

Pandas

```
In [2]: import pandas as pd
```

```
In [3]: df_infy = pd.read_csv("INFY.csv")
```

```
In [4]: df_infy.head()
```

```
Out[4]:
```

	Date	Open	High	Low	Close	Adj Close	Volume
0	2021-09-17	22.950001	22.990000	22.719999	22.850000	22.401333	8070400
1	2021-09-20	22.750000	22.950001	22.570000	22.760000	22.313101	7148800
2	2021-09-21	23.129999	23.230000	23.100000	23.120001	22.666033	4785000
3	2021-09-22	23.120001	23.379999	23.070000	23.270000	22.813086	6613800
4	2021-09-23	23.680000	23.790001	23.600000	23.730000	23.264053	5323800

```
In [5]: df_infy.sort_values(by = "Volume", ascending = False)
```

```
Out[5]:
```

	Date	Open	High	Low	Close	Adj Close	Volume
251	2022-09-16	17.200001	17.459999	17.090000	17.379999	17.379999	42686600
126	2022-03-18	24.389999	25.040001	24.209999	25.040001	24.763243	29035600
154	2022-04-28	20.410000	20.660000	20.350000	20.500000	20.273422	24229100
20	2021-10-15	22.900000	23.400000	22.760000	23.379999	22.920927	22295600
250	2022-09-15	17.700001	17.809999	17.510000	17.790001	17.790001	21950100
...
232	2022-08-19	20.040001	20.040001	19.750000	19.780001	19.780001	2868500
132	2022-03-28	24.730000	24.770000	24.520000	24.719999	24.446779	2834400
68	2021-12-23	24.690001	24.799999	24.500000	24.730000	24.456669	2488600
71	2021-12-29	25.190001	25.379999	25.139999	25.379999	25.099483	2330400
72	2021-12-30	25.540001	25.600000	25.389999	25.410000	25.129152	2293900

252 rows × 7 columns

```
In [6]: df_dict = {'player_id': {0: 401, 1: 401, 2: 401, 3: 401, 4: 401, 5: 401, 6: 401, 7
```

```
In [7]: df = pd.DataFrame(df_dict)
df
```

Out[7]:

	player_id	match_date	match_result
0	401	2021-05-04 00:00:00	W
1	401	2021-05-09 00:00:00	L
2	401	2021-05-16 00:00:00	L
3	401	2021-05-18 00:00:00	W
4	401	2021-05-22 00:00:00	L
5	401	2021-06-15 00:00:00	L
6	401	2021-06-16 00:00:00	W
7	401	2021-06-18 00:00:00	W
8	401	2021-07-06 00:00:00	L
9	401	2021-07-13 00:00:00	L
10	402	2021-05-14 00:00:00	L
11	402	2021-05-23 00:00:00	L
12	402	2021-05-24 00:00:00	W
13	402	2021-06-01 00:00:00	W
14	402	2021-06-02 00:00:00	W
15	402	2021-07-01 00:00:00	W
16	402	2021-07-11 00:00:00	W
17	402	2021-07-20 00:00:00	L
18	402	2021-07-26 00:00:00	L
19	402	2021-07-30 00:00:00	L
20	403	2021-05-03 00:00:00	L
21	403	2021-05-11 00:00:00	W
22	403	2021-05-12 00:00:00	W
23	403	2021-05-13 00:00:00	W
24	403	2021-05-20 00:00:00	W
25	403	2021-05-25 00:00:00	W
26	403	2021-07-06 00:00:00	L
27	403	2021-07-15 00:00:00	L
28	403	2021-07-22 00:00:00	W
29	403	2021-07-23 00:00:00	W
30	404	2021-05-10 00:00:00	W
31	404	2021-05-16 00:00:00	W
32	404	2021-05-20 00:00:00	W

	player_id	match_date	match_result
33	404	2021-05-22 00:00:00	W
34	404	2021-05-28 00:00:00	L
35	404	2021-06-06 00:00:00	L
36	404	2021-06-14 00:00:00	W
37	404	2021-07-25 00:00:00	W
38	404	2021-07-26 00:00:00	L
39	405	2021-05-07 00:00:00	L
40	405	2021-05-25 00:00:00	L
41	405	2021-06-06 00:00:00	L
42	405	2021-06-07 00:00:00	L
43	405	2021-06-14 00:00:00	L
44	405	2021-07-01 00:00:00	L
45	405	2021-07-02 00:00:00	L
46	405	2021-07-14 00:00:00	W
47	405	2021-07-16 00:00:00	L
48	405	2021-07-30 00:00:00	L

```
In [8]: # with open("csv/sample.json", "w") as file :  
#     json.dump(df_dict,file) #javascript
```

```
In [9]: df[df["player_id"] == 401]
```

