

# A Systematic Review on Auditory Parameters Analysis and Recommendation for Songs/Podcasts

Prof.Srushti Gunthe, Prasad Nalawade, Suchita Chavan, Rutuja Jadhav, Prajwal Korade

*Department of Information Technology, SITS,Pune*

Date of Submission: 20-02-2024

Date of Acceptance: 04-03-2024

## ABSTRACT-

In recent times almost, everything around us is being automated in some way just to make life easier for us. With the use of machine learning and deep learning we can help, respond and ease the tasks of human. The re-search generally focuses on recognition of audio parameters and recommending the best audio signals for songs and podcasts. The re-search genarally focuses on recognition of audio parameters and recommending the best audio signals for songs and podcasts. The voice recognition component utilizes MFCC to extract distinctive features from vocal inputs, enabling precise identification and differentiation. KBFE, through K-means clustering, refines the feature set, optimizing the model's ability of noticing small differences in what users like. Our recommendation system, focused on smart learning like RNN, learns from how each user acts and changes over time. KNN helps by finding users who are alike, making the system better at guessing what users might like. GMM and HMM are like the brain of the system, figuring out patterns in lots of audio data. They help smoothly put together the recommendation system. These fancy algorithms make sure the system not only hears voices well but also gives suggestions that really match what users like. This project is a mix of high-tech learning methods, making it easy for users to enjoy audio content that fits them. By combining MFCC, KBFE (knowledge based front end), RNN, KNN, GMM, HMM, and more, our project will recognize the audio parameters and will recommend which audio signals are suitable for songs and podcasts.).

Keywords-MFCC (Mel-Frequency Cepstral Coefficient), KNN, GMM, HMM, KBFE(Knowledge based front end)

## I. INTRODUCTION

Emotion recognition in speech, speech separation, music recommendations, and machine learning algorithm comparisons are pivotal research areas with widespread applications in enhancing human-computer interactions, voice processing systems, and recommendation systems [9]. Understanding and processing emotions in speech is essential for improving the responsiveness of computers to human feelings, thereby benefitting customer service, health-care, and more[1]. The use of Mel Frequency Cepstral Coefficients (MFCCs) and deep learning models like Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks have contributed to more accurate emotion recognition [8][3]. To further advance these systems, access to diverse databases of real-life emotional speech is crucial, and researchers are exploring the amalgamation of various emotion classifiers for superior results[2]. The future holds the promise of real-time applications, such as empathetic chatbots that can sense and respond to human emotions naturally.

Speech separation through Sparse Non-Negative Matrix Factorization (SNMF) addresses the challenge of isolating multiple voices in a single recording, with potential implications for voice recognition systems [10]. Researchers are utilizing techniques that dissect audio into individual components, complemented by deep learning methods like CNNs and Recurrent Neural Networks (RNNs) to enhance the clarity and naturalness of separated voices. These advancements are augmenting voice processing systems across various domains [6]. In the realm of music recommendations, sophisticated

methodologies are employed to tailor song suggestions to individual preferences. Music properties are analyzed using techniques like Shortest Time Fourier Transform (STFT), and dynamic K-means clustering is utilized to discern users' music tastes [3]. With the increasing availability of music online, real-time recommendations have gained significance, particularly in the context of music streaming. The future envisions systems that can amalgamate multiple recommendation approaches to provide users with the ultimate musical experience [4].

In the context of machine learning algorithms, this survey delves into a comparative analysis of K-Nearest Neighbor (KNN), Genetic Algorithm (GA), Support Vector Machine (SVM), Decision Tree (DT), and Long Short Term Memory (LSTM) [12]. These algorithms function as diverse tools applicable across various domains, including fraud detection and image recognition. The objective is to ascertain the most suitable algorithm for specific tasks, often relying on performance metrics [11]. Findings often highlight the superior performance of algorithms such as LSTM and SVM. The future is expected to bring about smarter systems with the integration of AI into various tasks, optimizing their functionality and efficiency [7].

This survey paper comprehensively reviews these research areas, offering a detailed examination of the methodologies, techniques, and potential future developments in emotion recognition, speech separation, music recommendations, and machine learning algorithm selection. Each section will provide an in-depth discussion, examining the significance of these domains and the research progress made within them.

The rest of the paper is organized as follows: Section 2 reviews literature survey of different research papers, Section 3 provides an overview of methodologies, Section 4 discusses challenges related to technologies, Section 4 provides an overview of methodologies, Section 5 reviews research directions, Section 6 concludes the paper.

