

# LELEC2870 - Project 2022

Predicting movie's revenues with machine learning models

Rousseau Mathieu, 67001800



UCLouvain  
Belgium  
23/12/2022

## 1 Introduction

Starting with a dataset containing data from *IMDb*<sup>1</sup> as well as *BoxOfficeMojo*, the goal of this project was to predict a movie's revenue in the USA with the help of machine learning models.

## 2 Exploratory data analysis and preprocessing

The first part of this project was to perform some data analysis to understand better our dataset and then preprocess it in order to enhance prediction results.

Our dataset does not contain that much features at a first glance. However, the principal difficulty was to manage the so-called **embeddings** that consist, for every observation, in large vectors of features (represented as float numbers). With *img\_embeddings* feature being a vector of size 2048 and *text\_embeddings* a vector of size 768. Keeping all these features would give us a very high dimensional dataset with a number of features greater than the number of observations. This situation should be avoided as it is prone to overfitting<sup>2</sup>.

The first thing was to remove obvious unnecessary features like *title* which contains only distinct features or *img\_url* and *description* that we judged unnecessary as we already have the embeddings. We also noticed there were duplicated observations and removed them.

We then took a look at the missing values. There were only 4 values missing for the column *genre* consisting in only a tiny portion both datasets (*X1.csv* and *X2.csv*) therefore we simply dropped the observations with missing genres. However, there were respectively 264 and 106 values missing for the column *runtime* in *X1.csv* and *X2.csv*. These counting for more than 5% of our dataset, we did not take the risk of removing the observations containing these to avoid introducing a bias in our datasets. Looking at the dataset, it seemed like these values were missing randomly. We noticed that the shape of the distribution of the *runtime* feature is not too far from a Normal. As a consequence, we decided to impute the missing values with the mean of this feature.

Looking at the distributions of the other variables, we noticed that *n\_votes* and *revenues* were heavily right skewed. Because of that, simply removing outliers outside  $1.5 * IQR$ <sup>3</sup> could lead to removing a lot of datas. To fix the skewness, we took the log of these variables resulting in a much more homogenous distribution. We also noticed that the variable *is\_mature* only contained a unique value so we deleted this column. Then, we removed the outliers for the *runtime*, *production\_year* and *release\_year* features. Models like the linear regression and k-nearest neighbors are really sensible to outliers. Because the distribution of these features were more or less normally distributed, we removed the observations having a z-score<sup>4</sup> above the range  $[-3, 3]$ . We only did that for the training set so that the testing set stays representative of the reality.

The *genres* feature consists for every observation in a list of the maximum 3 genres that most accurately represent the given movie. Because there were not that much unique genres, we decided to **one-hot encode** them. However, the *studio* feature suffer of high-cardinality. Indeed, this feature consists in more than 300 unique studios so to avoid exploding the dimension of our dataset (and therefore avoid the curse of dimensionality), we replaced the less representative studios (the ones that represent less than 1% of the dataset) by a special name : "other".

To manage the embeddings, we tried to run a PCA and Kernel PCA (with a radial based function kernel) against them to keep the most significant features (those that explain the most the variance of the data points). The later has the advantage of recognizing non-linear pattern and gave slightly better results for our model. At first, we wanted to keep a certain amount of variance but we noticed

---

1. International Movie Database.

2. Situation where the model learned results of training data instead of generalizing and therefore cannot predict any new data.

3. Interquartile range.

4. Number of standard deviation with respect to the mean of the distribution.

we got better results for our model by being more drastic. We then chose to keep only the five most important principal components for each of the embeddings.

Finally, we standardized our dataset to put every features on the same scale and to have better results with the linear regression and multi-layer perceptron models (MLP). For example, in the MLP, it allows the gradient descent to converge more quickly. We ensured to only compute the transformation on the training set and applying it on the testing set to avoid data leakage. We end up with a dataset of 62 features.

