

Predicting User Engagement on Social Media: A Comparative Study of Machine Learning-Based Modeling Approaches

Nagendra Singh¹, Manish Dixit¹, Saurabh Pachauri^{1,*}, Shobhit M. Sharma¹, Anil Singh¹

Abstract

Even while mobile social apps are becoming increasingly important in people's daily lives, little is known about the elements that motivate users to interact with these apps. The goal of this work is to analyze and define the transition patterns of individual users' in-app actions as a temporally changing action graph. According to our investigation, action graphs are a useful tool for characterizing user behavior patterns and offering ideas for further interaction. To record in-app usage patterns, extract multiple high-order graph elements. To further improve prediction power by developing an end-to-end, multichannel neural model that encodes activity sequences, temporal action graphs, and other macroscopic features. Machine learning (ML) has emerged as a powerful technique in several predictive analytics domains, including social media data analysis. Social media platforms have grown into massive user-generated content databases that offer valuable insights into the interests, behaviors, and trends of their users. To apply machine learning algorithms to anticipate user behavior based on social media data and detect important trends. The study makes use of a substantial dataset that comprises user profiles, blog posts, comments, and engagement metrics acquired from well-known social networking sites. Predictive models are developed using a variety of machine learning algorithms, including ensemble approaches, neural networks, decision trees, and support vector machines. Future work in this area must concentrate on resolving privacy and data quality concerns related to social media data to improve machine learning's prediction powers.

Keywords: Graph convolutional network, predictive analytics, support vector machines, long short-term memory, receiver operating characteristic

INTRODUCTION

There has been a noticeable upsurge in interest lately in understanding the elements that maintain user engagement with online programs (apps), as well as how people adopt these apps on mobile and internet platforms [1–10]. Even though a lot of research focuses on user profiling [11–21] and creating sophisticated models to show users better tailored content [17], there is still a knowledge vacuum about the crucial elements that influence user engagement with apps. Asking this question is crucial to keeping users interested and engaged.

*Author for Correspondence

Saurabh Pachauri
E-mail: saurabhme.pachauri@gmail.com

¹Assistant Professor, Department of Mechanical Engineering, Institute of Engineering and Technology, Khandari Campus, Agra, Uttar Pradesh, India

Received Date: September 24, 2025

Accepted Date: November 06, 2025

Published Date: December 24, 2025

Citation: Nagendra Singh, Manish Dixit, Saurabh Pachauri, Shobhit M. Sharma, Anil Singh. Predicting User Engagement on Social Media: A Comparative Study of Machine Learning-Based Modeling Approaches. International Journal of Computer Aided Manufacturing. 2025; 11(2): 41–61p.

The user's activities while utilizing the app are represented in Figure 1 by nodes in the graph, and the chance of switching between actions throughout a session (such as opening and shutting the app) is indicated by edges. Therefore, it is critical to understand, classify, and relate user behaviors to user engagement, for example, by predicting future engagement trends. In earlier attempts to analyze user interaction, the focus has been on extracting macro-level features from

activity metrics [1, 16, 21]. This study investigates the possibility of predicting users' future app engagement, such as active days over a given period, based on their in-app activity patterns. Previous work on user app engagement forecasting has mostly focused on various macroscopic indicators such as time-series data on activity frequency. A large amount of current research focuses on forecasting whether a user will return to the platform using different methodologies to uncover important characteristics of user behaviors. On social media platforms, predicting return rates is a prominent and important issue as businesses actively seek to understand what makes them successful. Some research focuses on predicting user retention or churn rates, which indicate a stop in participation on platform [2, 14]. Some tackle the problem by projecting when a user will return [13], if they will return [16], or how long a new user will live [22–30].

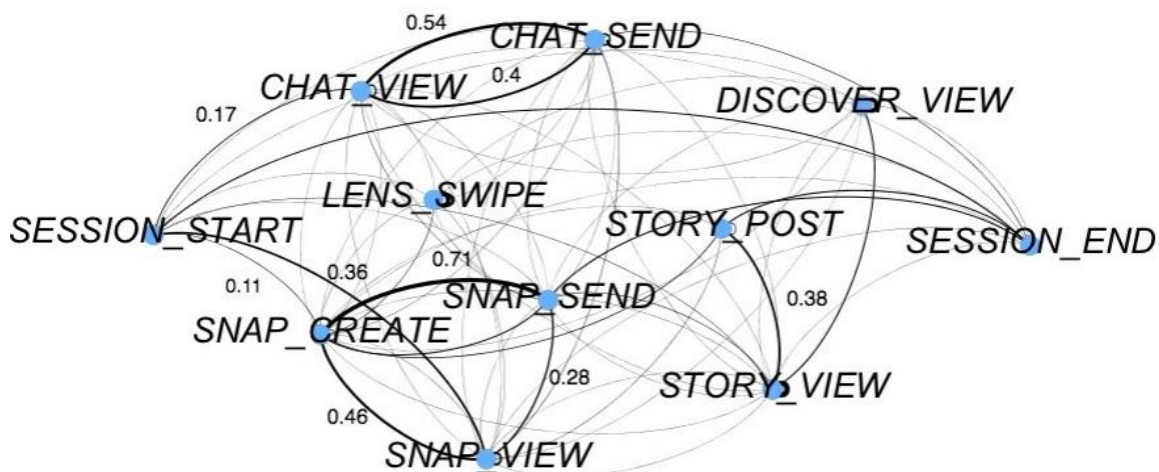


Figure 1. This is an example of an action graph made with a particular Snapchat user's in-app activity data [3].

In Figure 2, building on earlier studies, we present a novel method for capturing the transition patterns of individual users' in-app actions using action graphs, a weighted, directed graph [31–34]. Our first method introduces a feature-based model for interpretable user engagement prediction. This approach combines basic elements, such as explainable graph features and macroscopic user-level qualities, which are accomplished in two steps. Our deep neural network model uses the feature-based model as a point of reference, and it provides valuable insights and explanations for a more complex graph network [35].

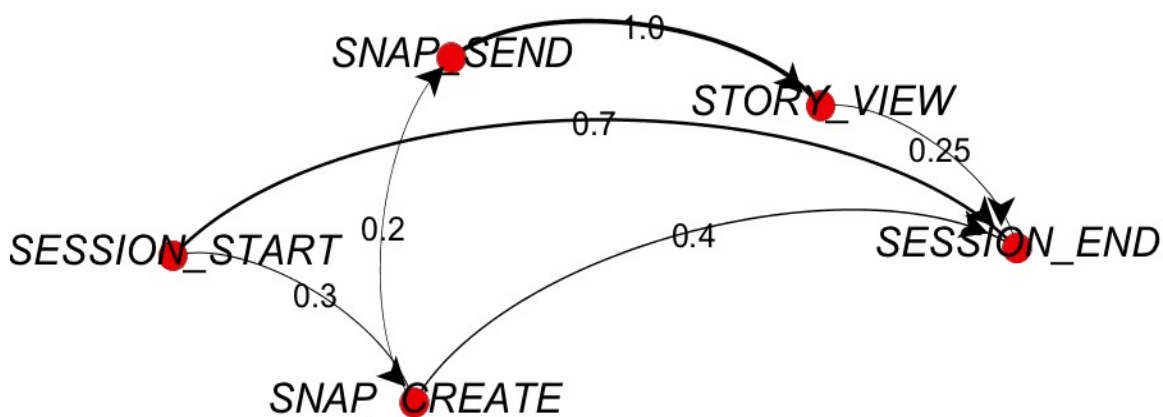


Figure 2. In contrast, an action graph has fewer actions and sparser transitions between actions [6].

To introduce many models to improve the prediction of user involvement. These include a static graph model that uses Graph Convolutional Network (GCN), an activity sequence model that uses

Long Short-Term Memory (LSTM) as a baseline [21], and a temporal action graph model that combines GCN and LSTM. Of these, the temporal graph model performs better than the prior baseline and other models in terms of capturing time-related dependencies in action graphs. A deep multichannel end-to-end model for enhanced performance that combines the activity sequence model, the temporal action graph model, and macroscopic aspects. Optimized performance can be achieved by training all these models together in an end-to-end arrangement, which allows for a more thorough knowledge of their combined effects. A new data model called the action graph, which is used to forecast user involvement and is intended to characterize user in-app behaviors. To investigate through our analysis how different user activity-related signals can help forecast a user's future engagement. By embedding high-order information modelling into action graphs [15]. A multichannel end-to-end forecasting system that integrates other valuable signals using a Long Short-Term Memory (LSTM) model trained on temporal action graphs. According to extensive research on both static and temporal action graphs. Our ablation study shows that easily interpreted graph characteristics are a good way to predict user engagement. Here is a screenshot of the Snapchat app's main interface, which shows off its ten main features.

The first step in the research process is data preparation, which involves tasks like cleaning, normalization, and feature extraction. Raw social media data is transformed into a machine learning algorithm-readable format using feature engineering techniques. Subsets of the dataset are then created for training and testing to accurately evaluate the performance of the models. Several machine learning techniques are applied after the dataset has been preprocessed. Decision trees are helpful for deriving rules from data because they generate models that are simple to understand [22]. Neural networks are better at capturing complicated linkages and nonlinear patterns than support vector machines are at classification. To increase prediction accuracy, ensemble techniques like gradient boosting and random forests use many models. Performance metrics, including F1-score, recall, accuracy, and precision, are used to evaluate prediction models. Additionally, the Receiver Operating Characteristic (ROC) and Area Under the Curve (AUC) curves are used to assess how well the models perform in classification tasks. These criteria are used to assess the models and identify the most effective method for using social media data for predictive analytics. The paper also addresses some challenges this sector encounters. Due to the noisy and unstructured nature of social media data, robust pretreatment techniques are essential for managing missing values, outliers, and irrelevant information. The study also covers the ethical concerns surrounding data privacy and user consent, as well as the ongoing need for model improvements due to social media platforms' dynamic nature [19]. The study also demonstrates the potential application of machine learning to predictive analytics using data from social media. The results demonstrate that ensemble approaches, and random forests in particular, outperform alternative algorithms in user behavior prediction.

