

▼ Student ID: 19522348

Name: Le Duc Tin

```
import numpy as np
import pandas as pd
```

▼ Creating a Series

You can convert a list, numpy array, or dictionary to a Series:

```
labels = ['a','b','c']
my_list = [10,20,30]
arr = np.array([10,20,30])
d = {'a':10, 'b':20, 'c':30}
```

**** Using Lists ****

```
pd.Series(data=my_list)
```

```
0    10
1    20
2    30
dtype: int64
```

```
pd.Series(data=my_list,index=labels)
```

```
a    10
b    20
c    30
dtype: int64
```

```
pd.Series(my_list,labels)
```

```
a    10
b    20
c    30
dtype: int64
```

**** NumPy Arrays ****

```
pd.Series(arr)
```

```
pd.Series(arr)

0    10
1    20
2    30
dtype: int64
```

```
pd.Series(arr, labels)
```

```
a    10
b    20
c    30
dtype: int64
```

**** Dictionary****

```
pd.Series(d)
```

```
a    10
b    20
c    30
dtype: int64
```

▼ Data in a Series

A pandas Series can hold a variety of object types:

```
pd.Series(data=labels)
```

```
0    a
1    b
2    c
dtype: object
```

Even functions (although unlikely that you will use this)

```
pd.Series([sum, print, len])

0    <built-in function sum>
1    <built-in function print>
2    <built-in function len>
dtype: object
```

▼ Using an Index

The key to using a Series is understanding its index. Pandas makes use of these index names or numbers by allowing for fast look ups of information (works like a hash table or dictionary).

Let's see some examples of how to grab information from a Series. Let us create two series, ser1 and ser2:

```
ser1 = pd.Series([1,2,3,4],index = ['USA', 'Germany', 'USSR', 'Japan'])
```

```
ser1
```

```
USA      1
Germany  2
USSR     3
Japan    4
dtype: int64
```

```
ser2 = pd.Series([1,2,5,4],index = ['USA', 'Germany', 'Italy', 'Japan'])
```

```
ser2
```

```
USA      1
Germany  2
Italy    5
Japan    4
dtype: int64
```

```
ser1['USA']
```

```
1
```

Operations are then also done based off of index:

```
ser1 + ser2
```

```
Germany  4.0
Italy    NaN
Japan    8.0
USA      2.0
USSR     NaN
dtype: float64
```

▼ DataFrames

DataFrames are the workhorse of pandas and are directly inspired by the R programming language. We can think of a DataFrame as a bunch of Series objects put together to share the same index. Let's use pandas to explore this topic!

```
import pandas as pd
import numpy as np
```

```
from numpy.random import randn
np.random.seed(101)
```

```
df = pd.DataFrame(randn(5,4),index='A B C D E'.split(),columns='W X Y Z'.split())
```

```
df
```

	W	X	Y	Z
A	2.706850	0.628133	0.907969	0.503826
B	0.651118	-0.319318	-0.848077	0.605965
C	-2.018168	0.740122	0.528813	-0.589001
D	0.188695	-0.758872	-0.933237	0.955057
E	0.190794	1.978757	2.605967	0.683509

▼ Selection and Indexing

Let's learn the various methods to grab data from a DataFrame

```
df['W']
```

```
A    2.706850
B    0.651118
C   -2.018168
D    0.188695
E    0.190794
Name: W, dtype: float64
```

```
# Pass a list of column names
df[['W','Z']]
```

	W	Z
A	2.706850	0.503826
B	0.651118	0.605965
C	-2.018168	-0.589001
D	0.188695	0.955057
E	0.190794	0.683509

```
# SQL Syntax (NOT RECOMMENDED!)
```

```
df.W
```

```
A    2.706850
B    0.651118
C   -2.018168
D    0.188695
E    0.190794
Name: W, dtype: float64
```

DataFrame Columns are just Series

```
type(df['W'])
```

```
pandas.core.series.Series
```

