

What is "Old Town Road"?



Arnav Deshwal | Pedro Ivo Rivas | Michael Sparkman He Wang | Rui Ying

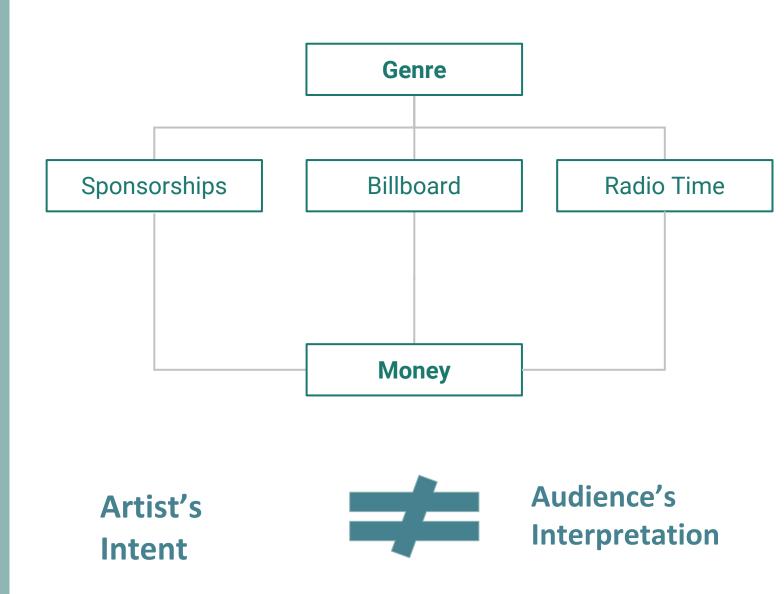


Content

- 1. Main Problem
- 2. Data Transformation
- 3. Descriptive Statistics
- 4. Model Selection
- 5. Results
- 6. Conclusion & Recommendations



Main Problem

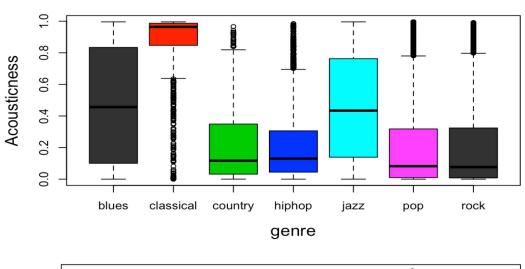


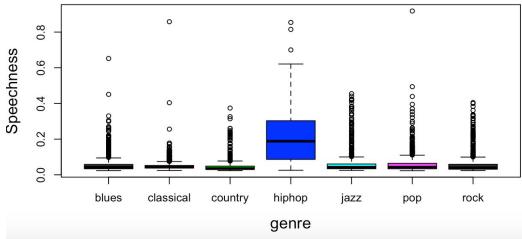
Data Transformation

Filtered Genres Cut Short Songs **135,000** Songs Cut songs less than 2 Chose 7 main genres and Dataset was found dropped sub-genres that containing 135,000 songs minutes to eliminate pulled from Spotify's API crossed between them (ie potential "dirty" data(ie songs that lasted one blues rock or pop rock) second) Final data set was

~**35,000** songs

Descriptive Statistics





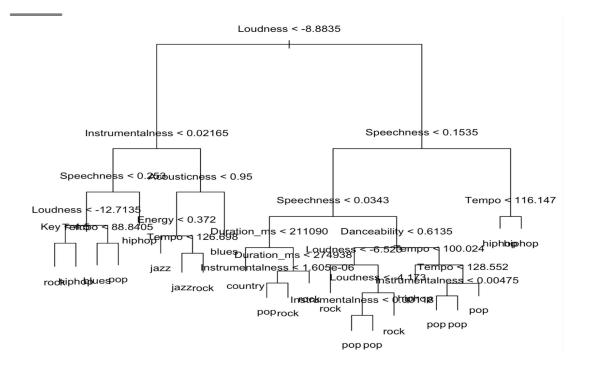


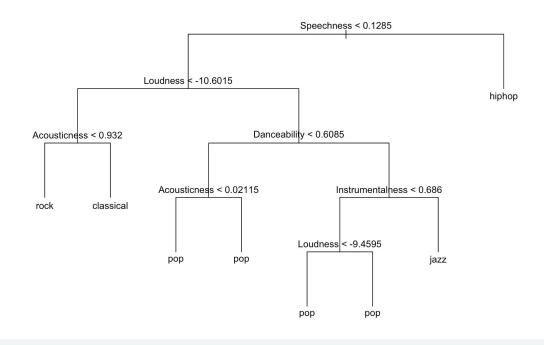
12 independent variables (e.g. loudness, key, energy, valence) ~ dependent variable: genre

02

Pop and rock are similar to each other; Classical, blues and jazz are very different from other categories.

