

Project: Other Analytics

Course: Massive Data Analytics

Professor: Dr. Iman Gholampoor

Mohammad Mehdi Zare

- **Introduction:**

In this report, I aim to address four questions utilizing the available data, including Twitter data and Persian news data, with the help of algorithms covered in the course. In the provided notebook, various libraries have been used, but the algorithms have also been implemented manually. Overall, the PySpark library and RDD API have been employed for data processing.

- **The questions that I will be answering next:**

1. If Twitter decides to verify some users for free, which users are candidates for verifying, and which users have lost their verified label?
2. How many distinct words are present in news data? (batch processing and stream processing approach)
3. Each user, based on their tweets and behavior on Twitter, is interested in a specific set of categories. News also belongs to various categories. If Twitter intends to suggest daily news to its users, which news does it recommend to each user?
4. Some news agencies copy news from others, which is a form of Plagiarism. I want to identify these instances in the news data.

- **Sections:**

1. **Dataset preparation and preprocessing:**

- a) At first, I clean news data. I remove punctuation and stopwords. Then, remove unnecessary columns from the data and keep columns: uid, body, source, date\_published, key\_words, and categories. The result dataset is:

```
[8]: [{"68fee4bbdc2d54adb2369",
      "varzesh3.com",
      "10/10/2023, 07:16:46",
      "سررمی ملوان وایرالی دیدار نسامی پست چاب منتشر گزارش ورزش مهدی تارنار خوشحالی عجیب غریب نیکه چهره چاب هله لیگ برتر اختصاص موزیتوچه هواداران فوتبال قرار تارنار انتشار پست اینستاگرامی تصویر چن خوشحالی نوشته جانوی طرهداران ازلی سررمی ملوان وده توشی ادمه فصل توان خواهد جنگی رضایت هواداران جلب نوشته ادمه میخواید عکسها بگه میکم متنی میوید پدم نوایم هرگز اینطور شوری بشم جانوی شملت سگو انجان حق انرژی میدید لحظه پیروزی برآیدمان وسفالتانی سگوها لب خط پای تویزیونیا مستطیل سبز هسل یکسدا میروید باهم میجگیم باهم میخندیم هرروز تانتش میکنیم خنده صورت دگه پیروزی شملت "ورزش سه", "varzesh3.com", "10/10/2023, 07:16:46", ["سررمی ملوان", "لیگ برتر", "مستطیل سبز", "ازلی", "sports"]], ("434e9837dd761812773fa4627", "سندگوی کیسبون تلقی برنامه هتم موافقت ضمنی دولت اجرای همنسازای برنامه خبر گزارش ایلا محسن زبگه سندگوی کیسبون تلقی برنامه هتم توسعه حشیه علی امروز بهشینه مهر ماه مجلس خبرندگان اشا زده وصیعت همنسازای حقوق بازنشگان شب گاشته جلس ساعته رئیس مجلس شورای اسلامی رئیس کیسبون تلقی برنامه هتم توسعه رئیس سازمان برنامه بودجه برگزار نت جمعیتی اولیه برنامه هتم توسعه بتوایم بحث متناد ب حقوق بازنشگان برانجام برسانیم هدف پایان برنامه هتم شاد تلقی موضوع بتایم سندگوی کیسبون تلقی برنامه هتم توسعه عنوان موضوع اصلاحات ساختاری صندوق بازنشگانی انجام مشکل ورنشگانی صندوق بازنشگانی "ای مرتفع", "خبرگزاری ایلا", "ilna.ir", "10/10/2023, 07:17:05"}]
```

