# Do You Feel The Music?

Kevin Kang, Raymond Sutanto
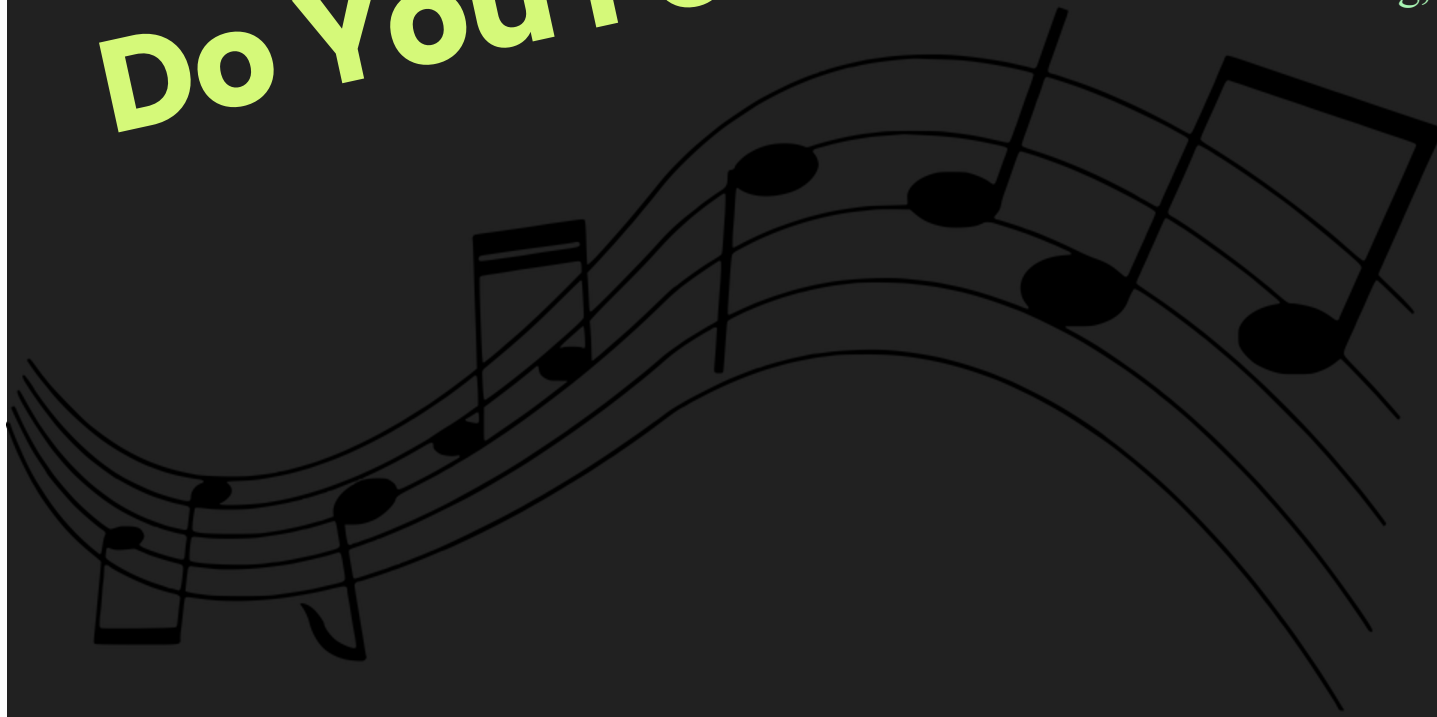
# Table of contents

