



Sharif University of Technology

Masoud Tahmasbi Fard



Student ID: 402200275

EE577: Massive Data Analytics

Assignment #3- Twitter User Engagement Recommender System

February 4, 2024

Table of Contents

1	Task Description	1
2	Algorithm Description	1
2.1	Theoretical Model	1
2.2	Implementation	2
3	Results	3
3.1	User #1	3
3.2	User #2	6
3.3	User #3	8
3.4	User #4	11
3.5	User #5	13
4	Future works	16
5	Conclusion	17

1 Task Description

In this assignment, we implement a Twitter User Engagement Recommender System using a dataset of Persian users' tweets. Each tweet contains multiple elements describing the tweet, including the text, user id, generator uid and tweet id (for "generated"-type tweets), quoted uid and tweet id (for "quoted"-type tweets), replied uid and tweet id (for "replied"-type tweets), retweeting uid and tweet id (for "retweeted"-type tweets), and more. Leveraging these fields, we develop a recommender system in PySpark that can recommend potentially similar users and, based on that, potentially similar tweets to engage users. The implementations utilize PySpark to analyze the large-scale tweet dataset and generate personalized user and tweet recommendations. Specifically, we extract user similarity features and feed them into a collaborative filtering algorithm to produce user-based recommendations. The goal is to recommend accounts and content that are relevant and engaging to Twitter users.

2 Algorithm Description

2.1 Theoretical Model

The recommender system is implemented using a user-based collaborative filtering algorithm, which makes recommendations by finding users with similar interests and recommending tweets related with those similar users.

To find similar users, we utilize Pearson correlation (1), a technique that measures the linear relationship between two users' tweet engagement patterns. It returns a value between -1 and 1 indicating how strongly two users' engagement tendencies align. A correlation of 1 means perfect alignment where users engage with precisely the same tweets, 0 is no correlation, and -1 is perfect misalignment.

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x) (r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}} \quad (1)$$

The algorithm first calculates a Pearson similarity score between the tweet engagement profile of

a given user and all other users. It selects other most similar users. Then from those similar users, the algorithm examines the tweets they engaged with that the given user has not engaged with. The highest scoring of these tweets, based on the aggregate of the similar users' engagements, are recommended to the given user.

2.2 Implementation

To implement the recommender system, a tweet scoring scheme is designed to quantify a user's engagement with each tweet. The scheme assigns points as follows:

- 3 points if the user generated the original tweet
- 2 points if the user quoted or replied to the tweet
- 1 point if the user retweeted the tweet

A scoring function encapsulates this logic, taking in default parameters for generated, quoted, reply, and retweet scores that can be customized as needed.

The implementation utilizes a map-reduce framework over the user RDD. Five target users are chosen from different communities and viewpoints. For each target user, the Pearson Correlation is calculated between them and all other users to measure tweet engagement profile similarity. Specifically, the correlation examines the subset of tweets rated by both user x and user y , comparing their assigned engagement scores r_{xs} and r_{ys} . In our Pearson correlation calculation, we omit subtracting the average rating values for r_x and r_y for two reasons. First, averaging is typically done to account for user rating biases; however, in our domain user averages do not represent systematic biases. Second, some users have no rating variance (e.g. a user only retweets), causing their average to equal every rating and produce undefined correlation values. By removing averages, we avoid invalid divisions while not needing to correct biases. The raw rating vectors still capture relative tweet engagement differences between users necessary for measuring similarity. Omitting averages retains computational stability without affecting the

utility of the resulting correlations for modeling tweet engagement tendencies. As a consequence, the calculated similarity is not necessarily between -1 and 1.

The most similar users to the target are selected. The tweets those similar users engaged with that the target user did not are extracted. The tweets are scored according to:

$$\text{tweet score}(t) = \text{sim}(x, y) * r_{yt} \quad (2)$$

Where r_{yt} is similar user y 's rating of tweet t . The top scoring tweets are recommended to the target user. This personalized recommendation pipeline allows suggesting potentially relevant tweets tailored to users based on those with historically analogous engagement tendencies.

3 Results

In this section the target user's tweets, similar users id's, and recommended tweets are shown for 5 different type of users.

3.1 User #1

```
target_user = grouped_rdd.takeSample(False,1,55)[0]
target_user

('1697206791715635200',
 [('1710302619371982848', '1736250426230493213', 1),
  ('1710302619371982848', '1736250178363883872', 1),
  ('1710302619371982848', '1736249928018764235', 1),
  ('1578454116107517962', '1735951927584682444', 2),
  ('1633911157755879425', '1736108738921033809', 1),
  ('1654409094076268546', '1736287140860829901', 1),
  ('1205723159984574464', '1736622445505741159', 2),
  ('1654179648761102337', '1734655981412573264', 1),
  ('1654179648761102337', '1734657820010991962', 1),
  ('1654179648761102337', '1734661707908120774', 1),
  ('1654179648761102337', '1734661220601352349', 1),
  ('1654179648761102337', '1734659777710498122', 1),
  ('1654179648761102337', '1734658201302577473', 1),
  ('1398658879827886085', '1734656185226424457', 2),
  ('1710302619371982848', '1734799518305034308', 2),
  ('1710302619371982848', '1734798971577536976', 1),
  ('1710302619371982848', '1734798971577536976', 2),
  ('1710302619371982848', '1734798173476950425', 1),
  ('1697206791715635200', '1734826498190717030', 3),
  ('1147519752740491266', '1734531657623859526', 1),
  ('1659343792422395907', '1734526333823176854', 1),
  ('1489683639180869634', '1734555309257982198', 1),
  ('1697206791715635200', '1734830893259239623', 3),
```

