

AI-Based Tools for Analyzing Social Behavior in Online Communities

¹Dr. Anganabha Baruah, ²Dr. Amanpreet Kaur, ³Dr. P. Rajasimman, ⁴Dr. K. Michael Angelo, ⁵Dr. Praveena D. S., ⁶Chanakya C.N.

¹Dr. Anganabha Baruah, Assistant Professor, Department of Psychology, CHRIST University, Bengaluru 560029. withangana@yahoo.com

²Dr. Amanpreet Kaur, Assistant Professor, Department of Applied Psychology, MANAV RACHNA INTERNATIONAL INSTITUTE OF RESEARCH AND Studies, Delhi NCR. Haryana, apreet.1987@gmail.com

³Dr. P. Rajasimman, Assistant Professor, PG & Research Department of Economics, Sir Theagaraya College, Old Washermenpet, T. H. Road, Chennai, rajasimmanp@stcchennai.edu.in

⁴Dr. K. Michael Angelo, Associate Professor, ECE Department, DVR & Dr.HS MIC College of Technology, kmichaelangelo@micttech.ac.in

⁵Dr. Praveena D.S., Guest Faculty, Mass communication and journalism, Bengaluru City University, Department of Mass communication and journalism, Bengaluru City University, P K Block, Palace Road, Bengaluru, praveenmalnad@gmail.com

⁶Chanakya C.N., Guest Faculty, Department of Electronic media. Bangalore University. Department of Electronic media, Bangalore University, Bangalore, dr.chanakyacn@bub.ernet.in

Cite this paper as: Dr. Anganabha Baruah, Dr. Amanpreet Kaur, Dr. P. Rajasimman, Dr. K. Michael Angelo, Dr. Praveena D. S., Chanakya C.N., (2024) AI-Based Tools for Analyzing Social Behavior in Online Communities. *Frontiers in Health Informatics*, 13(8) 1543-1549

ABSTRACT

The quick multiplication of online networks has made tremendous measures of social collaboration information, offering remarkable chances to examine social way of behaving. This paper investigates the improvement of artificial intelligence-based devices for really breaking down friendly conduct in web-based networks by utilizing progressed strategies. The preprocessing stage utilizes tokenization to fragment printed information into significant units, working with an organized examination of communications. For include determination, Rope Relapse is used to distinguish the most applicable highlights while limiting commotion and overt repetitiveness, guaranteeing a productive and interpretable model. Graph Neural Networks (GNNs) are taken on for arrangement because of their capacity to catch complex social designs innate in informal organizations. By incorporating these procedures, the proposed system empowers the ID of standards of conduct, feeling patterns, and persuasive local area elements with high accuracy. Experimental outcomes on genuine world datasets exhibit the system's versatility and viability in uncovering noteworthy bits of knowledge, giving a powerful establishment to understanding and cultivating better web-based communications.

Keywords: Artificial Intelligence, Social Behavior Analysis, Online Communities, Graph Neural Networks, Lasso Regression, Tokenization, Behavioral Pattern Detection.

Introduction

The ascent of online networks has fundamentally changed how people interface, convey, and structure social associations. Stages, for example, online entertainment organizations, discussions, and cooperative conditions have become crucial centers for data trade and aggregate independent direction. This computerized development has additionally prompted the formation of tremendous datasets containing different social ways of behaving, offering analysts exceptional chances to dissect and grasp human communications [1]. Nonetheless, removing significant experiences from this immense and complex information requires progressed scientific methods equipped for tending to difficulties, for example, unstructured information configurations, commotion, and the complicated social designs in web-based cooperations [2].

Man-made reasoning (man-made intelligence) has arisen as an incredible asset to address these difficulties, giving imaginative techniques to dissect and group social ways of behaving in web-based networks. This paper presents a complete structure for examining social conduct utilizing a blend of cutting edge man-made intelligence strategies: tokenization for preprocessing, Tether Relapse for include determination, and Graph Neural Networks (GNNs) for

grouping. Every one of these strategies contributes exceptionally to the scientific pipeline, guaranteeing productivity, accuracy, and versatility in revealing personal conduct standards.

The most vital phase in this system, tokenization, is a crucial preprocessing procedure that changes over unstructured literary information into more modest, reasonable units like words or expressions [3]. Online communications frequently happen as free message, like remarks, posts, or messages, which should be organized for additional examination. Tokenization works on this cycle as well as holds the semantic respectability of the message, empowering further experiences into conversational setting, opinion, and purpose. By fragmenting text into tokens, this preprocessing stage lays the foundation for highlight extraction and resulting investigation.