	player_id	match_date	match_result
0	401	2021-05-04 00:00:00	W
1	401	2021-05-09 00:00:00	L
2	401	2021-05-16 00:00:00	L
3	401	2021-05-18 00:00:00	W
4	401	2021-05-22 00:00:00	L
5	401	2021-06-15 00:00:00	L
6	401	2021-06-16 00:00:00	W
7	401	2021-06-18 00:00:00	W
8	401	2021-07-06 00:00:00	L
9	401	2021-07-13 00:00:00	L

```
In [10]: df["player_id"] == 401
```

```
Out[10]: 0      True
         1      True
         2      True
         3      True
         4      True
         5      True
         6      True
         7      True
         8      True
         9      True
        10     False
        11     False
        12     False
        13     False
        14     False
        15     False
        16     False
        17     False
        18     False
        19     False
        20     False
        21     False
        22     False
        23     False
        24     False
        25     False
        26     False
        27     False
        28     False
        29     False
        30     False
        31     False
        32     False
        33     False
        34     False
        35     False
        36     False
        37     False
        38     False
        39     False
        40     False
        41     False
        42     False
        43     False
        44     False
        45     False
        46     False
        47     False
        48     False
Name: player_id, dtype: bool
```

```
In [11]: df_infy
```

Out[11]:

	Date	Open	High	Low	Close	Adj Close	Volume
0	2021-09-17	22.950001	22.990000	22.719999	22.850000	22.401333	8070400
1	2021-09-20	22.750000	22.950001	22.570000	22.760000	22.313101	7148800
2	2021-09-21	23.129999	23.230000	23.100000	23.120001	22.666033	4785000
3	2021-09-22	23.120001	23.379999	23.070000	23.270000	22.813086	6613800
4	2021-09-23	23.680000	23.790001	23.600000	23.730000	23.264053	5323800
...
247	2022-09-12	19.230000	19.410000	19.190001	19.240000	19.240000	3607400
248	2022-09-13	18.760000	18.900000	18.490000	18.559999	18.559999	15227800
249	2022-09-14	18.370001	18.440001	17.969999	18.080000	18.080000	16757300
250	2022-09-15	17.700001	17.809999	17.510000	17.790001	17.790001	21950100
251	2022-09-16	17.200001	17.459999	17.090000	17.379999	17.379999	42686600

252 rows × 7 columns

In [12]: df_infy["Volume"]

```
Out[12]: 0      8070400
         1      7148800
         2      4785000
         3      6613800
         4      5323800
         ...
        247     3607400
        248     15227800
        249     16757300
        250     21950100
        251     42686600
Name: Volume, Length: 252, dtype: int64
```

In [13]: df_infy["Volume"].sum()

Out[13]: 2381188600

how many unique player id are there

In [14]: df["player_id"].unique()

Out[14]: array([401, 402, 403, 404, 405], dtype=int64)

In [15]: df

	player_id	match_date	match_result
0	401	2021-05-04 00:00:00	W
1	401	2021-05-09 00:00:00	L
2	401	2021-05-16 00:00:00	L
3	401	2021-05-18 00:00:00	W
4	401	2021-05-22 00:00:00	L
5	401	2021-06-15 00:00:00	L
6	401	2021-06-16 00:00:00	W
7	401	2021-06-18 00:00:00	W
8	401	2021-07-06 00:00:00	L
9	401	2021-07-13 00:00:00	L
10	402	2021-05-14 00:00:00	L
11	402	2021-05-23 00:00:00	L
12	402	2021-05-24 00:00:00	W
13	402	2021-06-01 00:00:00	W
14	402	2021-06-02 00:00:00	W
15	402	2021-07-01 00:00:00	W
16	402	2021-07-11 00:00:00	W
17	402	2021-07-20 00:00:00	L
18	402	2021-07-26 00:00:00	L
19	402	2021-07-30 00:00:00	L
20	403	2021-05-03 00:00:00	L
21	403	2021-05-11 00:00:00	W
22	403	2021-05-12 00:00:00	W
23	403	2021-05-13 00:00:00	W
24	403	2021-05-20 00:00:00	W
25	403	2021-05-25 00:00:00	W
26	403	2021-07-06 00:00:00	L
27	403	2021-07-15 00:00:00	L
28	403	2021-07-22 00:00:00	W
29	403	2021-07-23 00:00:00	W
30	404	2021-05-10 00:00:00	W
31	404	2021-05-16 00:00:00	W
32	404	2021-05-20 00:00:00	W