## II. LITERATURE REVIEW

Recognizing emotions in speech is crucial for applications like making computers more responsive to our feelings, improving customer service, and health-care. Researchers often use a technique called Mel Frequency Cepstral Coefficients (MFCCs) to analyze speech features and deep learning methods, like Convolutional Neural Networks (CNNs) and Long Short-Term

Memory (LSTM) networks, to make emotion recognition more accurate. To build better systems, they also need larger and more diverse databases of real-life emotional speech. A clever approach is to combine different emotion classifiers for better results. The future looks promising for real-time applications like chatbots that can sense emotions, making human-computer interactions more natural and empathetic. [2] [3].

Time-Frequency Matrix (TFM) presents a novel approach for audio feature extraction and classification, with a particular focus on environmental audio signals. These signals are challenging due to their non-stationarity and discontinuities. The primary objective of the research is to develop a feature extraction technique that effectively quantifies the non-stationarities in audio signals and improves the accuracy of audio classification. [4]. For speech separation using Sparse Non-Negative Matrix Factorization (SNMF), it's about separating multiple voices when there's only one recording. This can help improve things like voice recognition systems. The researchers had used techniques that break down audio into its individual [10]. They had also added deep learning methods like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) to make the separated voices sound clearer and more natural. These advances are making voice processing systems better, which is helpful in many areas [6]. An effective pitch estimation method using Canonical Correlation Analysis (CCA) for speech processing. The proposed method combines various techniques, including Dominant Harmonic Model (DHM), Empirical Mode Decomposition (EMD), and Normalized Autocorrelation Function (NACF) for accurate pitch estimation. It utilizes CCA to select Intrinsic Mode Function (IMF) components from reference signals, which are then used to reconstruct the signal, resulting in improved pitch estimation accuracy [1].

In the world of music recommendations, researchers are using fancy methods to suggest songs you might like. They analyze music properties using things like Shortest Time Fourier Transform (STFT) and figure out your music tastes with dynamic K-means clustering. It's all about making music recommendations more tailored to you. As more and more music is available online, real-time recommendations are becoming important, especially when you're streaming music. The future might hold systems that can mix different recommendation methods to give you the best musical experience [9]. In the paper comparing machine learning algorithms like K-Nearest Neighbor

bor (KNN), Genetic Algorithm (GA), Support Vector Machine (SVM), Decision Tree (DT), and Long Short Term Memory (LSTM), it's about finding out which algorithm is the smartest for different jobs. These algorithms are like different tools, and they're used in various areas, from catching fraud to recognizing images [8]. Researchers want to know which one works best, and they used measurements to figure it out. They find that algorithms like LSTM and SVM often give the best results. The future is all about making things smarter and using AI to do a lot of different tasks [12].

The middle-high frequency range, particularly between 3,000 Hz and 4,000 Hz, contains key components that contribute to the glossiness of a professional singer's voice. Through neural networks, the system can analyze the difference between an amateur singer's voice and a professional opera singer's voice. Parallel data is created by aligning the amateur and professional singing voices using phoneme information from lyrics [7]. When it comes to recognizing emotions in speech, the review focuses on identifying different feelings people express when they speak. It's essential for making devices like virtual assistants more aware of how you feel. Researchers use features from speech, like Mel Frequency Cepstral Coefficients (MFCCs), and machines like Gaussian Mixture Models (GMM) and k-Nearest Neighbour (k-NN) to figure out these emotions [5]. They use fancy measurements like Gross Pitch Error to check how well the system does. They're always looking for more real-life emotional data to make their systems work better. The future could bring real-time apps that sense your emotions and make your interactions with computers and robots more human-like and caring [11].

The survey paper presents various methodologies and approaches employed in the reviewed literature, focusing on four key areas: Emotion Recognition in Speech, Speech Separation, Music Recommendations, and Machine Learning Algorithm Comparison. Each section highlights the methods used to collect and analyze data in these domains.

### III. METHODOLOGY

1. Mel-Frequency Cepstral Coefficients (MFCCs) Mel-Frequency Cepstral Coefficients (MFCCs) are a critical feature extraction technique used in voice recognition and automatic speech recognition (ASR) systems. MFCCs are particularly effective in representing the

spectral characteristics of audio signals, making them a valuable tool in preprocessing and analyzing spoken language.

