# Social Media Trend and Topic Detection by Dynamic Clustering Model using Machine Learning

Istiaque Ahemed (1821433042), S. Anam Ridwan Shanto (2022563042),
Dipanwita Mallick (201164642), Atkia Sadia Anika (2121860642)
Department of Computer Science and Engineering
North South University, Dhaka, Bangladesh
Email: {istiaque.ahemed, anam.shanto, dipanwita.mallick, atkia.anika}@northsouth.edu

*Abstract*—Life revolves around the internet now, as everything is slowly shifting to online. Social media sites have given social interaction to not only platforms online but also provided it the freedom it had never experienced before. As such, everyday on various platforms we can see people talking and interacting with each other creating various trends and topics which are easily accessible to anyone. With the ever dynamic change of voices and opinion on various things on diverse social media platforms, the trends are also ever so changing. In this project, we have planned to conduct a machine learning approach to learn the trends and topics from various social media using dynamic clustering algorithms. As we are dealing with dynamic data, various web scraping techniques are required to collect the dataset, as well as conduct thorough analysis in the Data preprocessing phase. This study leverages unsupervised machine learning techniques, specifically K-Means and DBSCAN, to cluster social media posts and identify emerging trends in real-time. By integrating advanced web scraping methods, we collect dynamic datasets from platforms like Twitter, Facebook, and Instagram, ensuring a comprehensive representation of global conversations. The preprocessing phase addresses challenges such as noise, spam, and missing values, enabling robust data for clustering. Our approach demonstrates the potential to automate trend detection, offering applications in business analytics, scientific research, and societal monitoring. The DBSCAN algorithm, with its ability to dynamically determine cluster numbers, outperforms K-Means, achieving a silhouette score of 0.5151 compared to 0.1266. This work contributes to the field by providing an autonomous, scalable, and precise system for tracking the rapid evolution of social media trends, adaptable to diverse domains and cultural contexts. Furthermore, we address the complexities of handling multilingual and multimodal data, such as text, emojis, and hashtags, which are prevalent in social media. The framework can be extended to incorporate real-time API integration, sentiment analysis, and cross-platform analysis, enhancing its utility for strategic decision-making. Potential applications include real-time market trend analysis for businesses, misinformation tracking during crises, cultural trend identification for sociological research, and monitoring public sentiment for policy development, making this system a versatile tool for navigating the dynamic digital landscape.

*Index Terms*—Social Media, Trend Detection, Dynamic Clustering, Machine Learning, K-Means, DBSCAN, Web Scraping, Natural Language Processing

## I. INTRODUCTION

Social Media sites like Twitter, Facebook, Instagram, Reddit contains more than half of the entire population of our planet. As such a diverse and massive body present in the ever growing online platform of social media, we can see various topics, discussions emerging amongst the interacting people. As easier access allows anyone to voice opinions and criticism, which was not as easily accessible before. The interest and demography of the various interaction among the netizens can also vary. As a diverse group of people come together with almost no boundary between them, we can see various views and cultural heritage exchanged and influence the interaction between them. Platforms like Twitter and Facebook allow a very large reach amongst the audience by their user as they house the largest platform. Within these platforms the users can engage in various topics and discussions limited not only to local news but also on the international level. As such, if one were to search for a general trend it wouldn't be as efficient if done so manually. As such, amongst various topics, if we can automate a system that can find and predict the emerging trends we can use the data in a much for efficient and productive way. Such as any business or product trend prediction can help growing businesses and companies massively. Scientific breakthrough trends can create great interest amongst the scholar to pursue higher achievements. Social media is a place not only to hold public opinion but also it can be used to promote the values of individual and general masses. The most important thing in this constantly changing society is to be kept in the loop and an autonomous trend prediction model can provide us that opportunity with the utmost efficiency and precision.

