# SPOTIFY SONG ANALYSIS

A CLASSIFICATION PROBLEM OF SONG POPULARITY

- SULAIMON OYELEYE

# OUTLINE

❖ Motivation

❖ About the dataset

❖ Insights and Visualizations

❖ Model used

❖ Results & Comparison

❖ Model Evaluation on generated sample data

❖ Conclusion

# MOTIVATION

XYZ Music Production company understands that making popular song is very important for its business to stay competitive.

Hence, they have contacted the NLB Team to help in building a model to classify what kind of songs will be popular and figure out a way to make help the business stay competitive.

# ABOUT THE DATASET

➢ Web scrap using spotipy api.

➢ Datasets gotten in 2 parts: the track IDs and audio features. I combined both to make a more complete dataset.

➢ Data cleaning and duplicates removal was done.

➢ Saved to HDFS and Local Disk for future usage.

# INSIGHTS AND VISUALIZATIONS

## MOST POPULAR SONGS

| artist_name | popularity | track_name |
|---|---|---|
| 24kGoldn | 95 | Mood (feat. iann dior) |
| The Weeknd | 95 | Blinding Lights |
| KAROL G | 94 | BICHOTA |
| Bad Bunny | 94 | DÁKITI |
| Ariana Grande | 94 | positions |
| Bad Bunny | 93 | LA NOCHE DE ANOCHE |
| The Kid LAROI | 93 | WITHOUT YOU |
| CJ | 93 | Whoopty |
| Billie Eilish | 92 | Therefore I Am |
| SZA | 92 | Good Days |

## SONGS MOST SUITED TO DANCING

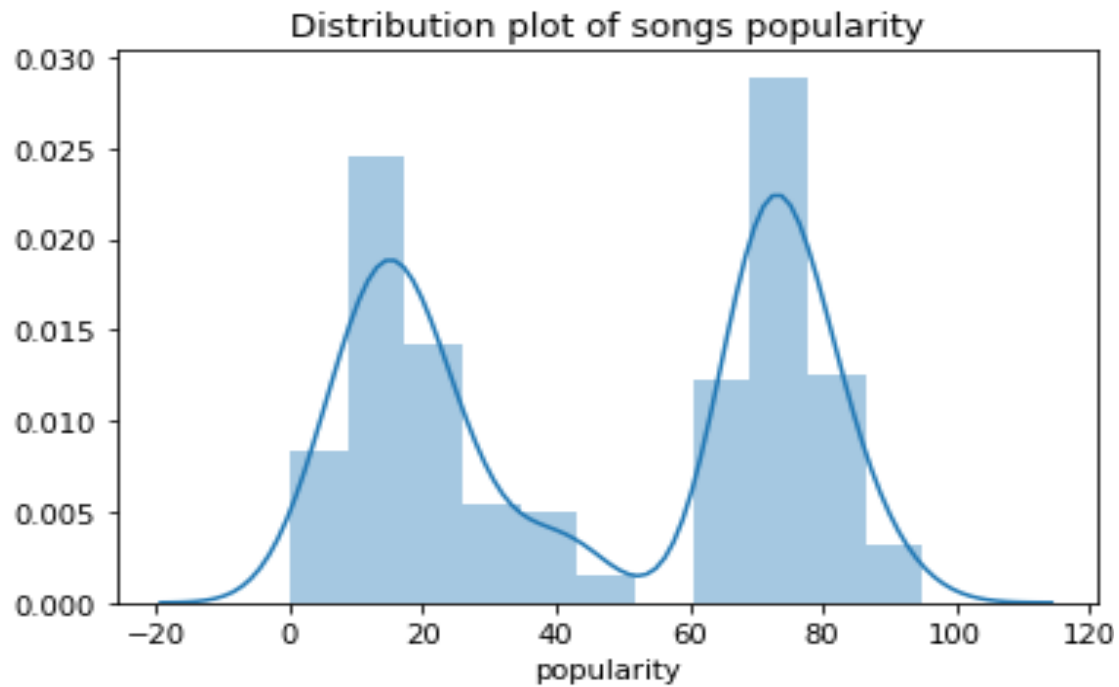| artist_name | track_name | danceability |
|---|---|---|
| 347aidan | Dancing in My Room | 0.980 |
| Championxiii | BOO! | 0.974 |
| Black Eyed Peas | GIRL LIKE ME | 0.965 |
| Young T & Bugsey | Don't Rush (feat. DaBaby) | 0.959 |
| Boosie Badazz | Mop Wit It | 0.958 |
| Kenndog | Drip Like ME | 0.956 |
| Erica Banks | Buss It | 0.956 |
| Sada Baby | Whole Lotta Choppas (Remix) [feat. Nicki Minaj] | 0.956 |
| T.I. | Pardon (feat. Lil Baby) | 0.955 |
| Saweetie | Tap In | 0.954 |

