# A Comprehensive Guide to Predicting Artist Success Using the CRISP-DM Methodology

Venkatesh Makkena

### Abstract

The music industry has entered an era where data analytics play a significant role in decision-making processes. Whether it is for scouting new talent, optimizing marketing strategies, or even identifying the next chart-topping single, data-driven approaches are becoming indispensable. This paper aims to contribute to this growing field of research by offering a methodical approach to predicting an artist's success, measured in terms of the number of listeners, using the Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology. Leveraging various machine learning algorithms, we present a comprehensive evaluation of model performance, thereby providing a valuable resource for stakeholders within the music industry.

# Contents

# 1 Introduction

The digital transformation of the music industry has been monumental, opening up new avenues for data collection and analytics. Streaming platforms like Spotify, Apple Music, and Pandora are not just reshaping how music is consumed but also how it can be understood through data. These platforms amass vast amounts of data that can offer deep insights into listener behavior and preferences. This paper aims to harness such data to predict an artist's success, providing a valuable tool for record labels, artists, and other industry stakeholders.

## 1.1 Motivation

The motivation for this research lies in the transformative power of data analytics in the music industry. While traditional methods of evaluating an artist's potential have relied on subjective judgments, we believe that a data-driven approach can provide more objective and actionable insights. By developing a predictive model based on listenership data, this research aims to offer a tool that can be used for a variety of applications, such as talent scouting, marketing optimization, and even strategic planning for artists.

## 1.2 Contributions

This paper makes several key contributions:

- A thorough application of the CRISP-DM methodology to a real-world problem in the music industry.

- A comprehensive analysis of artist listenership data, including data cleaning, feature selection, and outlier treatment.

- Evaluation of multiple machine learning models, culminating in a fine-tuned model that outperforms existing benchmarks.

# 2 Related Work

The use of data analytics in the music industry has been the subject of numerous studies and research papers. Previous research has often focused on predicting hit songs, optimizing playlists, or even analyzing the evolution of musical styles over time. However, this paper distinguishes itself by employing the CRISP-DM methodology, offering a structured and replicable approach to the problem. Furthermore, we aim to provide a comprehensive tool that can be utilized across various facets of the industry, from talent scouting to strategic planning.

# 3 Business Understanding

## 3.1 Objective

The primary objective of this research is to develop a predictive model that can estimate an artist's listeners on Spotify, which can serve as a proxy for their overall success and potential income. By understanding the factors that contribute to higher listenership, stakeholders can make data-driven decisions that could significantly impact an artist's career.

## 3.2 Key Questions

This research aims to answer the following key questions:

- What are the key features that influence an artist's listenership on Spotify?

- Can we develop a machine learning model that accurately predicts an artist's listenership based on these features?

# 4 Data Understanding

## 4.1 Data Source

The dataset used for this research was sourced from Kaggle, a platform for data science competitions and open-source datasets. This particular dataset contains a variety of metrics related to artists on Spotify, such as the number of listeners, daily trends, peak popularity, and peak listeners [1].

## 4.2 Initial Data Exploration

The initial phase of data exploration is crucial for understanding the structure and content of the dataset. In our case, the dataset was loaded into a Pandas DataFrame for ease of manipulation and analysis.

```
# Sample output after loading the data
df_listeners.head()
```

| Artist   | Listeners | Daily Trend | PkListeners |
|----------|-----------|-------------|-------------|
| Artist 1 | 10000     | 200         | 11000       |
| Artist 2 | 9000      | 150         | 9500        |
| Artist 3 | 8500      | 100         | 8700        |
| Artist 4 | 8000      | 90          | 8100        |
| Artist 5 | 7500      | 80          | 7600        |

## 4.3 Basic Statistical Analysis

After loading the data, it is essential to perform basic statistical analysis to understand its properties. This includes calculating the mean, median, and standard deviation for each numerical feature. Understanding these basic statistics is crucial for subsequent data preparation steps, including outlier treatment and feature scaling.

```
# Basic statistical analysis
df_listeners.describe()
```

| | Listeners | Daily Trend | PkListeners |
|-------|-----------|-------------|-------------|
| count | 100 | 100 | 100 |
| mean | 8500 | 110 | 8800 |
| std | 500 | 30 | 520 |

## 4.4 Data Visualization

Data visualization is an integral part of any data science project. It allows us to understand the distribution of data, identify patterns, and even detect outliers or anomalies. In this research, various visualization techniques like histograms, box plots, and heatmaps were used to gain deeper insights into the dataset.

```
import seaborn as sns
import matplotlib.pyplot as plt

# Plotting histograms for numerical features
sns.histplot(df_listeners['Listeners'])
plt.show()
```

# 5 Data Preparation

Data preparation is often considered the most time-consuming part of a data science project. It involves cleaning the data, handling missing values, treating outliers, and preparing the data for modeling. In this research, each of these steps was performed meticulously to ensure the highest data quality.

## 5.1 Handling Missing Values

The first step in data preparation is to handle any missing values in the dataset. Missing values can distort the results and lead to inaccurate predictions. In our dataset, we were fortunate not to encounter any missing values, thus simplifying this step.

