

Pandas

```
import pandas as pd
```

Creating Data

There are 2 core objects in pandas.

↳ DataFrame

↳ Series

Data Frame

```
=> pd.DataFrame({'Yes': [50, 21], 'No': [131, 2]})
```

	Yes	No
0	50	131
1	21	2

```
=> pd.DataFrame({'Bob': ['I liked it.', 'It was awful'],
                  'Sue': ['Pretty good!', 'Bland!']})
```

Bob	Sue
0 I liked it	Pretty good
1 It was awful	Bland

Up until this point index was by default starting with 0 and was increasing value.

Assigning customized index

`pd.DataFrame({ 'Bob': ['I liked it.', 'It was awful'],
 'Sue': ['Pretty good', 'Bland'],
 index=['Product A', 'Product B']})`

Product A	Bob I liked it Pretty good It was awful Pretty good	Sue Pretty good Bland
Product B		

Index

Series

=> pd.Series([1, 2, 3, 4, 5])

series consists of only

0 1
1 2
2 3
3 4
4 5

one column

dtype: int64

Also part of output

Assigning customized Index & an overall name

=> pd.Series([2100, 2200, 2300], index=

["2021 sales", "2015 sales",

"2027 sales"]
, name="Sales")

Reading data files

=> `wine_reviews = pd.read_csv("...")` address liut.csv'

only assigning .part is done

Shape

=> `wine_reviews.shape`

(129971, 14)

↓
row number
of dataset

↳ column number

Examining the DataFrame

=> `wine_reviews.head()`

(first 5 to row ashbe)

Index define w/ automatic 42m7

index (cont.)

Pointing specific column of the table
as Index.

wine-reviews = pd.read_csv("... address link...
, index_col=0)

Other column will be index.

Saving DataFrame as CSV file in hard disk

→ animals.to_csv("animals-in-Rangpur.csv")
DF

→ It'll save it.

Native Accessors

→ suppose reviews is a DataFrame, and country is one of its objects (columns).

⇒ reviews.
 |
 | country
 |
 | Native Accessor

O/P:

0	Italy
1	Portugal
..	China
129969	Nepal
129970	

turned a DF into
a series

→ columns can also be treated as keys
of dictionary as in Python.

⇒ reviews['country']

O/P:

same

$\Rightarrow \text{reviews}[\text{'country'}][0]$

O/P

'Italy'

This can't be done by reviews.country

Indexing in pandas

Pandas indexing works in one of two paradigms.

- Index-based selection
- Label-based "

Index-based Selection

→ Selecting data based on its numerical position, iloc operator is used.

$\Rightarrow \text{reviews}.iloc[0]$

$\downarrow \uparrow$
index location

0th row vs 0th data country

vow first column second

O/P:

country Italy

Description blank blank

Variety white

Flag WR

Name: 0, Length: number_of_vows, dtype: object
=> reviews.iloc[:, 0]

0th column 0th element COUNTRY

0 Italy

1 France

129969 Spain

129970 Portugal

Name: Length

dtype

=> reviews.iloc[:, 3, 0]

0 Italy

1 France

2 USA

~~2024-01-16 08:57 3rd element~~

0 element WR (3-1) element
= 2

\Rightarrow reviews.iloc [1:3,0]

OP!

1 France
2 USA

1st row $(3-1)$ th element
 $= 2^{\text{nd}}$ row

\Rightarrow reviews.iloc [[0,1,2],0]

at element show empty

\Rightarrow reviews.iloc [-5:]

129966

~~Q~~ v

awtcbgv

129967

pw

bldfnars

129968

~~Q~~

tulmactb

129969

~~Q~~

rmatcbjnv
nkwbwbf

129970

~~Q~~

cont cont at element empty

Label-based selection

loc operator is used

$\Rightarrow \text{reviews.loc[0, 'country']}$

O/P:

'Italy'

0 row country column can select
0 column ~ 50

ans

$\Rightarrow \text{reviews.loc[:, ['A', 'B', 'C']]}$

O/P:
(0-n) all row and specific 'A', 'B', 'C' column

data show ans

size = ~~n × 3~~ $(n+1) \times 3$

\Rightarrow reviews,iloc [0:10]

O/P

0

1

2

3

4

5

6

7

8

⑨

end

\Rightarrow reviews, loc [0:10]