2. Recurrent Neural Networks (RNNs) Recurrent Neural Networks (RNNs) are a class of deep learning models commonly used in voice recognition systems, particularly in the automatic speech recognition (ASR) component. RNNs are well-suited for tasks that involve sequential data, making them a natural choice for processing audio signals, which are essentially sequential in nature. Here's how RNNs can be employed in a voice recognition system.
3. K-Nearest Neighbors (KNN) K-Nearest Neighbors (KNN) is not typically used as a primary algorithm for voice recognition systems, especially in automatic speech recognition (ASR) tasks. KNN is more commonly applied in classification and clustering tasks. However, it can still play a role in certain aspects of voice recognition systems.
4. Hidden Markov Models (HMMs) Hidden Markov Models (HMMs) have historically played a crucial role in voice recognition systems, particularly in the domain of automatic speech recognition (ASR). HMMs are a statistical modeling technique that can be used to model both the acoustic and linguistic aspects of speech.

### IV. CHALLENGES

The challenges in auditory parameters analysis and recommendation for songs and podcasts are multi-faceted. One significant challenge lies in personalization, as creating recommendation systems that truly understand individual tastes and moods remains complex. Maintaining user privacy and data ethics is another hurdle, ensuring that recommendations are accurate without compromising personal information. Content licensing and copyright issues add a layer of complexity, as respecting intellectual property rights while recommending content is essential. Scalability is also a challenge, especially for streaming platforms with vast libraries and millions of users. Integrating auditory analysis with emerging technologies like virtual and augmented reality requires addressing compatibility and user experience concerns. Lastly, providing adaptive recommendations that evolve with changing user preferences in real-time remains an ongoing challenge.

Voice recognition and recommendation systems encounter their own set of challenges. These include improving accuracy in recognizing

emotions in speech and integrating real-life emotional data into systems. Striking the right balance between providing personalized recommendations and safeguarding user data and privacy is a constant challenge. Overcoming copyright and licensing constraints in audio content poses a challenge, particularly when recommending songs and podcasts. The scalability of such systems to accommodate a large user base and extensive content libraries is also a significant obstacle. Ensuring compatibility and seamless integration with various devices and platforms adds to the complexity. Lastly, keeping recommendations adaptive and responsive to evolving user preferences presents an ongoing challenge, demanding continuous research and innovation in the field.

## V. RESEARCH DIRECTIONS

In the field of auditory parameter analysis and recommendation for songs and podcasts, several promising research directions are emerging.

1. **Enhanced Personalization:** Investigating more advanced methods to tailor recommendations to individual preferences, considering factors beyond music genre and exploring mood-based or context-aware recommendations.
2. **Data Privacy and Ethics:** Research into safeguarding user data and privacy while developing recommendation systems, addressing ethical concerns, and ensuring responsible data usage.
3. **Content Licensing and Copyright:** Exploring solutions for content recommendation systems to respect copyright and intellectual property rights while providing valuable recommendations.
4. **Scalability:** Developing methods to handle the vast and ever-growing volume of audio content available online and making recommendations scalable for millions of users.
5. **Integration with Emerging Technologies:** Exploring how auditory analysis and recommendations can be integrated into emerging technologies like virtual and augmented reality for more immersive audio experiences.
6. **Adaptive Recommendations:** Researching ways to adapt recommendations in real-time as user preferences evolve, offering more dynamic and up-to-date suggestions.

## VI. CONCLUSION

The research we've explored shows great progress in how computers understand and interact with our voices. They're becoming better at recognizing our emotions, which can make computer interfaces more responsive and improve areas like customer service and healthcare. But to make these systems even better, we need bigger and more diverse databases of real-life emotional speech. We're also using advanced techniques like deep learning to make these systems more accurate and empathetic. In the field of music recommendations, we're personalizing music suggestions using techniques that analyze your music taste. This makes your music experience more tailored, especially when you're using online streaming services. Moreover, we're working on finding the best algorithms for various tasks. Some, like Long Short-Term Memory (LSTM) and Support Vector Machine (SVM), are proving to be very effective, promising a future where intelligent algorithms make various processes smarter and more efficient. All in all, these studies are leading us toward a future where technology understands our emotions and preferences, making our interactions with computers more natural and caring.

## REFERENCES

- [1]. Subrata Kumer Paul, Rakhi Rani Paule "Effective pitch estimation using Canonical Correlation Analysis" In 2020 Dept. of Computer Science and Engineering University of Rajshahi, Rajshahi, Bangladesh, doi: 110.1109/ICAICT51780.2020.9333460. "Effective pitch estimation using Canonical Correlation Analysis"
- [2]. Zhandos Yessenbayev "Robust Segmentation of Speech Signal Using MFCC and Acoustic Parameters" In 2012 Department of Information Technology, L.N. Gumilev Eurasian National University, Astana, Kazakhstan DOI: 10.1109/AMS.2012.26. "Robust Segmentation of Speech Signal Using MFCC and Acoustic Parameters"
- [3]. Om Deshmukh, Carol Y. Espy-Wilson and Amit Juneja "ACOUSTIC PHONETIC SPEECH PARAMETERS FOR SPEAKER-INDEPENDENT SPEECH RECOGNITION" In 2002 University of Maryland, Department of Electrical and Computer Engineering, A. V. Williams Bldg., College Park, MD 20752 DOI: 10.1109/ICASSP.2002.5743787.

