

Project

Summary for the Papers:

Paper 1: Spotify Data Analysis and Song Popularity Prediction

Authors: Sivasai Bhavanasi, Sahil Malla, V Manichetan , CVNJ Dhanush, Dr B Prakash

DOI: 10.35629/5252-0505296304

Summary: The paper "**Spotify Data Analysis and Song Popularity Prediction**" explores how machine learning can be leveraged to analyze the characteristics of songs and predict their success on Spotify. The authors collected a dataset of 30,000 songs from Spotify, which consisted of metadata such as track ID, artist, album, genre, and various audio features like danceability, energy, loudness, valence, and tempo. Data preprocessing involved handling missing values using forward and backward fill techniques and eliminating duplicates to ensure data integrity. Authors used Lasso Regression, they identified that acousticness, valence, and loudness positively impacted a song's popularity, whereas energy, instrumentalness, and liveness negatively correlate with the success of the song. Features such as key, mode, speechiness, valence, and tempo had little to no impact on determining popularity. To predict song popularity, various machine learning models were tested, including Random Forest, Gradient Boosting, Bagging Classifier, XGBoost, and Decision Tree Classifier, with the Decision Tree model emerging as the most accurate (93% accuracy). The findings suggest that machine learning can play a pivotal role in optimizing music recommendation algorithms. The authors suggest that future research could explore deep learning models to refine these predictions further, as the evolving music industry increasingly relies on data analytics to anticipate trends and audience preferences. The authors suggest that future work could explore deep learning models to refine these predictions further, as the evolving music industry increasingly relies on data analytics to anticipate trends and audience preferences.

Paper 2: Dance Hit Song Prediction

Authors: Dorien Herremans, David Martens, Kenneth Sörensen

DOI: 10.1080/09298215.2014.881888

Summary: The paper "**Dance Hit Song Prediction**" explores how machine learning can be employed to predict whether a dance song will become a Top 10 hit or remain lower ranked. This research is crucial for the music industry. The study builds on the emerging field of Hit Songs and aims to address prior inconsistencies in predicting song popularity by focusing particularly on dance music. To achieve this, the researchers compiled a dataset of dance hits spanning 1985 to 2013 from Billboard and the Official Charts Company, gathering detailed musical and temporal features. The features were analyzed into three categories: meta-information, basic audio features, and temporal features. The authors trained and tested multiple machine learning models using three datasets,

each differentiating between Top 10 hits and lower-ranked songs to minimize genre bias. Many classifiers were applied, like Decision Trees, RIPPER rule sets, Naive Bayes, Logistic Regression, and Support Vector Machines (SVMs). The best-performing model was Logistic Regression, achieving an AUC (area under the curve) of 0.67, indicating that certain audio characteristics strongly correlate with chart success. Another test was done on an out-of-time dataset which showed that models performed better on recent songs, implying that hit prediction might also depend on contemporary music trends. The study concludes that machine learning can successfully predict dance music hits based purely on audio features, and the model has even been implemented into an online application that allows users to test their songs' hit potential. Future research may integrate lyrics, social media engagement, and streaming data to enhance prediction accuracy and expand the model to other musical genres. The authors suggest that future work could integrate lyrics, social media engagement, and streaming data to enhance prediction accuracy and expand the model to other musical genres.

Paper 3: Interactive Music Learning Model Based on the RBF Algorithm

Authors: Fengqin Liu

DOI: 0.1155/2022/5759986

Summary: The paper, **Interactive Music Learning Model Based on the RBF Algorithm** explores how AI can evolve the music industry by making it more interactive and personalized for students. With the emergence of digital music tools and AI-driven applications, traditional teaching methods are evolving to include intelligent systems that enhance the learning process. The study uses an interactive teaching model using the Radial Basis Function algorithm, which is a type of neural network that processes music learning patterns and student performance feedback to create an adaptive and personalized learning experience. The RBF algorithm used in this study plays a huge role in analyzing music performance evaluations, identifying key learning patterns, and dynamically adjusting teaching methods to suit each student. The system consists of five layers, including input processing, normalization, and an output layer that evaluates student performance based on their music skill assessments. The model mimics the way a real teacher would guide students, adjusting lessons in real time to match their progress and learning needs. But beyond technical skills, the study highlights something even more important which is the emotional side of music education. Learning music isn't just about hitting the right notes but it's about feeling and expressing emotions through sound, and this AI-driven system recognizes that. By integrating multimedia tools, AI-powered emotional analysis, and interactive teaching techniques, the model creates a rich, immersive learning experience that keeps students engaged and motivated. Whether through group activities, hands-on classroom sessions, or instant feedback, this system helps both students and teachers connect more deeply with music. The authors suggest that future advancements in AI and deep learning could further refine the model, making it more responsive to individual learning styles and expanding its application beyond music education.

