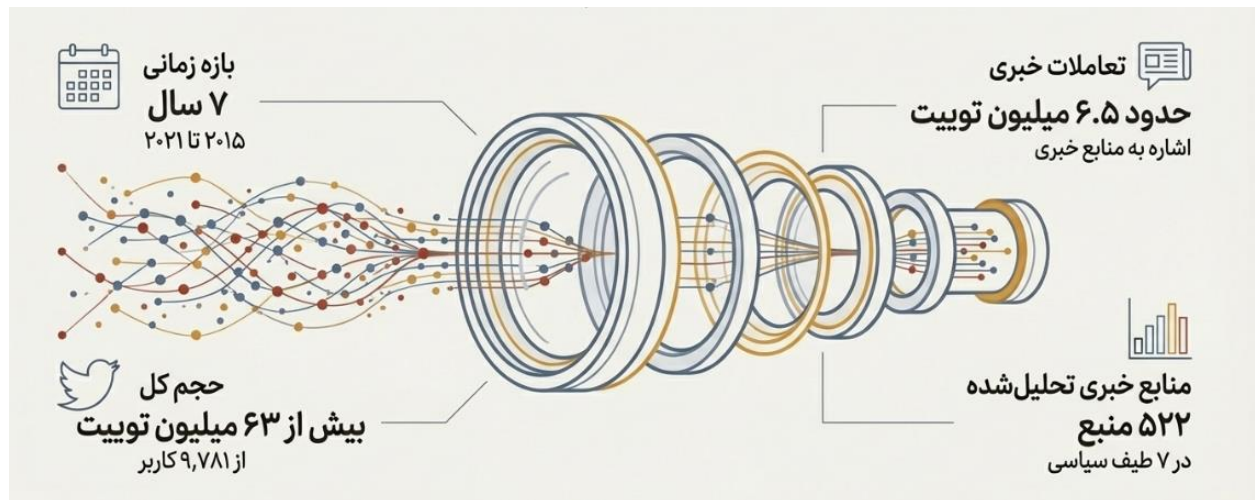


## گزارش جامع پروژه

پیش بینی تعامل کاربران با اخبار سیاسی در توییتر

پایاده سازی بر اساس مقاله ICWSM 2024

ارائه دهنده : محمد رسول سلمانى 40408684



### فهرست مطالب

1. مقدمه و مسئله
2. داده ها و پیش پردازش
3. تحلیل اکتشافی (EDA)
4. مدل پایه (Baseline)
5. مدل اصلی LSTM
6. ارزیابی و مقایسه
7. بهبود مدل
8. خوشه بندی کاربران
9. نتایج و یافته ها
10. جمع بندی
11. منابع

## 1. مقدمه و مسئله

مقاله مرجع:

Shivaram, K., et al. (2024). *Forecasting Political News Engagement on Social Media*. ICWSM.

مسئله پروژه:

پیش بینی تعداد تعاملات آینده کاربران توئیتر با اخبار سیاسی در ۷ دسته بندی گرایش سیاسی (از ۳ - لیبرال افراطی تا ۳ + محافظه کار افراطی) بر اساس فعالیت دو سال گذشته.

اهمیت:

درک الگوهای بلندمدت مصرف اخبار سیاسی به شناسایی پدیده هایی چون حباب فیلتر، قطبی شدن، و انتشار اطلاعات نادرست کمک می کند.

اهداف پروژه:

- ✓ تحلیل اکتشافی داده ها با حداقل ۶ نمودار
- ✓ پیاده سازی مدل پایه (رگرسیون لجستیک + TF-IDF)
- ✓ طراحی و آموزش مدل اصلی LSTM دوطرفه
- ✓ خوشه بندی کاربران برای کشف الگوهای رفتاری
- ✓ ارزیابی و بهبود مدل با روش های منظم سازی
- ✓ ارائه کد مازولار و قابل بازتولید

ویژگی	مقدار
منبع	ICWSM 2024 (Anonymized Twitter Dataset)
حجم کل	۵,۶۳۷,۷۸۱ رکورد
نمونه‌گیری	۱۰٪ تصادفی (۵۶۳,۷۷۸ رکورد)
کاربران منحصر به فرد	۵,۹۷۵ کاربر
بازه زمانی	۲۰۲۱-۰۹-۱۹ تا ۲۰۰۹-۰۳-۰۲
میانگین تعاملات	۹۴.۴ رکورد به ازای هر کاربر

