

## Pandas important concepts

- \* machine can only understand numbers and not the strings
- \* Here the species column is string so we need to convert them to numerical.  
so for this we can assume

$$\text{Iris-setosa} = 0$$

$$\text{Iris-Versicolor} = 1$$

$$\text{Iris-Virginica} = 2$$

Here we will learn advanced functions like:

- 1) map
- 2) apply
- 3) apply map

map(): It is a series function to do mapping.

- \* cannot be used on DataFrames.
- \* This is different from map in Python.
- \* we need to access the series and map on it.

`data['species'].unique()`

↳ `array(['Iris-setosa', 'Iris-Versicolor', 'Iris-Virginica'],  
 dtype=object)`

- \* Now if we want give the species their respective numbers we use `map()` function.

```
datac["species"].map({'Iris-setosa':0, 'Iris-Versicolor':1,  
                      'Iris-Virginica':2})
```

0 0  
1 0  
2 0  
⋮

Name: species, Length: 150, dtype=int64

This is a temporary change so if we want to use it we can save it and use it.

```
datac["species"] = datac["species"].map({'Iris-setosa':0,  
                                         'Iris-Versicolor':1,'Iris-Virginica':2})
```

```
datac["species"].unique()  
array([0, 1, 2], dtype=int64)
```

**apply():** It is both series and DataFrame function.

How apply() works on series

\* apply() takes any function as parameter

{  
 user defined  
 inbuilt  
 anonymous } input parameter

apply() applies this function elementwise to the series.

\* The function used inside apply() will not have '()' as apply calls the function for it.

```
datac[“species”].apply(type)
```

```
↳ 0 <class ‘int’>  
| <class ‘int’>  
| :  
| :  
| :
```

```
Name: species, length: 150, dtype: object
```

# convert SepalLengthcm to complex.

```
datac[“SepalLengthcm”].apply(complex)
```

```
↳ 0 5.1+0.0j  
| 4.9+0.0j  
| :  
| :
```

```
Name: SepalLengthcm, length: 150, dtype: Complex128
```

\* we can create our own function and use it

```
ex:- def sq(x):  
    return x**2
```

\* return is mandatory as otherwise we cannot use result.

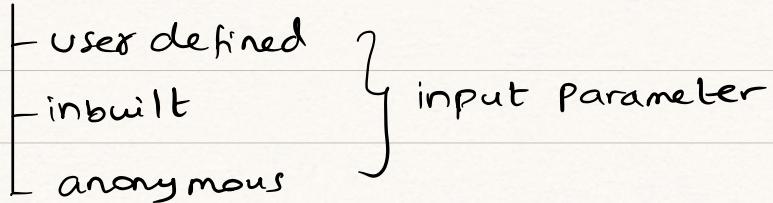
```
datac[“SepalLengthcm”].apply(sq)
```

```
↳ 0 26.01  
| 24.01  
| :  
| :
```

```
Name: SepalLengthcm, length: 150, dtype= float64
```

## apply() function for DataFrame

\* apply() takes any function as parameter



\* apply() applies the function axis wise

[ that is one row or one column at a time.  
this is not element wise.

\* So we need to use the functions that we can apply to the complete row or column.

ex:- `datac[['SepalLengthcm']].apply(int)`

↓  
X error as int function needs to be applied element wise not to entire row or column.

\* double square brackets indicate DataFrame.

`datac[['SepalLengthcm']].apply(sum)`

↓ SepalLengthcm 876.5

dtype: float64

\* default for sum is columnwise (axis=0)

SepalLengthcm has data only in single column so even if we do sum rowwise we get original values.

`datac[['SepalLengthcm']].apply(sum, axis=1)`

↓ 0 5.1

length: 150, dtype: float64

`datac[['SepalLengthcm', 'SepalWidthcm']].apply(sum)`



"SepalLengthcm" 876.5

"SepalWidthcm" 458.1

dtype: float64

\* adding "SepalLengthcm" with "SepalWidthcm" when  
axis=1 does not make sense logically.

applymap(): It is a DataFrame function.

\* This is not a Series function.

\* When we pass any function to DataFrame it is  
applied element wise.