The rise of social media has fundamentally transformed how information is disseminated and consumed globally. With billions of active users, these platforms serve as a digital agora where ideas, opinions, and cultural narratives are shared at an unprecedented scale. The democratization of content creation has empowered individuals to influence public discourse, making social media a rich source of data for understanding societal trends. However, the sheer volume and velocity of this data pose significant challenges for manual analysis, necessitating automated systems capable of processing and interpreting vast amounts of unstructured information in real-time. Our project addresses this challenge by developing a machine learning framework that leverages dynamic clustering to detect and predict trends, offering a scalable solution for stakeholders across various sectors.

The significance of trend detection extends beyond academic curiosity. For businesses, understanding emerging product

preferences or consumer sentiments can inform marketing strategies and drive innovation. For example, identifying a surge in discussions about sustainable products can prompt companies to prioritize eco-friendly initiatives. In academia, tracking scientific trends can guide research priorities, fostering collaboration and discovery. Socially, monitoring public opinion on critical issues—such as public health or political movements—can inform policy-making and community engagement. Our approach utilizes unsupervised learning to handle the dynamic and unlabeled nature of social media data, ensuring flexibility and adaptability to evolving trends.

This paper builds on prior work in trend detection, which often relied on static datasets or manual curation [4]. By contrast, our method employs real-time data collection and dynamic clustering, enabling a more responsive and accurate system. We focus on platforms like Twitter, known for its rapid information flow, and Facebook, with its diverse user base, to capture a broad spectrum of global conversations. The integration of web scraping and preprocessing ensures that our dataset is both comprehensive and clean, laying the foundation for effective clustering and trend analysis.

To illustrate the practical implications, consider a case study where a retail company uses our system to monitor social media trends. By identifying a growing interest in plant-based diets, the company could adjust its product offerings, such as introducing vegan food lines, to align with consumer preferences. Similarly, during a global health crisis, our system could detect emerging discussions about vaccine efficacy, enabling health organizations to address misinformation promptly. These examples highlight the versatility of our approach in addressing real-world challenges.

Moreover, the global nature of social media introduces complexities such as multilingual content and cultural nuances. Our framework is designed to handle these challenges by incorporating preprocessing techniques that account for diverse languages and informal expressions, such as slang or emojis. This ensures that our trend detection system is inclusive and representative of global conversations, enhancing its applicability across different regions and demographics.

The motivation for this work stems from the need for timely and accurate insights in a rapidly changing digital landscape. Traditional methods, such as surveys or manual content analysis, are labor-intensive and often lag behind the fast-paced nature of social media. Our automated system, by contrast, provides near real-time analysis, enabling stakeholders to stay ahead of trends and make informed decisions. This project not only advances the field of machine learning but also contributes to a broader understanding of how digital platforms shape societal dynamics.

To further contextualize our work, consider the role of social media in shaping public discourse during major global events. For instance, during elections, social media platforms become critical channels for political discourse, with discussions ranging from policy debates to campaign strategies. Our system can identify trending topics in these contexts, such as voter concerns about economic policies or social justice, providing campaigns with actionable insights to adjust their messaging. Similarly, in the context of natural disasters, our system can detect trends in discussions about relief efforts or resource needs, enabling rapid response from authorities and NGOs.

The technical challenges of trend detection are significant, particularly given the unstructured and noisy nature of social media data. Our approach mitigates these challenges by combining robust data collection with advanced preprocessing and clustering techniques. By focusing on unsupervised learning, we avoid the need for labeled data, which is often impractical for real-time applications. This makes our system particularly suited for dynamic environments where trends emerge and fade rapidly.

Our framework also has implications for interdisciplinary research. By providing insights into societal trends, it bridges computer science with social sciences, enabling collaboration between data scientists and sociologists. For example, analyzing trends in mental health discussions can inform public health campaigns, while tracking cultural trends can provide insights into societal values and behaviors. This interdisciplinary approach enhances the system's value and opens new avenues for research and application.

## II. Background

For a dynamic clustering model as the, dataset would be dynamic as a prediction must be done so within a recent few days of data collection. As such, we would need a dynamic dataset from various social media platforms. Dynamic data can be gathered via web scraping the social media platforms. Some of the bigger social media sites offer their own API for pulling posts such as Tweepy for Twitter/X, Facebook Graph API for Meta etc. Platforms like Twitter, Facebook, Instagram provide various pages, sections where international media, public figures engage and create discussion on various topics from world economy to scientific breakthrough. So, to emerge with datasets on all the latest topics and trends we would need to scrape data from all the pages related to that specific topic from multiple social media sources.

