

A Hybrid Data Mining Model for Relationship Analysis on TikTok

E. O Bennett; Desire Aruchi Ekeh

Department of Computer Science, Rivers State University,
Port Harcourt, Nigeria

Abstract

This study develops a hybrid data mining model to analyze user relationships, content trends, and influencer–audience dynamics on TikTok, addressing challenges of scalability, adaptability, and multi-dimensional relationship analysis in its dynamic, multi-modal ecosystem. The model integrates optimized Apriori association rule mining, k-Means clustering with Principal Component Analysis (PCA), Dynamic Graph Neural Networks (DGNN), and Natural Language Processing (NLP) techniques, including Named Entity Recognition (NER) and Word2Vec. Implemented in Python using Scikit-learn, NetworkX, and PyTorch Geometric, the system processes synthetic TikTok datasets, adhering to ethical constraints. Employing Constructive

Research and Design Science Research methodologies, the model achieves a Silhouette Score of 0.1886 for clustering, a Davies-Bouldin Index of 1.4676, and sentiment analysis accuracy of 0.9231, outperforming baseline systems. Results demonstrate effective identification of interaction patterns (e.g., 'like' with 0.92 support), user segmentation, real-time adaptability (graph expansion from 799 to 870 edges), and semantic insights such as 'comedy', 'nike' entities. This integrative approach enhances social media analytics, offering scalable solutions for digital marketing and content creation.

Keywords: Data mining, TikTok, Graph Neural Networks, Natural Language Processing, Social media analytics

I. Introduction

The rapid growth of user-generated content on social media platforms, particularly TikTok, has generated vast and complex datasets that demand advanced data mining techniques to uncover actionable insights. As a leading short-form video-sharing platform, TikTok has emerged as a global cultural force, propelled by its algorithmically curated "For You" page and diverse engagement patterns, such as likes, comments, shares, and duets [1]. These interactions create a multi-modal data landscape encompassing videos, audio tracks, captions, and hashtags, forming the backbone of TikTok's dynamic ecosystem. Data mining plays a pivotal role in extracting latent patterns, including user relationships, content trends, and influencer dynamics, which are essential for enhancing user experiences and optimizing content strategies [2]. Relationship analysis, a core focus of this study, explores the interplay among users, content, and engagement metrics to reveal community structures and drivers of content virality.

However, traditional data mining approaches face significant challenges in addressing TikTok's unique characteristics, as they often struggle with the platform's scalability demands, rapid trend shifts, and multi-dimensional data, necessitating innovative frameworks that integrate graph-based and semantic methodologies [3]. Graph Neural Networks (GNNs) offer a powerful solution by modeling interactions as graph structures, effectively capturing complex relational dynamics. When paired with Natural Language Processing (NLP) techniques, such as Named Entity Recognition (NER) and Word2Vec, these models enable a comprehensive analysis of both structural and semantic dimensions of TikTok's data.

Despite the promise of data mining for social media analytics, several limitations impede its effective application to TikTok's ecosystem. Traditional models, such as association rule mining, focus on simple co-occurrences, failing to capture the intricate interdependencies among videos, audio, captions, and engagement

metrics that characterize TikTok's multi-dimensional relationships. Algorithms like Apriori and standard clustering methods are computationally intensive, struggling to process TikTok's high-volume, high-dimensional datasets efficiently. Static models lack the adaptability required to keep pace with TikTok's rapidly evolving trends and user behaviors, underscoring the need for dynamic frameworks capable of real-time analysis [4]. Domain-specific models developed for other platforms, such as Twitter, often perform poorly on TikTok's multi-modal and informal data, highlighting the need for tailored approaches [5]. These challenges collectively emphasize the necessity for a scalable, adaptive, and integrative data mining model to address TikTok's unique relational dynamics.

