

PANDAS

classmate

Date _____
Page _____

↳ library used for data manipulation and analysis

+
cleaning
Handling missing data
Merging or joining data

- Data Manipulation

- Cleaning
- Handling missing data
- Merging, joining data

- Data Analysis

- Read Data
- Mean, Median, variance

- Handling different Data Format:

- CSV → JSON
- Excel → SQL databases

PANDAS DATASTRUCTURE

classmate

Date _____
Page _____

① Series (1 D labelled Array)

- Can hold any datatype → integers, float, strings, Python objects

② Data Frame (Table)

- Can store different type data
e.g. → integer, float, strings.
- Tabular Spread Sheet like str. containing rows, each of which have one or more columns

Series e.g.

↳ `data = [10, 20, 30, 40]`

`series = pd.Series(data,
index=['a', 'b', 'c',
'd'])`

→ a 10
b 20
c 30
d 40
} stored vertically

X ~~`pd.Series(data, index=False)`~~ X

→ By using
⇒ 1 D Data

Dataframe E.g

data = {

'Name': ['Alice', 'Bob', 'Charlie']

'Age': [25, 30, 35]

'Salary': [50000, 60000, 70000]

}

pd.DataFrame(data)

	Name	Age	Salary
0	Alice	25	50,000
1	Bob	30	60,000
2	Charlie	35	70,000

Q How you create 2D data for Dataframe?

- By using list in Dictionary.
- ⇒ 1D Data => list

Data

Series

- | | | | |
|---|--|---|--|
| ① | 1 D | ① | 2 D |
| ② | Single column or row of data | ② | Multiple rows or columns |
| ③ | Access element directly by their index | ③ | Access element using <code>.iloc[]</code> or <code>.loc[]</code> |

Q Difference in `.loc` & `.iloc[]` used in Dataframe of Pandas.

↳ can also be used in `.loc[]` in Series

Index	Name	Age
100	Alice	25
101	Bob	30
102	charlie	35

→ `df.iloc[1, 0]` → Bob.
 row ↑ column ↑

PANDAS

① → -

⑪

N

d.

→ df.loc[0, 'Name'] → Bob

Q Okay that's all right Vinayak
what's the actual difference?

→ is how they select data (access data)

iloc. → select data based on native integer location

loc → select data according to index of rows & labels of column

PANDAS (Code with Harry)

PANDAS FXN

① → To Excel (Save)

pd df.to_csv('Abc.csv')

② No index (Save to excel)

df.in
df.to_csv('Abc.csv', index=False)

(III) `df.head(2)` \Rightarrow Two values starting

(IV) `df.tail(5)` \Rightarrow Last 5 rows.

(V) `df.describe()`,

i) Work on Numerical column
ii) calculate

- \hookrightarrow I Count
- \hookrightarrow II Mean
- \hookrightarrow III Std
- \hookrightarrow IV Min
- \hookrightarrow V Max
- \hookrightarrow VI 25%
- \hookrightarrow VII 50%
- \hookrightarrow VIII Max
- \hookrightarrow IX 75%

E.g \Rightarrow

Name

Hello

Yellow

Mellow

Semi-Numerical

\hookrightarrow count

\hookrightarrow freq

\hookrightarrow top

\hookrightarrow unique

Q What happens if we use `describe()` fxn on not number data or semi numerical data?

$\star\star$ when applied on Table which have numeric column & Non-Numeric + Semi-Numeric \rightarrow only do

\rightarrow Can be applied on Non-Numerical data

\hookrightarrow But we usually don't do
 \hookrightarrow So, say it is only used on Numeric data

for
Numerical

Mixed

100

200

abc

200

- Non-Numeric Data →
 - i) count
 - ii) unique
 - iii) top (most occurring)
 - iv) freq ⇒ freq. of most occurring

E.g ⇒

Name	City
Hello	Path
Yellow	Amsr
Mellow	Cith

} object
↓
datatype

- Semi-Numeric

- ↳ count
- ↳ freq
- ↳ top
- ↳ unique

Mixed	count
100	unique
200	top
abc	freq
200	

- Mixed DataFrame

Name	Age	City
Alice	10	Pkt
Bob	20	Chd
Cily	40	AmSr

\Rightarrow op \Rightarrow

	Age
Count	3
mean	30
std	
min	

(VI) pd.read_csv('ABC.csv') \Rightarrow To Read CSV.

