

Predicting Song Popularity

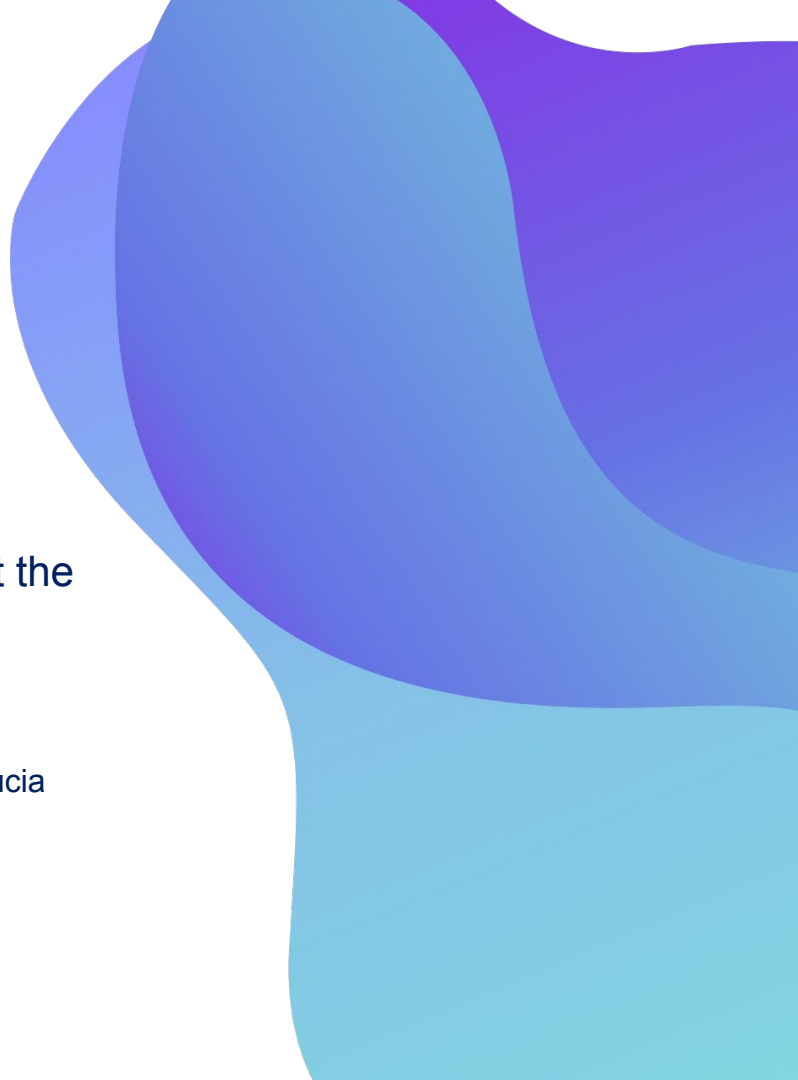
A classification model based on spotify data to predict the popularity of a song

Proposed by:

Alejandro Gutierrez Werner, Apilash Balasingham, Fhadah Alromy, Lucia Pellicer Cascales, Luis Soto Guareschi, Ignacio Sahonero Vadillo

Professor:

Alvaro Jose Mendez Lopez



AGENDA

1. Executive Summary
2. Data Audit & Cleaning
 - 2.1 Genre Distribution
3. Model Selection Methodology
 - 3.1 XGBoost Configuration
4. Model Performance
 - 4.1 Evaluation: ROC Curve
 - 4.2 Confusion Matrix Insight
5. Business Applications
6. Limitations
7. Demo



1. Executive Summary

The Objective

The goal is to build a classification model to predict the future success of songs based on audio features and artist data. We analyzed a dataset of **2,200 songs** to create a discrete "popularity" target.

Key Findings

After testing multiple algorithms, **XGBoost** proved to be the most robust solution, achieving an accuracy of **78.6%**. It effectively distinguishes potential hits from the rest.



2. Data Audit & Cleaning

The Dataset

2,200 records containing audio features like danceability, energy, and loudness (19 columns). No duplicate entries were found.