Figure 1: Targer user 1

```
[ ] target_user_tweets = parsed_rdd.flatMap(lambda x: target_user_tweets_parser(target_user,x))
target_user_tweets.collect()

[[
وصیت نامه یکی از شهادت بزرگوار : خدايا نشود که در زندگی دو دستی به دوخت زندگی بچشم و همچون منوه ای یوسیده بپاشم خدايا تا ایملی هست مرا برای خود قریبی')
شکر_تلویان_n\گنم_تل_ملر_n\کن https://t.co/CzBw8KokMm',
'generated'),
شکر_تلویان_n\گنم_تل_ملر_n\کن, اینکه قهرمان زندگی چه شخصی باشد در عقبیت بخیری خیلی مهمه\u200cزندگی یر از انتخاب هایی است که مسیر ها رو عوض می')
https://t.co/QGfEckpZ14',
'generated'),
https://t.co/JWvgiokZzn', گنم_تل_ملر_n\ها\u200cاداع_با_لا\u200cها\u200cا\u200cحضور مردم قدر نفس دینر تلویان در مراسم وداع با لاله',
'generated'),
https://t.co/1NNofp1AGe', گنم_تل_ملر_n\ها\u200cاداع_با_لا\u200cها\u200cا\u200cدینر تلویان و مردم همیشه در صحنه',
'generated'),
https://t.co/QgF5gIXk19', ها\u200cاداع_با_لا\u200cها\u200cا\u200cگنم_تل_ملر_n\...در بی اسم و',
'generated'),
می زندگنوهن می کنند و به آب هم میزنن \u200cاست قلم\u200cقیلا بلور نمی کردم-ولی باچتمان خود می بینم همنلیکه دم از وحدت و برادری می زندم-اسفاده اینجا علیه مقدمات اهل')
', همه توصیه ونگهدشتون اطاعت نمی کنید\u200cانتگش!نمی دوئم فرتون چیه!شماره خونتون رو بپرو رهبری قزاقه مان میونیدجیجور از آن',
'generated'),
و ایک رهروان مکتب حاج قاسم در بلبل فرزندان خونتلمان\ن\خسته نیاشی بیلون
https://t.co/w8Tds48qlM', شکر_تلویان_n\گنم_تل_ملر_n\❤ولمان را در آغوش میکشد
'generated'),
', گنم_تل_ملر_n\ها\u200cاداع_با_لا\u200cها\u200cا\u200cبوی ملر میدهد\آمیتت چه آمیتیت که عجب هوای شیرن\شید گنم سلام',
'generated'),
https://t.co/Zj4jzA9tCG', شکر_تلویان_n\گنم_تل_ملر_n\❤خطرات مرداتی که پروای نام ندارند و در کیف گنم میهند\.. در زیست',
'generated'),
', گنم_تل_ملر_n\شکر_تلویان_n\تو نشان از مادر داری گنم بمان\بخدا گنم یودن زیست\گنم بمان شید عزیز',
'generated'),
https://t.co/1oQCzwZBqC', شکر_تلویان_n\مهماتی که #گنم_تل_ملر هستند\های محصله و نفاع مقص\u200cمهماتی ویژه از سل\دینر تلویان این روزها مهمان دارن',
'generated'),
', شکر_تلویان_n\گنم_تل_ملر_n\شید گنم مله برکت برای نظام و باعث وحدت و هملی میتود',
'generated'),
```

Figure 2: Targer user tweets

```
[ ] similar_users = sorted_similars_rdd.top(20)
similar_users

[(1.6035674514745464, '1724483603600486400'),
 (1.6035674514745464, '1632351425383522312'),
 (1.6035674514745464, '1487097029990002688'),
 (1.3416407864998738, '774470366676877312'),
 (1.3416407864998738, '1725652340043100160'),
 (1.3416407864998738, '172582061237300965378'),
 (1.3416407864998738, '1721135891308007424'),
 (1.3416407864998738, '1719916251404021760'),
 (1.3416407864998738, '1719773177503682560'),
 (1.3416407864998738, '171969619692292096'),
 (1.3416407864998738, '1718837061237530624'),
 (1.3416407864998738, '1710302619371982848'),
 (1.3416407864998738, '1708866964020514816'),
 (1.3416407864998738, '1704606183858880513'),
 (1.3416407864998738, '1704531218174251012'),
 (1.3416407864998738, '1692155672807325699'),
 (1.3416407864998738, '1689309272180482048'),
 (1.3416407864998738, '1685928050708418560'),
 (1.3416407864998738, '1675154446966312962'),
 (1.3416407864998738, '1648401324608827392')]
```

Figure 3: Similar users and corresponding similarity score

```
[ ] recommend_tweets_ids = recommend_tweets_ids[0:20]
recommend_tweets_ids

[(4.810702354423639, '1734904359844982845'),
 (4.810702354423639, '1734599807044349982'),
 (4.810702354423639, '1734204787200004101'),
 (4.024922359499621, '1735389039786668038'),
 (4.024922359499621, '1734604136702935127'),
 (4.024922359499621, '1735901407872475393'),
 (4.024922359499621, '1735901649871253757'),
 (4.024922359499621, '1736285402858688631'),
 (4.024922359499621, '1736285577358573735'),
 (4.024922359499621, '1736348680079233107'),
 (4.024922359499621, '1736349049442210119'),
 (4.024922359499621, '1736349214286688597'),
 (4.024922359499621, '1736349372965655032'),
 (4.024922359499621, '1734182827418095904'),
 (4.024922359499621, '1736000149338325236'),
 (4.024922359499621, '1734227232418075035'),
 (4.024922359499621, '1734234103837327546'),
 (4.024922359499621, '1735690223877501163'),
 (4.024922359499621, '1734273745986994506'),
 (4.024922359499621, '1735714482779750596')]
```

