

DATE:

Global Transformation of Series by using map() method:-

Syntax: Series.map(arg, na_action=None)

eg: S = pd.Series(['cat', 'dog', np.nan, 'rat'])

S1 = S.map({'cat': 'kitten', 'dog': 'puppy'})

print(S1)

→ Iterating Elements of the Series:- Iterating means getting element one by one.

S = pd.Series([10, 20, 30, 40])

for i, v in S.items():

print(i, '-->', v)

Chapter-2 Data Frame

Data-Frame is a labeled 2-D array. It Contains multiple rows & columns.

eg: eno = [10, 20, 30, 40]

ename = ['A', 'B', 'C', 'D']

esal = [1000, 2000, 3000, 4000]

eadr = ['Hyd', 'Bang', 'Chennai', 'Ker']

df = pd.DataFrame({'eno': eno, 'ename': ename, 'esal': esal, 'eadr': eadr})

print(df)

print(type(df))

print(df.ndim) print(df.shape)

DATE: [.] [] [] [] []

Conclusion:

- ① Series is a labeled 1-D array where as DataFrame is labeled 2D array.
- ② Each Column of DataFrame acts as Series. Hence we can consider DataFrame is nothing but Collection of Series objects.
- ③ Series object holds only homogeneous values. But DataFrame contains multiple columns & each column values may be of different types.
Hence DataFrame holds heterogeneous data.

Syntax: pd.DataFrame(data=None, index=None, columns=None, dtype=None, copy=True)

Data-Frame is a Two-dimensional, Size-Mutable, Potentially heterogeneous tabular data.

Various ways to Create dataframe object:

- ① From dict of List objects. (Common way of creating data-frames)
- ② From dict of Tuple objects.
- ③ From dict of Series object.
- ④ From dict of dicts.
- ⑤ from NumPy ndarray.

Note: 'dtype' attribute only applicable for Series object but not for DataFrame. We have to use 'dtypes' attribute for DataFrame object.

DATE:

→ Creating DataFrame from numpy ndarray :-

```
df = pd.DataFrame(np.array([[1, 2, 3], [4, 5, 6]]),  
                  columns=['a', 'b', 'c'])  
print(df)
```

→ Common Methods & Attributes of Series & DataFrame :-

Most of methods & attributes of Series are applicable for DataFrame also.

eg: ① df.head() ③ df.shape() ⑤ df.index() ⑦ df.dtypes
 ② df.tail() ④ df.ndim() ⑥ df.values() ⑧ df.dtypes.value_counts()

Returns one row & column wise
no data type.

Attributes & Methods applicable only for DataFrame but not for Series:-

① info() method :- It generates summary information of DataFrame.

Syntax: DataFrame.info(verbose=None, buf=None, max_col=1000,
show_counts=None, null_counts=None)

② sample() Method :- We can use sample() method to get a random records from the DataFrame.

Syntax: DataFrame.sample(n=None, frac=None, replace=False, weights=None,
random_state=None, axis=None, ignore_index=False)

Returns a random sample of elements (or) items from a axis object.

eg: print(S.sample())

print(S.sample(n=2))

print(S.sample(frac=0.1))

DATE: .

(3) columns attribute: Returns column labels of the DataFrame.
eg: `print(df.columns)` → It is Index object contain column names.

(4) axes attribute: Returns a list representing the axes of the DataFrame.
eg: `print(df.axes)`

→ How to Select Only one required column of Data-Frame

1st way → `df.Column_name`

2nd way → `df['Column_name']` → 2nd way is recommended.

eg: `print(s['Names'])` `print(s[['Names', 'Marks']])` → Returns DataFrame
`print(s['Marks'])`

→ Arithmetic operations for DataFrame :-

We Cannot apply arithmetic operations directly on DataFrame Object because it contains multiple columns.

We can perform arithmetic operations by selecting a particular column of the DataFrame, because each column of a DataFrame is a Series object.

Eg: `print(s['Marks']+500)` → 500 will be broadcasted to every value of the Series. ↗ (or) `print(s['Marks'].add(500))`

DATE:

→ Adding New Columns to the DataFrame

1st way : df['Column'] = value → It will be added at last.

2nd way : df.insert(1, 'Country', 'India') → It will be added at specified row.

