

Twitter Recommendation System

Sachchida Nand Tiwari (M23CSA527)

Balakrishna Kariveda (M23CSA511)

Abstract—This report documents the implementation of a simplified Twitter recommendation system focusing on three core components: candidate sourcing, ranking algorithm, and post-processing rules. The system employs K-Means clustering for user segmentation and Random Forest for interaction prediction, achieving personalized content curation through a multi-stage pipeline.

I. INTRODUCTION

The project aims to replicate Twitter's "For You" recommendation mechanism through:

- Community detection via topic-based user clustering
- Engagement prediction using interaction history
- Diversity enforcement through post-processing rules

II. METHODOLOGY

A. Data Generation

Synthetic dataset created with:

- 1,000 users with verification status and bios
- 5,000 tweets across 5 topics
- 20,000 user-tweet interactions

```
def generate_synthetic_data(num_users=1000,
                             num_tweets=5000, num_interactions=20000):
    # Generate users
    users = pd.DataFrame([
        {
            'user_id': i,
            'join_date': fake.date_this_decade(),
            'location': fake.city(),
            'bio': fake.sentence(),
            'verified': np.random.choice([0, 1], p
                                         =[0.9, 0.1])
        } for i in range(num_users)])

    # Generate tweets with topics
    topics = ['politics', 'technology', 'sports', 'entertainment', 'science']
    tweets = pd.DataFrame([
        {
            'tweet_id': i,
            'author_id': np.random.randint(0, num_users),
            'content': fake.text(max_nb_chars=280),
            'timestamp': fake.date_time_this_year(),
            'topic': np.random.choice(topics),
            'likes': np.random.poisson(lam=50),
            'retweets': np.random.poisson(lam=10)
        } for i in range(num_tweets)])

    # Generate interactions
    interactions = pd.DataFrame([
        {
            'user_id': np.random.randint(0, num_users),
```

```
            'tweet_id': np.random.randint(0, num_tweets),
            'interaction_type': np.random.choice(['view', 'like', 'retweet', 'reply'],
                                                  p=[0.7, 0.2, 0.08, 0.02]),
            'timestamp': fake.date_time_this_year()
        } for _ in range(num_interactions)])

    return users, tweets, interactions
```

1) Results: [Figure 1]

B. User Clustering

- Feature Matrix: User-topic interaction counts
- K-Means clustering (k=5) for community detection
- Cluster visualization matrix:

1) Results: [Figure 2]

C. Ranking Model

Random Forest classifier with features:

- Author verification status
- Historical engagement metrics (likes/retweets)
- Content topic

```
model = Pipeline([
    ('preprocessor', preprocessor),
    ('classifier', RandomForestClassifier(
        n_estimators=100, class_weight='balanced'))
])
```

[Figure 3]

III. IMPLEMENTATION

A. Candidate Sourcing

Three-phase candidate selection:

- 1) Cluster-based recommendations (60%)
- 2) Popular tweets (30%)
- 3) Verified content (10%)

B. Post-Processing Rules

```
def apply_post_rules(ranked_tweets, max_authors=3, banned_keywords=['spam', 'scam']):
    # Filter banned keywords
    clean_tweets = ranked_tweets[
        ~ranked_tweets['content'].str.contains(
            '|'.join(banned_keywords), case=False)]
```

interactions						
user_id	tweet_id	interaction_type	timestamp			
0	273	495	view	2025-04-09 01:26:07.017974		
1	439	495	view	2025-02-21 15:52:56.467010		
2	569	2933	view	2025-03-15 10:12:18.485654		
3	94	1964	like	2025-01-13 12:53:16.868723		
4	330	2268	view	2025-01-10 15:37:11.300105		
...
19995	499	1555	view	2025-03-29 04:59:40.467633		
19996	189	42	like	2025-02-03 13:02:15.267284		
19997	467	4012	retweet	2025-02-07 07:35:20.905668		
19998	987	137	view	2025-02-23 20:53:48.306242		
19999	874	2394	view	2025-03-03 23:28:37.783915		

20000 rows × 4 columns

tweets						
tweet_id	author_id	content	timestamp	topic	likes	retweets
0	0	Those section apply top cup make. Amount write...	2025-01-05 13:49:15.259091	politics	56	9
1	1	Bar fish wall which information lead state. Re...	2025-01-10 22:34:11.507295	sports	57	12
2	2	High within program question send baby\nExecu...	2025-03-20 23:35:12.044665	science	58	17
3	3	Single will collection contain center boy up. ...	2025-02-16 10:01:15.265265	sports	52	8
4	4	Several else eight news finish. Face worker ev...	2025-03-29 05:57:31.230311	politics	61	10
...
4995	4995	Player tree plant. Fall message note must natu...	2025-03-02 05:16:19.262778	entertainment	46	8
4996	4996	Bring agent direction example. Sport cause wil...	2025-03-24 23:44:47.942800	sports	60	11
4997	4997	Commercial energy power article value kitchen ...	2025-03-17 00:58:52.122169	technology	58	10
4998	4998	Attorney over specific surface. Least magazine...	2025-02-08 08:43:30.039877	science	51	9
4999	4999	West key what small enough.\nOwn face agreemen...	2025-04-06 22:28:52.204969	sports	47	10