Creating a new column:

```
df['new'] = df['W'] + df['Y']
```

```
df
```

	W	X	Y	Z	new
A	2.706850	0.628133	0.907969	0.503826	3.614819
B	0.651118	-0.319318	-0.848077	0.605965	-0.196959
C	-2.018168	0.740122	0.528813	-0.589001	-1.489355
D	0.188695	-0.758872	-0.933237	0.955057	-0.744542
E	0.190794	1.978757	2.605967	0.683509	2.796762

**** Removing Columns****

```
df.drop('new',axis=1)
```

```

      W      X      Y      Z
-----
# Not inplace unless specified!
df

```

```

      W      X      Y      Z      new
-----
A  2.706850  0.628133  0.907969  0.503826  3.614819
B   0.651118 -0.319318 -0.848077  0.605965 -0.196959
C  -2.018168  0.740122  0.528813 -0.589001 -1.489355
D   0.188695 -0.758872 -0.933237  0.955057 -0.744542
E   0.190794  1.978757  2.605967  0.683509  2.796762

```

```
df.drop('new',axis=1,inplace=True)
```

```

df
      W      X      Y      Z
-----
A  2.706850  0.628133  0.907969  0.503826
B   0.651118 -0.319318 -0.848077  0.605965
C  -2.018168  0.740122  0.528813 -0.589001
D   0.188695 -0.758872 -0.933237  0.955057
E   0.190794  1.978757  2.605967  0.683509

```

Can also drop rows this way:

```
df.drop('E',axis=0)
```

```

      W      X      Y      Z
-----
A  2.706850  0.628133  0.907969  0.503826
B   0.651118 -0.319318 -0.848077  0.605965
C  -2.018168  0.740122  0.528813 -0.589001
D   0.188695 -0.758872 -0.933237  0.955057

```

**** Selecting Rows****

```
df.loc['A']
```

```
W    2.706850
X    0.628133
Y    0.907969
Z    0.503826
Name: A, dtype: float64
```

Or select based off of position instead of label

```
df.iloc[2]
```

```
W    -2.018168
X     0.740122
Y     0.528813
Z    -0.589001
Name: C, dtype: float64
```

**** Selecting subset of rows and columns ****

```
df.loc['B', 'Y']
```

```
-0.8480769834036315
```

```
df.loc[['A', 'B'], ['W', 'Y']]
```

	W	Y
A	2.706850	0.907969
B	0.651118	-0.848077

▼ Conditional Selection

An important feature of pandas is conditional selection using bracket notation, very similar to numpy:

```
df
```

	W	X	Y	Z
A	2.706850	0.628133	0.907969	0.503826

```
df>0
```

	W	X	Y	Z
A	True	True	True	True
B	True	False	False	True
C	False	True	True	False
D	True	False	False	True
E	True	True	True	True

```
df[df>0]
```

	W	X	Y	Z
A	2.706850	0.628133	0.907969	0.503826
B	0.651118	NaN	NaN	0.605965
C	NaN	0.740122	0.528813	NaN
D	0.188695	NaN	NaN	0.955057
E	0.190794	1.978757	2.605967	0.683509

```
df[df['W']>0]
```

	W	X	Y	Z
A	2.706850	0.628133	0.907969	0.503826
B	0.651118	-0.319318	-0.848077	0.605965
D	0.188695	-0.758872	-0.933237	0.955057
E	0.190794	1.978757	2.605967	0.683509

```
df[df['W']>0]['Y']
```

```
A    0.907969
B   -0.848077
D   -0.933237
E    2.605967
Name: Y, dtype: float64
```

```
df[df['W']>0]['Y']
```



```
df[['W', 'X', 'Y', 'Z']]
```

	Y	X
A	0.907969	0.628133
B	-0.848077	-0.319318
D	-0.933237	-0.758872
E	2.605967	1.978757

For two conditions you can use | and & with parenthesis:

```
df[(df['W']>0) & (df['Y'] > 1)]
```

	W	X	Y	Z
E	0.190794	1.978757	2.605967	0.683509

▼ More Index Details

Let's discuss some more features of indexing, including resetting the index or setting it something else. We'll also talk about index hierarchy!

```
df
```