**01**

**Motivation**

**02**

**Our Dataset**

**03**

**Cleaning & Exploration**

**04**

**Analysis**
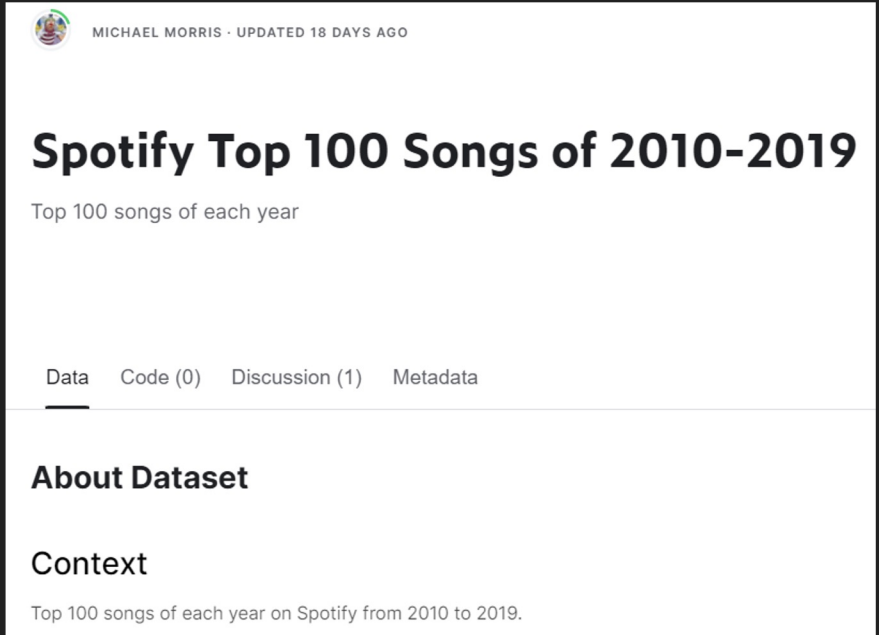
**05**

**Recommendations**

# Motivation

- Correlation or Causation?
- Finding The Formula
  - *Analysis*
  - *How variables affect the level of valence*
- "Aiding an algorithm"
  - *Accuracy*
  - *Personalization*

# Our Dataset

- Spotify Top 100 Songs:
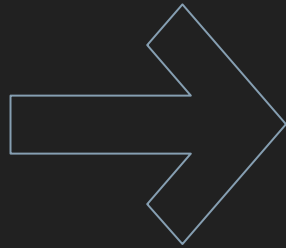  - *Timeframe: 2010-2019*

- Dataset URL (from KAGGLE):

*https://www.kaggle.com/datasets/muhmores/*
*spotify-top-100-songs-of-20152019*

kaggle

MICHAEL MORRIS · UPDATED 18 DAYS AGO

## Spotify Top 100 Songs of 2010-2019

Top 100 songs of each year

Data   Code (0)   Discussion (1)   Metadata

**About Dataset**

Context

Top 100 songs of each year on Spotify from 2010 to 2019.

# Cleaning (Overview)

**1000 records (100 songs x 10 years)**

*17* Variables

- Song_Name
- Artist_Name
- Song_Genre
- Year_Released
- Date_Added
- BPM
- Energy
- Danceability
- Decibel
- Liveness
- Valence
- Duration
- Acousticness
- Speechness
- Popularity_Score
- Top_Year
- Artist_Type

- Song_Genre
- BPM
- Energy
- Danceability
- Decibel
- Liveness
- Valence
- Duration
- Acousticness
- Speechness
- Popularity_Score
- Top_Year
- Artist_Type

*13* Variables

- Song_Genre
- Danceability
- Duration
- Acousticness
- Speechness
- Top_Year
- Popularity_Score

*7* Variables

# Cleaning

- Check & change necessary column data types to numerical
- Remove irrelevant variables
- Use python libraries in exploration & analysis