(VII) Harry

(VIII) pd.Series(np.random.rand(34))

34 ele in series

(IX) pd.DataFrame(np.random.rand(334, 5))

→ Create 334 rows of with 5 columns

`df.rows` \Rightarrow X

CLASSMATE

Date _____
Page _____

- (IX) `df.index` \Rightarrow rows name index.
 $\hookrightarrow 0, 1, 2, 3$
- (X) `df.columns` \Rightarrow Tells column name
 $\hookrightarrow \text{col1}, \text{col2}, \dots$
- (XI) `df.to_numpy()` \Rightarrow DataFrame converted to Numpy.
- (XII) $\text{Axis} = 0 \Rightarrow \text{Rows}$
 $= 1 \Rightarrow \text{Columns}$
- (XIII) `df.drop(0, axis=1)`
 \hookrightarrow Remove column
 \hookrightarrow will point by dropping \hookrightarrow But don't change that in original
- (XIV) `df.loc[[1, 2], ['C', 'D']]`
 $\underbrace{\text{rows}}$ $\underbrace{\text{columns}}$
- (XV) `df.loc[(condition)]`

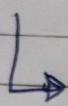
E.g. \Rightarrow `df.loc[newdf]`

\hookrightarrow `df.loc[(df['A'] < 0.3)]`
 \hookrightarrow Return all rows with column A satisfying this

(XVI)

`df.loc[(cond1) & (cond2)]`

E.g. `df.loc[(df['A'] < 0.3) & (df['C'] > 0.1)]`



Return all rows with col. satisfying that condit.

★★★
(XVII)

`df = df.drop(['A', 'D'], axis=1)`

OR

`df.drop(['A', 'D'], axis=1, inplace=True)`

will change in original Table

★★★
(XVIII)

`pd.read_json('ABC.json')`

★★★
(XIX)

`df["salary"].tolist()`

↳ convert to list

Q

create
Array

data

series

----> NumPy

★★★
①

0
1
2
3
4
5
6
7

→ df.

→ df.ilo

PANDAS (GFG)

classmate

Date _____

Page _____

= Q

Create Pandas Series from Numpy Array Data?

```
data = np.array(['g', 'e', 'k',  
                 's'])
```

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

= 1)

---> Numpy \Rightarrow converts list to array.

①	COL 1	COL 2
0	52	53
1	21	54
2	99	97
3	88	99
4	79	100
5	62	101
6	44	107
7	49	108

$\rightarrow df.loc[3:6] \Rightarrow$ will print row 3, 4, 5, 6.

$\rightarrow df.iloc[3:6] \Rightarrow$ will print 3, 4, 5.

Q You have two Series d1 & d2 . How to add subtract.

→ d1.add(d2)

→ d2.sub(d1)

d1

a	5
b	2
c	3
d	7

d2

a	1
b	6
c	4
d	9

→

a	6
b	8
c	3
d	11
e	9

• unique()

• value_count

the r

★ Series
array

→ dict

pd.Series

index

Gee

Boo

Cooo

Q Suppose you want to change the datatype of one column . How to do that?

→ astype()

e.g. → data["Salary"] = data["Salary"]
· astype(int)

★ DataFr
from

→ But
S

- `unique()` → See unique value in a particular column
- `value_counts()` → Method to count the number of times each

★ series can be created from Numpy array, list or dictionary

↳ dict = { 'Geeks': 10, 'Books': 20, 'Crooks': 30 }

→ pd. Series(dict)

Index.

O/p ⇒ Geeks 10
Books 20
Crooks 30.