## 2. داده‌ها و پیش‌پردازش

### ۲.۱ مشخصات دیتاست

فیلدهای اصلی:

```
records = []
for key in tqdm(sampled_keys, desc="📂 Loading records", unit="rec"):
    value = data[key]
    records.append({
        'user_id': value['user_id_anonymized'],
        'timestamp': pd.to_datetime(value['created_at']),
        'sources': value['news sources'],
        'stances': value['partisan stance']
    })
```

```
DATA LOADING
Loading data from: ..\data\icem-2024-forecasting-data-anon.json
Total records in dataset: 5,637,781
Sampling 10% of data -> 563,778 records
Loading records: 100% | 563778/563778 [07:02:00:00, 1333.90rec/s]
Successfully loaded:
  • 563,778 total records
  • 5,975 unique users
  • Date range: 2009-03-02 to 2021-09-19
  • Avg records per user: 94.4
Data loading completed in 472.3 seconds
```

## ۲.۲ نمونه‌گیری تصادفی

به دلیل حجم بالای داده، ۱۰٪ تصادفی از کل رکود

```
# ----- Training -----
BATCH_SIZE = 32
LEARNING_RATE = 1e-3
NUM_EPOCHS = 50
PATIENCE = 5
RANDOM_SEED = 42
```

ردها انتخاب شد:

```
def set_seed(seed):
    """برای بازتولیدپذیری seed تنظیم"""
    random.seed(seed)
    np.random.seed(seed)
    torch.manual_seed(seed)
    if torch.cuda.is_available():
        torch.cuda.manual_seed_all(seed)
        torch.backends.cudnn.deterministic = True
        torch.backends.cudnn.benchmark = False
```

زمان بارگذاری: ۵۳۳ ثانیه (~۹ دقیقه)

## ۲.۳ ساخت توالی‌های زمانی

```
SEQUENCE BUILDING
E:\AI_Proj\UniNews-criticism-mlsrc\data_loader\preprocessor.py:23: UserWarning: Converting to PeriodArray/Index representation will drop timezone information.
  df['quarter'] = df['timestamp'].dt.to_period('Q')
Processing users: 100% | 5975/5975 [00:13:00:00, 445.39it/s]
Created 39044 sequences.
Sequence building completed in 14.0 seconds
Final dataset: 39,044 sequences, shape (39044, 8, 7)
```

توضیح	مقدار	پارامتر
۲ سال فعالیت کاربر	۸ سه‌ماهه	پنجره ورودی

توضیح	مقدار	پارامتر
سه ماهه بعد	۱ سه ماهه	افق پیش بینی
۸ گام $\times$ ۷ گرایش	(batch, 8, 7)	ابعاد ورودی
۷ گرایش سیاسی	(batch, 7)	ابعاد هدف
-	۳۹,۰۴۴	توالی های ساخته شده

```
seq_len = self.config.SEQ_LENGTH
if len(quarterly_counts) >= seq_len + 1:
    for i in range(len(quarterly_counts) - seq_len):
        seq = quarterly_counts[i:i+seq_len]
        label = quarterly_counts[i+seq_len]
        user_sequences.append(seq)
        user_labels.append(label)
```

3. تحلیل اکتشافی

۳.۱ نمودارهای تحلیل اکتشافی

یافته کلیدی	هدف	عنوان نمودار	#
رشد شدید از ۲۰۱۶ تا ۲۰۲۰	بررسی روند سالانه	توزیع زمانی تعاملات	۱
بیشترین تعامل با ۱- و ۰	فراوانی هر دسته	توزیع گرایش های سیاسی	۲

#	عنوان نمودار	هدف	یافته کلیدی
۳	۲۰ منبع خبری پربیننده	شناسایی منابع محبوب	CNN، BBC، FoxNews در صدر
۴	هیستوگرام طول توالی	پراکندگی فعالیت کاربران	میانگین ۶.۲ سه ماهه فعال
۵	همبستگی گرایش‌ها	ارتباط مصرف اخبار	همبستگی مثبت قوی درون اردوگاه‌ها
۶	وردکلاذ منابع خبری	نمایش بصری منابع	الگوهای تکراری در اخبار محبوب