↳ Syntax → DataFrame.insert(loc, column, value, allow_duplicated=False)

eg: ① S.insert(1, 'ESAL', 999)

② S.insert(1, 'Country', ['A', 'B', 'C', 'D'])

→ How to drop DataFrame rows with missing values/null values?

We have to use dropna() method.

Syntax: DataFrame.dropna(axis=0, how='any'; thresh=None, subset=None, inplace=False)

Remove missing values. Removes rows from the DataFrame where at least one missing value.

eg: S = pd.read_csv('emp.csv')

S1 = S.dropna()

S = S.dropna(how='all')

default = any

Print(S1)

Any → If any NA values are present, drop that row (or) column

All → If ~~any~~ values are NA, drop that — / —

Subset → To delete rows if a particular column has missing values.

S1 = S.dropna(subset=['esal'])

axis → To delete column which have missing values.

DATE: [] [] [] [] []

→ How to fill missing/Null Values with our required values in DataFrame.
we have to use `fillna()` method for this requirement.

Syntax :- `DataFrame.fillna(value=None, method='None', axis=None, inplace=False, limit=None, downcast=None)`

eg: `df = s.fillna(0)`

`print(df)` → Replacing missing values of string with '0' is not meaningful.
So for different data types we have to fill the default value.

`df['Ename'] = s['Ename'].fillna('suhas')` → `df['Esal'] = df['Esal'].fillna(0)`

If column with different types :-

We can provide different replacement values for different columns based on column data type. We have to use `astype()` method.

For this we have to `dropna()` → to remove null values.

eg: `s = pd.read_csv('emp.csv')`

`print(s)`

`s.dropna(inplace=True)`

`s['Esal'] = s['Esal'].astype('int')`

`print(s)`

→ How to Sort values of the DataFrame :-

We can use `sort_values()` method.

Syntax :- `DataFrame.sort_values(by, axis=0, ascending=True, inplace=False, key=None, kind='quicksort', na_position='last', ignore_index=False)`

DATE:

We can use 'by' parameter to specify column based on which sorting has to be done.

eg:- `s.sort_values(by='Ename', inplace=True)`
`print(s)`

eg:- `s.sort_values(by='Ename', inplace=True, ascending=False, na_position='first')`

eg:- `s.sort_values(by=['Ename', 'Esal'], ascending=[False, True])`

→ How to sort DataFrame by index:-

Syntax: `DataFrame.sort_index(axis=0, level=None, ascending=True, inplace=False, key=None, kind='quicksort', na_position='last', ignore_index=False)`

eg:- `s.sort_index(inplace=True, ascending=False)`

→ How to rank Series of values by using rank() method:-

Syntax: `Series.rank(axis=0, method='average', ascending=True)`

`df = pd.read_csv('emp.csv').dropna(how='all')`

`df['Esal'] = df['Esal'].fillna(0).astype('int')`

`print(df)`

`df['Salary Rank'] = df['Esal'].rank(ascending=False).astype('int')`

`print(df)`

`df.sort_values(by='Salary Rank', inplace=True)`

`print(df)`

→ Default way to
remove na & fill
result at na to 0.

DATE: . .

→ How to filter data from DataFrame :-

- We can filter data based on some condition, ie by using subset of rows based on some selection criteria. Internally we can use Boolean masking for this filtering purpose.

Syntax: `[df [boolean masked array]]`

Case-1: Select employees belong to Bangalore

`print(df['Eaddrs'])` → `(df['Eaddrs'] == 'Bangalore')` → `print(df[df['Eaddrs'] == 'Bangalore'])`

Case-2: Select salary > 5000

`print(df[df['Esal'] > 5000])` similarly we can use `<, >, =, <=, >=`

→ How to filter data from DataFrame based on Multiple Conditions :-

We can use `&` and `|` operation to combine multiple conditions together.

Syntax: `df [C1 & C2]`

`df [C1 | C2]`

Case-1: To Select employee belongs to Bangalore & Salary > 5000

`C1 = df['Eaddrs'] == 'Bangalore'`

`C2 = df['Esal'] > 5000`

`print(df[C1 & C2])`

Case-2: Select employee who should be either Bangalore base who salary > 3000

(or) Mumbai based employee.

`C1 = df['Eaddrs'] == 'Bangalore' C2 = df['Esal'] > 3000`

`C3 = df['Eaddrs'] == 'Mumbai'`

`print(df[(C1 & C2) | C3])`

DATE: . .

