

School of Engineering & Technology

Affiliated to : University of Mumbai, Recognised by : DTE (Maharashtra) & Approved by : AICTE (New Delhi)

Name: Farhan Shikalgar	Course Name: SMA Lab
Class: BECO	Roll no: 19CO56

Experiment No: 3

Aim: Data Cleaning and Storage- Preprocess, filter and store social media data for business (Using Python, MongoDB, R, etc).

Theory:

Social Media Scraping of Myntra using Python:

1. <u>Instagram:</u>

```
2. # Instagram
3. from instagramy import InstagramUser
4. user = InstagramUser("Myntra")
5. print(f"Username: {user.fullname}")
6. print(f"Biography: {user.biography}")
7. print(f"Verified User: {user.is_verified}")
8. print(f"Website: {user.website}")
9. print(f"Followers: {user.number of followers}")
10.print(f"Following: {user.number_of_followings}")
11.print(f"No. Of Posts: {user.number of posts}")
12.posts = user.posts
13.print(posts[0])
14.instaPosts = []
15. for i in range(10):
16.
       post = \{\}
17.
       post["Likes"] = posts[i].likes
18.
       post["Comments"] = posts[i].comments
19.
       post["post_source"] = posts[i].post_source
       post["post_url"] = posts[i].post_url
20.
21.
       post["time"] = posts[i].taken_at_timestamp
       instaPosts.append(post)
23.insta_df = pd.DataFrame(instaPosts)
24.insta df.info()
```

```
25.print(insta_df.isna())
26.print(f"Latest posts:\n {insta df}")
```

OUTPUT

```
python Exp3-SMA.py
 Username: MYNTRA
Biography: If you never try, you'll never know ♥ 

#MyMyntraLook get featured on Myntra Studio
 Visit Myntra Studio for more fashion & style advice.
 Verified User: True
 Website: https://myntra.onelink.me/eWcy/gnozbjm9
 Followers: 3168336
 Following: 262
 No. Of Posts: 10804
No. Of Posts: 10804
Post(likes=283, comments=9, caption=None, is_video=True, timestamp=1675679094, location={'id': '109524955741121', 'has_pu blic_page': True, 'name': 'India', 'slug': 'india'}, shortcode='CoUUcXvJzT6', post_url='https://www.instagram.com/p/CoUUc XvJzT6', display_url='https://instagram.fbom56-1.fna.fbcdn.net/v/t51.2885-15/329543412_1342790036515421_8823910709516276
393_n.jpg?stp=dst-jpg_e15&_nc_ht=instagram.fbom56-1.fna.fbcdn.net&_nc_cat=1&_nc_ohc=sEZNfq2LUdkAX_mCWQ_&edm=ABfd0MgBAAAA&ccb=7-5&oh=00_AfAurul.bFXhitwlnp0VLfbppHaUgxzq70FfNuv1o9Up58w&oe=63E2D626&_nc_sid=7bff83', video_url='https://instagram.fbom56-1.fna.fbcdn.net/v/t66.30100-16/54642810_578501280529533_5147923975634763362_n.mp4?_nc_ht=instagram.fbom56-1.fna.fbcdn
.net&_nc_cat=109&_nc_ohc=rkgHinplNFIAX_mxMkR&edm=ABfd0MgBAAAA&ccb=7-5&oh=00_AfBJTQ0v6StFY_MRmGFFz3XpgEmJhP9rYSz38cukYINb-w&oe=63E274A7&_nc_sid=7bff83', video_view_count=4455, post_source='https://instagram.fbom56-1.fna.fbcdn.net/v/t66.30100-
16/54642810_578501280529533_5147923975634763362_n.mp4?_nc_ht=instagram.fbom56-1.fna.fbcdn.net&_nc_cat=109&_nc_ohc=rkgHinplNFIAX_mxMkR&edm=ABfd0MgBAAAA&ccb=7-5&oh=00_AfBJTQ0v6StFY_MRmGFFz3XpgEmJhP9rYSz38cukYINb-w&oe=63E274A7&_nc_sid=7bff83', taken_at_timestamp=datetime.datetime(2023, 2, 6, 15, 54, 54))
<class 'pandas.core.frame.DataFrame'></tl>
 <class 'pandas.core.frame.DataFrame'>
 RangeIndex: 10 entries, 0 to 9
 Data columns (total 5 columns):
     #
                  Column
                                                             Non-Null Count
                                                                                                               Dtype
     0
1
                  Likes
                                                             10 non-null
                                                                                                                int64
                                                                                                                int64
                                                             10 non-null
                  Comments
```

```
time
                                  10 non-null
                                                                  datetime64[ns]
dtypes: datetime64[ns](1), int64(2), object(2)
memory usage: 384.0+ bytes
      Likes Comments post_source
                                                               post_url
                                                                                    time
      False
                         False
                                                  False
                                                                      False
                                                                                  False
     False
                         False
                                                   False
                                                                      False
                                                                                  False
                         False
                                                  False
      False
                                                                      False
                                                                                   False
                                                  False
                                                                      False
      False
                         False
                                                                                  False
                         False
                                                  False
      False
                                                                      False
                                                                                   False
      False
                         False
                                                  False
                                                                      False
                                                                                  False
      False
                         False
                                                  False
                                                                      False
                                                                                   False
      False
                         False
                                                  False
                                                                      False
                                                                                  False
      False
                         False
                                                  False
                                                                      False
                                                                                   False
                         False
                                                  False
      False
                                                                      False False
Latest posts:
                                                                                                                post_url
        Likes Comments
                                               post_url time
https://www.instagram.com/p/CoUUcXv2T6/ 2023-02-06 15:54:54
https://www.instagram.com/p/CoTuff9D-xu/ 2023-02-06 10:27:51
https://www.instagram.com/p/CoURNTXJ9zQ/ 2023-02-06 15:29:32
https://www.instagram.com/p/CoUWRXxJGzV/ 2023-02-06 16:48:33
https://www.instagram.com/p/CoUWRXxJGx/ 2023-02-06 16:20:11
https://www.instagram.com/p/CoUWRPxs5la/ 2023-02-06 16:01:13
https://www.instagram.com/p/CoSwXmrjct9/ 2023-02-06 01:23:10
https://www.instagram.com/p/CoSEvupvhTz/ 2023-02-05 18:59:02
https://www.instagram.com/p/CoR3iSesDlg/ 2023-02-05 17:03:36
https://www.instagram.com/p/CoRj7jrpNPI/ 2023-02-05 14:12:17
          283
                               15
          790
1
2
3
4
5
6
7
          327
                               10
          505
                               30
          652
        2156
                               12
                                 9
          876
          796
                               21
8
                               55
        1320
        1092
[10 rows x 5 columns]
```