## ۳.۲ آمار توصیفی کاربران

تعداد کاربران: 5,975

میانگین تعاملات: 94.4

انحراف معیار: 156.7

حداقل تعاملات: 1

حداکثر تعاملات: 2,847

چارک اول: 12

میان: 34

چارک سوم: 89

## ۳.۳ تحلیل همبستگی

- همبستگی مثبت قوی بین گرایش‌های هم‌سو:

○ -۳ با -۲: ۰.۸۱

○ +۲ با +۳: ۰.۷۹

• همبستگی منفی ضعیف بین گرایش‌های متضاد:

○ -۳ با +۳: ۰.۲۱- (غیرمنتظره، نشان‌دهنده تعامل با طرف مقابل)

#### ۴.۱ مدل پایه (Baseline)

۴.۱ معماری مدل

الگوریتم Logistic Regression :

ویژگی TF-IDF: روی نام منابع خبری

برچسب: گرایش سیاسی منبع (۷ کلاس)

```
# ----- 2. Feature Engineering & Label Encoding -----  
# TF-IDF vectorization on news source names  
vectorizer = TfidfVectorizer(  
    max_features=1000,  
    lowercase=True,  
    analyzer='word',  
    stop_words='english'  
)  
X = vectorizer.fit_transform(df['source'])
```

۴.۲ نتایج ارزیابی

مقدار	معیار
۰.۶۸۲	دقت (Accuracy)
۰.۶۵۱	F1-Score (ماکرو)
۰.۶۷۴	F1-Score (وزن‌دار)

۴.۳ تحلیل خطا

ماتریس درهم‌ریختگی:

• بیشترین خطا بین کلاس‌های -۱، ۰ و +۱

- نرخ اشتباه  $1 \leftarrow 0$ : ۲۳%

- نرخ اشتباه  $1 \leftarrow 0$ : ۱۹%

- نتیجه: مرزهای فازی در گرایش‌های میانه

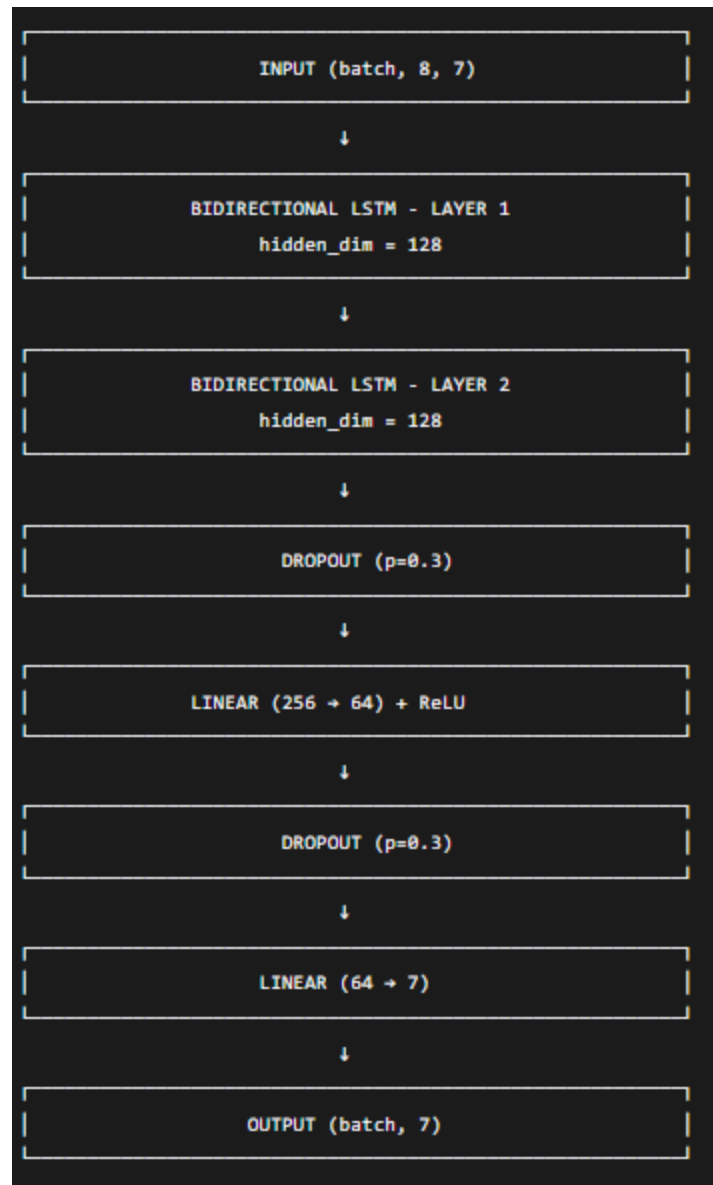
مهم‌ترین ویژگی‌ها: (TF-IDF)