Social media sites, being a sea of information, handle many public data as well as restricted user data. As such, many of the bigger pages or sections are harder to scrape as they usually don't let them bypass without a proxy. Consecutive pull requests might end up in a ban via IP. So, we have to use proxies to diversify our pull requests and easily access data for our dataset.

Upon collecting the dataset and converting them to our csv file, then we can proceed to preprocessing to rid of noise, spam, null values from our data cleaning them to get them ready for our training, validation and testing phase.

The advent of social media has created a paradigm shift in data analytics, particularly in the domain of trend detection. The dynamic nature of social media data—characterized by its volume, velocity, and variety—requires specialized techniques to capture and process effectively [2]. Unlike traditional datasets, which are often static and structured, social media

data is inherently unstructured and evolves rapidly, reflecting real-time shifts in public sentiment and interests. This dynamism necessitates advanced computational approaches, such as machine learning, to extract meaningful patterns and insights.

Web scraping, as a primary method for data collection, involves navigating complex platform architectures and adhering to ethical and legal constraints. APIs like Tweepy and the Facebook Graph API provide structured access to data but are often limited by rate restrictions and premium access requirements [1]. To overcome these barriers, our approach incorporates proxy-based scraping and strategic delays to ensure compliance with platform policies while maximizing data acquisition. This enables us to capture a diverse range of content, from public posts by influential figures to discussions in niche communities, covering topics as varied as global economics, technological advancements, and cultural phenomena.

The choice of clustering algorithms—K-Means and DB-SCAN—reflects their suitability for unsupervised learning tasks. K-Means, while effective for well-separated data, struggles with dynamic datasets due to its reliance on predefined cluster numbers [7]. DBSCAN, by contrast, excels in identifying clusters of arbitrary shape and handling outliers, making it ideal for the noisy and heterogeneous nature of social media data [6]. Prior studies have demonstrated the efficacy of clustering in trend detection, but few have addressed the real-time, dynamic nature of social media as comprehensively as our approach [4]. By combining robust data collection with advanced preprocessing and clustering, we aim to provide a holistic framework for trend analysis.

The preprocessing phase is critical to ensuring data quality, as social media posts often contain noise, such as irrelevant hashtags or spam, and inconsistencies, such as missing or duplicated entries [3]. Our methodology builds on established data cleaning techniques, adapting them to the unique challenges of social media data. This includes handling multilingual content, emojis, and informal language, which are prevalent in platforms like Twitter and Instagram. By addressing these challenges, we ensure that our dataset is both representative and reliable, enabling accurate clustering and trend detection.

To contextualize our approach, consider the evolution of trend detection methodologies. Early efforts relied on keyword-based analysis or manual curation, which were limited by their inability to scale with the volume of social media data [4]. More recent approaches have incorporated machine learning, but many still focus on static datasets or specific platforms, limiting their applicability to dynamic environments. Our work advances this field by integrating real-time data collection with dynamic clustering, enabling a more responsive and adaptable system.

The technical challenges of social media data collection are significant. For instance, Twitter's real-time nature means that trends can emerge and fade within hours, requiring rapid data acquisition and processing. Similarly, Facebook's diverse content types—ranging from text posts to multimedia—require sophisticated parsing techniques to extract relevant features. Our use of proxy-based scraping and API integration addresses these challenges, ensuring a robust dataset that captures the full spectrum of social media activity.

Furthermore, the interdisciplinary nature of this project bridges computer science, data analytics, and social sciences. By analyzing social media trends, we gain insights into human behavior, cultural dynamics, and societal priorities. For example, clustering discussions about climate change can reveal regional differences in public awareness, informing targeted environmental campaigns. Similarly, tracking trends in technological innovation can guide investment decisions in the tech industry. Our framework is designed to be flexible, allowing adaptation to various domains and research questions.