Project Report:

Aim

The goal of this project was to analyze a music dataset containing information such as tempo, loudness, energy, danceability, and speechiness to predict song popularity and classify songs into different genres or categories. Understanding music trends is essential for improving recommendation systems, aiding music producers, and enhancing user experience on streaming platforms.

Importance of the Chosen Data

The dataset we used in this project contains music information from different decades, from the 1960s to the 2010s. This data is important because it helps us understand how music trends have changed over time. By analyzing the features of songs, we can group similar songs, predict song popularity, and classify music styles. This information can be useful for music streaming platforms, researchers, and even musicians who want to understand what makes a song popular.

Preprocessing Details

The preprocessing involved several key steps to ensure data quality:

- **Data Integration:** Multiple CSV files representing different decades were merged into a single dataset.
- **Duplicate Removal:** Duplicate entries were identified and removed to maintain unique records.
- **Missing Values Handling:** A heatmap visualization confirmed that there were no missing values in the dataset.
- **Feature Selection:** Numeric and non-numeric features were identified for further analysis.
- **Normalization:** Some features were normalized to ensure consistency in machine learning models.

Feature Selection and Cleaning

Important features considered in the analysis included attributes such as tempo, loudness, energy, danceability, and speechiness. These features were chosen based on their impact on clustering and classification tasks. Unnecessary columns were dropped, and data distributions were examined to detect any anomalies.

Prediction Task and Best R^2 Score

The primary prediction focused on danceability using various regression models. The best-performing model was **Random Forest**, achieving an R^2 score of **0.481**, indicating its moderate effectiveness in capturing the variance in danceability. The model captured non-linear relationships better than linear regression. However, some features with moderate correlation, such as loudness and instrumentality, did not perform as expected in the prediction.

Clustering Methods and Evaluation

- K-Means Clustering: Provided well-separated clusters in the "energy-loudness" feature space.
- Mini-Batch K-Means: Showed similar results but introduced minor variations.
- Mean Shift C Agglomerative Clustering: Evaluated but produced less distinct clusters.

Best Clustering Method:

Among the clustering methods applied, **Agglomerative clustering** provided the most well-defined and meaningful clusters. The visualization shows distinct groupings with smooth transitions, indicating that the **energy, loudness, and danceability** features naturally form relationships. Mini-Batch K-Means followed closely, offering a better cluster quality, though some boundary noise was visible. Mean-Shift clustering failed to detect meaningful clusters.

Classification and Evaluation

Models Implemented:

- Logistic regression
- Decision Tree Classifier
- Support Vector Machines (SVM)
- Neural Network (MLP classifier)

Model Evaluation:

Among the tested classifiers, **the Neural Network (MLP Classifier) provided the best performance**, achieving **72.56% accuracy**, closely followed by **SVM at 71.65%**. Logistic Regression and Decision Tree models performed worse, struggling to distinguish between Medium and other categories. Across all models, **the Medium category was never correctly predicted**, suggesting a **class imbalance issue**.

Challenges Faced

- Handling missing values: Deciding between mean and median replacements was context dependent.
- Data inconsistency: Some features required additional preprocessing before analysis.
- Computational limitations: Processing large datasets efficiently.
- Feature selection complexity: Identifying the most impactful features.
- Parameter tuning: Optimizing clustering and classification algorithms.
- Model interpretability: Random Forest, despite its accuracy, was challenging to interpret at an individual prediction level.

Conclusion and Takeaways

This project successfully implemented preprocessing, clustering, and classification techniques on a large-scale music dataset. The findings demonstrated that **Agglomerative clustering** was the most effective for grouping songs based on their characteristics, and **Neural Network (MLP Classifier)** was the best model for classification tasks. Correlation analysis helped identify key features, but regression modeling validated their actual importance.