	player_id	match_date	match_result
33	404	2021-05-22 00:00:00	W
34	404	2021-05-28 00:00:00	L
35	404	2021-06-06 00:00:00	L
36	404	2021-06-14 00:00:00	W
37	404	2021-07-25 00:00:00	W
38	404	2021-07-26 00:00:00	L
39	405	2021-05-07 00:00:00	L
40	405	2021-05-25 00:00:00	L
41	405	2021-06-06 00:00:00	L
42	405	2021-06-07 00:00:00	L
43	405	2021-06-14 00:00:00	L
44	405	2021-07-01 00:00:00	L
45	405	2021-07-02 00:00:00	L
46	405	2021-07-14 00:00:00	W
47	405	2021-07-16 00:00:00	L
48	405	2021-07-30 00:00:00	L

```
In [16]: df = pd.DataFrame(df_dict)
df
```

	player_id	match_date	match_result
0	401	2021-05-04 00:00:00	W
1	401	2021-05-09 00:00:00	L
2	401	2021-05-16 00:00:00	L
3	401	2021-05-18 00:00:00	W
4	401	2021-05-22 00:00:00	L
5	401	2021-06-15 00:00:00	L
6	401	2021-06-16 00:00:00	W
7	401	2021-06-18 00:00:00	W
8	401	2021-07-06 00:00:00	L
9	401	2021-07-13 00:00:00	L
10	402	2021-05-14 00:00:00	L
11	402	2021-05-23 00:00:00	L
12	402	2021-05-24 00:00:00	W
13	402	2021-06-01 00:00:00	W
14	402	2021-06-02 00:00:00	W
15	402	2021-07-01 00:00:00	W
16	402	2021-07-11 00:00:00	W
17	402	2021-07-20 00:00:00	L
18	402	2021-07-26 00:00:00	L
19	402	2021-07-30 00:00:00	L
20	403	2021-05-03 00:00:00	L
21	403	2021-05-11 00:00:00	W
22	403	2021-05-12 00:00:00	W
23	403	2021-05-13 00:00:00	W
24	403	2021-05-20 00:00:00	W
25	403	2021-05-25 00:00:00	W
26	403	2021-07-06 00:00:00	L
27	403	2021-07-15 00:00:00	L
28	403	2021-07-22 00:00:00	W
29	403	2021-07-23 00:00:00	W
30	404	2021-05-10 00:00:00	W
31	404	2021-05-16 00:00:00	W
32	404	2021-05-20 00:00:00	W

	player_id	match_date	match_result
33	404	2021-05-22 00:00:00	W
34	404	2021-05-28 00:00:00	L
35	404	2021-06-06 00:00:00	L
36	404	2021-06-14 00:00:00	W
37	404	2021-07-25 00:00:00	W
38	404	2021-07-26 00:00:00	L
39	405	2021-05-07 00:00:00	L
40	405	2021-05-25 00:00:00	L
41	405	2021-06-06 00:00:00	L
42	405	2021-06-07 00:00:00	L
43	405	2021-06-14 00:00:00	L
44	405	2021-07-01 00:00:00	L
45	405	2021-07-02 00:00:00	L
46	405	2021-07-14 00:00:00	W
47	405	2021-07-16 00:00:00	L
48	405	2021-07-30 00:00:00	L

```
In [17]: df["player_id"].unique() #
```

```
Out[17]: array([401, 402, 403, 404, 405], dtype=int64)
```

```
In [18]: df["player_id"].nunique() # count of unique columns
```

```
Out[18]: 5
```

```
In [19]: spotify_dict = {'id': {0: 303651, 1: 85559, 2: 1046089, 3: 350824, 4: 776822, 5: 46
```

```
In [20]: df_spotify = pd.DataFrame(spotify_dict)
```

```
In [21]: df_spotify
```

Out[21]:

	id	position	trackname	artist	streams	url	date	region
0	303651	52	Heart Won't Forget	Matoma	28047	https://open.spotify.com/track/2of2DM5LqTh7ohm...	2017-02-04 00:00:00	nc
1	85559	160	Someone In The Crowd - From "La La Land" Sound...	Emma Stone	17134	https://open.spotify.com/track/7xE4vKvjqUTtHyJ...	2017-02-26 00:00:00	fi
2	1046089	175	The Greatest	Sia	10060	https://open.spotify.com/track/7xHWNBfM6ObGEQP...	2017-03-06 00:00:00	c
3	350824	25	Unforgettable	French Montana	46603	https://open.spotify.com/track/3B54sVLJ402zGa6...	2017-10-01 00:00:00	nc
4	776822	1	Bad and Boujee (feat. Lil Uzi Vert)	Migos	1823391	https://open.spotify.com/track/4Km5HrUvYTaSUfi...	2017-01-27 00:00:00	us
...
95	792423	2	DNA.	Kendrick Lamar	3013496	https://open.spotify.com/track/6HZILIRieu8S0iq...	2017-04-15 00:00:00	us
96	792223	2	DNA.	Kendrick Lamar	3643231	https://open.spotify.com/track/6HZILIRieu8S0iq...	2017-04-14 00:00:00	us
97	793422	1	HUMBLE.	Kendrick Lamar	3144482	https://open.spotify.com/track/7KXjTSCq5nL1LoY...	2017-04-20 00:00:00	us
98	793622	1	HUMBLE.	Kendrick Lamar	3172718	https://open.spotify.com/track/7KXjTSCq5nL1LoY...	2017-04-21 00:00:00	us
99	793022	1	HUMBLE.	Kendrick Lamar	3394456	https://open.spotify.com/track/7KXjTSCq5nL1LoY...	2017-04-18 00:00:00	us

100 rows × 8 columns

In [22]: df_spotify["streams"]

Out[22]:

```
0    28047
1    17134
2    10060
3    46603
4    1823391
      ...
95   3013496
96   3643231
97   3144482
98   3172718
99   3394456
Name: streams, Length: 100, dtype: int64
```

```
In [23]: df_spotify["streams"].sum()
```

```
Out[23]: 39673624
```

```
In [24]: df_spotify["streams"].sum()/100
```

```
Out[24]: 396736.24
```

```
In [25]: host_dict = {'host_id': {0: 0, 1: 1, 2: 2, 3: 3, 4: 4, 5: 5, 6: 6, 7: 7, 8: 8, 9: 9}}
```

```
In [26]: df_host = pd.DataFrame(host_dict)
df_host
```

```
Out[26]:
```

	host_id	nationality	gender	age
0	0	USA	M	28
1	1	USA	F	29
2	2	China	F	31
3	3	China	M	24
4	4	Mali	M	30
...
171	7	Luxembourg	F	25
172	6	Luxembourg	M	25
173	7	Luxembourg	F	25
174	6	Luxembourg	M	25
175	7	Luxembourg	F	25

176 rows × 4 columns

```
In [27]: appart_dict = {'host_id': {0: 0, 1: 0, 2: 0, 3: 1, 4: 1, 5: 2, 6: 3, 7: 3, 8: 4, 9: 9}}
```

```
In [28]: df_apprt = pd.DataFrame(appart_dict)
df_apprt
```

Out[28]:

	host_id	apartment_id	apartment_type	n_beds	n_bedrooms	country	city
0	0	A1	Room	1	1	USA	New York
1	0	A2	Room	1	1	USA	New Jersey
2	0	A3	Room	1	1	USA	New Jersey
3	1	A4	Apartment	2	1	USA	Houston
4	1	A5	Apartment	2	1	USA	Las Vegas
5	2	A6	Yurt	3	1	Mongolia	-
6	3	A7	Penthouse	3	3	China	Tianjin
7	3	A8	Penthouse	5	5	China	Beijing
8	4	A9	Apartment	2	1	Mali	Bamako
9	5	A10	Room	3	1	Mali	Segou
10	5	A11	Room	2	1	Mali	Segou
11	6	A12	Penthouse	6	6	Luxembourg	Luxembourg
12	7	A13	Room	4	1	Luxembourg	Luxembourg
13	8	A14	Apartment	2	1	Australia	Perth
14	9	A15	Apartment	2	1	Australia	Perth
15	9	A16	Apartment	2	1	Australia	Perth
16	10	A17	Room	4	1	Brazil	Rio De Janeiro
17	10	A18	Room	4	1	Argentina	Mendoza
18	10	A19	Room	4	2	Uruguay	Mercedes
19	10	A20	Room	4	2	Brazil	Brasilia
20	11	A21	Apartment	2	2	Mexico	Mexico City

In [29]: df_host_apprt = pd.merge(df_host,df_apprt, on = "host_id")

In [30]: df_host_apprt

Out[30]:

	host_id	nationality	gender	age	apartment_id	apartment_type	n_beds	n_bedrooms	count
0	0	USA	M	28	A1	Room	1	1	U
1	0	USA	M	28	A2	Room	1	1	U
2	0	USA	M	28	A3	Room	1	1	U
3	0	USA	M	28	A1	Room	1	1	U
4	0	USA	M	28	A2	Room	1	1	U
...
347	10	Brazil	M	39	A20	Room	4	2	Bra
348	11	Brazil	F	42	A21	Apartment	2	2	Mexi
349	11	Brazil	F	42	A21	Apartment	2	2	Mexi
350	11	Brazil	F	42	A21	Apartment	2	2	Mexi
351	11	Brazil	F	42	A21	Apartment	2	2	Mexi