Figure 4: Recommended tweets and corresponding tweet score

[illegible]

Figure 5: Recommended tweets

3.2 User #2

```
[ ] target_user = grouped_rdd.takeSample(False,1,1000)[0]
target_user

('1036885525326647296',
 [('1591547142136627202', '1736619970027462949', 1),
 ('1227221292685938688', '1734241138918175049', 1),
 ('1227221292685938688', '1733918649746919740', 1),
 ('1591547142136627202', '1734596007038517749', 1),
 ('1591547142136627202', '1735659470351433916', 1),
 ('1591547142136627202', '1735310756118155394', 1),
 ('1591547142136627202', '1736442052454310094', 1),
 ('1591547142136627202', '1736076770128851061', 1),
 ('1591547142136627202', '1736096711972475174', 1)])
```

Figure 6: Targer user 2

```
[ ] target_user_tweets = parsed_rdd.flatMap(lambda x: target_user_tweets_parser(target_user,x))
target_user_tweets.collect()
```

'روحانی#n'. در ملک‌ای که شورای نگهبان مصلحت‌اندیشی کند و مجمع تشخیص مصلحت‌های فقهی هیچ چیز درست نخواهد شد'.

'generated'.

فی‌المثل آیا با گروه‌های دیگر حد البته نظر شما برای حقیر مهم و تعیین‌کننده خواهد بود؟ لطفاً بفرمایید از بنده انتظار دارید در انتخابات پیش رو چه کشتگری داشته باشیم؟ اینبار فرق می‌کند؟

و... چه مسایلی برایتان اولویت دارد؟ خطوط قرمز شما چیست؟ اگر به چه اشخاص، گروه‌ها و جریان‌هایی؟ انجمن کیم یا خیر؟

'generated'.

باید با کمبود باید تکلیف‌مان را با #سريلزای اجباری مشخص کنیم؟ یک بار برای همیشه باید تکلیف‌مان را با بی‌تفاوتی و در مواردی شبیه‌ت دولتی پاکستان در برقراری امنیت مرزها مشخص کنیم؟

باید تکلیف‌مان را با مسئولین بی‌عرضه‌ای که در تهران نشسته و با عیبه اینفلوئنسر تروریست‌ها در مرزهای داخلی! مشخص کنیم؟ تجهیزات و نبود آموزش نیروهای مرزی مشخص کنیم؟

'ایران تسلیت#n#راسک#n' برای مرزها نسخه غلط می‌پیچد مشخص کنیم

'generated'.

'@M_Afsari2 @RezaGholami1354 در باب گفتگو هم قرار بود دیروز ایشان تشریف بیآورند دقت ما که ظاهر کاری برایشان پیش آمد و از جانب ایشان موکول شد به بعد از ایام شهادت

'replied'.

زای.ن۲! بحث تبدیل، بحثی انحرافی است، حق پدرم تیرنه است. n۲\اند u200c\گاه دیلمت نبوده و از سال ۷۱ که از دانستنی استغنا می‌دهند، کارگردانی داشته u200c\پدرم هیچ. n۱:چندنگاه مهم'.

'https://t.co/DBxSHJlQv1' n#FreeHamidNouri!! حکم جیس اید هنوز نهایی نشده است؟ شنبه هفته قی حکم تجدیدنظر صادر می‌شود u200c\بدوی، ۲۳ تیر ۱۴۰۱ صادر شده و سه

'generated'.

طبق صحیح‌ترین منابع - n۲ منابع شیعیه و سنی معتقد و معتزله هرکس بدون امام و بیعت او می‌رسد به مرگ جاهلیت مرده است لذا جهنمی است - n۱\! من امام فاطمه س را امام می‌دانم'.

دوباره طبق منابع معتبر سنی حضرت زهرا س بدون بیعت با خلیفه اول از دنیا رفت در حالی که از او خشمگین و ناراضی - n۲\ اهل سنت حضرت زهرا سلام الله علیها سید زنان بهشت است

سوال مهم: یا حضرت زهرا س المایه الله جهنمی است که نیست یا خلیفه رسول الله ایوبی نبوده و حضرت فاطمه س در بیعت با امام و خلیفه‌ای دیگر (امام علی ع) از دنیا رفته! بود

'یا زهراء#n'. است

'generated'.

برادران! فاطمه سلام الله علیها بنیادگذار این مکتب است n... آن مظلوم باید مقدر باشد نه ذلیل: منتها یک شرط دارد. هیچ چیز مانند مظلومیت، رسواکننده ظلم و ظالم نیست'.

'عزم-تلاطم#n' مکتب-مظلوم#n\، خواهرانمان در عزم خواه بداند یا نداند پیر این مکتب هستند u200c\و

'generated'.

دوما: در متون روایی و تاریخی به این مهم اشاره شده است مانند اصول کافی و من لایحضره الفقیه به سند صحیح n\ اولاً: تقویم جمهوری اسلامی که مبنای تشخیص امر تاریخی نیست'.