- ”ACOUS- TIC -PHONETIC SPEECH PARAMETERS FOR SPEAKER-INDEPENDENT SPEECH RECOGNITION ”
- [4]. Behnaz Ghoraani,Sridhar Krishnan “Time–Frequency Matrix Feature Extrac-tion and Classification of Environmental Audio Signals” In 2011 Dept. of Elec-trical Computer and Biomedical Engi- neering, Toronto Metropolitan University, Toronto DOI:10.1109/TASL.2011.2118753. ” Time–requency Matrix Feature Extraction and Classification of Environmental Audio Signals”
- [5]. Hyan-Soo BaeHo-Jin Lee; Suk-Gyu Lee “Voice recognition based on adaptive MFCC and deep learning” In 2017 IEEE 11th, fer-ence on Industrial Electronics and Applications (ICIEA) DOI: 10.1109/ICIEA.2016.7603830. ” Voice Recognition Based on Adaptive MFCC and Deep Learning ”
- [6]. Ashesh Jain<sup>1,2</sup>, Amir R. Zamir<sup>2</sup>, Silvio Savarese<sup>2</sup>, and Ashutosh Saxena<sup>3</sup> “Structural-RNN: Deep Learning on Spatio-Temporal Graphs” In 2016 Cornell University<sup>1</sup>, Stanford University<sup>2</sup>, Brain Of Things Inc.<sup>3</sup> DOI: 10.1109/CVPR.2016.573 . “ Structural-RNN: Deep Learning on Spatio-Temporal Graphs”
- [7]. Ryuka Nanzaka, Tsuyoshi Kitamura, Yuji Adachi, Kiyoto Tai, Tetsuya Takiguchi “Spec-trum Enhancement of Singing Voice Using Deep Learning” In 2018 IEEE International Symposium on Multimedia (ISM) DOI: 10.1109/ISM.2018.00-18. ”Spectrum Enhance-ment of Singing Voice Using Deep Learning”
- [8]. Shiqing Zhang, Xiaoming Zhao and Qi Tian “Spontaneous Speech Emotion Recognition Using Multiscale Deep Convolutional LSTM” In 2019 nstitute of Intelligent Information Processing Taizhou University Taizhou, China. DOI: 10.1109/TAFFC.2019.2947464. “Spon-taneous Speech Emotion Recognition Using Multiscale Deep Convolutional LSTM”
- [9]. Dong-Moon Kim, Kun-su Kim, Kyo-Hyun Park , Jee-Hyong Lee and Keon Myung Lee “A Music Recommendation System with a Dynamic K-means Clustering Algorithm” In 2007 Department of Electrical and Electronic Engineering, SungkyunKwan University, South Korea. DOI: 10.1109/ICMLA.2007.97. ”A Mu-sic Recommendation System with a Dynamic K-means Clustering Algorithm ”
- [10]. Mikkel N. Schmidt and Rasmus K. Ols-son “Single-Channel Speech Separation us-ing Sparse Non-Negative Matrix Factoriza-tion ” In 2006 Informatics and Mathematical Modelling, Technical University of Denmark, doi: 10.21437/Interspeech.2006-655. ” Single-Channel Speech Separation using Sparse Non-Negative Matrix Factorization”
- [11]. Rahul B. Lanjewar, Swarup Mathurkar, Nilesh Patel “Implementation and Comparison of Speech Emotion Recognition System using Gaussian Mixture Model(GMM) and K-Nearest Neighbor(KNN) techniques.” ”Implementation and Comparison of Speech Emotion Recognition System using Gaussian Mixture Model(GMM) and K-Nearest Neighbor(KNN) techniques. ”
- [12]. Malti Bansal, Apoorva Goyal, Apoorva Choud-hari “A comparative Analysis of K-Nearest Neighbor, Genetic, Support Vector Machine, Decision Tree, and Long Short Term Memory algorithms in machine learning” ”A compara-tive Analysis of K-Nearest Neighbor, Genetic, Support Vector Machine, Decision Tree, and Long Short Term Memory algorithms in ma-chine learning. ”
- [13].