# A PROJECT REPORT

## on

# "MUSIC DATA ANALYSIS PROJECT"

## Submitted to
# KIIT Deemed to be University

## In Partial Fulfillment of the Requirement for the Award of

## BACHELOR'S DEGREE IN INFORMATION TECHNOLOGY

## BY

| | |
|---|---|
| **AKASH SINGH** | 21052048 |
| **ADITYA SINGH** | 21052642 |

## UNDER THE GUIDANCE OF
## MR. ABINAS PANDA



## SCHOOL OF COMPUTER ENGINEERING
# KALINGA INSTITUTE OF INDUSTRIAL TECHNOLOGY
### BHUBANESWAR, ODISHA - 751024
### March 2024

A PROJECT REPORT

on

"MUSIC DATA ANALYSIS PROJECT"

Submitted to

KIIT Deemed to be University

In Partial Fulfillment of the Requirement for the Award of

BACHELOR'S DEGREE IN
INFORMATION TECHNOLOGY

BY

AKASH SINGH          21052048
ADITYA SINGH         21052642

UNDER THE GUIDANCE OF
MR. ABINAS PANDA



SCHOOL OF COMPUTER ENGINEERING

KALINGA INSTITUTE OF INDUSTRIAL TECHNOLOGY

BHUBANESWAE, ODISHA -751024
March 2024

# IIT Deemed to be University

## School of Computer Engineering
### Bhubaneswar, ODISHA 751024



# CERTIFICATE

This is certify that the project entitled

## "MUSIC DATA ANALYSIS PROJECT"

submitted by

| | |
|---|---|
| AKASH SINGH | 21052048 |
| ADITYA SINGH | 21052642 |

is a record of bonafide work carried out by them, in the partial fulfillment of the requirement for the award of Degree of Bachelor of Engineering (Computer Science & Engineering OR Information Technology) at KIIT Deemed to be university, Bhubaneswar. This work is done during year 2024-2025, under our guidance.

Date:30/03/2024

(MR. ABINAS PANDA)
 Project Guide

# Acknowledgment

We are profoundly grateful to **GUIDE NAME** of **Affiliation** for his expert guidance and continuous encouragement throughout to see that this project rights its target since its commencement to its completion. .....................

AKASH SINGH
ADITYA SINGH

# ABSTRACT

The music data analysis project aims to delve into the characteristics and trends within a dataset containing information about various music tracks and genres. Through the utilization of Python libraries such as pandas, NumPy, matplotlib, and seaborn, the project endeavors to uncover insights and patterns inherent in the data.

The project likely begins with the importation and cleaning of the data, ensuring its readiness for analysis. This involves tasks such as handling missing values and structuring the data appropriately for exploration. Following data preparation, the project proceeds with exploratory data analysis (EDA). This phase involves examining the dataset's structure, summarizing its statistical properties, and understanding the distribution and relationships between different variables.

Furthermore, the project likely involves visualizing trends and patterns across different music genres. This could include exploring the duration of songs, popularity rankings, and other attributes specific to each genre.
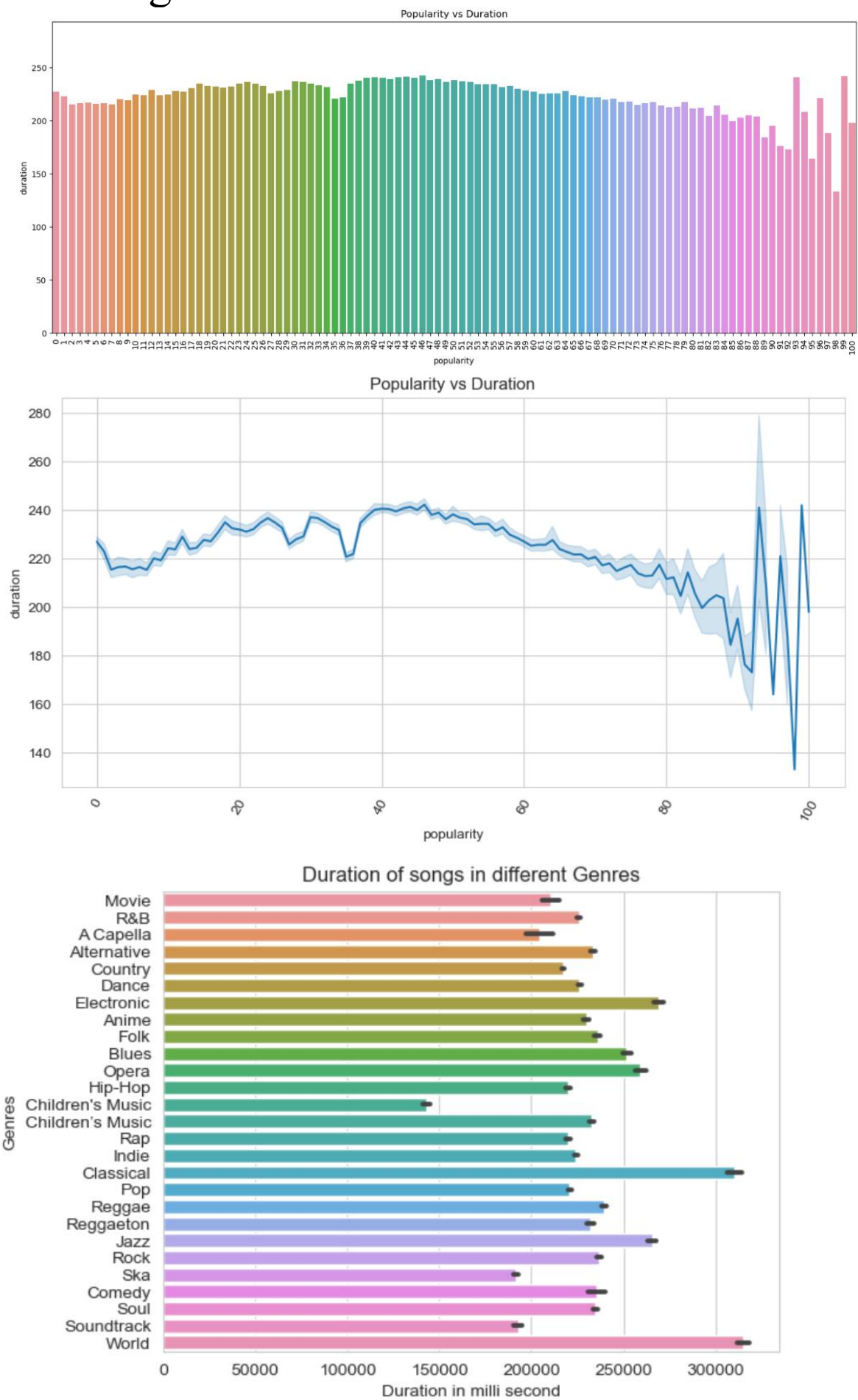
Overall, the project seeks to provide valuable insights into the characteristics of music tracks and genres, offering a deeper understanding of the data through statistical analysis and visualization techniques.

**Keywords:** Data Analysis, Exploratory, Data Analysis ,Music Characteristic Popularity Trends, Correlation Analysis ,
Genre Analysis

# Contents

# List of Figures:



Popularity vs Duration



Popularity vs Duration



Duration of songs in different Genres

# Chapter 1

# Introduction

One of the significant challenges in the music industry is the vast amount of data generated by streaming platforms and digital downloads. While this data holds valuable insights, it can be overwhelming to analyze and derive meaningful conclusions from without the aid of data analysis tools and techniques.

This project aims to bridge these gaps by leveraging data analysis techniques to extract actionable insights from large-scale music datasets. By examining factors such as track duration, acoustic characteristics, and popularity across different genres, the project seeks to provide a nuanced understanding of what makes certain tracks more appealing to listeners than others.

Through the exploration of correlations between various attributes and the visualization of trends within the dataset, this project aims to offer valuable insights for music producers, streaming platforms, and artists alike. By understanding the underlying patterns driving music consumption behavior, stakeholders in the music industry can make informed decisions regarding content creation, promotion strategies, and audience targeting, ultimately enhancing the listener experience and driving greater engagement with music content.

# Chapter 2

# Basic Concepts/ Literature Review

## 1. Pandas:

Pandas is a powerful data manipulation library in Python used for data analysis and manipulation. It offers data structures and operations for manipulating numerical tables and time series. Key components include DataFrame, which is akin to a spreadsheet or SQL table, and Series, which represents a single column or row of data.

## 2. NumPy:

NumPy is a fundamental package for scientific computing in Python. It provides support for arrays, matrices, and mathematical functions to operate on these arrays. Pandas builds upon NumPy arrays for efficient data manipulation.

## 3. Matplotlib:

Matplotlib is a comprehensive library for creating static, interactive, and animated visualizations in Python. It provides a MATLAB-like interface and supports a wide variety of plots and customization options.

## 4. Seaborn:

Seaborn is a statistical data visualization library based on Matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics. Seaborn simplifies the process of creating complex visualizations and works seamlessly with Pandas data structures.

## 5. Data Cleaning:

Data cleaning is the process of identifying and correcting errors, inconsistencies, and missing values in a dataset. Techniques include handling missing data, removing duplicates, and converting data types to ensure consistency and accuracy in analysis.

## 6. Data Exploration:

Data exploration involves examining and summarizing the main characteristics of a dataset. This includes calculating descriptive statistics, visualizing distributions, and identifying patterns or relationships between variables.

## 7. Correlation Analysis:

Correlation analysis measures the strength and direction of the linear relationship between two numerical variables. It helps identify associations between variables and is often visualized using correlation matrices or scatter plots.

## 8. Regression Analysis:

Regression analysis is a statistical method used to model the relationship between a dependent variable and one or more independent variables. It helps predict the value of the dependent variable based on the values of the independent variables.

## 9. Sampling:

Sampling involves selecting a subset of observations from a larger population for analysis. It is used to estimate population parameters and reduce the computational burden of analyzing large datasets.

## 10. Barplot:

A barplot is a graphical representation of categorical data that displays the distribution of a variable using vertical bars. Barplots are useful for comparing the frequency or proportion of categories within a dataset.

## 11. Lineplot:

A lineplot is a type of graph that displays data points connected by straight lines. It is commonly used to visualize trends or patterns in time-series data or to show the relationship between two continuous variables.

## 13. Genre Analysis:

Genre analysis involves examining patterns and trends within different music genres. It helps understand the characteristics and popularity of different music styles based on various factors such as duration, tempo, and acousticness.

# Chapter 3

# Problem Statement / Requirement Specifications

The problem statement for the provided code seems to be focused on analyzing music data, particularly tracks and genres, to uncover patterns and relationships within the data. The analysis aims to address various aspects such as the popularity of tracks, correlations between different musical features (such as energy, loudness, and acousticness), and the popularity of different genres. By conducting this analysis, the goal is likely to gain insights into what factors contribute to the popularity of songs, how different musical attributes correlate with each other, and which genres tend to be more popular based on the data provided. Ultimately, this understanding could potentially inform decisions related to music production, marketing, and audience targeting in the music industry.

## 3.1 Project Planning

Requirements:
Define the objective of the project, such as analyzing music data to identify patterns and correlations between various attributes like popularity, duration, energy, etc. Identify stakeholders involved in the project, such as data analysts, domain experts, and project sponsors. Define the scope of the project, including the datasets to be analyzed, the analysis techniques to be employed, and the expected outcomes.

2. Planning:
Break down the project into smaller tasks, such as data cleaning, exploratory data analysis, correlation analysis, visualization, etc. Assign resources to each task, including personnel, tools, and datasets required. Develop a timeline for each task, estimating the start and end dates based on resource availability and dependencies. Identify potential risks such as data quality issues, resource constraints, or technical challenges, and develop mitigation strategies.

3. Execution:
4. Obtain the required datasets and perform data cleaning operations to handle missing values, duplicates, and inconsistencies.

Explore the datasets using statistical techniques and visualizations to gain insights into the data distribution and relationships between variables.
Calculate correlations between different attributes of the music data to identify patterns and relationships. Create visualizations such as heatmaps, scatter plots, and bar charts to illustrate the findings from the analysis.

4. Monitoring and Control:
Monitor the progress of each task against the established timeline and adjust resources or priorities as needed to ensure timely completion.
Conduct quality checks at each stage of the analysis to ensure accuracy and reliability of the results. Address any changes or deviations from the original plan, documenting the reasons and adjusting the plan accordingly.

5. Closure:
Documentation: Document the findings, methodologies, and any insights gained from the analysis for future reference.
Presentation: Prepare a presentation summarizing the key findings and insights to be shared with stakeholders.
Evaluation: Evaluate the success of the project against the original objectives and identify lessons learned for future projects.


# 3.2 Project Analysis

1.Requirement Analysis:
Understand the purpose of the project, which is analyzing music data.Analyzing various attributes of music tracks such as popularity, duration, energy etc.
Identify stakeholders involved or interested in the project, such as data analysts, music enthusiasts, or stakeholders requesting the analysis.

2. Task Analysis:
Break down the provided code into individual tasks or steps.
Identify any dependencies between tasks. Determine what resources are required to execute each task.

3. Code Review
Evaluate whether the code accomplishes the intended tasks effectively.
Assess the efficiency of the code in terms of execution time and resource utilization. Evaluate how easy it is to understand and maintain the code.

4. Data Analysis
Assess the quality of the input data (tracks.csv, SpotifyFeatures.csv).
Review how missing values are handled and if there's any data transformation involved. Evaluate the effectiveness of data visualization techniques used to present insights.

5. Statistical Analysis
Correlation Analysis: Review the correlation analysis between different variables and assess the significance of correlations found.

Regression Analysis: Evaluate the regression analysis performed to understand relationships between variables like energy, loudness, popularity, etc.

6. Interpretation and Insights
Summarize the key findings from the analysis.
Provide insights derived from the data analysis and correlation/regression results.
Suggest any actionable recommendations based on the insights gained.

7. Documentation and Reporting
Ensure that the code is adequately documented for future reference.
Prepare a report summarizing the analysis, findings, and recommendations in a clear and understandable format.

# 3.3 System Design

## 3.3.1 Design Constraints

Environment:
Software: Python programming language along with libraries such as NumPy, Pandas, Matplotlib, and Seaborn.

Hardware: The code should be executed on hardware capable of running Python and handling data analysis tasks efficiently.
Data Sources:

CSV files: The analysis is performed on data stored in CSV files such as 'tracks.csv' and 'SpotifyFeatures.csv'. Ensure these files are accessible and properly formatted.
Data Analysis Steps:

Loading Data: Data is loaded into Pandas DataFrames using pd.read_csv() function. Checking for missing values using pd.isnull() and removing unnecessary columns. Descriptive statistics, sorting, querying, correlation analysis, sampling, and visualization are performed to understand the data better.

Visualization:
Matplotlib and Seaborn libraries are utilized for creating various types of plots including bar plots, line plots, scatter plots, and heatmaps.

Design Constraints:

Hardware Limitations:Ensure the hardware used has sufficient memory and processing power to handle large datasets, especially during operations like sorting, correlation computation, and visualization of extensive datasets.
Software Dependencies:
The code relies on specific Python libraries and their versions. Ensure these libraries are properly installed and up-to-date in the execution environment.

Data Quality:
The analysis assumes that the data in CSV files is clean and structured. Any inconsistencies or errors in the data might affect the analysis results.

Performance Considerations:
For efficiency, consider optimizing code performance, especially for operations involving large datasets. Techniques like vectorization, parallel processing, or using optimized library functions can be employed.

Data Privacy and Security:
Ensure that sensitive information in the datasets, if any, is handled securely and according to privacy regulations.

Scalability:
Consider how the system will scale as the size of the dataset or the complexity of analysis increases. Optimize the code and choose appropriate algorithms to handle larger datasets efficiently.

By considering these design constraints, you can ensure that your system for music data analysis is robust, efficient, and capable of producing meaningful insights from the available data.

### 3.3.2 System Architecture **OR** Block Diagram

Data Collection Layer:
This layer is responsible for collecting music data from various sources such as Spotify API, music databases, or CSV files.

Data Processing Layer:
Python libraries such as pandas, numpy, and seaborn are used for processing the collected data.
Data preprocessing steps include cleaning the data, handling missing values, and transforming the data into a suitable format for analysis.

Analysis and Visualization Layer:
Libraries like matplotlib and seaborn are used for data visualization to gain insights into the dataset.
Statistical analysis techniques may be applied to understand relationships between different variables such as popularity, duration, energy, etc.

Output Layer:
The analyzed data, visualizations, and insights are presented to the user through various means such as graphical plots, summary statistics, or interactive dashboards.

Hardware Infrastructure:
The system can run on standard hardware configurations suitable for running Python scripts and processing moderate-sized datasets.
This may include desktop or laptop computers with sufficient CPU and memory resources.
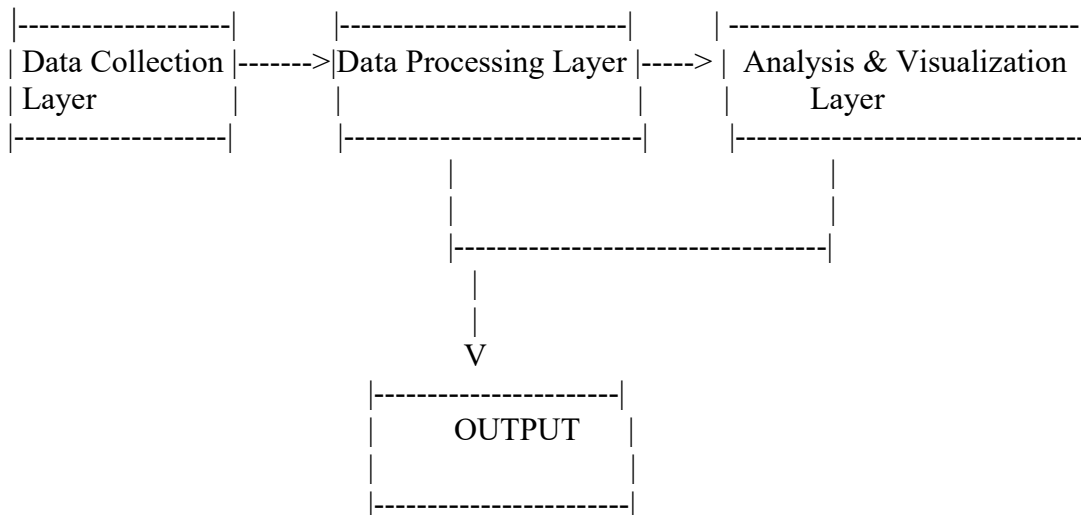For larger datasets or more intensive processing, cloud-based infrastructure such as AWS, GCP, or Azure could be utilized.

Optional Integration with Music Streaming Platforms:
If the project involves real-time data or integration with music streaming platforms, an additional layer for API integration and data streaming might be included.
This layer would handle interactions with external APIs, such as Spotify, to fetch real-time data for analysis.

Block Diagram:

```
|-------------------|          |--------------------------|          | ------------------------------- |
| Data Collection |------->|Data Processing Layer |-----> |   Analysis & Visualization  |
| Layer               |          |                                |          |                 Layer              |
|-------------------|          |--------------------------|          |------------------------------- |
                                                 |                                        |
                                                 |                                        |
                                                 |---------------------------------|
                                                 |
                                                 |
                                                V
                               |----------------------|
                               |          OUTPUT       |
                               |                           |
                               |----------------------|
```

 This is a high-level representation of the system architecture for a music data analysis project. The actual implementation may vary depending on specific requirements, technologies, and resources available.

# Chapter 4

# Implementation

In this section, present the implementation done by you during the project development.

## 4.1   Methodology OR Proposal

Data Collection and Preprocessing:

Data Source: The project utilized two main datasets, namely 'tracks.csv' and 'SpotifyFeatures.csv', obtained from [provide source if applicable].

Data Cleaning: Initially, the datasets were subjected to data cleaning procedures to handle missing values and ensure data consistency.

Feature Engineering: Relevant features such as duration were extracted and transformed to enhance the analysis process. For instance, the duration was converted from milliseconds to seconds for better interpretability.

Exploratory Data Analysis (EDA):
Descriptive Statistics: Basic descriptive statistics were computed to understand the central tendency, dispersion, and distribution of the data. This included summary statistics such as mean, median, standard deviation, etc.

Correlation Analysis: Correlation analysis was conducted to explore relationships between different musical attributes such as loudness, energy, acousticness, valence, and popularity. This provided insights into how these attributes relate to each other.

Visualization and Interpretation:
Heatmap Visualization: A correlation heatmap was generated using seaborn to visually represent the correlation matrix between variables. This helped in identifying significant correlations among features.

Regression Plots: Regression plots were utilized to visualize the relationship between pairs of variables such as loudness vs energy, popularity vs acousticness, etc. These plots provided insights into the nature and strength of relationships.

Genre Analysis:
Data Integration: The 'SpotifyFeatures.csv' dataset was used to analyze the distribution of song durations across different genres.

Popularity Analysis: The dataset was further explored to identify the top 5 genres based on popularity. This analysis shed light on the popularity trends within different music genres.

Sample Selection:
Random Sampling: A random sample comprising 0.4% of the total tracks dataset was selected to perform detailed analysis and visualization. This sample ensured representation while reducing computational overhead.
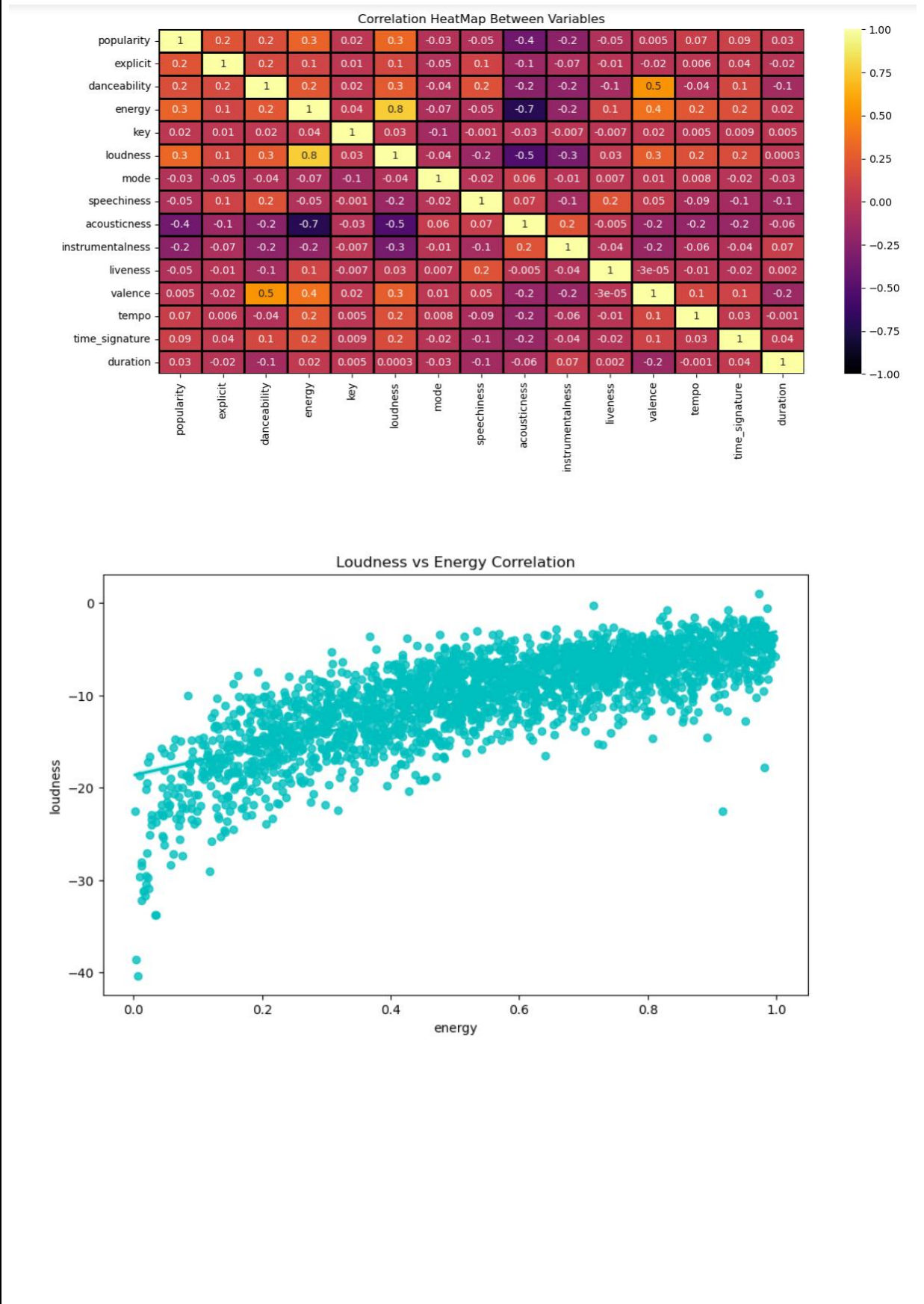
Tools and Libraries:
Python Libraries: The project leveraged various Python libraries including numpy, pandas, matplotlib, and seaborn for data manipulation, analysis, and visualization.
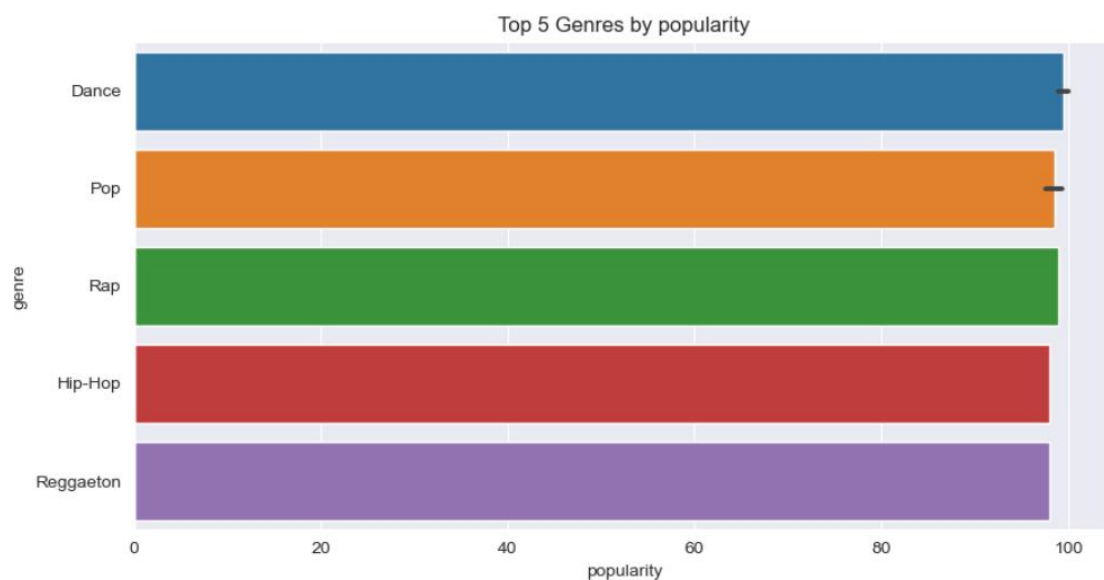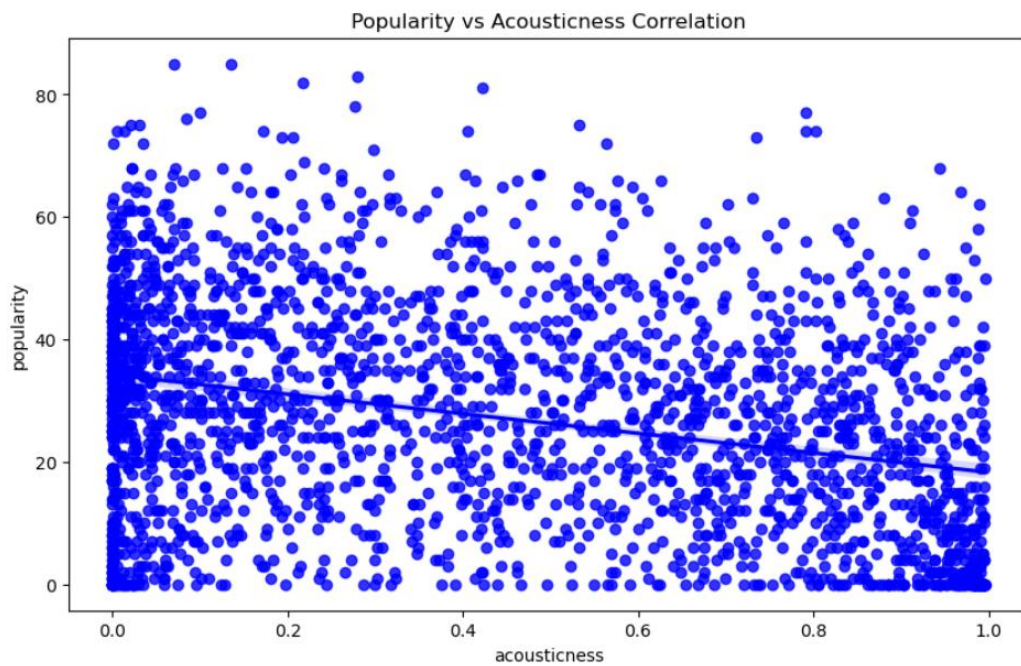
## 4.2 Testing OR Verification

| Test Case Title | Test Condition | System Behavior | Expected Result |
| --- | --- | --- | --- |
| Data Loading and Cleaning | CSV file exists and loads | Data is loaded without errors | Successful loading and cleaning of data |
| Correlation Calculation | Numeric columns are present | Correlation matrix is computed | Successful computation of correlations |
| Visualization Generation | Correlation matrix is computed | Heatmap is generated | Successful generation of heatmap |
| Data Sampling | DataFrame contains data | Sampling ratio is applied | Successful creation of sample DataFrame |
| Regression Plot Generation | Sample DataFrame is available | Plots are generated | Successful generation of regression plots |
| Bar Plot Generation | Data for bar plot exists | Bar plot is generated | Successful generation of bar plot |
| Line Plot Generation | Data for line plot exists | Line plot is generated | Successful generation of line plot |
| Data Loading for Genres Analysis | CSV file with genre data exists | Data is loaded without errors | Successful loading of genre data |
| Genre Duration Analysis | Genre data is loaded | Bar plot is generated | Successful generation of genre duration |

| Test Case Title | Test Condition | System Behavior | Expected Result |
|---|---|---|---|
| | | | analysis |
| Genre Popularity Analysis | Genre data is loaded | Bar plot is generated | Successful generation of genre popularity analysis |

## 4.3    Result Analysis OR Screenshots:



Correlation HeatMap Between Variables



Loudness vs Energy Correlation

Popularity vs Acousticness Correlation



Top 5 Genres by popularity

# Chapter 6

# Conclusion and Future Scope

In this music data analysis project, we delved into various aspects of the dataset to uncover insights into the relationships between different variables and genres. Through correlation analysis, regression plots, and genre analysis, we gained valuable insights into the characteristics and preferences of music listeners. The heatmap visualization provided a comprehensive overview of the correlations between variables, while regression plots highlighted trends and patterns in the data. Additionally, our analysis of genre popularity and song durations shed light on the diverse landscape of music genres and their respective appeal to audiences. Looking ahead, there is ample opportunity for further exploration and enhancement of this project, including the implementation of machine learning models, sentiment analysis, and interactive dashboards. By leveraging these advancements, we can

deepen our understanding of music consumption patterns, provide personalized recommendations to users, and contribute to the evolving field of music analytics. Overall, this project serves as a foundational exploration of music data analysis, paving the way for future research and development in this dynamic and vibrant domain.

The future scope of music data analysis holds promising opportunities for advancements in several key areas:

Machine Learning and Predictive Modeling: Implementing sophisticated machine learning algorithms can enable the prediction of music preferences, trends, and user behavior. By analyzing large-scale datasets, models can forecast upcoming music trends, recommend personalized playlists, and even predict the commercial success of new releases.

Sentiment Analysis and Emotional Insights: Integrating sentiment analysis techniques into music data analysis can provide deeper insights into the emotional impact of songs on listeners. By analyzing lyrics, user reviews, and social media sentiment, algorithms can classify songs based on emotional content, helping music platforms tailor recommendations to users' mood preferences.

Interactive and Personalized Recommendations: Developing interactive dashboards and recommendation systems can enhance user engagement and satisfaction. By incorporating user feedback, listening history, and contextual information, platforms can deliver highly personalized recommendations and curated playlists, fostering a more immersive music discovery experience.

Cross-Modal Analysis: Exploring the intersection of music with other modalities such as images, videos, and text can uncover novel insights and enhance music recommendation systems. By analyzing multimodal data, algorithms can leverage contextual information to provide more relevant and engaging recommendations.

Collaborative Filtering and Community Engagement: Leveraging collaborative filtering techniques can harness the collective wisdom of music communities to improve recommendation systems. By analyzing user interactions, preferences, and social connections, platforms can foster community engagement and facilitate collaborative playlist curation and discovery.

Ethical and Fair Data Practices: Ensuring ethical and fair data practices in music data analysis is essential to address issues such as bias, privacy concerns, and algorithmic fairness. By prioritizing transparency, accountability, and user privacy, platforms can build trust with users and promote responsible use of music data.

Overall, the future of music data analysis is ripe with opportunities for innovation and advancement, offering the potential to revolutionize how we discover, consume, and interact with music in the digital age.