352 rows × 10 columns

In [31]: df_host_apprt[df_host_apprt["host_id"]==0]

Out[31]:

	host_id	nationality	gender	age	apartment_id	apartment_type	n_beds	n_bedrooms	count
0	0	USA	M	28	A1	Room	1	1	U
1	0	USA	M	28	A2	Room	1	1	U
2	0	USA	M	28	A3	Room	1	1	U
3	0	USA	M	28	A1	Room	1	1	U
4	0	USA	M	28	A2	Room	1	1	U
...
211	0	USA	M	28	A2	Room	1	1	U
212	0	USA	M	28	A3	Room	1	1	U
213	0	USA	M	28	A1	Room	1	1	U
214	0	USA	M	28	A2	Room	1	1	U
215	0	USA	M	28	A3	Room	1	1	U

216 rows × 10 columns

In [32]: df_host_apprt[df_host_apprt["nationality"] != df_host_apprt["country"]]

Out[32]:

	host_id	nationality	gender	age	apartment_id	apartment_type	n_beds	n_bedrooms	cou
232	2	China	F	31	A6	Yurt	3	1	Mong
233	2	China	F	31	A6	Yurt	3	1	Mong
234	2	China	F	31	A6	Yurt	3	1	Mong
235	2	China	F	31	A6	Yurt	3	1	Mong
333	10	Brazil	M	39	A18	Room	4	1	Arger
334	10	Brazil	M	39	A19	Room	4	2	Uruç
337	10	Brazil	M	39	A18	Room	4	1	Arger
338	10	Brazil	M	39	A19	Room	4	2	Uruç
341	10	Brazil	M	39	A18	Room	4	1	Arger
342	10	Brazil	M	39	A19	Room	4	2	Uruç
345	10	Brazil	M	39	A18	Room	4	1	Arger
346	10	Brazil	M	39	A19	Room	4	2	Uruç
348	11	Brazil	F	42	A21	Apartment	2	2	Me
349	11	Brazil	F	42	A21	Apartment	2	2	Me
350	11	Brazil	F	42	A21	Apartment	2	2	Me
351	11	Brazil	F	42	A21	Apartment	2	2	Me

In [33]: df_host_apprt[df_host_apprt["nationality"] == df_host_apprt["country"]]

Out[33]:

	host_id	nationality	gender	age	apartment_id	apartment_type	n_beds	n_bedrooms	count
0	0	USA	M	28	A1	Room	1	1	U
1	0	USA	M	28	A2	Room	1	1	U
2	0	USA	M	28	A3	Room	1	1	U
3	0	USA	M	28	A1	Room	1	1	U
4	0	USA	M	28	A2	Room	1	1	U
...
339	10	Brazil	M	39	A20	Room	4	2	Bra
340	10	Brazil	M	39	A17	Room	4	1	Bra
343	10	Brazil	M	39	A20	Room	4	2	Bra
344	10	Brazil	M	39	A17	Room	4	1	Bra
347	10	Brazil	M	39	A20	Room	4	2	Bra

336 rows × 10 columns

In [34]: `df_host_apprt[df_host_apprt["nationality"] != df_host_apprt["country"]]["host_id"].`

Out[34]: 3

Joins are very important : https://www.w3schools.com/sql/sql_join.asp

rename

In [35]: `df_host_apprt.rename(columns = {"host_id" : "renamed_host_id"}, inplace = True)`

In [36]: `df_host_apprt`

Out[36]:

	renamed_host_id	nationality	gender	age	apartment_id	apartment_type	n_beds	n_bedroom
0	0	USA	M	28	A1	Room	1	
1	0	USA	M	28	A2	Room	1	
2	0	USA	M	28	A3	Room	1	
3	0	USA	M	28	A1	Room	1	
4	0	USA	M	28	A2	Room	1	
...
347	10	Brazil	M	39	A20	Room	4	
348	11	Brazil	F	42	A21	Apartment	2	
349	11	Brazil	F	42	A21	Apartment	2	
350	11	Brazil	F	42	A21	Apartment	2	
351	11	Brazil	F	42	A21	Apartment	2	

352 rows × 10 columns

In [37]: df_host_apprt.columns

Out[37]: Index(['renamed_host_id', 'nationality', 'gender', 'age', 'apartment_id', 'apartment_type', 'n_beds', 'n_bedrooms', 'country', 'city'], dtype='object')