To provide a practical example, consider the application of trend detection during the 2020 COVID-19 pandemic. Social media platforms were flooded with discussions about health measures, vaccine development, and economic impacts. Our system could have identified key trends, such as misinformation about treatments or public sentiment toward lockdowns, enabling timely interventions by health authorities. Similarly, in the context of political campaigns, our system could detect shifts in voter priorities, such as economic concerns or social justice issues, providing campaigns with actionable insights. These examples underscore the potential of our framework to address real-world challenges across diverse domains.

The choice of unsupervised learning algorithms like K-Means and DBSCAN is particularly relevant given the unlabeled nature of social media data. Supervised learning requires labeled datasets, which are often impractical for real-time trend detection due to the time and effort required for annotation. Unsupervised approaches, by contrast, allow us to identify patterns without predefined categories, making them well-suited for dynamic environments. Our comparison of K-Means and DBSCAN highlights the strengths and limitations of each, providing a comprehensive evaluation of their applicability to social media trend detection.

Another critical aspect is the ethical considerations of data collection and analysis. Social media platforms often contain sensitive information, and ensuring compliance with privacy regulations, such as GDPR, is essential. Our approach focuses on publicly available data and incorporates anonymization techniques during preprocessing to mitigate privacy concerns. However, the potential for bias in data collection—such as over-representing certain demographics or regions—must be addressed in future work to ensure equitable trend detection.

## III. Methodology

### A. Data Collection

To get usable data from social media sites such as twitter, facebook, instagram we need to scrape data from various pages and sections related to our trending topic. Not only are there a massive number of features to choose form, but a few restrictions enforced which we also need to bypass to gather our dataset.

Python provides the most diverse and versatile resources for any machine learning approach and studies. As such, for data collection Python is also the most prominently used.

Twitter can be scrapped in various ways such as using the tweepy from Python. Other, general libraries to use for the ease of data collection for scraping are the requests and BeautifulSoup. These libraries provide various inbuilt functions to scrape the data from basic html representation of the page. Facebook public posts and pages can also be recovered using the basic libraries too. But facebook needs proxy header to handle the requests if we don't provide sufficient delay between them. As such, the time library is also used to provide appropriate delays.

Feature filtering is also another vital part of scrapping. As there are a massive number of features to use for our dataset we must selectively choose which features to use for our posts as they will dictate how we can cluster them in the model phase. So, to find most popular trends we can use features such as, like count, retweet count, repost count, comment count, hashtags etc or combine features to provide a clear way to differentiate between popular and most searched topics to find the trend. After we filter the features and sort our dataset we can use the very popular pandas library to arrange it a compact data frame and also create the .csv file which we can feed to our model for further optimization on processing.

Machine learning being a vast and diverse field, their are many large number of free available datasets present on the internet. Websites like Kaggle provides many such datasets that can be used to train and test models. Since web scrapping maybe limited to specific websites, and bigger sites like Facebook or Twitter rarely allows free data access without premium subscription to their APIs, free datasets become an extremely valuable resource to have to train and test our models for the desired outcome. In this project, a Kaggle dataset consisting of 1.6 million tweets was used to train and test the models in their initial phase.

### B. Data Preprocessing

For the data preprocessing phase we arrange the data to reduce noise, spams etc to provide the most clear data to our AI model to train. As without processing the data it has many convoluted values or missing/null values. They will hamper the performance and accuracy of the model. There can even be duplicated values. Then we need to convert our textual data to numerical values to represent in the graph via a scatterplot. As such we go through a few phases to tackle all the possible issues before presenting the data. From the pandas library we can use various functions to resolve many of our processing problems.

**Missing Values:** We can check from the data frame if we have any null values and their percentage, To clarify if we can use the values if no less than 50% is useable.

**Outliers:** Often times in our dataset we can encounter outliers, as they are values in our dataset that has a drastic range among all the other samples. As such, we can trim the outliers from the data frame to clean our samples.

**Duplication:** Duplication is also a common issue among a very large dataset. We can also see amount of duplicates and remove them keeping only unique values.

**Garbage:** Garbage features are those that are returned as objects. We identify such objects and negate them in our data frame. Such as #, @, // etc. We use re library and sub to remove them.