O/P

0

1

2

3

4

5

6

7

8

9

Manipulating the index:

=> reviews.set_index("title")

O/P

title name ↗, ↘ column ↗
 ↗ index ↗, ↘ table ↗
 ↗ ↗ ↗ ↗ ↗ ↗

Conditional selection:

=> reviews.country == 'Italy'

O/P:

'Italy' ↗, ↘ row ↗, country
 ↗, ↘ column ↗, ↘ show ↗, ↘ True &
 ↗, ↘ False ↗, ↘ show ↗, ↘

=> reviews.loc[reviews.country == 'Italy']

=> reviews.loc[0]

2025.01.16 08:57

→ 'Italy' ↗, ↘ column ↗, ↘

=> reviews.loc[(reviews.country == 'Italy') &
(reviews.points >= 90)]

=> reviews.loc[(reviews.country == 'Italy') |
(reviews.points >= 90)]

Selecting data in a list of values.

=> reviews.loc[reviews.country.isin(['Italy',
'France'])]

Selecting Null or not Null values

=> reviews.loc[reviews.Price.notnull()]

=> reviews.loc[reviews.Price.isnull()]

Assigning data

→ reviews['critic'] = 'everyone'

O/P:

critic name → കരിക്ക് column create 22
സ്റ്റോറ നേരിലെ ഫോറ്മാറ്റ് സൗഖ്യം
സ്റ്റോറ കുടുംബം എല്ലാവർക്ക്
സ്റ്റോറ എല്ലാവർക്ക് സൗഖ്യം

critic
0 everyone
1 everyone
2 everyone
3 everyone
...
n everyone

→ reviews['block'] = range(len(reviews), 0, -1)

0 n
1 n-1
...
n-1 2
n 1

2025.01.16 08:58

n ഒരു പദ്ധതിയും
-1 ഒരു പദ്ധതിയും
0 ഒരു പദ്ധതിയും

Summary

Functions

.describe()

reviews.points.describe()

column
of
reviews
table

O/P

count @ 129971.00

mean 88.447 ...

75% 91.00

max 100.00

Name: points, Length: 8, dtype: float.

=> reviews.taster_name.describe

O/P

count 103727

unique 19

top chofan

freq 25514

.describe() is
data-type aware,
that's why we got
different summary

for different numerical & string column.

=> reviews.points.mean()

O/P

88.447

=> reviews.taster_name.unique()

O/P

taster_name column → unique value

गुरु विजय

⇒ => reviews.taster_name.value_counts()

O/P

Unique name : ~~frequent~~ how many times

Roger. Voss

25519

Michael Schachner

15139

- for different numerical & string column.

=> reviews.Points.mean()

O/P

88.447

=> reviews.taster_name.unique()

O/P

taster-name column vs unique value

गुरु विजय

=> reviews.taster_name.value_counts()

O/P

Unique name frequent how many times

Roger. Voss

25514

Michael Schachner

15134

Maps

Takes one set of values and "maps" them to another set of values.

There are two mapping methods.

→ `map()` [slightly simpler]

→ `apply()`

map()

⇒ `review_points_mean = reviews.points.mean()`

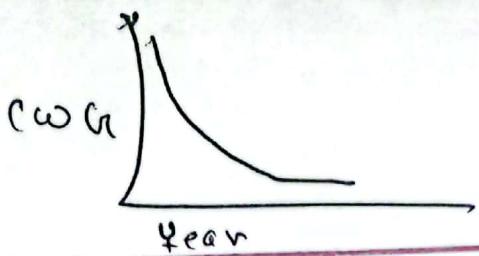
`reviews.points.map(lambda p: p - review_points_mean)`

P \Rightarrow A function.
P \Leftrightarrow P = points - review_points_mean

maps ~~are~~ return a new series.

apply()

To return a whole DataFrame `apply()` is used.