Project

**Audio-Based Songs for Personalized Playlist
Recommendations On Spotify**

SAMPLE DATASET

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S |
|---------------------------------|-------------------|---------------------------------------|--------------|--------|-----|----------|------|-------------|--------------|------------------|----------|---------|---------|-------------|----------------|------------|----------|--------|
| track | artist | uri | danceability | energy | key | loudness | mode | speechiness | acousticness | instrumentalness | liveness | valence | tempo | duration_ms | time_signature | chorus_hit | sections | target |
| Lucky Man | Montgomery Gentry | spotify:track:4GiXBCUF7H6YfNQsnBRlzl | 0.578 | 0.471 | 4 | -7.27 | 1 | 0.0289 | 0.368 | 0 | 0.159 | 0.532 | 133.061 | 196707 | 4 | 30.88059 | 13 | 1 |
| On The Hotline | Pretty Ricky | spotify:track:1zyqZONW985C4osz9wlsu | 0.704 | 0.854 | 10 | -5.477 | 0 | 0.183 | 0.0185 | 0 | 0.148 | 0.688 | 92.988 | 242587 | 4 | 41.51106 | 10 | 1 |
| Clouds Of Dementia | Candlemass | spotify:track:6cH2f7RbxXCKwEkgAZT4mY | 0.162 | 0.836 | 9 | -3.009 | 1 | 0.0473 | 0.000111 | 0.00457 | 0.174 | 0.3 | 86.964 | 338893 | 4 | 65.32887 | 13 | 0 |
| Heavy Metal, Raise Hell! | Zwartketterij | spotify:track:2lJBpP2vMeX7LggzRN3iSX | 0.188 | 0.994 | 4 | -3.745 | 1 | 0.166 | 7.39E-06 | 0.0784 | 0.192 | 0.333 | 148.44 | 255687 | 4 | 58.59528 | 9 | 0 |
| I Got A Feelin' | Billy Currington | spotify:track:1fF370eYXUcWwklvaq3lGz | 0.63 | 0.764 | 2 | -4.353 | 1 | 0.0275 | 0.363 | 0 | 0.125 | 0.631 | 112.098 | 193760 | 4 | 22.62384 | 10 | 1 |
| Dantzig Station | State Of Art | spotify:track:5Z3nrC0JbJmXaOGiXTuNFk | 0.726 | 0.837 | 11 | -7.223 | 0 | 0.0965 | 0.373 | 0.268 | 0.136 | 0.969 | 135.347 | 192720 | 4 | 28.29051 | 10 | 0 |
| Divorced | Blacklisted | spotify:track:0iAdSLiQBliZTAilUP7p5E | 0.365 | 0.922 | 1 | -2.644 | 1 | 0.071 | 0.00285 | 0 | 0.321 | 0.29 | 77.25 | 89427 | 4 | 45.77202 | 4 | 0 |
| Where I Come From | Alan Jackson | spotify:track:6ej1QJ8elYmhsyTlvgDajy | 0.726 | 0.631 | 11 | -8.136 | 0 | 0.0334 | 0.22 | 0 | 0.193 | 0.746 | 124.711 | 239240 | 4 | 35.59732 | 10 | 1 |
| Nothin' To Die For | Tim McGraw | spotify:track:3iRSz6HujrSy9b3LXg2Kg9 | 0.481 | 0.786 | 10 | -5.654 | 1 | 0.0288 | 0.0538 | 0 | 0.0759 | 0.389 | 153.105 | 253640 | 4 | 19.65701 | 11 | 1 |
| I Want to Know Your Plans | Say Anything | spotify:track:3pinCLIHbRczUjenW0Eo56 | 0.647 | 0.324 | 7 | -9.679 | 1 | 0.0377 | 0.354 | 0 | 0.115 | 0.344 | 124.213 | 314286 | 3 | 32.66343 | 16 | 0 |
| F.U.R.B. (F U Right Back) | Frankee | spotify:track:7jElrCgQJBcVLsbMRKni2t | 0.787 | 0.632 | 8 | -3.487 | 1 | 0.137 | 0.103 | 6.78E-06 | 0.388 | 0.612 | 141.026 | 198173 | 4 | 18.091 | 9 | 1 |
| Amarillo Sky | Jason Aldean | spotify:track:0axUHkhMMY0YSC1jFBVwqv | 0.491 | 0.776 | 2 | -3.887 | 1 | 0.0393 | 0.314 | 0 | 0.146 | 0.428 | 154.988 | 202547 | 4 | 37.86861 | 9 | 1 |
| Gin And Juice | Hot Rod Circuit | spotify:track:1xthH0Ze4FYo2y99quJUJj | 0.455 | 0.737 | 1 | -6.206 | 1 | 0.0272 | 0.00114 | 0.000487 | 0.195 | 0.206 | 94.028 | 224053 | 4 | 62.77759 | 9 | 0 |
| Six-Pack Summer | Phil Vassar | spotify:track:23qDMWnwf8p0pr5slIjB6i | 0.725 | 0.733 | 0 | -6.66 | 1 | 0.0242 | 0.508 | 0 | 0.173 | 0.843 | 100.311 | 219827 | 4 | 27.0663 | 8 | 1 |
| Hatho Pai Kariya Na Kar | Kartar Ramla | spotify:track:7KA8tQVcGqCHHdGwSCRi6v | 0.497 | 0.421 | 2 | -14.059 | 0 | 0.187 | 0.985 | 0.839 | 0.148 | 0.804 | 92.19 | 232719 | 4 | 33.3613 | 15 | 0 |
| Serious Hardcore - Original Mix | Ham | spotify:track:6mGnxmegYYJA2TkhXCbOkb | 0.509 | 0.942 | 11 | -6.899 | 1 | 0.0628 | 0.000639 | 0.882 | 0.0788 | 0.163 | 169.96 | 460080 | 4 | 34.78597 | 12 | 0 |
| Daddy Won't Sell The Farm | Montgomery Gentry | spotify:track:2Ww173KM9i97KSRE2nCiini | 0.708 | 0.778 | 7 | -7.039 | 1 | 0.0298 | 0.011 | 2.26E-05 | 0.0556 | 0.704 | 119.881 | 258640 | 4 | 33.81641 | 15 | 1 |

