

Pandas

1/1

- pandas is a python package providing fast, flexible, and expressive data structures designed to make working with relational or labeled data both easy and intuitive

- Data structures in pandas

↳ Series (similar to 1 dimensional numpy array)

↳ DataFrame (similar to 2 dim. numpy array)

- Installation command

- ! pip install pandas (Jupyter Notebook)

- pip install Pandas (Command prompt)

- Importing pandas

- import pandas as pd

Series (ndarray-like)

| s = pd.Series(data, index=index) |

list → data = [1, 2, 3, 4] ↗ python list

dictionary → data = {'a': 1, 'b': 2} ↗ dictionary

numpy array ↗ data = np.array([1, 2, 3])

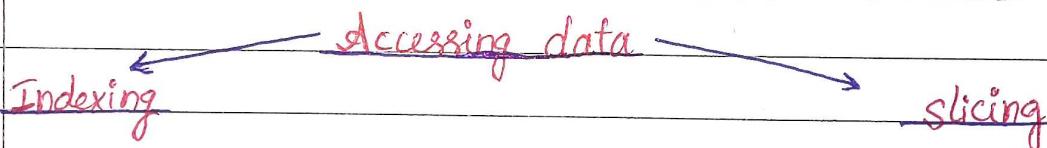
scalar → data = 5, index = ['a', 'b', 'c']

- Homogenous
- One dimensional labeled array
- Holding any data type (int, str, float, date, objects etc)
- Row labels = index
- Value-mutable (the values they contain can be altered) but not size-mutable (immutable)



- | | |
|------------------------------------|--|
| - Data Type : <code>S.dtype</code> | - Head : <code>S.head()</code> |
| - Shape : <code>S.shape</code> | - Tail : <code>S.tail()</code> |
| - Values : <code>S.values</code> | - Info : <code>S.info()</code> |
| - Array : <code>S.array</code> | - To numpy : <code>S.to_numpy()</code> |

`s = pd.Series(data= , index= , dtype= , name=)`



single Index

`s[2]` ← using index value

`s['a']` ← using row label

`s[start : stop : step]`

To slice the original string object from each row in a Series.

Multiple Indexes

`s[[1, 4]]` ② `s[['a', 'c']]`

list of indexes

`s.str.slice(start= , stop= , step=)`

`s.str : [start, stop, step]`

DataFrame

- Heterogeneous typed columns
- Value-mutable as well as size-mutable.
- Each column in df is series.
- Working with Tabular data. (spread sheets stored in database also)

`df = pd.DataFrame(data, index = idxs, columns = cols)`

Create DataFrame

*using
dictionary*

```
[ data = { 'Name' : ['Ravi', 'Shahees'],
      'Age' : [26, 23] } ]
```

*using
tuple*

```
[ data = ((('Ravi', 26), ('Shahees', 23)),  
        columns = ['Name', 'Age']) ]
```

*using
list*

```
[ data = ([['Ravi', 26], ['Shahees', 23]]), columns = ['Name', 'Age']]
```

*using
NumPy
array*

```
[ array = np.array([1, 2, 3], [4, 5, 6]) ]
```

```
[ data = array, columns = [col-1, col-2] ]
```

DataFrame

Attributes / Properties

- Shape of df : `df.shape`
- columns : `df.columns`
- data types : `df.dtypes`
- Axes : `df.axes`
- Values : `df.values`

Methods

- Head : `df.head()`
- Tail : `df.tail()`
- Info : `df.info()`

Read data to pandas

- `pd.read_*` (* = csv, excel, sql, json, parquet etc ...)

Eg :- `pd.read_csv('path')`

Export / Save from pandas

- `df.to_*` (* = csv, excel, sql, json, parquet. etc ...)

Eg :- `df.to_excel('path', sheet-name='name', index=False)`

Note : In pandas default Axis = 0

Axis = 0 → columnwise

Axis = 1 → Rowwise.

Non-visual Data Analysis (Statistical)

- min df.min()
- max df.max()
- mean df.mean()
- median df.median()
- var df.var()
- std df.std()
- skew df.skew()
- kurt df.kurt()
- sum df.sum()
- count df.count()
- describe → df.describe(include = ['object']) ← categorical
 (include = ['numbers']) ← numerical
 (include = 'all') ← All.
- num_df = df.select_dtypes(include = ['int32', 'float64'])
- cat_df = df.select_dtypes(include = ['object'])
- Aggregation :- df[['col-1']].agg(['min', 'max', 'mean', 'median', 'std'])
- df.agg({'col-1': ['min', 'max', 'mean', 'median'],
 'col-2': ['count', 'nunique', 'unique', 'value_counts']})

Indexing and slicing.

