# Predicting the success factors of a song

Sophie De Becker
Yao Di
Mattia Gallese
Giacomo Martiriggiano

11th May, 2020

# PRESENTATION OVERVIEW

| 1 | Business Context |

| 2 | Data Collection and Exploration |

| 3 | Classification Problem |

| 4 | Regression Approach |

| 5 | Conclusions |

# The Global Music Industry
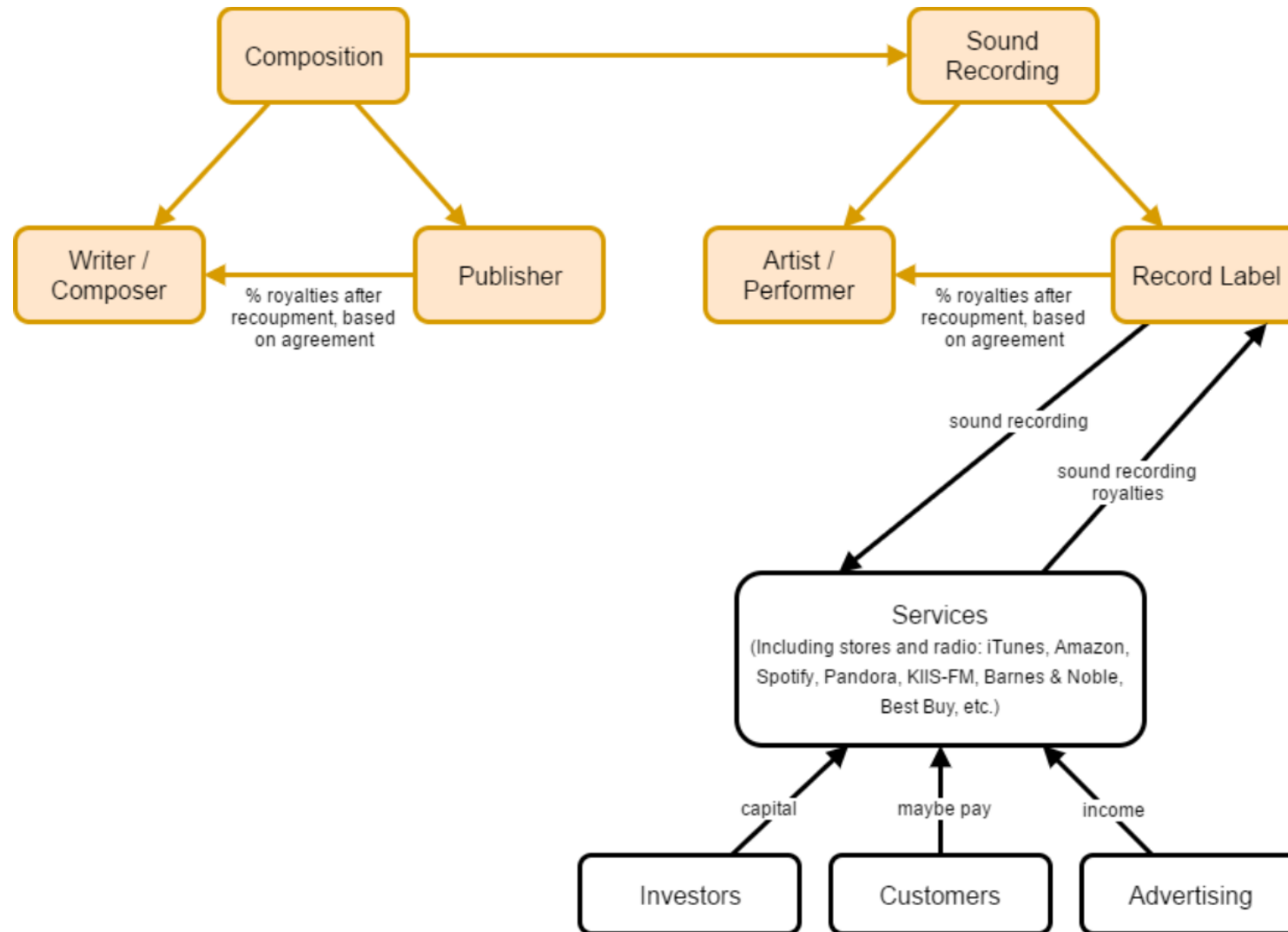
19.1 B USD in 2018
Highest Value Ever

59% of Value from
Digital (11.2 B USD)

Large Social Impact

# Stakeholder's Interactions



© 2017 TheMusicMaze.com

# Data Collection & Error Handling

# Success Definition

| Distinct songs | Best Position Ever Reached |
|:---:|:---:|
| Top | 1-20 |
| Not Top | |

# Exploratory Data Analysis – Continuous Variables

# Exploratory Data Analysis – Discrete Variables

Key has an **impact**



boxplot of best position grouped by key

Mode has
**no significant** impact



boxplot of best position grouped by mode

# Classification Problem Process Flow

# Classification Problem Results



Random Forest

**Final model selection**

SVM

Logistic Regression

F1 score:  0.097

F1 score: 0.289

F1 score: 0.276

# Success Prediction - Another Approach

Classification Model: no satisfactory performance

# Success Prediction - Another Approach



New Approach: Regression for Classification

# Feature Engineering for Regression Model

Features for Classification

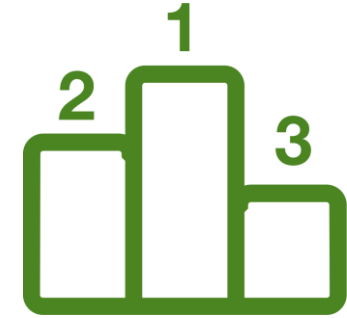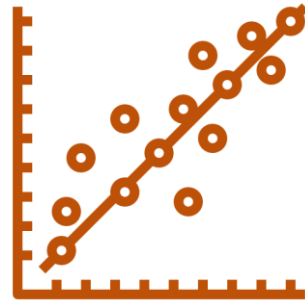Add Dynamic Features