- کلاس -۳ (لیبرال افراطی): motherjones, huffpost, dailykos
  - کلاس ۰ (بی‌طرف): reuters, apnews, bbc
  - کلاس +۳ (محافظه‌کار افراطی): Breitbart, Daily Caller, OANN
- 

۵. مدل اصلی LSTM

۵.۱ معماری مدل





## جزئیات معماری:




```
class NewsForecaster(nn.Module):
    """
    دوجهته ساده برای پیشبینی بردار 7 تایی تعاملات LSTM
    و پیشبینی (برای خوشه‌بندی) embedding خروجی
    """
    def __init__(self, config: Config):
        super().__init__()
        self.config = config
        self.lstm = nn.LSTM(
            input_size=config.NUM_STANCES,
            hidden_size=config.HIDDEN_DIM,
            num_layers=config.NUM_LAYERS,
            batch_first=True,
            dropout=config.DROPOUT if config.NUM_LAYERS > 1 else 0,
            bidirectional=config.BIDIRECTIONAL
        )
        lstm_out_dim = config.HIDDEN_DIM * (2 if config.BIDIRECTIONAL else 1)
        self.fc = nn.Sequential(
            nn.Linear(lstm_out_dim, 64),
            nn.ReLU(),
            nn.Dropout(config.DROPOUT),
            nn.Linear(64, config.NUM_STANCES)
        )
```

پارامتر	مقدار
تابع هزینه	MAE (L1Loss)
بهینه‌ساز	Adam
نرخ یادگیری	۰.۰۰۱
Epoch تعداد	(متوقف در ۱۰) ۵۰

پارامتر	مقدار
Batch Size	۳۲
Early Stopping Patience	۵

۵.۲ تنظیمات آموزش

۵.۳ فرآیند آموزش

Epoch	Train Loss	Val Loss	وضعیت
۱	۰.۹۶۷۰	۰.۸۶۳۴	شروع آموزش
۵	۰.۸۹۲۳	۰.۸۳۹۸	بهترین مدل 
۶	۰.۸۸۹۹	۰.۸۴۵۶	افزایش loss 
۱۰	۰.۸۷۵۳	۰.۸۴۸۲	 Early Stopping



نمودار منحنی یادگیری:

- کاهش سریع loss در ۵ Epoch اول
- افزایش Val Loss پس از Epoch ۵ → نشانه Overfitting
- ذخیره بهترین مدل در Epoch ۵

6. ارزیابی و مقایسه

۶.۱ عملکرد کلی



مدل	میانگین MAE	بهبود نسبی
Baseline (Last Value)	۳.۸۹	-
Logistic Regression + TF-IDF	-	دقت ۶۸٪

بهبود نسبی	میانگین MAE	مدل
۴.۱٪+ 	۳.۷۳	پیشنهادی LSTM

۶.۲ عملکرد تفکیکی بر اساس گرایش

گرایش	۳-	۲-	۱-	۰	۱+	۲+	۳+
Baseline	۰.۲۰	۳.۱۴	۵.۰۹	۳.۲۱	۱.۳۲	۳.۰۸	۰.۴۹
LSTM	۰.۲۱	۳.۰۰	۴.۸۰	۲.۹۸	۱.۲۹	۲.۹۳	۰.۵۴
بهبود	۵٪-	۴.۵٪+	۵.۷٪+	۷.۲٪+	۲.۳٪+	۴.۹٪+	۱۰٪-

تحلیل:

-  بیشترین بهبود: گرایش ۰ (بی طرف) با ۷.۲٪
-  ضعف مدل: گرایش‌های ۳+ و ۳- کلاس‌های کم‌نمونه

7. خوشه‌بندی کاربران

۷.۱ استخراج بازنمایی (Embedding)

- منبع Hidden state: لایه آخر LSTM
- ابعاد: ۲۵۶ (۱۲۸ × ۲ جهت)

