

عنوان : Predicting Emerging Trends: A Machine Learning Approach to Topic Popularity on Social Media
نویسنده‌ان : Zhe Wu et al.
محله : PeerJ Computer Science (2025)

1. مسئله اصلی مقاله و اهمیت آن

رشد سریع شبکه‌های اجتماعی باعث تولید روزانه حجم عظیمی از محتوا و موضوعات (Topics) شده است. تنها بخش کوچکی از این موضوعات به «ترند» یا موضوعات محبوب تبدیل می‌شوند، در حالی که بسیاری دیگر به سرعت فراموش می‌شوند. مسئله اصلی مقاله این است:

آیا می‌توان در مراحل اولیه شکل‌گیری یک موضوع در شبکه‌های اجتماعی، پیش‌بینی کرد که آیا آن موضوع در آینده محبوب (Popular) خواهد شد یا خیر؟

اهمیت این مسئله از چند جنبه قابل توجه است:

- کنترل و شناسایی زودهنگام اطلاعات نادرست و شایعات که می‌توانند افکار عمومی را تحت تأثیر قرار دهند.
- کاربردهای تجاری و بازاریابی مانند زمان‌بندی کمپین‌ها و تشخیص موضوعات داغ بازار.
- تحلیل افکار عمومی و سیاست‌گذاری در بحران‌ها یا رویدادهای اجتماعی.

برخلاف بسیاری از کارهای قبلی که به مدل‌های سری زمانی متکی‌اند (و نیازمند داده‌های بلندمدت هستند)، این مقاله تلاش می‌کند محبوبیت نهایی یک موضوع را تنها با استفاده از ویژگی‌های اولیه آن (مثلاً ۱۲ ساعت اول) پیش‌بینی کند. این تغییر دیدگاه، مسئله را از «پیش‌بینی روند» به «مسئله طبقه‌بندی زودهنگام» تبدیل می‌کند که از نظر مهندسی و عملیاتی بسیار کاربردی‌تر است.

2. ورودی‌ها و خروجی‌های مدل / سیستم

ورودی‌ها

ورودی مدل، مجموعه‌ای از ویژگی‌های استخراج شده در مراحل اولیه انتشار هر Topic است. این ویژگی‌ها در چهار دسته کلی قرار می‌گیرند:

1. ویژگی‌های آماری پایه (طول متن، تعداد لینک، ایموجی و ...)

2. ویژگی‌های متنی و احساسی (Sentiment)، نوع موضوع (Topic)

3. ویژگی‌های زمانی (زمان انتشار، روز هفته، تعطیلات)

4. ویژگی‌های انتشار اولیه (تعامل کاربران در ۱۲ ساعت اول)

در مجموع، ۳۸ ویژگی برای هر Topic تعریف شده است.

خروجی‌ها

خروجی سیستم یک برچسب دودویی (Binary Classification) است:

• ۱محبوب (Popular) → Topic

• ۰غیرمحبوب (Non-popular) → Topic

معیار محبوبیت بر اساس یک امتیاز نهایی محاسبه شده از تعداد پست‌ها، لایک‌ها، کامنت‌ها و اشتراک‌گذاری‌ها تعیین می‌شود. اگر این امتیاز از یک آستانه مشخص (مثلاً ۲۵۰۰) عبور کند، Topic محبوب تلقی می‌شود.

3.داده مورد استفاده (نوع، منبع، اندازه)

نوع داده

• داده‌های متنی پست‌های Facebook

• داده‌های تعاملی کاربران (Likes, Comments, Shares)

• داده‌های زمانی انتشار

منبع داده

• داده‌ها از مجموعه داده عمومی Facebook استخراج شده‌اند.

• منبع خام داده در GitHub و Zenodo منتشر شده است.

اندازه داده

• حدود ۳۵.۳۱ میلیون پست Facebook

• بازه زمانی: فوریه ۲۰۲۱ تا ژانویه ۲۰۲۲

• حجم داده خام: حدود ۴۶.۵ گیگابایت

- تعداد Topic های استخراج شده اولیه: ۱۰,۹۳۲
- پس از پاکسازی دادهها: Topic ۳,۵۷۰ معتبر
- توزیع کلاس‌ها:

 - Topic ۱۵۷ محبوب ○
 - Topic ۳,۴۱۳ غیرمحبوب ○

عدم توازن شدید دادهها (۱ به ۲۲) یکی از چالش‌های اصلی مقاله است.

4. روش پیشنهادی مقاله (به زبان ساده)

گام 1: استخراج Topic ها

- پست‌های روزانه Facebook با استفاده از TF-IDF تمايش عددی می‌شوند.
- الگوریتم K-Means برای خوشه‌بندی متن‌ها استفاده می‌شود.
- هر خوشه به عنوان یک Topic در نظر گرفته می‌شود.
- های مشابه در روزهای مختلف با Cosine Similarity ادغام می‌شوند.

گام 2: محاسبه محبوبیت Topic

محبوبیت نهایی Topic به صورت زیر محاسبه می‌شود:

$$P = 0.46 \times \text{Posts} + 0.27 \times \text{Shares} + 0.17 \times \text{Comments} + 0.10 \times \text{Likes}$$

این وزن‌ها با ترکیب روش AHP و نظر خبرگان تعیین شده‌اند.

گام 3: استخراج ویژگی‌ها (۱۲ ساعت اول)

- استخراج ۳۸ ویژگی از متن، زمان و تعامل کاربران در ۱۲ ساعت اول انتشار Topic
- تحلیل اهمیت ویژگی‌ها نشان می‌دهد که:

 - Shares و Comments Likes اولیه بیشترین نقش را دارند.

گام 4: آماده‌سازی داده

- حذف Topic های کوتاه‌عمر یا ناقص

رفع عدم توازن داده با:

SMOTE (Oversampling) ○

Class Weighting ○

گام 5: مدل سازی و یادگیری

مدل های زیر آموزش و مقایسه شده اند:

Random Forest •

XGBoost •

LightGBM •

CatBoost •

SVM •

شبه کد کلی روش

For each day:

Collect Facebook posts

Preprocess text

Cluster posts → Topics

For each Topic:

Extract first 12h features

Compute final popularity score

Assign label (Popular / Non-popular)

Balance dataset (SMOTE)

Train ML classifiers

Evaluate using Accuracy, Precision, Recall, F1, AUC

5. نتایج اصلی، محدودیت‌ها و ایده‌های ادامه

نتایج اصلی

- دقت کلی مدل‌ها بالاتر از ۹۰٪ است.

XGBoost و Random Forest عملکرد متعادلی در Precision و Recall دارند.

ویژگی‌های انتشار اولیه مهم‌ترین عامل پیش‌بینی هستند.

روش پیشنهادی برای پیش‌بینی زودهنگام بسیار مناسب است.

محدودیت‌ها

استفاده از داده‌های یک پلتفرم \rightarrow (Facebook) تعمیم‌پذیری محدود

عدم استفاده از مدل‌های عمیق مانند LSTM یا Transformer

وابستگی به آستانه محبوبیت

احتمال overfitting ناشی از SMOTE

ایده‌های ادامه تحقیق

استفاده از Transformer‌ها (BERT) برای نمایش متنی قوی‌تر

تحلیل کاربران اولیه منتشرکننده Topic

یادگیری وزن‌های محبوبیت به صورت داده‌محور

گسترش مدل به چند پلتفرم (...), Instagram, Twitter

ترکیب تشخیص ناهنجاری (Anomaly Detection) با طبقه‌بندی