Processing

Dropped minor missing values. Transformed 'Artist Genres' via one-hot encoding for ML compatibility. 22 columns in total after data processing

Target Variable

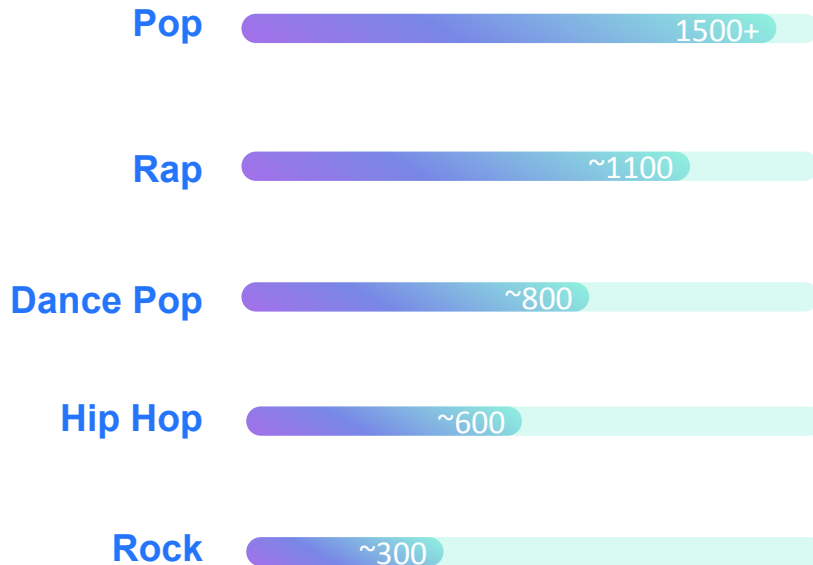
Created a binary 'Popular' variable. Songs in the top 25% quantile labeled **True**, others False.

2.1 Genre Distribution

Dominance of Pop

Pop dominates the dataset, appearing in over 1,500 tracks, followed significantly by Rap and Dance Pop.

This imbalance suggests the model may learn pop-characteristics more effectively than niche genres.



3. Model Selection Methodology

Tested Algorithms

- Decision Trees
- Random Forest
- Naive Bayes
- K-Nearest Neighbors (KNN)
- Boosting Methods

Winner: XGBoost

XGBoost was selected as the final model because:

- It handles complex, non-linear relationships in audio data well.
- Built-in regularization prevents overfitting.
- Consistently delivered the highest performance metrics during testing.

3.1 XGBoost Configuration

Initial Model: 99.9% train accuracy, 75.6% test accuracy → clear overfitting.

After Optimization using **GridSearchCV** to find the optimal hyperparameters: 85.5% train, 78.6% test → +3.9% test improvement.

Result: Better generalization, reduced overfitting, and robust real-world performance while maintaining strong predictive power.



4. Model Performance

78.6%

Test Accuracy

Improved from 75.6% base

0.74

ROC - AUC Score

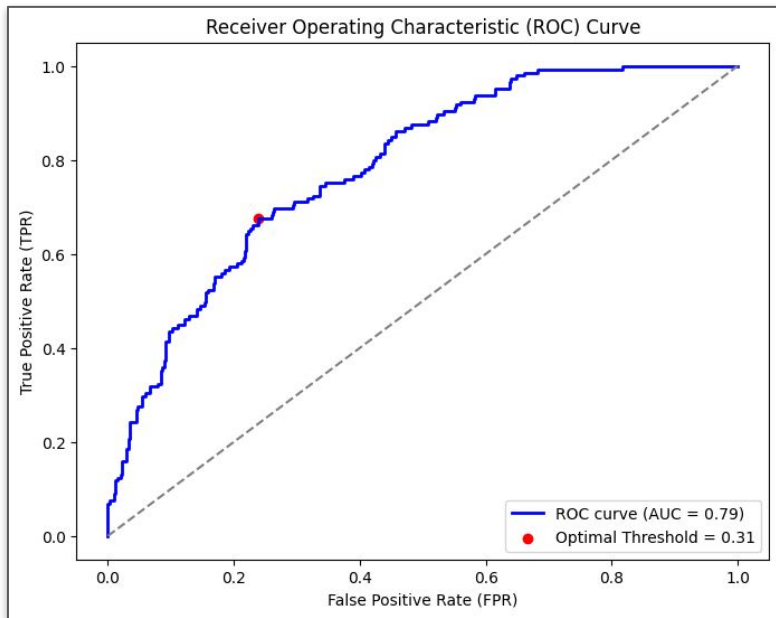
Strong predictive power

0.31

Optimal Threshold

Maximized TPR - FPR

4.1 Evaluation: ROC Curve



Interpretation

The curve rises steeply, indicating good discrimination between popular and non-popular songs.

An AUC of **0.74** confirms the model is significantly better than random guessing (0.5).