This research has implications for sentiment analysis, targeted advertising, personalized suggestions, social network analysis, and more. Through the application of machine learning techniques to leverage the rich insights obtained from social media data, organizations can optimize marketing tactics, boost customer engagement, and enhance decision-making processes. It can be difficult to extract useful insights from this massive amount of data, though. This is where machine learning approaches applied to predictive analytics are useful [13]. Machine learning, a subfield of artificial intelligence, offers methods and tools for comprehending and analyzing vast amounts of data. Important patterns and trends are found in social media data by using statistical models and algorithms. This feature has completely changed how companies and organizations use predictive analytics in the context of social media. To predict future events or actions, predictive analytics makes use of both historical and present data. Predictive analytics looks for buried correlations and patterns in the data to produce precise forecasts. Its use in the social media space may be quite beneficial for businesses, helping them to better understand consumer behavior, predict trends, find influencers, and improve marketing tactics. Another significant application is the prediction of user behavior, wherein machine learning algorithms scan social media data to identify patterns in user behavior such as engagement levels, favorite content, and frequency of posts. By identifying these trends, businesses

may enhance their social media strategy, target certain user demographics, and personalize their content to boost user engagement and motivate desired actions. In addition, relationships and hidden groupings inside social networks can be found through social media network research employing machine learning techniques [22]. Businesses can refine their targeting strategies and boost the efficacy of their social media campaigns by looking into user interactions, finding influencers, and identifying interest groups. However, even if machine learning has great promises for predictive analytics in social media data, there are still problems that need to be fixed. The dynamic nature of social media networks and issues with data privacy and quality are major barriers. Additionally, evaluating and comprehending machine learning models that are applied to social media data may be challenging. A delicate balance between ethical issues, legislative frameworks, and technological advancements must be achieved in order to effectively handle these challenges.

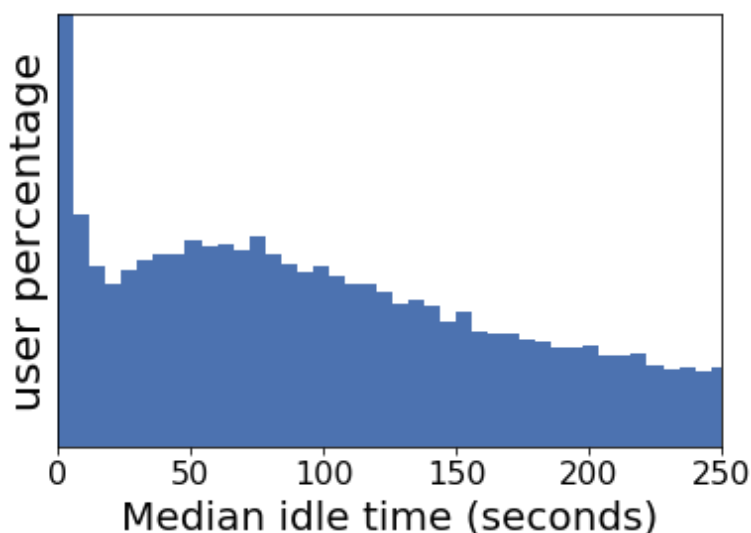
LITERATURE SURVEY

The various machine learning techniques used to analyze sentiment in social media data are examined in this review of the literature. It assesses multiple algorithms, showing the advantages and disadvantages of each in terms of sentiment extraction from user-generated content. Examples of this include support vector machines, ensemble methods, and recurrent neural networks [23, 24]. This paper provides an overview of deep learning techniques used in predictive analytics for social media data. It covers the use of attention processes, recurrent neural networks, and convolutional neural networks for various applications such as event detection, trend prediction, and user behavior tracking [25, 26]. This review looks at machine learning algorithms that are used for social media user profiling. It examines several techniques, like topic modeling, clustering, and classification, with the aim of extracting meaningful data from user-generated content for the purpose of creating user profiles. These profiles are then utilized by recommendation engines, targeted advertising, and tailored services [27, 28]. This work investigates sentiment analysis and opinion mining in social media data from a machine learning standpoint. It explores further into machine learning techniques for sentiment and opinion mining and includes topics including feature selection, sentiment lexicons, and model evaluation. Additionally, the study clarifies the difficulties and encouraging prospects in this field of study [29–30]. This research investigates the application of machine learning techniques to detect fraudulent information in social media data. It investigates several methods for identifying and stopping the spread of false information such as feature engineering, graph analysis, and deep learning [31, 32]. Another research effort investigates the machine learning methods used in recommender systems on social media platforms. This review of the literature investigates machine learning techniques for social network analysis in social media data, with a focus on user preferences and social interactions. It evaluates collaborative filtering, content-based filtering, and hybrid ways to recommend relevant products, people, and businesses [33, 34].

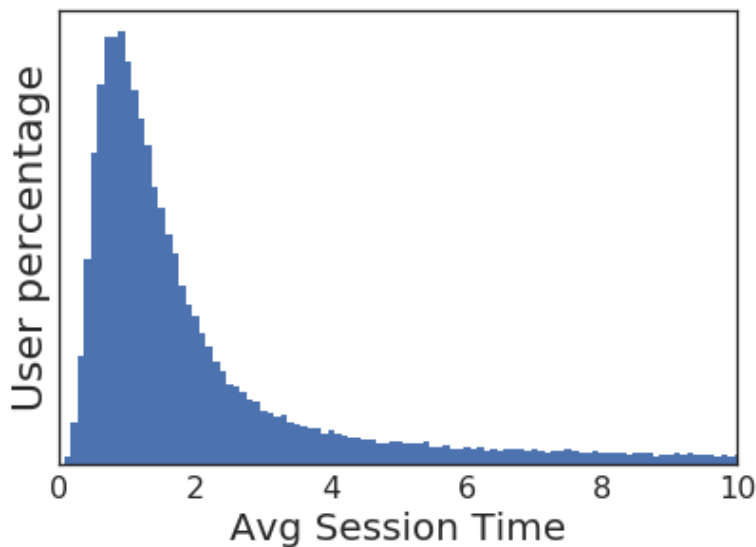
The paper also looks at machine learning methods that use data from social media platforms to analyze social networks. It examines techniques for identifying communities, identifying relationships, and evaluating impact, demonstrating how machine learning algorithms can extract valuable data from user networks and the connections that bind them [35, 36]. This work explores the use of machine learning algorithms for event recognition in social media data. It assesses various methods, including supervised and unsupervised techniques, for automatically detecting and tracking events from user-generated content. Real-time monitoring and situational awareness are made feasible by this [37–47]. This review of the literature examines machine learning techniques for analyzing user behavior in data from social media platforms in an alternative context. This study provides information about how machine learning can help comprehend the behaviors and preferences of users. It includes methods for predicting user involvement, identifying influential users, and segmenting users [39, 40, 46]. This study also investigates the use of machine learning techniques to predict user behavior in social media data. Predictive analytics can improve user engagement on social media platforms and marketing methods, as the article highlights. It examines elements including content qualities, network architecture, and user demographics.

ACTION GRAPH FOR CHARACTERIZING USER ENGAGEMENT

This section explains the action graph concept, introduces the dataset we used from Snapchat for our study, and goes into detail on how we created these action graphs using user activity data. A sample of this information represented more than 25 million new users in a given nation throughout the given time frame. After a user registers, we collect activity data for 4 weeks, concentrating on 10 important in-app features (shown in Figure 3) that we thought were very important to represent user behavior on Snapchat [46]. We used data from the first 2 weeks for our study and to create user action graphs. The data from the final 2 weeks was set aside to assess projections about users' future involvement. An in-app session should ideally consist of a sequence of continuous operations with distinct breaks for the user to disconnect. Signals signaling the start of interaction (e.g., App Open) might originate from both the initial app opening and the app returning from background mode in real-world mobile usage scenarios. Disengagement can be ambiguous, in fact, because in the aforementioned instances, users frequently pick up where they left off and begin immediately thereafter [47]. For session splitting, it would be insufficient to accurately capture the flow of user actions and their intent for using the app by relying solely on recorded disengagement signals such as App Open and App Close. Consequently, a more thorough explanation of in-app sessions is required. As an alternative to merely employing heuristic methods to segment sessions, previous work has explored session identification by fitting interactivity time into various distributions [11]. A combination of the long-tail and normal distributions characterizes the distribution of the median idle time per user between the start and end engagement signals in Figure 3(a) because average idle time is prone to severe outliers, median idle time was chosen. It is noteworthy to note that the 25-second point corresponds with the lower 10th percentile of the idle time distribution, as well as the intersection of two distributions. As a result, we set a threshold idle time of 25 seconds to divide sessions along with the disengagement signals (App Open/Close) in our tests. Figure 3(b) illustrates the distribution of time spent on Snapchat inside each session using our definition of in-app sessions. We provide a novel concept for modeling user behavior in this study [18]. An action graph, which is a directed graph with ten in-app actions, has two nodes representing Session Start and Session End (Figure 4). Like a Markov chain model, the edges of the action graph indicate the probability of transitions between action nodes. Only outgoing and incoming edges are present on the nodes that symbolize the beginning and conclusion of a session, called Session Start and Session conclusion, respectively. For the duration of the monitored period, every user has their own action graph. Every user has a unique action graph created at every time interval for temporal graph modelling. Action graph samples are shown in the above graphics; Figure 1 shows a more engaged user, and Figure 2 shows a less engaged user.