`datac[['SepalLengthcm', 'SepalWidthcm']].applymap(type)`

↓ SepalLengthcm    SepalWidthcm

0 <class 'float'> <class 'float'>

1

⋮

⋮

150 rows x 2 columns

`datac[['SepalLengthcm', 'SepalWidthcm']].applymap(lambda x: str(x))`

↓ SepalLengthcm    SepalWidthcm

0 5.0 3.7

1

⋮

150 rows x 2 columns

## String methods in pandas

- \* Pandas has a module called 'str' using which we can apply string methods.
- \* These string methods can only be applied on Series.
- \* All these are applied elementwise.

ex:- Some of the supported methods:

- 1) str.upper()
- 2) str.lower()
- 3) str.capitalize()
- 4) str.title()
- 5) str.isalpha()
- 6) str.isalnum()
- 7) str.isdigit()
- 8) str.startswith()
- 9) str.endswith()
- 10) str.swapcase()
- 11) str.replace()
- 12) str.index()
- 13) str.find()
- 14) str.strip()
- 15) str.rstrip()
- 16) str.lstrip()
- 17) str.partition()

18) str.contains( ) - checks if given string is present or not.

19) str.replace( ) - This cannot be directly applied to data we need to first get translate table assign it to variable and use this in translate method.

20) str.translate( )

ex:- datac["species"].str.upper()

0 IRIS-SETOSA  
1 :  
:

Name: species, length: 150, dtype: object

datac["species"].str.isalpha()

0 False  
1 False  
:  
:

Name: species, length: 150, dtype: object

datac["species"].str.startswith("I")

0 True  
1 True  
:  
:

Name: species, length: 150, dtype: object

\* we can use these methods in conditions

data[.loc[data["species"] . str. ends with ("a") == False]]

↓ Sepallengthcm Sepalwidthcm Petallengthcm petalwidthcm Species  
50 7.0 3.2 4.7 1.4 Iris-Vericolor

:

data[["species"] . str.replace ("i", "I")]

↓ 0 Iris-setosa

:

Name: species, length: 150, dtype: object

data[["species"] . str.partition ("-")]

↓ 0 1 2  
0 Iris - setosa

:

150 rows x 3 columns

\* a = str.maketrans ("In", "123")

data[["species"] . str.translate(a)]

↓ 0 123S-setosa

:

Name: species, length: 150, dtype: object

`datac[“species”].str.contains(“virginica”)`

↳ 0 False

⋮ ⋮

145 True

Name: species, length: 150, dtype: object

Assignment:

`import pandas as pd`

`pd.DataFrame({“name”: [“Subhadra Bhupathiraju”, “pranavi  
kondraguntla”, “Raghav kondraguntla”]})`

↳ name

0 Subhadra Bhupathiraju

1 Pranavi kondraguntla

2 Raghav kondraguntla

`data = pd.DataFrame({“name”: [“Subhadra Bhupathiraju”,  
“pranavi kondraguntla”, “Raghav kondraguntla”]})`

`type(data)`

↳ `pandas.core.frame.DataFrame`

`data[“name”].str.split()`

↳

0 [Subhadra, Bhupathiraju]

1 [Pranavi, kondraguntla]

2 [Raghav, kondraguntla]

# write code such that first column has first name and second column has last name.

hint: Can use apply(lambda(split))

elementwise  
✓ so we use  
python split

data["first"] = data["name"].apply(lambda x:x.split()[0])

data["second"] = data["name"].apply(lambda x:x.split()[1])

data	name	First	Second
0	Subhadra Bhupathiraju	Subhadra	Bhupathiraju
1	Pranavi Kondraguntla	Pranavi	Kondraguntla
2	Raghav Kondraguntla	Raghav	Kondraguntla

## Advanced Concepts in Python Pandas

### Groupby :-

let us assume we are given a DataFrame consisting of Height, weight and Gender.

height	weight	gender
H1	w1	m
H2	w2	f
H3	w3	m
H4	w4	f
H5	w5	f
H6	w6	m

3 feature Variables  
= D6 (data)

when we are doing analysis on this data . If we want to calculate things like.

