

# Dokumentacja Skryptu Analizy Regresji Danych Spotify

Michał Herka, Daniel Brzezicki

## Wprowadzenie

Ten skrypt przeprowadza analizę regresji na danych Spotify, wykorzystując bibliotekę scikit-learn. Celem jest przewidywanie cechy 'streams' - ilości odtworzeń na podstawie wybranych cech wejściowych.

## Zależności

- pandas
- numpy
- scikit-learn (sklearn)
- Matplotlib

## Enumy

### RegressionMethod

- SINGLE: Analiza pojedynczej regresji
- TRAIN\_TEST\_SPLIT: Analiza regresji z podziałem na zestaw treningowy i testowy
- KFOLD: Analiza regresji z użyciem K-krotnego podziału krzyżowego

### RegressorType

- DECISION\_TREE: Drzewo decyzyjne
- RANDOM\_FOREST: Las losowy

## Funkcje

### `singleRegressorMAE(regressor, x_train, y_test)`

Dopasowuje dostarczony regresor do danych treningowych (`x_train`, `y_train`) i zwraca średni błąd bezwzględny (MAE) na danych testowych.

**trainTestSplitRegressorMAE(regressor, df\_features,  
df\_main\_feature, test\_size=0.5)**

Dzieli dane na zestaw treningowy i testowy, dopasowuje regresor i zwraca MAE na zestawie testowym.

**kfoldRegressorMAE(regressor, df\_features,  
df\_main\_feature, n\_splits=5)**

Przeprowadza K-krotny podział krzyżowy, dopasowuje regresor i zwraca maksymalne MAE we wszystkich foldach.

**getMaeValues(dataMethod: RegressionMethod,  
regressorType: RegressorType, X, y, depth=5)**

Generuje wartości MAE dla różnych głębokości regresora na podstawie określonej metody regresji i typu regresora.

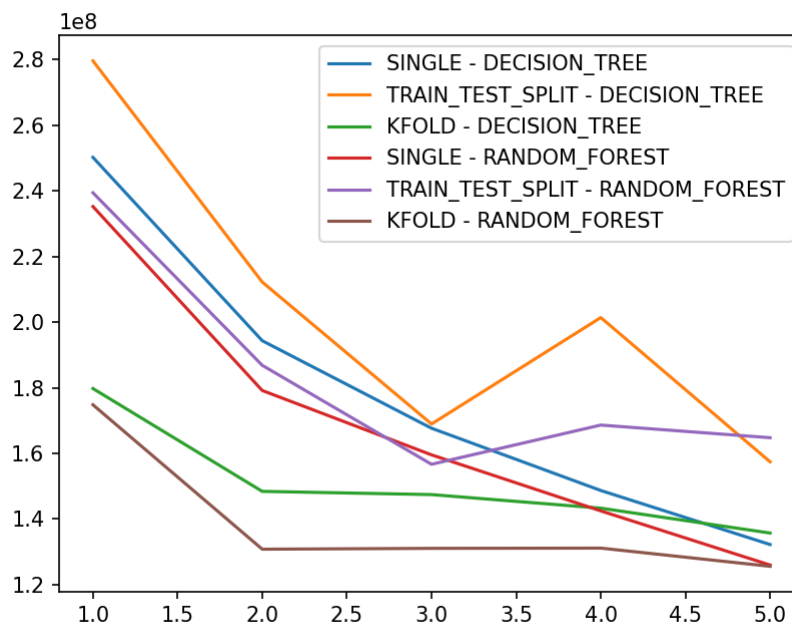
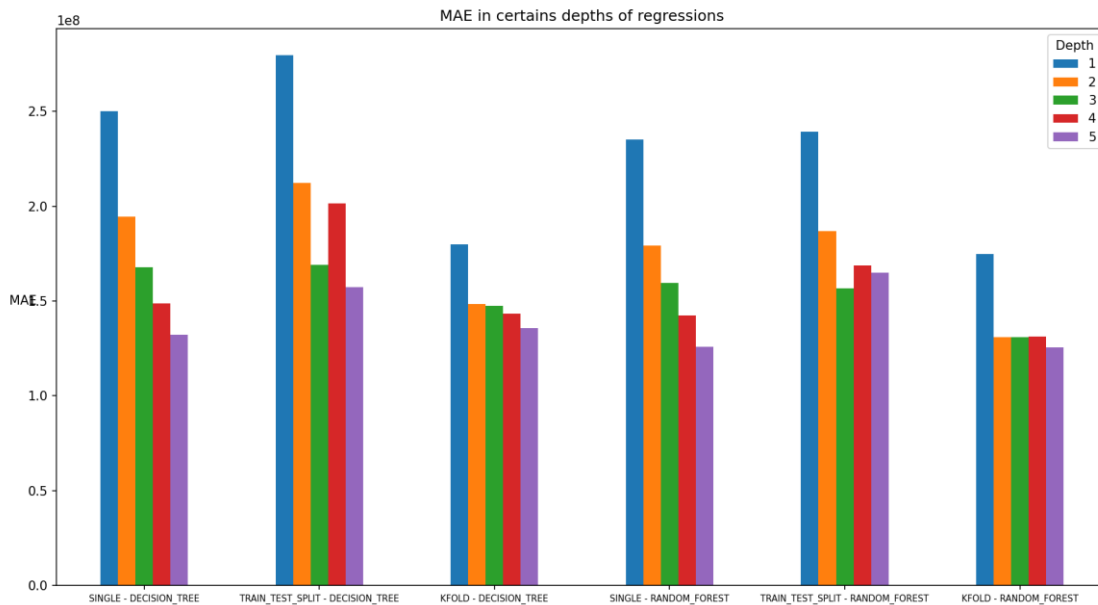
**predictValueByFeatures(dataFrame, features,  
mainFeature)**