(a) Median idle time distribution.



(b) Average session time.

Figure 3. (a) The distribution of the user's idle time between sessions and the average amount of time spent in each session are shown in Graph (b). Notably, both y-axes in graph (b) are hidden to prevent the presentation of absolute values, and the x-axis has been rescaled [10].

Only keep users who have five or more valid sessions left over for our study to make sure the transition probabilities in the action graphs are more meaningful for everyone [29]. After filtering, we get about 150,000 sampled users for our research. Users participate in about seven sessions a day on average, and throughout the first 2 weeks of the observation period, they participate in 98 sessions. More specifically, each action graph is created using the activity data from all the legitimate sessions that a user has participated in. This probability is aggregated over all user sessions to find the transitional probability between two activities in an action graph. The maximum number of nodes in an action graph is 12, which includes the session start and session end nodes.

Proposed System

The field of predictive analytics in social media data is mostly dependent on machine learning algorithms, which use patterns found in the data to forecast future events. These algorithms continuously learn from fresh data and user feedback to increase their accuracy and efficiency over time. Their predictive powers, which quickly analyze enormous volumes of social media data, are useful for guiding strategy creation and decision-making processes (Figure 4).

This is examining textual data from social media sites to ascertain the tone or emotional background of the information. Sentiment analysis automatically classifies social media postings or comments as neutral, negative, or positive using machine learning algorithms. Understanding consumer sentiment, identifying brand evangelists and skeptics, and improving marketing tactics can all be greatly aided by this data [37]. Through numerous interactions, including posts, comments, likes, and shares, social media platforms have completely changed how people communicate and exchange information. This has resulted in the generation of vast volumes of data. However, extracting valuable knowledge from this massive amount of data can be difficult. To overcome this challenge, the proposed system applies predictive analytics using machine learning algorithms on social media data.

System Architecture and Expected Benefits

In this stage, methods including text tokenization, feature scaling, and feature encoding are used. Predictive analytics jobs are carried out using multiple machine learning techniques [17]. For sentiment analysis tasks, for example, algorithms, such as Naive Bayes, Support Vector Machines

(SVM), or Recurrent Neural Networks (RNNs), might be used. Trends can be predicted by using time series forecasting methods like ARIMA or Long Short-Term Memory (LSTM) networks. Techniques for examining user behavior include clustering techniques such as DBSCAN and K-means. Predictive analytics data are presented and communicated using an intuitive interface. Users can more easily comprehend patterns, trends, and sentiment distributions in social media data by using visualizations like word clouds, graphs, and charts. Performance indicators, such as F1-score, recall, accuracy, and precision, are used to evaluate how effective the suggested method is. To find any weaknesses in the system and implement the required fixes, user input and comparisons with current methods are utilized.



Figure 4. Domains of predictive analysis [13].

Predictive analytics on social media data can help businesses make wise judgments. This entails seeing new trends, figuring out what customers want, and modifying marketing plans accordingly. Businesses may interact with customers more successfully by using sentiment analysis and an understanding of user behavior [25]. Businesses can adjust their offerings and messages to meet client expectations by analyzing data from social media. The suggested method gives businesses instant access to social media trends and conversations, allowing them to quickly respond to questions, complaints, and emerging issues from customers. Businesses can obtain a competitive edge in the market by utilizing predictive analytics' capabilities. This entails spotting unexplored prospects, gauging the mood of the market, and modifying plans as necessary [41]. The suggested approach uses predictive analytics and machine learning to extract insightful information from social media data. Businesses that gather, preprocess, and use sophisticated algorithms can improve their comprehension of customer preferences, forecast trends, and make wise judgments. Technology can greatly improve marketing efforts, increase customer interaction, and provide businesses with a competitive edge in the ever-changing social media landscape.

LARGE-SCALE DATA ANALYSIS

Action Type

Figure 5 shows how users interact with Snapchat’s fundamental capabilities in three main categories like broadcasting, narrowcasting, and content creation. Narrowcasting actions, which include Chat Send, Chat View, Snap Send, and Snap View, are one-to-one conversations for direct engagement with users and their friends. Consuming and disseminating content to one’s whole network are examples of broadcasting behaviors. Understanding user involvement and content consumption with the intention of increasing or maintaining the level of engagement and consumption is the core objective of most of the social media research [12]. A highly engaged user should ideally switch between narrowcast and broadcast activities, as well as between creation and consumption activities, more frequently. It is better to move from narrowcast to broadcast actions since watching broadcast content brings in advertising while consumers use the app to chat with friends. Analyzing changes in behavior from broadcasting to narrowcasting might indicate that a user is very involved in their social network and is more likely to continue participating. According to our analysis, almost one-third of the time a session moves on to broadcasting duties after beginning with narrowcasting tasks. On the other hand, about 25% of the time, a session moves from broadcasting chores to narrowcasting actions at the start. It is possible that users only complete the activities they set out to do, such as communicate or consume content, but this does not imply that they are not actively participating [18]. Here, the main goal is to encourage less engaged users to explore the app’s cross-functionality to increase their level of engagement. According to Figure 5, “Chat View” and “Snap Create” are the two activities that encourage involvement the most frequently. To prevent absolute values from being displayed, the Y-axis is concealed.

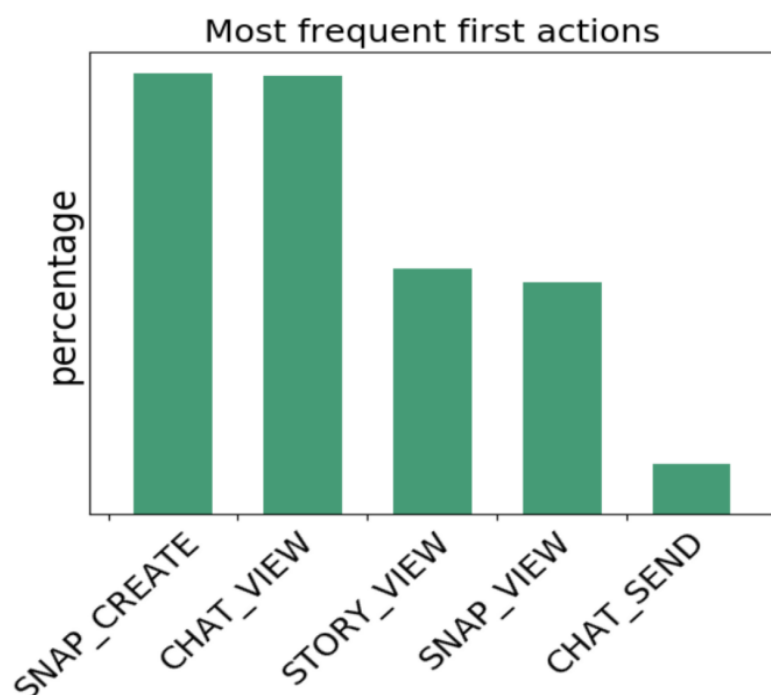


Figure 5. First things that happen most often after a session starts [17].

To gain further understanding of what drives users to engage with Snapchat, we look at the first action that takes place in each session following Session Start. Remarkably, we discover that the most favored initial actions for users are “Chat View” (a narrowcasting action) and “Snap Create” (a content generation action), as illustrated in Figure 5. We then investigate how long users spend on the app during each session after narrowcasting actions, such as “Snap View” and “Story View,” which are broadcasting-consuming activities [46]. Figure 4(b) displays the average session duration per user, which has a normal distribution with peaks occurring within a few minutes.

Session Gene

To gain a better understanding of session sequences and to establish stronger connections with users' intentions while using the app (i.e., content generation, narrowcasting, and broadcasting), we segment sessions into session genes using the soft clustering technique Latent Dirichlet Allocation (LDA) [6]. We regard every action as a word and every session as a document, much like topic modelling. The subjects found using LDA are referred to as "session genes," and they are fundamental elements that provide a more thorough description of a session's functions (Figure 6). Using LDA, we treated every session sequence as a document and every action as a word in our dataset. The study yielded five interpretable topics, which are determined by looking at the terms that have the highest weights within each topic [19]. These subjects are combined to form each session; these are known as "session genes." Chat, Snap, Story View, Discover View, and Content Creation are the five genes that have been found. Figure 6 shows the sequence of events for each gene. The fact that the content creation gene covers a wide range of activities raises the possibility that users are playing with the app and producing content that they may not share. Our action graphs benefit from the accurate description of each session's makeup provided by session genes. Among the five genes, communication-related genes, such as "Chat" and "Snap", are the most common and probable, as shown in Figure 7(a). This suggests that broadcasting actions, like Story View, are used more frequently than narrowcasting communication features. Going ahead, we will calculate the mean for all users by considering each user's sessions [23]. This will allow us to aggregate data at the user level. The topic probabilities in Figure 7(b) are more uniformly distributed, suggesting that most Snapchat users utilize the program in a balanced manner across all features. A comparison between the user-level aggregate in Figure 7(b) and the session-level aggregation in Figure 7(a) reveals that users who use Snapchat's narrowcasting communication features, such as "Chat" and "Snap," more frequently tend to have more sessions on average.