- Avg height of men
- Avg height of women
- Avg weight of men
- min height of men

for this we can use `data.loc[condition]` but we have advanced concept called Groupby.

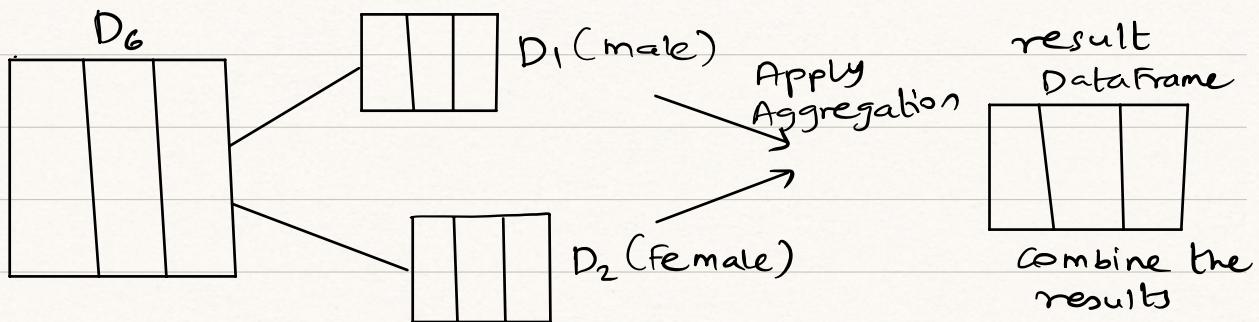
Split - apply - Combine : This refers to chain of three steps :

- 1) split a table (DataFrame) in to groups .
- 2) Apply operations to each of the smaller tables .
- 3) Combine the results .

## Group by :-

we can form groups based on conditions.

we can divide the given data into 2 groups based on a condition (here we can use gender based on the problem statement)



\* How many small DataFrames depends on the unique value of that feature variable.

D1 (DataFrame)		
H1	W1	M
H3	W3	M
H6	W6	M

D2 (DataFrame)		
H2	W2	f
H4	W4	f
H5	W5	f

\* we can apply any aggregations on this like:  
sum, mean, median, ...

\* The result is combined to give one DataFrame.

\* Groupby is mostly done as categorical data  
as we will get a subset of DataFrames.

The concept is groupby.

```
import pandas as pd
```

```
pd.read_csv(r"C:\....\Book4.csv")
```

	Name	Gender	Height	weight
0	P1	m	123	45
1	P2	m	156	67
2	P3	f	125	67
3	P4	f	145	89
4	P5	f	12	67
5	P6	f	133	43
6	P7	m	157	90

data = pd.read\_csv(r"C:\....\Book4.csv")

# **groupby()** function takes the condition as input  
# This is a DataFrame function.

data.groupby("Gender")

→ Groupby object

g = data.groupby("Gender")

\* There are various methods that we can apply on this Groupby object - we will learn a few.

# **n\_groups** gives how many groups are present

g.ngroups

→ 2

# **groups** gives the groups that are there and their row index.

`g.groups`

↳  $\{ 'f': [2, 3, 4, 5], 'm': [0, 1, 6] \}$

# `g` actually internally behaves as dictionary

# `get_group()` we can get that particular group.

`g.get_group("m")`

↳ Name Gender Height weight

0 P1 m 123 45

1 P2 m 156 67

6 P7 m 157 90

`g.get_group("f")`

↳ Name Gender Height weight

2 P3 f 125 67

3 P4 f 145 89

4 P5 f 12 67

5 P6 f 133 43

\* Since GroupBy object acts as dictionary we can use for loop on it.