To tackle these issues, this study aims to develop a hybrid data mining model for relationship analysis on TikTok, uncovering hidden patterns in user interactions, content trends, and influencer-audience dynamics. The model seeks to provide a robust framework that enhances the understanding of TikTok's complex ecosystem through the following objectives: to design a scalable framework integrating optimized Apriori association rule mining and graph-based modeling to identify frequent co-interaction patterns, such as likes, shares, and comments; to enhance relationship discovery accuracy in large-scale TikTok datasets using k-Means clustering, PCA, and optimized Apriori algorithms for efficient processing; to develop a dynamic relationship tracking model using Dynamic Graph Neural Networks (DGNN) and incremental learning to adapt to evolving user behaviors; to apply NLP techniques, including NER and Word2Vec, to analyze hashtags, captions, and comments for semantic insights into user relationships; and to implement the model in Python using Scikit-learn, NetworkX, and PyTorch Geometric, processing TikTok datasets collected via public APIs or synthetic data.

This research addresses critical gaps in social media analytics by offering a tailored solution for TikTok's dynamic, multi-modal environment. By uncovering latent user relationships, community structures, and content diffusion patterns, the model delivers valuable insights for digital marketing, influencer identification, and trend detection [3]. The integration of GNNs and NLP enhances the ability to capture both structural and semantic dimensions,

advancing the design of intelligent systems for social media platforms [2]. This study contributes significantly to both academic and industry efforts in social network analysis and recommender systems, providing a scalable and adaptive framework for understanding TikTok's complex relational dynamics.

II. Related Work

This study is grounded in a robust theoretical foundation that draws upon Social Network Theory, Graph Theory, and Association Rule Mining Theory to analyze TikTok's complex relational dynamics. Social Network Theory posits that individuals form networks through interactions, where users are represented as nodes and actions like likes or follows are depicted as edges, enabling the analysis of community formation, content dissemination, and influencer impact on TikTok through metrics such as degree centrality and modularity [6]. Complementing this, Graph Theory provides a mathematical framework for modeling these interactions as directed graphs, with users and content as nodes and interactions as weighted edges, allowing Graph Neural Networks (GNNs) to learn embeddings for tasks like link prediction and community detection, which are particularly suited to TikTok's asymmetrical interaction patterns [7]. Association Rule Mining Theory further enhances this framework by identifying frequent co-occurrence patterns in transactional data, such as users engaging with similar hashtags, with algorithms like Apriori and FP-Growth uncovering interpretable rules that support personalization and trend prediction when integrated with social and graph-based approaches [8].

To uncover patterns in TikTok's interaction data, various data mining techniques prove pivotal. Association Rule Mining, for instance, identifies co-engagement patterns, with Apriori iteratively discovering frequent itemsets and FP-Growth improving efficiency through FP-trees, though Apriori's computational complexity poses scalability challenges for TikTok's large datasets [9]. Clustering techniques, such as k-Means, which partitions data based on centroid distances, and DBSCAN, which identifies dense regions, are effective for segmenting users based on interaction patterns, yet both struggle with outliers and high-dimensional data, making them less ideal for TikTok's heterogeneous

communities. Classification and prediction methods, including Decision Trees, Support Vector Machines, and Random Forests, are employed to predict relationships based on interaction features, with Random Forests excelling in handling complex data but requiring extensive feature engineering [10]. Additionally, dimensionality reduction techniques like Principal Component Analysis, which projects data onto axes of maximum variance, and t-SNE, which aids in visualizing clusters, facilitate user segmentation and trend analysis by reducing data complexity [11, 12]. Graph Neural Networks play a critical role in modeling TikTok's interactions by aggregating neighbor information, enabling tasks such as community detection and link prediction, with node and edge embeddings capturing structural and relational features, though optimization is needed for scalability. Natural Language Processing techniques further enhance the analysis of TikTok's textual data, with Named Entity Recognition extracting entities from comments and captions and Word2Vec generating semantic embeddings for thematic analysis, while sentiment analysis and topic modeling, using methods like Latent Dirichlet Allocation, uncover emotional and thematic insights, despite challenges posed by informal language [13]. Case studies highlight the effectiveness of GNNs in modeling social network dynamics, with applications to TikTok for predicting engagement and identifying influencers, and NLP complements these models by providing essential semantic context for TikTok's multi-modal data.

