

# 統計, 機械学習, AIを用いた楽曲のヒット予測

兵庫県立姫路西高等学校 木村優介 吉田隼輔

## 1. 研究動機・目的

毎年, 3000以上の曲がリリースされているが, その中でヒットする曲には法則があると考えた。ヒット曲の予測に関しては様々な先行研究がある。本研究では, 音楽の重要な要素である「コード進行」だけで, ヒット曲の予測モデルを構築する。AIが機械学習し, より精度の高いヒット曲予測ができる仕組みを作成することを目的とした。

## 2. 研究手法

【使用データ】 Billboard Japan Year End Hot 100

- ・2010~2019年の各年上位20曲の全200曲
- ・サビの部分だけの「コード進行」のデータを収集した。
- 全コード数: 約4000(主音を長調はC, 短調はAに統一)

【分析手法】

### 第一分析

「コード」から曲の複雑性をデータ化

活用した統計手法

主成分分析

### 第二分析

「コード」の組み合わせをデータ化

Ngram解析・決定木分析

### 第三分析

「コード」から曲全体の特徴をデータ化

トピックモデル分析 (LDA)

3つのデータを説明変数, 順位を目的変数

ヒット曲(順位)を予測する回帰式の作成 (予測モデルの作成)

重回帰分析

モデル再学習

曲のコードを入力する

予想順位が出る!

コードと順位の入力を自動化  
1年毎に予測モデルの更新

## 3. 分析

### 第一分析 (曲の複雑性を分析・数値化)

手順①

- a サビのコード種類数
- b ダイアトニックコードを含む割合
- c サビの長さに対するコードの固まりの繰り返し割合の3つのデータを2次元化(主成分分析法の利用)

手順② 主成分得点をクラスタリングする。

- ・エルボー法によりクラスタ数を4に決定(図1)
- ・クラスタごとに分類(図2)

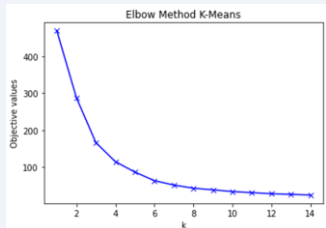


図1 エルボー法からクラスタ数の決定

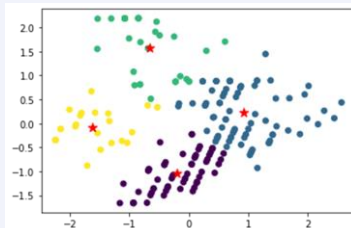


図2 4つのクラスタに分類

結果 (所属したクラスタの個数) ÷ (代表点との距離) を

第一分析によるデータ **A** とする。

### 第二分析 (コードの組の特徴を分析・数値化)

2010年~2019年の曲のサビコードにTri-gramを実施し, 3つ区切りのコードのtf-idf値を算出した。そのうち, 全データの上位100個を決定木分析し, 8グループに分類した。

各グループにはグループ内のtf-idf値平均値とTri-gramの個数の和の数値を与えることに設定する。

このグループに所属した場合  
tf-idf値0.55  
個数17

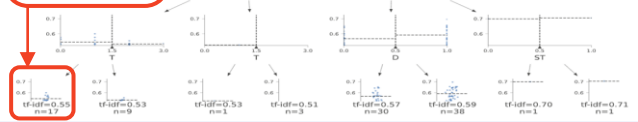


図3 pythonによる決定木分析結果

例 このグループに所属した場合

$$(\text{その曲tf-idf値})^2 \times (\text{tf-idf値}0.55 + \text{個数}17 + 2) = B_1$$

※その曲tf-idf値とは, 1曲あたり(コード数-2(個))現れるため, 特徴を表す大きい値から3つを採用する。

※「tf-idf値」+「個数」では「個数」の影響が大きいため標準化した。

※定数項2は, 全ての値を正にするために設定した。

### 結果

1曲につき  $B_1$  の値が3つ現れるため, 3つの値の積を **B** とする。

### 第三分析 (コード進行全体の特徴を分析・数値化)

コード進行の特徴データ化するため, 自然言語処理し, グループ分けをする。

手順① 各曲のサビコードから7割以上出現・10個以下しか出現しないコードを削除した。

手順② トピック数は評価指標である

PerplexityとCoherencelによってpythonの処理から7グループに分けることを決定した(右図)。

手順③ LDAを実行した。

結果: トピック(0~6)ごとに現れるコードの割合が示される。



トピック0 図4 LDAの結果の可視化

※例えば, トピック0は, G7が0.134, A#が0.113, Cm7が0.112, D#が0.093の割合で出現するグループと判断できる。

結果  $\sum_{k=0 \rightarrow 6} P_k \times S_k = C$  とする。

( $P_k$ : 演算する曲がトピックkに所属する確率)

( $S_k$ : 分析データがトピックkに所属する確率の総和)

## 4. 研究結果

2020年のランキング上位20曲のサビの部分のコードを用いて, 各分析結果から  $A \cdot B \cdot C$  の値を算出する。

説明変数:  $A \cdot B \cdot C$  の値 目的変数: 各曲の順位として, 重回帰分析法により回帰式を作成した。

回帰式 (重決定係数  $R^2=0.32$ )

$$y = -0.04 \times A - 3.25 \times B - 0.15 \times C + 21.0$$

|     | A     | B     | C     |
|-----|-------|-------|-------|
| p 値 | 0.096 | 0.235 | 0.183 |

## 5. 考察・今後の展望

予測モデル構築のためのフローは完成した。しかしながら, 有意水準10%としても, まだよりよいモデルの改善が見込まれる。特に, 第二分析が不十分である。

今後, 各年の上位20曲と, それ以外の曲との比較を行い, コードの情報だけでなく, AIが機械学習し精度の高いヒット曲を予測するシステムを構築する。人間が感覚で意思決定している要素を曲のデータから見出していく。