For keys, values in `g`:

`print(g.keys())`

↳ f  
m

for keys in g :

    print(keys)

    ↳ ('f', Name Gender Height weight)

        2    P3        f        125      67

        3    P4        f        145      89

        4    P5        f        12        67

        5    P6        F        133      43 )

    ('m', Name Gender Height weight)

        0    P1        m        123      45

        1    P2        m        156      67

        6    P7        m        157      90 )

for keys, values in g :

    print(keys)

    print(values)

    print("----" \* 25)

    ↳ f

Name Gender Height weight

        2    P3        f        125      67

        3    P4        f        145      89

        4    P5        f        12        67

        5    P6        F        133      43

        -----

        m

Name Gender Height weight

        0    P1        m        123      45

        1    P2        m        156      67

6 P7 m 157 90

If we want to see only female data

for keys, values in g :

print(1keys)



f

Name Gender Height weight

2	P3	f	125	67
3	P4	f	145	89
4	P5	f	12	67
5	P6	f	133	43

# first data from the sub-dataset

g.first()



Name Height weight

Gender

f	P3	125	67
---	----	-----	----

m

P1

123

45

`g.last()`

↳ Name Height weight

Gender

f	P6	133	43
m	P7	157	90

`g.min()`

↳ Name Height weight

Gender

f	P3	12	43
m	P1	123	45

`g[["Height"]].min()`

↳

Gender

f	12
m	123

Name: Height, dtype: int64

# we can count non-null values in variables

`g.count()`

↳ Name Height weight

Gender

f	4	4	4
m	3	3	3

# max height and weight in male and female  
# these can be of different people.

`g[['Height', 'Weight']].max()`

↓ Height weight

Gender

f	145	89
m	157	90

# Average height and weight in male and female

\* Groupby is very very powerful in pandas.

`g[['Height', 'Weight']].mean()`

↓ Height weight

Gender

f	103.750...	66.50...
m	145.833...	67.33...

# describe() will give details for female and male based on height & weight as these are numerical values.

`g.describe()`

↓ Height

weight

↓ Count mean std min 25% 50% 75% max Count mean std min 25% 50% 75% max

Gender

f	4.0	103.0...	61.71...	12.0	96.75	-	-	-	-	-	-	-
m	3.0	145.33...	19.34...	12.0	139.50	-	-	-	-	-	-	-

```
data.describe()
```

	<u>Height</u>	<u>weight</u>
Count	7.00...	7.00...
mean	121.8714	66.857
min	-	-
25%	-	-
50%	-	-
75%	-	-
max	-	-

# describe by default only gives info for numerical columns.

# if we want to see describe for Categorical columns we can use include\_object.

```
data.describe(include='object')
```

	<u>Name</u>	<u>Gender</u>
Count	7	7
unique	7	2
top	P1	f
freq	1	4

# top is which has highest frequency. If equal frequency then topmost value is used.

\* we can also apply aggregation directly.

`g.agg([{"mean", "median", "count"}])`



Height

weight

Gender	Height			Weight		
	mean	median	count	mean	median	count
f	103.75..	129.0	4	66.50...	67.0	4
m	145.33..	156.0	3	67.33..	67.0	3

\* If we only want to apply aggregations on one particular column.

`g[["Height"]].agg([{"mean", "median", "count", "min", "max"}])`



mean median count min max

Gender	mean	median	count	min	max
f	103.75..	129.0...	4	12	145
m	145.33..	156.0..	3	123	157

\* If we want to apply aggregation column wise differently for different columns.

# aggregation column wise

`g.agg({ "Height": [{"mean", "median"}], "weight": [{"min", "max"}] })`



Height

weight

Gender	mean	median	min	max
f	103.750...	129.0	43	89
m	145.33..	156.0	45	90

Concatenation: Concatenation means combining 2 or more DataFrames into a single DataFrame

- row wise (axis=1)

- column wise (axis=0) (default)

\* Default concatenation is columnwise and uses column labels.

\* When we are doing row wise we will use row indexes, it will check for each row index and concatenate, missing data is marked as NaN.