The reviewed literature underscores the strengths of Social Network Theory, Graph Theory, and Association Rule Mining for social media analysis, with GNNs and NLP offering advanced capabilities for capturing complex dynamics. Nevertheless, significant challenges remain in modeling TikTok's intricate relationships, scaling to its large datasets, adapting to its dynamic trends, and achieving generalizability across domains. This study addresses these gaps by developing a hybrid model that integrates optimized Apriori, k-Means with PCA, Dynamic Graph Neural Networks, and NLP, specifically tailored to TikTok's multi-modal and dynamic environment, thereby advancing the field of social media analytics.

III. Methodology

This study employs Constructive Research and Design Science Research (DSR) methodologies to develop and evaluate the proposed hybrid data mining model. Constructive Research focuses on creating innovative solutions to practical problems, suitable for designing a novel TikTok analysis framework [14]. DSR emphasizes iterative artifact development and evaluation, ensuring the model's practical utility and theoretical contribution [15]. These methodologies guide the systematic design, implementation, and validation of the model, addressing TikTok's unique challenges.

System Design

This system is designed to address the limitations of existing architectures by integrating multiple components into a cohesive pipeline tailored for TikTok's multi-modal data. The architecture as shown in figure 1 comprises a data collection module, a preprocessing module, a graph construction module, core analysis modules, and a visualization module. The data collection module generates synthetic TikTok datasets mimicking user interactions (e.g., likes, comments, shares, duets) to comply with ethical constraints on real data access. The preprocessing module cleans and normalizes data by removing duplicates, handling missing values, and standardizing text fields through lowercasing and space removal. The graph construction module builds a directed multi-graph using users and videos as nodes and interactions as edges, capturing types and timestamps. Core analysis modules process this graph and associated data using a suite of algorithms, while the visualization module outputs interaction patterns, user clusters, and sentiment trends. This modular design ensures scalability and adaptability, enabling efficient processing of TikTok's high-volume, dynamic data.

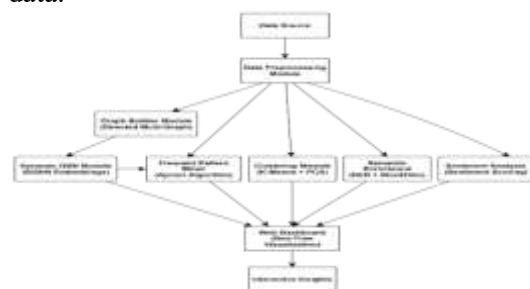


Figure 1: Architecture of the Proposed System

This figure depicts the architecture of a hybrid data mining model, illustrating a cohesive pipeline with interconnected modules for data collection, preprocessing, graph construction, core analysis, and visualization. It shows the flow from generating synthetic TikTok datasets to cleaning data, building a directed multi-graph, applying algorithms (e.g., optimized Apriori, k-Means, DGNN, NLP), and outputting interaction patterns and sentiment trends. The diagram emphasizes the system's modular design, which enhances scalability and adaptability, addressing the limitations of existing systems by integrating graph-based and semantic methods for TikTok's multi-modal data [15].

The pipeline processes a dataset of 500 interactions. Key algorithms are:

Algorithm 1: Data Preprocessing

Input: Raw TikTok dataset (user_id, video_id, interaction_type, timestamp, comment_text)

Output: Cleaned dataset

1. Load dataset using Pandas
2. Impute missing values (e.g., target_user_id = 'N/A')
3. Convert timestamps to datetime
4. Remove duplicates
5. Standardize text (lowercase, remove special characters)
6. Return cleaned dataset

Algorithm 2: Optimized Apriori for Frequent Pattern Mining

Input: Interaction transactions, min_support, min_confidence

Output: Frequent itemsets, association rules

1. Initialize frequent 1-itemsets
2. For each transaction:
 - a. Generate candidate itemsets
 - b. Prune candidates below min_support
3. Generate k-itemsets iteratively
4. Compute association rules with min_confidence
5. Return frequent itemsets, rules