## MOST CHEERFUL SONGS

| artist_name | track_name | valence |
|---|---|---|
| Simon Patterson | Close My Eyes (Mixed) | 0.978 |
| Workout Music | Memories (Remix) | 0.977 |
| Greg Sletteland | I Have a Dream (Deep House Dance Party Remix) | 0.971 |
| Armin van Buuren | A State Of Trance (ASOT 991) - ASOT Tune Of Th... | 0.969 |
| Shawn Mendes | There's Nothing Holdin' Me Back | 0.968 |
| The Cog is Dead | Let Me Be Your Man (Remastered 2020) | 0.968 |
| DaBaby | JUMP (feat. YoungBoy Never Broke Again) | 0.966 |
| Workout Music | RITMO (Bay Boys For Life) [Remix] | 0.965 |
| Camilo | BEBÉ | 0.965 |
| Lele Pons | Se Te Nota (with Guaynaa) | 0.963 |

## MOST DEPRESSING SONGS

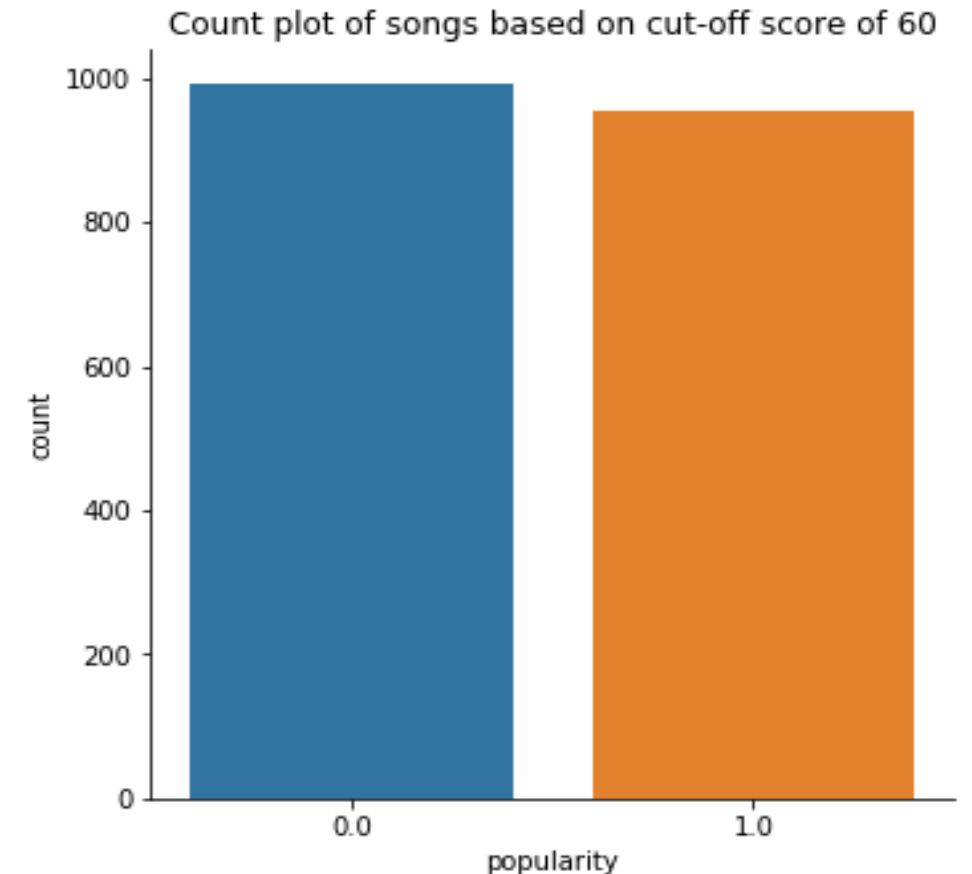| artist_name | track_name | valence |
|---|---|---|
| Water Sound Natural White Noise | Deep Sleep Recovery Noise | 0.00000 |
| No Spirit | It's Been A Year | 0.00000 |
| Epic Soundscapes | Heavy Rain | 0.00001 |
| Creatress | Steady Forest Rain | 0.00001 |
| Om Zone | Reaching Zen (Tibetan Bells & Soft Rain) | 0.02790 |
| Hammock | Longest Year - 2020 | 0.03120 |
| White Noise Baby Sleep Music | Feel the Rhythm of and Yet More Pink Soft Comp... | 0.03220 |
| Jai Wolf | Indian Summer - 2020 Encore Mix | 0.03260 |
| Arkham Knights | Closing In (Year in Review 2020) | 0.03280 |
| Above & Beyond | Surrender (Year in Review 2020) - Genix Remix | 0.03310 |

Correlation plot of features

Shapiro-Wilk Normality Test confirms that popularity is not a Gaussian distribution.

Linear regression won't be used.