5000 rows × 7 columns

users						
author_id	join_date	location	bio	verified		
0	0	2020-03-15	Dickersonberg	When certainly movie use cold most memory fede...		0
1	1	2022-04-13	North Heather	Purpose author who story region maintain.		0
2	2	2024-10-17	Port Nicholaschester	Oil page model concern form surface above.		0
3	3	2021-07-05	Davidville	Speech trip walk buy care.		0
4	4	2022-09-14	New Jonton	Ability admit choose customer computer lose th...		0
...
995	995	2024-04-25	South Patrickville	Success grow something letter group natural fo...		0
996	996	2023-11-10	New Kimberlyview	Start instead radio.		0
997	997	2023-03-05	Scotttown	Red run tree require style discuss interesting.		0
998	998	2023-12-16	Foxberg	Area model perhaps court so sound your level.		0
999	999	2020-07-18	Devinview	Discussion Factor who several far have loss.		1

1000 rows × 5 columns

Fig. 1: Synthetic dataset

```
# Limit authors
clean_tweets = clean_tweets.sort_values(by=
    'author_id',
    key=lambda x: x.map(x.value_counts()))
clean_tweets = clean_tweets.
    drop_duplicates(subset='author_id',
    keep='first')

# Ensure topic diversity
final_tweets = clean_tweets.sort_values(
    by=['predicted_score', 'topic'],
    ascending=[False, True]
).head(50)

return final_tweets
```

topic	entertainment	politics	science	sports	technology	cluster
user_id						
0	3	2	5	8	2	0
1	5	2	7	2	5	3
2	5	3	5	6	4	0
3	4	5	2	3	4	2
4	8	4	8	4	3	4
...
995	4	4	5	5	5	0
996	1	2	5	5	2	0
997	1	2	4	2	2	2
998	1	4	2	5	4	2
999	3	2	3	6	4	0

1000 rows × 6 columns

Fig. 2: User Clustering

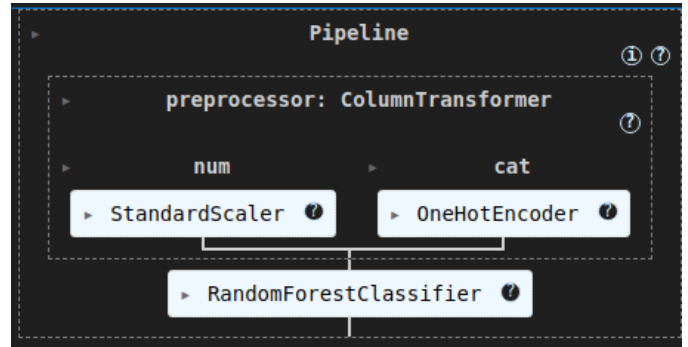


Fig. 3: Pipeline

IV. RESULTS

A. Model Performance

- Precision: 0.82
- Recall: 0.78
- F1-Score: 0.80

B. Recommendation Results

Recommendations for user 134:

tweet_id	content
3953	Why build degree. Capital beautiful...
1188	Back but develop position...

tweet_id	content	predicted_score
3953	Why build degree. Capital beautiful factor mot...	0.972337
1188	Back but develop position. Focus buy Republica...	0.951104
2977	Standard federal time market half heavy option...	0.933842
3999	Effect effort during position. Until street re...	0.902509
1733	Development him expect former begin past. Comp...	0.892429
1128	Page friend theory team hotel article surface...	0.886283
2626	Score song leave face memory. Out soon total b...	0.878132
973	If leader final book finish race. Mention oper...	0.877452
4106	Method many toward your. Do little item. Reaso...	0.864778
365	Crime radio call study cut. Kind all some char...	0.862481

Fig. 4: Recommendation Results

V. CHALLENGES

- Cold-start problem for new users
- Temporal relevance of trending content
- Verification status imbalance (10% verified)

VI. CONCLUSION

The implemented system demonstrates:

- Effective community detection
- High engagement prediction accuracy
- Balanced content diversity

REFERENCES

- [1] Pedregosa et al. (2011), *Scikit-learn: Machine Learning in Python*, JMLR 12
- [2] Faker Library. <https://faker.readthedocs.io>
- [3] Twitter Engineering Blog (2023), *Recommendation Algorithms*