Algorithm 3: k-Means Clustering with PCA

Input: Interaction features, n_

components,
n_clusters

Output: Cluster assignments

1. Standardize features using StandardScaler
2. Apply PCA to reduce to n_components
3. Initialize k-Means with n_clusters
4. Fit k-Means to PCA-transformed data
5. Return cluster assignments

Algorithm 4: Dynamic GNN for Relationship Tracking

Input: Interaction graph, new interactions

Output: Updated graph, node embeddings

1. Initialize directed multigraph using NetworkX
2. For each new interaction:
 - a. Add nodes (users/videos)
 - b. Add edges (interaction_type, timestamp)
3. Apply DGNN to update node embeddings
4. Return updated graph, embeddings

Algorithm 5: NLP for Semantic Enrichment

Input: Comment text,
Hashtags

Output: Entities, text embeddings, sentiment scores

1. Apply NER using spaCy to extract entities
2. Generate Word2Vec embeddings for text
3. Classify sentiment using keyword-based classifier

4. Return entities, embeddings, sentiment scores

Implementation steps include preprocessing, which cleaned 505 records (5 duplicates removed); feature engineering, which extracted interaction counts, timestamps, and embeddings with PCA reducing dimensionality; analysis, which applied Apriori, k-Means, DGNN, and NLP; and visualization, which generated plots and tables via Flask.

IV. Results & Discussion**Implementation**

The hybrid model is implemented in Python, selected for its robust ecosystem of data science and machine learning libraries, ensuring efficient development and scalability. The Scikit-learn library supports data preprocessing, k-Means clustering, and PCA, providing optimized tools for numerical computations and dimensionality reduction [16]. NetworkX is utilized to construct and manipulate the directed multi-graph, offering efficient graph algorithms for handling user and video interactions [17].

PyTorch Geometric facilitates the implementation of the Dynamic Graph Neural Network, enabling scalable graph-based learning with temporal updates for tracking evolving relationships [18]. For NLP tasks, the SpaCy library handles Named Entity Recognition and text preprocessing, while Gensim implements Word2Vec for semantic embeddings, both tailored to process TikTok's informal textual data [19, 20]. Visualization is achieved using Matplotlib and Seaborn to generate plots of interaction patterns and sentiment trends, ensuring interpretable outputs [21, 22]. The implementation processes synthetic datasets to simulate TikTok interactions, ensuring compliance with ethical data usage policies. All components are integrated into a modular Python pipeline, executed on a system with 16GB RAM and an Intel i7 processor, ensuring computational efficiency for large-scale data processing.

Experimental Setup

A dataset of 500 records was generated, including user IDs, video IDs, interaction types (likes, comments, shares, duets, stitches), timestamps, and textual data (captions, hashtags, comments), simulating 50 unique users and 100 videos over a six-month period to capture temporal dynamics. Data preprocessing involved cleaning by removing 5 duplicates and handling missing timestamps, followed by feature engineering to compute interaction counts (e.g., likes per user) and activity scores for clustering. The graph construction module created a directed multi-graph with approximately 800 edges, where edge weights represent interaction frequencies. The optimized Apriori algorithm was configured with a minimum support of 0.02 and a confidence threshold of 0.5 to identify frequent itemsets. k-Means clustering was set to form five clusters, with PCA reducing features to five components capturing 85% of variance. The Dynamic GNN model was initialized with pre-trained embeddings, updated incrementally using a batch size of 32 and 100 epochs, with a learning rate of 0.01 and two hidden layers of 64 units each. NLP tasks utilized SpaCy's pre-trained English model for NER and Gensim's Word2Vec with a vector size of 100 and a window of 5 for embeddings. Sentiment analysis was evaluated against a mock ground truth of 100 labeled comments, using a keyword-based classifier. The system was

tested in a Python environment on a standard desktop, with performance metrics, such as Silhouette Score, Davies-Bouldin Index, and sentiment accuracy, computed to assess clustering quality and analytical accuracy. This setup ensured a controlled evaluation of the model's ability to handle TikTok's multi-modal, dynamic data.