- b) Second, I clean the tweet dataset and remove incomplete rows in the data. I remove 214 rows:

```
0]: # In data, some tweets have incomplete data
incomp=tweets_json_rdd.filter(lambda x: x['tweet_type']=='quoted' and x.get('quoted_status').get('id')==0).collect()
print("sample of incomplete data: ")
print(incomp[1])
print("number of incomplete data : "+str(len(incomp)))

# filter them
tweets_json_rdd_cleaned=tweets_json_rdd.filter(lambda x: not (x['tweet_type']=='quoted' and x.get('quoted_status').get('user')==0))
tweets_json_rdd_cleaned=tweets_json_rdd_cleaned.filter(lambda x: not x.get('nlp')==0)

sample of incomplete data:
{'publish_source': None, 'in_reply_to_status_id_str': None, 'in_reply_to_user_id_str': None, 'in_reply_to_status_id': None, 'in_reply_to_user_id': None, 'in_reply_to_screen_name': None, 'truncated': False, 'is_quote_status': True, 'retweet_count': 1, 'reply_count': 0, 'quote_count': 1, 'favorite_count': 3, 'favorited': False, 'retweeted': False, 'possibly_sensitive': None, 'lang': 'fa', 'geo': None, 'view_count': 235, 'tweet_type': 'quoted', 'emoji': No ne, 'text': 'با ازدواج نکند یا از مهر بالا فرار کند\انی که به مهر بالا اصرار دارد رنگه به چه زبونی بهتر بگه میخواد بعدا بنشینون کنه\اینو اقا اون مهریه کوئی رو تعداد بالا بنویسید: "id": "1363594775539159045", "id_str": "1363594775539159045", "created_at": "2023-09-16T11:36:16+03:30", "name": "دیکل ارمیا", "screen_name": "md_s_law", "location": "tehran", "description": "حسبنا الله و نعم الوکیل", "verified": False, "followers_count": 1124, "friends_count": 1431, "listed_count": 1, "favourites_count": 3955, "statuses_count": 1334, "profile_image_url": "https://pbs.twimg.com/profile_images/1729131836213800960/VbZfy7s-normal.jpg", "profile_banner_url": "https://pbs.twimg.com/profile_banners/1363594775539159045/1647984275", 'default_profile': True, 'default_profile_image': False, 'entities': {'urls': [], 'hashtags': [], 'symbols': [], 'user_mentions': []}, 'quoted_status': {}, 'engagement': 4, 'impression': 235, 'reach': 0.2091, 'influence_v2': 0.0036, 'impression_rate_v2': 0.1762, 'engagement_rate_v2': 0.003, 'twitter_engagement_rate_v2': 1.7021, 'timestamp': 1702393654, 'index_version': 1, 'created_at': 1698831343, 'id': "1719649410987028916", 'tweet_source_type': 'tweet_timeline_v2', 'text_lang': 'fa', 'nlp': {'offensive': 'non_offensive', 'sentiment': 'Negative', 'emotion': ['anticipation'], 'keywords': [{'keyword': 'مهریه'}, 'question': 'question', 'hashtag_polarity': [], 'event': 'NotEvent', 'classification': ['social', 'persona']}, 'ner_polarity': []}]
number of incomplete data : 214
```

## 2. Exploration

Let's go to answer the above questions:

- a) **If Twitter decides to verify some users for free, which users are candidates for verifying, and which users have lost their verified label?**

To answer this question, I use PageRank and TrustRank algorithms. users are linked to each other by their interaction with tweets. So with the PageRank algorithm, we can find a weight for the user. Also, verified users can be used as trust sets and then use the TrustRank algorithm. I use Networkx library to create a directed graph.

Steps:

- Find verified users as trust set. In data, some users have verified labels, so I filter them. You can see the number of users and verified users in the data, below picture:

```
[19]: # find source user
user_id_list = tweets_data['user_id'].tolist()
user_id_list = list(set(user_id_list))

[20]: # find target user
target_id_list = tweets_data['in_reply_to_user_id'].tolist()
target_id_list = list(set(target_id_list))

[21]: # find user who is verified
verified_users = []
for user_id in user_id_list:
    user = users[user_id]
    if user['verified']:
        verified_users.append(user_id)

[22]: # find user who is not verified
not_verified_users = []
for user_id in user_id_list:
    user = users[user_id]
    if not user['verified']:
        not_verified_users.append(user_id)

[23]: # count number of users and verified users
number_of_users = len(user_id_list)
number_of_verified_users = len(verified_users)

[24]: # print number of users and verified users
print("Number of users: ", number_of_users)
print("Number of verified users: ", number_of_verified_users)
```