\* During concatenation index values are also concatenated so we use ignore\_index=True.

```
import pandas as pd
```

```
data = pd.read_csv(r"C:\.....\Book4.csv")
```

```
data
```

	Name	Gender	Height	Weight
0	P1	m	123	45
1	P2	m	156	67
2	P3	f	125	67
3	P4	f	145	89
4	P5	f	12	67
5	P6	f	133	43
6	P7	m	157	90

`dataac = data.copy()`

`Pd.concat([data, dataac])`

	Name	Gender	Height	Weight
0	P1	m	123	45
1	P2	m	156	67
2	P3	f	125	67
3	P4	f	145	89
4	P5	f	12	67
5	P6	f	133	43
6	P7	m	157	90
0	P1	m	123	45
1	P2	m	156	67
2	P3	f	125	67
3	P4	f	145	89
4	P5	f	12	67
5	P6	f	133	43
6	P7	m	157	90

`Pd.concat([data, dataac], ignore_index=True)`

`Name Gender Height Weight`

0	P1	m	123	45
1	P2	m	156	67
2	P3	f	125	67
3	P4	f	145	89

4	P5	f	12	67
5	P6	f	133	43
6	P7	m	157	90
7	P1	m	123	45
8	P2	m	156	67
9	P3	f	125	67
10	P4	f	145	89
11	P5	f	12	67
12	P6	f	133	43
13	P7	m	157	90

Pd.concat ([data, datac], axis=1)

	Name	Gender	Height	Weight	Gender	Height	Weight
0	P1	m	123	45	m	123	45
1	P2	m	156	67	m	156	67
2	P3	f	125	67	f	125	67
3	P4	f	145	89	f	145	89
4	P5	f	12	67	f	12	67
5	P6	f	133	43	f	133	43
6	P7	m	157	90	m	157	90

data.drop ([ "weight" ], inplace=True, axis=1)

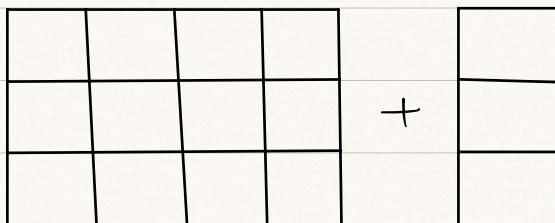
Pd.concat([data, datac], ignore\_index=True)

	Name	Gender	Height	Weight
0	P1	m	123	Nan
1	P2	m	156	Nan
2	P3	f	125	Nan
3	P4	f	145	Nan
4	P5	f	12	Nan
5	P6	f	133	Nan
6	P7	m	157	Nan
7	P1	m	123	45
8	P2	m	156	67
9	P3	f	125	67
10	P4	f	145	89
11	P5	f	12	67
12	P6	f	133	43
13	P7	m	157	90

# if we want data weight to data

# Since we are combining DataFrame with Series

we do combine row wise ... axis=1

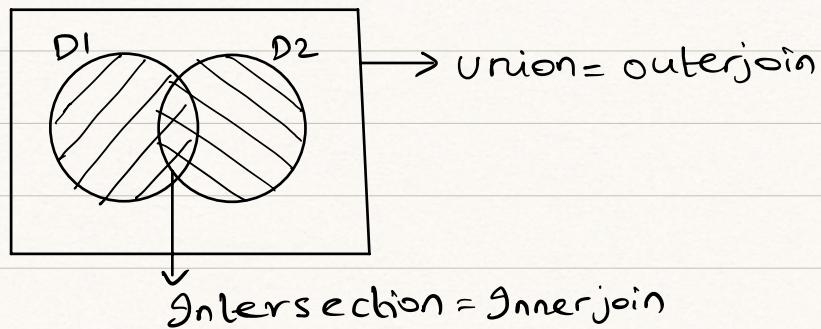


Pd. concat ([data, datac[ "weight" ]], axis=1)

	Name	Gender	Height	weight
0	P1	m	123	45
1	P2	m	156	67
2	P3	f	125	67
3	P4	f	145	89
4	P5	f	12	67
5	P6	f	133	43
6	P7	m	157	90

Merge: we can only merge 2 Dataframes.