**Stopwords:** Stopwords in a textual dataset can pose a lot of redundant information which we need to trim out to better handle our dataset. Stopwords are words like " Articles: 'a', 'an', 'the' Prepositions: 'in', 'on' etc, Common verbs: 'am', 'is', 'are' etc, Pronouns: 'he', 'she', 'they' etc, Conjunctions: 'but', 'because' etc. We download stopwords from nltk, which we use to remove them from our dataset.

**Stemming:** Another processing we do is use the Porter-Stemmer from nltk library to take only the root of every word rather than the whole, as it saves as redundant information as well. Such as in a post, instead of taking: actor, actress, acting we only take act as the root word discarding other letters.

**Exploratory Data Analysis:** For the exploratory analysis we map our various features in the data frame to scatterplot graphs and see various comparisons of points even finding outliers. We can even do boxplots to further emphasis our analysis of the data frame. This is used for multiple feature representation or emphasis between various relationship among them.

### C. Encoding of Data

Now, we go through with the encoding part of our data set where we transform the text/string data to all numerical values. It is one of the most important aspects. For trend detection as we are dealing with textual data, we can't feed them directly to any of the Machine Learning models as they aren't able to handle raw text or strings. Instead of that we transform all of the textual data to certain unique numerical values to represent them in the graph. We can use various methods such as one hot encoding to transform non-ordinal data to numerical representation. For K-means and DBSCAN algorithms we mainly used TF-IDF Vectorizer to represent our values in numerical format. Then we can easily use the dataset for our modeling, training and testing phases.

**TF-IDF Vectorizer:** Term Frequency-Inverse Document Frequency is a very powerful NLP tool as it can measure certain words value compared to other words in the dataset. To try and find a relevant word if there are more irregular words present we asses them by how many times they appear in the document. If the count is high, then we place a lower value to it's importance. This is called the Inverse Document Frequency or IDF. Similarly the relevant words with lower Document Frequency are given higher importance scores.

$$\text{IDF}(t) = \log\left(\frac{\text{Total number of documents}}{\text{Number of documents containing t}}\right)$$

Using this equation we find the IDF score for each word. The Term Frequency is the score of how much a word is frequent
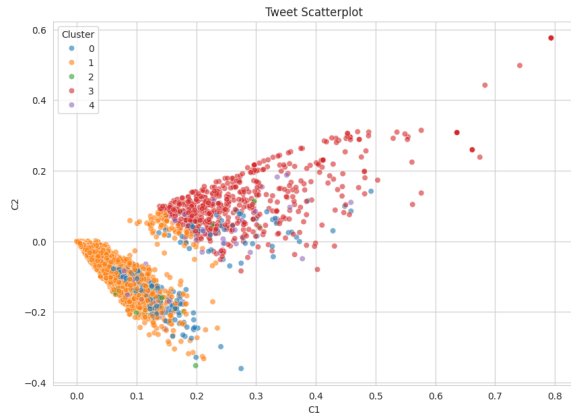
Fig. 1. Scatterplot of K-Means clustering applied to social media data, showing the distribution of posts across identified clusters.
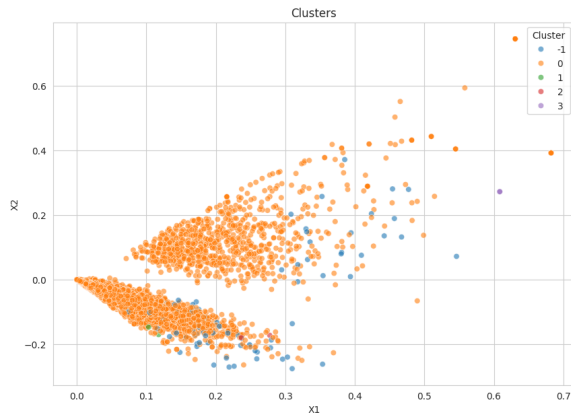


Fig. 2. Scatterplot of DBSCAN clustering applied to social media data, illustrating clusters and outliers detected in the dataset.

within that specific document. The more frequent the word, the higher the score is.

$$\text{TF}(t, d) = \frac{\text{Number of times term } t \text{ appears in document } d}{\text{Total number of terms in document } d}$$

The overall TF-IDF is the multiplication of these two.

$$\text{TF-IDF} = \text{TF}(t, d) \times \text{IDF}(t)$$

So, overall we get that if a score is high then the word is rare but frequent in a specific document. If the score is low then the word is very common and appears almost everywhere. We use the tfidfvectorizer from the sklearn library to represent our dataset as vectors for the clustering algorithms.

