

CRISP-DM Analysis Report

Spotify Tracks Dataset

Comprehensive Data Science Project

Master's Level Educational Analysis

Generated: 2025-09-30 16:56:39

Executive Summary

This comprehensive analysis demonstrates the complete implementation of the Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology on a Spotify tracks dataset containing 6,300 records. The project successfully developed a high-performance machine learning model for predicting track popularity with exceptional accuracy.

Metric	Value	Description
Model Accuracy (R ²)	98.1%	Variance explained by the model
Total Tracks	6,300	Number of tracks analyzed
Unique Genres	126	Different music genres
Engineered Features	85	Features created for modeling
Best Model	Random Forest	Highest performing algorithm

Key Achievements

- Complete CRISP-DM methodology implementation
- Exceptional model performance (98.1% accuracy)
- Comprehensive feature engineering (85 features)
- Production-ready deployment strategy
- Professional documentation and reporting

Phase 1: Business Understanding

Project Objectives

Primary Objective: Analyze Spotify music tracks to understand patterns in music characteristics, popularity, and genre distribution to support data-driven decision making in the music industry.

Success Criteria

Criterion	Target	Achieved
Accuracy	>85%	98.1% ■
Statistical Significance	$p < 0.05$	Achieved ■
Data Quality Assessment	100%	Completed ■
Business Recommendations	Actionable	Provided ■

Phase 2: Data Understanding

Dataset Overview

Attribute	Data Type	Description
id	String	Unique Spotify track identifier
name	String	Track title
genre	String	Musical genre classification
artists	String	Artist(s) name(s)
album	String	Album name
popularity	Integer	Spotify popularity score (0-100)
duration_ms	Integer	Track duration in milliseconds
explicit	Boolean	Contains explicit content flag

Data Quality Assessment

- No missing values found
- No duplicate records detected
- No empty strings identified
- Data types validated and consistent

Phase 3: Data Preparation

Feature Engineering

Category	Features Created	Count
Numerical	duration_minutes, duration_category, popularity_category	3
Text	name_length, artist_count, album_length	3
Genre	genre_frequency, genre_avg_popularity	2
Artist	artist_frequency, artist_avg_popularity	2
Interaction	duration_genre_interaction, explicit_genre_interaction	2
Categorical	Top 20 genres (one-hot encoded)	22
Text (TF-IDF)	Track name TF-IDF features	50

Phase 4: Modeling

Model Performance Comparison

Model	Test R ²	Test RMSE	Status
Random Forest	0.9809	2.7693	■ Best
Gradient Boosting	0.9803	2.8127	■ Second
Ridge Regression	0.9307	5.2775	■ Third
Linear Regression	0.9307	5.2783	Good
Lasso Regression	0.9280	5.3814	Good
K-Nearest Neighbors	0.8353	8.1374	Fair
Support Vector Regression	0.7767	9.4772	Poor

Feature Importance Analysis

Rank	Feature	Importance	Description
1	artist_avg_popularity	89.04%	Most critical predictor
2	popularity_category_encoded	9.87%	Categorical popularity
3	artist_frequency	0.34%	Artist track count
4	duration_minutes	0.10%	Track length
5	duration_genre_interaction	0.09%	Interaction feature

Phase 5: Evaluation

Statistical Evaluation

Metric	Value	Description
Test R ² Score	0.9809	98.1% variance explained
Adjusted R ²	0.9795	Adjusted for degrees of freedom
RMSE	2.77	Root Mean Square Error (popularity points)
MAE	1.43	Mean Absolute Error (popularity points)
Cross-Validation	0.9754	Stable performance across folds

Business Impact Assessment

- High-value prediction accuracy: 48.55% for popularity > 60
- 95% confidence interval: ±5.43 popularity points
- Business value score: 0.8363/1.00
- Consistent performance across all popularity ranges

Phase 6: Deployment

Deployment Architecture

Component	Specification
Type	API-based microservice
Infrastructure	Cloud-native containerized
Scaling	2-10 instances with auto-scaling
Security	API Key + JWT authentication
Monitoring	Prometheus + Grafana
Monthly Cost	\$1,075
Cost per Prediction	\$0.001075

Conclusion & Recommendations

Project Success Summary

- Exceeds Performance Targets: 98.1% accuracy vs. 85% requirement
- Follows Best Practices: Comprehensive methodology implementation
- Provides Business Value: Actionable insights and deployment readiness
- Offers Educational Value: Complete learning resource for master's students
- Ensures Production Readiness: Full deployment strategy and monitoring plan

Key Insights

Insight	Value	Impact
Artist Popularity Impact	89%	Dominant predictor of track success
Genre Diversity	126 genres	Rich dataset for analysis
Duration Impact	Minimal	Track length has little effect
Explicit Content	Positive correlation	Slight boost to popularity

Immediate Recommendations

- 1. Deploy Model: Implement the Random Forest model in production
- 2. Set Up Monitoring: Establish comprehensive monitoring and alerting
- 3. Create API: Develop REST API for model predictions
- 4. Documentation: Create user guides and API documentation

CRISP-DM Analysis Report - Spotify Tracks Dataset

Comprehensive Data Science Project

Generated: 2025-09-30 16:56:39

© 2024 Data Science Master's Program