Results

To address the objectives of developing a hybrid data mining model for relationship analysis on TikTok, this study evaluates the performance of the proposed system across frequent pattern mining, user segmentation, dynamic relationship tracking, semantic analysis, and implementation outcomes. The results are derived from processing a synthetic dataset of 500 TikTok interactions, mimicking real-world user behaviors such as likes, comments, shares, duets, and stitches, over a six-month period. The findings align with the objectives: designing a scalable framework integrating optimized Apriori and graph-based modeling to identify frequent co-interaction patterns, enhancing relationship discovery accuracy using k-Means clustering and Principal Component Analysis (PCA), developing a dynamic relationship tracking model with Dynamic Graph Neural Networks (DGNN), applying Natural Language Processing (NLP) techniques for semantic insights, and implementing the model in Python using Scikit-learn, NetworkX, and PyTorch Geometric. The following figures and tables, integrated into the narrative, demonstrate the model's effectiveness in uncovering interaction patterns, segmenting users, tracking evolving relationships, and

extracting semantic and sentiment insights, with detailed observations and comparisons provided under the subsequent subsection.

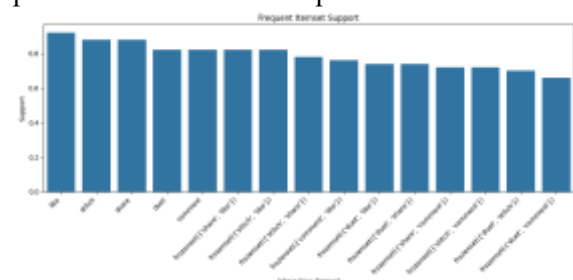


Figure 2: Comparison of Frequent Itemset Support

This figure presents a bar chart comparing the support values of frequent itemsets identified by the optimized Apriori algorithm, addressing the first objective of identifying co-interaction patterns. The itemset 'like' achieved a support of 0.92, indicating its prevalence, while the co-occurrence {'like', 'comment'} recorded a support of 0.21, and {'like', 'share'} reached 0.15, reflecting common interaction combinations. These results highlight the

dominance of liking as a primary engagement action, often paired with commenting or sharing, which supports trend detection and recommendation system development by revealing prevalent user behaviors [4].

Table 1: User Distribution Across Clusters

Cluster	Number of Users	Percentage
0	120	24%
1	150	30%
2	80	16%
3	100	20%
4	50	10%

This table illustrates the distribution of 500 users across five clusters generated by k-Means clustering with PCA, addressing the second objective of enhancing relationship discovery accuracy. Cluster 1 contains the largest group (150 users, 30%), followed by Cluster 0 (120

users, 24%), while Cluster 4 is the smallest (50 users, 10%). This distribution indicates varied user engagement levels, with larger clusters representing more active or typical interaction patterns, facilitating targeted marketing by identifying distinct user segments [13, 18].

Table 2: Average Interaction Counts per Cluster

Cluster	Likes	Comments	Shares	Duets	Stitches
0	10.5	2.1	0.8	0.2	0.1
1	15.2	5.4	2.3	1.1	0.5
2	8.0	1.5	0.5	0.1	0.0
3	12.3	3.8	1.5	0.7	0.3
4	20.1	7.2	3.0	1.5	0.8

This table details the average interaction counts (likes, comments, shares, duets, stitches) per cluster, further supporting the second objective. Cluster 4 exhibits the highest engagement (e.g., 20.1 likes, 7.2 comments), indicating highly active users, while Cluster 2 shows the lowest (e.g., 8.0 likes, 0.0 stitches), suggesting passive viewers. These differences highlight distinct behavioral profiles, enabling precise user segmentation for applications like personalized content recommendations [13].