- Filtering Single column vs Multiple columns :-

Syntax

- single column : `df['col-name']` ← Returns Series.
`df[['col-name']]` ← Returns DataFrame.

- multiple columns : `df[['col-1', 'col-2', 'col-3']]`

- Filtering Rows :-

- Way - 1 : `df[starting-row-index : ending-row-index : step]`
Eg : `df[1:6]`

- way - 2 : `df[condition]`

Eg : `df[df["max temp"] > 95]`

** `df[df['date'].isin(['10-5-2016', '10-4-2016'])]`

or

`df[(df['date'] == '10-5-2016') | (df['date'] == '10-4-2016')]`

- Filtering Specific Rows & columns from df :-

Syntax :

`df.loc[row-label, column-label]`

`df.iloc[row-index, column-index]`

Eg :-

`df.loc[100, ['date', 'snowfall']]`

`df.iloc[100, [0,5]]`

df.loc[10:15, 'max temp': 'precipitation'] ②
df.iloc[10:15, 1:5]

- Accessing rows based on condition :

df.loc[condition, column_labels]

Eg:- df.loc[df['maxtemp'] > 95, 'date']

- Accessing rows based on multiple conditions :

df.loc[(cond-1) & (cond-2) | (cond-3), column_labels]

Eg:- df.loc[(df['maxtemp'] > 95) & (df['mintemp'] > 80), 'date']

- Renaming columns :

Syntax → `df.rename(index=None, columns=None)`

Eg: → `df = df.rename(columns = {'description': 'product description', 'customers Id': 'cust id'})`

↑ ↑
previous name New Name.

- A very common Renaming strategy :-

- ① strip Extra spaces
- ② convert to lower case
- ③ Remove all the special characters including spaces .

`col-names = [col.strip().lower().replace(' ', '_').replace('-', '_') for col in df.columns]`

`df.columns = col-names`

(or)

`df.columns = df.columns.str.strip().str.replace(' ', '_').str.replace('-', '_').`

Note :

By doing the above column renaming we can access the columns like properties of the data frame like `df.columns, df.shape....`

`df.col-name`

- Modifying columns Datatype:

- using df.astype()

`df['quantity'] = df['quantity'].astype(int)`

(or)

`df = df.astype({'quantity': int, 'country': str})`

- using df.apply()

`df['stock'] = df['stock'].apply(pd.to_numeric)`

`pd.to_datetime` (or)
`pd.to_timedelta`

- Creating a Derived column:

`df['amount'] = df['quantity'] * df['unit-price']`

- Creating columns using apply() function:

data frame `df.apply(function, axis=1)`

Series `s.apply(function, axis=1)`

using Anonymous function {
`df['amount'] = df.apply(lambda row: row['quantity'] * row['unit.price'], axis=1)`

- Handling TimeSeries data:

pandas has Great Support for time series and has an extensive set of tools for working with dates, times and time-indexed etc.

<u>datetime</u>		
<u>properties</u>		<u>methods</u>
- date - df['date'].dt.day		- day name - .dt.day_name()
- month -	• dt.month	- month name - .dt.month_name()
- year -	• dt.year	
- day_of_week -	• dt.day_of_week	
- day_of_year -	• dt.day_of_year	
- is_leap_year -	• dt.is_leap_year	
- week -	• dt.week	

<u>time delta</u>		
<u>properties</u>		<u>methods</u>
→ difference b/w two dates (OS)		
- days	• dt.days	- total seconds
- hours	• dt.hours	• dt.total_seconds()
- minutes	• dt.minutes	
- Seconds	• dt.seconds	
- components	• dt.components	

- days	• dt.days	- total seconds	• dt.total_seconds()
- hours	• dt.hours		
- minutes	• dt.minutes		
- Seconds	• dt.seconds		
- components	• dt.components		

- pd.to_datetime()

This function helps us to convert the 'date' column to its actual datatype datetime.

`df['date'] = df['date'].apply(pd.to_datetime())`

- using this function, we end up with many warnings because of format.

- These warnings are generated for a reason. Since dates can be specified in various formats for eg:
`DD/MM/YYYY` (or) `YYYY/MM/DD` (or) `MM/DD/YYYY` etc.