سوما: وفات در زبان عربی n\ قال: إِنَّ قَلْبِي عَ مَبِيتِي تَهَيَّأَ وَ إِنَّ بَيْتِي الْأَكْبَرُ لَا يَطْمَئِنُّ. [1] امام کاظم علیه السلام می‌فرماید: فاطمه (علیها السلام) صدیقه‌ای است که شهید شده است n\:

رایجاً: ملاحظه افرادی که تحلیل شان مانند شمامت باعث شد از دهه هفتاد n\ شمال مرگ و شهادت است و متحد در اسناد تاریخی و روایی درباره افراد مختلف که شهانستان قطعی است استفاده شده

'... به بعد از لفظ شهادت استفاده شود

'replied'.

Figure 7: Targer user tweets


```
[ ] similar_users = sorted_similars_rdd.top(20)
similar_users

[(1.3416407864998738, '988684923203710976'),
 (1.3416407864998738, '952463173868818432'),
 (1.3416407864998738, '951201559597191169'),
 (1.3416407864998738, '948126072603754496'),
 (1.3416407864998738, '941220112845164544'),
 (1.3416407864998738, '919323412195102721'),
 (1.3416407864998738, '918806712995348480'),
 (1.3416407864998738, '911453224049164288'),
 (1.3416407864998738, '90177775383332736'),
 (1.3416407864998738, '880353064477233152'),
 (1.3416407864998738, '878382611730751488'),
 (1.3416407864998738, '878217414323961856'),
 (1.3416407864998738, '874242781790818305'),
 (1.3416407864998738, '855438360730578946'),
 (1.3416407864998738, '834665971075907585'),
 (1.3416407864998738, '822564412997402624'),
 (1.3416407864998738, '819266108351856640'),
 (1.3416407864998738, '818514349782405120'),
 (1.3416407864998738, '818336090998374400'),
 (1.3416407864998738, '814591408862416897')]
```

Figure 8: Similar users and corresponding similarity score

```
[ ] recommend_tweets_ids = recomm_tweets_ids[0:20]
recommend_tweets_ids

[(4.024922359499621, '1736200632980103386'),
 (4.024922359499621, '1736200676332417325'),
 (4.024922359499621, '1736571240553283866'),
 (4.024922359499621, '1736571274892296333'),
 (4.024922359499621, '1735849256391475619'),
 (4.024922359499621, '1735849284082303219'),
 (4.024922359499621, '1735867134490607746'),
 (4.024922359499621, '1736640572977250357'),
 (4.024922359499621, '1736663226119512083'),
 (4.024922359499621, '1735958921079112138'),
 (4.024922359499621, '1736695367943561440'),
 (4.024922359499621, '1736703321509892379'),
 (4.024922359499621, '1736722218724274318'),
 (4.024922359499621, '1735994535237325134'),
 (4.024922359499621, '1736358408054296791'),
 (4.024922359499621, '1735647971692204432'),
 (4.024922359499621, '1735648652553560540'),
 (4.024922359499621, '1736760514624016488'),
 (4.024922359499621, '1735654953446736373'),
 (4.024922359499621, '1735672954925515228')]
```

Figure 9: Recommended tweets and corresponding tweet score

Figure 10: Recommended tweets

Figure 11: Targer user 3

Figure 13: Similar users and corresponding similarity score

```
[ ] recommend_tweets_ids = recomm_tweets_ids[0:20]
recommend_tweets_ids

[(4.024922359499621, '1734638622920298899'),
 (4.024922359499621, '1736829679644868705'),
 (4.024922359499621, '1735674240941035753'),
 (4.024922359499621, '1735703356704727448'),
 (4.024922359499621, '173430308772683235'),
 (4.024922359499621, '1735743796418150653'),
 (4.024922359499621, '1735431247524999614'),
 (4.024922359499621, '1736551527198658808'),
 (4.024922359499621, '1735489509817532581'),
 (4.024922359499621, '1735490187638034513'),
 (4.024922359499621, '1735491111907393642'),
 (4.024922359499621, '1735492658921923071'),
 (4.024922359499621, '1735493682059452843'),
 (4.024922359499621, '1734512533652857019'),
 (4.024922359499621, '1734908880897929434'),
 (4.024922359499621, '1734665362262671482'),
 (4.024922359499621, '1734665362262671482'),
 (4.024922359499621, '1735713588088607109'),
 (4.024922359499621, '1736168545619231053'),
 (4.024922359499621, '1736169799418331389')]
```

Figure 14: Recommended tweets and corresponding tweet score

[illegible]

Figure 15: Recommended tweets

3.4 User #4

```
target_user = grouped_rdd.takeSample(False,1,777)[0]
target_user

('4702453171',
 [('4702453171', '1734001719539810691', 3),
 ('4702453171', '1734471154646049090', 3),
 ('4702453171', '1734471452450005121', 3),
 ('1653467502972682270', '1734079495177461916', 2),
 ('25710094', '1734490763394798040', 2),
 ('60171166', '1736513161631588535', 2),
 ('17532913', '1734073641908981882', 2),
 ('4702453171', '1733813004280459569', 2),
 ('4702453171', '1734086631366311958', 3),
 ('4702453171', '1734086631366311958', 2),
 ('17532913', '1734087759646871761', 2),
 ('4702453171', '1734096096526291087', 3),
 ('4702453171', '1734099408688488815', 3),
 ('4702453171', '1734099878626705917', 3),
 ('4702453171', '1734100866527817896', 3),
 ('4702453171', '1735916939380212209', 3),
 ('4702453171', '1734103110279102546', 3),
 ('4702453171', '1726253167547924649', 2),
 ('15090494', '1734104725874713032', 2),
 ('360461731', '1734106135055958309', 2),
 ('1271067471760568322', '1734105477301117383', 2),
 ('4702453171', '1734108334402228647', 3),
 ('767126934', '1734108466426380588', 2),
 ('302678816', '1736273645201182727', 2),
 ('4702453171', '1736290273133490433', 3),
 ('302678816', '1736289193905594796', 2),
 ('767126934', '1734118849220374966', 2),
 ('712752292194357249', '1735382946834518426', 2),
 ('767126934', '1734119784206172308', 2),
```