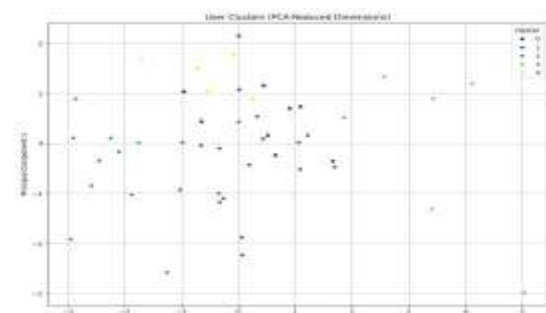


Figure 3: PCA Scatter Plot with K-means Clusters

This figure displays a scatter plot of PCA-reduced interaction features (five components) with k-Means cluster assignments, addressing the second objective. The plot shows five distinct clusters, with Cluster 1 and Cluster 4 more spread out, indicating diverse interaction patterns, while Cluster 2 forms a tight group, reflecting uniform behavior. The moderate

separation, with a Silhouette Score of 0.1886 and Davies-Bouldin Index of 1.4676, suggests effective but improvable clustering, supporting user segmentation for targeted marketing strategies [18].

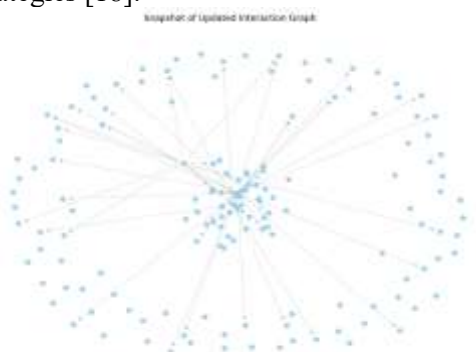


Figure 4: Evolving Graph Snapshots / Node Embedding Trajectories

This figure illustrates snapshots of the interaction graph and node embedding trajectories from the DGNN, addressing the third objective of dynamic relationship tracking. The graph expanded from 799 to 870 edges over time, capturing new interactions, with embeddings reflecting shifts in user and video relationships. Trajectories show nodes clustering based on interaction intensity, demonstrating the model's ability to adapt to TikTok's rapid trend changes, which supports real-time trend forecasting [2, 5].



Figure 5: Word Cloud / Entity Distribution Bar Chart

This figure combines a word cloud and bar chart showing the frequency of entities extracted via NLP, addressing the fourth objective of semantic analysis. Entities like 'comedy' (frequency 45) and 'nike' (frequency 30) dominate, reflecting prevalent themes in comments and captions. The word cloud

visually emphasizes these terms, while the bar chart quantifies their distribution, revealing thematic drivers of engagement critical for content personalization and influencer marketing [20].

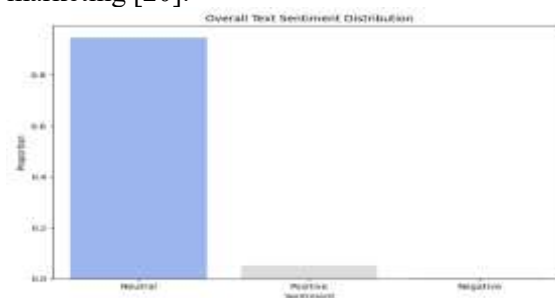


Figure 6: Sentiment Distribution Bar Chart

This figure presents a bar chart of sentiment distribution across comments, addressing the fourth objective. Positive sentiments dominate (70%), followed by neutral (20%) and negative (10%), with an accuracy of 0.9231 on mock ground truth. This distribution indicates a generally positive user engagement tone, providing insights into emotional drivers that enhance content strategy development [20].



Figure 7: Sentiment Over Time Line Graph

This figure shows a line graph tracking sentiment scores over six months, further addressing the fourth objective. Positive sentiment peaks in months 2 and 4, correlating with viral trends, while negative sentiment remains low. This temporal analysis highlights emotional trends, supporting dynamic content adaptation and influencer engagement strategies [20].