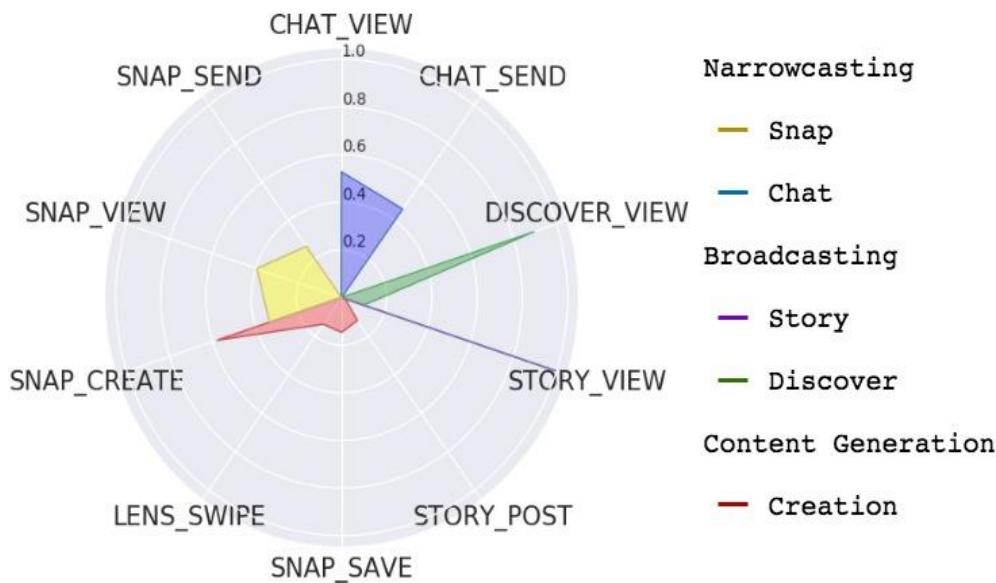
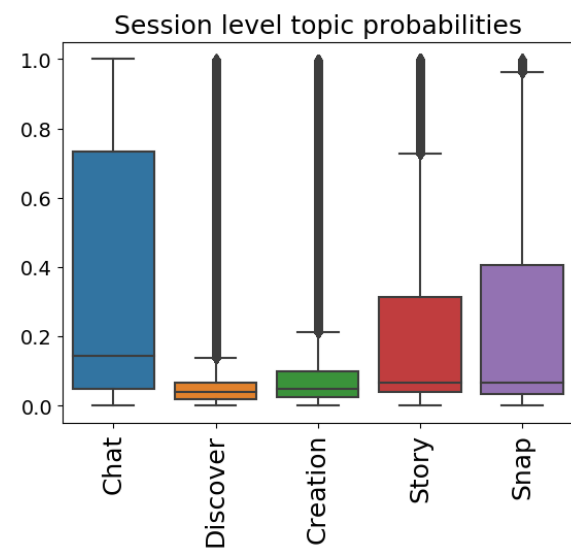
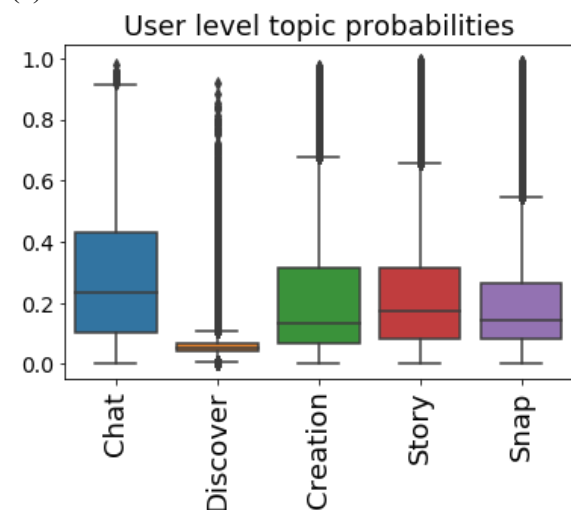


Figure 6. We discovered five session genes; each gene's probability composition is shown in the chart [20].

These paths are the foundation for creating action graphs because the transition probabilities of graph edges are obtained from common pathways. By removing bi-grams, tri-grams, and so on from session sequences, these pathways can be found. Common routes preserve the sequence order of action sequences, in contrast to session genes, where each session sequence is viewed as a bag of words in an LDA model fit. Common pathways, which are the fundamental building blocks of action graphs, facilitate comprehension of how action graphs illustrate patterns of user behavior [23]. Different levels of aggregation are used for comparison and to provide analytical insights.



(a) Session level.



(b) User level.

Figure 7. Distribution of session genes at both the session and user levels [24].

User Clusters and Engagement

Our goal in this part is to investigate the connection between user involvement and the previously described behavioral patterns. To demonstrate the significance of modelling higher-order action graphs, we first look at the relationship between user engagement and lower-level features. At each level, we use feature sets to do user clustering, and we apply silhouette analysis to determine the number of clusters [42]. For comparability, we keep the silhouette score constant at 4; therefore, we choose the cluster number with the highest silhouette score. Based on the probability of session genes, we cluster lower-level graph characteristics using K-means clustering. We can divide users into groups with different dominant genes by employing 4 clusters. As a result, distinct user groups emerge that are primarily characterized by their story, creation, snap, and chat genes (Figure 8). Though these user clusters can distinguish between users with different levels of interaction, as Figure 10(a) shows, a clear differentiation is not immediately visible.

Based on the makeup of users' session genes, we create four user clusters such as the Story Viewer, Creator, Snap-er, and Chatter clusters. Our process for cluster users is the same as it is for higher-order features like shared pathways [17]. We find four user clusters using the probability normalized by the number of sessions and the average session duration per user. With user groups displaying

preferences for specific paths, each cluster displays a unique distribution over the top paths shown in Figure 9. It is unclear, therefore, how the engagement rates for each user cluster differ in Figure 10(b). Consequently, we move on to characteristics at a higher level.

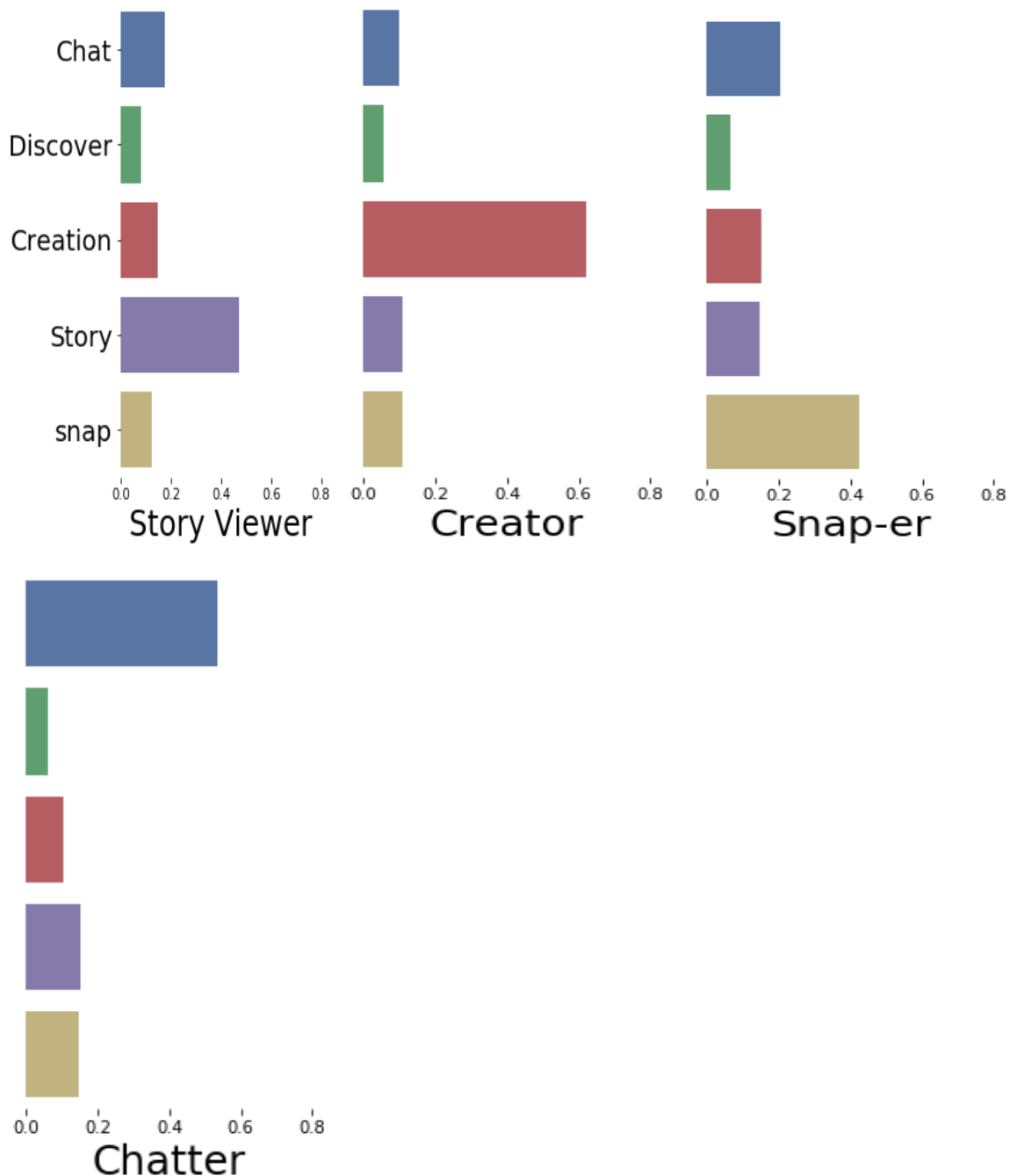


Figure 8. Session gene profiles [27].