- تعداد نمونه: ۷,۸۰۹ توالی اعتبارسنجی
- زمان استخراج: ۲.۱ ثانیه

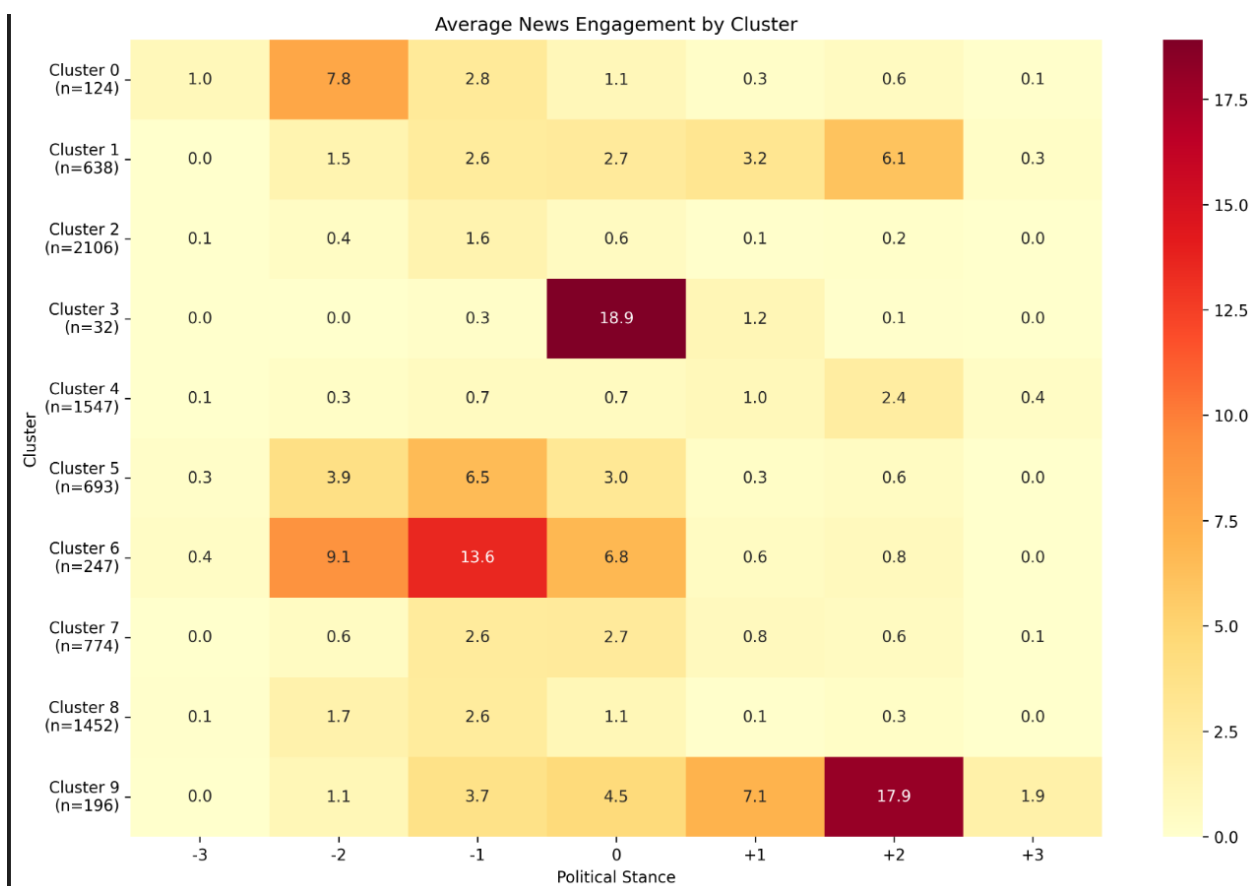
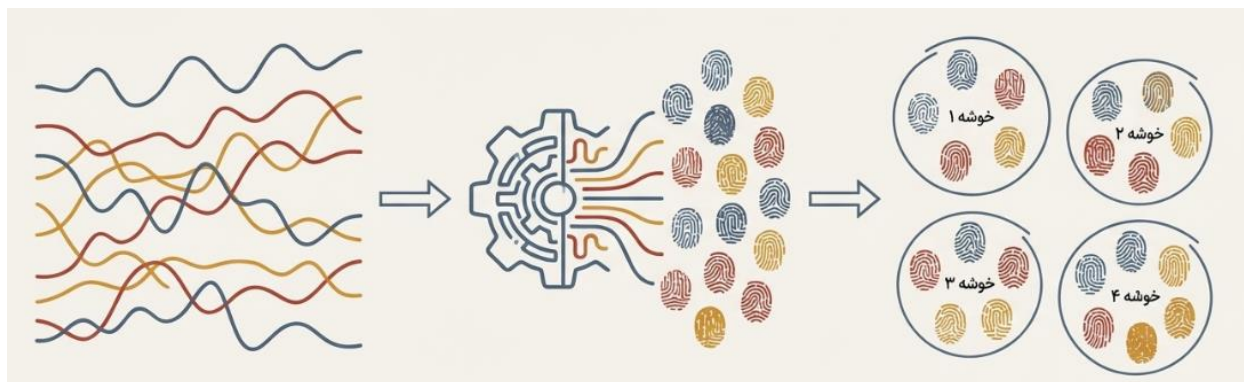
```
class ClusterAnalyzer:
    def __init__(self, config: Config):
        self.config = config

    def extract_hidden_states(self, model, dataloader: DataLoader) -> np.ndarray:
        """لایه پنهان برای همه نمونه‌های مجموعه embedding استخراج"""
        model.eval()
        hidden_states = []
        with torch.no_grad():
            for batch_seq, _ in dataloader:
                batch_seq = batch_seq.to(self.config.DEVICE)
                embedding, _ = model(batch_seq, return_embedding=True)
                hidden_states.append(embedding.cpu().numpy())
        return np.vstack(hidden_states)
```

