



PANDAS

Library



Pandas

import Pandas as pd

pd is object

pd. methods() ← we gonna use for

To Read a any file: CSV, json, Savl

CSV - pd.read_csv("path")

json - pd.read_json("path")

Savl - pd.read_Sav("path", compression)

To Know the type

df = pd.read_csv("data.csv")

type(df) = DataFrame

In excel, row and column are called series.

Accessing any columns:-

df[["City"]] → type(df[["City"]]) → Series

This will shows the City rows data

df['City'] → type(df['City']) → Series

Skip the row:-

df.read_csv('data.csv'), skiprows=4
(4 rows are skipped)

Dimension of DataFrame

df.shape = (Total row, total colm)

df.shape[0] = column

first and last top & bottom of DataFrame:-

→ df.head() → top few data

→ df.tail() → last few data

Datframe. information

df.info() - return information about data

value_counts()

gives the analysis of most occurring data frequencies.

df.city.value_counts()

ans condy = True / False

dropna = drop empty values

sort_values()

→ give the sorted series

sort_values(by=)

→ sort the multiple series at once.

new series = df.city.sort_values()



Because its inplace = False means it will

~~create~~ create a new series instead of overwriting the old one.

Boolean Indexing:-

AND = &

OR = |

NOT = ~

df. Metal == 'Gold' → Returns Tom/Roh

df [df. Metal == "Gold"] → Returns Rohan which has Gold

df [(df. Metal == "Gold") & (df. City == "Delhi")]

This return DataFrame which consists
having Gold and living in Delhi

String handling:- (Str. contains(" "))

→ df. name. str. contains(" Kushnora")

→ Return Tom / Roh in all DataFrame

→ df [df. name. str. contains (" Kushnora")]

→ Return DataFrame which contains
Kushnora

challenge

- 1) In which events did Jesse Owens win a medal?

sol-

df[(df.Athlete == "Jesse Owens")
 & (df.Medal == "Gold" or
 df.Medal == "Silver" or
 df.Medal == "Bronze")]

OR

new = df[df.Athlete == "Jesse Owens"]

→ give all Jesse Owens Events

stv in new variable

→ new.events.value_counts()

→ give all events.

2) Most men gold medal in badminton
Sort by Player name.

→ $\text{df}[(\text{df. metal} == \text{"Gold"}) \& (\text{df. gender} == \text{"men"}) \& (\text{df. sport} == \text{"badminton"})]$

→ Store in variable.

→ $\text{new_sort_values}(\text{by} = \text{"Athlete,"})$

3) Which three countries have won the most medals in recent years (1984 to 2008)

$\text{new} = \text{df}(\text{df. year} >= 1984)$

$\text{new.NOC.value} \text{conts}().head(3)$

give result of three country
having high medals.

y) Rishabh made gold medal for 100m from
most recent with City, Bdivision, Athlete
and country -

new = df[(df['gender' == "Male"] & (df['medal' == "Gold"]) & (df['sport' == "100m"])]]

new. sort values(['Bdivision'], axis = 0, ascending = False)

to print extra series thus gives

{['City', 'Bdivision', 'Athlete', 'NotC']}

BASIC PLOTTING

→ Matplotlib

Import matplotlib.pyplot as plt
%matplotlib inline

→ merger command → helps to
view graphs in notebook

Plot types → plot()

- plot(kind = 'line') → Start time period
- plot(kind = 'bar') → change on large
- plot(kind = 'barh') → some
- plot(kind = 'pie') → Show categories

f0 = df[df.event == 1846]

f0.head(5)

f0['sport'].value_counts().plot(); → to remove
default = line
the object
the

color

fo-sport.value_count().plot (kind='bar',)
color = 'monotone')

figsize

plot (figsize = (width, height))

fo-sport.value_count().plot (kind='bar',
figsize=(10,3))

CountPlot () — for complex data

Seaborn.countplot () :

data - the dataset

hue - categorical variables

order - sequence when using categorical level

palette - colors to use for different levels of
hue variable

import seaborn as sns

sns.countplot (x = 'medal', data=fo)



Sns. countplot (x = 'medal', data=fo, hue=
'Gender')



Indexing:-

1) type (fo.index) → RangeIndex

fo.index(10) → give some value

Indx are immutable

2) Set_Index (keys, inplace = True)