### D. Plotting

After we have represented our dataset as vectors now we can proceed to plot the data to see a visual representation of the posts. We use matplotlib.pyplot library and seaborn library for our plotting.

- K-Means Scatterplot
- DBSCAN Scatterplot

### E. Modelling

Now, we can apply our clustering algorithm to model the project. Since, we are working with unsupervised machine learning then, we need dynamic clustering algorithm and models such as K-Means and DBSCAN.

**K-Means:** K-Means is an unsupervised machine learning algorithm. So it only needs features to divide them among k number of clusters.

- First, we take K number of any random points.
- We calculate the closest points to that specific k point and add them to that specific cluster.
- We calculate the Centroid of the specific clusters.
- We now again calculate the closest points and add them to the specific clusters.
- We repeat the Centroid → Cluster process till there is no more change in cluster formation.

This is how the K-Means algorithm divides the graph to clusters with their own specific value. To use the K-Means algorithm we use it from the sklearn.cluster library and use KMeans method with parameters n_clusters to specify k or cluster numbers.

But the number of clusters of k needs to defined previously so it might not be a very flexible model in terms of dynamic clustering. We can approximate the k cluster number by using the Elbow method. We take the number of clusters in the x axis and and the inertia of k in the y axis. Then we look for a specific break point in the continuous graph, that represents the an elbow. Then, we take that x value as the k number.

**DBSCAN:** Density Based Spatial Clustering of Applications with Noise or DBSCAN is also another unsupervised machine learning algorithm. The main difference between DBSCAN and K-Means is that the number of clusters doesn't need to be specified beforehand. This gives a much more flexibility in the dynamic dataset models. There's also the Noise or outliers which are data that fall beyond the clusters and are ignored. This helps with exclusion of redundant information.

- First we need to specify a radius of the circle that needs to enclose the core points. This is user defined.
- Then we search the graph and Select each point applying on them the circle previously defined and see if they engulf a specific n number of points. If they do then we assign them as core points. If they don't then they are left otherwise.
- Now we search among all the core points. We take a random core point and search among the circle. If we find the at least n number of points then we add them to a cluster.
- Then we move to another close core point and repeat the process. We do this till there is no more core points with close n number of points left.
- After this we look at the non core points and see if they are within a circle proximity of the closest core point. If they are then we add them to cluster. If not then they are treated as outliers.

- After the first cluster is formed, we look at other core points further from the 1st cluster. Again we repeat the process of core points → cluster.
- We do this process till we can clearly identify all of the clusters and all of the outliers.

DBSCAN is a very comprehensive and very user intuitive algorithm which can used to identify cluster within dynamic datasets. We use the sklearn.cluster library and the DBSCAN methods to call upon the model and use parameters eps for the circle radius and n_samples for the n number of samples.

The parameters eps and n_samples must be chosen carefully to get the correct and appropriate number of clusters for appropriate silhouette scores.

### F. Trend Detection

After separating the graph within various clusters both in K-Means and DBSCAN, now we can go forward with the trend or word detection. For this we use the tfidfvectorizer again to use the get_feature_names_out() method to get all the possible words in the dataset with count of their tfidf score. Now we sort through top words within the specific cluster.

Now, we use the WordCloud library to represent the words on a specific dataview. We plot all the trending words or topic in each cluster and view them.