Date wide data is called timeseries data we can use this to predict something

19th sept 2022

In [38]: df_infy

Out[38]:

	Date	Open	High	Low	Close	Adj Close	Volume
0	2021-09-17	22.950001	22.990000	22.719999	22.850000	22.401333	8070400
1	2021-09-20	22.750000	22.950001	22.570000	22.760000	22.313101	7148800
2	2021-09-21	23.129999	23.230000	23.100000	23.120001	22.666033	4785000
3	2021-09-22	23.120001	23.379999	23.070000	23.270000	22.813086	6613800
4	2021-09-23	23.680000	23.790001	23.600000	23.730000	23.264053	5323800
...
247	2022-09-12	19.230000	19.410000	19.190001	19.240000	19.240000	3607400
248	2022-09-13	18.760000	18.900000	18.490000	18.559999	18.559999	15227800
249	2022-09-14	18.370001	18.440001	17.969999	18.080000	18.080000	16757300
250	2022-09-15	17.700001	17.809999	17.510000	17.790001	17.790001	21950100
251	2022-09-16	17.200001	17.459999	17.090000	17.379999	17.379999	42686600

252 rows × 7 columns

In [41]: `mean_open = df_infy["Open"].mean()`In [46]: `df_infy["far_apart_open"] = df_infy["Open"]**2 - mean_open**2`In [47]: `df_infy`

Out[47]:

	Date	Open	High	Low	Close	Adj Close	Volume	far_apart_open
0	2021-09-17	22.950001	22.990000	22.719999	22.850000	22.401333	8070400	59.995105
1	2021-09-20	22.750000	22.950001	22.570000	22.760000	22.313101	7148800	50.855059
2	2021-09-21	23.129999	23.230000	23.100000	23.120001	22.666033	4785000	68.289413
3	2021-09-22	23.120001	23.379999	23.070000	23.270000	22.813086	6613800	67.827006
4	2021-09-23	23.680000	23.790001	23.600000	23.730000	23.264053	5323800	94.034959
...
247	2022-09-12	19.230000	19.410000	19.190001	19.240000	19.240000	3607400	-96.914541
248	2022-09-13	18.760000	18.900000	18.490000	18.559999	18.559999	15227800	-114.769841
249	2022-09-14	18.370001	18.440001	17.969999	18.080000	18.080000	16757300	-129.250504
250	2022-09-15	17.700001	17.809999	17.510000	17.790001	17.790001	21950100	-153.417405
251	2022-09-16	17.200001	17.459999	17.090000	17.379999	17.379999	42686600	-170.867406

252 rows × 8 columns

In [48]: `summation_of_far_apart_open = df_infy["far_apart_open"].sum()`

```
In [49]: summation_of_far_apart_open
```

```
Out[49]: 1367.4751837390095
```

```
In [50]: var = summation_of_far_apart_open/25  
var
```

```
Out[50]: 5.448108301749042
```

```
In [55]: type(df_infy.shape)
```

```
Out[55]: tuple
```

```
In [56]: df_infy.head()
```

```
Out[56]:
```

	Date	Open	High	Low	Close	Adj Close	Volume	far_apart_open
0	2021-09-17	22.950001	22.990000	22.719999	22.850000	22.401333	8070400	59.995105
1	2021-09-20	22.750000	22.950001	22.570000	22.760000	22.313101	7148800	50.855059
2	2021-09-21	23.129999	23.230000	23.100000	23.120001	22.666033	4785000	68.289413
3	2021-09-22	23.120001	23.379999	23.070000	23.270000	22.813086	6613800	67.827006
4	2021-09-23	23.680000	23.790001	23.600000	23.730000	23.264053	5323800	94.034959

we want make a dict with date and far apart open

```
In [57]: def convert_to_dict(row):  
    return{row["Date"] : row["far_apart_open"]}
```

```
In [59]: df_infy["date_spread_dict"] = df_infy.apply(convert_to_dict, axis = 1)
```

```
In [63]: df_infy.drop(["Date", "far_apart_open"], axis = 1, inplace = True)
```

```
In [64]: df_infy
```

Out[64]:

	Open	High	Low	Close	Adj Close	Volume	date_spread_dict
0	22.950001	22.990000	22.719999	22.850000	22.401333	8070400	{'2021-09-17': 59.995105167695044}
1	22.750000	22.950001	22.570000	22.760000	22.313101	7148800	{'2021-09-20': 50.855059267694}
2	23.129999	23.230000	23.100000	23.120001	22.666033	4785000	{'2021-09-21': 68.28941300769503}
3	23.120001	23.379999	23.070000	23.270000	22.813086	6613800	{'2021-09-22': 67.82700550769488}
4	23.680000	23.790001	23.600000	23.730000	23.264053	5323800	{'2021-09-23': 94.03495926769398}
...
247	19.230000	19.410000	19.190001	19.240000	19.240000	3607400	{'2022-09-12': -96.91454073230597}
248	18.760000	18.900000	18.490000	18.559999	18.559999	15227800	{'2022-09-13': -114.76984073230597}
249	18.370001	18.440001	17.969999	18.080000	18.080000	16757300	{'2022-09-14': -129.25050399230503}
250	17.700001	17.809999	17.510000	17.790001	17.790001	21950100	{'2022-09-15': -153.417405332305}
251	17.200001	17.459999	17.090000	17.379999	17.379999	42686600	{'2022-09-16': -170.867406332305}

252 rows × 7 columns

how to separate the above column date_spread_dict --> apply

In [68]:

```
def extract_col(row):
    key = list(row["date_spread_dict"].keys())[0]
    value = list(row["date_spread_dict"].values())[0]
    row["date"] = key
    row["far_apart_open"] = value
    return row
```

In [69]:

```
df_infy.apply(extract_col, axis = 1)
```

Out[69]:

	Open	High	Low	Close	Adj Close	Volume	date_spread_dict	
0	22.950001	22.990000	22.719999	22.850000	22.401333	8070400	{'2021-09-17': 59.99510516769504}	2021-1
1	22.750000	22.950001	22.570000	22.760000	22.313101	7148800	{'2021-09-20': 50.855059267694}	2021-1
2	23.129999	23.230000	23.100000	23.120001	22.666033	4785000	{'2021-09-21': 68.28941300769503}	2021-1
3	23.120001	23.379999	23.070000	23.270000	22.813086	6613800	{'2021-09-22': 67.82700550769488}	2021-1
4	23.680000	23.790001	23.600000	23.730000	23.264053	5323800	{'2021-09-23': 94.03495926769398}	2021-1
...
247	19.230000	19.410000	19.190001	19.240000	19.240000	3607400	{'2022-09-12': -96.91454073230597}	2022-1
248	18.760000	18.900000	18.490000	18.559999	18.559999	15227800	{'2022-09-13': -114.76984073230597}	2022-1
249	18.370001	18.440001	17.969999	18.080000	18.080000	16757300	{'2022-09-14': -129.25050399230503}	2022-1
250	17.700001	17.809999	17.510000	17.790001	17.790001	21950100	{'2022-09-15': -153.417405332305}	2022-1
251	17.200001	17.459999	17.090000	17.379999	17.379999	42686600	{'2022-09-16': -170.867406332305}	2022-1

252 rows × 9 columns

TQDM--> check

In [70]: `!pip install tqdm`

```
Collecting tqdm
  Downloading tqdm-4.64.1-py2.py3-none-any.whl (78 kB)
----- 78.5/78.5 KB 4.3 MB/s eta 0:00:00
Requirement already satisfied: colorama in c:\users\mause\appdata\local\programs\python\python310\lib\site-packages (from tqdm) (0.4.5)
Installing collected packages: tqdm
Successfully installed tqdm-4.64.1

WARNING: You are using pip version 22.0.4; however, version 22.2.2 is available.
You should consider upgrading via the 'C:\Users\mause\AppData\Local\Programs\Python\Python310\python.exe -m pip install --upgrade pip' command.
```

In [77]: `from tqdm import tqdm`In [78]: `tqdm.pandas()`In [79]: `df_infy = df_infy.progress_apply(extract_col, axis = 1)`