	W	X	Y	Z
A	2.706850	0.628133	0.907969	0.503826
B	0.651118	-0.319318	-0.848077	0.605965
C	-2.018168	0.740122	0.528813	-0.589001
D	0.188695	-0.758872	-0.933237	0.955057
E	0.190794	1.978757	2.605967	0.683509

```
# Reset to default 0,1...n index
df.reset_index()
```

	index		W	X	Y	Z
0	A		2.706850	0.628133	0.907969	0.503826
1	B		0.651118	-0.319318	-0.848077	0.605965

```
newind = 'CA NY WY OR CO'.split()
```

```
df['States'] = newind
```

4	E		0.190794	1.978757	2.605967	0.683509
---	---	--	----------	----------	----------	----------

df

		W	X	Y	Z	States
A		2.706850	0.628133	0.907969	0.503826	CA
B		0.651118	-0.319318	-0.848077	0.605965	NY
C		-2.018168	0.740122	0.528813	-0.589001	WY
D		0.188695	-0.758872	-0.933237	0.955057	OR
E		0.190794	1.978757	2.605967	0.683509	CO

```
df.set_index('States')
```

		W	X	Y	Z
States					
CA		2.706850	0.628133	0.907969	0.503826
NY		0.651118	-0.319318	-0.848077	0.605965
WY		-2.018168	0.740122	0.528813	-0.589001
OR		0.188695	-0.758872	-0.933237	0.955057
CO		0.190794	1.978757	2.605967	0.683509

df

		W	X	Y	Z	States
A		2.706850	0.628133	0.907969	0.503826	CA
B		0.651118	-0.319318	-0.848077	0.605965	NY
C		-2.018168	0.740122	0.528813	-0.589001	WY
D		0.188695	-0.758872	-0.933237	0.955057	OR
E		0.190794	1.978757	2.605967	0.683509	CO

```
df.set_index('States', inplace=True)
```

```
df
```

	W	X	Y	Z
States				
CA	2.706850	0.628133	0.907969	0.503826
NY	0.651118	-0.319318	-0.848077	0.605965
WY	-2.018168	0.740122	0.528813	-0.589001
OR	0.188695	-0.758872	-0.933237	0.955057
CO	0.190794	1.978757	2.605967	0.683509

▼ Multi-Index and Index Hierarchy

Let us go over how to work with Multi-Index, first we'll create a quick example of what a Multi-Indexed DataFrame would look like:

```
# Index Levels
outside = ['G1', 'G1', 'G1', 'G2', 'G2', 'G2']
inside = [1,2,3,1,2,3]
hier_index = list(zip(outside,inside))
hier_index = pd.MultiIndex.from_tuples(hier_index)
```

```
hier_index
```

```
MultiIndex([('G1', 1),
            ('G1', 2),
            ('G1', 3),
            ('G2', 1),
            ('G2', 2),
            ('G2', 3)],
           )
```

```
df = pd.DataFrame(np.random.randn(6,2), index=hier_index, columns=['A', 'B'])
df
```

		A	B
G1	1	0.302665	1.693723
	2	-1.706086	-1.159119
	3	-0.134841	0.390528

Now let's show how to index this! For index hierarchy we use `df.loc[]`, if this was on the columns axis, you would just use normal bracket notation `df[]`. Calling one level of the index returns the sub-dataframe:

```
df.loc['G1']
```

	A	B
1	0.302665	1.693723
2	-1.706086	-1.159119
3	-0.134841	0.390528

```
df.loc['G1'].loc[1]
```

```
A    0.302665
B    1.693723
Name: 1, dtype: float64
```

```
df.index.names
```

```
FrozenList([None, None])
```

```
df.index.names = ['Group', 'Num']
```

```
df
```

```

                A      B
df.xs('G1')

                A      B
Num
1    0.302665  1.693723
2   -1.706086 -1.159119
3   -0.134841  0.390528