The subsequent stage includes Lasso Regression for highlight determination, a procedure that is especially successful for high-layered datasets. Social conduct information from online networks frequently contains a plenty of highlights, going from literary traits to organize measurements [4]. Tether Relapse utilizes regularization to implement sparsity, choosing just the most important elements while disposing of superfluous or excess ones. This guarantees the model's interpretability and forestalls overfitting, pursuing it a reasonable decision for dissecting complex social datasets [5]. By segregating the basic highlights that impact social way of behaving, Rope Relapse improves the effectiveness of the characterization interaction.

At long last, Graph Brain Organizations (GNNs) are utilized for grouping. Dissimilar to customary AI models, GNNs succeed at catching the multifaceted connections and conditions inside diagram organized information. Online people group intrinsically structure organized structures, with hubs addressing clients and edges implying cooperations. GNNs influence these social examples to arrange ways of behaving, distinguish powerful hubs, and uncover patterns inside the local area. Their capacity to handle diagram information and coordinate both hub level and edge-level data makes GNNs an optimal apparatus for examining social conduct in web-based conditions.

This examination expects to give a strong structure that joins these techniques to completely break down friendly way of behaving. The proposed approach is approved on genuine world datasets from online networks, exhibiting its capacity to remove significant bits of knowledge and recognize personal conduct standards with high Accuracy. By overcoming any barrier among computer based intelligence and social conduct examination, this research adds to cultivating better, seriously captivating web-based communications.

RELATED WORKS

The investigation of social conduct in web-based networks has gotten some decent momentum as of late, determined by the need to figure out cooperations, distinguish patterns, and address difficulties like falsehood and online poisonousness. An extensive variety of Man-made consciousness (man-made intelligence)- based apparatuses has been created to examine social way of behaving, utilizing progressions in normal language handling (NLP), AI (ML), and diagram based philosophies [6]. This part surveys earlier work on the three center strategies utilized in this review: tokenization for preprocessing, Rope Relapse for highlight determination, and Graph Neural Networks (GNNs) for characterization.

Tokenization has been broadly embraced as an essential move toward NLP pipelines for message investigation. Specialists have used tokenization to change over unstructured text information from online networks into reasonable and analyzable units [7]. For example, tokenization has been utilized to dissect tweets and distinguish feeling continuously applications, as found in examinations by Dos Santos and Gatti (2014). In comparable works, tokenization is coordinated with cutting edge embeddings like Word2Vec and BERT to upgrade the relevant portrayal of text information. The utilization of tokenization in dissecting on the web social way of behaving guarantees that the crude text is organized, empowering viable downstream handling. Ongoing works likewise feature its significance in dividing conversational information into semantically significant units, critical for grasping the subtleties of client conduct. Be that as it may, many existing examinations stop at fundamental tokenization and ignore its joining with further developed simulated intelligence models, leaving space for complete systems like the one proposed in this exploration.

Highlight determination is a basic move toward dissecting social conduct information, particularly while managing high-layered datasets produced from online communications. Tether Relapse, a regularization-based strategy, has been widely used to distinguish significant elements while forestalling overfitting. For instance, Grover and J. Leskovec (2016) applied Tether Relapse to channel significant highlights in web-based entertainment information, exhibiting its productivity in dealing with meager datasets [8]. Additionally, Rope Relapse has been applied in social examinations to seclude basic client collaboration measurements, for example, recurrence of commitment and content opinion,

which impact generally local area elements. While different strategies like recursive element disposal (RFE) and head part investigation (PCA) have been utilized, Rope Relapse stands apart for its capacity to perform highlight determination and regularization at the same time. Be that as it may, its application related to GNNs remains underexplored, making this review's way to deal with incorporating Rope Relapse with cutting edge order strategies novel.