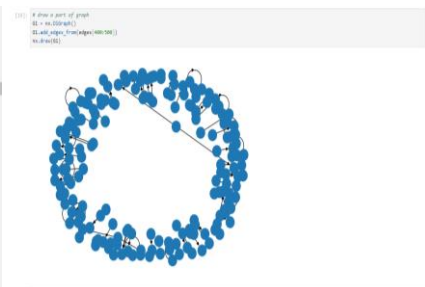
- Make a directed graph. Users interact with each other by replying, retweeting, and quoting. Therefore, I make edges from user interaction in data. And make a directed graph. In the below pictures you can see some edges and a part of the graph:

```
[25]: # create a directed graph
graph = nx.DiGraph()

# add nodes
for user_id in user_id_list:
    graph.add_node(user_id)

# add edges
for tweet in tweets_data:
    user_id = tweet['user_id']
    in_reply_to_user_id = tweet['in_reply_to_user_id']
    retweet_of_tweet_id = tweet['retweet_of_tweet_id']
    quote_tweet_id = tweet['quote_tweet_id']

    if in_reply_to_user_id:
        graph.add_edge(user_id, in_reply_to_user_id)
    if retweet_of_tweet_id:
        graph.add_edge(user_id, retweet_of_tweet_id)
    if quote_tweet_id:
        graph.add_edge(user_id, quote_tweet_id)
```



- Run page Rank algorithm. The result is:

```
[19]: [{"user_id": "38154763", "page_rank": 0.019669219532820137},
      {"user_id": "2682820352", "page_rank": 0.00505240408135082},
      {"user_id": "810859229267378176", "page_rank": 0.005047751465343615},
      {"user_id": "4075582205", "page_rank": 0.004502849214073731},
      {"user_id": "1440806264255045638", "page_rank": 0.003923641660580251},
      {"user_id": "867126568451280896", "page_rank": 0.003609715762625503},
      {"user_id": "9458863960310592", "page_rank": 0.002927343671582436},
      {"user_id": "1612794343911297025", "page_rank": 0.002476016227295738},
      {"user_id": "2739506010", "page_rank": 0.002474665074495339},
      {"user_id": "2331584635", "page_rank": 0.002310059354181786},
      {"user_id": "1408312021", "page_rank": 0.0022640175312704947},
      {"user_id": "1290359738585288704", "page_rank": 0.002155140438080757},
      {"user_id": "122797681901399681", "page_rank": 0.0019409359004997234},
      {"user_id": "49127938", "page_rank": 0.001904467540991717},
      {"user_id": "46482865", "page_rank": 0.0018732180562105003},
      {"user_id": "7227576277335040", "page_rank": 0.001818383265052587},
      {"user_id": "122277330244148416", "page_rank": 0.0017797144975011},
      {"user_id": "146270073786371801215", "page_rank": 0.001633263606166075}
```

- Run trust Rank algorithm. The result is:

```
[22]: [{"user_id": "38154763", "trust_rank": 0.06050678180503591},
      {"user_id": "49127938", "trust_rank": 0.0107245931254231},
      {"user_id": "2682820352", "trust_rank": 0.01015200715102240},
      {"user_id": "3414505683", "trust_rank": 0.00956065942950009},
      {"user_id": "900512912", "trust_rank": 0.008712373671555865},
      {"user_id": "122797681901399681", "trust_rank": 0.007860097748656215},
      {"user_id": "1440806264255045638", "trust_rank": 0.007367714883802768},
      {"user_id": "2739506010", "trust_rank": 0.006894261694057001},
      {"user_id": "742451244028039168", "trust_rank": 0.006762484959106495},
      {"user_id": "160940788", "trust_rank": 0.00673242143866629},
      {"user_id": "1458337727536144386", "trust_rank": 0.0060605019641884956},
      {"user_id": "4075582205", "trust_rank": 0.00576468497048296},
      {"user_id": "1408312021", "trust_rank": 0.00568109402041521},
      {"user_id": "1575005035501735936", "trust_rank": 0.005270689332861157},
      {"user_id": "9458863960310592", "trust_rank": 0.00517939983047827},
      {"user_id": "72296083385173760", "trust_rank": 0.0050401659163541185},
      {"user_id": "867126568451280896", "trust_rank": 0.005038288959481802},
      {"user_id": "146270073786371801215", "trust_rank": 0.00476033336276260}
```