# MODEL USED

➢ Classification Model

○ Logistic Regression Vs Decision Tree

➢ Three Categorical columns: Key, mode and time_signature
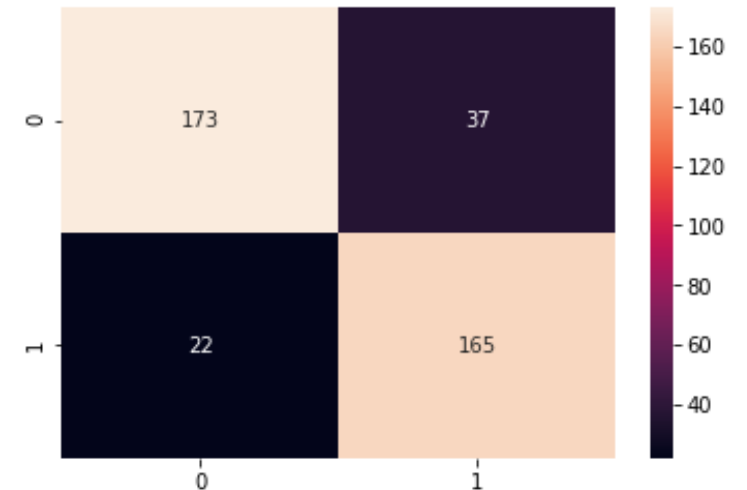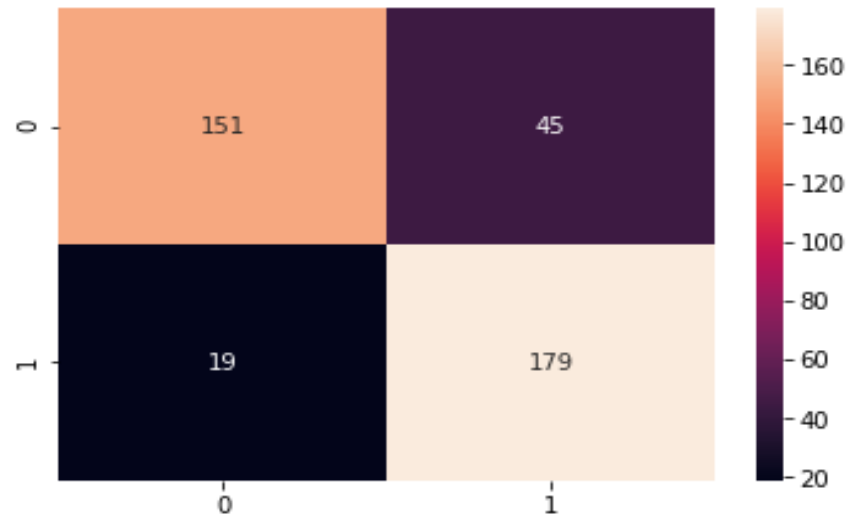
➢ Feature's engineering and scaling


Count plot of songs based on cut-off score of 60

# RESULTS & COMPARISON

| MEASURE | LOGISTIC REGRESSION | DECISION TREE |
|---|---|---|
| Accuracy (Train, Test) | (86.8 %, 83.8 %) | (90%, 85.1 %) |
| Precision (Not Pop., Pop) | (0.89, 0.80) | (0.89, 0.82) |
| Recall (Not Pop., Pop) | (0.77, 0.9) | (0.82, 0.88) |
| F1-Score (Not Pop., Pop) | (0.83, 0.85) | (0.85, 0.85) |
| Confusion Matrix | | |

# MODEL EVALUATION ON GENERATED SAMPLE DATA

## LOGISTIC REGRESSION

## DECISION TREE

| probability | prediction |
|---|---|
| [0.663617749194562,0.336382250805438] | 0.0 |
| [0.9915718467369371,0.008428153263062975] | 0.0 |
| [0.23819355954507418,0.761806440549258] | 1.0 |
| [0.9996709015962023,3.290984037976616E-4] | 0.0 |
| [0.605710619867777,0.3945289301322236] | 0.0 |
| [0.970576046716037,0.029423953528396215] | 0.0 |
| [0.041834858428763014,0.95816141571237] | 1.0 |
| [0.9986920818030733,0.001307918196926535] | 0.0 |
| [0.03712617064289648,0.962873829571035] | 1.0 |
| [0.8070134412944134,0.1929865587055656] | 0.0 |

| probability | prediction |
|---|---|
| [0.8095238095238095,0.19047619047619047] | 0.0 |
| [0.9747899159663865,0.025210084033613446] | 0.0 |
| [0.13071895424836602,0.86928104575163] | 1.0 |
| [0.9747899159663865,0.025210084033613446] | 0.0 |
| [1.0,0.0] | 0.0 |
| [0.32075471698113206,0.6792452830188679] | 1.0 |
| [0.32075471698113206,0.6792452830188679] | 1.0 |
| [0.9747899159663865,0.025210084033613446] | 0.0 |
| [0.13071895424836602,0.86928104575163] | 1.0 |
| [0.32075471698113206,0.6792452830188679] | 1.0 |

# CONCLUSION

➢ Model can be improved upon on both cases in the future

➢ Decision Tree is a better model for the most part

➢ Popular songs predicted were found to have high danceability value, high energy value, low speechiness value and above 0.5 loudness value.

# THANKS