- Two ways to get rid of warnings:

Way - 1 add parameter `dayfirst = True`

Way - 2 add parameter `format = "%d/%m/%Y"`

Note: Best Suggested way to convert 'date' to it's actual data type ↴ (Recommended : Kanav sir)

`df['date'] = df['date'].apply(pd.datetime())`

*** ✓ `df['date'] = pd.to_datetime(df['date'], dayfirst=True)`

(or)
`format = "%d/%m/%Y"`

- while Reading / Loading data

** ✓ `df = pd.read_csv(PATH, parse_dates = ['cols'])`

today = pd.to_datetime('today')
df['age'] = today.year - df['DOB'].dt.year

+ Improve performance by Setting date column as the Index:

df = df.set_index(['date']).
modifying the index inplace

df = set_index(['date'], inplace=True)

Select data with a specific year & perform aggregation

df.loc['2018'] → Select data with Specific year
df.loc['2018-5-1'] → Specific date
df.loc['2018-5-1':'2018-5-4'] → slicing
df.loc['2018-5-1':'2018-5-5', ['sales']].mean() →
agg. with slice

④ sort_index :

df.sort_index(ascending=False)

④ sort_values :

df.sort_values(by='col-name')

④ reset_index :

df.reset_index()

Intermediate pandas

1 / 1

Note: `groupby()` and `pivot_table()` are very powerful in analyzing and summarizing the data.

- Group by :

`grouped_df = df.groupby('category')` ← splits based on category.

`grouped_df = df.groupby(['category', 'sub-category'])`

It splits the

df based on 'category' first
and then 'sub-category'.

* Groups :

`grouped_df.groups`

- The "groups" attribute is a dictionary whose keys are the computed unique Groups and corresponding values are the axis labels belonging to each group.
- Here we can apply dictionary Methods.

`grouped_df.groups.keys()` → Returns unique Group keys.

Filter data on the basis of Group Keys :

`grouped_df.get_group("Group-key-name")` ↴

Grouped on category only.

`grouped_df.get_group(("Group-category", "Group-sub-category"))` ↴

Grouped based on category and Sub-Category .

+ Returning First row, Last row and nth row for each group

- First row of each group :

`grouped_df.first()`

- Last row of each group :

`grouped_df.last()`

- nth row of each group :

`grouped_df.nth(n)` $n = 1, 2, 3, 4, 5, \dots$

Note : The groupby method is used to support this type of operations.

This fits in the more general,
split - apply - combine

- Split the data into groups.
- apply a function to each group independently.
- combine the results into a data structure.

- In the apply step, we might wish to do one of the following :

- Aggregation : Compute a Summary statistic for each group

- Filtration: Discard some groups, according to a group-wise computation that evaluates to True or False.
- Transformation: perform some group-specific computations and return a like-indexed object

Aggregation:

Built-in Aggregation methods

Eg: `grouped_df['Sales'].min()`

- `any()` - compute whether any of the values in the groups are truthy.
- `all()` - compute whether all of the values in the groups are truthy.
- `count()` - computes the non-NA values in the groups.

- `cov()`* - compute the co-variance of the Groups
- `first()` - compute the first occurring value in the Group
- `last()` - compute the last occurring value in each Group
- `nth()` - compute the nth occurring value in each Group

`idxmax()*` - compute the index of maximum value in each group

`idxmin()*` - index of minimum value in each group.

`min()` - minimum value in Each Group

`max()` - maximum value in Each Group

`mean()` - mean of Each Group.

`median()` - median of Each Group.

- `nunique()` - No. of unique values in Each Group.
- `prod()` - product of values in Each Group.
- `quantile()` - a Given quantile of the values in each group.
- `Sem()` - The standard errors of the mean of the values in each group.
- `size()` - The no. of values in each group.
- `Skew()*` - The skew of the values in each group.
- `sum()` - The sum of the values in each group.
- `Vars()` - The variance of the values in each group.
- `Std()` - The Std of the values in each group.

User-defined Aggregation

Aggregation on Single Column:

```
grouped_df['sales'].agg(['min', 'max', 'mean'])
```

Different Aggregations on different columns :

```
grouped_df.agg({'sales': ['mean', 'std'],
                 'order-date': ['min', 'max']})
```

Filtration:Built-in Filtration

- `head()` - Select the top row(s) of each group
- `tail()` - bottom row(s)
- `nth()` - nth row(s)

Eg: `grouped_df.head()`

User-defined Filtration

`grouped_df.filter(lambda group: group['sales'].mean > 20)`

Transformation:

unlike aggregations, the groupings that are used to split the original object are not included in the result.