Figure 16: Targer user 4

```
[ ] target_user_tweets = parsed_rdd.flatMap(lambda x: target_user_tweets_parser(target_user,x))
target_user_tweets.collect()
```

((: ولی من یا ۳۰۰ لیر می‌تونم بخرمش حالت عالی ۱۰۰ دلارمه واقعا. رو برای استور ترکیه ببینید #Blinkist شما فقط قیمت یک آب مثل'))
<https://t.co/snE0lCFuKD>,
 'generated'),
 زده ۲۵ گیگ اینترنت با ۳ ماه یوتیوب یرمیوم و دیتای نامحدود برای واتس‌اپ و یوتیوب و اسپاتیفای. همه اینا برای ۳ ماه اول ۳۷۰۰ لیر. آقا چرا بسته های اینترنتی این کنار اینترنت خون آخه
 لیر <https://t.co/xoSKCGtk1k> Vodafone 🇮🇷 چرا تو ایران همچین چیزایی نداریم:/:
 'generated'),
 قدم بعدی حتما خرید یک روتر n/n 🇮🇷 چوشتی اندروید که واقعا وضیعت و خیمه داره رو بعدا میدم تمیز ولی تا اون موقع به عنوان روتر دارم استفاده میکنم ازش و خیلی خوب شده
 واقعه
 'generated'),
 برای گوگل n/ولی چرا سیستم پرداخت ایل اینترنت بی دردم تر از گوگل n/کلی هم یلیند به قوانین مالی آمریکا هستند ایل هم بزرگه n/گوگل بزرگه n/من درک نمیکم این قضیه رو
 همین موردم باعث میشه بازم بخوام این گوشتیو نگه n/n. باید کلی سند و مدرک و ایپی رزینشنال جور کنی. ولی ایل میگه فقط پولو بده بهم و به روش پرداخت تو اون کشور داشته باش و تمام
 دارم.
 'generated'),
 (@Nikolaona ها، ولی حالا دیگه مجبورم داخل فولدرشم بزنم u200cمارک u200cکل یوک
 'replied'),
 (@Nikolaona ها، ولی حالا دیگه مجبورم داخل فولدرشم بزنم u200cمارک u200cکل یوک
 'replied'),
 امیدوارم مدت زیادی n/n نگهتیم خیلی خوب بلطری نگه میداره و منی که از یک گشتی با اسکرین تلم نهایت ۲ ساعت میام. ایپی که کل روز شارژ نگه میداره برام خیلی خوشحال کنتم
 رو همینطوری نمونه
 'generated'),
 ملت از دید آزادی اینترنت و گرون بودنش میگویند که خب هر دو n/n اون از دید اقتصادی میگه آقا وضیعت درآمدی خرابه n/ولی هر دلی که منبرعامل شاتل توپیت میزنه بهش حمله میکنم
 ()))): اینا دست حکومت و نه یک شخص و مدیریش. دیگه در این حد آزادی نداریم
 'generated'),
 ملت از دید آزادی اینترنت و گرون بودنش میگویند که خب هر دو n/n اون از دید اقتصادی میگه آقا وضیعت درآمدی خرابه n/ولی هر بار که منبرعامل شاتل موبایل توپیت میزنه بهش حمله میکنم
 ()))): دو اینا دست حکومت و نه یک شخص و مدیریش. دیگه در این حد آزادی نداریم
 'generated'),
 سال با به گشتی سر کردم که اسکرین n/5 ()): فریادی بلطری بد دارم n/یکی باید یک توضیح عملی و درست حسابی بده من چطوری این آیفون رو نگه دارم که بلطریش به فنا نره
 🇮🇷 این الان به ۷ ساعتی میده خیلی خرابه n/نیم ۲ ساعت میداد
 'generated'),

Figure 17: Targer user tweets

```
[ ] similar_users = sorted_similars_rdd.top(20)
similar_users

[(1.386750490563073, '855095592'),
 (1.386750490563073, '2564151662'),
 (1.386750490563073, '1572479734335815684'),
 (1.386750490563073, '1206642366830002179'),
 (1.386750490563073, '1152293110212263937'),
 (1.386750490563073, '1151820549653446656'),
 (1.386750490563073, '1141720339'),
 (1.3416407864998738, '941203650818560000'),
 (1.3416407864998738, '927327889'),
 (1.3416407864998738, '880550727680184320'),
 (1.3416407864998738, '860039158370316288'),
 (1.3416407864998738, '818518646968422400'),
 (1.3416407864998738, '801430424'),
 (1.3416407864998738, '75530244'),
 (1.3416407864998738, '720019748780699649'),
 (1.3416407864998738, '712037416345989123'),
 (1.3416407864998738, '36924087'),
 (1.3416407864998738, '36024428'),
 (1.3416407864998738, '30186494'),
 (1.3416407864998738, '2421831930')]
```