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 13 columns):
 #   Column            Non-Null Count   Dtype
---  ------            --------------   -----
 0   Popularity_Score  1000 non-null    int64
 1   Song_Genre        1000 non-null    string
 2   BPM               1000 non-null    int64
 3   Energy            1000 non-null    int64
 4   Danceability      1000 non-null    int64
 5   Decibel           1000 non-null    int64
 6   Liveness          1000 non-null    int64
 7   Valence           1000 non-null    int64
 8   Duration          1000 non-null    int64
 9   Acousticness      1000 non-null    int64
 10  Speechness        1000 non-null    int64
 11  Top_Year          1000 non-null    int64
 12  Artist_Type       987 non-null     float64
dtypes: float64(1), int64(11), string(1)
memory usage: 101.7 KB
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1000 entries, 392 to 182
Data columns (total 7 columns):
 #   Column            Non-Null Count  Dtype
---  ------            --------------  -----
 0   Popularity_Score  1000 non-null   int64
 1   Song_Genre        1000 non-null   int32
 2   Danceability      1000 non-null   int64
 3   Duration          1000 non-null   int64
 4   Acousticness      1000 non-null   int64
 5   Speechness        1000 non-null   int64
 6   Top_Year          1000 non-null   int64
dtypes: int32(1), int64(6)
memory usage: 58.6 KB
```
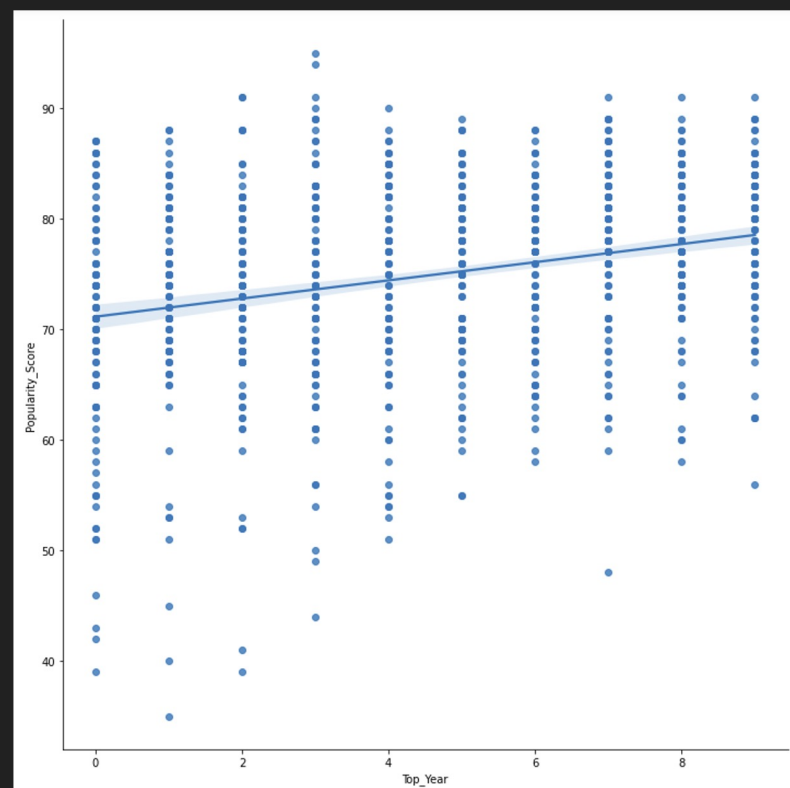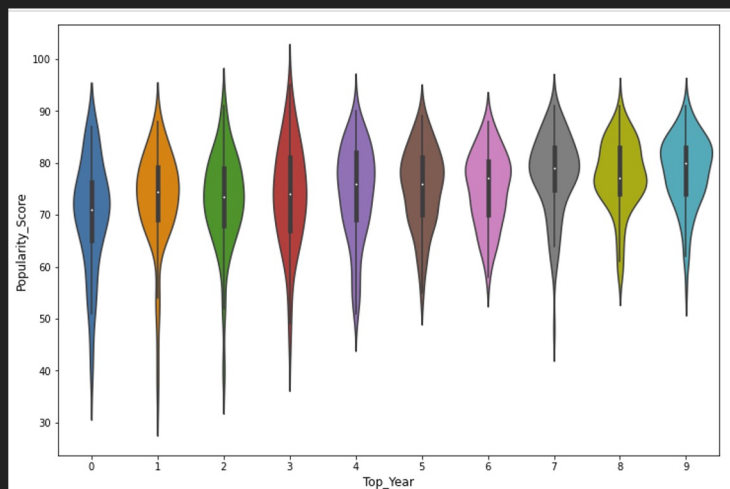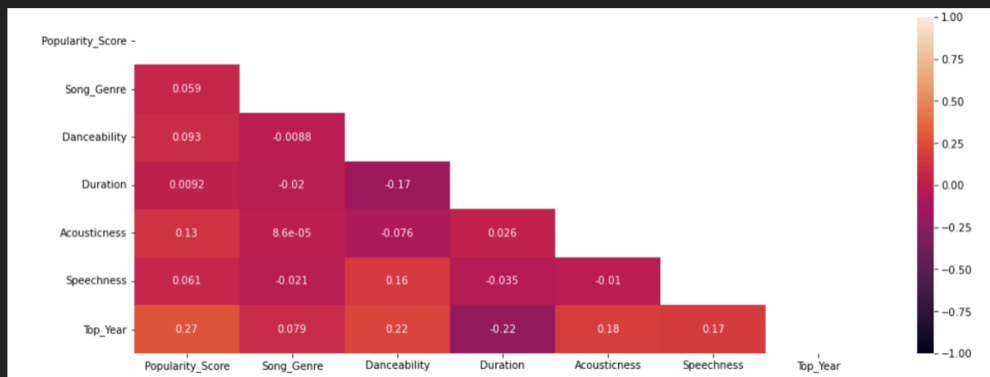
# Dataset Preparation (Cleaning)



| | Popularity_Score | Song_Genre | BPM | Energy | Danceability | Decibel | Liveness | Valence | Duration | Acousticness | Speechness | Top_Year | Artist_Type |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 70 | dance pop | 140 | 81 | 61 | -6 | 23 | 23 | 203 | 0 | 6 | 0 | 2.0 |
| 1 | 68 | dance pop | 138 | 89 | 68 | -4 | 36 | 83 | 192 | 1 | 8 | 0 | 2.0 |
| 2 | 72 | pop soul | 95 | 48 | 84 | -7 | 9 | 96 | 243 | 20 | 3 | 0 | 1.0 |
| 3 | 80 | atl hip hop | 93 | 87 | 66 | -4 | 4 | 38 | 180 | 11 | 12 | 0 | 1.0 |
| 4 | 79 | atl hip hop | 104 | 85 | 69 | -6 | 9 | 74 | 268 | 39 | 5 | 0 | 1.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 995 | 86 | hip hop | 155 | 73 | 83 | -4 | 12 | 45 | 313 | 1 | 22 | 9 | 1.0 |
| 996 | 85 | hip hop | 80 | 50 | 55 | -9 | 80 | 41 | 190 | 23 | 7 | 9 | 1.0 |
| 997 | 68 | grime | 103 | 77 | 89 | -5 | 9 | 46 | 177 | 1 | 7 | 9 | 1.0 |
| 998 | 67 | afroswing | 138 | 58 | 53 | -6 | 10 | 59 | 214 | 1 | 10 | 9 | 2.0 |
| 999 | 75 | atl hip hop | 98 | 59 | 80 | -7 | 13 | 18 | 200 | 2 | 15 | 9 | 1.0 |

Chart Pre-Cleaning

Heatmap Pre-Cleaning

# Dataset Exploration

# Dataset Analysis (Machine Learning)

## PRE