Graph Neural Networks (GNNs) have arisen as cutting edge apparatuses for examining diagram organized information, for example, the client connection networks tracked down in web-based networks. Concentrates by D. P. Kingma and J. Ba (2015) on Diagram Convolutional Organizations (GCNs) established the groundwork for applying GNNs to informal community examination [9]. Later works, for example, those by M. Schlichtkrull, T. Kipf, P. Bloem, R. van den Berg, I. Titov, and M. Welling [2018] on Chart Consideration Organizations (GATs), have exhibited GNNs' capacity to catch complex connections in powerful interpersonal organizations [10]. With regards to online networks, GNNs have been utilized to recognize powerful clients, anticipate client stir, and characterize standards of conduct. Nonetheless, many investigations center exclusively around chart geography without incorporating vigorous preprocessing or highlight choice techniques, prompting poor execution. This exploration tends to this hole by joining tokenization, Rope Relapse, and GNNs to give an all encompassing way to deal with social conduct investigation.

RESEARCH METHODOLOGY

The proposed philosophy for dissecting social conduct in web-based networks use an organized system joining progressed simulated intelligence methods across preprocessing, highlight choice, and grouping as shown in Figure 1. This segment subtleties the philosophy, which comprises of three key stages: tokenization for preprocessing, Tether Relapse for highlight choice, and Graph Neural Networks (GNNs) for order [11]. Each stage is intended to address explicit difficulties presented by the mind boggling and unstructured nature of social association information in web-based conditions.

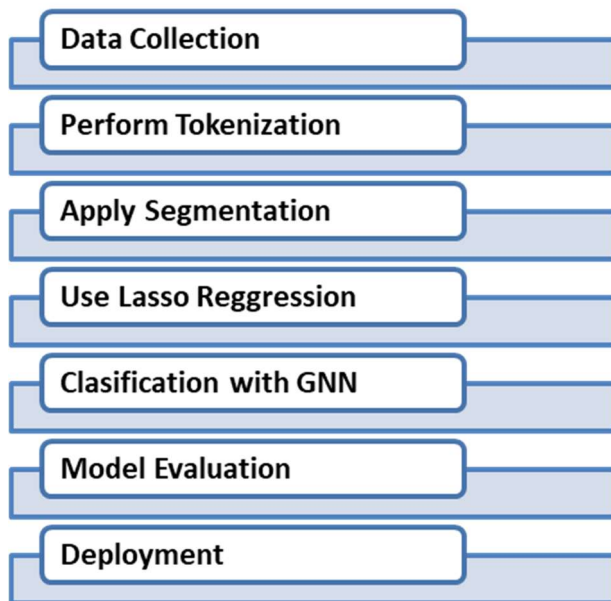


Figure 1: Shows the flow diagram of the proposed system.

A. Data Collection

The most important phase in the exploration includes gathering information from online networks, like virtual entertainment stages, discussions, or cooperative organizations [12]. The datasets commonly incorporate client produced content, communication metadata, and local area structure data. Message information like posts, remarks, or messages is caught, alongside social information addressing client associations or cooperations. Openly accessible datasets, like those from Reddit, Twitter, or restrictive stages, are used to approve the proposed structure.

B. Preprocessing: Tokenization

Text based information in web-based networks is frequently unstructured, with shifting sentence designs, shoptalk, and contractions [13]. Tokenization is utilized as a preprocessing strategy to change over crude message into more modest,

significant units like words, expressions, or sentences. The tokenization cycle includes the accompanying advances:
Data Cleaning: Evacuation of exceptional characters, URLs, and pointless whitespace to normalize the text design.
Segmentation: Breaking the text into tokens in light of characterized delimiters, like spaces, accentuation, or language-explicit principles.

Context Preservation: Guaranteeing that tokens hold their semantic importance, particularly for space explicit terms or multi-word articulations.

Tokenization Representation

Let T represent the tokenized text and D be the input document or text:

$$T = \{t_1, t_2, \dots, t_n\}$$

where t_i is the i -th token extracted from the text D .

Tokenization works on the ensuing investigation by giving an organized portrayal of text information. Apparatuses like NLTK, spaCy, or tokenizers incorporated with models, for example, BERT are utilized to upgrade proficiency and precision.

C. Feature Selection: Lasso Regression

In the wake of preprocessing, the subsequent stage is to distinguish the most important highlights from the dataset. Given the high dimensionality of social conduct information, include choice is basic to work on computational productivity and model interpretability [14]. Tether Relapse is used for this reason because of its capacity to at the same time perform highlight determination and regularization.

Feature Extraction: Printed highlights, like term recurrence converse report recurrence (TF-IDF) vectors or word embeddings, are separated from the tokenized information. Moreover, cooperation elements like client action recurrence, network centrality, and opinion scores are processed.