```

```
df.xs(['G1',1])
```

```

<ipython-input-109-c549ee06ce91>:1: FutureWarning: Passing lists as key for xs is depre
df.xs(['G1',1])
A    0.302665
B    1.693723
Name: (G1, 1), dtype: float64

```

```
df.xs(1,level='Num')
```

```

                A      B
Group
G1    0.302665  1.693723
G2    0.166905  0.184502

```

▼ Missing Data

Let's show a few convenient methods to deal with Missing Data in pandas:

```
import numpy as np
import pandas as pd
```

```
df = pd.DataFrame({'A':[1,2,np.nan],
                    'B':[5,np.nan,np.nan],
                    'C':[1,2,3]})
```

```
df
```

	A	B	C
0	1.0	5.0	1
1	2.0	NaN	2
2	NaN	NaN	3

```
df.dropna()
```

	A	B	C
0	1.0	5.0	1

```
df.dropna(axis=1)
```

	C
0	1
1	2
2	3

```
df.dropna(thresh=2)
```

	A	B	C
0	1.0	5.0	1
1	2.0	NaN	2

```
df.fillna(value='FILL VALUE')
```

	A	B	C
0	1.0	5.0	1
1	2.0	FILL VALUE	2
2	FILL VALUE	FILL VALUE	3

```
df['A'].fillna(value=df['A'].mean())
```

0	1.0
1	2.0
2	1.5

Name: A, dtype: float64

V. Merging, Joining and Concatenating

1. Concatenation

```
import pandas as pd
```

```
df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'],
                    'B': ['B0', 'B1', 'B2', 'B3'],
                    'C': ['C0', 'C1', 'C2', 'C3'],
                    'D': ['D0', 'D1', 'D2', 'D3']},
                    index=[0, 1, 2, 3])
```

```
df2 = pd.DataFrame({'A': ['A4', 'A5', 'A6', 'A7'],
                    'B': ['B4', 'B5', 'B6', 'B7'],
                    'C': ['C4', 'C5', 'C6', 'C7'],
                    'D': ['D4', 'D5', 'D6', 'D7']},
                    index=[4, 5, 6, 7])
```

```
df3 = pd.DataFrame({'A': ['A8', 'A9', 'A10', 'A11'],
                    'B': ['B8', 'B9', 'B10', 'B11'],
                    'C': ['C8', 'C9', 'C10', 'C11'],
                    'D': ['D8', 'D9', 'D10', 'D11']},
                    index=[8, 9, 10, 11])
```

df1

	A	B	C	D
0	A0	B0	C0	D0
1	A1	B1	C1	D1
2	A2	B2	C2	D2
3	A3	B3	C3	D3

df2

	A	B	C	D
4	A4	B4	C4	D4
5	A5	B5	C5	D5
6	A6	B6	C6	D6
7	A7	B7	C7	D7

df3

	A	B	C	D
8	A8	B8	C8	D8
9	A9	B9	C9	D9
10	A10	B10	C10	D10
11	A11	B11	C11	D11

```
pd.concat([df1,df2,df3])
```

	A	B	C	D
0	A0	B0	C0	D0
1	A1	B1	C1	D1
2	A2	B2	C2	D2
3	A3	B3	C3	D3
4	A4	B4	C4	D4
5	A5	B5	C5	D5
6	A6	B6	C6	D6
7	A7	B7	C7	D7
8	A8	B8	C8	D8
9	A9	B9	C9	D9
10	A10	B10	C10	D10
11	A11	B11	C11	D11

```
pd.concat([df1,df2,df3],axis=1)
```


	A	B	C	D	A	B	C	D	A	B	C	D
0	A0	B0	C0	D0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
1	A1	B1	C1	D1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	A2	B2	C2	D2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	A3	B3	C3	D3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN	A4	B4	C4	D4	NaN	NaN	NaN	NaN