As illustrated in Figure 10(c), user engagement rates clearly reveal a divide in contrast to the user groups based on lower-level attributes. It is clear that activity graphs and their fundamental characteristics are a superior way to identify different kinds of users [25]. In summary, action graphs provide detailed information and are far more instructive when it comes to determining user engagement. This emphasizes how crucial it is to use action graph modelling to forecast future interaction.

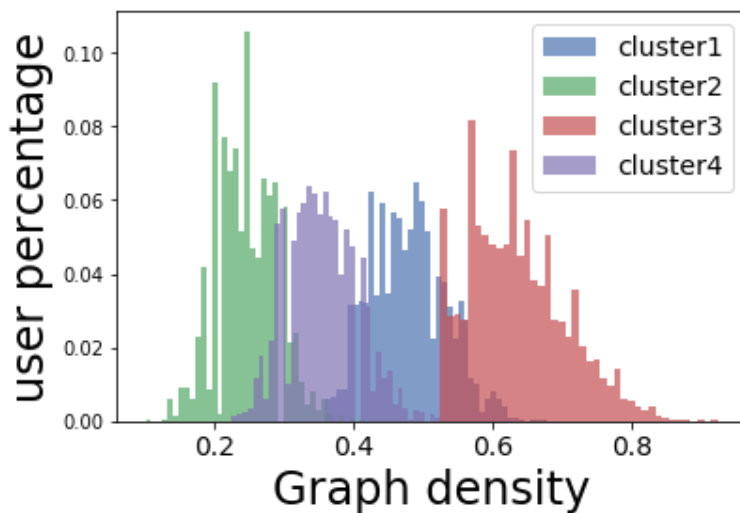


Figure 9. Graph metrics are utilized to determine the density distribution of clusters [31].

FORECASTING USER ENGAGEMENT

In this part, we outline two prediction objectives pertaining to the prediction of user interaction and suggest various forecasting techniques to address these requirements. These models blend human-crafted, interpretable properties with deep neural networks. Initially, we offered a feature-based approach that utilizes several feature detail levels [26]. To show how action graphs can be used to predict future user participation levels, we then go on to a more intricate end-to-end prediction model.

Task Formulation

We present a novel prediction target designed to predict a user’s future interaction. This assignment emphasizes the importance of action graphs in improving prediction accuracy and expands upon previous findings. There are two possible structures for the prediction task is regression or classification. First, we establish a three-class categorization task. We record the days that a user is actively using the app throughout the next 2 weeks, from week 3 to week 4, which is the observation period after registration. When a user has at least one legitimate session and at least one activity that we record during that session, that day is considered active [31]. Our classification of a user’s engagement pattern is “Increases” if their number of active days throughout week 3 and week 4 is higher than it was during the prior period; “Decreases” if not. If there is no difference in the number of active days between the first 2 weeks and the next 2 weeks, the user will be placed in a different class that indicates nothing has changed. We also observe that active days appear to show more noticeable variations from week 2 to week 4 compared to the first 2 weeks. We characterize the challenge of active rate prediction as a regression, with the goal of forecasting not only the shift in engagement trend between the two periods but also the active rate for the next several weeks.

Feature-based Forecasting Model

We first propose a feature-based model that combines explainable graph properties with macroscopic features. This model provides interpretability and acts as a foundation for models based on deep neural graphs. A simple forecast is based on the macroscopic features, which are derived from user profiles and include measures like session count and average session time [33]. As seen in the figs, clustering based on graph-level attributes has the strongest relationship with the active rate.

The macroscopic features include information about (1) the average number of daily sessions per user, (2) the average length of a Snapchat session, (3) the gender of the user, (4) the user’s maximum age during the observed time, and (5) the total number of friends. These characteristics provide clear insights into potential effects on the user’s rate of activity. On the other hand, an action graph captures a multitude of user behavior patterns, many of which are complex and difficult to see from a high-

level viewpoint. To show how action graphs may be used to capture user behavior data and provide insights into activity patterns and app interaction, we extract several interpretable graph features from them [35]. This shows what happened right after the session started, revealing the user’s desired first interaction or the app trigger. These probabilities indicate the possibility that different activities will take place, offering more information on user preferences and behavior. Most likely the last activities before the end of the session. These paths start at the beginning of the session and show the most likely paths of action that a user could take. They are calculated by taking the joint probability of the transition probabilities T_i of the edges that make up the path. These paths trace from the start node to the end node using BFS search [39]. A limit on path length of six is incorporated to prevent endless searches within loops and self-loops. For every path, we evaluate its potency as $P \cdot 1/N$, where P is the joint probability of edges in the path and N is the total number of edges. By utilizing Johnson’s technique, we can extract every elementary cycle present in a graph. $P \cdot 1/N$, where N is the number of edges in the cycle and P is the joint probability $\prod(T_i)$ of edges, represents the cycle’s strength. We will showcase a few graph elements that provide information on user involvement in section 5.6. Simple feature vectors are integrated into the feature-based model and used as the classifiers’ input. To choose the best performance, we investigated SVM and SoftMax classifiers in addition to Ridge and Linear regression techniques for the regression problem. While the SoftMax classifier performs best when all features are included, SVM serves as a baseline.

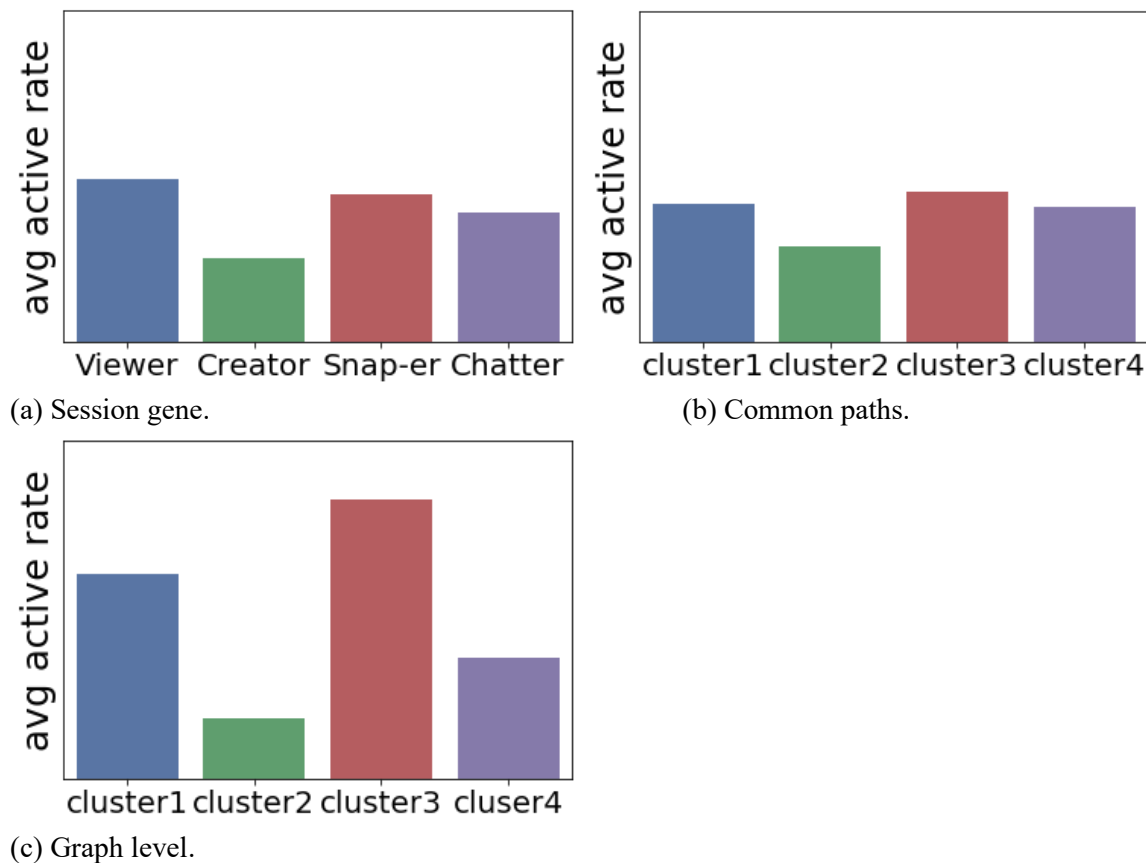


Figure 10. Average active rate of different clustering [35].

Deep Neural Forecasting Models

More intricate neural models that can process a variety of input signals are incorporated to improve the forecast accuracy of user involvement. As a baseline model, we first use LSTM to capture the temporal patterns in activity sequences. Next, we present our most successful engagement prediction model – a multichannel end-to-end training framework – along with temporal and static graph modelling methodologies. Motivated by previous work on user sequential behavior data modelling,

our solution combines LSTM sequence-to-sequence learning techniques to encode and predict churn rate [21]. We use the activity sequences shown in Figure 3 as a 10-dimensional time series input to create 2-layer LSTMs. The sequential character of behavioral data is well captured by the LSTM design, which enables us to comprehend how user behaviors change over time. We intend to combine LSTMs' superior ability to capture temporal dependencies in behavioral data with other features to produce a prediction model that is more robust. Our explainable graph features are good at portraying the graph, but they could miss some dependencies and patterns. This contrasts with static graphs, which are combined over the course of the 2-week observation period and provide a single unique graph for every user. Therefore, we choose a more powerful approach to make predictions and encode static graphs. For node classification, the Graph Convolutional Network (GCN) may work in semi-supervised and supervised environments by encoding each node as an embedding vector. Using message passing at each layer to leverage neighbor information, a GCN learns the node representation and updates the node embedding [23].