- Verified users whose rank difference in the two algorithms is greater than .03 are fake users. And users that rank 1 to 100 are new verified users.

```
23]: dict_g1=dict(pagerank_sorted)
    dict_g2=dict(pagerank_sorted1)
    result_dict = {key: dict_g1[key] - dict_g2[key] for key in dict_g1}
    alpha = .03
    fake_user = [key for key, value in result_dict.items() if abs(value) > alpha]
    fake_user
```

```
24]: # fake verified user
      set(verified_nodes).intersection(fake_user)
```

```
25]: # user that can verify
first_100_user=list(zip(*pagerank_sorted1))[0][1:100]
new_verified_user=set(first_100_user).difference(verified_nodes)
new_verified_user
```

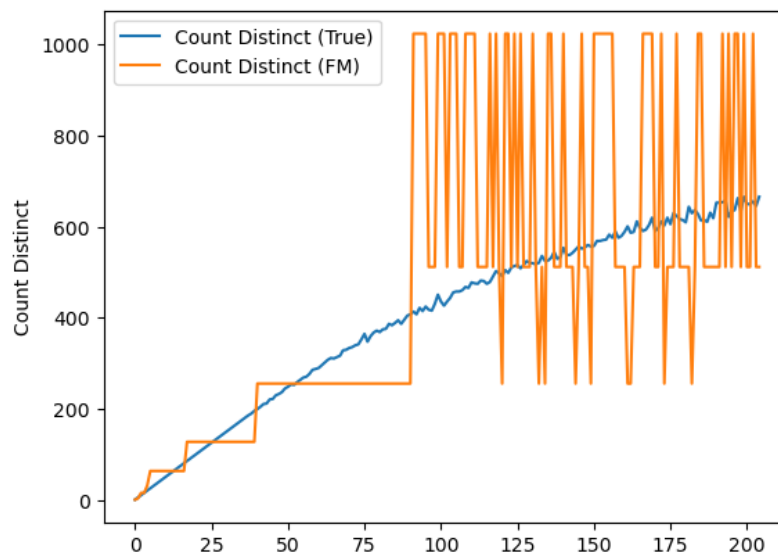
```
25]: {'1065658304884547585',
      '1222773302441148416',
      '1515956552',
      '1588740471546847238',
      '1612794343911297025',
      '2509360928',
      '3189754680',
      '46482865',
      '722757627877335040',
      '802011472964546560',
      '940962097130889216'}
```

To answer this question, I implement the Flajolet Martin algorithm. Steps:

- ```
52]: print(tokens[1:100])
```

[illegible]

- In each partition, I start with defining a closed hash range, big enough to hold the maximum number of unique values possible ( $2^{64}$ ). Every word of the tokens is passed through a hash function that permutes the elements in a uniform distribution. For this hash value, I find the position of the rightmost 1s bit and mark the corresponding position in the bit vector as 1. Once all the words of the partition are processed, the bit vector will have 1s at all the positions corresponding to the position of every rightmost 1s bit for all words in the partition. Now I find the position,  $b$ , of the leftmost 1 in this bit vector. This position  $b$  corresponds to the max length of unset bits that I have seen while processing the words. Then find the median of the estimated value for each hash function. Finally, estimate distinct words by averaging values of partitions. In the below picture, you can see the fluctuation of estimated values for generated data.