Figure 18: Similar users and corresponding similarity score

```
[ ] recommend_tweets_ids = recomm_tweets_ids[0:20]
recommend_tweets_ids

[(4.160251471689219, '1736313883789312092'),
 (4.160251471689219, '1734302996685582511'),
 (4.160251471689219, '1736543598210908218'),
 (4.160251471689219, '1736191078921826671'),
 (4.160251471689219, '1736194209722352050'),
 (4.160251471689219, '1736196909767745970'),
 (4.160251471689219, '1735911029438083239'),
 (4.160251471689219, '1734100975072207178'),
 (4.160251471689219, '1736660798875881627'),
 (4.160251471689219, '1734851686357782537'),
 (4.160251471689219, '1734118177137717730'),
 (4.160251471689219, '1735933230006558790'),
 (4.160251471689219, '1734597062904140077'),
 (4.160251471689219, '1734599210136154571'),
 (4.160251471689219, '1736693297865806316'),
 (4.160251471689219, '1736693537289228516'),
 (4.160251471689219, '1735961908656689410'),
 (4.160251471689219, '1736321041838879073'),
 (4.160251471689219, '1734158029661077888'),
 (4.160251471689219, '1736336769023385988')]
```

Figure 19: Recommended tweets and corresponding tweet score


```
[ ] recommended_tweets = parsed_rdd.flatMap(lambda x: recommended_tweets_parser(recommend_tweets_ids,x))
recommended_tweets.collect()
```

((:)) فی شبکه بیتکوین به نفسی کشیدن\آخریتش\شش\
 و\اولین هست در نوع خودش\اینسکریپت میته CBRC20 و BRC20 همزمان و بصورت دوگانه روی هردو پروتکل\خیلی چیز جالبیه \$SYMM توکن\
 خب من الان ساعت ۵ صبح یا چه حالی بلند بشم اتی بریزم تو ولت\بخدا این اصناف نیست نکومی که ما باید بخوایم\آمریکایا بیدارن همه چو یامب میکن\
 و\بعد ما میریم تسک لینیا میزنیم\عجیب و غریب\!چنان حجم پولی تزریق شده نوی اکوسیستم بیتکوین که باور نمیکند\
 و\ولی اگر مینت کردید الان یراحتی میتونید ۱۸۰۰ دلار بفروشید\بون invalid تایید شد و \$UTXO متاسفانه تراکتهای خورم روی\
 bitcoin\...بیبتیم کف، قیمت رو حمایت میکنه یا باید بریم یلینش\ولکتش به میدلین بود (قبلا چارکتو تو کتال گذاشته بودم)\ن: چه عجب بالاخره به کتال قرمز هم دیدیم\
<https://t.co/97we1i01xc>,
 و\ارزش خرید و یا نگهداری ندارن\اینسکریپشن روی هر بلاکچینی بجز بیتکوین\
 و\سخنهای خوبی براش در راه\ SRC20\NFT روی پروتکل توکنهای\باشه #SRC721 حواسون به\
 و\دامینشن نتر جای خیلی حساسی هست\ <https://t.co/UB81eBubl3>,
 و\روی شبکه بیتکوین\ \$STAMP 35x\N#SRC20 48x\N\$UTXO 14x\N\$SATO 2.5x\N\میترونه\ اینم که جدید ولی از همین الان داره میترونه\N#AVAX 3.5x\N#ATOM 1.6x\N\$TIA 1.7x\N\اگر به یلین شکسته بشه، آلتسیزنی در بیش داریم که بیا و ببین\این ترند صعودیه که یکسال هست دامینشن رو نگه داشته\
 و\خوبی داشتید، دویلت هم میتونید انجام بدید\ (به دایرکت توینتر مراجعه کنید)
 و\اگر به یلین شکسته بشه، آلتسیزنی در بیش داریم که بیا و ببین\این ترند صعودیه که یکسال هست دامینشن رو نگه داشته\ <https://t.co/p8StFdHubZ>,
 و\این ترند صعودیه که یکسال هست دامینشن رو نگه داشته\ <https://t.co/uHTbfNbm3Q>,
 ای که تحویل\u200c\بسته\عدم آیدیت وضیعت مرسله\بسته های له و لورده\تاخیر در تحویل\قرار بر این بود یست هوشمند بشه ولی الان تبدیل شده به رقیب مخبرات نوی قتل بودن\
 و\!تشد رو یسجی محترم میزنه تحویل گردیده\
 و\همش تصویر شما مردم که نمیزارید کارتونو بکنند دیگه\!حالا اگه مثل کره شمالی اینترنت و ارتباطات خارجی رو قطع کرده بودند، این مشکلات هک پیش نمیومد\
 و\همش تصویر شما مردم که نمیزارید کارتونو بکنند دیگه\!حالا اگه مثل کره شمالی اینترنت و ارتباطات خارجی رو کاملاً قطع کرده بودند، این مشکلات هک پیش نمیومد\
 و\جرا یلک بیتب قتلگم رو زود فروختم\حسرت امروز\!
 و\روی تمام بلاکچینها بودن inscription بلکه لایه های انگل وارن\بود و نه لایه ۲ NFT ترند جدید مارتک ته\
 و\مخالفت\اعتراض مننی قسط به قضیه\صادق یوقی تلقی میگرده\
 و\بهتر نبود این انرژی رو صرف اعتراض به چیزهای مهمتری میکردیم\!نداشت\؟
 و\بازی دو سر بلکت\!بخاطر فی بالا، معامله ها بسیار کم شده و قیمتها در حال ریزشه\!شبکه بیتکوین زایلیده\
 و\کد\c\200u\های عاشقانه منتشر می\c\200u\ای هرکی به ویلگ داشت و توش شس و متن\c\200u\ایش بخیر به دوره\