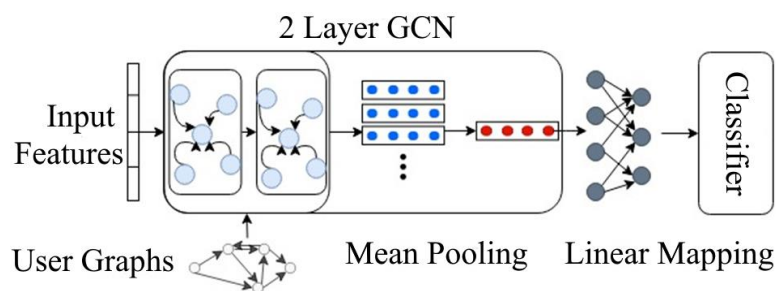


Figure 11. GCN model over static graph for engagement forecasting [38].

The inputs of a two-layer GCN structure with mean-pooling applied after the last layer are the user's static action graph and its node features. After undergoing a linear mapping process, the single graph embedding vector is sent into the classifier. As shown in Figure 11, we modified a two-layer GCN structure to fit our situation of directed graph categorization. We train all these little graphs collectively as a batched giant graph, which enables us to categorize many small graphs from multiple users. We can efficiently batch these graphs together for training by utilizing the DGL library. A mean-pooling transformation (v_i) is applied to the two-layer GCN's output at the end of each iteration [45]. This transformation creates a unique embedding vector for each graph that captures its structural properties by combining all the node embeddings of the graph into a single representation. For our classification prediction, we then apply a SoftMax layer to the linearly projected graph embedding vector v_G . Because the GCN graph embedding is trained with our target engagement data and incorporates neighboring information from the action graph on each iteration, it turns out to be a good representation of action graphs. By adding more features to the graph embedding, we can improve the prediction capability. Action graphs in real-world situations change over time. In the past, action graphs were handled as static objects, with each session over the observation period being combined into a single graph. We now provide temporal graph modelling, an adaptation of action graph modelling that is time dependent. While temporal graphs record the dynamic changes in user patterns at each time step, static graphs ignore the evolution of user behavior [47]. To create an action graph, we aggregate all sessions within each time step. We use a day as the time-step unit in our arrangement. Consequently, we produce 14 different temporal graphs during a 14-day observation period. Because LSTM networks are excellent at capturing temporal correlations, they are a good fit in this situation. Within our model, every graph is transformed by a Graph Convolutional Network (GCN) to produce a unique graph embedding vector. For each step, we encode all graphs using a single, unified GCN. A sequential sequence is formed by the graph embeddings produced by the GCN at every time step. For prediction, this sequence is then fed into a single LSTM network. In short, we leverage each user's daily sessions to compute an action graph at the end of each day. We take the graph embeddings out of every graph and feed them chronologically into the LSTM network. For the classification task, a SoftMax layer is then applied to the LSTM's final hidden state output.

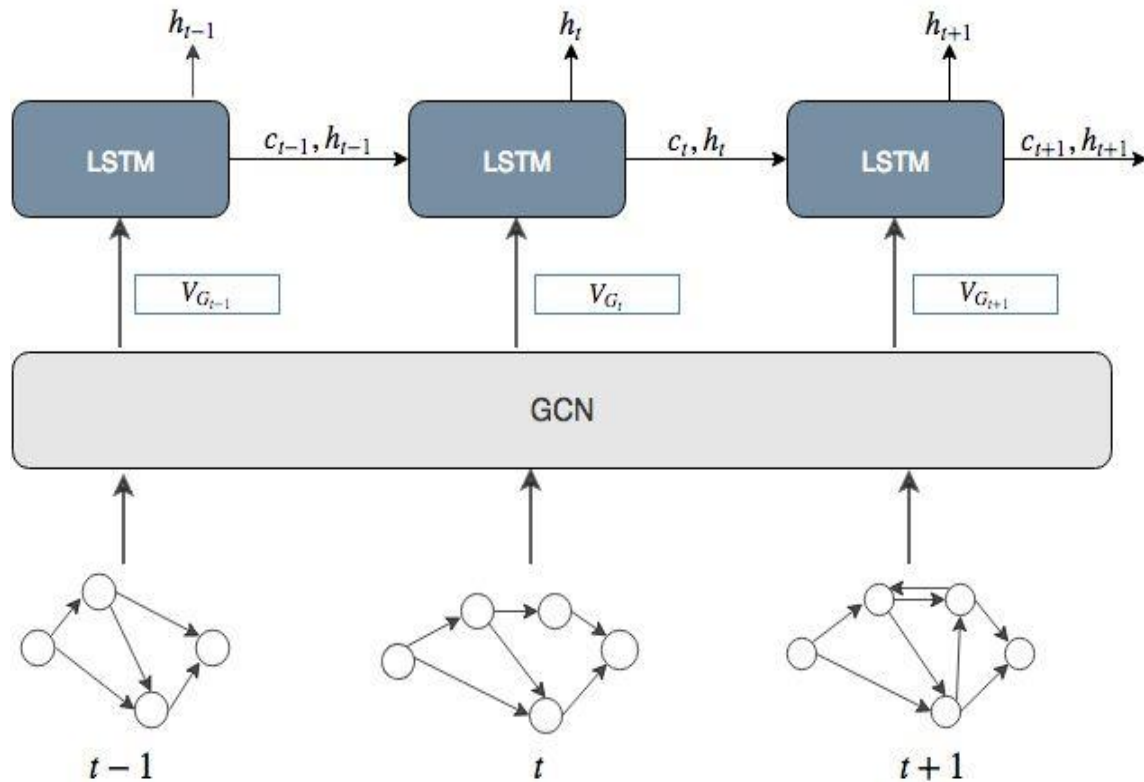


Figure 12. Using a temporal action graph made up of several action graph snapshots, we apply the GCN-LSTM model to predict engagement [41].

We feed the current action graph snapshot into a GCN module for every time step, and after mean pooling, we transmit the output embedding to an LSTM for prediction. We propose a simpler method to improve the training process, such as combining the GCN action graph model with the LSTM for activity sequences. In the last step, we also incorporate macro features to compute a single loss that is backpropagated to all the models at the same time [37]. It is possible to calculate the cross-entropy loss for every class in the SoftMax classifier for every user and the matching binary ground truth. In particular, the loss is computed in accordance with class c if the user's ground truth matches it. Alternatively, the SoftMax classifier can be swapped out for a linear SVM classifier, which would result in the model parameters indicated by w being used to compute a multi-class hinge loss shown in Figure 12. A linear regression layer is employed for the regression task and Mean Squared Error (MSE) is used to calculate the loss. To finish the training process, this loss value, represented by l , is then backpropagated across the two models in the end-to-end structure. When the models are trained end-to-end, they can collaborate to learn and achieve optimal performance [39]. The integration of the temporal GCN-LSTM model, LSTM activity sequences, and macroscopic characteristics is shown in Figure 13. According to our research, end-to-end joint training of these models produces the greatest results for forecasting user involvement.

EXPERIMENTS ON USER ENGAGEMENT PREDICTION

We test our hypothesis by gradually adding features and evaluating each one separately. To keep consistency for comparison, we combine 2 weeks of user graph data with 2 weeks of activity sequence data, which covers the whole time of our observation. Anticipating future participation in two areas is our aim: (1) user engagement trends and (2) future active rates. We conduct experiments in which we feed our model with different feature combinations to predict a single class or value for each of the two tasks. We predict intuitively that graph embedding will perform better than explainable graph features, offering an improvement over the existing activity sequence modelling

method [35]. The best results are obtained when temporal graph embedding, activity sequence embedding, and macroscopic characteristics are combined. A subsampled dataset of newly added users taken from the Snapchat database was used for the trials. The initial new user data was subsampled into a 1:1:1 split for each of the three types of engagement trends to avoid a skewed distribution. Based on the comparison of active days between the first 2 weeks and the next 2 weeks, engagement trend labels were assigned. The active days of the next 2 weeks are compared to the active days of the previous 2 weeks to determine the labels for the regression job of forecasting activity rate. By dividing the data into 80:20 training and testing sets at random, repeating this procedure 10 times for 10-fold cross-validation, and summing the results, we were able to assess the models. These tests were carried out on a Google Cloud Engine system that was just outfitted with CPUs.