### 3 Feature selection

Even if we ran a PCA on the embeddings, we still have a lot of features in our dataset and not enough data in order for our models to generalize well. We performed feature selection in order to decrease the number of variables in our dataset. We investigated multiple technique to select features : correlation, mutual information and recursive feature elimination. The problem with the first one is that it is only relevant for linear models since a low correlation does not mean the absence of relationship for non-linear models. We could have selected features highly correlated with the target for the linear regression but we found that even for this model we had slightly better results with mutual information. The last one on the other hand is a wrapper method and therefore use a model to compute a score and remove unnecessary features accordingly. The problem is that it is computationnaly intensive and very dependent on the chosen model.

We first tried to remove redundant features. The idea was to check the features that are highly correlated together (*above a certain threshold*) and then keep the one that had the most mutual information with the target. However we didn't find any redundant features in our dataset (at least after the PCA).

#### 3.1 Mutual information

We chose to use the mutual information to select the features. Mutual information has the advantage over correlation to detect nonlinear relationships between variables and therefore is a filter of choice for nonlinear models. Our strategy was to minimize the redundancy between inputs variables and a maximize relevancy between inputs variables and the target variable (*minimum redundancy, maximum relevance*). The idea is that even if 2 features are highly relevant, we shouldn't add both of them to our subset of features if they are highly correlated since it would increase the model complexity and could cause overfitting.

To minimize the redundancy, we computed a normalized mutual information matrix (fig. 5) and decided to remove for each group of redundant features, the feature that share the less information with the target. As you can see on the heatmap (fig. 5) in the appendix, *release\_year* shares a lot information with *production\_year*. *img\_feature1* shares also a significant amount of information with *text\_feature0*. Looking at the mutual information of each feature with the target variable (fig. 6), we decided to remove *production\_year* and *text\_feature0* since they share less information with the target.

### 4 Model selection

We tried severall models : linear regression, K-nearest neighbors (KNN), multi-layer perceptron and a random forest. For each model except the linear regression, we choosed to tweak severall well chosen hyperparameters and used a searching algorithm along with cross-validation (5 folds) to find the best hyperparameters for each model. We used the well known **Grid Search** despite having tested other searching algorithms like **Random Search** and **Bayesian Search**. These last two have lower computational cost due to the fact they make less iterations (the number of iterations is chosen by

the user) while finding a good model. However the results were not convincing so we decided to stick with Grid Search.

For each model, we run a 5-fold cross-validation to avoid overfitting and we compute the validation score as well as the test score that's been computed on a part of the  $X_1$  dataset that has not been used to train any of our models.

## 4.1 Linear regression

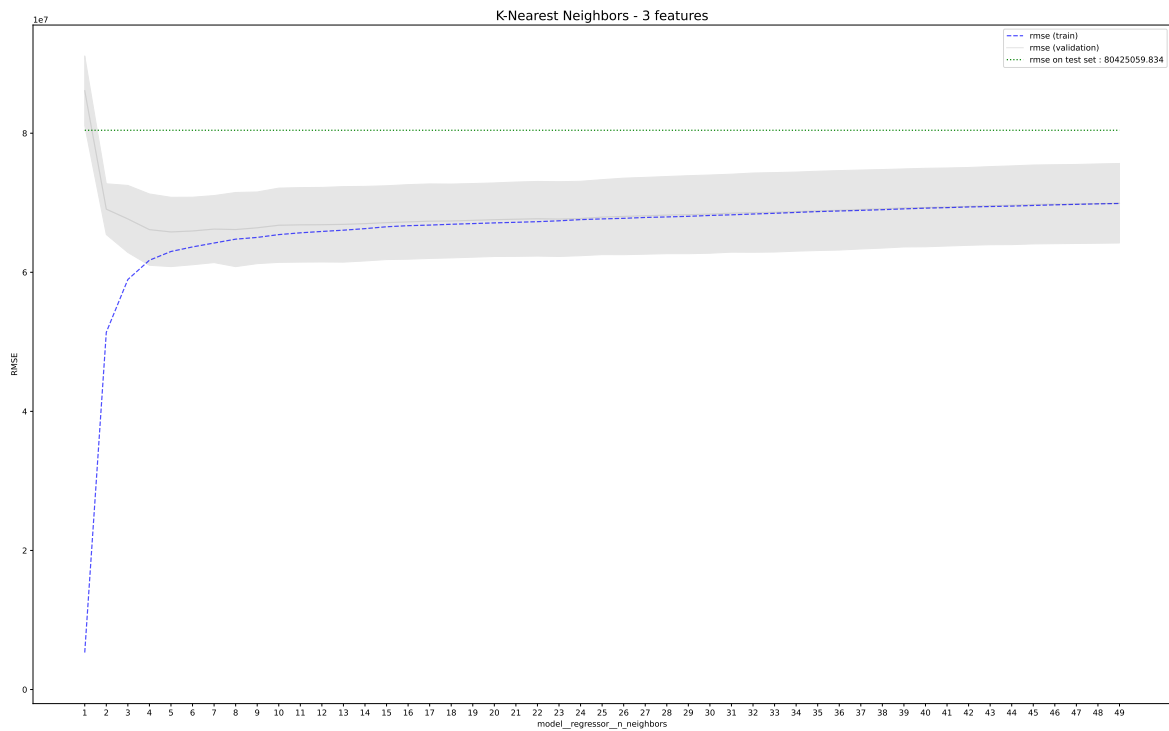
The linear regression has no hyperparameters to tune so we just ran a cross-validation for each set of features we decided to keep. The results show us that the best linear model is the one with the **30** best features selected.