Built-in Transformation

- `bfill()` - Back fill NA values within each group.
- `cumcount()` - The cumulative count within each group.
- `cummax()` - The cumulative max within each group.
- `cummin()` - The cumulative min within each group.
- `cumprod()` - The cumulative product within each group.
- `cumsum()` - The cumulative sum within each group.
- `diff()` - The difference b/w adjacent values within each group.
- `ffill()` - Forward fill NA values within each group
- `fillna()` - Fill NA values within each group.

- pct_change() - compute the percent change b/w adjacent values within each group.
- rank() - compute the rank of each value within each group.
- shift() - shift values up ⬆ down ⬇ within each group.

user-defined Transformation:
grouped_df.transform (lambda x: x+1)

Reshaping Layout of Tables:

Q: How to reshape the layout of tables?

A: - Change the structure of your data table in multiple ways.

- you can `melt()` your data table from wide to long/tidy form or `pivot_table()` from long to wide format.

+ `pivot_table()` - used to Summarize and Analyse data.

① `df.pivot_table(values = "Sales", index = ["Region"], aggfunc = ['sum', 'count'])`

② `df.pivot_table(values = 'sales', index = ['Region'], aggfunc = ['sum'], columns = ['Category']).apply(lambda val: val * 100 / sum(val))`

category	color	Region	units	sales	category		
					Region	Tech	fancy
Tech	Red	West	1	\$11	West	\$11	\$26
Tech	Blue	South	8	\$96	South	\$96	\$22
fancy	Green	West	2	\$26	North	\$84	\$104
Tech	Blue	North	7	\$84	East	\$156	0
fancy	Green	North	8	\$104			
fancy	Red	South	2	\$22			
Tech	Blue	East	5	\$60			
Tech	silver	East	8	\$96			

pivot-table

Delete column(s) in a DataFrame :

Syntax-1 : Dropping columns by using column name:

dropping two columns by passing column names.

`inplace=True` parameter performs the operation saves the result back to the data frame.

```
df.drop(['col1', 'col3'], axis=1, inplace=True)
```

Syntax-2 : Removing columns by using columns name using .loc[]:

Removing columns b/w col2 and col4

```
df.drop(df.loc[:, 'col2':'col4'], inplace=True)
```

Syntax-3 : Removing columns based on index:

Remove three columns as index base.

```
df.drop(df.columns[[0, 4, 2]], axis=1, inplace=True)
```

Syntax-4 : Removing columns based on index using .iloc[]:

Removing two columns b/w column index 1 to 3.

```
df.drop(df.iloc[:, 1:3], inplace=True)
```

Syntax-5 : DataFrame.pop() method:

using `pop()` we can delete single column at a time

```
df.pop('col4', inplace=True)
```

Adding / inserting Rows ↴

Syntax-1: Inserting a single Row

create a new record using Dictionary.

- new_record = pd.DataFrame([{'day': '1/7/2017', 'temp': 36, 'windspeed': 4, 'event': 'sunny'}])

inserting row at the end ↴

- df = pd.concat([df, new_record], ignore_index=True)

inserting row at the top ↴

- df = pd.concat([new_record, df], ignore_index=True)

Syntax-2: Inserting multiple rows (i.e; a batch of data)

create a new record(s) using dictionary.

- batch_records = pd.DataFrame([{'day': '1/8/2017', 'temp': 30, 'windspeed': 3, 'event': 'Rain'}, {'day': '1/9/2017', 'temp': 27, 'windspeed': 4, 'event': 'snow'}])

inserting at the end ↴

- df = pd.concat([df, batch_records], ignore_index=True)

inserting at the top ↴

- df = pd.concat([batch_records, df], ignore_index=True)

Syntax-3 : Inserting a row using df.loc[]

`df.loc[len(df)] = ['1/12/2017', 28, 2, 'Rain']`

Note : Using `df.iloc[]` we can't Insert a row.

It Generates Error - you can't use `iloc[]` to Enlarge the target object (i.e; `iloc[]` can't be used to add new records), it can be used for updating and overwriting the Existing row.

Syntax-4 : Inserting a row at Specific index

adding at row label 8.5

`df.loc[8.5] = ['1/11/2017', 30, 3, 'Rain']`

Sort & Reset Index

`df = df.sort_index().reset_index(drop=True)`

combining data from Multiple tables using pd.merge() function:

- Remember:
- ① Multiple tables can be concated both column-wise and row-wise using the concat function
 - ② For data-base like merging/joining of tables, use the merge function.

Syntax: pd.merge(left, right, how='inner', on=None, left_on=None, right_on=None, left_index=False, right_index=False, sort=False)

Note:

- * By default **how = 'inner'**
- 'left'
- 'right'
- 'outer'
- 'cross'

pandas

SQL

- Inner join
- left outer join
- right outer join
- full outer join
- Cross join

* cross-join don't need '**On**' parameter

Brief primer on merge methods (Relational Algebra):

- Experienced users of relational databases like SQL will be familiar with the terminology used to describe join operations b/w two SQL-table like structures (dataframes).
- There are several cases to consider which are very important to understand:

① One-to-one joins: for example when joining two data frame objects on their indexes (which must contain unique values)

② many-to-one joins: for example when joining an index (unique) to one or more columns in a different data frame.

③ many-to-many joins: joining columns on columns.

Visual Data Analysis using pandas (plotting) :

Syntax:

`df.plot(kind = "kind")`

The kind of plot to produce :

- ① line : line plot (default 'kind')
- ② bar : vertical bar plot
- ③ barh : horizontal bar plot
- ④ hist : Histogram
- ⑤ box : Box plot
- ⑥ kde : kernel density Estimation plot
- ⑦ density : same as "kde"
- ⑧ area : area plot
- ⑨ pie : pie plot
- ⑩ scatter : scatter plot
- ⑪ hexbin : hexbin plot

Note: ① whenever we draw 'bar' plot / categorical column follow the syntax below ↴

`df['Species'].value_counts().plot(kind = "bar")`

combining data from multiple tables using pd.merge()

Examples :

	left				right			
	key 1	key 2	A	B			C	D
0	K0	K0	A0	B0			C0	D0
1	K0	K1	A1	B1			C1	D1
2	K1	K0	A2	B2			C2	D2
3	K2	K1	A3	B3			C3	D3

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

Innner Join :

	key 1	key 2	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

Left Join : how='left'

	key 1	key 2	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

Right Join : how='right'

	key 1	key 2	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
3	K2	K0	Nan	Nan	C3	D3

11

Outer Join : how = 'outer'

	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

Cross Join : how = 'cross'

Note: for cross join "on" parameter is not required

	key1-x	key2-x	A	B	key1-y	key2-y	C	D
0	K0	K0	A0	B0	K0	K0	C0	D0
1	K0	K0	A0	B0	K1	K0	C1	D1
2	K0	K0	A0	B0	K1	K0	C2	D2
3	K0	K0	A0	B0	K2	K0	C3	D3
4	K0	K1	A1	B1	K0	K0	C0	D0
5	K0	K1	A1	B1	K1	K0	C1	D1
6	K0	K1	A1	B1	K1	K0	C2	D2
7	K0	K1	A1	B1	K2	K0	C3	D3
8	K1	K0	A2	B2	K0	K0	C0	D0
9	K1	K0	A2	B2	K1	K0	C1	D1
10	K1	K0	A2	B2	K1	K0	C2	D2
11	K1	K0	A2	B2	K2	K0	C3	D3
12	K2	K1	A3	B3	K0	K0	C0	D0
13	K2	K1	A3	B3	K1	K0	C1	D1
14	K2	K1	A3	B3	K1	K0	C2	D2
15	K2	K1	A3	B3	K2	K0	C3	D3

- pd.melt():

```
data = {'Name': ['Ravi', 'shahees'],
        'math': [90, 95],
        'English': [85, 88]}
```

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

df

Name	math	English
Ravi	90	85
shahees	95	88

```
melted_df = pd.melt(df, id_vars=['Name'],
                     var_name=['subject'],
                     value_name='Score')
```

print(melted_df)

	Name	Subject	Score
0	Ravi	Math	90
1	Shahees	Math	95
2	Ravi	English	85
3	Shahees	English	88

- pd.explode() :

```
data = {'Name': ['Ravi', 'Abhi'],
        'Subjects': [['Math', 'English'], ['Science', 'History']]}
```

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

```
df
```

Name	Subjects
Ravi	['Math', 'English']
Abhi	['Science', 'History']

```
exploded_df = df.explode('Subjects')
```

```
print(exploded_df)
```

	Name	Subjects
0	Ravi	Math
1	Ravi	English
2	Abhi	Science
3	Abhi	History.