| A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S |
|---------------------------------|-------------------|--------------------------------------|--------------|--------|-----|----------|------|-------------|--------------|------------------|----------|---------|---------|-------------|----------------|------------|----------|--------|
| track | artist | uri | danceability | energy | key | loudness | mode | speechiness | acousticness | instrumentalness | liveness | valence | tempo | duration_ms | time_signature | chorus_hit | sections | target |
| Lucky Man | Montgomery Gentry | spotify:track:4GiXBCUF7H6YfNQsnBRlZl | 0.578 | 0.471 | 4 | -7.27 | 1 | 0.0289 | 0.368 | 0 | 0.159 | 0.532 | 133.061 | 196707 | 4 | 30.88059 | 13 | 1 |
| On The Hotline | Pretty Ricky | spotify:track:1zyqZONW985Cs4osz9wlsu | 0.704 | 0.854 | 10 | -5.477 | 0 | 0.183 | 0.0185 | 0 | 0.148 | 0.688 | 92.988 | 242587 | 4 | 41.51106 | 10 | 1 |
| Clouds Of Dementia | Candlemass | spotify:track:6chZl7RbxXCKwEkgAZT4mY | 0.162 | 0.836 | 9 | -3.009 | 1 | 0.0473 | 0.000111 | 0.00457 | 0.174 | 0.3 | 86.964 | 338893 | 4 | 65.32887 | 13 | 0 |
| Heavy Metal, Raise Hell! | Zwartketterij | spotify:track:2lJBPP2vMeX7LggzRN3iSX | 0.188 | 0.994 | 4 | -3.745 | 1 | 0.166 | 7.39E-06 | 0.0784 | 0.192 | 0.333 | 148.44 | 255687 | 4 | 58.59528 | 9 | 0 |
| I Got A Feelin' | Billy Currington | spotify:track:1fF370eYXUcWwklvaq3lGz | 0.63 | 0.764 | 2 | -4.353 | 1 | 0.0275 | 0.363 | 0 | 0.125 | 0.631 | 112.098 | 193760 | 4 | 22.62384 | 10 | 1 |
| Dantzig Station | State Of Art | spotify:track:5Z3nrC0JbJmXaOGiXTuNfK | 0.726 | 0.837 | 11 | -7.223 | 0 | 0.0965 | 0.373 | 0.268 | 0.136 | 0.969 | 135.347 | 192720 | 4 | 28.29051 | 10 | 0 |
| Divorced | Blacklisted | spotify:track:0iAdSLiQBliZTAilUP7p5E | 0.365 | 0.922 | 1 | -2.644 | 1 | 0.071 | 0.00285 | 0 | 0.321 | 0.29 | 77.25 | 89427 | 4 | 45.77202 | 4 | 0 |
| Where I Come From | Alan Jackson | spotify:track:6ej1QJ8elYmhsyTlvgDaj | 0.726 | 0.631 | 11 | -8.136 | 0 | 0.0334 | 0.22 | 0 | 0.193 | 0.746 | 124.711 | 239240 | 4 | 35.59732 | 10 | 1 |
| Nothin' To Die For | Tim McGraw | spotify:track:3lRSz6HujrSy9b3LXg2Kq9 | 0.481 | 0.786 | 10 | -5.654 | 1 | 0.0288 | 0.0538 | 0 | 0.0759 | 0.389 | 153.105 | 253640 | 4 | 19.65701 | 11 | 1 |
| I Want to Know Your Plans | Say Anything | spotify:track:3pinCLIHbRczUjenW0Eo56 | 0.647 | 0.324 | 7 | -9.679 | 1 | 0.0377 | 0.354 | 0 | 0.115 | 0.344 | 124.213 | 314286 | 3 | 32.66343 | 16 | 0 |
| F.U.R.B. (F U Right Back) | Frankee | spotify:track:7JElrCgQJBcVLsbMRKni2t | 0.787 | 0.632 | 8 | -3.487 | 1 | 0.137 | 0.103 | 6.78E-06 | 0.388 | 0.612 | 141.026 | 198173 | 4 | 18.091 | 9 | 1 |
| Amarillo Sky | Jason Aldean | spotify:track:0axUHkhMMY0YSC1fBVWqv | 0.491 | 0.776 | 2 | -3.887 | 1 | 0.0393 | 0.314 | 0 | 0.146 | 0.428 | 154.988 | 202547 | 4 | 37.86861 | 9 | 1 |
| Gin And Juice | Hot Rod Circuit | spotify:track:1xthH0Ze4FYo2y99QuJUJj | 0.455 | 0.737 | 1 | -6.206 | 1 | 0.0272 | 0.00114 | 0.000487 | 0.195 | 0.206 | 94.028 | 224053 | 4 | 62.77759 | 9 | 0 |
| Six-Pack Summer | Phil Vassar | spotify:track:23qDMWnwf8p0pr5sljB6i | 0.725 | 0.733 | 0 | -6.66 | 1 | 0.0242 | 0.508 | 0 | 0.173 | 0.843 | 100.311 | 219827 | 4 | 27.0663 | 8 | 1 |
| Hatho Pai Kariya Na Kar | Kartar Ramla | spotify:track:7KA8tQVcGqCHdGwSCRi6v | 0.497 | 0.421 | 2 | -14.059 | 0 | 0.187 | 0.985 | 0.839 | 0.148 | 0.804 | 92.19 | 232719 | 4 | 33.3613 | 15 | 0 |
| Serious Hardcore - Original Mix | Ham | spotify:track:6mGnxmegYYJA2TkhXCbOkb | 0.509 | 0.942 | 11 | -6.899 | 1 | 0.0628 | 0.000639 | 0.882 | 0.0788 | 0.163 | 169.96 | 460080 | 4 | 34.78597 | 12 | 0 |
| Daddy Won't Sell The Farm | Montgomery Gentry | spotify:track:2Wt173KM9i97KSR5FnCuni | 0.708 | 0.728 | 7 | -7.039 | 1 | 0.0298 | 0.011 | 2.26E-05 | 0.0556 | 0.704 | 119.881 | 258640 | 4 | 33.81641 | 15 | 1 |