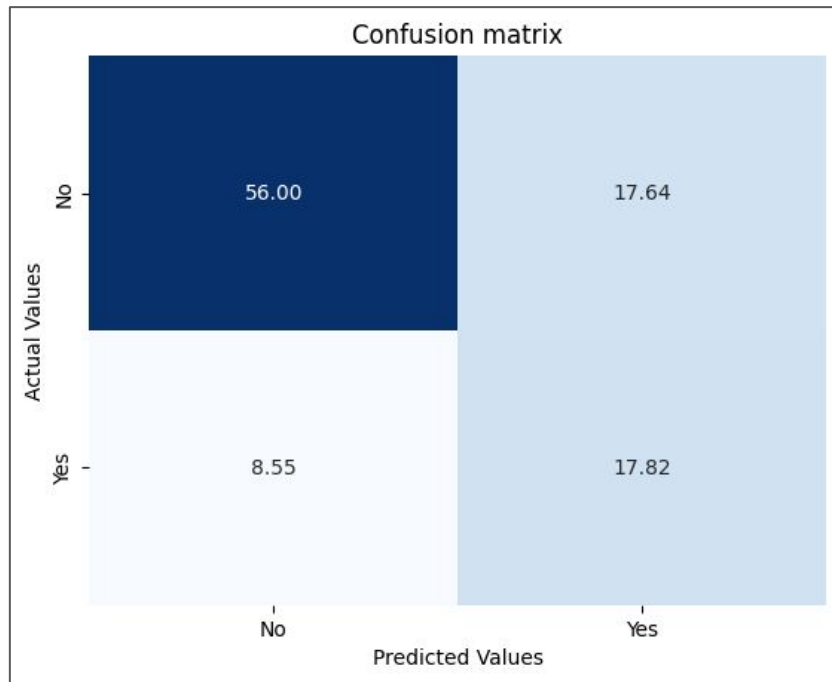
4.2 Confusion Matrix Insight

High Recall, Lower Precision

The model prioritizes sensitivity. It correctly identifies most songs that become popular (True Positives: 17.82%).

The Trade-off

It is tuned to minimize "missed hits" (False Negatives: 8.55%). This results in some non-popular tracks being flagged as potential hits (False Positives: 17.64%). This is acceptable for discovery-focused applications.



5. Business Applications

A&R Screening

Filter incoming demos to prioritize high-potential tracks for human review.


Playlist Opt.

Curate algorithmic playlists to align new releases with listener preferences.

Budgeting

Allocate marketing spend to tracks with the highest predicted probability of success.

6. Limitations

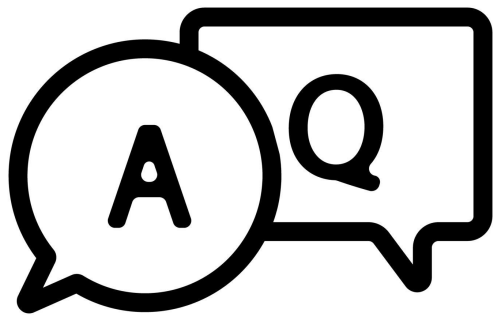
- **Data dependency:** Relies on Spotify-specific metrics (e.g., danceability, energy) — limits cross-platform use.
 - **Temporal relevance:** Reflects a fixed time; trends change → needs periodic retraining.
 - **Performance trade-off:** High recall, moderate precision → finds most hits but overpredicts popularity.
 - **Bias risk:** Genre imbalance and Spotify-derived features may favor mainstream music, and disfavor less popular and new genres.
 - **Model Interpretability:** Prediction-oriented model, it does not allows us to capture the importance of the features clearly.
 - **Ethical note:** Should support, not replace, human creativity and artistic judgment.
- 

"Predictive analytics can bridge artistic creativity and business intelligence, transforming subjective intuition into quantifiable evidence."