★ DataFrame can also be created from list, dictionary, list of list or dictionary

→ But prefer list or Dict:
{ 'Name': [..] },

NIGHT
thon.
Re

① List of List : DataFrame

lst = [[1, 'Geeks'], [2, 'for'],
[3, 'Geeks']]

pcl.DataFrame (lst, col=['Id', 'Data'])

Id	Data
1	Geeks
2	for
3	Geeks

② List of Dictionary : DataFrame

lst = [{1: 'Geeks', 2: 'for',
3: 'Geeks'},
, {1: 20, 2: 30, 3: 40}]

1	2	3
Geeks	for	Geeks
20	30	40

Q How to ?

df = df.rename

Q Assign

df = df.

Q Descr

Q set

dates

df = pd

Q How to rename column in a DataFrame
?

df = df.rename(columns={ 'col1' : 'Sach' ,
 'col2' : 'Piggy' })

Q Assign custom index to DataFrame

df = pd.DataFrame(d1, index=[0, 1, 2])

Q describe() → only numerical
column are selected.

Q set date as index

dates = pd.date_range('20250118',
 period=6)

df = pd.DataFrame(chita, index=dates,
 columns=[0, 1, 2])

MY
NIGHT
Python.
P.

CLASSMATE

Date _____
Page _____



Q = Access Method using Different Strategies in DataFrame

→ what allowed / what not

df.

df [:] ['A'] ✓

df [:, 'A'] X

df [:] or df [1:4] ✓

df [:][['A', 'B']] ✓

df [:][['A', 'B']] X

Syntax for Directly accessing

df [] []
 ↓ +
 col rows

→ Using loc & iloc

→ df3

`df.loc[:, :, ['A', 'B']]`

`df.loc` \Rightarrow Print memory location

`df.loc[:, :] == df.loc[:, :, :] == df.loc[:, :, :, :]`

`df.loc[:, :, :, :]` \Rightarrow Show all index

`df.loc[:, :, :]` \Rightarrow X

Concatenate Two DataFrame

df2

0	1	0	1	2
x	a	a	b	c
y	b	d	e	f
z	c	g	h	i
		j	k	l

`df3 = pd.concat([df1, df2], axis=0)`
rowwise

Table

	0	1	2
0	x	a	NaN
1	y	b	NaN
2	z	c	NaN
0	a	b	c
1	d	e	f
2	g	h	i
3	j	k	l

→ df4 = pd.concat([df1, df2], axis=1)

column wise

i) pd.mer

→ To correct index use

→ df4.reset_index(drop=True, inplace=True)

ii) pd.

iii) pd.

JOIN IN PANDAS

classmate

Date _____
Page _____

Q Types of Joins?

- => Inner
- < Outer
- < Left
- < Right

Q Left Join, Outer Join, Inner two Dataframe

DF1

A	B
1	11
2	12
3	13
4	14

DF2

A	C
1	21
2	22
3	23
5	24

=1)

wise

i) pd.merge(var1, var2, on = 'A')
↳ inner

ii) pd.merge(var1, var2, how = "left")

iii) pd.merge(var1, var2, how = "outer")

★ ★

①

How to access first row?

→ df[column][row] X

↳ can't be used
in this

df.head(1)
 df.iloc[0]
 df[:, 0]

✓
✓

X

★ ★ Need to practice indexing ⇒ Accessing
element from Kaggle

MEAN, MEDIAN, UNIQUE,
VALUE_COUNTS

↳ Review ⇒ Table

① reviews.description.mean()

② reviews.points.median()

↳ middle value when data is
set in order

③

O/P ⇒

⑤

review

→ Litera
(deep)

CH

→ Un
ce

★

⑥

rev

reviews.points.mode()

↳ most frequently occurring.

⑤ reviews.points.unique()

↳ only shows the unique values every value once (col name)

→ Like Like if there are two name (duplicate) 'Roger', 'Roger', 'Hello'.

→ unique will only show it once a name

↳ 'Roger'
'Hello'

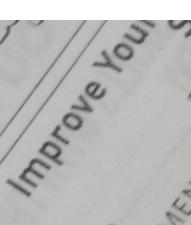
★

⑥ reviews.points.value_counts()

↳ How many times each value occurs.

↳ all unique() ones

Op =>	Roger	10
	Hello	5
	Kane	1
	:	:



Improve Your

ACID
ARGUMENT
BEAK

MAPS

- ↳ are the way to modify the data in Series or DataFrame.
- ↳ Perform operation element-wise.