DATA CLEANING

Merging Datasets : Combined multiple CSV files (1960s–2010s) into a single dataset.

Removing Duplicates : Checked and removed duplicate rows to avoid redundancy.

Handling Missing Values : Verified missing values using `df.isnull().sum()`. No missing values were found in the dataset.

Data Type Verification & Conversion : Identified numeric and non-numeric columns. Converted potential numeric columns where applicable.

Checking Data Distribution : Used histograms and boxplots to analyze distributions and detect anomalies.

Feature Selection Based on Correlation : Applied Pearson correlation to determine relationships.

Dropped weakly correlated features: key, mode, liveness, tempo, duration_ms, time_signature, chorus_hit, and sections.

Standardization of Features : Used StandardScaler to normalize loudness for better model performance.

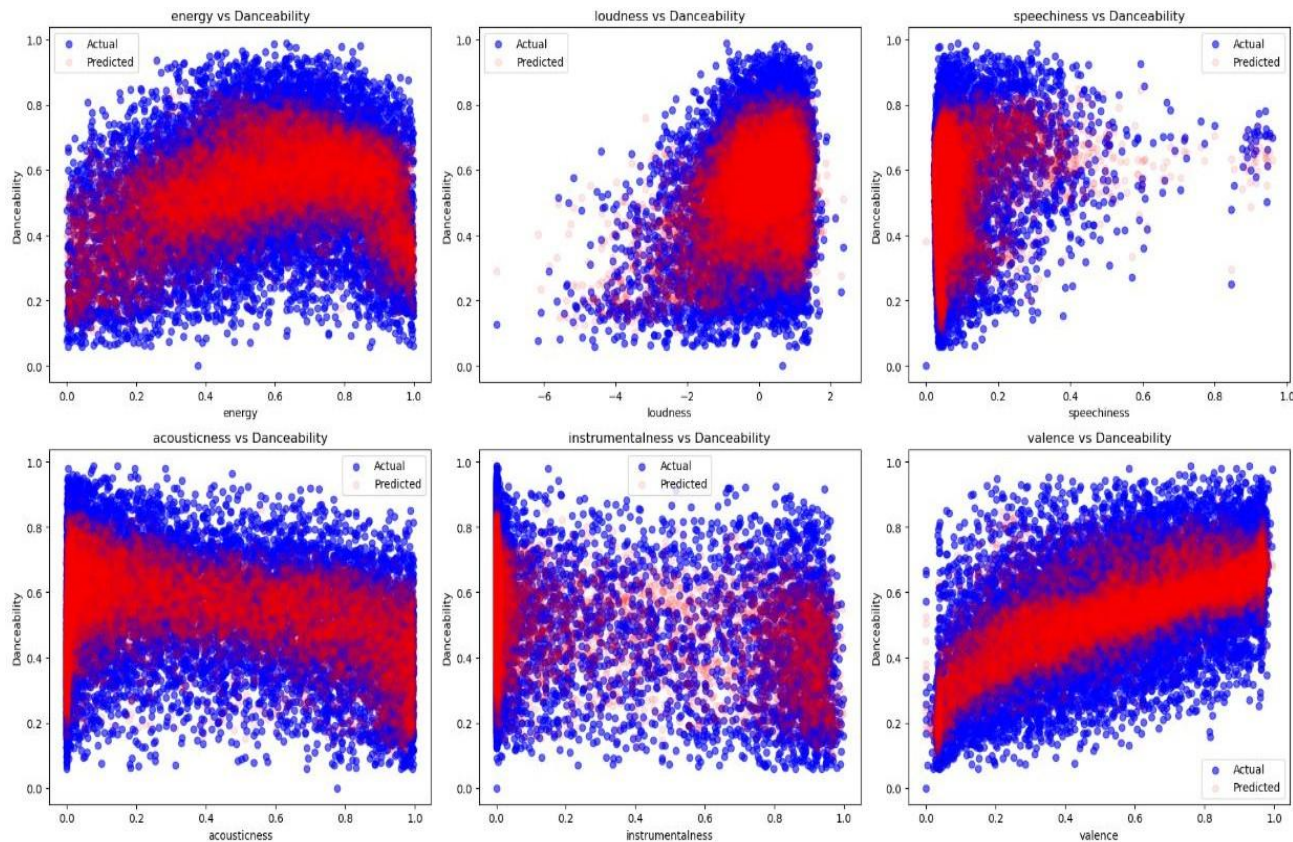


REGRESSION


Goal: Predict danceability of a song based on other audio features.

Feature used: Energy, Loudness, Speechiness, Acousticness, Instrumentalness, Valence.

Models Tested: Linear Regression (Baseline model) , Random Forest Regressor (Best performing model)



CLUSTERING

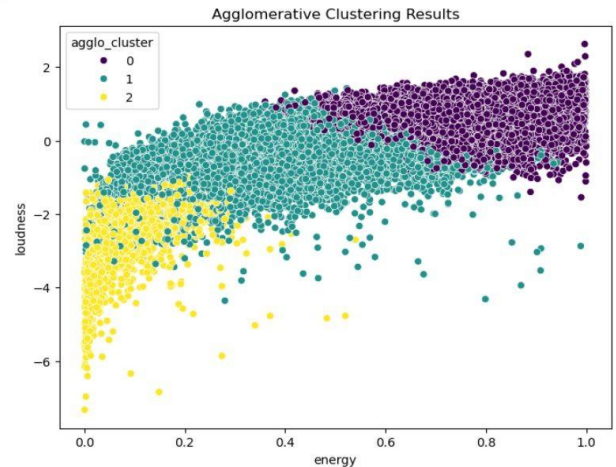
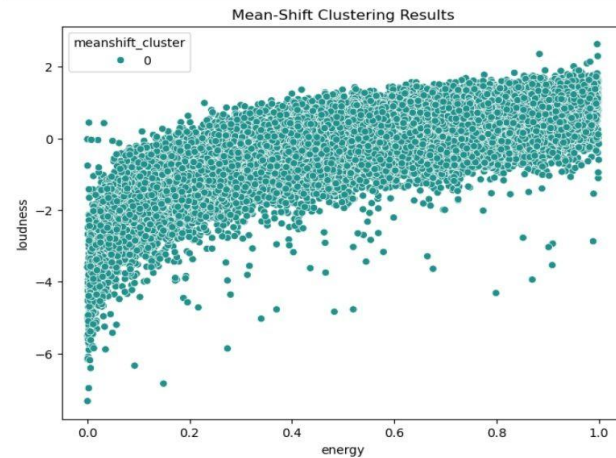
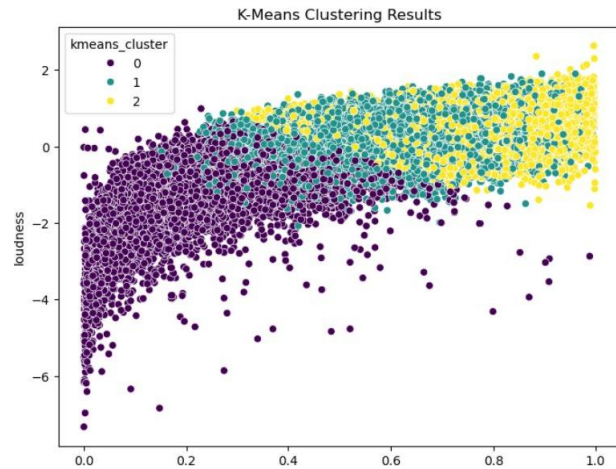


K-Means Clustering: Identified three clusters based on energy, loudness, and danceability. Good separation of groups.

