**PG Certificate Course in Data Science, AI/ML and Data Engineering by IIT Roorkee**

**Final Project Submission –Hemant Agarwal**

# Agenda

# Agenda (Conti.)

# Agenda (Conti.)
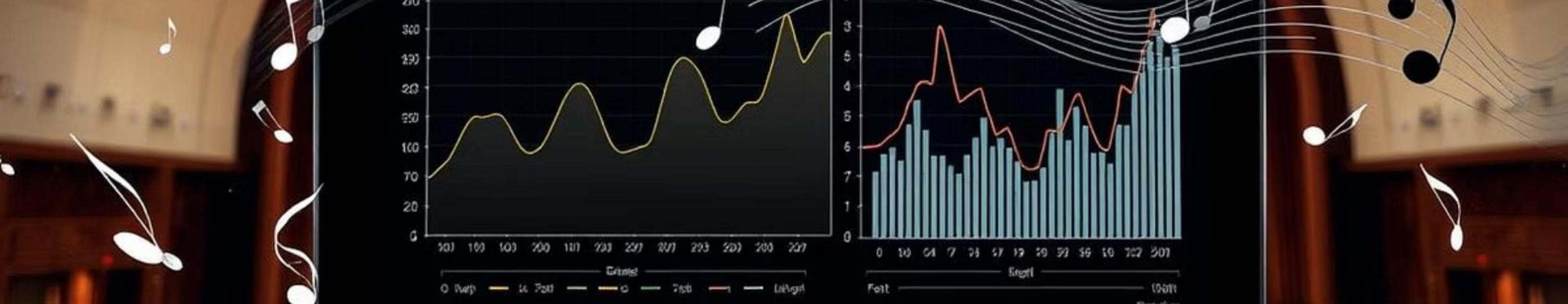
# Predicting Song Popularity with Machine Learning

Forecast song popularity pre-release to boost marketing and engagement. Leverage audio features to identify potential hits early.

# Project Objective & Importance

**Project Goal**

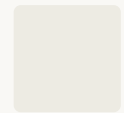Predict song popularity score via musical features

**Key Features**

- Danceability
- Energy
- Tempo

**Business Value**

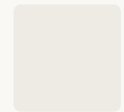Optimize marketing, identify hits early, enhance user engagement

# Proof of Concept Results

### Model Used

Random Forest Regressor to predict popularity

### Key Validation

Strong correlations between audio features and popularity

# Business Impact

♫ **Curate Trends**

Enhance playlist recommendations

◎ **Optimize Marketing**

Pre-launch promotional tactics

☆ **Increase Reach**

Engage larger listener bases

# Machine Learning Problem & Alternatives

## Problem Type

Regression: Predict continuous popularity score 0–100

## Alternative Approach

Classification considered: Hit vs Flop

Rejected due to loss of score granularity

Residials over 0.08

# Key Technical Metrics

### R² Score

Explained variance metric of prediction quality

### RMSE

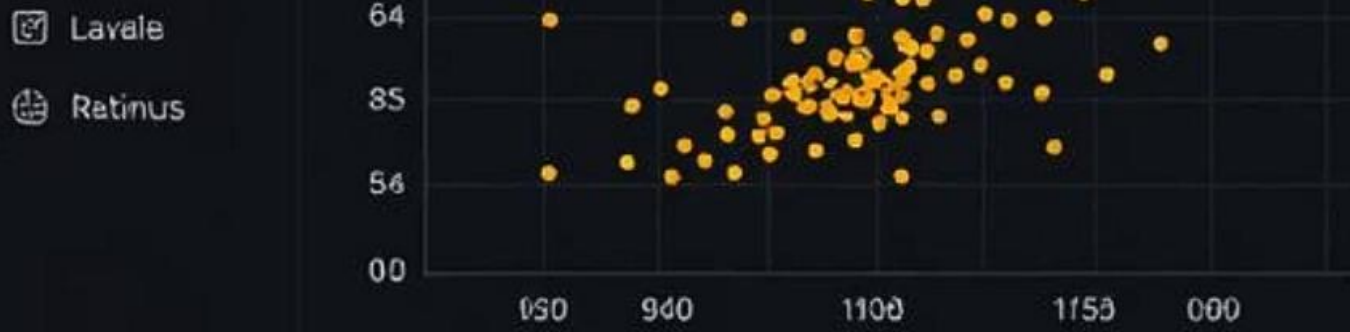Average prediction error, lower is better