MAP

Two Method to Apply

(I) map & lambda

(II) apply()

* → map - s a values

(II) apply()

Map

(I) reviews.points.map(lambda p: p

- reviews.points
- mean())

(I) works only

(II) Element

(III) No axis

(II) reviews.apply(

↳ I can't understand

O/p =	0	-1.447
	1	-1.499
	2	1.55
	3	
	.	
	.	

* → map & lambda combination return
- s a new series where all the
values are transformer

② Apply()

Map()

- ① works on series only
- ② Element wise
- ③ No axis

Apply()

- ① works on series & DataFrame
- ② Element wise
- ③ axis = 0 (columns)
axis = 1 (rows)

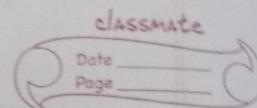
sum is a built-in function.

A B

Apply



1	4
2	5
3	6



= Q

Return a
many to
"fruity"
app de

① result = df.apply(sum, axis=0)



O/P

A	5
B	15

① →

n-trop =
(lambda
un

② result = df.apply(lambda
row: row * 2, axis=1)

A	B
2	8
4	10
6	12

---> True
Sum

---> Built-in Pandas

→ sum()

↳ a - b → a - b.mean()

↳ a + b → a + b.mean()



② n-fruity

(lambda

→ sum()

classmate

Date _____
Page _____

Q Return a series counting how many times word "Tropical" or "fruity" appears in column app description

① $n_trop = reviews.description.map(\lambda desc: "tropical" \text{ in } desc).sum()$

0	True
1	False
2	True
3	False

...> True is treated as 1 in sum() & False as 0.

→ sum() $\Rightarrow 2$.

② $n_fruity = reviews.description.map(\lambda desc: "fruity" \text{ in } desc).sum()$

0	True	2	False
1	True	3	True

→ sum() $\Rightarrow 3$.

des. count = pd.Series([n_trop, n_fruity], index=['tropical', 'fruity'])

Q Why not use here count()?

↳ Because statement in bracket returns Boolean value True or False, count will consider all value in count.

GROUPBY

↳ Similar to GROUP BY in DBMS (SQL)

A	B	C	D
10	60	70	1
10	70	90	2
15	90	80	3
20	100	77	5
30	190	67	4

→ df.groupby('A')

A	B
10	130
15	90
20	100
30	190

→ df.groupby('B')

A	B
10	15
15	20
30	30

→ df.groupby('C')

A	B
10	15
20	30

classmate
Date _____
Page _____

open
in fruit
[tropical]
]
L?
bracket
True or
either all

BY in

D

1
2
3
5
4

→ df.groupby('A').sum()

A	B	C	D
10	130	160	1
15	90	1	1
20	100	1	1
30	190	1	1

→ df.groupby('A').B

A	B
10	[60, 70]
15	[40]
20	[100]
30	[190]

→ df.groupby('A').B.min()

A	B
10	60
15	40
20	100
30	190

★ we can also use aggregate function with groupby

↳ `reviews.groupby('country').price.agg([len, min, max])`

Output → will calculate count of prices
min price
max price

SORT EXN()

↳ `Table-name.sort_values(by='col-name')`

↳ `Table-name.sort_values(by='col-name', ascending=False)`

↳ `Table-name.sort_values(by=['country', 'len'])`

first sorted by country after that in that sorted by len.

★ Column consists don't get the instead given

Q change the reviews point

HANDLING VALUES

Q check all a column

→ reviews

913 NAN

450 NAN

1011 NA

★ Column consisting entirely of strings
don't get their own type; they are
instead given the object type.

Q change the datatype of column?

reviews.points.astype('float64')

HANDLING MISSING

VALUES

Q check all the missing value in
a column?

→ reviews[pd.isnull(reviews.country)]

913 NAN

450 NAN

1011 NAN

Q Fill Non-Null (NaN) value with something else?

→ reviews.region.fillna("Unknown")

= Q Calculate total no. of entries is null in a column?