#features	scores	
	validation rmse (\$)	test rmse (\$)
5	1.629e08	1.154e08
10	1.009e08	1.041e08
15	7.850e07	9.070e07
20	7.166e07	8.836e07
25	6.813e07	8.738e07
<b>30</b>	<b>6.740e07</b>	<b>8.539e07</b>
35	6.649e07	8.827e07
40	6.984e07	8.892e07

**TABLE 1** – *Cross-validation results of linear regression model*

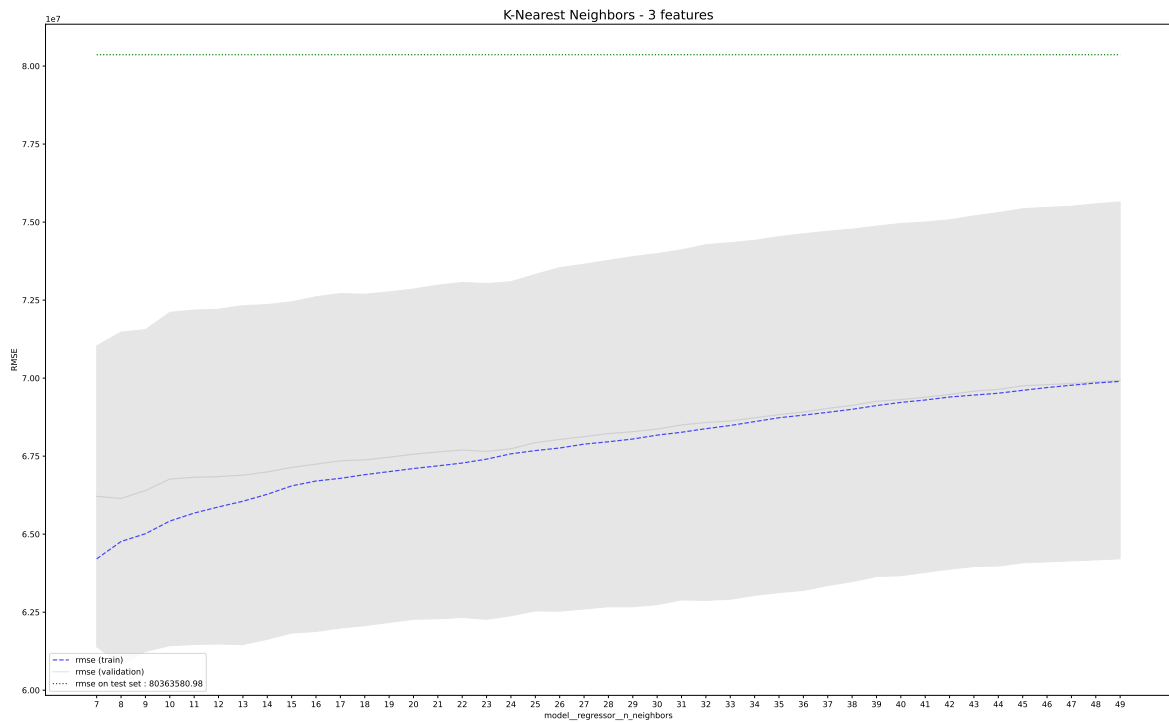
## 4.2 K-Nearest Neighbors

For the K-Nearest Neighbors model, we need to take care of adjusting the number of neighbors hyperparameter. Indeed, a value too low will induce overfitting and a value too high will induce underfitting. Moreover, because this algorithm suffers a lot from the curse of dimensionality, we ran it with a smaller subset of features : [3, 5, 7, 9, 11, 13, 15]. We grid searched for a number of neighbors between 1 and 50 first but as you can see on the figure below, this model overfits with a number of neighbors approximatively less than 7. Then, its performance slightly decrease as this number increases.



**FIGURE 1** – RMSE score with respect to the number neighbors for the best K-Nearest Neighbors model

We chose to run this algorithm again but this time starting from 7 neighbors up to 50



**FIGURE 2** – RMSE score with respect to the number neighbors for the best K-Nearest Neighbors model

You can see the results of this second pass on the following table,

#features	scores	
	validation rmse (\$)	test rmse (\$)
<b>3</b>	<b>6.614e07</b>	<b>8.036e07</b>
5	6.642e07	8.166e07
7	6.572e07	8.124e07
9	6.997e07	8.481e07
11	6.980e07	8.541e07
13	7.691e07	9.440e07
15	7.677e07	9.440e07

TABLE 2 – Cross-validation results of K-Nearest Neighbors model

The best model is the one with the **3** best features and has 8 neighbors.

### 4.3 Multi-Layer Perceptron

Since the multi-layer perceptron has the ability to set weights to zero, we can make the hypothesis that the feature selection will have less influence. Therefore, we wanted to test that model on a greater number of features. We also tested with all the features. We tested different number of layers and neurons per layers as well as different activation functions. We can observe that two layers gives the best results and adding more layers to the network tends to decrease the model performance.