After running the algorithm in news data, the result is:

```
0]: tokens=news_cleaned1.flatMap(lambda x : x[1].split()).collect()
part_tokens=np.array_split(tokens, 3)
res=[]
for i in part_tokens:
    tok=list(i)
    tokens_hashed = [int(hashlib.md5(token.encode()).hexdigest(), 16) % (2 ** 64) for token in tok]
    res.append(count_distinct_fm(tokens_hashed,params))

print(len(tokens), len(set(tokens)))
print("estimated words: ")
print(statistics.mean(res))

1017002_47456 True distinct words
estimated words:
49152
```

We can see that the estimated count-distinct using the Flajolet-Martin algorithm is very close to the actual deterministic words.

- c) **Each user, based on their tweets and behavior on Twitter, is interested in a specific set of categories. News also belongs to various categories. If Twitter intends to suggest daily news to its users, which news does it recommend to each user?**

In this part, I use stream processing. users like some topics. news has topics too. so we can use the bloom filter to recommend news to users each day. Steps:

- Make data as a stream data.
- Implement bloom filter
- Make categories in news data and tweet data the same. Twitter and Wikipedia use different words for categories. For example, economy and economical. You can see the categories in tweets and news:

```
1]: # list of categories in new and tweets
categories=news_cleaned1.flatMap(lambda x : x[6]).distinct().collect()
classifications=tweets_rdd.ok().flatMap(lambda x : x['nlp']).get('classification').distinct().collect()

2]: print(categories)
print(classifications)

['politics', 'social', 'economy', 'security', 'culture', 'sports', 'science_and_technology', 'religious', 'health', 'military']
['entertainment', 'security', 'military', 'health', 'offensive', 'culture', 'religious', 'others', 'politics', 'science_and_technology', 'sports', 'economy', 'personal', 'advertising', 'social']
```

- I use categories that exist in the intersection set.

```
8]: print(set(categories).intersection(classifications))
print(set(categories).difference(classifications))

{'politics', 'security', 'social', 'health', 'economy', 'sports', 'culture', 'military', 'religious', 'science_and_technology'}
set()
```

- As we know, The bloom filter has false positive. I calculate the number of hash functions and the size of the bit array. According to the number of categories and False positive Probability.

```
[40]: recom_newsUser["1004052664882118657","list_news_day,tweets_rdd_ok"]

Size of bit array:19
False positive Probability:0.0001
Number of hash functions:13
category for day 10/09/2023 : ['politics']
In day :10/09/2023
user 1004052664882118657 maybe like news :
{"0184a2f6b9d8f482d03247d9f91f": "آنجایی واتس‌گرم تلاش میکند مانع افزایش قیمت نفت در جهان پالایشگاه چین مطمئن صادرات نفت ایران ادامه خواهد گرازش گروه اقتصاد بینالملل فارس نل بلومبرگ آنجیکه واتس‌گرم و واتس‌گرم
نفت میکند مانع افزایش قیمت نفت در جهان پالایشگاه چین مطمئن صادرات نفت ایران ادامه خواهد گرازش پالایشگاه مستلزم چن موسوم نتایجها مشتق خرید نفت ایران بلومبرگ مدعی افزایش ایران مدخلان افزایش ادعای حمله فکلفکاره میوات
تفاوت نفت کشور موضع سبتر آمریکا لک ترمیر ایران معرض تهدید راف حاضر چنین پالایشگاه، چنین نفت خام ایران میبردند هیچگونه اختلال مسدودگری پیشبینی نمی‌نمودند نتایجها مشتق نفت ایران بلومبرگ آنجیکه واتس‌گرم و واتس‌گرم
دولت ایران فشار تحریم آمریکا از آید بررسی خواهد گشته تحریمها بدون اشتیاق چن خرید نفت خام مدعی مزاری تأمین ایران شرکت می‌نمودند نتیجه این تحریمها
ت ایران چن ماه گرویزت دولته بلومبرگ می‌نمودند پیشه روز رسیده پیشترین رقم گشته میروند با ادعای شبکه خبری تهدیدات درحالیکه رسانه آمریکایی ادعا میکند ایران حملات موقعت نفت ائتلافی بلکان وزیر خارجه آمریکا شبکه خبر
گروه الفمل فارس." وزیر خارجه آمریکا" و"التوی بلکان" و"التشعات برای مجه" "10/09/2023", "farsnews.ir", "ی ای نیس" جانسی شواهدی این ائتلاف مستقیم دمات پایان پیام" "خبرگزاری فارس"
category : ['politics']]
category for day 10/10/2023 : ['security']
In day :10/10/2023
user 1004052664882118657 maybe like news :
{"a7e44da4f2f282902cebed90bf": "گزارش ایفا امر سربط طریحه‌ها الهامی جاشتن فرمانده نیروی پادافه هواری ارتش دومین گرازش فراز شهید اسانم مدخلان اظهاراتات خلواده شهبا تقدیم ارتش کشور مردم امتیاز"
رامش جی‌تیست عرب نیس انتفاع انگیزه حرکت شهبا نل امروزی وایولارها کنه‌گرها جزو ادعای اصلی قرار فرموده درحضر معتمد انقلاب منتظله‌هایی از ثواب اقامات امر شهبا کنتر اثره گشتان دفاع قضای امتیاز اصل سیمیتت ۲
```