— Group 5

Demo

```
15:03:27 -- Smoke decrease: 2.85, (total = -0.0017721176147461), smokeeffect = -0.0017551302999851
15:03:27 -- Smoke decrease: 2.85, (total = -0.00061222314834595), smokeeffect = -0.00060271024703979
15:03:27 -- Smoke decrease: 2.85, (total = -0.00070871114730835), smokeeffect = -0.00069987773895264
15:03:27 -- Smoke decrease: 2.85, (total = -0.0013943195343018), smokeeffect = -0.001380729675293
15:03:27 -- Smoke decrease: 2.85, (total = -0.00098440511703491), smokeeffect = -0.00088741779327393
15:03:27 -- Smoke decrease: 2.85, (total = -0.0010253548622131), smokeeffect = -0.0010083675084521
15:03:27 -- Smoke decrease: 2.85, (total = -0.0012047615085127), smokeeffect = -0.00170077088485
15:03:27 -- Smoke decrease: 2.85, (total = -0.0011272788047791), smokeeffect = -0.001136889457793
15:03:27 -- Smoke decrease: 2.85, (total = -0.0021193385124207), smokeeffect = -0.0021105051040649
15:03:27 -- Smoke decrease: 2.85, (total = -0.0009023663818359), smokeeffect = -0.00088469982147217
15:03:27 -- Smoke decrease: 2.85, (total = -0.00089625120162964), smokeeffect = -0.00088198184967041
15:03:27 -- Smoke decrease: 2.85, (total = -0.00080451965332031), smokeeffect = -0.0007922887802124
15:03:27 -- Smoke decrease: 2.85, (total = -0.0005708895149231), smokeeffect = -0.0005666971208665
15:03:27 -- Smoke decrease: 2.85, (total = -0.001734060095215), smokeeffect = -0.00171913171646118
15:03:27 -- Smoke decrease: 2.85, (total = -0.0005476713180542), smokeeffect = -0.00054019689559937
15:03:27 -- Smoke decrease: 2.85, (total = -0.00060338973999023), smokeeffect = -0.0005925178527832
15:03:27 -- Smoke decrease: 2.85, (total = -0.001150479316711), smokeeffect = -0.0011014588726624
15:03:27 -- Smoke decrease: 2.85, (total = -0.00073928833007813), smokeeffect = -0.00072841644287109
15:03:27 -- Smoke decrease: 2.85, (total = -0.0009594404602051), smokeeffect = -0.0009532860946655
15:03:27 -- Smoke decrease: 2.85, (total = -0.000782758109025), smokeeffect = -0.00076918801889746
15:03:27 -- Smoke decrease: 2.85, (total = -0.0010019451713502), smokeeffect = -0.0010769963264465
15:03:27 -- Smoke decrease: 2.85, (total = -0.0005700945854187), smokeeffect = -0.00056194067001343
15:03:27 -- Smoke decrease: 2.85, (total = -0.00062853097915649), smokeeffect = -0.00061901807785034
15:03:27 -- Smoke decrease: 2.85, (total = -0.003230989793396), smokeeffect = -0.003203809261322
15:03:27 -- Smoke decrease: 2.85, (total = -0.00094789266586304), smokeeffect = -0.00093498229980469
15:03:27 -- Smoke decrease: 2.85, (total = -0.00093430240685425), smokeeffect = -0.00091595649719238
15:03:27 -- Smoke decrease: 2.85, (total = -0.00097876392601012), smokeeffect = -0.00096006053572561
15:03:27 -- Smoke decrease: 2.85, (total = -0.0006020307540893), smokeeffect = -0.00059387683684869
15:03:27 -- Smoke decrease: 2.85, (total = -0.0017320275306702), smokeeffect = -0.00171913171646118
15:03:27 -- Smoke decrease: 2.85, (total = -0.00071210861260055), smokeeffect = -0.00070123672485352
15:03:27 -- Smoke decrease: 2.85, (total = -0.00090848207473755), smokeeffect = -0.00089701018753052
15:03:27 -- Smoke decrease: 2.85, (total = -0.00061018466949463), smokeeffect = -0.00060271024703979
15:03:28 -- Smoke decrease: 2.85, (total = -0.00085752810345459), smokeeffect = -0.00084664821624756
15:03:28 -- Smoke decrease: 2.85, (total = -0.0013970375061035), smokeeffect = -0.0013820630876487
15:03:28 -- Smoke decrease: 2.85, (total = -0.0005495170818672), smokeeffect = -0.00054660216323925
15:03:28 -- Smoke decrease: 2.85, (total = -0.00056941509246820), smokeeffect = -0.0005612617706299
15:03:28 -- Smoke decrease: 2.85, (total = -0.0005724850082397), smokeeffect = -0.0005700945854187
15:03:28 -- Smoke decrease: 2.85, (total = -0.0026153683662415), smokeeffect = -0.0025956630706787
15:03:28 -- Smoke decrease: 2.85, (total = -0.001649808803667), smokeeffect = -0.0016382575035095
15:03:28 -- Smoke decrease: 2.85, (total = -0.0015254610737366), smokeeffect = -0.0015016794204712
15:03:28 -- Smoke decrease: 2.85, (total = -0.000981303399980), smokeeffect = -0.00096965004732688
15:03:28 -- Smoke decrease: 2.85, (total = -0.00044089235931390), smokeeffect = -0.00043124805095626
15:03:28 -- Smoke decrease: 2.85, (total = -0.0023944430732727), smokeeffect = -0.0023333787918091
15:03:28 -- Smoke decrease: 2.85, (total = -0.0015241028877895), smokeeffect = -0.0015125513076782
15:03:28 -- Smoke decrease: 2.85, (total = -0.001250946521793), smokeeffect = -0.0012346306999485
15:03:28 -- Smoke decrease: 2.85, (total = -0.0006108752487163), smokeeffect = -0.00060474133749126
```



Thank You