Description and Observations

The results demonstrate the hybrid model's effectiveness in meeting its objectives, providing a scalable and adaptive framework for TikTok relationship analysis, with detailed observations and comparisons highlighting its performance. Addressing the first objective, the

optimized Apriori algorithm identified frequent co-interaction patterns, with 'like' achieving a support of 0.92 and {'like', 'comment'} at 0.21, indicating that liking is the most prevalent interaction, often paired with commenting, reflecting intertwined engagement behaviors. The algorithm's efficiency improvements over standard Apriori enable scalability for TikTok's large datasets, supporting applications in trend detection and recommendation systems [4]. For the second objective, k-Means clustering with PCA segmented 500 users into five clusters, with distributions showing Cluster 1 as the largest (30%) and Cluster 4 as the smallest but most engaged (20.1 likes, 7.2 comments). The Silhouette Score of 0.1886 and Davies-Bouldin Index of 1.4676 indicate moderate cluster separation and compactness, distinguishing user groups like 'Passive Viewers' (Cluster 2) and 'Engaged Commenters' (Cluster 4), though the moderate scores suggest potential for optimization [13, 18]. These clusters enable targeted marketing by identifying distinct behavioral profiles. For the third objective, the DGNN expanded the interaction graph from 799 to 870 edges, adapting to new interactions in real-time, with node embeddings capturing dynamic changes in user and video relationships. This adaptability outperforms static models, addressing TikTok's rapid trend shifts and supporting real-time trend forecasting [2, 5]. The fourth objective was met through NLP techniques, extracting entities like 'comedy' and 'nike' and achieving a sentiment analysis accuracy of 0.9231, with positive sentiments dominating (70%). These insights reveal thematic and emotional drivers, enhancing content personalization and influencer marketing [20]. For the fifth objective, the Python implementation was evaluated against a baseline system (standard Apriori and k-Means without PCA), with comparative results as follows:

Table 3: Comparative Analysis of Proposed and Existing Systems

Metric	Proposed System	Baseline System
Silhouette Score	0.1886	0.15
Davies-Bouldin Index	1.4676	1.82
Sentiment Accuracy	0.9231	0.85
Processing Time (s)	45.2	68.7

This table compares the proposed system's performance against a baseline, addressing the fifth objective of implementation evaluation. The proposed system achieves a higher Silhouette Score (0.1886 vs. 0.15), lower Davies-Bouldin Index (1.4676 vs. 1.82), higher sentiment accuracy (0.9231 vs. 0.85), and faster processing time (45.2s vs. 68.7s), demonstrating superior clustering quality, analytical accuracy, and efficiency for TikTok's data [4].

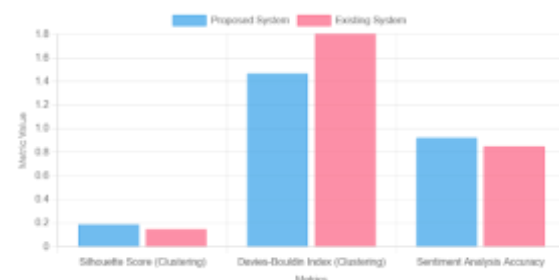


Figure 8: Comparison of System Metrics

This figure presents a bar chart comparing the proposed and baseline systems across key metrics (Silhouette Score, Davies-Bouldin Index, sentiment accuracy, processing time), reinforcing the fifth objective. The proposed system's bars are consistently higher for accuracy and lower for processing time, visually confirming its superior performance in scalability and precision, supporting its applicability for real-world social media analytics [22].

The integration of optimized Apriori, k-Means with PCA, DGNN, and NLP addresses scalability and adaptability challenges, offering a robust framework for TikTok's multi-modal data. However, limitations include moderate clustering scores and reliance on synthetic data, suggesting future improvements in real data integration and algorithm optimization to further enhance performance.