Przetwarza dane, przeprowadza analizę regresji dla różnych metod i typów, a następnie przedstawia wyniki za pomocą wykresów słupkowych i liniowych.

## **Przykład Użycia**

```
result = predictValueByFeatures(pd.read_csv('data/spotify-2023.csv',  
encoding='latin-1'),  
['artist_count', 'released_year', 'in_apple_playlists', 'in_spotify_playlists', 'in_spotify_charts', 'danceability_%', 'streams'], 'streams')
```

## **Wyjście z dobraniem powyższych cech**



## Decision Tree

W przypadku pojedynczego drzewa decyzyjnego, obserwuje się tendencję do poprawy dokładności modelu wraz z zwiększaniem głębokości drzewa. Warto jednak zauważyć, że głębokość drzewa równa 5 wydaje się być optymalnym wyborem, gdyż dalsze zwiększanie głębokości nie przynosi już tak znaczącej poprawy. Wyniki uzyskane na zbiorze testowym (TRAIN\_TEST\_SPLIT) różnią się od wyników na zbiorze treningowym, co może sugerować pewne przetrenowanie modelu. Optymalna

głębokość drzewa dla tego zbioru wydaje się być mniejsza niż dla pojedynczego drzewa.

## Random Forest

Random Forest, będący ensemblem drzew decyzyjnych, wykazuje ogólnie lepszą zdolność do generalizacji niż pojedyncze drzewo. Otrzymane wyniki są bardziej stabilne i mniej podatne na przetrenowanie. W przypadku Random Forest również zauważa się tendencję do poprawy dokładności wraz z zwiększaniem głębokości drzewa, ale różnice między kolejnymi głębokościami są mniejsze niż dla pojedynczego drzewa.

## Porównanie między pojedynczym drzewem a Random Forest:

Random Forest uzyskuje niższe wartości błędów średnich bezwzględnych (MAE) niż pojedyncze drzewo decyzyjne dla wszystkich głębokości drzewa i dla wszystkich rodzajów podziałów danych (SINGLE, TRAIN\_TEST\_SPLIT, KFOLD). Wyniki na zbiorze KFOLD są najbardziej stabilne, co potwierdza, że Random Forest jest bardziej odporny na różnice między podziałami zbioru danych.

## Wnioski ogólne:

Random Forest wydaje się być bardziej wszechstronnym modelem do tego zadania, oferując lepszą zdolność do generalizacji niż pojedyncze drzewo. Optymalna głębokość drzewa może zależeć od specyfiki danych i zadania. Warto przeprowadzić dalsze eksperymenty, aby zoptymalizować parametry modelu. Przy analizie wyników ważne jest również zwrócenie uwagi na ewentualne przetrenowanie modelu, co można zaobserwować na zbiorze testowym w przypadku pojedynczego drzewa decyzyjnego. Regularyzacja modelu lub ograniczenie głębokości drzewa może pomóc w tym przypadku.