27. Facebook:

OUTPUT

```
False
                                                                     False
                                                                                                False
                                                                                                                                                  False
32
                  False
                                      False
                                                                     False
                                                                                                False
                                                                                                                              False
                                                                                                                                                  False
                                                                                                                                                                       False
33
                  False
                                      False
                                                                     False
                                                                                               False
                                                                                                                             False
                                                                                                                                                 False
                                                                                                                                                                       False
34
                  False
                                      False
                                                                     False
                                                                                               False
                                                                                                                              False
                                                                                                                                                 False
                                                                                                                                                                       False
35
                  False
                                     False
                                                                     False
                                                                                               False
                                                                                                                              False
                                                                                                                                                 False
                                                                                                                                                                       False
Latest posts:
             post_id
10161169058823221
                                                                                                                                                                                                              text ...
                                                                                                                                                                                                                                              likes
                                                                    26.0 2023-02-06 16:20:21
10.0 2023-02-06 15:55:05
9.0 2023-02-06 15:29:42
11.0 2023-02-06 12:31:15
14.0 2023-02-06 11:15:04
            10161169030738221
10161168987678221
             10161168767013221
             10161168690158221
                                                                                                                                                                                                                                            . 14.0 2023-02-06 11:15:04
29.0 2023-02-06 01:23:22
82.0 2023-02-05 14:12:22
28.0 2023-02-04 23:38:35
29.0 2023-02-04 19:46:55
70.0 2023-02-04 18:43
24.0 2023-02-04 13:36:57
20.0 2023-02-04 11:55:11
27.0 2023-02-04 11:07:25
34.0 2023-02-03 23:55:34
29.0 2023-02-03 23:59:64
             10161167893718221
             10161166649103221
             10161165440628221
             10161162761178221
                                                                     Jump into an unparalleled listening experience...
            10161164921723221
10161164580113221
                                                                     Shhhh...something extraordinary coming your way, ...
10
                                                                     Shhhh...something extraordinary coming your way, ...
             10161164466943221
11
                                                                     Shhhh...something extraordinary coming your way, ...
             10161164422793221
                                                                     Shhhh...something extraordinary coming your way, ...
12
13
             10161163521578221
                                                                     Denim game strong with pantaloons!\nCheck out
                                                                                                                                                                                                                                               29.0 2023-02-03 23:09:04
14
            10161163439843221
                                                                    A race before the actual race!\nThe Early Bird...
                                                                   Keep swiping ...

Wibe check \( \cap{\n}\) \( \text{ophoolfilledvish} \cap{\n}\) \( \text{eep swiping } \( \text{ophoolfilledvish} \cap{\n}\) \( \text{ophoolfilledvish} \) \( \
                                                                                                                                                                                                                                              121.0 2023-02-03 21:12:09
62.0 2023-02-03 00 :12:56
19.0 2023-02-02 23 :41:28
45.0 2023-02-02 22:38:48
60.0 2023-02-02 20:30:45
15
             10161163209568221
16
             10161161257373221
17
             10161161191548221
             10161161063803221
18
             10161160787473221
                                                                                                                                                                                                                                                        NaN 2023-02-02 00:01:45
20
             10161158857943221
             10161158772288221
                                                                                                                                                                                                                                                 81.0 2023-02-01 23:23:33
```