Figure 20: Recommended tweets

3.5 User #5

```
[ ] target_user = grouped_rdd.takeSample(False,1,91)[0]
target_user

('1693973461306109952',
 [('1574904470835798041', '1734442494039826873', 2),
 ('1574904470835798041', '1734883099647914204', 2),
 ('1574904470835798041', '1735771172472492519', 2),
 ('1693973461306109952', '1736853602969280923', 3),
 ('1693973461306109952', '1736855151917334935', 3),
 ('1693973461306109952', '1735069630308266029', 3),
 ('1693973461306109952', '1735401606160785610', 3),
 ('1693973461306109952', '1735071498032058504', 3),
 ('1693973461306109952', '1735071498032058504', 2),
 ('1574904470835798041', '1735403820757426674', 2)])
```

Figure 21: Targer user 5

```
[ ] target_user_tweets = parsed_rdd.flatMap(lambda x: target_user_tweets_parser(target_user,x))
target_user_tweets.collect()
```

'افتین خمتند\n ما دیوانه شده ها\n و این دنیای موبد جای «ما» نیست\n شده اید\n حتی «شو» نیستی\n تو «دیگر»
'generated'),

کلم قرمز ، بنفش و اکون... از آن روز به بد شدم باغبان باغچه ی کوچک با یک گل... تمام سیاره ها را گشتم تا به جایی رسیدم که خالی از هر کسی بود و آنگاه دلم آرام گریست
'generated'),

نه! هوجی، چیزی نفهمیدم، فقط خوانیدم.\n...! امروز خواستم از عشق بگم، دیدم عشقی نیست! خواستم از عاشقونه ها حرف بزنم دیدم چیزی نیست، یادم اومد از دل تنگ بنویسم که فهمیدم
'generated'),

اگره گوشت هوشیار شد، تیزش کردی و صدایش رو نشنیدی.\n اگره چشمت صدادت کرد، سرتو بلند کردی ندیدیش، بلختی\n اگره دستت بی وزن شد، حرکتش دادی و نتوانستی لمسش کنی، بلختی
'generated'),

گفت یعنی چی؟! گفتم دوس داشتن یا هست یا نه... خیلی زود جواب دادم دوس داشتن که اندازه نداره رفیق\n به روز یکی از آملای اطرافم که دنیاش خیلی رنگیه بهم گفت: چقدر دوسم داری؟!
'generated'),

و... گفت باز نفهمیدم چی میگی!!! گفتم هر وقت کسی رو دوس داشتی متوجه میشی. بهر داد زد و گفت آهان.\n نیست، اندازه نداره
'generated'),

و... عشقت روزی به خاله ی ما خواهد رسید؟\n اما مگر قانون این نبود که هر آنچه دیر می آید\n شاید خیلی دیر رسیدم خیلی دیر\n از بوی تنگ نفهمیدم
'generated')

Figure 22: Targer user tweets

```
[ ] similar_users = sorted_similars_rdd.top(20)
similar_users
```

[(1.0000000000000002, '1693973461306109952'),
(1.0, '1574904470835798041'),
(0, '999996903185100800'),
(0, '99999538'),
(0, '999995055090929664'),
(0, '999988080819015680'),
(0, '999970696435523584'),
(0, '999964377922134016'),
(0, '999963219425931266'),
(0, '999959604556697602'),
(0, '999956594686349317'),
(0, '999953030060675072'),
(0, '999952212783783939'),
(0, '999951354620727296'),
(0, '999951276'),
(0, '999928626987569152'),
(0, '999925604957937665'),
(0, '999924334402531328'),
(0, '999922321816080386'),
(0, '999918484485128195')]

Figure 23: Similar users and corresponding similarity score


```
[ ] recommend_tweets_ids = recomm_tweets_ids[0:20]
recommend_tweets_ids

[(3.0000000000000001, '1736853602969280923'),
 (3.0000000000000001, '1736855151917334935'),
 (3.0000000000000001, '1735069630308266029'),
 (3.0000000000000001, '1735401606160785610'),
 (3.0000000000000001, '1735071498032058504'),
 (3.0, '1734042842803110258'),
 (3.0, '1736261559083491431'),
 (3.0, '1735326293720228136'),
 (3.0, '1736124195354890571'),
 (3.0, '1735771172472492519'),
 (3.0, '1735406823128125835'),
 (3.0, '1735410144727478743'),
 (3.0, '1735422143100244168'),
 (2.0000000000000004, '1734442494039826873'),
 (2.0000000000000004, '1734883099647914204'),
 (2.0000000000000004, '1735771172472492519'),
 (2.0000000000000004, '1735071498032058504'),
 (2.0000000000000004, '1735403820757426674'),
 (0, '1735767625915859133'),
 (0, '1735872482198860152')]
```

Figure 24: Recommended tweets and corresponding tweet score

```
[ ] recommended_tweets = parsed_rdd.flatMap(lambda x: recommended_tweets_parser(recommend_tweets_ids,x))
recommended_tweets.collect()
```