# Model Architecture & Workflow

| 1 | **Step 1: Data Ingestion**<br><br>Import CSV dataset |
|---|---|
| 2 | **Step 2: Preprocessing**<br><br>Clean, normalize, filter data |
| 3 | **Step 3: Modeling**<br><br>Train Random Forest on features |
| 4 | **Step 4: Evaluation**<br><br>Assess metrics, visualize results |

# Data Sources & Preparation Steps

## Data Origin

Kaggle dataset with song metadata + audio features

## Preparation Tasks

- Missing value imputation

- Drop irrelevant columns

- Scaling and normalization

# Feature Engineering Highlights

| Feature Selection | Top Features | Outcome |
|---|---|---|
| Removed low variance and multicollinear features | • Energy<br>• Danceability<br>• Valence | Improved accuracy and interpretability |

# Hyperparameter Tuning & Tech Stack

## Hyperparameter Tuning

- GridSearchCV for n_estimators, max_depth
- Improved stability and reduced overfitting

## Tech Stack

- Python, Pandas, NumPy, Scikit-learn
- Django for web integration
- Jupyter Notebook, SQLite database

# Core Concepts in Music Popularity Prediction

Key techniques and metrics for modeling song popularity.

# Model Training & Evaluation

## Training Strategy

- 80-20 train-test split
- Random Forest on engineered features

## Evaluation Metrics

- High $R^2$ score
- Low RMSE
- Error distribution visualization

# Robust Testing Strategy

### Manual Prediction Testing

Test realistic input scenarios

### Django Unit Testing

Validate forms and backend logic

### Cross-Validation

Ensure statistical reliability

# Integration & Deployment Overview

## User Interface

Simple and intuitive song feature input

## Backend Integration

Random Forest model seamlessly integrated

## Real-Time Predictions

Instantaneous popularity output

# Project Requirements

## Languages & Libraries

- Python 3.x
- Pandas, NumPy, Scikit-learn

## Tools & Interface

- Jupyter Notebook
- SQLite database
- Django web app via browser

# Key Learnings & Tips

**Data Quality**

Critical for stable
predictions

**GridSearchCV**

Powerful yet
computationally expensive

**Visual Analytics**

Useful for feature selection

# Impact & Future Roadmap

**Current Use**

Music trend forecasting & curation

**Lyrics Analysis**

Apply NLP for deeper insights

**Real-Time Data**

Integrate streaming data sources

**Model Enhancements**

Incorporate deep learning architectures

1

2

3

4

# Trade-Offs & Decisions

## Model Choice

Random Forest preferred for interpretability

## Database

SQLite chosen for lightweight setup

## Framework

Django for admin & UI convenience

# File Structure Summary

- README.md

- dataset/MusicDataset.csv

- notebooks/Song_popularity_prediction

- models/scaler.pkl

- models/ridge.pkl

- application.py

- templates/home.html

- templates/index.html

- .ebextensions/python.config

- requirements.txt

- .vscode/settings.json, extensions.json, tasks.json

# STAR Story – Technical Challenge

**1** **Situation**

Inconsistent predictions despite similar inputs

**2** **Task**

Identify instability source and refine model

**3** **Action**

- Dropped irrelevant columns

- Handled multicollinearity

- Applied GridSearchCV tuning

**4** **Result**

- 20% RMSE reduction

- Stable, interpretable predictions

- Increased model trustworthiness

# Link – Code Repo, Dataset, PDF

[Predicting-Song-Popularity](Predicting-Song-Popularity)

# Thank You & Questions

### Appreciate your attention
Thank you for joining today

### Open Q&A
Any questions or feedback?

# Appendix

## Confusion Matrix

|        | Predicted 0 | Predicted 1 |
|--------|-------------|-------------|
| Actual 0 | 1759 | 4 |
| Actual 1 | 230 | 8 |

## Top 15 Feature Importances

| Feature | Importance Score |
|---------|------------------|
| artist_hotttnesss | ~0.096 |
| tempo | ~0.090 |
| artist_familiarity | ~0.088 |
| loudness | ~0.086 |
| duration | ~0.082 |
| start_of_fade_out | ~0.078 |
| mode_confidence | ~0.077 |
| key_confidence | ~0.076 |
| end_of_fade_in | ~0.064 |
| time_signature_confidence | ~0.060 |
| year | ~0.053 |
| key | ~0.046 |
| artist_longitude | ~0.040 |
| artist_latitude | ~0.037 |
| time_signature | ~0.021 |