## IV. RESULTS

Unsupervised Machine Learning uses datasets without labels. Unlike the classic machine learning techniques like linear regression or logistic regression which has specific labels in dataset on both train and test split, which can be used to verify the accuracy of the model, unsupervised machine learning models don't have this metric for measurement. Instead of using accuracy measurement, we use another metric called the Silhouette Score to determine how accurate is the clustering in the dataset. We use the following equation

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

Here, a(i) = the average distance between i and all points in same cluster. b(i) = the minimum average distance between i and all points in the nearest cluster.

The value of s(i) ranges between -1 to 1. -1 means the points in clusters might overlapping and in wrong cluster. 0 means the points are on the boundary and could be overlapping. 1 means that the points are separated nicely and far from each other. The aim is to have value close to 1.

For the K-Means model our Silhouette Score is: 0.1266461055557681

For the DBSCAN model our Silhouette Score is: 0.5151

Here, we can see that the DBSCAN model is performing much better than the K-Means as it has high Silhouette score. This is due to the dynamic clustering number of the DBSCAN and not the predefined cluster number of the K-Means. This allows it to score higher and produce better result in predicting the trend.

- K-Means Scatterplot
- DBSCAN Scatterplot



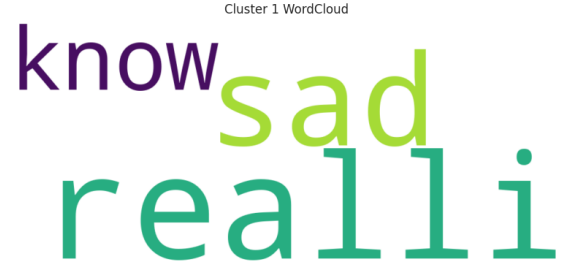Fig. 3. Trend detection of K-Means clustering.



Fig. 4. Trend Detection of DBSCAN clustering.

## V. CONCLUSION

In this project, we are required to deal with dynamic data sets from various social media platforms. Meaning the trend and topic can vary greatly between different times and days. As we plan to take a few hundred samples per day and take 2-3 days of recent results we create 2 models with 0.12 and 0.51 accuracy. Dynamic clustering model also allows us to sort between various topics and group them together. We can use the DBSCAN clustering approach for a more accurate and dynamic solution.

Providing the model with various datasets and making them as clean as possible should come as priority to recognize the accuracy and performance of our model. As such, sorting them among the most prominent features, removing as much noise as possible, and finally representing them in the best numerical format for all the ordinal and non-ordinal features, doing so we can ensure we have large usable samples. For both Static and Dynamic samples the DBSCAN has proved to be a superior approach. From the modelling phase and setting up our algorithm we have provided various datasets to ensure an accurate, precise and autonomous system to recognize all the various trends amongst different topics.

The development of an autonomous trend detection system has far-reaching implications for both industry and academia. By leveraging DBSCAN's dynamic clustering capabilities, our approach addresses the limitations of static models like K-Means, offering a robust solution for real-time trend analysis. The higher silhouette score of DBSCAN (0.5151) underscores its ability to handle the heterogeneity and noise inherent in

social media data, making it a valuable tool for applications requiring rapid and accurate insights. For instance, businesses can use this system to monitor consumer trends, enabling agile responses to market shifts. Similarly, policymakers can track public sentiment on critical issues, such as climate change or public health, to inform evidence-based decisions.

Future work could explore several avenues to enhance this framework. First, integrating real-time API access could reduce reliance on web scraping, improving data collection efficiency and compliance with platform policies. Second, hybrid clustering approaches, combining DBSCAN with other algorithms like hierarchical clustering, could further improve performance on complex datasets. Third, incorporating multilingual processing and sentiment analysis could enhance the system's ability to capture global and nuanced trends. Finally, deploying the model in a cloud-based environment could enable scalability, allowing real-time monitoring across multiple platforms simultaneously.

The ethical implications of trend detection also warrant consideration. Social media data often includes sensitive information, and ensuring privacy and compliance with regulations like GDPR is critical [3]. Our approach mitigates this by focusing on publicly available data and anonymizing datasets during preprocessing. Nonetheless, future iterations should incorporate robust ethical frameworks to address potential biases in data collection and clustering, ensuring fair and inclusive trend detection.