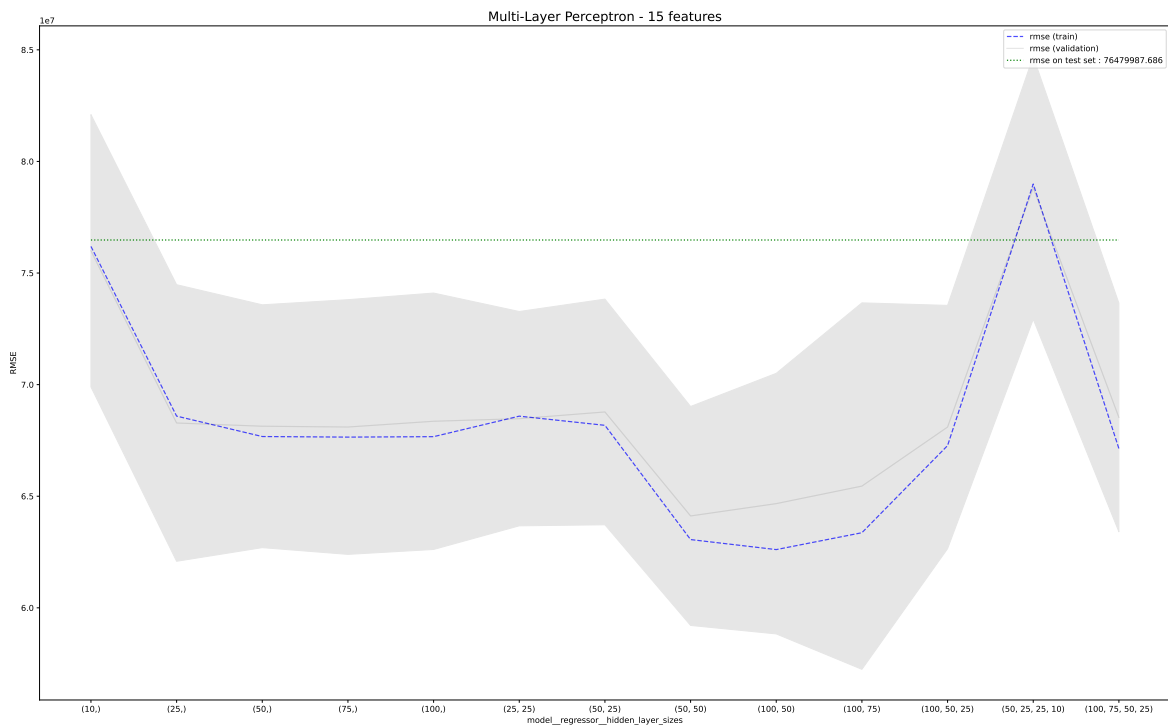


FIGURE 3 – RMSE score with respect to the number of hidden layers and their sizes for the best multi-layer perceptron model

The results are summarised in the following table,

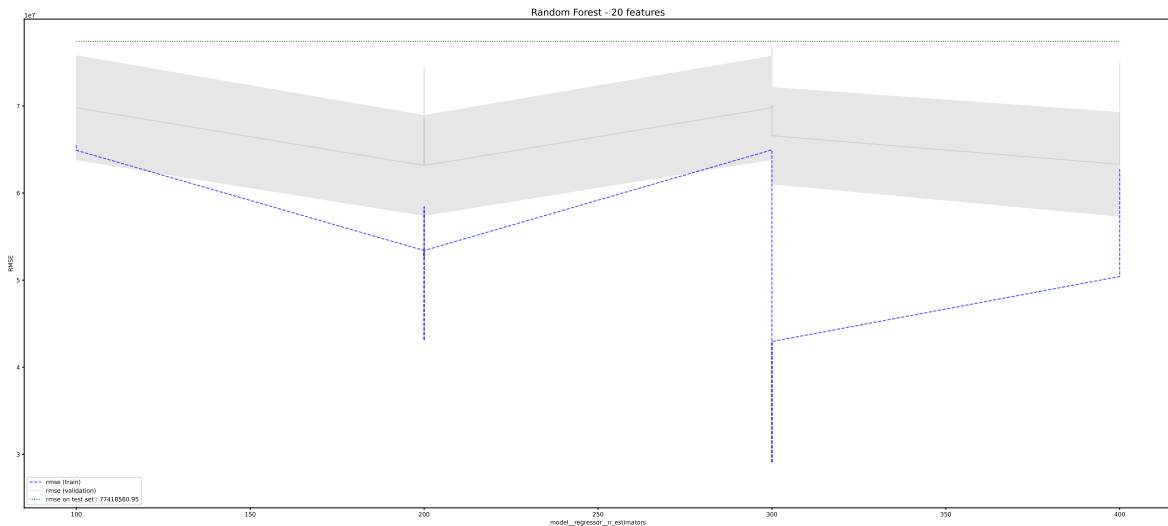
#features	scores	
	validation rmse (\$)	test rmse (\$)
5	6.674e07	8.096e07
<b>15</b>	<b>6.367e07</b>	<b>7.647e07</b>
25	6.716e07	8.272e07
35	6.674e07	8.637e07
all	6.739e07	8.289e07

**TABLE 3** – Cross-validation results of Multi-Layer perceptron model

The best model is the one with the **15** best features and has its first layer with 100 neurons and its second layer with 50 neurons. It uses a tanh · activation function and has a starting learning rate of 0.001.

#### 4.4 Random Forest

A random forest is an ensemble technique that combines multiple decision trees. As a consequence, it has a better generalization performance than a single decision tree due to randomness which decrease the model's variance. Another thing to note is that it is little sensitive to outliers and do not require much hyperparameters tuning. For this model, we performed a **Bayesian Search** with 30 iterations. Maybe we didn't find the best model but it still should be a good model. Moreover, considering the number of hyperparameters we wanted to try, it allowed use to win much computational time.



**FIGURE 4** – RMSE score with respect to the number trees for the best random forest model

#features	scores	
	validation rmse (\$)	test rmse (\$)
5	6.390e07	8.131e07
10	6.333e07	7.785e07
15	6.267e07	7.816e07
<b>20</b>	<b>6.316e07</b>	<b>7.741e07</b>
25	6.271e07	7.774e07

**TABLE 4** – *Cross-validation results of Random Forest model*

The best model is the one with the **20** best features. It uses 200 estimators.