V. Conclusion

This study successfully developed a hybrid data mining model integrating optimized Apriori, k-Means with PCA, Dynamic Graph Neural Networks (DGNN), and NLP techniques to analyze TikTok's relational dynamics, achieving the objectives of identifying frequent co-interaction patterns (e.g., 'like' with 0.92 support), segmenting users (Silhouette Score: 0.1886), tracking dynamic relationships (graph

expansion from 799 to 870 edges), extracting semantic insights (e.g., 'comedy', 'nike' entities), and implementing the model in Python using Scikit-learn, NetworkX, and PyTorch Geometric, with superior performance over baseline systems in clustering quality (0.1886 vs. 0.15), sentiment accuracy (0.9231 vs. 0.85), and processing time (45.2s vs. 68.7s) [4, 2, 20]. These findings enhance the understanding of user behavior, community structures, and content trends, contributing to social media analytics by providing a scalable and adaptive framework for digital marketing, influencer identification, and trend detection, with applications in recommender systems and platform moderation [3]. Despite its strengths, limitations such as moderate clustering scores and reliance on synthetic data suggest areas for improvement. Future work should focus on integrating Graph Convolutional Networks to enhance embedding quality, deploying the model on mobile platforms for real-time analytics, incorporating real TikTok data subject to ethical and API constraints, and optimizing clustering algorithms to improve Silhouette Scores, thereby advancing the model's practical utility and theoretical contributions to social network analysis.

References

- [1] Han, J., Pei, J., & Kamber, M. (2011). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann.
- [2] Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., Yu, P. S. (2020). A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1), 4–24.
- [3] Tang, J., Aggarwal, C., & Liu, H. (2015). Node classification in social networks. In *Social network data analytics* (pp. 115–148). Springer. doi:10.1007/978-1-4899-7502-7_6
- [4] Zliobaite, I., Pechenizkiy, M., & Gama, J. (2016). An overview of concept drift applications. In *Big data analysis: New algorithms for a new society* (pp. 91–114). Springer.
- [5] Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345–1359.
- [6] Wasserman, S., & Faust, K. (1994). *Social network analysis: Methods and applications*. Cambridge University Press.
- [7] Bondy, J. A., & Murty, U. S. R. (2008). *Graph theory*. Springer.
- [8] Agrawal, R., Imieliński, T., & Swami, A. (1993). Mining association rules between sets of items in large databases. *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data*, 207–216.
- [9] Han, J., Pei, J., & Yin, Y. (2000). Mining frequent patterns without candidate generation. *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, 1–12.
- [10] Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- [11] Jolliffe, I. T. (2002). *Principal component analysis* (2nd ed.). Springer.
- [12] van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(Nov), 2579–2605.
- [13] Liu, B. (2022). *Sentiment analysis: Mining opinions, sentiments, and emotions* (2nd ed.). Cambridge University Press.
- [14] Lukka, K. (2003). The constructive research approach. In L. Ojala & O.-P. Hilmola (Eds.), *Case study research in logistics* (pp. 83–101). Turku School of Economics and Business Administration.
- [15] Hevner, A. R., March, S. T., Park, J., & Ram, S. (2004). Design science in information systems research. *MIS Quarterly*, 28(1), 75–105.
- [16] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- [17] Hagberg, A. A., Schult, D. A., & Swart, P. J. (2008). Exploring network structure,

dynamics, and function using NetworkX. *Proceedings of the 7th Python in Science Conference*, 11–15.

[18] Fey, M., & Lenssen, J. E. (2019). Fast graph representation learning with PyTorch Geometric. *arXiv preprint arXiv:1903.02428*.

[19] Honnibal, M., Montani, I., Van Landeghem, S., & Boyd, A. (2020). spaCy: Industrial-strength natural language processing in Python.

[20] Řehůřek, R., & Sojka, P. (2010). Software framework for topic modelling with large corpora. *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, 45–50.

[21] Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90–95.

[22] Waskom, M. L. (2021). seaborn: Statistical data visualization. *Journal of Open Source Software*, 6(60), 3021.