`> def remean_points(row):`

`row.points = row.points - review_points_mean`

`return row`

`reviews.apply(remean_points, axis=1)`

O/P

`reviews.apply func 92 rows 0 rows`

~~apply~~ remean_points ^{func} a call 22(2) 1

~~func~~ 22(2) row.points in change 92

~~func~~

Result ~~func~~ point column modified

~~func~~ dataframe return 92 1

~~func~~ variable

$\Rightarrow \text{review_points_mean} = \text{reviews.points.mean}$

reviews.points - review_points_mean
2025.01.16 08:59

\Rightarrow reviews.country + " - " + reviews.region - 1

Grouping

Q3
=> reviews.groupby('points').price.min()

O/P

points

ପ୍ରକାଶକ

group

၅၇၃

୪୮

price

price

۱۷۴

price. min() ରୁ ଫରା ହାତେ price. (ନୟାପାତ୍ର) ,

6

```
=> reviews.groupby('winery').apply(lambda df:
```

df.title.iloc[0])

OP

~~WINE TASTING~~ unique winery so set up first

element

፩፻፲፭ 2025.01.16 08:59

Multiple Groupby

=> reviews.groupby(['country', 'province']).
apply(lambda df: df.loc[dt.points.idxmax()])

O/P:

Country અનુયાયી group ૨૮૦ રાજ્યોને province
ગુજરાતી group ૨૮૧ Then ૨૮૧ loc વડ.
point max ગુજરાતી ૨૮૧ data ગુજરાતી

Aggregate Functions using groupby

=> reviews.groupby(['country']).price.agg([len, min, max])

O/P

કાર્યોને group વડ. len, min price, max price
wise ~~sort~~ O/P 2025-08-08 08:59

Multi-indexes

=? countries_reviewed = reviews.groupby([['country', 'province']]).description.agg(len)

O/P

len doesn't count the number of char instead how many descriptions there are.

Converting back to regular index

=? countries_reviewed.reset_index()

Sorting

=? countries_reviewed.sort_values(by='len')

O/P

2025.01.16 08:59

groupby vs no sort in 'len' vs Refactor

Ascending Order ০. $10 \rightarrow 0$

Descending Order

countries_reviewed.sort_values(by='len',
ascending=False)

Sort by sorting by index

countries_reviewed.sort_index()

N.B: By default index set sort অর্থাৎ এই
বিষয়ে কোনো sorting technique use.
অর্থাৎ এই

Sorting by Multiple columns

=>countries_reviewed.sort_values(by=['country',
'len'])

Data Types & Missing Values

`@ dtype('float64')`

=> reviews.price.dtype

O/P

`dtype('float64')`

=> reviews.dtypes

O/P

country	object
description	object
...	...
variety	object
winery	object

Missing data

NaN → Not a Number

NaN is always float64 type

⇒ reviews[pd.isnull(reviews, country)]
↓
Σ DF ↓
Σ column → Σ NaN
Σ Σ Σ

Replacing missing values with fillna()

⇒ reviews.region_2.fillna("Unknown")

O/P

~~region_2~~

0 Unknown

1 Unknown

..

m-1

"

n

"

2025.01.16 08:59

Replacing non-null value

⇒ reviews[^{↑ DF} & taster-^{↑ column} twitter-handle].replace
(" @balsal ", "@bal2sal2")

Changing dtype

=> reviews.points.astype('float64')

Renaming

⇒ reviews.rename(columns = {'points': 'score'})
previous changed

Renaming entries

⇒ reviews.rename(index = {0: 'firstEntry',
1: 'secondEntry'})

which entry to change

I

change to what

Renaming index

→ reviews.rename_axis("wines", axis='rows')
· rename_axis("fields", axis='columns')

By Default

columns	A	B	C	D	E	F
Rows	E	M	F	T	X	
0	-	-	-	-	-	
1	-	-	-	-	-	,
2	-	-	-	-	-	.
3						

This will change to:

titles
wines

Rest is same

Combining

concat()

=> canadian_youtube = pd.read_csv("link~~1~~.csv")
bangladeshi_youtube = pd.read_csv("link2.csv")
pd.concat([canadian_youtube, bangladeshi_youtube])

In this case the ~~row~~ columns of both
~~CSV~~ files dataFrames or series must be same.

left = can

join()

=> left = can-yt.set_index(['title', 'trending-date'])

right = bd-yt.set_index(['title', 'trending-date'])

left.join(right, lsuffix=' - CAN', rsuffix=' - BD')
কোন DF যোগ করা হবে
কোন DF উপর যোগ করা হবে 2025/01/15 10:00:00
left এর পার্ট অসম্ভব হবে - CAN এর পার্ট
প্রতিকর্তা হবে - BD এর পার্ট