## 5 Chosen model and predictions



## Appendix

### Mutual information matrix

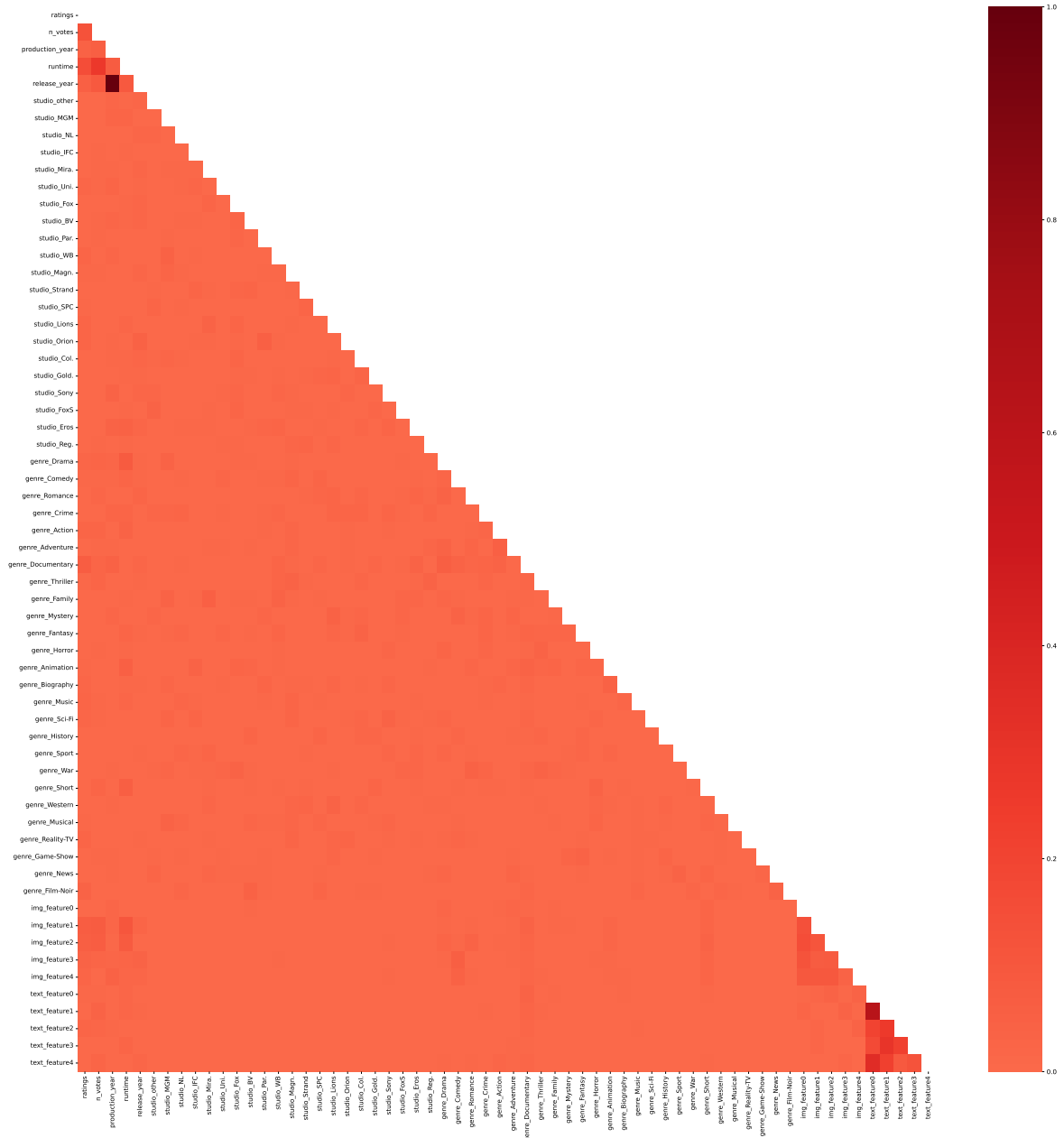
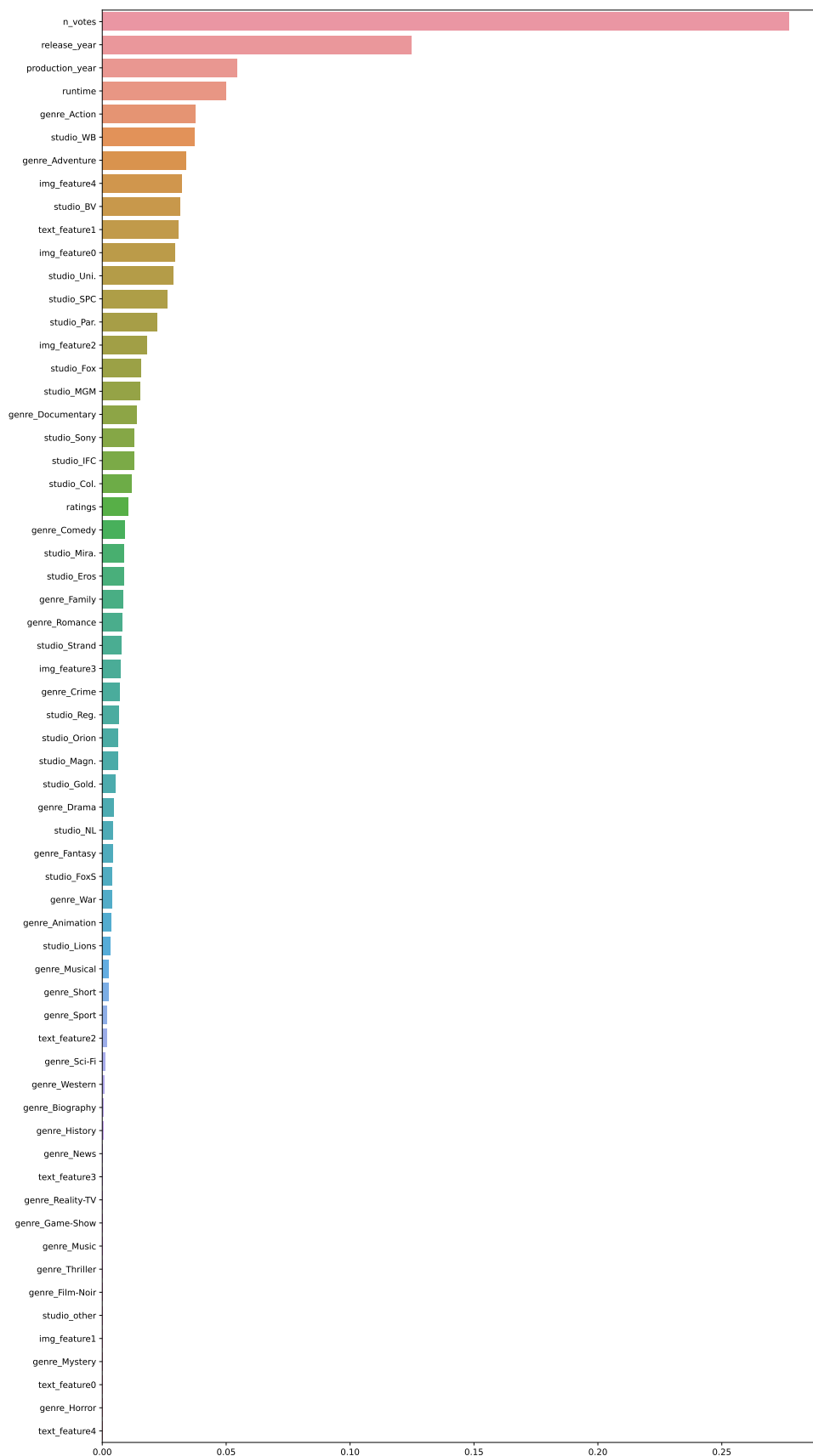


FIGURE 5 – *Mutual information matrix*



## Mutual information with the target variable



2022-2023 FIGURE 6 – Mutual information of each feature with the target variable (revenues)