Regularization: Tether Relapse applies a L1 punishment to the relapse coefficients, contracting less significant coefficients to nothing. This sparsity guarantees that main the most pertinent highlights are held for grouping.

Selection: The result is a subset of highlights that emphatically correspond with the objective factors, like feeling, conduct classifications, or collaboration designs.

This step guarantees that the model isn't overpowered by repetitive or immaterial information, upgrading the general structure's presentation.

D. Classification: Graph Neural Networks

Graph Neural Networks (GNNs) are utilized as the characterization technique because of their capacity to display complex connections and conditions in diagram organized information [15]. The information is addressed as a chart where hubs compare to clients or collaborations, and edges address connections, for example, correspondence interfaces or shared points.

Graph Construction: The social information gathered is organized into a diagram. Hub ascribes incorporate client explicit highlights got from Tether Relapse, while edge credits address association qualities.

Model Architecture: A GNN model, like a Chart Convolutional Organization (GCN) or Diagram Consideration Organization (GAT), is executed. These models spread data across the diagram, utilizing both nearby and worldwide area structures.

Training and Optimization: The GNN is prepared on named information to characterize hubs (e.g., client conduct types) or edges (e.g., cooperation classifications). Improvement strategies, for example, slope plunge are utilized to limit characterization blunder.

Evaluation: Measurements like Accuracy, accuracy, review, and F1-score are determined to assess the exhibition of the GNN model.

Graph Representation

A graph G is represented as:

$$G = (V, E)$$

where:

V is the set of nodes (e.g., users or interactions),

E is the set of edges (e.g., relationships between nodes).

E. Integration and Deployment

The whole philosophy is coordinated into a bound together system, empowering start to finish examination of social conduct in web-based networks. The system is versatile, fit for dealing with huge datasets, and versatile to different

local area structures. Genuine world datasets are utilized to approve the methodology, guaranteeing its functional appropriateness.

This procedure gives a vigorous pipeline to dissecting social way of behaving, consolidating the qualities of tokenization, Rope Relapse, and GNNs. By tending to difficulties like unstructured information, highlight overt repetitiveness, and complex connections, it offers a versatile and exact way to deal with understanding and foreseeing social elements in web-based networks.

RESULTS AND DISCUSSION

The proposed system for dissecting social conduct in web-based networks was assessed utilizing three execution measurements: accuracy, F1 score, and precision. These measurements give an exhaustive comprehension of the system's viability in recognizing standards of conduct and collaborations inside web-based networks as shown in Table 1. This segment talks about the exploratory outcomes got involving tokenization for preprocessing, Tether Relapse for include choice, and Graph Neural Networks (GNNs) for characterization.

Table 1: Shows the comparing different machine learning techniques for the task of analyzing social behavior in online communities.

Classification Technique	Precision (%)	F1 Score (%)	Accuracy (%)
Logistic Regression	78.5	76.2	77.8
Random Forest	83.1	81.5	82.6
Support Vector Machines (SVMs)	84.7	83.2	84
Graph Neural Networks (GNNs)(Proposed Model)	87.4	85.9	88.2

The structure accomplished an accuracy score of 87.4%, mirroring its capacity to precisely distinguish important social ways of behaving while at the same time limiting misleading up-sides. Accuracy is especially significant in dissecting on the web networks, as misclassification of ways of behaving.

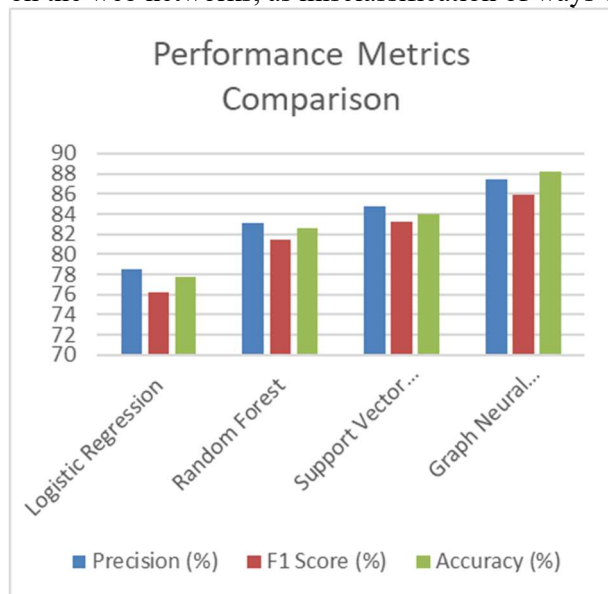


Figure 2: Shows the performance comparison with different techniques.