\* Merge and concatenation are different  
merging is similar to set theory



\* inner join: intersection (common elements in both DataFrames)

\* outer join: Union (all the elements in 2 DataFrames)  
\* left join: left DataFrame values and intersection  
\* Right join: intersection and right DataFrame values.

import pandas as pd

```
data1 = pd.DataFrame({ "state": ["FL", "KT", "AP", "UP", "TM"],  
                      "temp": [34, 53, 23, 45, 29],  
                      "humd": [34, 54, 34, 56, 57]})
```

```
data2 = pd.DataFrame({ "state": ["FL", "KT", "AP", "UP", "KR"],  
                      "temp": [42, 56, 41, 55, 38],  
                      "humd": [24, 44, 54, 66, 67]})
```

data1

	state	temp	humd
0	tl	34	34
1	kt	53	54
2	ap	23	34
3	up	45	56
4	tm	29	57

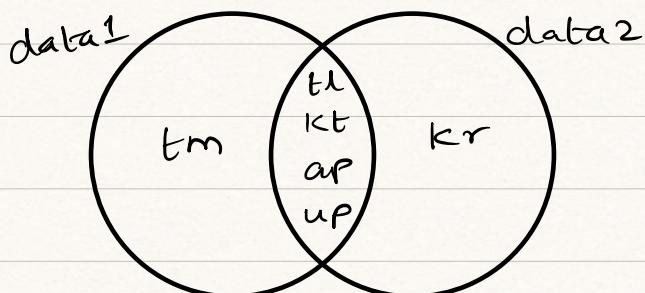
data

	state	temp	humd
0	tl	42	24
1	kt	56	44
2	ap	41	54
3	up	55	66
4	kr	38	67

\* For merge we have to mention based on which column label we have to merge.

\* default is inner join.

\* If we are merging on state.



innerjoin: tl, kt, ap, up

outerjoin: tl, kt, ap, up, tm, kr

pd.merge(data1, data2, on="state") # default inner join

↳ state temp\_x humd\_x temp\_x humd\_x

0	tl	34	34	42	24
1	kt	53	54	56	44
2	ap	23	34	41	54
3	up	45	56	58	66
4	tm	29	57	38	67

\* 'how' means how to join, we have to mention this if we want other joins. 'outer', 'left', 'right'.

pd.merge(data1, data2, on="state", how="outer")

↳ state temp\_x humd\_x temp\_x humd\_x

0	tl	34.0	34.0	42.0	24.0
1	kt	53.0	54.0	56.0	44.0
2	ap	23.0	34.0	41.0	54.0
3	up	45.0	56.0	58.0	66.0
4	tm	29.0	57.0	nan	nan
5	kr	nan	nan	38.0	67.0

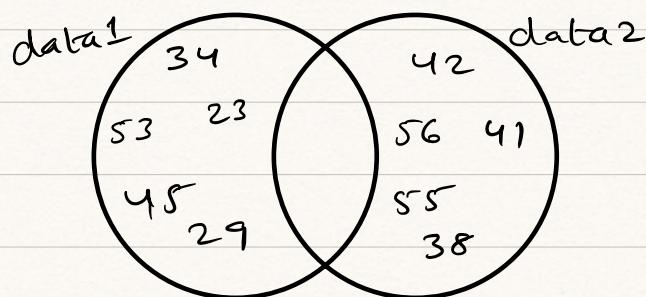
`pd.merge(data1, data2, on="state", how="left")`

	state	temp_x	humd_x	temp_x	humd_x
0	tl	34	34	42.0	24.0
1	kt	53	54	56.0	44.0
2	ap	23	34	41.0	54.0
3	up	45	56	55.0	66.0
4	tm	29	57	nan	nan

`pd.merge(data1, data2, on="state", how='right')`

	state	temp_x	humd_x	temp_x	humd_x
0	tl	34.0	34.0	42	24
1	kt	53.0	54.0	56	44
2	ap	23.0	34.0	41	54
3	up	45.0	56.0	55	66
4	kr	nan	nan	38	67

\* Merge on temp



`pd.merge(data1, data2, on="temp")`

	state_x	temp	humd_x	state_y	humd_y
					#Empty

#Empty

`pd.merge(data1, data2, on="temp", how="outer")`

	state_x	temp	humd_x	state_y	humd_y
0	tl	34	34.0	NaN	NaN
1	kt	53	54.0	NaN	NaN
2	ap	23	34.0	NaN	NaN
3	up	45	56.0	NaN	NaN
4	tm	29	67.0	NaN	NaN
5	nan	42	NaN	tl	24.0
6	nan	56	NaN	kt	44.0
7	nan	41	NaN	ap	54.0
8	nan	55	NaN	up	66.0
9	nan	38	NaN	kr	67.0