100%|██████████| 252/252 [00:00<00:00, 907.26it/s]

```
In [80]: df_infy.tail(10)
```

```
Out[80]:      Open    High     Low   Close  Adj Close  Volume  date_spread_dict  
242  18.299999  18.440001  18.000000  18.090000  18.090000  5248800  {'2022-09-02':  
                           -131.81747733230503}  2022-  
243  18.139999  18.219999  17.930000  17.980000  17.980000  7158700  {'2022-09-06':  
                           -137.647877012305}  2022-  
244  18.139999  18.450001  18.110001  18.430000  18.430000  4769300  {'2022-09-07':  
                           -137.647877012305}  2022-  
245  18.350000  18.530001  18.280001  18.530001  18.530001  3936900  {'2022-09-08':  
                           -129.98494073230597}  2022-  
246  18.770000  19.100000  18.750000  19.059999  19.059999  5513600  {'2022-09-09':  
                           -114.39454073230604}  2022-  
247  19.230000  19.410000  19.190001  19.240000  19.240000  3607400  {'2022-09-12':  
                           -96.91454073230597}  2022-  
248  18.760000  18.900000  18.490000  18.559999  18.559999  15227800  {'2022-09-13':  
                           -114.76984073230597}  2022-  
249  18.370001  18.440001  17.969999  18.080000  18.080000  16757300  {'2022-09-14':  
                           -129.25050399230503}  2022-  
250  17.700001  17.809999  17.510000  17.790001  17.790001  21950100  {'2022-09-15':  
                           -153.417405332305}  2022-  
251  17.200001  17.459999  17.090000  17.379999  17.379999  42686600  {'2022-09-16':  
                           -170.867406332305}  2022-
```

```
In [84]: df_infy_aj_close = df_infy[["Adj Close","date"]]
```

```
In [85]: df_infy_aj_close.dtypes
```

```
Out[85]: Adj Close    float64  
          date        object  
          dtype: object
```

```
In [86]: date = "2021-02-21"
```

```
In [87]: date.split("-")
```

```
Out[87]: ['2021', '02', '21']
```

```
In [88]: date.split("-")[1]
```

```
Out[88]: '02'
```

```
In [92]: df_infy_aj_close.date.str.split("-")
```

```
Out[92]: 0      [2021, 09, 17]
         1      [2021, 09, 20]
         2      [2021, 09, 21]
         3      [2021, 09, 22]
         4      [2021, 09, 23]
         ...
        247     [2022, 09, 12]
        248     [2022, 09, 13]
        249     [2022, 09, 14]
        250     [2022, 09, 15]
        251     [2022, 09, 16]
Name: date, Length: 252, dtype: object
```

what is stringmethod?

```
In [94]: df_infy_aj_close["month"] = df_infy_aj_close.date.str.split("-").str[1]
```

```
C:\Users\mause\AppData\Local\Temp\ipykernel_9784\660655713.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead
```

```
See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/
user_guide/indexing.html#returning-a-view-versus-a-copy
df_infy_aj_close["month"] = df_infy_aj_close.date.str.split("-").str[1]
```

```
In [95]: df_infy_aj_close
```

```
Out[95]:   Adj Close      date  month
0    22.401333 2021-09-17      09
1    22.313101 2021-09-20      09
2    22.666033 2021-09-21      09
3    22.813086 2021-09-22      09
4    23.264053 2021-09-23      09
...
247  19.240000 2022-09-12      09
248  18.559999 2022-09-13      09
249  18.080000 2022-09-14      09
250  17.790001 2022-09-15      09
251  17.379999 2022-09-16      09
```

252 rows × 3 columns

```
In [96]: df_infy_aj_close.groupby("month")["Adj Close"].mean() # this is a series
```

```
Out[96]: month
01    23.894947
02    22.542275
03    23.791064
04    21.559053
05    19.106501
06    18.640476
07    18.665000
08    19.644348
09    20.294808
10    22.369523
11    22.804656
12    23.589092
Name: Adj Close, dtype: float64
```

```
In [97]: df_infy_aj_close.groupby("month")["Adj Close"].mean().reset_index() # this is a dat
```

```
Out[97]:   month  Adj Close
0      01  23.894947
1      02  22.542275
2      03  23.791064
3      04  21.559053
4      05  19.106501
5      06  18.640476
6      07  18.665000
7      08  19.644348
8      09  20.294808
9      10  22.369523
10     11  22.804656
11     12  23.589092
```

to get sort by month and year

```
In [98]: df_infy_aj_close["year"] = df_infy_aj_close.date.str.split("-").str[0]
```

C:\Users\mause\AppData\Local\Temp\ipykernel_9784\2795584799.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df_infy_aj_close["year"] = df_infy_aj_close.date.str.split("-").str[0]

```
In [99]: df_infy_aj_close
```

Out[99]:

	Adj Close	date	month	year
0	22.401333	2021-09-17	09	2021
1	22.313101	2021-09-20	09	2021
2	22.666033	2021-09-21	09	2021
3	22.813086	2021-09-22	09	2021
4	23.264053	2021-09-23	09	2021
...
247	19.240000	2022-09-12	09	2022
248	18.559999	2022-09-13	09	2022
249	18.080000	2022-09-14	09	2022
250	17.790001	2022-09-15	09	2022
251	17.379999	2022-09-16	09	2022

252 rows × 4 columns

In [105...]

```
def extract_month_year(row):
    row["year"] = row["date"].split("-")[0]
    row["month"] = row["date"].split("-")[1]
    return row
```

In []: