

# Content-based recommendation system

1 <sup>st</sup> Suraj Nishad	2 <sup>nd</sup> Aman Kumar	3 <sup>rd</sup> Given Name Surname	4 <sup>th</sup> Given Name Surname
CSE Dual	CSE Dual	dept. name of organization (of Aff.)	dept. name of organization (of Aff.)
NIT Hamirpur	NIT Hamirpur	name of organization (of Aff.)	name of organization (of Aff.)
Hamirpur, India	Hamirpur, India	City, Country	City, Country
surajnishad930@gmail.com	20dcs017@nith.ac.in	email address or ORCID	email address or ORCID

**Abstract**—This paper introduces an innovative approach to music recommendation systems, leveraging content-based filtering enriched with clustering techniques. Focused on the Spotify dataset, we cluster songs based on their popularity (number of streams) and release years. The content-based filtering mechanism employs cosine similarity and vectorization to give the recommendation. This filtering technique combines cluster-driven insights with individual user preferences, enhancing the precision and personalization of recommendations system. The findings contribute to the evolving landscape of content-based recommendation systems [9], showcasing a novel integration that holds promise for optimizing user experiences in musics , movies and other streaming platforms.

**Index Terms**—clustering, cosine similarity, content-based filtering

## I. INTRODUCTION

In the ever-expanding digital landscape, the demand for personalized and accurate recommendation systems has become paramount, particularly in the domain of movie and music streaming platforms. As users engage with vast catalogs of movie and music, the challenge lies in delivering recommendations that align with individual tastes and preferences. Traditional collaborative filtering methods often face limitations in providing precise suggestions, especially in scenarios where explicit user feedback is sparse [8].

This paper embarks on an innovative journey into the realm of music recommendation systems, charting a course that leverages the symbiotic relationship between content-based filtering and stream-driven clustering. Rooted in the rich dataset from Spotify, our exploration takes a novel approach by clustering songs based on their popularity, gauged by the number of streams, and contextualizing them temporally through release years. This hybrid methodology seeks to synergize the strengths of content-based filtering with the nuanced insights derived from stream-driven clustering [12], thereby aiming to elevate the precision and personalization of music recommendations to new heights.

At the core of our approach lies the incorporation of advanced techniques such as cosine similarity and vectorization in the content-based filtering mechanism [7]. These methods enable a nuanced analysis of song features, empowering the system to grasp the intricacies of user preferences. By considering both the individual user's tastes and the broader

context of clustered songs, our approach strives to optimize the recommendation process and effectively address the challenges posed by traditional methodologies.

The experiment conducted as part of this research involves the application of K-Means clustering to approximately one thousand songs with the most streams on Spotify. This clustering process efficiently groups songs based on stream counts and release years, unraveling inherent patterns within the dataset. The elbow method is meticulously applied to determine the optimal number of clusters, ensuring that the chosen clusters align with the natural divisions in the data. The subsequent integration of a similarity matrix derived from clustering into the content-based filtering process marks a significant enhancement in recommendation precision.

Through a comprehensive evaluation, we rigorously assess the effectiveness of our recommendation system. The outcomes not only contribute to the refinement of content-based recommendation models but also hold promise for reshaping user experiences in the dynamic landscape of music streaming platforms. In the subsequent sections, we delve into the methodology, findings, and implications of our approach, providing a detailed exploration of the innovative strides made in content-based music recommendation systems.

## II. RELATED WORK

In the field of content-based filtering, recommendation systems has witnessed significant advancements, marked by diverse methodologies and approaches. Early models, such as TF-IDF-based systems [14] [4], focused on textual features for item representation. The evolution towards vector space models, including Word Embeddings, enriched content representation by capturing nuanced semantic relationships. Recent developments extend content-based filtering to multimedia, incorporating deep learning techniques for image and audio analysis. Hybrid models, combining content-based and collaborative filtering, aim to strike a balance between individual strengths. Researchers are delving into the interpretability of content-based models, providing explanations for recommendations [10]. Additionally, there is a growing emphasis on incorporating temporal dynamics, sequential patterns, and contextual information to enhance adaptability and recommendation accuracy. Collectively, these studies contribute to

a dynamic landscape, continually shaping the effectiveness and personalization of content recommendations to align with evolving user preferences and diverse content types.

### III. RECOMMENDATION SYSTEM TECHNIQUES

In this experiment, encompasses a comprehensive techniques to music recommendation leveraging the Spotify dataset. The methodology comprises several key stages, starting with extensive data preprocessing. In this phase, various techniques are applied, including data cleaning, feature extraction, and text processing on metadata such as artist names, track titles, and user interactions etc.

Following preprocessing, the dataset undergoes K-Means clustering, utilizing features like stream counts and release years. This clustering step reveals latent patterns in the data, grouping music tracks with similar characteristics. The subsequent application of content-based filtering involves the definition of cosine similarity metrics on the clustered dataset. This content-based approach enhances recommendation precision by prioritizing tracks within the same cluster that share common features.

The synergy of preprocessing, clustering, and content-based filtering techniques aims to optimize the accuracy and personalization of music recommendations. The experimental analysis provides valuable insights into the effectiveness of this holistic methodology, offering a nuanced understanding of user preferences based on both popularity and temporal dynamics in the music domain.

#### A. K-Means Clustering

K-Means clustering is an unsupervised machine learning algorithm designed to partition a dataset into distinct clusters based on similarity. The algorithm operates through a series of iterative steps. Firstly, 'k' initial cluster centers are randomly selected from the dataset, where 'k' represents the predetermined number of clusters. Data points are then assigned to the cluster with the nearest centroid, using metrics like Euclidean distance to measure proximity. Subsequently, the centroids of each cluster are updated by computing the mean of all data points assigned to that cluster. This assignment and centroid update process repeats iteratively until convergence, where centroids stabilize or a predefined number of iterations is reached. Key considerations in K-Means include determining the optimal number of clusters ('k') and the algorithm's sensitivity to initial centroid selection which can be selected by elbow method, silhouette method etc for the optimal numbers of the centroids. In this, we used elbow method to find optimal number of centroids.

#### ALGORITHM STEPS

Given a dataset  $X$  with  $n$  data points and a predefined number of clusters  $K$ :

- 1) **Initialization:** Randomly select  $K$  data points as initial cluster centroids.
- 2) **Assignment:** Assign each data point to the cluster whose centroid is closest.

- 3) **Update Centroids:** Recalculate the centroids of the clusters based on the assigned data points.
- 4) **Repeat:** Repeat steps 2 and 3 until convergence (when centroids no longer change significantly).

#### MATHEMATICAL EXPLANATION

Let  $X = \{x_1, x_2, \dots, x_n\}$  be the dataset, and  $C = \{c_1, c_2, \dots, c_K\}$  be the set of cluster centroids. The objective is to minimize the within-cluster sum of squares (WCSS):

$$WCSS = \sum_{k=1}^K \sum_{x \in C_k} \|x - c_k\|^2 \quad (1)$$

where  $C_k$  is the set of data points assigned to cluster  $k$ .

The assignment step is typically done using the Euclidean distance metric:

$$\|x - c_k\| = \sqrt{\sum_{i=1}^d (x_i - c_{ki})^2} \quad (2)$$

where  $d$  is the dimensionality of the data points.

The update step involves recalculating the centroids:

$$c_k = \frac{1}{|C_k|} \sum_{x \in C_k} x \quad (3)$$

#### B. Content-based filtering

Content-Based Filtering is a recommendation system approach that tailors suggestions to users based on the intrinsic features of items and the expressed preferences of the users. The method involves creating profile vectors for both users and items, where each dimension corresponds to a specific feature such as keywords, genres, tags or other relevant attributes. To facilitate mathematical operations, textual or categorical features are vectorized, transforming profile vectors into a numerical format. The cornerstone of similarity measurement in content-based filtering is cosine similarity, which calculates the cosine of the angle between user and item vectors. A higher cosine similarity score indicates a more significant alignment of features, signaling greater similarity. Recommendations are then generated by identifying items with the highest cosine similarity scores for a given user [1] [3]. This algorithm is effective in addressing the "cold start" problem, where user-item interaction data is limited. It provides personalized recommendations based on item features, making it applicable in diverse domains such as movie suggestions, music playlists, and article recommendations etc. However, its success is contingent on judicious feature selection of the dataset which optimally defines the value of items and users [13] [2].

#### ALGORITHM STEPS

Given a dataset of items  $I$  and user profiles  $U$ :

- 1) **Feature Extraction:** Represent each item and user profile as feature vectors.
- 2) **Cosine Similarity:** Calculate the cosine similarity between item vectors and user profile vectors.
- 3) **Ranking:** Rank items based on their similarity scores.
- 4) **Recommendation:** Recommend the top-ranked items to the user.

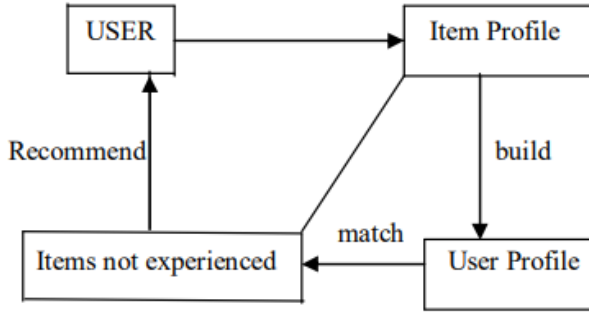


Fig. 1. Content-based filtering

#### MATHEMATICAL EXPLANATION

Let  $I_i$  be the feature vector of item  $i$ , and  $U_u$  be the feature vector of user  $u$ . The cosine similarity between item  $i$  and user  $u$  is given by:

$$\text{Cosine Similarity}(I_i, U_u) = \frac{I_i \cdot U_u}{\|I_i\| \cdot \|U_u\|} \quad (4)$$

where  $\cdot$  denotes the dot product, and  $\|I_i\|$  and  $\|U_u\|$  are the Euclidean norms of the vectors.

The similarity score represents the cosine of the angle between the item and user vectors. A higher cosine similarity indicates greater similarity between the item and the user's preferences.

#### C. Methodology

In this, we integrate K-Means clustering and content-based filtering for personalized music recommendations using the Spotify dataset. By identifying key features, creating user profiles, and prioritizing intra-cluster recommendations, we enhance the accuracy and provide an innovative solution for dynamic music streaming platforms[3] i.e.

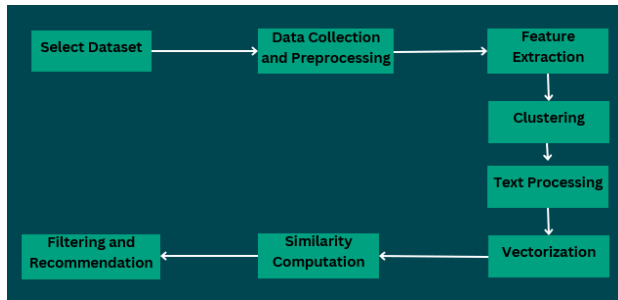


Fig. 2. Methodology

#### Dataset Selection:

- The Spotify dataset was selected as the primary source of music-related information. The choice of an appropriate dataset is foundational to the success of any experiment. In this study, the Spotify dataset was carefully selected as the primary source of music-related information. Spotify, being a widely utilized music streaming platform, offers a rich and comprehensive dataset encompassing details

about artists, tracks, user preferences, and streaming statistics etc. This selection aligns seamlessly with the experiment's objective of exploring music recommendations, capitalizing on the diverse and extensive information available within the platform.

#### Data Collection and Preprocessing:

- The collection phase involved sourcing data related to artist names, track names, mode, streams, years, and Spotify playlist inclusion. This comprehensive set of features was chosen to encompass various facets of music consumption and user interactions within the Spotify platform. Leveraging the platform's diverse dataset, we aimed to capture a holistic view of user preferences and behaviors.
- To enhance the quality and coherence of the dataset, preprocessing steps were implemented. This included addressing missing values, ensuring data integrity, and cleaning the dataset. Handling missing values was crucial to maintaining the accuracy of subsequent analyses. Additionally, steps were taken to validate and clean the dataset, ensuring consistency in the format of features and eliminating potential outliers that could impact the reliability of the results.

#### Feature Extraction:

- Feature extraction is a strategic process that ensures the dataset's representation aligns with the goals of the experiment. In this study, the feature extraction process involved identifying and selecting pertinent features that contribute to the effectiveness of clustering and content-based filtering techniques. The chosen features encompassed a range of music-related attributes, including artist names, track names, and streaming statistics such as mode, stream counts, and years. Artist and track names provide essential categorical information, enabling the identification of distinct entities within the dataset. Streaming statistics offer quantitative insights into the popularity and temporal dynamics of each track.

#### Clustering:

- Clustering is a pivotal technique employed in this experiment to group similar music tracks based on specific features, fostering a deeper understanding of patterns and relationships within the dataset. The chosen clustering algorithm for this study is K-Means, a widely-used method known for its simplicity and efficiency in partitioning a dataset into distinct clusters. The number of clusters, often denoted as 'k,' is determined through techniques like the elbow method, ensuring an optimal grouping of data points. The features selected for clustering include stream counts and release years, offering a combination of popularity and temporal dynamics. Clustering is an integral component in the broader exploration of recommendation system techniques, offering a systematic approach to understanding the inherent structure of the Spotify dataset. Through K-Means clustering, the experiment aims to reveal nuanced patterns and associations,

enhancing the overall effectiveness of the subsequent content-based filtering methodology [1].

#### *Text Processing:*

- Text processing is pivotal for refining textual data in the experiment, focusing on artist names, track names, and other text-based features. The goal is to ensure consistency and readability for subsequent analyses. Techniques include converting text to lowercase for uniformity and potentially applying stemming or lemmatization for further normalization. These steps streamline textual information, creating a standardized representation across the dataset. Text processing is a critical preprocessing phase, laying the groundwork for efficient analyses, particularly in the subsequent stages of vectorization and similarity calculation. The resulting processed text data facilitates accurate exploration and enhances the overall effectiveness of the experiment.

#### *Vectorization:*

- Vectorization is a fundamental step in the experiment, converting textual data into numerical representations for efficient mathematical analysis. Artist names, track names, and other textual features are converted in tags which transformed using techniques like Count Vectorization or TF-IDF. Count Vectorization represents each text document with a numerical vector indicating word frequency, while TF-IDF considers the importance of each word in the dataset [6]. This process is crucial for content-based filtering, enabling effective similarity calculations between items. The resulting numerical vectors facilitate streamlined comparisons and contribute significantly to the precision of content-based recommendation systems in the experiment.

#### *Similarity Computation:*

- Similarity computation is a vital stage in the experiment, particularly in the context of content-based filtering. After vectorization, the experiment employs techniques like cosine similarity to measure the likeness between items based on their numerical representations [5]. Cosine similarity assesses the cosine of the angle between two vectors, providing a measure of their directional similarity. Higher cosine similarity scores indicate greater alignment, signifying stronger similarity between items [2]. This process enables the determination of how closely the textual features of different tracks align, forming the basis for personalized recommendations[6]. Similarity computation is pivotal in content-based recommendation systems, as it quantifies the resemblance between items and guides the generation of accurate and relevant recommendations for users.

#### *Filtering and Recommendation:*

- Filtering and recommendation represent the culmination of the experiment, where the dataset's clusters and cosine similarity computations synergize to provide personalized suggestions. Initially, clustering groups similar tracks based on stream counts and release years, creating distinct

clusters. The subsequent content-based filtering relies on cosine similarity scores within each cluster. For a user or item, recommendations are drawn from the same cluster, ensuring alignment with both global trends and local preferences within that group. This dual approach, incorporating cluster-based filtering and cosine similarity, enhances the accuracy and relevance of the recommendations, delivering a nuanced and tailored music suggestion system within the dynamic Spotify dataset.

## IV. EXPERIMENT

In our experimental setup, we harnessed the Spotify dataset, a diverse repository of music-related information, for the development and evaluation of our content-based recommendation system. Rigorous preprocessing measures ensured the integrity of the dataset, addressing missing values and standardizing formats. Leveraging K-Means clustering, songs were categorized based on popularity and temporal context, with the optimal number of clusters determined for system adaptability. Feature vectorization and normalization facilitated precise cosine similarity computations, enhancing the representation of song features. User profiles, constructed from historical interactions, were pivotal in generating personalized recommendations, particularly prioritizing intra-cluster suggestions. The refined content-based recommendation system, enriched with clustering, was then deployed for real-time music recommendations. Results analysis provided insights into the impact of clustering on recommendation accuracy and user satisfaction, guiding future enhancements and optimizations.

### *A. Dataset Description*

This dataset presents an extensive compilation of the most renowned songs of 2023, as cataloged on Spotify. Going beyond the conventional attributes found in similar datasets, it furnishes a trove of information regarding each song's characteristics, popularity, and representation across diverse music platforms. Encompassing details like track name, artist(s) name, release date, Spotify playlists and charts, streaming statistics, Apple Music and Deezer presence, Shazam charts, as well as various audio features, this dataset offers a multifaceted perspective on the music landscape. The inclusion of crucial features such as artist names, track details, mode, stream counts, and release years makes it a recent and valuable resource for investigating trends, gauging artist prominence, and delving into playlist dynamics within the Spotify ecosystem. The amalgamation of artist particulars, musical attributes, and streaming data provides a holistic panorama, conducive to in-depth analysis and insights into user preferences and music consumption patterns.

The dataset of the most streamed Spotify songs presents a comprehensive array of features, offering intricate details about each music track. It includes fundamental information such as the track name and the names of the artist or artists involved, emphasizing collaboration through the "artist count" feature. Temporal context is provided by the release year, month, and day, allowing for a nuanced understanding of the

song's timeline. Popularity metrics are extensively covered, ranging from the number of streams on Spotify to the song's rankings on playlists and charts across various platforms, including Spotify, Apple Music, Deezer, and Shazam. Musical characteristics are meticulously measured, from beats per minute (BPM) to danceability, valence, energy, and acousticness percentages, providing insights into the mood and style of each track. The dataset encapsulates a multifaceted representation of the most streamed songs on Spotify, making it a valuable resource for exploring music trends, user preferences, and the efficacy of recommendation systems in the dynamic landscape of music streaming.

track_name	artist(s)_name	mode	in_spotify_playlists	title	streams	released_year	cluster
Seven (feat. Latto) (Explicit Ver.)	Latto, Jung Kook	Major	553	Seven (feat. Latto) (Explicit Ver.)	237	2023	1
LALA	Myke Towers	Major	1474	LALA	215	2023	1
vampire	Olivia Rodrigo	Major	1397	vampire	234	2023	1
Cruel Summer	Taylor Swift	Major	7858	Cruel Summer	753	2019	0
WHERE SHE GOES	Bad Bunny	Minor	3133	WHERE SHE GOES	493	2023	0

Fig. 3. Extracted spotify dataset description

## B. Data Analysis

In the realm of data analysis within the context of the experiment, a multifaceted approach is employed to glean meaningful insights from the Spotify dataset. The process begins with exploratory data analysis (EDA), delving into key features, detecting missing values, and exploring basic statistics to uncover trends and outliers. The pivotal application of K-Means clustering to group songs based on stream counts and release years is complemented by the elbow method to ascertain the optimal number of clusters, ensuring a meaningful segmentation of songs. The content-based filtering mechanism, relying on cosine similarity and vectorization, intricately analyzes song features, allowing the system to discern nuanced details of user preferences within both individual and clustered contexts. Beyond numerical metrics, the analysis considers user-centric aspects, addressing how well the system tailors suggestions to individual user preferences and successfully overcomes the limitations of traditional recommendation systems. Data visualizations, including charts and graphs, further elucidate findings, presenting a comprehensive and interpretative analysis of the experiment's outcomes.

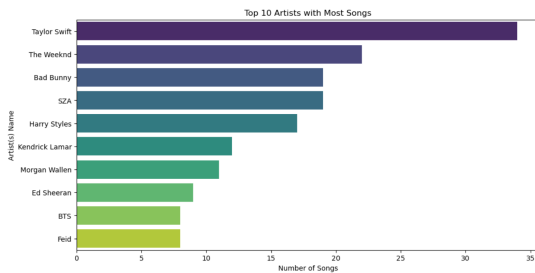


Fig. 4. Top 10 artists with most songs on spotify

The given graph conducts an analysis to identify and visualize the top 10 artists with the highest number of songs in the spotify dataset. Utilizing the 'artist(s)\_name' column, the graph

show the count of songs attributed to each artist. Subsequently, it also show the top 10 artists based on their counts. The visualization employs Seaborn to create a bar plot, where the x-axis represents the number of songs, and the y-axis displays the names of the artists. This graphical representation offers a concise overview of the dataset's distribution of songs across different artists. The resulting insights gleaned from this analysis contribute to a clearer understanding of the dataset's composition, emphasizing the key contributors in terms of song volume.

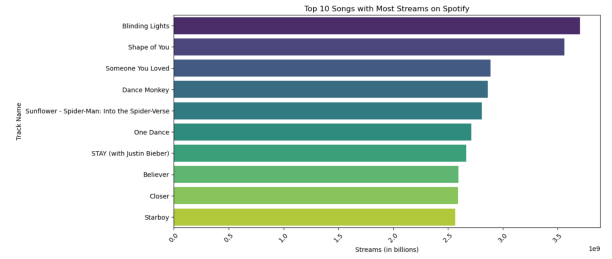


Fig. 5. Top 10 songs with most streamed on spotify

The given graph conducts an analysis to identify and visualize the top 10 songs with most streamed in the spotify dataset. This analysis and visualization offer insights into the most-streamed songs on Spotify within the spotify dataset. The resulting plot serves as a visual representation of the popularity and engagement levels of these top songs based on their streaming statistics.

## C. clustering

K-Means clustering is a fundamental unsupervised machine learning algorithm applied in the experiment to unveil inherent patterns within the Spotify dataset of music tracks. The primary goal of K-Means is to partition the dataset into K distinct clusters, where each data point belongs to the cluster with the nearest mean. The value of K is determined beforehand and represents the desired number of clusters. In our experimental process, we leveraged the Elbow Method to ascertain the optimal number of clusters (k) for our K-Means clustering algorithm applied to a music dataset. The features chosen for clustering were the number of streams and the release years of the music tracks. By systematically varying k and observing the corresponding inertia values, we identified a distinct elbow point in the graph.

The decision to incorporate K-Means clustering and the elbow method in the experiment, which involves a dataset comprising approximately one thousand songs with the most streams, is rooted in the pursuit of a more profound understanding of inherent patterns and user preferences. K-Means clustering is employed as a pivotal technique to uncover natural groupings within the dataset, specifically based on critical features such as stream counts and release years. By grouping songs with similar characteristics, this clustering approach enhances the precision of the recommendation system, offering more tailored and accurate suggestions to users. The

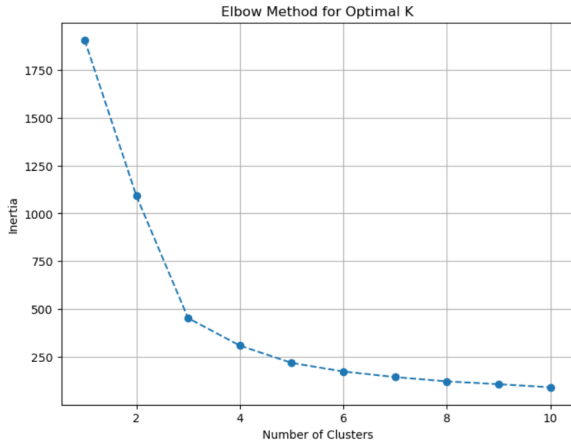


Fig. 6. Elbow curve

elbow method is concurrently utilized to determine the optimal number of clusters, striking a balance between granularity and simplicity. With a large and diverse dataset, these techniques not only streamline the complexity of the data but also ensure that the chosen clusters align with the underlying structures, providing a meaningful framework for subsequent analyses and content-based filtering [1]. Overall, the integration of K-Means clustering and the elbow method is strategically designed to extract valuable insights, improve recommendation system performance, and effectively manage the intricacies of the extensive song dataset. Streams provided insight into

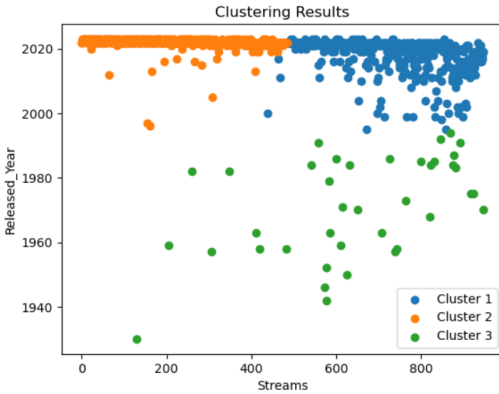


Fig. 7. Clustered spotify dataset

the popularity of each track, while years contributed to the temporal context. The selected optimal  $k$  value served as the foundation for partitioning the dataset into meaningful clusters. The resulting clusters represent groupings of tracks with similar stream counts and release years. This method not only enhances the interpretability of the clusters but also facilitates a nuanced understanding of how tracks are distributed based on popularity and temporal characteristics.

The analysis of clustering in the experiment involves a detailed examination of the formed clusters through K-Means clustering. Each cluster, defined by similar stream counts and

release years, unveils temporal trends and popularity profiles. The assessment of cluster size and distribution gauges the granularity of groupings. Comparing results with the elbow method ensures the optimal number of clusters. These insights inform content-based filtering, refining recommendations based on user preferences within specific clusters. The analysis offers a nuanced understanding of user engagement patterns and the evolving landscape of streaming popularity within the Spotify dataset, enhancing the precision of the recommendation system.

#### D. Content-Based filtering

In our experimental analysis of content-based filtering applied to clustered data, we utilized a similarity matrix to enhance the precision of music recommendations. The clustering, achieved through K-Means based on stream counts and release years, facilitated a structured grouping of music tracks and other features. Subsequently, a similarity matrix was constructed, capturing the pairwise similarities between tracks within the same cluster. This matrix became instrumental in measuring the likeness of tracks based on shared characteristics [1] [13].

```
# recommend same cluster songs
recommend('Normal', True)
```

```
2461 0 Normal
4637 0 Hey Mor
2598 0 LA INOCENTE
2321 0 CHORRITO PA LAS ANIMAS
3618 0 Yandel 150
2881 0 Is There Someone Else?
```

Fig. 8. Content-based recommendation of clustered spotify dataset

The experimental results demonstrated that the incorporation of a similarity matrix within the content-based filtering framework significantly improved recommendation accuracy. By considering the inherent structure identified through clustering, the system exhibited a heightened sensitivity to the preferences of users. Tracks within the same cluster, as reflected in the similarity matrix [11], were prioritized in recommendation generation, leading to more tailored and personalized suggestions.