→ How to filter data from DataFrame :

- We can filter data based on some condition, i.e. by using subset of rows based on some selection criteria. Internally we can use Boolean masking for this filtering purpose.

Syntax : `df[boolean masked array]`

Case-1 : Select employee belong to bangalore

`print(df['Eaddr'])` → `(df['Eaddr'] == 'Bangalore')` → `print(df[df['Eaddr'] == 'Bangalore'])`

Case-2 : Select salary > 5000

`print(df[df['Esal'] > 5000])` similarly we can use `<`, `=`, `>`, `<=`,

→ How to filter data from DataFrame based on Multiple Conditions :-

We can use `&` and `|` operation to combine multiple conditions together.

Syntax : `df [C1 & C2]`

`df [C1 | C2]`

Case-1 : To Select employee belongs to bangalore & salary > 5000.

`C1 = df['Eaddr'] == 'Bangalore'`

`C2 = df['Esal'] > 5000`

`print(df[C1 & C2])`

Case-2 : Select employee who should be either Bangalore base who salary > 5000 or Mumbai based employee.

`C1 = df['Eaddr'] == 'Bangalore' C2 = df['Esal'] > 5000`

`C3 = df['Eaddr'] == 'Mumbai'`

`print(df[(C1 & C2) | C3])`

DATE:

→ How to check inclusion by using isin() method :-

Eg:- To select employees who are belongs to Hyderabad, Mumbai & Delhi?
Without using isin method

$C_1 = df['Eaddr'] == 'Hyderabad'$

$C_2 = \dots == 'Mumbai'$

$C_3 = \dots / == 'Bangalore'$

`print(df[C1 | C2 | C3])`

Syntax `df.isin(values)`

Note: To ignore Case :- $C_1 = df['Eaddr'].str.lower().isin(['Hyd', 'Bang'])$

→ How to check for inclusion by using between() method :-

Without using between() method

$C_1 = df['Esal'] \geq 2500$

$C_2 = df['Esal'] \leq 10000$

`print(df[C1 & C2])`

With using between() method

$C = df['Esal'].between(2000, 5000)$

case-1: inclusive : { 'both', 'neither', 'left', 'right' }

$C = df['Esal'].between(2000, 5000, inclusive='both')$

→ How to check missing & non-missing Values by using isnull() & notnull() methods :-

`Series.isnull()` → Detect missing values.

`Series.notnull()` → Detect existing (non-missing) values.

DATE:

--	--	--	--	--	--

eg: C = df['Ename'].isnull()
print(df[C])

C = df['Ename'].notnull()
print(df[C])

→ How to check duplicate dataframe rows by using duplicated() method

eg: C = df['Ename'].duplicated()
print(df[C])

eg.2: C = df['Ename'].duplicated(keep=False)

eg.3: C = ~df['Ename'].duplicated(keep=False)

→ How to drop duplicate rows in DataFrame by using drop_duplicates()

Syntax: df.drop_duplicates(subset=None, keep='first', inplace=False, ignore_index=False)

eg: df.drop_duplicates(subset=['Ename'], keep=True, inplace=True)

→ How to identify & count unique values in DataFrame :-

S.unique() → returns unique values.

S.nunique() → number of values unique, by ignoring NaN.

S.nunique(dropna=False) → number of unique values without ignoring NaN.

eg: df['Ename'].nunique()

→ Extracting Values from the DataFrame :-

① How to customize Column as index

1st way:- df = pd.read_csv('emp.csv', index_col = 'Ename')

2nd way:- df.set_index(by='Ename', inplace=True)

DATE:

→ How to reset index from custom index to default numerical index:-
`df.reset_index(inplace=True)`

→ How to select rows by index label with `loc[]` indexer:-

Syntax- `df.loc[index_label]` ∵ `df = pd.read_csv('emp.csv', index_col='Ename')`

① `df.loc['Suhar']` → if only one row matched & returns in the form of Series.

If multiple rows matched returns Datframe

If specified index not matched we get KeyError.

② Using slice operator :- `df.loc[index_label1:index_label2:step]`

Here both Begin & end index labels are inclusive.

Eg: `print(df.loc['Suhar':'Ramesh':2])`

③ With multiple index labels of `df.loc[[index1, index2, index3]]`

→ How to Access Rows of DataFrame by using Index position:-

We have use by `iLoc[]` indexer.