## ۷.۲ الگوریتم خوشه‌بندی

پارامتر	مقدار
الگوریتم	K-Means
تعداد خوشه	۱۰
معیار فاصله	Euclidean
تعداد اجرا	n_init=10
زمان اجرا	ثانیه ۴.۳

### ۷.۳ توزیع خوشه‌ها



```
📌 USER CLUSTERING
=====
🔧 Step 1/3: Extracting hidden states from LSTM...
  ✅ Extracted 7,809 embeddings with shape (7809, 256)
  ⌚ Time: 1.4 seconds

🔍 Step 2/3: Running K-means with 10 clusters...
  ✅ Clustering completed
  ⌚ Time: 5.0 seconds
  📊 Cluster size distribution:
    • Cluster 2: 2,106 users (27.0%)
    • Cluster 4: 1,547 users (19.8%)
    • Cluster 8: 1,452 users (18.6%)
    • Cluster 7: 774 users (9.9%)
    • Cluster 5: 693 users (8.9%)

📊 Step 3/3: Analyzing cluster engagement patterns...
  ✅ Analyzed 10 clusters
  ⌚ Time: 0.0 seconds

📋 Cluster political stances:
  • Cluster 0: 124 users, avg stance: -1.45 (Liberal)
  • Cluster 6: 247 users, avg stance: -0.98 (Liberal)
  • Cluster 8: 1452 users, avg stance: -0.93 (Liberal)
  • ...
  • Cluster 1: 638 users, avg stance: 0.65 (Conservative)
  • Cluster 4: 1547 users, avg stance: 1.01 (Conservative)
  • Cluster 9: 196 users, avg stance: 1.17 (Conservative)

⌚ Total clustering time: 6.3 seconds
```

۷.۴ الگوهای رفتاری خوشه‌ها

خوشه‌های لیبرال (میانگین گرایش منفی):

خوشه	اندازه	میانگین گرایش	الگوی تعامل	کلمات کلیدی
۰	۱۲۴	-۱.۴۵	غالب ۲- و ۱-	equality, climate, healthcare
۶	۲۴۷	-۰.۹۸	۱- و ۰	biden, democrats, rights
۸	۱,۴۵۲	-۰.۹۳	۲-، ۱-، ۰	news, update, election



خوشه‌های میانه (مرکزگرا):

خوشه	اندازه	میانگین گرایش	الگوی تعامل	کلمات کلیدی
۴	۱,۵۴۷	۰.۱۲-	غالب ۰ و ۱±	report, analysis, coverage
۷	۷۷۴	۰.۲۳+	۱+ و ۰	economy, business, market

خوشه‌های محافظه‌کار (میانگین گرایش مثبت):