Classification Tree





A Big Tree

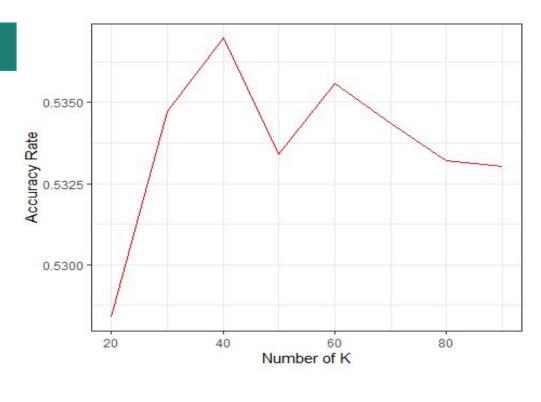
Predict using test dataset: Accuracy_bigtree = **45.57%**

After pruning it

Set the number of tree nodes = 8
Predict using test dataset:
Accuracy_prunetree = **48.06**%

K-Nearest Neighbor

Method 1	Method 2
Using all the 12 predictors Cross validation:	Variable Selection: ended up with 9 variables • Danceability
K = 63 Accuracy: 53.80%	 Energy Loudness Speechness Acouticness Instrumentalness Valence Liveness Duration_ms



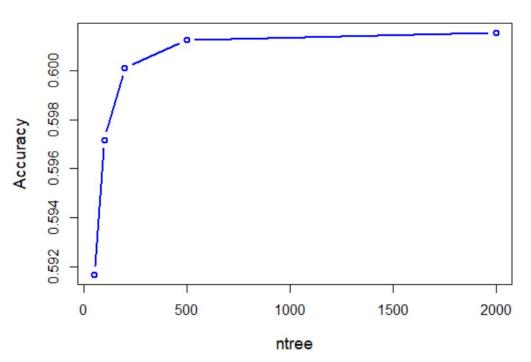
Accuracy:

K=43

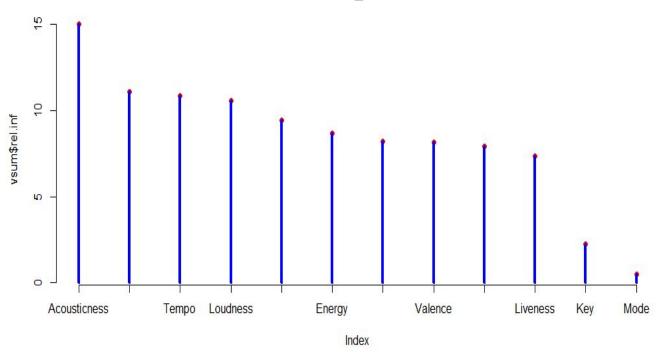
53.88%

Bagging and Boosting

Number of Trees vs. Accuracy



Variable Importance



Bagging:

Accuracy: **60.15%** with **2000** trees

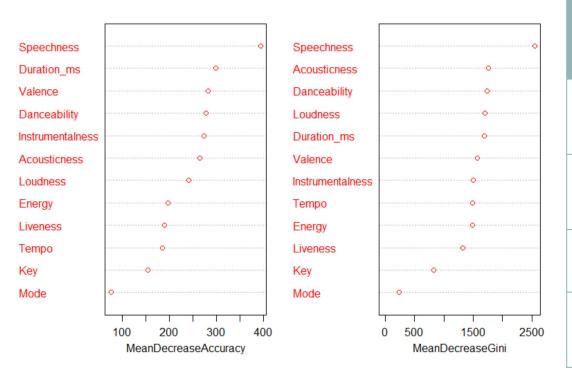
Boosting:

Accuracy: **56.52%** with 1000 trees

Acoustiness and Speechness shown to be most important; Key and Mode by far least important