① `df = pd.read_csv('emp.csv', index_col='Ename')` ② by using slice operator:

`print(df.iLoc[0:7])`

`df.iLoc[2:5:2]`

③ with multiple indexpositions ∵ `df.iLoc([index_1, index_2, index_3])`

→ How to get values of the required cells:-

Syntax:- `df.loc[index_label, Column_name]` → Eg: `print(df.loc['Suhar', 'ESal'])`

`df.iLoc(index_position, column_name)`

DATE:

Case study: ① `df.loc['sunny', 'Esal'].iloc[0]`

→ How to set ~~value~~ new value of the specified cells:

Syntax: `df.loc[index_label, col-name] = new value`

`df.iloc[index-position, col-name] = new value`

Eg: ① `df.loc['Suhar', 'Esal'] = 777` ② `df.loc[['A', 'B', 'C'], 'Eaddr'] = 'pu'`
`print(df) or print(df.loc['Suhar'])`

Q: Replace every occurrence of Hyderabad to Cyberabad

`C = df['Eaddr'] == 'Hyderabad'`

`df.loc[C, 'Eaddr'] = 'cyberabad'`

→ How to rename index labels & column names in a DataFrame:

By using `rename()` method for index labels (rows) → axis 0

For column names → axis 1

Renaming index labels in two ways: ① By using `mapper` & `axis` parameters

② By using `index parameter`

Syntax: `DataFrame.rename(mapper=None, index=None, columns=None, axis=None, copy=True, inplace=False)`

Eg: ① `df.rename(mapper={'sunny': 'leone', 'Bunny': 'Arjun'}, axis=0, inplace=True)`

② `df.rename(index={'sunny': 'leone', 'Bunny': 'Arjun'}, inplace=True)`

DATE:

→ Renaming Column Names :- ① By using mapper & axis parameters

② By using columns parameter

- ① eg: `df.rename(mapper={'Esal': 'salary', 'Eaddr': 'city'}, axis=1, inplace=True)`
- ② `df.rename(columns={' ': ' '}, inplace=True)`

→ How to delete rows & columns from the DataFrame :-

By using drop() method.

Syntax: `df.drop(labels=None, axis=0, index=None, column=None, inplace=False)`

eg: `df.drop(labels='Sunny', inplace=True)` → For rows default axis=0

`df.drop(labels='Esal', axis=1, inplace=True)` → To delete columns

→ How to filter DataFrame with query() method :-

We can use query() method to filter rows from the DataFrame. It is the most commonly used method in Pandas.

eg: `df.query(expr, inplace=False, **kwargs)`

While using query() method the following 2 points are important:-

- ① expr parameter value should be a valid string.
- ② The column names should not contain spaces.

eg: ① `df.query("Eaddr == 'Hyderabad'", inplace=True)`
`print(df)`

② `df.query("Esal > 1500", inplace=True)`

DATE: []

③ To handle multiple conditions :- we can use & or |

df. query ("Eaddr == 'Bangalore' and Esal > 1500", inplace=True)

④ in, not in operators in query() method :-

df. query ("Eaddr in ['Hyderabad', 'Bangalore']", inplace=True)
not in

Note: If column names contain space then query() method won't work.

How to resolve this? → Replace space with underscore (-) symbol

& then apply query() method.

df.columns = [col_name.replace(' ', '_') for col_name in df.columns]

df.query ("Total_Salary > 1500", inplace=True)

print(df)

→ How to use apply() method for DataFrame rows :-

It is exactly same as Series class apply() method except that it is applicable row wise.

df.apply(func, axis=0, raw=False, result_type=None, **kwargs)

Q: Double the Salary of each employee:

df[['Esal']].apply(lambda x: 2*x, inplace=True)

DATE:

→ How to get random sample of rows from DataFrame() :-

By using sample() method.

df.sample(n=None, frac=None, replace=False)

eg: ① print(df.sample(2)) ② print(df.sample(n=2, axis=1))

→ How to use nlargest() method of the DataFrame

eg: df.nlargest(n=3, columns=['Esal'])

→ How to use nlargest() method:-

eg: df.nlargest(n=3, columns=['Esal'])

→ How to filter DataFrame with where() method:-

Syntax: df.where(cond., other=None, inplace=False, axis=None)

Replaces values where Condition is False, i.e. it returns Full DataFrame where original data will be Considered if Condition is True & Nan will be Considered if Condition is False.

eg: To select rows where Eaddr is Hyderabad ?

df = pd.read_csv('emp.csv', index_col='Ename')

C = df['Eaddr'] == 'Banglore'

print(df.where(C)) (or) print(df.where(C, other='Not usfull'))

Note: Usage of where() method on the Series object is more meaningful than DataFrame.