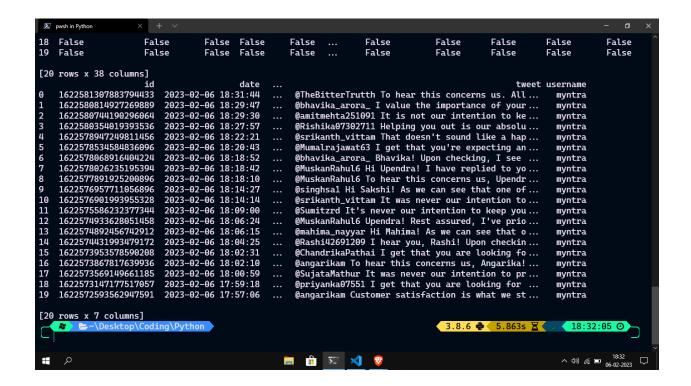
39. Twitter:

```
40.# Twitter
41.import twint
42.c = twint.Config()
43.c.Lang = "en"
44.c.Username = "Myntra"
45.c.Pandas = True
46.# Run
47.twint.run.Search(c)
48.Tweets_df = twint.storage.panda.Tweets_df
49.Tweets_df.info()
50.print(Tweets_df.isna())
51.Tweets_df = Tweets_df[['id', 'date', 'nlikes', 'language', 'nreplies', 'tweet', 'username']]
53.print(Tweets_df)
```

Output

```
1622578947249811456 2023-02-06 18:22:21 +0530 <myntra> @srikanth_vittam That doesn't sound like a happy situation, Vittam I I've prioritized this concern &amp; our case manager will contact you at the earliest. -KB https://t.co/WsLuXFYSSb 1622578534584836096 2023-02-06 18:20:43 +0530 <myntra> @Mumalrajawat63 I get that you're expecting an update, and it was never our objective to keep you waiting, Mumal! Upon checking, I see that an email was sent to you under reference number 'IN230203165210822064U53' on 06-Feb-2023 at 11:46 AM. Kindly check (cont) https://t.co/CbJqJhBR20 1622578068916404224 2023-02-06 18:18:52 +0530 <myntra> @bhavika_arora_Bhavika! Upon checking, I see that one of our case managers tried calling you today, 12:40 PM, but unfortunately, could not connect. A notification was sent to you via the Myntra App with the reference number IN23020515312830651385. You (cont) https://t.co/Tjl3hTmwJq 162257806235105394 2023-02-06 18:18:42 +0530 <myntra> @MuskanRahul6 Hi Upendra! I have replied to you here https://t.co/qRsSJRKsim. Please check. - ST 16225778091925200896 2023-02-06 18:18:10 +0530 <myntra> @MuskanRahul6 To hear this concerns us, Upendra! I regret the hass ce caused. Rest assured, we're closely following up with the logistics partners to have the product picked up from you so on. One of our case managers will connect with you to share an update as early as possible. - ST 1622576957711056896 2023-02-06 18:14:27 +0530 <myntra> @singhsal Hi Sakshi! As we can see that one of our case managers ontacted you on 01 Feb 2023, 01:15 PM and ask for images of the product with bit more information. Also, send an in-app w ith reference dI IN23013111020216736311. kindly have a check and revert for further assistance. -AB 1622576901993955328 2023-02-06 18:10:14 Hb 2530 <myntra> @srikanth_vittam It was never our intention to cause any disruption to our customers, Srikanth! Please be assured that I've highlighted your concern, and we are looking into it with the u tmost priority. One of our case managers will co
```

≥ p	wsh in Python X	+ ~		-	ō
Data	columns (total 3	8 columns):			
#	Column	Non-Null Count	Dtype		
					
Θ	id	20 non-null	object		
1	conversation_id	20 non-null	object		
2	created_at	20 non-null	float64		
3	date	20 non-null	object		
4	timezone	20 non-null	object		
5	place	20 non-null	object		
6	tweet	20 non-null	object		
7	language	20 non-null	object		
8	hashtags	20 non-null	object		
9	cashtags	20 non-null	object		
10	user_id	20 non-null	int64		
11	user_id_str	20 non-null	object		
12	username	20 non-null	object		
13	name	20 non-null	object		
14	day	20 non-null	int64		
15	hour	20 non-null	object		
16	link	20 non-null	object		
17	urls	20 non-null	object		
18	photos	20 non-null	object		
19	video	20 non-null	int64		
20	thumbnail	20 non-null	object		
21	retweet	20 non-null	bool		
22	nlikes	20 non-null	int64		
23	nreplies	20 non-null	int64		
24	nretweets	20 non-null	int64		
25	quote_url	20 non-null	object		



Conclusion:

We have successfully cleaned Myntra's social media data from websites like Instagram, Facebook, and Twitter using various data preprocessing techniques.