## Dane wyjściowe z przetrenowanego datasetu:

```
[10 rows x 24 columns]
Top 10 Predictions:
1. Title: Blinding Lights, Actual Streams: 3703895074.0, Predicted Streams: 3477042232.46
2. Title: Shape of You, Actual Streams: 3562543890.0, Predicted Streams: 3390391773.93
3. Title: Dance Monkey, Actual Streams: 2864791672.0, Predicted Streams: 2853848666.42
4. Title: Someone You Loved, Actual Streams: 2887241814.0, Predicted Streams: 2851084972.08
5. Title: One Dance, Actual Streams: 2713922350.0, Predicted Streams: 2845959256.40
6. Title: Sunflower - Spider-Man: Into the Spider-Verse, Actual Streams: 2808096550.0, Predicted Streams: 2798196210.70
7. Title: STAY (with Justin Bieber), Actual Streams: 2665343922.0, Predicted Streams: 2702468026.08
8. Title: Believer, Actual Streams: 2594040133.0, Predicted Streams: 2614473093.72
9. Title: Starboy, Actual Streams: 2565529693.0, Predicted Streams: 2611986729.07
10. Title: Closer, Actual Streams: 2591224264.0, Predicted Streams: 2605190308.44
PS C:\Users\dbenze\Desktop\uczenie_maszynowe>
```

#### Blinding Lights:

- Rzeczywiste odsłuchy: 3 703 895 074
- Przewidziane odsłuchy: 3 477 042 232
- Wnioski: Model prognozowania wydaje się skuteczny, ale przewidywania są nieco niższe niż rzeczywiste wartości, co może oznaczać pewne niedoszacowanie.

#### Shape of You:

- Rzeczywiste odsłuchy: 3 562 543 890
- Przewidziane odsłuchy: 3 390 391 774
- Wnioski: Model ponownie wydaje się być skuteczny, ale przewidywania są nieco niższe niż rzeczywiste wartości.

#### Dance Monkey:

- Rzeczywiste odsłuchy: 2 864 791 672
- Przewidziane odsłuchy: 2 853 848 666
- Wnioski: Model dobrze przewiduje odsłuchy, zbliżając się do rzeczywistych wartości.

#### Someone You Loved:

- Rzeczywiste odsłuchy: 2 887 241 814
- Przewidziane odsłuchy: 2 851 084 972
- Wnioski: Model prawdopodobnie dokładnie przewiduje odsłuchy, choć istnieje niewielka różnica.

#### One Dance:

- Rzeczywiste odsłuchy: 2 713 922 350
- Przewidziane odsłuchy: 2 845 959 256
- Wnioski: Model wydaje się nieco przeszacowywać odsłuchy.

#### Sunflower - Spider-Man: Into the Spider-Verse:

- Rzeczywiste odsłuchy: 2 808 096 550
- Przewidziane odsłuchy: 2 798 196 210
- Wnioski: Model jest skuteczny, a prognozy są zbliżone do rzeczywistych wartości.

#### STAY (with Justin Bieber):

- Rzeczywiste odsłuchy: 2 665 343 922
- Przewidziane odsłuchy: 2 702 468 026
- Wnioski: Model jest skuteczny, ale przewidywania są nieco wyższe niż rzeczywiste wartości.

#### Believer:

- Rzeczywiste odsłuchy: 2 594 040 133
- Przewidziane odsłuchy: 2 614 473 093
- Wnioski: Model jest skuteczny, a prognozy są zbliżone do rzeczywistych wartości.

Starboy:

- Rzeczywiste odsłuchy: 2 565 529 693
- Przewidziane odsłuchy: 2 611 986 729
- Wnioski: Model jest skuteczny, a prognozy są zbliżone do rzeczywistych wartości.

Closer:

- Rzeczywiste odsłuchy: 2 591 224 264
- Przewidziane odsłuchy: 2 605 190 308
- Wnioski: Model jest skuteczny, a prognozy są zbliżone do rzeczywistych wartości.

Podsumowanie:

Model jest skuteczny, ale istnieją pewne różnice między przewidywaniami a rzeczywistymi wartościami. Dla niektórych utworów prognozy są niższe, a dla innych wyższe. Może to wynikać z wielu czynników, takich jak brakujące informacje, które mogą wpływać na jakość modelu. Warto także zauważyć, że różnice te są stosunkowo niewielkie w porównaniu z ogólnymi liczbami odsłuchań.