2. Merging

⌵ NaN NaN NaN NaN A0 B0 C0 D0 NaN NaN NaN NaN

```
left = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],
                      'A': ['A0', 'A1', 'A2', 'A3'],
                      'B': ['B0', 'B1', 'B2', 'B3']})

right = pd.DataFrame({'key': ['K0', 'K1', 'K2', 'K3'],
                      'C': ['C0', 'C1', 'C2', 'C3'],
                      'D': ['D0', 'D1', 'D2', 'D3']})
```

left

	key	A	B
0	K0	A0	B0
1	K1	A1	B1
2	K2	A2	B2
3	K3	A3	B3

right

	key	C	D
0	K0	C0	D0
1	K1	C1	D1
2	K2	C2	D2
3	K3	C3	D3

```
pd.merge(left,right,how='inner',on='key')
```

	key	A	B	C	D
0	K0	A0	B0	C0	D0
1	K1	A1	B1	C1	D1
2	K2	A2	B2	C2	D2

```
left = pd.DataFrame({'key1': ['K0', 'K0', 'K1', 'K2'],
                     'key2': ['K0', 'K1', 'K0', 'K1'],
                     'A': ['A0', 'A1', 'A2', 'A3'],
                     'B': ['B0', 'B1', 'B2', 'B3']})

right = pd.DataFrame({'key1': ['K0', 'K1', 'K1', 'K2'],
                     'key2': ['K0', 'K0', 'K0', 'K0'],
                     'C': ['C0', 'C1', 'C2', 'C3'],
                     'D': ['D0', 'D1', 'D2', 'D3']})
```

```
pd.merge(left, right, on=['key1', 'key2'])
```

	key1	key2	A	B	C	D
0	K0	K0	A0	B0	C0	D0
1	K1	K0	A2	B2	C1	D1
2	K1	K0	A2	B2	C2	D2

```
pd.merge(left, right, how='outer', on=['key1', 'key2'])
```

	key1	key2	A	B	C	D
0	K0	K0	A0	B0	C0	D0
1	K0	K1	A1	B1	NaN	NaN
2	K1	K0	A2	B2	C1	D1
3	K1	K0	A2	B2	C2	D2
4	K2	K1	A3	B3	NaN	NaN
5	K2	K0	NaN	NaN	C3	D3

```
pd.merge(left, right, how='right', on=['key1', 'key2'])
```

	key1	key2	A	B	C	D
0	K0	K0	A0	B0	C0	D0

```
pd.merge(left, right, how='left', on=['key1', 'key2'])
```

	key1	key2	A	B	C	D
0	K0	K0	A0	B0	C0	D0
1	K0	K1	A1	B1	NaN	NaN
2	K1	K0	A2	B2	C1	D1
3	K1	K0	A2	B2	C2	D2
4	K2	K1	A3	B3	NaN	NaN

3. Joining

```
left = pd.DataFrame({'A': ['A0', 'A1', 'A2'],
                     'B': ['B0', 'B1', 'B2']},
                    index=['K0', 'K1', 'K2'])
```

```
right = pd.DataFrame({'C': ['C0', 'C2', 'C3'],
                     'D': ['D0', 'D2', 'D3']},
                    index=['K0', 'K2', 'K3'])
```

```
left.join(right)
```

	A	B	C	D
K0	A0	B0	C0	D0
K1	A1	B1	NaN	NaN
K2	A2	B2	C2	D2

```
left.join(right, how='outer')
```

	A	B	C	D
K0	A0	B0	C0	D0
K1	A1	B1	NaN	NaN
K2	A2	B2	C2	D2
K3	NaN	NaN	C3	D3

VI. Operations

```
import pandas as pd
df = pd.DataFrame({'col1':[1,2,3,4], 'col2':[444,555,666,444], 'col3':['abc','def','ghi','xyz']}
df.head()
```

	col1	col2	col3
0	1	444	abc
1	2	555	def
2	3	666	ghi
3	4	444	xyz

1. Info on Unique Values

```
df['col2'].unique()

array([444, 555, 666])
```

```
df['col2'].nunique()

3
```

```
df['col2'].value_counts()

444    2
555    1
666    1
Name: col2, dtype: int64
```

