','a','b','b','b','a'],
                        'keys2': ['one','two','one','two','one'],
                        'data': np.random.randn(6)})
df

Out[85]:
   keys  keys2  data    data2
0      a      one  1.254414  1.149078
1      a      two  1.439302  -1.339778
2      b      one -0.010864  1.142461
3      b      two  2.527437  1.509448
4      a      one -1.507006  1.007775

In [86]: #dates = pd.date_range('20131001',periods=5)
df = pd.DataFrame(np.random.randn(5,4),index=dates,columns=list('ABCD'))
df

Out[86]:
              A          B          C          D
2013-01-01  0.680595  0.014873 -0.375660 -0.030234
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484
2013-01-03  0.750006  0.350003  1.076121  0.102141
2013-01-04  0.367974 -0.044724 -0.302375 -0.220484
2013-01-05  0.461100  0.712678  1.113340 -0.350979

In [88]: #get first n rows from the data frame
df[:3]

Out[88]:
              A          B          C          D
2013-01-01  0.680595  0.014873 -0.375660 -0.030234
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484
2013-01-03  0.750006  0.350003  1.076121  0.102141

In [89]: #get first n rows from the data frame
df[:3]

Out[89]:
              A          B          C          D
2013-01-01  0.680595  0.014873 -0.375660 -0.030234
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484
2013-01-03  0.750006  0.350003  1.076121  0.102141

In [91]: #filter based on column names
df.loc[:,['A','B']]

Out[91]:
              A          B
2013-01-01  0.680595  0.014873
2013-01-02  0.367974 -0.044724
2013-01-03  0.750006  0.350003
2013-01-04  0.367974 -0.044724
2013-01-05  0.461100  0.712678

In [92]: #filter based on both row index labels and column names
df.loc['2013-01-01':'2013-01-03',['A','B']]

Out[92]:
              A          B
2013-01-01  0.680595 -0.044724
2013-01-02  0.367974 -0.044724
2013-01-03  0.750006  0.350003

In [93]: #filter based on index of columns
df.iloc[0].reset_index(drop=True)

Out[93]:
0      0.680595 -0.044724 -0.375660 -0.030234
1      0.367974 -0.044724 -0.302375 -0.220484
2      0.750006  0.350003  1.076121  0.102141
3      0.367974 -0.044724 -0.302375 -0.220484
4      0.461100  0.712678  1.113340 -0.350979

In [94]: #resample a specific range of rows
df.iloc[2:4, 0:2]

Out[94]:
              A          B
2013-01-03  0.750006  0.350003
2013-01-04  0.367974 -0.044724

In [95]: #create specific row and columns using lists containing columns or row indexes
df.iloc[[1,3],0:2]

Out[95]:
              A          C
2013-01-02  0.367974 -0.302375
2013-01-04  0.461100  0.223139

In [96]: df[df.A > 1]

Out[96]:
              A          B          C          D
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484
2013-01-04  0.461100  0.712678  1.113340 -0.350979

In [97]: #we can copy the data set
df2 = df.copy()
df2['E'] = ['one','one','two','three','four','three']
df2

Out[97]:
              A          B          C          D          E
2013-01-01  0.680595  0.014873 -0.375660 -0.030234  one
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484  one
2013-01-03  0.750006  0.350003  1.076121  0.102141  two
2013-01-04  0.367974 -0.044724 -0.302375 -0.220484  three
2013-01-05  0.461100  0.712678  1.113340 -0.350979  four

In [98]: #select rows based on column values
df2[df2['E'] != 'one' & df2['C'] > 1]

Out[98]:
              A          B          C          D          E
2013-01-03  0.750006  0.350003  1.076121  0.102141  two
2013-01-04  0.367974 -0.044724 -0.302375 -0.220484  three
2013-01-05  0.461100  0.712678  1.113340 -0.350979  four

In [99]: #select all rows except those with two and four
df2[~df2['E'].isin(['two','four'])]

Out[99]:
              A          B          C          D          E
2013-01-01  0.680595  0.014873 -0.375660 -0.030234  one
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484  one
2013-01-03  0.750006  0.350003  1.076121  0.102141  two
2013-01-04  0.367974 -0.044724 -0.302375 -0.220484  three
2013-01-05  0.461100  0.712678  1.113340 -0.350979  four

In [100]: #list all columns where A is greater than C
df[df['A'] > C]

Out[100]:
              A          B          C          D
2013-01-01  0.680595  0.014873 -0.375660 -0.030234
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484
2013-01-03  0.750006  0.350003  1.076121  0.102141
2013-01-04  0.367974 -0.044724 -0.302375 -0.220484
2013-01-05  0.461100  0.712678  1.113340 -0.350979

In [101]: #using OR condition
df.query('A < B | C > A')

Out[101]:
              A          B          C          D
2013-01-01  0.680595  0.014873 -0.375660 -0.030234
2013-01-02  0.367974 -0.044724 -0.302375 -0.220484
2013-01-03  0.750006  0.350003  1.076121  0.102141
2013-01-04  0.367974 -0.044724 -0.302375 -0.220484
2013-01-05  0.461100  0.712678  1.113340 -0.350979

In [102]: #create a data frame
data = pd.DataFrame({'group':['a','a','a','b','b','b','b','c','c','c'],
                    'ounces': [4, 3, 12, 6, 7, 8, 3, 5, 6, 5]})

Out[102]:
   group  ounces
0      a      4.0
1      a     3.0
2      a    12.0
3      b     6.0
4      b     7.5
5      b     8.0
6      c     3.0
7      c     5.0
8      c    12.0
```