```python
#split into feature and target variables

X = spo_final.drop(['Popularity_Score'],1)
y = spo_final['Popularity_Score']

#split into train and test sets, 80:20 ratio
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size = 0.8, random_state = 0)
```

## POST

```python
#split into feature and target variables

Xx = spy.drop(['Popularity_Score'],1)
yy = spy['Popularity_Score']

#split into train and test sets, 80:20 ratio
Xx_train, Xx_test, yy_train, yy_test = train_test_split(Xx, yy, train_size = 0.8, random_state = 0)
```

# Dataset Analysis (Machine Learning)

| | PRE | POST |
|---|---|---|
| Linear Regression | Cross Validation Score:  [-17.10  -4.29 -4.37]<br>Average CVS: -14.58<br>Mean Squared Error:  76.07 | Cross Validation Score =  [-17.82  -3.36  -4.48]<br>Average CVS: -14.17<br>Mean Squared Error =  79.67 |
| Logistic Regression | Cross Validation Score: [0.057 0.063 039]<br>Average CVS: 0.046<br>Mean Squared Error =  108.47<br>Accuracy =  0.065 | Cross Validation Score =  [0.054 0.06 0.05]<br>Average CVS: 0.056<br>Mean Squared Error =  92.19<br>Accuracy =  0.085 |
| Decision Tree | Cross Validation Score =  [0.056 0.06 0.069]<br>Average CVS: 0.072<br>Mean Squared Error =  129.155<br>Accuracy =  0.045 | Cross Validation Score =  [0.042 0.057 0.03]<br>Average CVS: 0.068<br>Mean Squared Error =  136.73<br>Accuracy =  0.065 |

# Dataset Analysis (Machine Learning)

|  | PRE | POST |
|---|---|---|
| KNN | Cross Validation Score = [0.012 0.045 0.039]<br>Average CVS: 0.036<br>Mean Squared Error = 201.07<br>Accuracy = 0.075 | Cross Validation Score = [0.045 0.06 0.036<br>Average CVS: 0.039<br>Mean Squared Error = 239.62<br>Accuracy = 0.03 |
| SVC | Cross Validation Score = [0.066 0.063 0.079]<br>Average CVS: 0.069<br>Mean Squared Error = 98.24<br>Accuracy = 0.05 | Cross Validation Score = [0.069 0.067 0.075]<br>Average CVS: 0.064<br>Mean Squared Error = 98.24<br>Accuracy = 0.05 |
| Neural Network | Cross Validation Score = [0.045 0.039 0.072]<br>Average CVS: 0.057<br>Mean Squared Error = 106.801<br>Accuracy = 0.055 | Cross Validation Score = [0.06 0.039 0.045]<br>Average CVS: 0.053<br>Mean Squared Error = 124.61<br>Accuracy = 0.06 |

# Dataset Analysis (Machine Learning)

|  | PRE | POST |
|---|---|---|
| Logistic Regression | Cross Validation Score: [0.057 0.063 039]<br>Average CVS: 0.046<br>Mean Squared Error = 108.47<br>Accuracy = 0.065 | Cross Validation Score = [0.054 0.06 0.05]<br>Average CVS: 0.056<br>Mean Squared Error = 92.19<br>Accuracy = 0.085 |
| SVC | Cross Validation Score = [0.066 0.063 0.079]<br>Average CVS: 0.069<br>Mean Squared Error = 98.24<br>Accuracy = 0.05 | Cross Validation Score = [0.069 0.067 0.075]<br>Average CVS: 0.064<br>Mean Squared Error = 98.24<br>Accuracy = 0.05 |

# Conclusion

- Recommended model after dropping variables: Logistic regression

- What are the important variables in creating a song?
  - Song Genre, Danceability, Duration, Acousticness, Speechness, Top Year

- Why are our indicators have low numbers?
  - Dropped variables
  - Other qualitative variables

- Business Recommendation:
  - A song producer should consider the above variables in creating a popular song