\* If we use `indicator=True` it shows where the values are present i.e left-only, both, right-only

`pd.merge(data1, data2, on="state", how="outer", indicator=True)`

	state	temp_x	humd_x	temp_x	humd_x	_merge
0	tl	34.0	34.0	42.0	24.0	both
1	kt	53.0	54.0	56.0	44.0	both
2	ap	23.0	34.0	41.0	54.0	both
3	up	45.0	56.0	58.0	66.0	both
4	tm	29.0	57.0	NaN	NaN	left-only
5	kr	NaN	NaN	38.0	67.0	right-only

## PIVOT & pivot table

\* Pivot and pivot table are different

Pivot	column1	column2	values
C <sub>1</sub>	A <sub>1</sub>	x <sub>1</sub>	
C <sub>2</sub>	A <sub>2</sub>	x <sub>2</sub>	
C <sub>1</sub>	A <sub>3</sub>	x <sub>3</sub>	

transform my column values in to row index  
and column index

	A <sub>1</sub>	A <sub>2</sub>	A <sub>3</sub>
C <sub>1</sub>	x <sub>1</sub>	NaN	x <sub>3</sub>
C <sub>2</sub>	NaN	x <sub>2</sub>	NaN

Pivot we will learn in more detail in statistics

import pandas as pd

import numpy as np

pd.read\_csv(r"C:\....\weather.csv")



	date	city	temperature	humidity
0	5/1/2017	newyork	65	56
1	5/2/2017	newyork	66	58
2	5/3/2017	newyork	68	60
3	5/1/2017	mumbai	75	80
4	5/2/2017	mumbai	78	83

5	5/3/2017	mumbai	82	85
6	5/1/2017	beijing	80	26
7	5/2/2017	beijing	77	30
8	5/3/2017	beijing	79	35

`df = pd.read_csv(r'c:\....\weather.csv')`

\* Here if we want to find what is temperature on a particular date we can use PIVOT.

\* Date will become row Index

\* City will become Column Index.

\* remaining columns will behave as values.

\* Pivot() is a Dataframe function

\* makes sense to use Categorical data in columns  
also index better to use categorical data.

# Sub index of date is created.

`df.pivot(index='city', columns='date')`

city	temperature			humidity		
	5/1/2017	5/2/2017	5/3/2017	5/1/2017	5/2/2017	5/3/2017
beijing	80	77	79	26	30	35
mumbai	75	78	82	80	83	85
new york	65	66	68	56	58	60

# if we want to see only humidity

df.pivot(index='city', columns='date', values='humidity')



humidity

date 5/1/2017 5/2/2017 5/3/2017

City

beijing	26	30	35
mumbai	80	83	85
newyork	56	58	60

(or) we can change the view we want to see

df.pivot(index='date', columns='city')



temperature

humidity

City beijing mumbai newyork beijing mumbai newyork

Date

5/1/2017	80	75	65	26	80	56
5/2/2017	77	78	66	30	83	58
5/3/2017	79	82	68	35	85	60

\* Pivot table: It is a powerful tool to calculate, summarize and analyze data that let's you see comparisons, patterns, and trends in your data.

\* pivot\_table() function creates a spread-sheet style Pivot table as DataFrame.

\* Pivot-table is mainly used to summarize the data.

\* Instead of actual values if we want aggregate values.

# margins = True will give row wise and column wise result, by default will calculate mean.

```
df.pivot_table(index='date', columns='city',  
               margins=True)
```

Date	humidity				temperature			
	beijing	mumbai	newyork	All	beijing	mumbai	newyork	All
5/1/2017	26.0...	80.0...	56	54.0	80.0...	75.0...	68.0...	73.33...
5/2/2017	30.0...	83.0...	58	57.0	77.0...	78.0...	66.0...	73.66...
5/3/2017	38.0...	85.0...	60	60.0	79.0...	82.0...	68.0...	76.33...
All	30.33..	82.66..	58	57.0	78.66..	78.33..	66.33..	74.44...