2. Selecting Data

```
#Select from DataFrame using criteria from multiple columns
newdf = df[(df['col1']>2) & (df['col2']==444)]
```

```
newdf
```

	col1	col2	col3
3	4	444	xyz

3. Applying Functions

```
def times2(x):
    return x*2
```

```
df['col1'].apply(times2)
```

```
0    2
1    4
2    6
3    8
Name: col1, dtype: int64
```

```
df['col3'].apply(len)
```

```
0    3
1    3
2    3
3    3
Name: col3, dtype: int64
```

```
df['col1'].sum()
```

```
10
```

Permanently Removing a Column

```
del df['col1']
```

```
df
```

	col2	col3
0	444	abc
1	555	def
2	666	ghi
3	444	xyz

Get column and index names:

```
df.columns
```

```
Index(['col2', 'col3'], dtype='object')
```

```
df.index
```

```
RangeIndex(start=0, stop=4, step=1)
```

Sorting and Ordering a DataFrame:

```
df
```

	col2	col3
0	444	abc
1	555	def
2	666	ghi
3	444	xyz

```
df.sort_values(by='col2') #inplace=False by default
```

	col2	col3
0	444	abc
3	444	xyz
1	555	def
2	666	ghi

Find Null Values or Check for Null Values

```
df.isnull()
```

	col2	col3
0	False	False
1	False	False
2	False	False
3	False	False

```
# Drop rows with NaN Values  
df.dropna()
```

	col2	col3
0	444	abc
1	555	def
2	666	ghi
3	444	xyz

Filling in NaN values with something else:

```
import numpy as np  
  
df = pd.DataFrame({'col1':[1,2,3,np.nan],  
                  'col2':[np.nan,555,666,444],  
                  'col3':['abc','def','ghi','xyz']})  
df.head()
```

	col1	col2	col3
0	1.0	NaN	abc
1	2.0	555.0	def
2	3.0	666.0	ghi
3	NaN	444.0	xyz

```
df.fillna('FILL')
```

	col1	col2	col3
0	1.0	FILL	abc
1	2.0	555.0	def
2	3.0	666.0	ghi
3	FILL	444.0	xyz

```
data = {'A': ['foo', 'foo', 'foo', 'bar', 'bar', 'bar'],
        'B': ['one', 'one', 'two', 'two', 'one', 'one'],
        'C': ['x', 'y', 'x', 'y', 'x', 'y'],
        'D': [1, 3, 2, 5, 4, 1]}
```

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

```
df
```

	A	B	C	D
0	foo	one	x	1
1	foo	one	y	3
2	foo	two	x	2
3	bar	two	y	5
4	bar	one	x	4
5	bar	one	y	1

```
df.pivot_table(values='D',index=['A', 'B'],columns=['C'])
```

		C	x	y
A	B			
bar	one	4.0	1.0	
	two	NaN	5.0	
foo	one	1.0	3.0	
	two	2.0	NaN	

VII. Data Input and Output

```
import numpy as np
import pandas as pd
```

```
df = pd.read_csv('example')
df
```


	a	b	c	d
0	0	1	2	3

▼ CSV Output

```
df.to_csv('example',index=False)
```

▼ Excel

Pandas can read and write excel files, keep in mind, this only imports data. Not formulas or images, having images or macros may cause this read_excel method to crash.

▼ Excel Input

```
pd.read_excel('/Excel_Sample.xlsx', sheet_name='Sheet1')
```

	a	b	c	d
0	0	1	2	3
1	4	5	6	7
2	8	9	10	11
3	12	13	14	15

```
df.to_excel('Excel_Sample.xlsx',sheet_name='Sheet1')
```

```
from sqlalchemy import create_engine
```

```
engine = create_engine('sqlite:///memory:')
```

```
df.to_sql('data', engine)
```

```
4
```

```
sql_df = pd.read_sql('data', con=engine)
```

```
sql_df
```

	index	a	b	c	d
0	0	0	1	2	3
1	1	4	5	6	7
2	2	8	9	10	11
3	3	12	13	14	15