DATE:

→ How to get Separate copy of the DataFrame :-

To create separate independent isolated copy of DataFrame (or) Series, we have to use `copy()` method.

`df1 = df.copy()` → If we perform any changes to the original object, then those changes won't be reflected to the copy.

→ How to handle text data with String methods :-

The most commonly used data type in any programming language is text data (Str type). To perform common requirement to handle text data like removing spaces, Converting into required case etc.

Pandas library provides the following methods to handle text data:-

- ① `lower()`
- ② `title()`
- ③ `replace((old-chr, new-chr))`
- ④ `strip()`
- ⑤ `lstrip()`
- ⑥ `rstrip()`
- ⑦ `upper()`
- ⑧ `len()`
- etc

While using these methods we should use 'str' prefix.

① ex : ① Convert all employees into uppercase :-

`df['Ename'] = df['Ename'].str.upper()`

② To remove \$ symbol in Esal :-

`df['Esal'] = df['Esal'].str.replace('$', '')`

③ To remove \$, ., , symbols in Esal & then convert values into float type :-

`df['Esal'] = df['Esal'].str.replace('$', '').str.replace(',', '').str.replace('.', '').astype(float)`

④ To add new column 'Number of Characters' which is evaluated from Ename value :-

`df['Number of Char'] = df['Ename'].str.len()`

DATE: []

Case Study: Convert Esal into proper numeric values:-

Esal

e.g.: df = pd.read_csv('emp2.csv')

15K

Conversion-rate = {'k': 1000, 'L': 100000, 'C': 10000000}

20K

df Convert(value):

30K

unit = value[-1] → ^{1 2}₁

40L

numeric = value[1:-1] ^{0 1 2 3 4}
\$200K

3C

updated-value = float(numeric) * Conversion-rate[unit]

return '\$' + str(updated-value)

df['Esal'] = [Convert(sal) for sal in df['Esal']]

Print(df)

→ How to filter DataFrame rows with String Methods :-

Contains(), startswith(), endswith().

① Contains(Sub-string) → Returns True if the string contains provided Sub-string irrespective of position.

② startswith(Sub-string) :- Returns True if the string starts with provided Sub-string

③ endswith(Sub-string) :- → → String ends with → →

e.g.: ① To filter all rows where ename starts with 'A'?

C = df['Ename'].str.startswith('A') → Case Sensitive.

C = df['Ename'].str.lower().str.startswith('a')

DATE:

② To filter all rows where ename contain 's'?

$C = df['Ename'].str.lower().str.contains('s')$

③ Ends with y :- $C = df['Ename'].str.lower().str.endswith('y')$

→ How to use String split() method :-

We can use split() method to split the given string based on some condition.

Eg: Consider the Ename column of DataFrame & split every name to two names - First Name & last name.

$df = pd.read_csv('emp2.csv')$

$df.dropna(inplace=True)$

$df.insert(loc=1, column='First Name', value=df['Ename'].str.strip().str.split(expand=True).get(0).str.title())$

$df.insert(loc=2, column='Last Name', value=df['Ename'].str.strip().str.split(expand=True).get(1).str.title())$

$def df['Ename']$

$print(df)$

Eg: Ename = 'Sunny Leone Pande' (n parameter)

$df['Ename'].str.split() \rightarrow ['Sunny', 'Leone', 'Pande']$

$df['Ename'].str.split(n=1) \rightarrow ['Sunny', 'Leone Leone Pande']$

DATE:

→ How to Concatenate, merge & join multiple DataFrames into a single DataFrame.

Syntax: `Pd.concat(objs, axis=0, join='Outer', sort=False)`

eg: `df1 = pd.read_csv('emp1.csv')`

`df2 = pd.read_csv('emp2.csv')`

`combined = pd.concat(objs=[df1, df2])`

`print(combined)`

→ How to remove duplicate index labels :- Two ways to remove duplicate index

1st way: by using `reset_index()` method

`df.reset_index(level=None, drop=False, inplace=False, col_fill=' ')`

eg: `combined.reset_index(inplace=True)`

`print(combined)`

2nd way: Ignore_index parameter of `concat()` method

eg: `combined = pd.concat(objs=[df1, df2], ignore_index=True)`

`print(combined)`

→ Defining MultiIndexing By using keys parameter.

eg: `combined = pd.concat(objs=[df1, df2], keys=['month1', 'month2'])`

→ How to access data if we are using MultiIndex :-

eg: `print(combined.loc['month1', 1])`

→ To enforce unique index labels →

By using `verify_integrity=True`

DATE:

→ How to concatenate Columns

Combined = pd.concat(objs=[df1, df2], axis=1)
print(Combined)

→ Merge Method: It is most powerful than concat() method. Its functionality is exactly same as SQL Join queries. Grouping can be done based on matched content.

Syntax :- pd.merge(left, right, how='inner', On=None, right_on=None, left_on=None)

eg:- To select all students who registered for the course in 1 month:-

df1 = pd.read_csv('month_1 reg. Csv')

df2 = pd.read_csv('students.Csv')

s = pd.merge(df1, df2)

print(s)

eg:- s = pd.merge(df1, df2, On='student_id', suffixes=('_month_1', '_month_2'))

eg:- If there is no Common Column name :-

If there is no common name, merge operation fails by default.

To Overcome This problem, we have to use left_on & right_on parameters

eg:- df1 = pd.read_csv('month_1 reg. Csv') → student id Course id

df2 = pd.read_csv('Courses.Csv') → CID, Course Name, Fee

s = pd.merge(df1, df2, left_on='Course id', right_on='CID')

s.drop(['CID'], axis=1, inplace=True)

In df2

DATE: [] [] [] [] []

→ Merge Outer Join :- how = 'outer' → All keys are selected. similar to set union operation. If data present in both DataFrames is important & if we don't want to miss any data then we should go for outer join.

eg:- df₁ = pd.read_csv('month_1-reg.csv')
df₂ = pd.read_csv('courses.csv')
S = pd.merge(df₁, df₂, how = 'outer')
print(S)

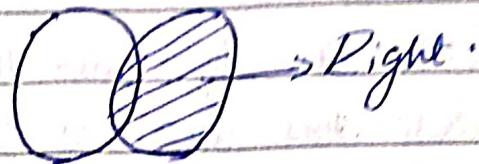
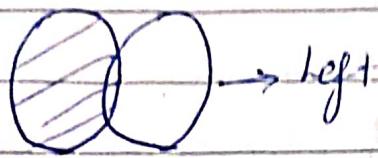
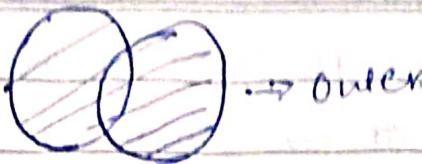
→ default is 'inner' join

Indicator Parameter :- indicator :- bool (or) str, default False
If True, adds a column to the output DataFrame called 'merge'
With information on the source of each row.
S = pd.merge(df₁, df₂, how = 'outer', indicator=True)
How to find source count of rows:
print(S['merge'].value_counts())

eg:- Select only rows present in either left (or) Right DataFrame but not both
W = S['merge'].isin(['right-only', 'left-only'])
print(S[W])

→ Left Join :- Data present in Left DataFrame is very important & we have to consider complete information, but right DataFrame is not important we have to consider only matched content, then we should go for left join.
eg:- df₁ = pd.read_csv('month_1-reg.csv')
df₂ = pd.read_csv('students.csv')
S = pd.merge(df₁, df₂, how = 'left')

DATE:



Right Join:

```
df1 = pd.read_csv('month-1 reg.csv')
```

```
df2 = pd.read_csv('students.csv')
```

```
S = pd.merge(df1, df2, how='right')
```

```
S.to_csv('complete.csv') → To write new csv file.
```

Syntax: df.merge(right, how='inner', on=None, right_on=None, left_on=None, left_index=False, right_index=False, suffixes=('_x', '_y'), copy=True, indicator=False, validate=None)

On → Based on Which Column , merge has to perform.

left-on → Based on Which Column in left DataFrame, merge has to perform

right-on → Based on Which Column in right DataFrame, merge has to perform

DATE:

→ join() method: If we want to perform merge operation based on indexes of both Dataframe then we should use join() method or Shorthand technique.

Pd.merge(df₁, df₂, left_index=True, right_index=True)

We can rewrite this line as :- df₁.join(df₂).

→ How to use Regular Expressions in Pandas:-

If we want to represent a group of characters according to a particular pattern then we should go for regular expressions.

prefixed character classes:-

I S → Space Character

I d → Any digit

I S → Any character except space character.

I D → Any character except digit

I W → [a-zA-Z0-9]

• → Any character.

I W → Any character except word character

Quantifiers:-

a → Exactly one 'a'

a{ⁿ} → Exactly n number of a's

a+ → Atleast one 'a'

a{^{m,n}} → minimum m & a's & max n

a* → Any number of a's including zero a{^{m,n}} → minimum m number & a's & max any of a's

a? → Atmost One 'a'

a{^{l,n}} → minimum any number of a's & max n number of a's

DATE:

replace(), contains() or split() methods with regular expressions:-

Syntax: Series.str.replace(pat, repl, n=-1, case=None, flags=0, regex=None)

eg: To remove \$, %, Symbols in Esal column values. Then convert values into float type.

eg: df['Esal'] = df['Esal'].str.replace(r'[%,\\$]', '', regex=True).astype(float)

Contains():

Syntax: Series.str.contains(pat, case=True, flags=0, na=None, regex=True)

eg: C = df['Ename'].str.contains(r'kmj', case=False, regex=True)
Print(df[C])

Split():

Syntax: Series.str.split(pat=None, n=-1, expand=False)

Split key column into 3 parts: k1, k2, k3

df[['k1', 'k2', 'k3']] = df['key'].str.split(r'_', expand=True)

del df['key']

Print(df)

DATE:

→ Working with MultiIndex :- Multi level Index (or) Hierarchical Index
Sometimes One level index is not enough to access data from the DataFrame. We should go for multiple levels of indexing which is nothing but MultiIndex.

Eg : df = pd.read_csv('multi_indexing.csv')

print(df)

Print(df.index) # RangeIndex (start=0, stop=8, step=1)

print(df.columns)

RangeIndex is the child class of Index. → If we are not providing explicit index Pandas will generate default numeric index which is of type RangeIndex. Eg it is immutable

→ How to set our own column as Index Column

① By using set_index() method → df.set_index(keys='Date', inplace=True)

② By using index_col parameter. → df = pd.read_csv('emp.csv', index_col='Date')

How to define Multi Indexing :- We have to pass list of Column names.

① df.set_index(keys=['Date', 'Domain'], inplace=True)

② df.read_csv('emp.csv', index_col=['Date', 'Domain'])

Eg : df = pd.read_csv('emp.csv', index_col=['Date', 'Domain'])

Note :- RangeIndex → Default Numeric Index
Index → Column based Index.

DATE: . .

→ How to access data from Multi-Indexed DataFrame

There are multiple ways available to extract data :-

① By using loc indexer with index labels :-

loc → Label based selection.

iloc → Position based selection.

ex:1: To get all batches info on '25-01-2023' ? → ie, access data by using first level index :-

```
df = pd.read_csv('multi_courses.csv', index_col=['Date', 'Subject'])  
print(df.loc['25-01-2023'])
```

ex-2: To get Java batch information on '25-01-23' ?

```
print(df.loc['25-01-2023', 'python'])
```

ex:3: To get Fee info of 'java' batch on 25-01-23 & Timing.

```
print(df.loc['25-01-2023', 'python'], 'Fee')
```

```
print(df.loc['25-01-2023', 'python'], ['Fee', 'Timing'])
```

→ How to access with Range of indexes :-

eg.: To get python batch information from '25-01-23' to '28-01-23' ?

```
print(df.loc[['25-01-2023', '26-01-2023']])
```

ex:2:

```
print(df.loc[['25-01-23', '26-01-2023'], ['SQL']])
```

The total tuple is considered as index labels only.
It represents all column names.

DATE:

Q3: print(df.loc[['2022-01-22', '2022-01-23'], ['Java', 'Python'], ['Fee', 'Timing']])

Syntax: df.loc[(), ()]

() → index labels

[] → column names

Note: It is very common & very easy to use for indexer when compared with iloc indexer.

Q4: To get Python & Java Batches Fee & Faculty information on 24-01-22 and 25-01-22?

print(df.loc[['24-01-22', '25-01-22'], ['Java', 'Python'], ['Fee', 'Faculty']])

→ How to perform slice operation?

→ How to get cross section from the DataFrame by using xs() method
xs → Cross Section.

Syntax: df.xs(key, axis=0, level=None, drop_level=True)

Returns Cross-Section from one Series / DataFrame.

Q5: df = pd.read_csv('multi_courses.csv', index_col=['Date', 'Subject'])
print(df.xs('25-01-2023'))

Q6 To get Java Batches info for all dates:

print(df.xs('Java', level=1, drop_level=False))

DATE: [] [] [] [] []

Q. To get Java info on '28-01-2023'?

```
print(df.xs('28-01-2023', 'Java'), level=(0, 1), drop_level=False)
```

→ The attributes of MultiIndex object:-

We can use multiple attributes on MultiIndex object. These are very helpful if we are working on very complex hierarchical DataFrames.

Q. df = pd.read_csv('new-batches.csv', index_col=['Date', 'Subject'])

```
print(df) print(type(df))
```

```
print(df.index) print(type(df.index))
```

① names: print(df.index.names) List of names for each of the index levels.

② nlevels: Integer number of levels in the multiindex

```
print(df.index.nlevels)
```

③ levels: The unique labels for each level:-

```
print(df.index.levels)
```

④ lenshape: A tuple with the length of each level.

```
print(df.index.lenshape) → (5, 4)
```

i. First level index contains 5 values, 2nd level index contains 4 levels.

⑤ size: (df.index.size)

⑥ values: (df.index.values)

DATE:

→ How to add another level of the Multi-Index :-

1st way:- Read again data & Create DataFrame.

```
df = pd.read_csv('emp.csv', index_col=['Date', 'Subject', 'Faculty'])
```

2nd way:- If DataFrame already available.

If DataFrame already available, then we should add new level of index by using set_index() method.

Syntax:- df.set_index(keys, drop=True, append=False, inplace=False)

e.g:- df = pd.read_csv('emp.csv', index_col=['Date', 'Subject'])

```
df.set_index(keys='Faculty', append=True, inplace=True)
```

```
print(df) print(df.index.names)
```

→ How to swap index levels :- In Two ways.

① By using index position:- To swap any two levels.

Syntax:- df.swaplevel(i=-2, j=-1, axis=0)

e.g:- df₁ = df.swaplevel() → innermost 2 level i.e last 2 levels will be swapped

```
df2 = df.swaplevel(0, 2) → we can use either
```

```
df3 = df.swaplevel('Date', 'Subject') index position (or) Index, axis
```

② By using reorder_levels():- To swap more than 2 levels.

Syntax:- df.reorder_levels(order, axis=0)

e.g:- df₁ = df.reorder_levels(order=[2, 0, 1])

```
df2 = df.reorder_levels(order=['Subject', 'Faculty', 'Date'])
```

DATE:

→ How to remove existing index level from MultiIndexed DataFrame

By two ways :- ① By using `droplevel()` method. (not recommended)
② By using `reset_index()` method.

① By using `droplevel()` :-

```
df = pd.read_csv('multi-course.csv', index_col=['Date', 'Subject'])
```

```
df1 = df.droplevel(level='Subject')
```

```
print(df1)
```

Not recommended to use `droplevel()` because there may be a chance of loss of information.

② By using `reset_index()` method :-

Syntax :- `df.reset_index(level=None, drop=False, inplace=False, col_level=0)`

e.g. `df = pd.read_csv('emp.csv', index_col=['Date', 'Subject', 'Faculty'])`

```
df.reset_index(level=[0, 2], inplace=True)
```

```
print(df) print(df.index.names)
```

To remove all indexes

```
df.reset_index(inplace=True)
```

DATE:

--	--	--	--	--	--

→ How to Sort Index by using sort_index() method:-

It is highly recommended to perform sorting of indices.

Syntax: df.sort_index(axis=0, level=None, key=None, ascending=True, inplace=False)

eg:- df=pd.read_csv('emp.csv', index_col=['Date', 'Subject', 'Faculty'])
df.sort_index(inplace=True)
print(df)

→ How to Sort a Particular levels :-

df.sort_index(level=2, inplace=True, ascending=False)

df.sort_index(level=(0,1), ascending=[False, True], inplace=True)

→ Important Methods of Multi-Index :-

Multi-Index is a class in pandas library.

- ① is_lexsorted() :- Returns True if index is lexicographically sorted.
- ② sortlevel() :- df.index.sortlevel(level=0, ascending=False)