# we can tell it to apply different aggregations

# we can use name for margins

```
df.pivot_table(index='date', columns='city', margins=True,  
               margins_name='sum', aggfunc=np.sum)
```

↓

Date	humidity					temperature			
	City	beijing	mumbai	newyork	Sum	beijing	mumbai	newyork	Sum
5/1/2017	26	80	56	162	80	75	65	220	
5/2/2017	30	83	58	171	77	78	66	221	
5/3/2017	35	85	60	180	79	82	68	229	
Sum	91	248	174	513	236	235	199	670	

# we can apply different aggregations to different numerical data.

Jupyter notebook giving errors when using one aggregation as std

```
df.pivot_table(index='date', columns='city', margins=True,
               margins_name='sum', aggfunc=[{"humidity": np.mean, "temperature": np.median}])
```

↓

Date	humidity					temperature			
	City	beijing	mumbai	newyork	All	beijing	mumbai	newyork	All
5/1/2017	26.0..	80.0..	56	54.0	80	75	65	75.0	
5/2/2017	30.0..	83.0..	58	57.0	77	78	66	77.0	
5/3/2017	35.0..	85.0..	60	60.0	79	82	68	79.0	
All	30.33..	82.66..	58	57.0	79	79	66	77.0	

Crosstab() function is used to compute a simple cross tabulation of two (or more) factors.

\* if we are given data like

Gender	Color
m	B
m	B
f	W
f	B
m	W

Cross tab applies 'and' operator

	B	W
m	2	1
f	1	1

↳ these are not actual values but counts.

```
import pandas as pd  
import numpy as np  
pd.read_excel(r"C:\....\survey.xls")
```



	Name	Nationality	sex	Age	Handedness
0	Kathy	USA	Female	23	Right
1	Linda	USA	Female	18	Right
2	Peter	USA	Male	19	Right
3	John	USA	Male	22	Left
4	Fatima	Bangladesh	Female	31	Left
5	Kadir	Bangladesh	Male	25	Left
6	Dhaval	India	Male	35	Left
7	Sudhir	India	Male	31	Left
8	Parthir	India	Male	37	Right
9	Yan	china	Female	52	Right
10	Juan	china	Female	58	Left
11	Liang	china	Female	43	Left

```
df = pd.read_excel(r"C:\....\survey.xls")
```

- \* Crosstab is visualization concept when we are having category vs category
- \* x-axis and y-axis should have Categorical data

\* result will be 'and' operator so we will get count

pd.crosstab(df.Nationality, df.Handedness)



Nationality	Handedness	Left	Right
Bangladesh		2	0
China		2	1
India		2	1
USA		1	3

# from this we can summarize that

- In Bangladesh there are 2 left handed people and 0 right handed people.
- In China there are 2 left handed people and 1 right handed person.
- In India there are 2 left handed people and 1 right handed person.
- In USA there are 1 left handed person and 3 right handed people.

pd.crosstab(df.Sex, df.Handedness)



Handedness		left	Right
Sex	Female	2	3
	male	5	2

# from this we can summarize that

- Females are more right handed
- male are more left handed

# we can add margins=True, we will see rowwise  
and column wise margins

pd.crosstab(df.Sex, df.Handedness, margins=True)



Handedness		left	Right	All
Sex	Female	2	3	5
	male	5	2	7
All		7	5	12

# we can have multiple index columns and rows.

pd.crosstab(df.sex, [df.Handedness, df.Nationality], margins=True)

↓

		Handedness			left			Right		All	
		Nationality			Bangladesh	china	India	USA	china	India	USA
		sex									
Female		1		1	0	0	1	0	2	5	
male		1		1	2	1	0	1	1	7	
All		2		2	2	1	1	1	3	12	

(or)

pd.crosstab([df.sex, df.Handedness], df.nationality, margins=True, margins\_name="total")

↓

		Nationality				Bangladesh	china	India	USA	total
Sex	Handedness	Left	Right	Left	Right	Left	Right	Left	Right	total
Female	Left	1		1	0	0	0	2		2
Female	Right	0		0	1	0	2	0	3	3
male	Left	1		1	2	2	1	1	5	5
male	Right	0		0	1	1	1	1	2	2
total		2		3	3	3	4	4	12	12