→ reviews.price.isnull().sum()

= Q Print correlation matrix

→ df.corr()

= Q Drop the rows with missing value in Table.

df = df.dropna();

= Q Drop the rows having more than two columns with missing value in Table.

df = df.dropna(thresh=2)

Q Drop the have missin

→ df.drop

= Q Removes the Table

→ df.drop

Q Removes value er

→ user[

Q How to DataFr

df['T

df

Q Drop the rows where age column have missing value

→ df.dropna(subset=['Age'])

Q Removes Duplicate rows from the Table?

→ df.drop_duplicates()

Q Removes Rows with Duplicate value in age column?

→ user['age'].drop_duplicates()

Q How to add new column in DataFrame?

df['Total_fare'] = 5 * df['price']

df

Total_fare

NUMPY

→ Numpy is a library which provides multidimensional array object & tools for working with these arrays.

→ LIST to Numpy Array

① $\text{lst} = [1, 2, 3, 4]$
 $\text{arr} = \text{np.array(lst)}$

O/P \Rightarrow $\text{array}(1, 2, 3, 4)$

$\text{type(arr)} \Rightarrow \text{numpy.ndarray}$.

$\text{arr.shape} \Rightarrow (4,)$

②

$\text{lst1} = [1, 2, 3, 4, 5]$
 $\text{lst2} = [2, 3, 4, 5, 6]$
 $\text{lst3} = [3, 4, 5, 6, 7]$

$\text{arr1} = \text{np.array}([\text{lst1}, \text{lst2}, \text{lst3}, \text{lst4}])$

O/P

$\text{arr1.shape} \Rightarrow (3, 5)$

↳ DSA
only

↳ IN NIGHT
↳ Python.
↳

0/p \Rightarrow array ([[1, 2, 3, 4, 5],
[2, 3, 4, 5, 6],
[3, 4, 5, 6, 7]])

\rightarrow [1, 2, 3, 4, 5]
 \rightarrow arr < 2
 \hookrightarrow True

★ If Numpy Array is 1D \Rightarrow
Retrieve it like LIST.

↳ In NDA
POSE T

1D Numpy \Rightarrow [1, 2, 3, 4, 5]
Array

\hookrightarrow like arr
arr[].

★ If Numpy Array is 2D or Multi
→ Retrieve it like iloc in Data
frame

2D \Rightarrow [1, 2, 3, 4, 5]
[2, 3, 4, 5, 6]
[3, 4, 5, 6, 7]

↳ (create)

arr[:, 1]

arr[2:4, :]

arr[2:5, :: -1]

np.array

\Rightarrow arr[]

\Rightarrow np.array

classmate
Date _____
Page _____

classmate
Date _____
Page _____

→ [1, 2, 3, 4, 5]

→ arr < 2
↳ True, False, False, False, False

Q In NDARRAY, WE CAN ALSO TRANSPOSE THAT ARRAY.

↳ If array is of shape => (3, 5)

arr1.reshape(5, 3)

↳

1	2	3
2	3	4
3	4	5
4	5	6
5	6	7

Q Create Numpy array using range.

np.arange(1, 10, 1)

=> array([1, 2, 3, 4, 5, 6, 7, 8, 9])

→ np.arange(1, 10, 1).reshape(2, 5)

① arr * arr \Rightarrow element wise multiplication
② arr * 2 \Rightarrow 2 multiplied element wise

③ np.random.randint(10, 50, 4)
 \hookrightarrow 15, 35, 40, 39
array([15, 35, 40, 39])

④ np.random.randint(10, 50, 4).reshape(2, 2)
 \hookrightarrow array([[20, 35],
[40, 49]])

⑤ np.random.random_sample((4, 7))
 \hookrightarrow 0.0 & 1
 \hookrightarrow array([[0.42, 0.12, 0.72, 0.32, 0.52, 0.22, 0.62],
[0.32, 0.22, 0.42, 0.52, 0.12, 0.62, 0.72],
[0.52, 0.42, 0.32, 0.22, 0.12, 0.62, 0.72],
[0.22, 0.12, 0.62, 0.72, 0.42, 0.32, 0.52]])