Features for Regression

Dynamic features: impact of historical data

# Feature Engineering for Regression Model

Feature augmentation:
Linear model for non-linear regression

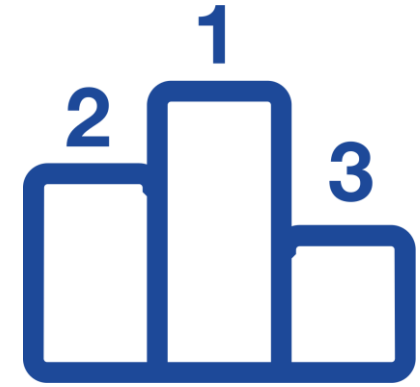# Result for Regression Model



| Metric Type | Metrics Score |
|-------------|---------------|
| Accuracy | 0.92 |
| Recall | 0.67 |
| Precision | 0.891 |
| F1 | 0.77 |

# Business Insights

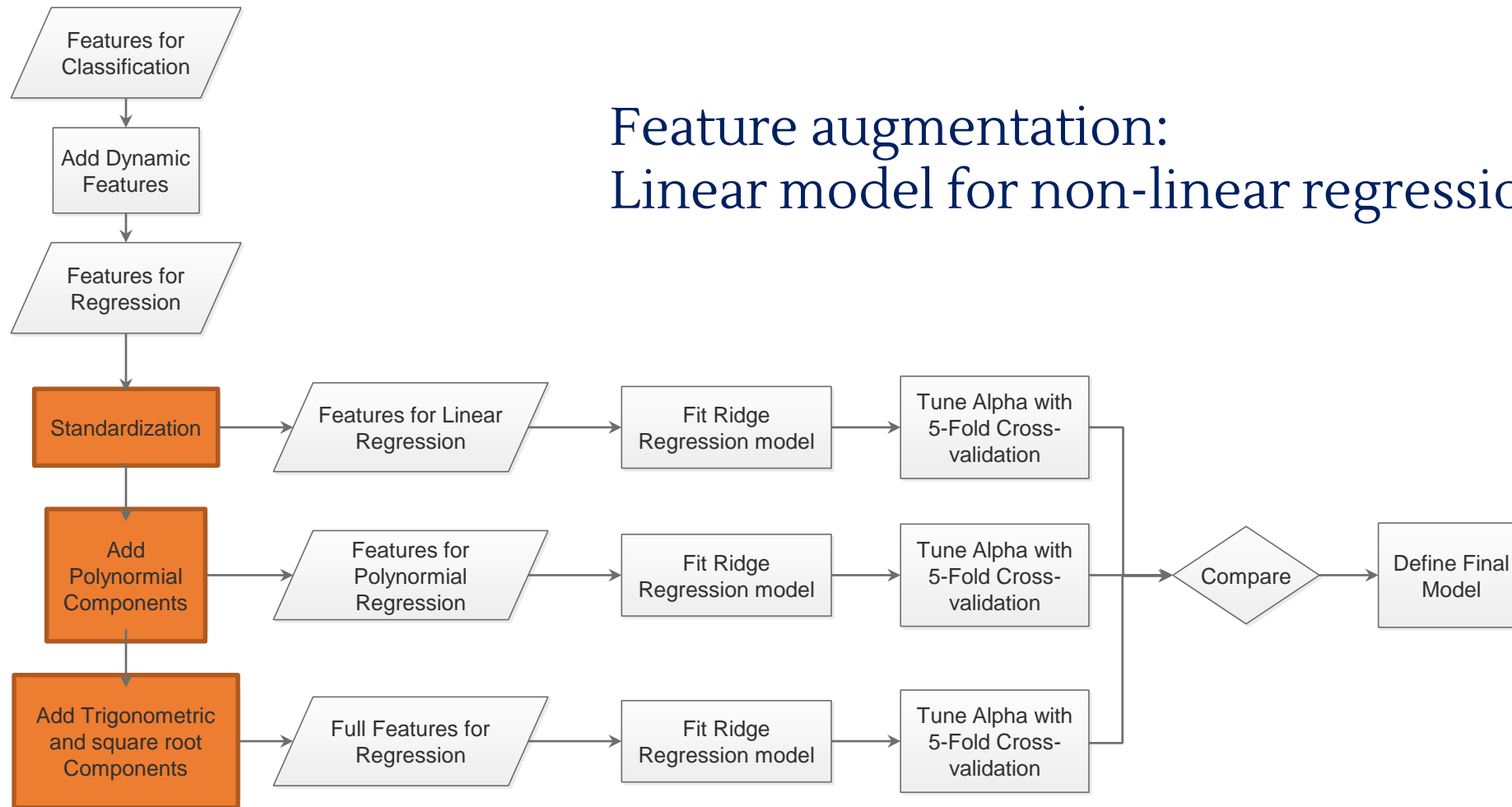| Inaccuracy Factors |
|---|

| No predictable effect | Time impact | Classification issues |
|---|---|---|

| Tv-show or movie could boost song streaming | Non-hit song could generate more money than a hit | Analysis on a small database |
|---|---|---|

# Conclusion

**1**    Difficult Industry to Model (As Expected)

**2**    Easier to Predict Evolution Based on Historical Data

**3**    There is Room for Improvement