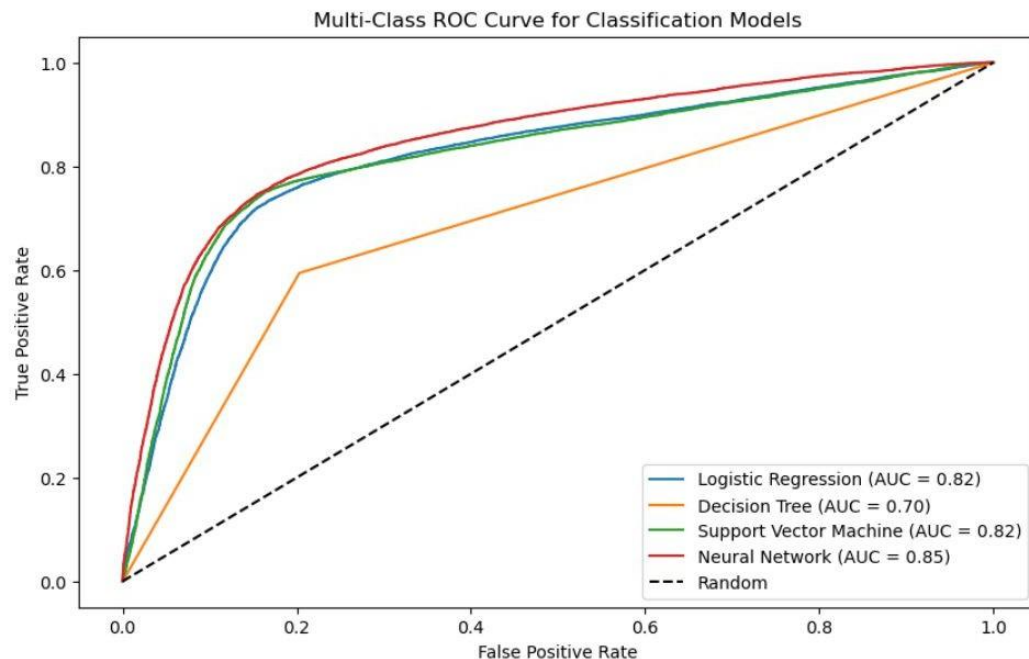
Mini-Batch K-Means: Faster version of K-Means but slightly noisier results.

Mean-Shift Clustering: Identified only one cluster, indicating low density variation.

Agglomerative Clustering: Best results among clustering techniques. Defined groups with smooth transitions.



CLASSIFICATION



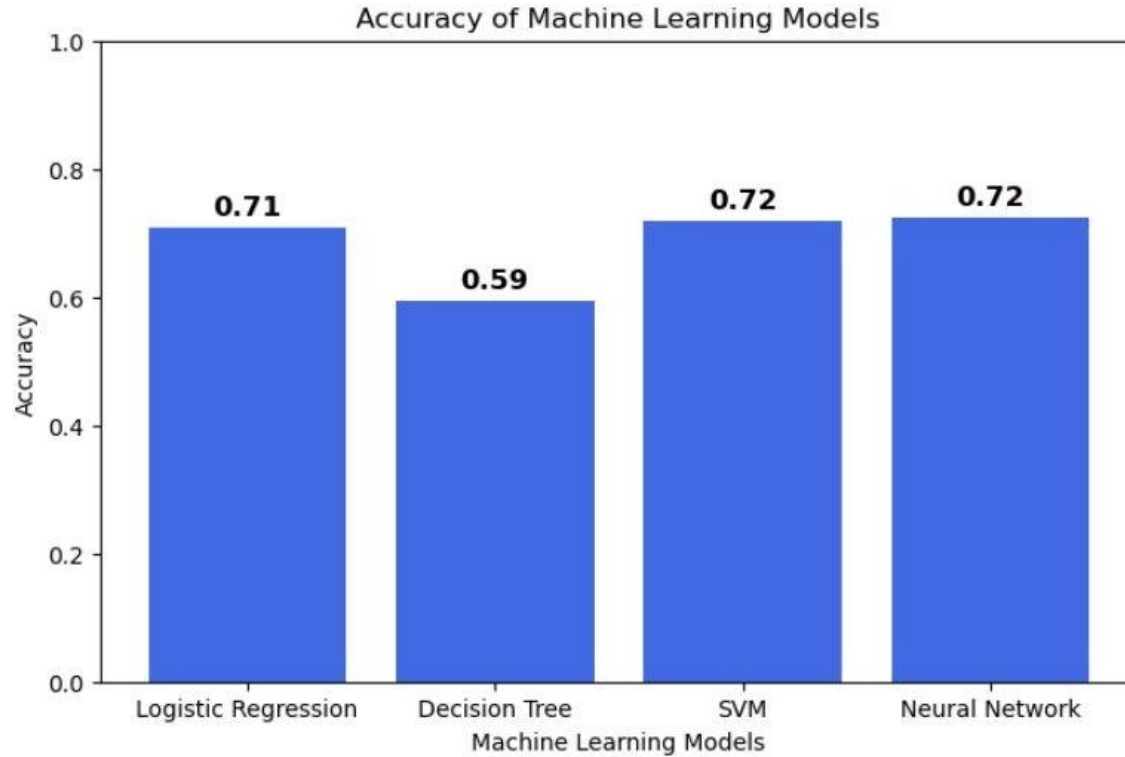
Goal : Classify instrumentalness into "Low," "Medium," "High."

Features Used : Energy, Speechiness, Acousticness, Danceability.

Models applied :

- **Logistic Regression :** 70.78% accuracy
Performed well but misclassified most "Medium" values.
- **Decision Tree :** 59.42% Accuracy : Balanced predictions but poor overall accuracy.
- **Support Vector Machine (SVM) :** 71.95% Accuracy : Better than Decision Tree, but still struggled with "Medium" values.
- **Neural Network (MLP Classifier) :** Best Model (72.30% Accuracy) : Outperformed others but still failed to predict "Medium" correctly.

CLASSIFICATION MODEL COMPARISON



IMPLEMENTING ReLU AND Softmax

Overall Accuracy: 72.59%

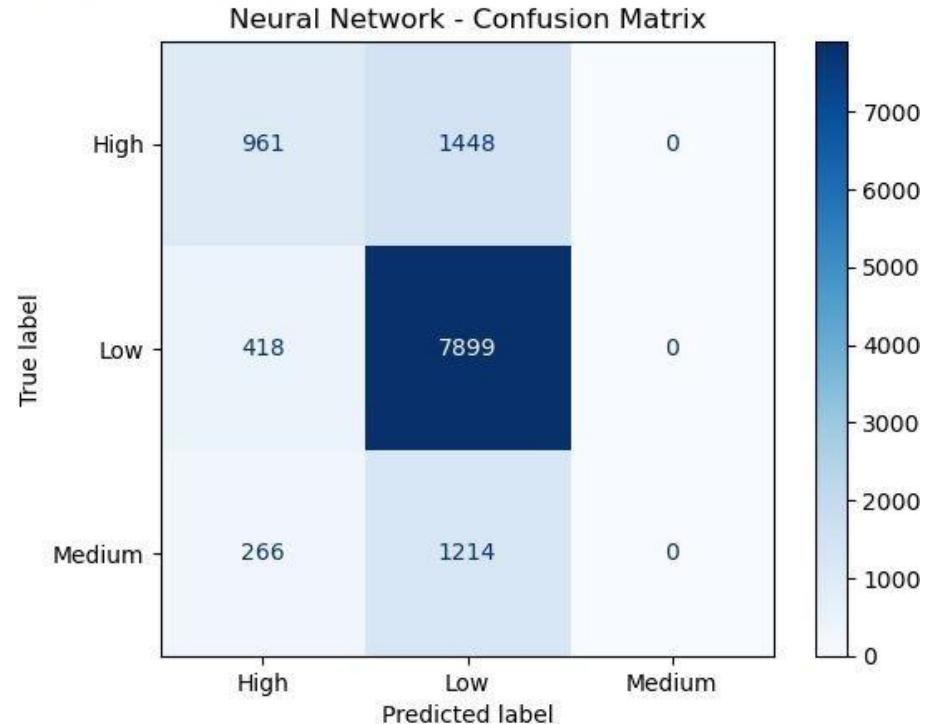
Insights:

- The model performs best in predicting Low class, with the highest correct classification count.
- High and Medium classes show misclassification, primarily into the Low category.
- The lack of Medium class predictions suggests potential class imbalance or insufficient feature separation.

Conclusion:

After using ReLU in hidden layers and Softmax in the output layer we still don't see any significant improvement in our model's accuracy.

Neural Network Accuracy: 0.7259




CONCLUSION

1. **Limitations Identified** – Class imbalance affected model predictions, and danceability variance wasn't fully explained ($R^2 = 0.481$).
2. **Result** : Neural Networks worked best for classification, Agglomerative Clustering was most effective, and Regression had moderate success in predicting danceability.
3. **Challenges Faced** : Class imbalance, feature selection, and computational limitations impacted model performance.
4. We learned that how important it is to deal with unbalanced data, scale features correctly, tune parameters carefully, and understand results accurately in machine learning. These basic things affect how well our models work, which is essential for making machine-learning projects successful.



RESEARCH PAPERS

- 
- **Spotify Data Analysis and Song Popularity Prediction** – This study explores how machine learning can predict a song's popularity on Spotify by analyzing audio features like danceability, loudness, and acousticness. Using models like Decision Trees and XGBoost, the study found that the Decision Tree classifier performed best with 93% accuracy, providing valuable insights for the music industry.
 - **Dance Hit Prediction** – This research focuses on predicting whether a dance song will be a Top 10 hit using machine learning models trained on audio features and chart history. The Logistic Regression model achieved an AUC of 0.67, suggesting that certain timbre-related features play a stronger role in determining hit potential than danceability or energy.
 - **Interactive Music Learning Model Based on RBF Algorithm** – The paper presents an AI-powered interactive music education system that leverages the Radial Basis Function (RBF) algorithm to personalize learning. By integrating multimedia tools, AI-driven emotional feedback, and adaptive learning, it creates a dynamic and immersive music education experience, making learning more engaging and effective.



Thank
You