- Some Twitter users do plagiarism, meaning they send out tweets that are essentially identical to another user's tweet, with slight modifications in sentence structure. Due to constraints in processing tweet text, I implement the algorithm on news data, leveraging the results of Exercise 1 in this section. Additionally, the first news agency to publish the text is considered trustworthy. note: I implemented tf-idf in hw1. To save time I used the Mllib library in the project.

I chose news from hw2 results that I have found similar news for that. After calculating tf-idf for all news and candidate news, calculate the cosine similarity between the candidate and other news and find the top five news that has high cosine similarity. According to their published date, we can say which ones are plagiarism.

```
for t in topFive:
    print("doc '%s' has score %.4f" % (t[0], t[1]))
```

```
doc 'a58188b693d8167ecf144685a' has score 1.0000
doc '6789ed5448fe3bcb8e9bf4bcb' has score 0.9947
doc '49a301cb98246ca6d89e33755' has score 0.9947
doc 'c672d416dfda66a9b979e1202' has score 0.9913
doc '35a62bfc87becf6a2cfa6357' has score 0.9909
```

News are:

```
[47]: x=yzip(*topFive)
y=list(x)
top_news=news_cleaned1.filter(lambda x: x[0] in y)
top_news.take(10)
```

[47]: [['a58188b693d8167ecf144685a',  
مراسم افتتاح ایستگاه مترو خطوط آذال صلیبات حطاری مکنیزه تولل خط حضور رئیسجمهوری انعام گزارشعزوه، دولت خبرگزاری تسنیم سید ابراهیم رئیسی رئیسجمهور صبح امروز سبته حضور ایستگاه شهید آرمان طویوری منط'  
فه گوهرام افتتاح رسمی ایستگاه خط مترو تهران قطار ایستگاه شهران ایستگاه شگل رسمی افتتاح مراسم افتتاح ایستگاه شهید طویوری حضور خانواده برگزار رئیس جمهور پیشهاد تابلویی زنگینامه شهید وایفام ووردی ایستگاه نصب  
همزمان مراسم ایستگاه شهیر زیبا خط میان کتاب خط متروی تهران بهرهبرداری رسیدند مراسم امروز صلیبات حطاری مکنیزه تولل خط متروی تهران کارگاه دریاچه شهبادی حطاری فارس آذال خط شرق پایتخت ادامه مییاد خط کیومرث  
'، ایستگاه جاجایی مسافران منتقل پوشتن خواهد  
'، اخبارگزاری تسنیم'  
'www.tasnimnews.com',  
'10/10/2023, 07:33:10',  
'کارگاه دریاچه شهبادی حطاری فارس',  
'ایستگاه شهید آرمان علی وردی',  
'ایستگاه شهید علی وردی',  
'ایستگاه شهران',  
'politics']],  
( '35a62bf8c7becf6a2cfca6357',  
گزارش خبرگزاری مهر تحت الاسلام سید ابراهیم رئیسی رئیسجمهور صبح امروز سبته سید حضور ایستگاه شهید آرمان طویوری منطقه کوسمار افتتاح رسمی ایستگاه خط مترو تهران قطار ایستگاه شهران ایستگاه شگل رسمی'  
افتتاح گزارش مراسم ایستگاه شهید طویوری حضور خانواده برگزار رئیس جمهور پیشهاد تابلویی زنگینامه شهید والاعام ووردی ایستگاه نصب ابراهیم رئیسی همزمان مراسم ایستگاه شهیر زیبا خط میان کتاب خط متروی تهران بهرهرد  