Random Forest

Variable Importance

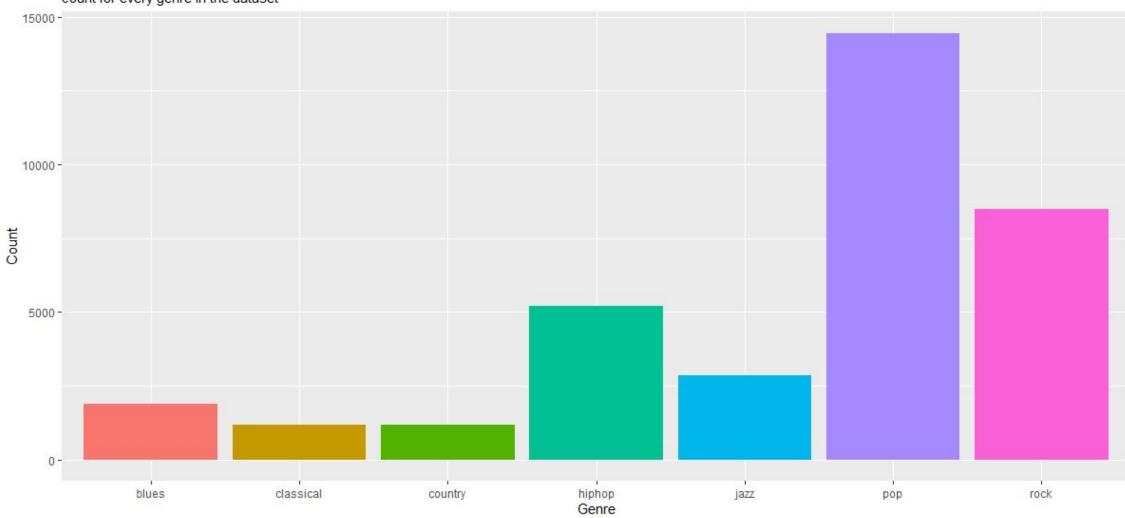


Number of Trees	m(Max Accuracy)	Accuracy
50	6	60.97%
200	6	61.30%
500	3	61.46%
2000	3	61.56%

Balancing Dataset

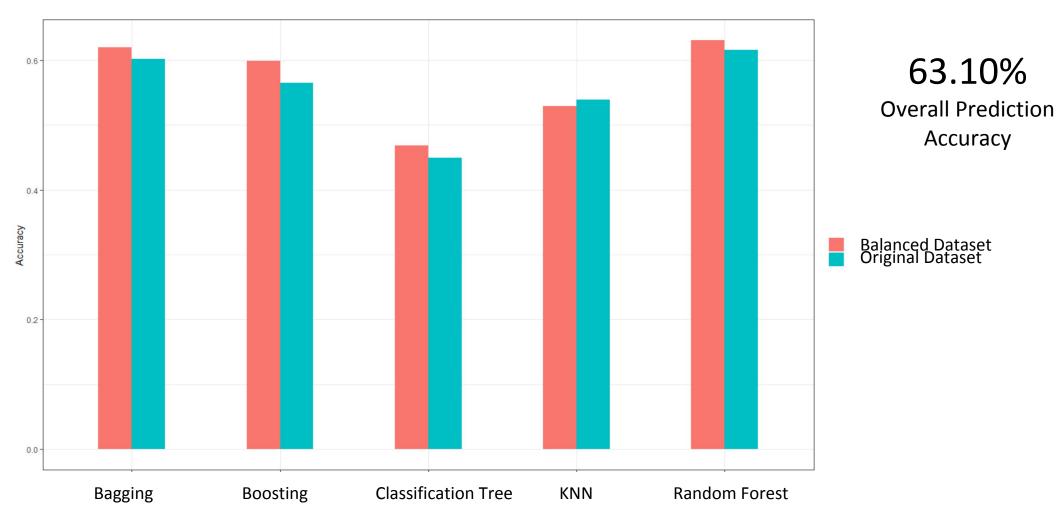
Genre Count

count for every genre in the dataset



Results

Accuracy(Balanced Dataset) vs. Accuracy(Original Dataset)

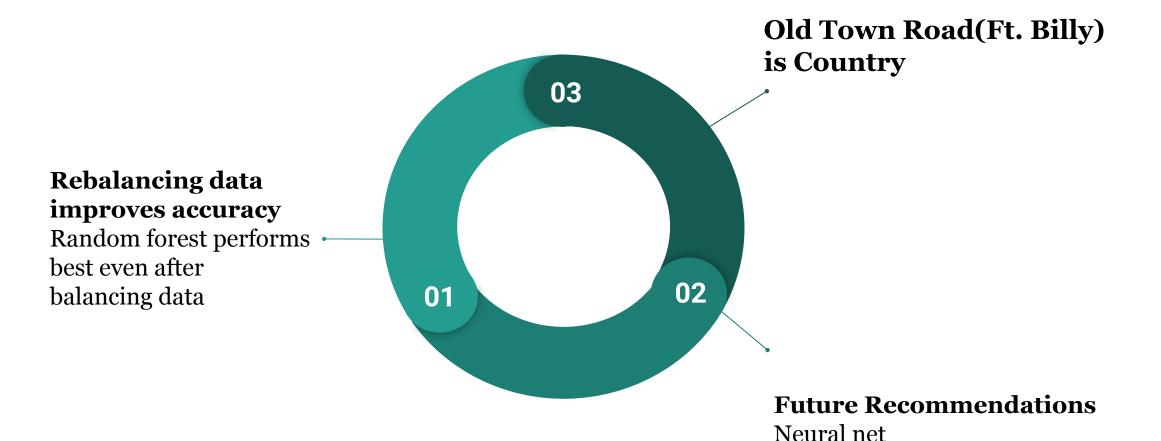


Prediction

Song	Actual Genre
Old Town Road	Hiphop
Old Town Road (Ft. Billy)	??
I want it that way	Pop
Metallica -Fade to black	Rock
Iron Maiden - Trooper	Rock

Predicted Genre		
Hiphop		
Country		
Pop		
Blues		
Pop		

Conclusion and Future Recommendations



Find a better way to rebalance data

Add more genres including cross-genre

(e.g. pop-rock)



Thank You!

Q&A