[[('!است!خوبتخت!چقدر!که! نمی داند! "آجا!n" \n\n!آجایی!آو!')
 اینا دیگه تمارف رو کنار!ها هست(منم بلاک کرده) امروز علنا گوشه ایون ۱۳(یرومکس تو برنامه زنده گرفته دستش u200c!مقدم که مجری محبوب آخونده و تو همه شبکه u200c!این شریفی.)
),('گشتن و به ریش ملت میخندن
),('...کسانی کسی که دوستش داری شروع میشه! بخیر ترین صبح قطعا اونیه که!
),('...آجا که تو نیستی، کجا خواهد بود!
),('افشین!خمتان!n ما بخواه شده ها!n و این دنیای موبد جای «ها» نیست!n!شما» شده ای!n!حتی «هو» نیستی!n تو « دیگر»!
 گلم قرمز ، بنفش و اکون...از آن روز به بعد شدم باغبان باغچه ی کوچک یا یک گل...تمام سیاره ها را گشتم تا به جایی رسیدم که خالی از هر کسی بود و آنگاه دلم آرام گریست!
),('...ای را در خود دارد
),('...و خدا همه چیز را درست میکند! گمان میکنی پایان است!
 نه! هیچی، چیزی نفهمیدم، فقط خوابیدم...!!n...امروز خواستم از عشق بگم، دیدم عشقی نیست! خواستم از عشقونه ها حرف بزنم دیدم چیزی نیست، یادم اومد از دل تنگ بنویسم که فهمیدم
),('...همین
 اگر گوشتا هوسیل شد، نیزش کردی و صداقت رو نشنیدی!n!اگر چشت صداقت کرد، سرتو بلند کردی ندیدیش، بلختی!n!اگر دستت بی وزن شد، حرکتش دادی و نتوانستی لمسش کنی، بلختی!
),('...و امشب من بازنده ترین!n!اگر دلت تنگ شد، موبایلت رو ورداشتی، زنگش زدی، جوابت رو نداد، بلختی!n!بلختی
 گفت یعنی چی؟! گفتم دوس داشتن یا هست...n!خیلی زود جواب دادم دوس داشتن که اندازه نداره رفیق!n!به روز یکی از آهای اطرافم که دنیاش خیلی رنگیه بهم گفت: چقدر نوسم داری?!
),('...گفت باز نفهمیدم چی میگویی!!! گفتم هر وقت کسی رو دوس دایمی متوجه میشی. یهو داد زد و گفت آهان!n!نیست، اندازه نداره
 بجز واحدهای نویو، اطلاعات، یگان ویژه و امنشون و اشخاص!n!استخدام بیمانی پلیس تنها جا دست خلیپها برای رسیدن به حناقلهاس، فاسد هستن اما اکثرا شر نیستن @Orwellian2017k
),('و...از قدیم گفتن مرگ جزئی از این نظام!n!اهداف آنست! اند اما تکتیری بر هسته اصلی سرکوب(سیاه/سرسبز/آتش به اختیار) نداره!n!سرسبزده
),('...عاقبت روزی به خاله ی ما خواهد رسید!n!اما مگر قاتون این نبود که هر آچه دیر می آید!n!شاید خیلی دیر رسیدم خیلی دیر!n!از بوی تنگت فهمیدم!
),('...شهر حسودی بکند!n!بظلم کن همه جا!
),('...در هستی پنهان من!n!ای هست تو پنهان شده!n!ای روی تو ایمان من!n!ای دیدن تو دین من!
 کجا!n!به او می گویم حالت چطور است؟ و به تو فکر میکنم!n!به آن زن تلخ میکنم و درعین حال به تو فکر میکنم!n!کرم می کنم و به تو فکر میکنم!
),('...به او ایخند می زنم!مشماتم را می بندم و به تو فکر می کنم!n!برویم گویه تو فکر میکنم

Figure 25: Recommended tweets

4 Future works

Here are some additional metrics and features that could be incorporated to improve user similarity analysis and tweet recommendations:

- **Hashtag usage**

tracking which hashtags users frequently use can help identify shared topical interests. Correlating common hashtags can boost similarity.

- **Tweet entity extraction**

extracting mentioned users, links, media, locations etc. can uncover more semantic similarities between users' tweets.

- **Follower/following overlap**

users sharing a lot of followers/followings may have closer alignment of interests. The social graph offers useful signals.

- **Tweet timestamp analysis**

users tweeting at similar times of day or days of week may have linguistic/interest similarities.

- **Tweet language modeling**

topic models and distributional semantic vectors derived from tweet text can reveal similarities.

- **Profile metadata**

shared descriptors in user profiles like location, website, bio keywords could indicate similarities.

- **Tweet propagation patterns**

how tweets spread can indicate influence and interest similarities between original posters.

5 Conclusion

The results verify the proper functionality of the implemented recommender system, which effectively recommends tweets highly similar to a user's related tweets. The system was tested on various users with differing tweet volume, viewpoints, and tweet subjects. In all cases, accurate recommendations were generated, with only a negligible amount of uncorrelated tweets. The implementations were demonstrated for 5 sample users due to resource constraints, but could be scaled up to all users given a more powerful system configuration.

A key advantage of this system is that it can operate in both online and offline settings. When new tweets or users are added, only the affected user recommendation lists need to be efficiently updated while all other users remain unchanged. This allows incorporating new data with minimal computation by incrementally updating tweet scores as needed, rather than in real-time. In summary, the personalized tweet recommendation engine developed provides relevant engagement suggestions by modeling user similarity based on tweet activity patterns. The map-reduce algorithm combined with Pearson correlation user profiling effectively produces customized recommendations. The system's scalability and updatability makes it well-suited for real-world usage on Twitter's firehose of ever-growing data.

Is the implemented model online or offline?

Implemented collaborative filtering algorithm is naturally more amenable to online updates:

- User tweet engagement profiles are directly compared via Pearson correlation without a full graph precompute
- New tweets or users can have their profiles appended to the system and correlated against existing ones
- Only correlations involving the new data would need recomputation

So as new tweets come in:

- Their engagements would update related user profiles

- Those user profile changes would incrementally adjust correlations and recommendations

This localized adjustment to impacted user profiles and correlations enables easy assimilation of streaming data.

Additionally, optimization strategies like approximate similarity search techniques could accelerate inbox recommendation generation over growing user bases.

The intrinsic design of profiling users, correlating their behaviors, and making recommendations avoids an upfront graph construction allowing innate online capabilities as data evolves.