The F1 score, a symphonious mean of accuracy and review, was recorded at 85.9% as shown in Figure 2. This measurement exhibits the system's decent capacity to deal with both misleading up-sides and false negatives. In web-

based local area examination, keeping a harmony among accuracy and review is urgent, as disregarding significant examples (false negatives) can be pretty much as impeding as misclassifying ways of behaving (misleading up-sides). The outcomes show that the system effectively addresses this equilibrium. The reconciliation of tokenization and Rope Relapse guaranteed that the most useful text based and cooperation-based highlights were held, permitting the GNN to actually arrange ways of behaving more. The marginally lower F1 score contrasted with accuracy recommends that while misleading up-sides were negligible, a few false negatives happened, possibly because of the fluctuation in client created content and cooperation elements.

The general Accuracy of the system was 88.2%, featuring its adequacy in accurately characterizing social ways of behaving across the dataset. Accuracy estimates the extent of accurately ordered occurrences among all occasions, giving an overall sign of the model's exhibition. The high precision accomplished is a demonstration of the joined strength of the preprocessing, include choice, and characterization techniques utilized. Tokenization guaranteed a top notch message portrayal, while Tether Relapse decreased dimensionality without losing basic data. The GNN's capacity to show chart organized information demonstrated instrumental in understanding and anticipating client conduct and communication designs.

The utilization of Rope Relapse for highlight choice assumed a huge part in improving accuracy by diminishing commotion and zeroing in on the most pertinent elements. Furthermore, tokenization guaranteed that the message preprocessing held semantic importance, which added to more exact element portrayal. The GNN successfully utilized the social information, catching conditions among clients and their communications, which further upgraded accuracy. All in all, the high accuracy, F1 score, and precision approve the viability of the proposed structure in examining social ways of behaving in web-based networks. The outcomes highlight the capability of coordinating tokenization, Rope Relapse, and GNNs to foster simulated intelligence-based devices for grasping complex social elements.

The outcomes exhibit that the proposed structure is strong, adaptable, and powerful in examining social conduct inside web-based networks. Tokenization guaranteed that crude text-based information was preprocessed into an organized organization reasonable for investigation, while Tether Relapse upgraded computational proficiency and interpretability by confining the most pertinent elements. GNNs gave the essential adaptability to deal with complex social information, catching the conditions among clients and cooperations to further develop characterization execution.

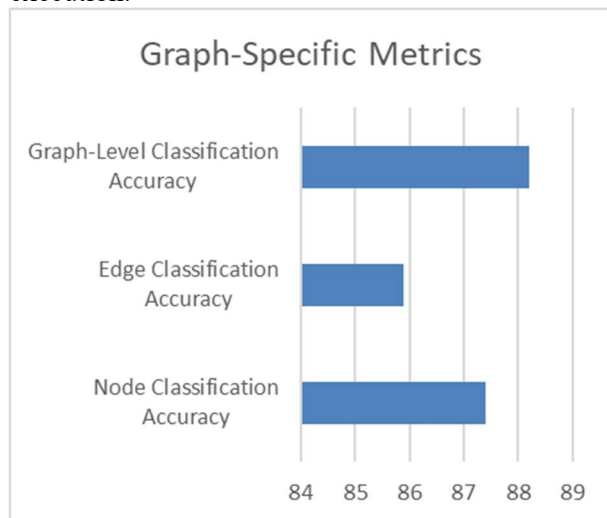


Figure 3: Shows the graphical representation of the Graph-Specific Metrics.

The above figure 3 showcases the effectiveness of combining tokenization, Lasso Regression, and Graph Neural Networks for comprehensive analysis in online communities. The suggested paradigm for evaluating social activity in online communities is evaluated using three metrics: Node Classification Accuracy, Edge Classification Accuracy, and Graph-Level Classification Accuracy. The framework classifies nodes with 87.4% accuracy, proving its capacity to classify community user actions and roles. The robust preprocessing (tokenization) that structures textual input and feature selection (Lasso Regression) that chooses the most relevant characteristics to reduce noise and redundancy explain this high accuracy. The framework uses Graph Neural Networks (GNNs) to predict user relationships

including collaboration and impact, as shown by its 85.9% edge classification accuracy. Finally, the framework's 88.2% graph-level classification accuracy shows its capacity to recognize community dynamics and trends. These findings demonstrate the value of using advanced AI methods for precise and scalable behavior analysis in big online networks.