0
1
2
3
4

Index

fo.set_index('Athlete')

first col become Athlete

3) reset_index () → Reset the set index
vsif inplace False
→ df.reset_index (inplace = True)

sort_index()

df. sort_index (i-place = True)

→ Replace and Sort the index

loc[] → Dataframe.loc[]

loc[] will raise a Key Error
when the items are not found

→ first set the index by df.set_index

→ df. [df. Atlete = "USA20 Bolt"]

give some result or error

→ df. loc[df. Atlete = "USA20. Bolt"]

iloc[] → gives specific result of list
of index.

df. loc[[102, 222, 1500]]

shows specific result of index of the
in list.

`df.loc[1:4]` → Show 1, 2, 3

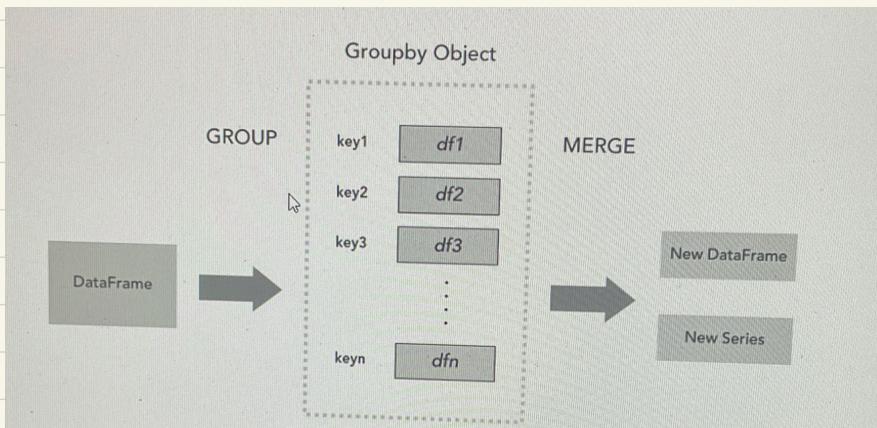
Index marks

Groupby → split backbone into groups

Returns
an object

→ Apply function to each group
independently -

→ Combine the result into DataFrame



`df.groupby('Edition')`

To view what's in group:-

`list(df.groupby('Edition'))`

Iterate through a group

for key, group in Dataframe.groupby()

print(key)

print(group)

for groupkey, groupvals in df.groupby('Edm')

print(groupkey)

print(groupvals) # type = Dataframe

groupby

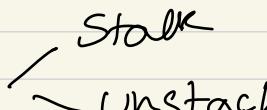
df.groupby('Edition').size()

df.groupby(['Edition', 'NC', 'Medl']).

agg('count')

`df.loc[df['Athlete'] == 'Lewis, Carl'].groupby`

`('Athlete').agg([{'Edition': ['min', 'max', 'last']}`

Reshaping:-  ?]

Stack - allow to move inner col to the rows for the ~~dataframe~~

→ Return ~~dataframe~~ and series

→ by default dropna = True

No NA value will be show

→ Pivot a level of the column labels,
returning a ~~dataframe~~ of series,
with a new innermost level of
row label.

Diagram illustrating the effect of `DataFrame.stack()`. An arrow points from the original wide DataFrame on the left to a taller, stacked DataFrame on the right. The code `DataFrame.stack()` is written above the arrow.

Original DataFrame (Wide Format):

	Discipline	Athletics	
		Event	100m
NOC	Gender		
JAM	Men	1.0	1.0
	Women	3.0	2.0
TRI	Men	1.0	NaN
	Men	1.0	2.0
USA	Women	NaN	1.0

Stacked DataFrame (Tall Format):

NOC	Gender	Discipline	Athletics
		Event	
JAM	Men	100m	1.0
		200m	1.0
JAM	Women	100m	3.0
		200m	2.0
TRI	Men	100m	1.0
		200m	1.0
USA	Men	100m	2.0
		200m	1.0
USA	Women	200m	1.0

df.stack()

Unstack:-

→ Return Series and DataFrame

→ missing data values are displayed

→ Unpack the Stack data

challenges

import pandas as pd

df = pd.read_csv('..')

print(df.head(10))