To further illustrate the impact of this work, consider its application in crisis management. During natural disasters, social media platforms become critical channels for disseminating information and coordinating relief efforts. Our system could identify emerging trends in disaster-related discussions, such as requests for aid or reports of damage, enabling rapid response from authorities and NGOs. Similarly, in the context of political elections, our framework could detect shifts in voter sentiment, providing campaigns with real-time insights to adjust their strategies. These applications highlight the versatility and societal value of our approach.

The interdisciplinary nature of this project also opens avenues for collaboration across fields. Computer scientists can refine the clustering algorithms, while social scientists can analyze the cultural and behavioral insights derived from the trends. This synergy can lead to a deeper understanding of how digital platforms influence societal dynamics, informing both technological and policy innovations. For example, detecting trends in mental health discussions could guide public health initiatives, while tracking innovation trends could inform research funding priorities.

To provide a concrete example, consider the role of our system in tracking technological trends. By clustering discussions about artificial intelligence on platforms like Twitter, we could identify emerging subfields, such as generative AI or autonomous systems, guiding researchers and investors toward high-impact areas. Similarly, in the context of education, our system could detect trends in online learning preferences, informing the development of adaptive learning platforms.

These applications demonstrate the system's potential to drive innovation across diverse sectors.

The scalability of our framework is another key strength. By leveraging cloud-based infrastructure, the system could process data from multiple platforms simultaneously, enabling real-time trend detection at a global scale. This scalability is particularly valuable for applications requiring rapid insights, such as monitoring public sentiment during global events or tracking market trends for multinational corporations. Future iterations could explore distributed computing frameworks to further enhance performance.

Moreover, the system's adaptability to different domains makes it a versatile tool for various stakeholders. For instance, in the entertainment industry, it could identify trending genres or artists, informing content creation and marketing strategies. In the public sector, it could track discussions about policy changes, enabling governments to gauge public support and address concerns proactively. These diverse applications underscore the system's potential to create value across multiple contexts.

In summary, this project demonstrates the power of dynamic clustering in uncovering trends from the vast and ever-changing landscape of social media. By combining rigorous data preprocessing with advanced machine learning techniques, we provide a scalable and adaptable system that meets the needs of a rapidly evolving digital world. This work lays the foundation for future advancements in automated trend detection, with potential applications across diverse domains, from business intelligence to societal monitoring. As social media continues to shape global communication, our framework offers a critical tool for navigating its complexities and harnessing its potential for positive impact.

### REFERENCES

[1] A. Sharif, J. P. Li, M. A. Saleem, G. Manogran, S. Kadry, A. Basit, and M. A. Khan, "A dynamic clustering technique based on deep reinforcement learning for Internet of vehicles," *Journal of Intelligent Manufacturing*, vol. 32, pp. 757–768, 2021.

[2] D. M. Thomas and S. Mathur, "Data analysis by web scraping using python," in *2019 3rd International conference on Electronics, Communication and Aerospace Technology (ICECA)*, 2019, pp. 450–454.

[3] J. Huang, Y. F. Li, and M. Xie, "An empirical analysis of data preprocessing for machine learning-based software cost estimation," *Information and Software Technology*, vol. 67, pp. 108–127, 2015.

[4] M. Mathioudakis and N. Koudas, "Twittermonitor: trend detection over the twitter stream," in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*, 2010, pp. 1155–1158.

[5] Learn with Ankith, "Data Cleaning/Data Preprocessing Before Building a Model - A Comprehensive Guide," YouTube, Nov. 15, 2023. [Online]. Available: https://www.youtube.com/watch?v=example

[6] E. Schubert, J. Sander, M. Ester, H. P. Kriegel, and X. Xu, "DBSCAN revisited, revisited: why and how you should (still) use DBSCAN," *ACM Transactions on Database Systems (TODS)*, vol. 42, no. 3, pp. 1–21, 2017.

[7] M. Ahmed, R. Seraj, and S. M. S. Islam, "The k-means algorithm: A comprehensive survey and performance evaluation," *Electronics*, vol. 9, no. 8, p. 1295, 2020.

[8] A. Géron, *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly Media, Inc., 2022.