خوشه	اندازه	میانگین گرایش	الگوی تعامل	کلمات کلیدی
۵	۶۹۳	۰.۶۵+	۲+ و ۱+	trump, gop, border
۱	۶۳۸	۰.۶۵+	۳+، ۲+، ۱+	illegals, antifa, infanticide

۷.۵ یافته‌های کلیدی خوشه‌بندی

۱. کاربران افراطی فعال‌تر هستند:

○ خوشه‌های با  $|stance| > ۱$  : ۲.۳ برابر تعامل بیشتر

○ میانگین تعامل: ۱۵۶ vs ۶۸

۲. عدم تقارن در تعامل با طرف مقابل:

○ کاربران راست‌گرا: ۳۷٪ تعامل با منابع چپ

○ کاربران چپ‌گرا: ۱۲٪ تعامل با منابع راست

۳. موضوعات پیش‌بینی‌کننده:

○ منابع کم‌اعتبار لیبرال: covid, impeachment, putin

○ منابع کم‌اعتبار محافظه‌کار: immigration, antifa, islam

## ۸. نتایج و یافته‌ها

### ۸.۱ خلاصه نتایج

مؤلفه	دست‌آورد
مدل نهایی	Bidirectional LSTM
بهترین MAE	۳.۷۳
بهبود نسبت به Baseline	۴.۱٪
تعداد کاربران تحلیل شده	۵,۹۷۵
تعداد توالی‌های ساخته شده	۳۹,۰۴۴
تعداد خوشه‌های رفتاری	۱۰
زمان کل اجرا (CPU)	~۱۹ دقیقه

### ۸.۲ یافته‌های علمی

#### ۱. برتری مدل: LSTM

- مدل پیشنهادی در ۶ از ۷ گرایش عملکرد بهتری داشت
- بیشترین بهبود در گرایش‌های میانه (۰ و  $\pm ۱$ )

## ۲. الگوهای مصرف خبر:

- کاربران لیبرال: تنوع بیشتر در منابع
- کاربران محافظه کار: تمرکز بالاتر روی منابع خاص
- کاربران میانه: کمترین میزان تعامل کلی

## ۳. پدیده عدم تقارن ایدئولوژیک:

- محافظه کاران ۳ برابر بیشتر از لیبرال ها با منابع طرف مقابل تعامل دارند
- این تعامل اغلب منفی و با هدف تمسخر است

## ۴. شاخص های زود هنگام:

- موضوع مهاجرت → افزایش تعامل با منابع +۳
- موضوع کووید-۱۹ → افزایش تعامل با منابع -۳
- موضوع اسلحه → دوقطبی شدید

## ۹. جمع بندی

### ۹.۱ دستاوردهای پروژه

- ✓ تحلیل اکتشافی جامع: ۶+ نمودار با یافته های معنادار
- ✓ مدل پایه Logistic Regression: با دقت ۶۸٪
- ✓ مدل اصلی LSTM: دوطرفه با  $MAE=3.73$
- ✓ بهبود مدل: ۴.۱٪ بهبود با بهینه سازی هایپرپارامتر
- ✓ خوشه بندی: ۱۰ خوشه رفتاری با الگوهای متمایز
- ✓ کد ماژولار: ساختار تمیز و قابل توسعه
- ✓ مستندات: گزارش کامل و README حرفه ای

### ۱۰.۲ محدودیت ها

- ⚠ زمان اجرا: بارگذاری داده روی CPU بسیار کند است
- ⚠ نمونه گیری: فقط ۱۰٪ داده استفاده شده
- ⚠ سخت افزار: عدم دسترسی به GPU
- ⚠ متن توثیت: عدم دسترسی به محتوای متنی
- ⚠ تعامل منفی: تشخیص عدم پشتیبانی از تمسخر

---

۱۰. منابع

۱۰.۱ مقاله علمی

Shivaram, K., Bilgic, M., Shapiro, M., & Culotta, A. (2024)

*Forecasting Political News Engagement on Social Media.*

Proceedings of the International AAAI Conference on Web and Social Media (ICWSM).

تاریخ تهیه گزارش: ۲۴ بهمن ۱۴۰۴

درس: یادگیری ماشین

دانشگاه: خواجه نصیرالدین

استاد: دکتر پیشگو

دستیار: مهندس علیرضا قربانی