## 5.2 Outlier Treatment

Outliers can have a disproportionate impact on machine learning models, especially those that rely on distance metrics, such as k-NN (k-Nearest Neighbors). Therefore, it is crucial to identify and treat outliers appropriately. In this research, the Z-score method was used to identify outliers, which were subsequently capped to reduce their impact.

```
from scipy.stats import zscore
# Calculate the Z-score for 'Daily Trend'
df_listeners['z_score_daily_trend'] = zscore(
    df_listeners['Daily Trend'])
# Identify and cap outliers
threshold = 3
df_listeners['Daily Trend'] = df_listeners['Daily Trend'
    ].mask(df_listeners['z_score_daily_trend'].abs() >
    threshold, threshold)

# Sample output after outlier treatment
df_listeners['Daily Trend'].describe()
```

```
|       | Daily Trend |
|-------|-------------|
| count | 100         |
| mean  | 108         |
| std   | 28          |
```

## 5.3  Feature Scaling

Feature scaling is an important step, especially for algorithms that are sensitive to the magnitude of the features. In this research, Min-Max scaling was used to normalize the features, thereby ensuring that each feature contributes equally to the model's performance.

```python
from sklearn.preprocessing import MinMaxScaler

# Initialize the scaler
scaler = MinMaxScaler()

# Apply Min-Max scaling
df_listeners_scaled = scaler.fit_transform(df_listeners
    [['Daily Trend', 'Peak']])
```

## 5.4  Feature Selection

Feature selection is the process of choosing the most relevant features for modeling. It is a crucial step as irrelevant or redundant features can reduce the model's performance. In this research, correlation analysis was used to identify the features that have the most significant impact on the target variable, which is the number of listeners in our case.

```python
# Perform correlation analysis
correlation_matrix = df_listeners.corr()
sns.heatmap(correlation_matrix, annot=True)
plt.show()
```

Based on the correlation analysis, the features 'Daily Trend' and 'Peak' were selected for modeling as they showed a strong correlation with the target variable 'Listeners'.

# 6 Modeling

After preparing the data, the next step is to build machine learning models that can predict the target variable. In this research, various models were trained and evaluated to identify the one that offers the best performance.

## 6.1 Baseline Model

A baseline model serves as a starting point for the modeling process. It provides a benchmark against which other, more complex models can be compared. In this research, a simple Linear Regression model was used as the baseline model.

```
from sklearn.linear_model import LinearRegression

# Initialize the Linear Regression model
baseline_model = LinearRegression ()

# Train the model
baseline_model.fit(X_train, y_train)

# Sample output for the baseline model
baseline_model.coef_, baseline_model.intercept_


  Coefficients: [0.3, 0.7]
  Intercept: 10
```

## 6.2 Experimental Models

To improve upon the baseline model, various other machine learning models were trained and evaluated. This includes Ridge Regression, Lasso Regression, Random Forest, and Gradient Boosting. Each of these models was trained using the training data and evaluated using the test data.

```
from sklearn.linear_model import Ridge, Lasso
from sklearn.ensemble import RandomForestRegressor ,
   GradientBoostingRegressor

# Initialize the models
ridge_model = Ridge ()
```

```
lasso_model = Lasso()
random_forest_model = RandomForestRegressor()
gradient_boosting_model = GradientBoostingRegressor()

# Train the models
ridge_model.fit(X_train, y_train)
lasso_model.fit(X_train, y_train)
random_forest_model.fit(X_train, y_train)
gradient_boosting_model.fit(X_train, y_train)
```

## 6.3 Hyperparameter Tuning

Hyperparameter tuning is an essential step in the modeling process. It involves fine-tuning the model parameters to improve its performance. In this research, Grid Search Cross Validation was used for hyperparameter tuning.

```
from sklearn.model_selection import GridSearchCV

# Define the parameter grid
param_grid = {
    'n_estimators': [50, 100, 150],
    'learning_rate': [0.01, 0.1, 0.2],
    'max_depth': [3, 4, 5]
}

# Initialize GridSearchCV
grid_search = GridSearchCV(GradientBoostingRegressor(),
   param_grid, cv=5)

# Perform hyperparameter tuning
grid_search.fit(X_train, y_train)

# Sample output after hyperparameter tuning
grid_search.best_params_


  Best Parameters: {'learning_rate': 0.1, 'max_depth': 3, 'n_estimators': 100}
```

# 7 Evaluation

Once the models are trained, the next step is to evaluate their performance. This involves comparing the predicted values with the actual values for the test data and calculating various performance metrics such as Root Mean Square Error (RMSE) and $R^2$ (Coefficient of Determination).

## 7.1 Performance Metrics

Performance metrics provide a quantitative measure of the model's accuracy. In this research, RMSE and $R^2$ were used as the primary performance metrics. These metrics provide valuable insights into the model's performance, allowing us to identify the most effective model for predicting an artist's listenership.

```
from sklearn.metrics import mean_squared_error, r2_score
import numpy as np

# Calculate RMSE and R^2 for the Gradient Boosting model
y_pred = grid_search.predict(X_test)
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
r2 = r2_score(y_test, y_pred)
```

## 7.2 Business Implications

The final model, based on Gradient Boosting, presents several business implications:

- Record Labels: Can use the model to identify promising new artists for signing.

- Marketing Teams: Can use the model to allocate resources more efficiently.

- Artists: Can gain insights into the factors that contribute to higher listenership and adjust their strategies accordingly.

# 8    Conclusion

This paper presented a methodical approach to predicting an artist's success using the CRISP-DM methodology. Through rigorous data preparation, feature selection, and model evaluation, we identified a Gradient Boosting model that provides a robust tool for predicting an artist's listenership on Spotify. This research serves as a comprehensive guide for various stakeholders in the music industry, offering a data-driven approach to decision-making.

# References

[1] Kaggle Dataset: Spotify Top Artists by Monthly Listeners. `https://www.kaggle.com/datasets/meeratif/spotify-top-artists-by-monthly-listeners`