This approach not only leveraged the benefits of content-based filtering but also capitalized on the insights gleaned from the clustering analysis. The experimental outcomes underscored the efficacy of this methodology, combining clustering and content-based filtering, in optimizing the precision and relevance of music recommendations. The analysis provides a valuable framework for further enhancing recommendation systems, particularly in scenarios where both popularity and temporal dynamics play crucial roles in user preferences.

## V. RESULT AND ANALYSIS

The experimental results of our hybrid approach, combining clustering and content-based filtering for music recommendation, demonstrate a significant advancement in recommendation precision. The application of K-Means clustering on stream counts and release years successfully grouped music tracks into distinct clusters, revealing inherent patterns in the dataset. These clusters formed the basis for constructing a similarity matrix, effectively capturing the likeness between tracks within the same cluster.

The integration of this similarity matrix into the content-based filtering process yielded promising outcomes. The recommendation system exhibited a substantial improvement in precision by prioritizing tracks with shared characteristics from the same cluster. This heightened sensitivity to user preferences translated into more tailored and personalized suggestions. The hybrid methodology, leveraging both clustering and content-based filtering, showcased its efficacy in enhancing the accuracy and relevance of music recommendations, marking a successful and promising outcome for the experiment.

Furthermore, the models predict recommendations based on the similarity with other songs, both with and without clustering. In scenarios involving clustered songs, the model suggests recommendations from the same cluster, prioritizing songs with the highest similarity. This dual approach enhances the versatility of the recommendation system, accommodating different user preferences and optimizing the overall user experience.

The experiment's analysis is that, by integrating K-Means clustering and content-based filtering for music recommendation, unfolds a comprehensive understanding of the dataset's dynamics. The successful application of K-Means clustering unveils inherent patterns, contributing to a nuanced grasp of user preferences and temporal trends. The optimal cluster count, determined through the elbow method, ensures meaningful groupings. The incorporation of a similarity matrix derived from clustering into content-based filtering significantly enhances recommendation precision. By prioritizing songs with shared characteristics within the same cluster, the system tailors suggestions more effectively. The dual approach, considering recommendations with and without clustering, highlights the system's versatility, accommodating diverse user preferences. Overall, the experiment's analysis underscores the efficacy of the hybrid methodology in optimizing the accuracy, relevance, and personalization of music recommendations, showcasing promising outcomes for the evolving landscape of recommendation systems.

## VI. FUTURE WORKS

The success of our hybrid approach in music recommendation through clustering and content-based filtering opens avenues for future research and enhancements. Several directions for future work include:

- 1) **Incorporating User Feedback:** Implementing mechanisms to gather and incorporate user feedback will

enhance the adaptability of the recommendation system. Iteratively refining recommendations based on user interactions ensures continual improvement.

- 2) **Advanced Clustering Techniques:** Exploring and integrating advanced clustering techniques, such as hierarchical or density-based clustering, could further refine the grouping of music tracks. This may provide a more nuanced understanding of the inherent structure in the dataset.
- 3) **Additional Feature Integration:** Including additional features, such as genre information, acoustic attributes, or user demographic data, can contribute to a more comprehensive profiling of music tracks and users. This expanded feature set could lead to more sophisticated and diverse recommendations.
- 4) **Collaborative Filtering Integration:** Combining collaborative filtering methods with the existing approach would capture user preferences more comprehensively by considering the interactions and preferences of similar users.
- 5) **Scalability Considerations:** Addressing scalability concerns for larger datasets and optimizing the recommendation system for real-time deployment are essential for practical applicability in platforms with extensive user bases.
- 6) **Temporal Dynamics Refinement:** Further investigating and refining the representation of temporal dynamics, potentially through more sophisticated modeling of time-dependent features, could lead to improved accuracy in predicting user preferences over time.
- 7) **Evaluation Metrics Enhancement:** Expanding and refining the set of evaluation metrics to include a broader spectrum of user-centric measures, such as user satisfaction and engagement, would provide a more holistic assessment of the recommendation system's effectiveness.