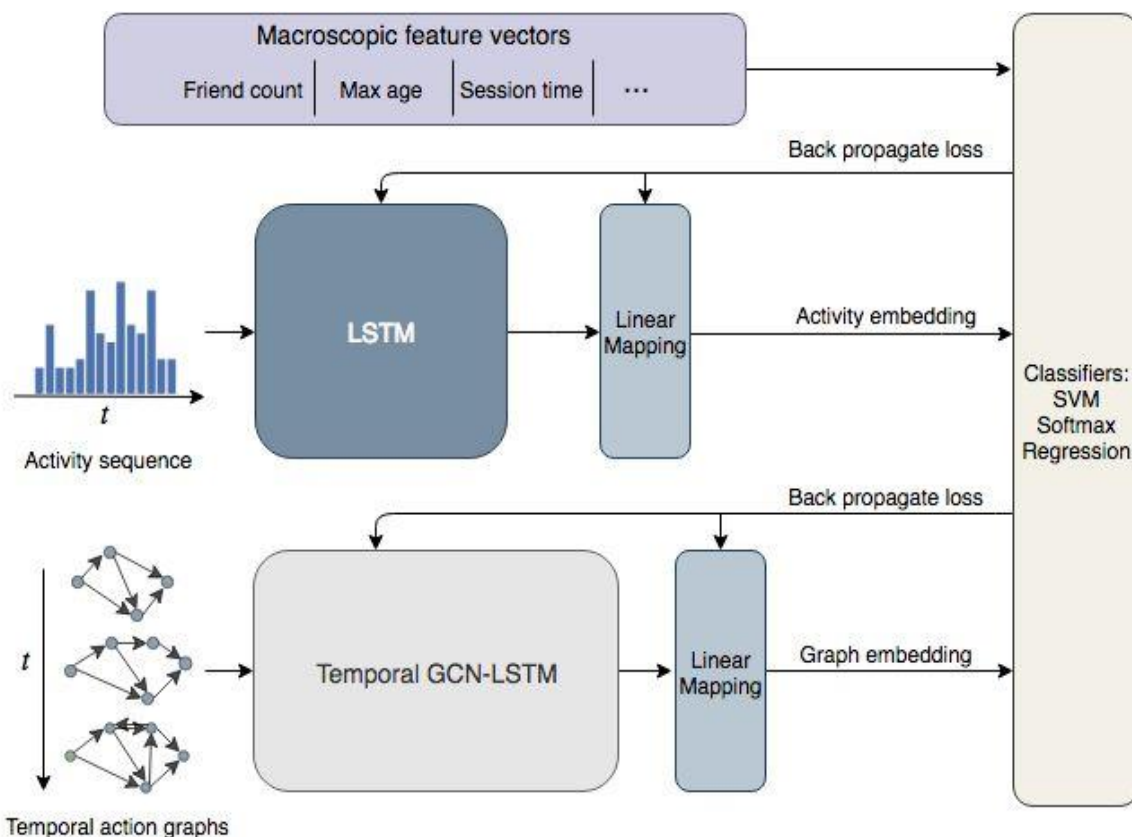


Figure 13. Deep multichannel forecasting framework [45].

Model Training and Hyperparameters

SVM, SoftMax, and linear regression were the last prediction layers we tested in our model. All variables were scaled using resilient-to-outlier approaches before training. We set up two levels in our GCN network, and in the final layer, we used mean pooling to get a single vector representation for every action graph. Each graph's adjacency matrix served as an input feature [32]. Our two-layer LSTM network is configured with an embedding size of 32 and a dropout rate of 0.5. These parameters were established empirically based on earlier research, and our study provides additional validation for them. We used PyTorch and the DGL library to create our end-to-end framework, the LSTM models, and GCN. Interestingly, we found that training the LSTM and GCN models end-to-end from scratch yielded comparable (or marginally better) results, thus we decided against pretraining them. We used a held-out validation set and grid search to find our model's hyperparameters. To avoid overfitting, we used L2 regularization ($1e-3$) on the classifier layer and a dropout rate of 0.5 on the LSTM layer in our trials. Given the comparatively short size of the action graph in our situation, the hidden state size of the GCN was adjusted to 20 to reduce overfitting. For

our GCN network, we also tried with various input configurations, such as random Xavier normal/uniform initialization, but we discovered that they produced results that were comparable, if not better.

Performance Study on Classification of User Engagement Trends

We evaluate graph features and graph embeddings against multiple baselines are SVM as our baseline classifier, a weaker baseline that uses macroscopic characteristics for prediction, and a stronger baseline that uses LSTM-based activity sequence prediction. The classification task's ultimate performance is displayed by the SoftMax classifier. Basic macroscopic characteristics, graph features, and other embedding features can all be used to train the feature-based model in two steps. On the other hand, the end-to-end model concurrently trains every model parameter in concert. Because it is simpler and usually produces a better fit for predictions, the end-to-end model is preferred [28]. While temporal graphs aggregate data daily for each user, static graphs aggregate data over the course of the 2-week observation period. Temporal graphs, in contrast to static graphs, consider how user behavior patterns change over time. Predictive performance is enhanced by modeling the time variation of action graphs. This means that by employing action graphs that change every day, we can efficiently collect user behavior trends. We investigate several feature combinations in our task of using 3-way classification to forecast engagement trends.

Building on earlier research, we take LSTM activity sequence embedding as a baseline [21]. It is interesting to see that our temporal graph model outperforms the LSTM activity sequence on its own, where we began with a basic model that only included macroscopic features and then progressively included additional graph features. Performance is improved by combining activity embedding and static graph embedding with macroscopic characteristics. Combining temporal graph embedding, activity embedding, and macroscopic characteristics yields the best results. Taking baselines in the form of macroscopic features, each graph feature adds distinct effects and improves performance above baseline. When all graph attributes are combined, it can be seen that they improve the capacity to anticipate user involvement and provide light on the reasons why a more intricate graph network performs better.

RESULTS ON ACTIVE RATE PREDICTION

Since it requires solving a regression problem over a 14-day period, predicting active rate is a more involved challenge than predicting engagement level. We experimented with the same feature combinations that are employed in the prediction of engagement trends. In this finer-grained challenge, temporal graph modelling works substantially better than the baseline activity sequence, yielding far superior prediction results. The result that such a huge performance jump can be achieved by modeling temporal graphs shows its Combining numerous feature combinations greatly improves our model's prediction of future engagement rates [27]. We obtain the highest result by combining temporal action graphs, macroscopic characteristics, and baseline activity sequences. Through the comprehensive data gathered by the action graphs, companies may decide on targeted incentives and user retention with knowledge.

Case Study on Graph Features

We are especially interested in finding the explainable graph traits that best distinguish between different user involvement levels. We may generate sparse coefficients that naturally function as a feature selection tool by incorporating a straightforward L1 penalty term into a linear kernel SVM classifier. After training every graph feature with the linear SVM, this method enables us to extract the features that have the greatest predictive power. The highest-order predictive graph traits that were shown to be most predictive are cycle strength, end-to-end paths, and k-hop paths. Almost all the top three paths from the start of the session involve conversation activity. This implies that chat activities – which entail conversing with friends – are reliable markers of higher levels of engagement. However, rather than being simple cycles, the top cycles are primarily composed of complicated structures [45]. This suggests that a more well-rounded user on Snapchat tends to be more engaged

with the platform's primary features. In a similar vein, longer and more complicated end-to-end paths were chosen rather than short and straightforward ones. This suggests that consumers show higher levels of engagement when they spend more time and use more functionalities during each session. The interpretability of these elementary graph properties provides the basis for comprehending the effectiveness of deeper neural models.

Related Work

Numerous studies use different approaches to simulate user behavior to forecast if a user would return to the site. Predicting return rates is one of the most common and important tasks on social media platforms as businesses look for the secrets that make them successful. Numerous studies try to forecast the churn rate, or the point at which a user would stop using the platform, as well as user retention [2, 14, 21]. Some studies take a different approach to the problem, estimating the lifetime of a new user [22], forecasting when a user will return [13], or identifying whether a user will return [16]. Finding and predicting user intention is another goal of user behavior modelling. Studies have been done to forecast a user's intention [8], purchase intent on Pinterest [17], and user intent and subsequent behavior. It is interesting to note that several research in this field concentrate on simulating consumer consumption patterns [5, 19]. Prediction tasks have been approached using a variety of techniques. Several studies that employ gradient boosted tree or basic logistic regression techniques using macroscopic features report successful outcomes [1, 5, 16]. Others use Long Short-Term Memory (LSTM) structures [21] or Cox proportional hazard models [13, 22] that take time information into account when making predictions. Many articles in the domains of social and search platforms model activity sequences as graphs. On the other hand, related literature models queries and clickstreams into activity graphs (Markov chain graphs), in contrast to our approach, which models user activities into activity graphs. Queries can be modelled into sequential graphs in search platform research to forecast search success [12] or provide recommendations [7]. Moreover, query graphs are used to simulate user search behavior [3]. Like this, on social and web platforms, clickstreams are used to simulate user behavior as Markov chain models or sequential graphs [4, 10, 18, 20]. Although research has been done on activity sequence graphs, including click streams and search queries, little has been done to model activity sequence graphs involving user actions on social media platforms.

CONCLUSION

This study examines and analyzes new Snapchat users' behavior using action graphs. Our investigation has led to the development of a strong and anticipatory action graph modelling framework. By means of our analysis, we provide significant insights into understanding post-registration trends of new users and their likelihood of future interaction. While Snapchat data is the primary focus of our analysis and models, other online or offline social media platforms can also benefit from the use of the same techniques and models. Snap Inc. is using our action graph prediction platform to give analytical insights and benefits for a range of business and production decisions, including user modeling, growth, and retention, among other things. Future studies could investigate the relationship between a user's ego-network and their behavioral tendencies. When it comes to applying machine learning for predictive analytics in social media data, social media platforms hold a lot of promises. Machine learning algorithms can be used by businesses to spot trends, predict outcomes, and make data-driven decisions. However, concerns, like data privacy, data quality, and interpretability, need to be addressed to guarantee the moral and proper use of these powerful tools in social media.