'، ارئ رسیدند مراسم امروز صلیبات حطاری مکنیزه تولل خط متروی تهران کارگاه دریاچه شهبادی حطاری فارس آذال خط شرق پایتخت ادامه مییاد گفتنی خط کیومرث ایستگاه جاجایی مسافران منتقل پوشتن خواهد  
'، Meh News Agency' | خبرگزاری مهر | اخبار ایران و جهان'  
'mehnews.com',  
'10/10/2023, 07:46:15',  
'کارگاه دریاچه شهبادی حطاری فارس',  
'ایستگاه شهید آرمان علی وردی',  
'ایستگاه شهید علی وردی',  
'مهر',  
'social']],  
( '6789ed544fe3bc8e9bf4bcb',  
مراسم افتتاح ایستگاه خطوط مترو حطاری مکنیزه تولل خط متروی تهران حضور رئیس جمهور آذال گزارشی هشتورئی آنتنن حضور سید ابراهیم رئیسی ایستگاه شهید آرمان طویوری ایستگاه مراسم افتتاح ایستگاه مترو آر'  
سما مراسم خانواده شهید علی وردی حضور رئیس جمهور پیشهاد تابلویی زنگینامه شهید ووردی ایستگاه نصب رئیس جمهور قطار مترو ایستگاه ایستگاه شهران همزمان مراسم ایستگاه شهیر زیبا خط میان کتاب خط متروی تهران بره  
بهرداری رسیدند امروز همزمان مراسم افتتاح ایستگاه مترو صلیبات حطاری مکنیزه تولل خط متروی تهران کارگاه دریاچه شهبادی حطاری فارس آذال خط شرق پایتخت ادامه مییاد خط کیومرث ایستگاه جاجایی مسافران منتقل پوشتن خواهد  
'،  
'هشتورئی آنتنن',  
'hamshahrionline.ir',  
'10/10/2023, 07:48:02',  
'کارگاه دریاچه شهبادی حطاری فارس',  
'آرمان علی وردی',  
'ابراهیم رئیسی',  
'ایستگاه شهران',  
'social']],

you can see that all the top five news are the same and about specific topics and published on specific dates.

```
top_news.map(lambda x : (x[0],x[4])).take(10)
```

Therefore, news with Id a58188b693d8167ecf144685a is a source, and others may copy from that.

The news agency of this news is :

```
59: top_news.filter(lambda x: x[0]=='a58188b693d8167ecf144685a').map(lambda x:(x[2],x[3])).collect()

59: [('خبرگزاری تسنیم', 'www.tasnimnews.com')]
```

### Conclusion:

In the surrounding datasets, there is a wealth of information that can be uncovered through various data analysis methods. We discussed and explored different approaches for processing stream and persistent data, each having its own advantages. In the project, we demonstrated that stream processing methods may have errors, but they still provide good estimates of information, enabling us to make strong assumptions backed by the science of statistics and probability regarding data behavior.

We also emphasized the importance of data and how it can enhance the user experience in applications. It can assist us in identifying fraudulent activities and plagiarism by individuals from other sources. Moreover, based on users' behavioral history, we can even suggest news or tweets that align with their preferences and interests.