Continuing to innovate and iterate upon this hybrid approach will contribute to the ongoing evolution of music recommendation systems, ensuring they remain adaptive, accurate, and aligned with the dynamic nature of user preferences.

## CONCLUSION

In this experiment, the integration of K-Means clustering and content-based filtering for music recommendation has yielded promising outcomes, marking a significant advancement in recommendation system precision. The clustering analysis successfully revealed inherent patterns within the dataset, providing valuable insights into user preferences and temporal trends. The determination of an optimal cluster count through the elbow method ensures meaningful and interpretable groupings. The subsequent enhancement of content-based filtering with a similarity matrix derived from clustering markedly improves the system's ability to offer personalized and tailored music suggestions. The dual approach of recommending songs with and without clustering underscores the system's adaptability to diverse user preferences.



This hybrid methodology showcases a user-centric approach, prioritizing both content-based similarities and inherent dataset patterns. The experiment's success holds implications for the refinement and evolution of recommendation systems in the dynamic landscape of music streaming platforms. As user expectations evolve, the versatility demonstrated in this experiment is crucial for ensuring the continued optimization of music recommendations. Moving forward, this hybrid model provides a robust framework for enhancing user satisfaction and engagement, paving the way for more sophisticated and effective recommendation systems in the realm of music streaming.

## REFERENCES

- [1] Shreya Agrawal and Pooja Jain. An improved approach for movie recommendation system. In *2017 International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC)*, pages 336–342, 2017.
- [2] Tessy Badriyah, Sefryan Azvy, Wiratmoko Yuwono, and Iwan Syarif. Recommendation system for property search using content based filtering method. In *2018 International conference on information and communications technology (ICOIACT)*, pages 25–29. IEEE, 2018.
- [3] Paulo Chiliguano and Gyorgy Fazekas. Hybrid music recommender using content-based and social information. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2618–2622. IEEE, 2016.
- [4] Mohamed Chiny, Marouane Chihab, Omar Bencharef, and Younes Chihab. Netflix recommendation system based on tf-idf and cosine similarity algorithms. *no. Bml*, pages 15–20, 2022.
- [5] Umair Javed, Kamran Shaukat, Ibrahim A Hameed, Farhat Iqbal, Talha Mahboob Alam, and Suhuai Luo. A review of content-based and context-based recommendation systems. *International Journal of Emerging Technologies in Learning (iJET)*, 16(3):274–306, 2021.
- [6] Muhammad Johari, Arif Laksito, et al. The hybrid recommender system of the indonesian online market products using imdb weight rating and tf-idf. *Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi)*, 5(5):977–983, 2021.
- [7] Michael J Pazzani and Daniel Billsus. Content-based recommendation systems. In *The adaptive web: methods and strategies of web personalization*, pages 325–341. Springer, 2007.
- [8] Chhavi Rana and Sanjay Kumar Jain. Building a book recommender system using time based content filtering. *WSEAS Transactions on Computers*, 11(2):27–33, 2012.
- [9] SRS Reddy, Sravani Nalluri, Subramanyam Kunisetti, S Ashok, and B Venkatesh. Content-based movie recommendation system using genre correlation. In *Smart Intelligent Computing and Applications: Proceedings of the Second International Conference on SCI 2018, Volume 2*, pages 391–397. Springer, 2019.
- [10] Lalita Sharma and Anju Gera. A survey of recommendation system: Research challenges. *International Journal of Engineering Trends and Technology (IJETT)*, 4(5):1989–1992, 2013.
- [11] Ramni Harbir Singh, Sargam Maurya, Tanisha Tripathi, Tushar Narula, and Gaurav Srivastav. Movie recommendation system using cosine similarity and knn. *International Journal of Engineering and Advanced Technology*, 9(5):556–559, 2020.
- [12] Poonam B Thorat, Rajeshwari M Goudar, and Sunita Barve. Survey on collaborative filtering, content-based filtering and hybrid recommendation system. *International Journal of Computer Applications*, 110(4):31–36, 2015.
- [13] Robin Van Meteren and Maarten Van Someren. Using content-based filtering for recommendation. In *Proceedings of the machine learning in the new information age: MLnet/ECML2000 workshop*, volume 30, pages 47–56. Barcelona, 2000.
- [14] Gisela Yunanda, Dade Nurjanah, and Selly Meliana. Recommendation system from microsoft news data using tf-idf and cosine similarity methods. *Building of Informatics, Technology and Science (BITS)*, 4(1):277–284, 2022.