REFERENCES

1. Althoff T, Leskovec J. Donor retention in online crowdfunding communities: A case study of DonorsChoose.org. Proc World Wide Web Conf. 2015;2015:34–44. doi: 10.1145/2736277.2741120.
2. Au WH, Chan KCC, Yao X. A novel evolutionary data mining algorithm with applications to churn prediction. IEEE Trans Evol Comput. 2003;7(6):532–545. doi: 10.1109/TEVC.2003.819264.

3. Baeza-Yates R, Hurtado C, Mendoza M, Dupret G. Modeling user search behavior. In: Web Congress, 2005. LA-WEB 2005. Third Latin American. 2005. doi: 10.1109/LAWEB.2005.23.
4. Benevenuto F, Rodrigues T, Cha M, Almeida V. Characterizing user behavior in online social networks. In: Proceedings of the 9th ACM SIGCOMM Conference on Internet Measurement 2009. Illinois: 2009. doi: 10.1145/1644893.1644900.
5. Benson AR, Kumar R, Tomkins A. Modeling user consumption sequences. In: WWW '16: Proceedings of the 25th International Conference on World Wide Web. 2016. pp. 519–529. doi: 10.1145/2872427.2883024.
6. Blei DM, Ng AY, Jordan MI. Latent Dirichlet allocation. *J Mach Learn Res.* 2003;3:993–1022.
7. Boldi P, Bonchi F, Castillo C, Donato D, Gionis A, Vigna S. The query-flow graph: Model and applications. In: Proc CIKM; 2008. doi: 10.1145/1458082.1458163.
8. Cheng J, Lo C, Leskovec J. Predicting intent using activity logs: How goal specificity and temporal range affect user behavior. In: WWW '17 Companion: Proceedings of the 26th International Conference on World Wide Web Companion. 2017. pp. 593–601 doi: 10.1145/3041021.3054198 2017.
9. Ciampaglia GL, Taraborelli D. MoodBar: Increasing new user retention in Wikipedia through lightweight socialization. In: Proc 18th ACM Conf Comput Supported Coop Work Soc Comput; 2015. pp. 734–742. doi: 10.1145/2675133.2675181.
10. Gündüz Ş, Özsü MT. A web page prediction model based on click-stream tree representation of user behavior. In: The ninth ACM SIGKDD international conference. 2003. doi: 10.1145/956750.956815.
11. Halfaker A, Keyes O, Kluver D, Thebault-Spieker J, Nguyen T, Shores K, et al. User session identification based on strong regularities in inter-activity time. In: WWW '15: Proceedings of the 24th International Conference on World Wide Web. 2015. pp. 410–418. doi: 10.1145/2736277.2741117.
12. Hassan A, Jones R, Klinkner KL. Beyond DCG: User behavior as a predictor of a successful search. In: WSDM '10: Proceedings of the third ACM international conference on Web search and data mining. 2010. pp. 221–230. doi: 10.1145/1718487.1718515.
13. Kapoor K, Sun M, Srivastava J, Ye T. A hazard based approach to user return time prediction. In: KDD '14: Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining. 2014. pp. 1719–1728. doi: 10.1145/2623330.2623348.
14. Kawale J, Pal A, Srivastava J. Churn prediction in MMORPGs: A social influence based approach. In: Proc CSE. IEEE; 2009;4:423–428.
15. Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks. arXiv preprint. 2016;arXiv:1609.02907. doi: 10.48550/arXiv.1609.02907.
16. Lin Z, Althoff T, Leskovec J. I'll be back: On the multiple lives of users of a mobile activity tracking application. Proc Int World Wide Web Conf. 2018;2018:1501–1511. doi: 10.1145/3178876.3186062.
17. Lo C, Frankowski D, Leskovec J. Understanding behaviors that lead to purchasing: A case study of Pinterest. In: KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016. pp. 531–540. doi: 10.1145/2939672.2939729.
18. Sadagopan N, Li J. Characterizing typical and atypical user sessions in clickstreams. In: WWW '08: Proceedings of the 17th international conference on World Wide Web. 2008. pp. 885–894. doi: 10.1145/1367497.1367617.
19. Trouleau W, Ashkan A, Ding W, Eriksson B. Just one more: Modeling binge watching behavior. In: KDD '16: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2016. pp. 1215–1224. doi: 10.1145/2939672.2939792.
20. Wang G, Zhang X, Tang S, Wilson C, Zheng H, Zhao BY. Clickstream user behavior models. *ACM Trans Web.* 2017;11(4):1–37. doi: 10.1145/3068332.
21. Yang C, Shi X, Jie L, Han J. I know you'll be back: Interpretable new user clustering and churn prediction on a mobile social application. In: The 24th ACM SIGKDD International Conference. 2018. doi: 10.1145/3219819.3219821.

22. Yang J, Wei X, Ackerman MS, Adamic LA. Activity lifespan: An analysis of user survival patterns in online knowledge sharing communities. In: Proceedings of the International AAAI Conference on Web and Social Media. 2010;4(1):186–193. doi: 10.1609/icwsm.v4i1.14010.
23. Baccianella S, Esuli A, Sebastiani F. SentiWordNet 3.0: An enhanced lexical resource for sentiment analysis and opinion mining. In: Proc Int Conf Lang Resour Eval (LREC); 2010. pp. 2200–2204.
24. Go A, Bhayani R, Huang L. Twitter sentiment classification using distant supervision. CS224N Proj Rep. 2009;1:12.
25. Li J, Cardie C, Ji H. Mining point-wise mutual information from tweet streams. In: Proc 50th Annu Meet Assoc Comput Linguist (Short Papers); 2012. pp. 381–385.
26. Zhang X, Zhao J, LeCun Y. Character-level convolutional networks for text classification. In: Adv Neural Inf Process Syst; 2015. pp. 649–657.
27. Kwak H, Lee C, Park H, Moon S. What is Twitter, a social network or a news media? In: WWW '10: Proceedings of the 19th international conference on World wide web. 2010. pp. 591–600. doi: 10.1145/1772690.1772751.
28. Pennacchiotti M, Popescu AM. Democrats, Republicans and Starbucks aficionados: User classification in Twitter. In: Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2011. pp. 430–438. doi: 10.1145/2020408.2020477.
29. Pang B, Lee L. Opinion mining and sentiment analysis. Found Trends Inf Retr. 2008;2(1–2):1–135. doi: 10.1561/15000000011.
30. Liu B. Sentiment analysis and opinion mining. Synth Lect Hum Lang Technol. 2012;5(1):1–167. doi: 10.2200/S00416ED1V01Y201204HLT016.
31. Shu K, Mahudeswaran D, Wang S, Lee D. Fake news detection on social media: A data mining perspective. ACM SIGKDD Explor Newsl. 2017;19(1):22–36. doi: 10.48550/arXiv.1708.01967.
32. Vosoughi S, Roy D, Aral S. The spread of true and false news online. Science. 2018;359(6380):1146–1151. doi: 10.1126/science.aap9559.
33. Koren Y, Bell R, Volinsky C. Matrix factorization techniques for recommender systems. Computer. 2009;42(8):30–37.
34. Ricci F, Rokach L, Shapira B. Recommender Systems Handbook. Berlin: Springer; 2015. doi: 10.1007/978-0-387-85820-3_1.
35. Newman MEJ. Networks: An Introduction. Oxford: Oxford University Press; 2010. doi: 10.1093/acprof:oso/9780199206650.001.0001.
36. Leskovec J, Krevl A. SNAP datasets: Stanford large network dataset collection. 2014. Available at <http://snap.stanford.edu/data>
37. Becker H, Naaman M, Gravano L. Beyond trending topics: Real-world event identification on Twitter. In: Proc Int AAAI Conf Weblogs Soc Media (ICWSM); 2011.
38. Petrović S, Osborne M, Lavrenko V. Streaming first story detection with application to Twitter. In: Proc NAACL-HLT; 2010. pp. 181–189.
39. Romero DM, Meeder B, Kleinberg J. Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on Twitter. In: Proc World Wide Web Conf (WWW); 2011. pp. 695–704.
40. Lerman K, Hogg T. Using a model of social dynamics to predict popularity of news. In: WWW '10: Proceedings of the 19th international conference on World wide web. 2010. pp. 621–630. doi: 10.1145/1772690.1772754.
41. Tsai CF, Wang SW, Huang YM, Tseng SS. Predicting user engagement on social media: A data perspective. Soc Netw Anal Min. 2014;4(1):1–18.
42. Yang J, Counts S, Hoff A. Predicting the speed, scale, and range of information diffusion in Twitter. In: Proc Int Conf Weblogs Soc Media (ICWSM); 2011.
43. Dinakar JR, Vagdevi S. Real-time streaming analytics using big data paradigm and predictive modelling based on deep learning. Int J Recent Innov Trends Comput Commun. 2023;11:161–165. doi: 10.17762/ijritcc.v11i4s.6323.
44. Bai VS, Sudha T. A systematic literature review on cloud forensics in cloud environment. Int J Intell Syst Appl Eng. 2023;11(4s):565–578.

45. López M, Popović N, Dimitrov D, Botha D, Ben-David Y. Efficient dimensionality reduction techniques for high-dimensional data. *Kuwait J Mach Learn*. 2023;1(4).
46. Dwarkanath Pande S, Hasane Ahammad DS. Cognitive computing based network access control system in secure physical layer. *Res J Comput Syst Eng*. 2022;3(1):14–20.
47. Dhabliya D. An integrated optimization model for plant diseases prediction with machine learning model. *Mach Learn Appl Eng Educ Manag*. 2021;1(2):21–26.