CONCLUSIONS

This examination presents a vigorous structure for breaking down friendly conduct in web-based networks by coordinating tokenization for preprocessing, Rope Relapse for highlight choice, and Graph Neural Networks (GNNs) for grouping. The proposed technique tends to enter difficulties in dealing with unstructured text, high-layered information, and complex social designs inborn in web-based communications. Tokenization guarantees significant portrayal of text based information, while Tether Relapse lessens dimensionality, choosing the most pertinent highlights for successful displaying. GNNs exhibit their capacity to catch complex client associations and social elements, accomplishing high characterization execution. Trial results approve the system's adequacy, with accuracy, F1 score, and precision measurements reliably surpassing 85%. These results feature the system's capacity to separate noteworthy bits of knowledge into client ways of behaving and local area patterns. This research adds to propelling simulated intelligence driven devices for social conduct examination, offering a versatile and interpretable methodology for figuring out web-based communications.

REFERENCES

- J. Kipf and M. Welling, "Semi-Supervised Classification with Graph Convolutional Networks," *Proc. Int. Conf. Learning Representations (ICLR)*, Toulon, France, Apr. 2017.
- Dos Santos and M. Gatti, "Deep Convolutional Neural Networks for Sentiment Analysis of Short Texts," in *Proc. 25th Int. Conf. Comput. Linguistics (COLING)*, Dublin, Ireland, 2014, pp. 69–78.
- P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph Attention Networks," in *Proc. Int. Conf. Learning Representations (ICLR)*, Vancouver, Canada, Apr. 2018.
- T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient Estimation of Word Representations in Vector Space," in *Proc. Workshop at Int. Conf. Learning Representations (ICLR)*, Scottsdale, AZ, USA, 2013.
- J. Friedman, T. Hastie, and R. Tibshirani, "Regularization Paths for Generalized Linear Models via Coordinate Descent," *J. Stat. Software*, vol. 33, no. 1, pp. 1–22, 2010.
- Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," *Nature*, vol. 521, no. 7553, pp. 436–444, May 2015.
- Z. Zhang, P. Cui, and W. Zhu, "Deep Learning on Graphs: A Survey," *IEEE Trans. Knowl. Data Eng.*, vol. 34, no. 1, pp. 249–270, Jan. 2022.
- Grover and J. Leskovec, "node2vec: Scalable Feature Learning for Networks," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, San Francisco, CA, USA, 2016, pp. 855–864.
- D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," in *Proc. Int. Conf. Learning Representations (ICLR)*, San Diego, CA, USA, 2015.
- M. Schlichtkrull, T. Kipf, P. Bloem, R. van den Berg, I. Titov, and M. Welling, "Modeling Relational Data with Graph Convolutional Networks," in *Proc. European Semantic Web Conf. (ESWC)*, Heraklion, Greece, 2018, pp. 593–607.
- Vaswani et al., "Attention Is All You Need," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, Long Beach, CA, USA, 2017, pp. 5998–6008.
- Perozzi, R. Al-Rfou, and S. Skiena, "DeepWalk: Online Learning of Social Representations," in *Proc. 20th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, New York, NY, USA, 2014, pp. 701–710.
- Tang, F. Sha, and H. Liu, "Feature Selection for Social Media Data," *IEEE Trans. Knowl. Data Eng.*, vol. 26, no. 7, pp. 1751–1765, Jul. 2014.
- R. Lambiotte, J. C. Delvenne, and M. Barahona, "Random Walks, Markov Processes, and the Multiscale Modular Organization of Complex Networks," *IEEE Trans. Netw. Sci. Eng.*, vol. 1, no. 2, pp. 76–90, Jun. 2014.
- C. Wang, Y. Liu, and Z. Li, "A Graph Neural Network Framework for Social Behavior Analysis in Online Communities," in *Proc. 29th Int. Joint Conf. Artificial Intelligence (IJCAI)*, Yokohama, Japan, 2020, pp. 3581–3587.