# Movie Genre Prediction by Its Synopsis

Artem Razdyakonov

January 2023

## Abstract

This document is the report for NLP course project and will provide you insights about different ways to classify movies as either comedies or dramas based on their synopsis. Here is a link to this project code: https://github.com/rtm-1dyakonov/nlp_course.

## 1 Introduction

The task of determining the genre of a film based on its description has long ceased to be a challenge for machine and deep learning algorithms. Even the straightforward BoW can reach high score. It is not surprising since key words, for instance, of horror genre are quite different from key words of family film.

The moment when the task is to identify similarly close genres, complexity appears. Specifically, if you try to google a comedy film, almost every film will contain both comedy and drama genre in its description. Examples of this are presented on Fig. 1 – 4. Therefore, this research aims to check how well do deep learning algorithms cope with the task of distinguish between comedies and dramas.



Figure 1: The Intouchables genres.
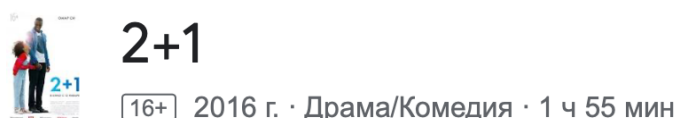


Figure 2: The Specials genres.



Figure 3: Two Is a Family genres.

Figure 4: Silver Linings Playbook genres.

## 1.1 Team

Artem Razdyakonov has performed research and prepared this document.

## 2 Related Work

The thing is there are very few works of such specifics. Most of the time the main goal of movie genre prediction projects is to train model, which will be able to distinguish several different genres, instead of a couple of close genres.

Quan Hoang [1] studied the problem of predicting genres from movie plot and used different methods like Recurrent Neural Networks and Word2Vec+XGBoost are used text classification, also K-binary transformation and probabilistic classification were employed to tackle the multi label problem. He attained a high F1-score of 0.56 along with a hit rate of 80%.

A. Blackstock and M. Spitz [2] attempt to classify movie scripts by building a logistic regression model using NLP-related features extracted from the scripts such as the ratio of descriptive words to nominals or the ratio of dialogues frames to non-dialogue frames. The experiment is done on a small dataset with only 399 scripts and the best subset of features achieves an F1-score of 0.56.

Ka wing Ho [3] investigates different methods to classify movies' genres based on synopsis. The methods examined include One-Vs-All Support Vector Machines (SVM), Multi-label K-nearest neighbor (KNN), Parametric mixture model (PMM) and Neural network. The experiment is limited to only predicting only 10 most popular genres, including action, adventure, comedy, crime, documentary, drama, family, romance, short films, and thrillers. Overall, SVM achieves the highest F1-score of 0.55.

## 3 Model Description

Before any experiment with neural network models all movies' plots were gone through preprocessing functions to make all words lowercase and get rid of stopwords. The code is represented on Fig. 5.

```
1 %%time
2
3 lower = news['Text'].str.lower()
4 noStops = removeStopWords(lower)
5 news['Text'] = noStops
6 df['OriginalPlot'] = df['Plot']
7 lower = df['Plot'].str.lower()
8 cleaned = removeStopWords(lower)
9 df['Plot'] = cleaned

CPU times: user 5min 45s, sys: 45.4 s, total: 6min 30s
Wall time: 6min 32s
```

Figure 5: Preprocessing of plots.

After the preprocessing stage all plots have been tokenized with Tokenizer from keras.preprocessing.text with padding each sequence up to the maximum length of plot vector, which is shown on Fig. 6.

```
 1  %%time
 2
 3  tokenizer = Tokenizer()
 4  tokenizer.fit_on_texts(list(df['Plot']))
 5  sequences = tokenizer.texts_to_sequences(list(df['Plot']))
 6  max_len = np.max([len(sequence) for sequence in sequences])
 7  print(f"Максимальная длина синопсиса: {max_len}")
 8  word_index = tokenizer.word_index
 9  print(f'Найдено {len(word_index)} уникальных токенов.')
10  data = pad_sequences(sequences, maxlen=max_len)
11
12  print(f'Размерность данных: {data.shape}')
```

```
Максимальная длина синопсиса: 1536
Найдено 81642 уникальных токенов.
Размерность данных: (10343, 1536)
CPU times: user 4.84 s, sys: 21.4 ms, total: 4.86 s
Wall time: 4.91 s
```

Figure 6: Using Tokenizer to create sequences.

A CNN-based approach, a hybrid CNN-LSTM approach, and an ensemble of the two previous ones are applied to the problem.

## 4   Dataset

For training and testing purposes two datasets were used, which are Wikipedia Movie Plots[1] and News Category Dataset[2]. The last one was used as supplemental training data taken from the news category dataset, which contains headlines and short descriptions of Huffington Post articles, in addition to categories that include comedy. News descriptions are quite short, so they were concatenated with headlines to add information. Examples of each dataset are presented on Fig. 7 and Fig. 8. You can see the distribution of classes in Wikipedia Movie Plots on the Tab. 1, and the distribution of classes in News Category Dataset on the Tab. 2

---

[1]The Wikipedia Movie Plots link: https://www.kaggle.com/datasets/jrobischon/wikipedia-movie-plots

[2]The News Category Dataset link:

https://huggingface.co/datasets/khalidalt/HuffPost/blob/main/News_Category_Dataset_v2.json

| Genre | Plots count |
|---|---|
| drama | 5964 |
| comedy | 4379 |

Table 1: Distribution of classes in Wikipedia Movie Plots.

| Genre | News count |
|---|---|
| comedy | 5175 |
| other | 5000 |

Table 2: Distribution of classes in News Category Dataset.



Figure 7: Wikipedia Movie Plots dataframe.



Figure 8: News Category Dataset dataframe.

On Fig. 9 you can see the result of using Tokenizer's dictionary to transform text to vector and vice.



Figure 9: Demonstration of Tokenizer.

The data is divided into training and validating sets as follows – 80% for training and 20% for validation.

## 5 Experiments

### 5.1 Metrics

To evaluate the quality of all algorithms, a standard Accuracy metric will be used since our dataset is fairly balanced.

### 5.2 Experiment Setup

Firstly, we split the data into training and validating sets. Additionally, separate training sets for training data enhanced with news data and training data consisting only of movie data were created. It is presented on Fig. 10. Moreover, the GloVE 6B 100d[3] word embeddings were loaded in hopes that they will enhance the accuracy and training speed of the neural network models. After that 2 Keras embedding layers were created to convert the texts to embeddings. One uses Glove, and the other does not.

```
 1 X_train, X_test, y_train, y_test = train_test_split(data, df['GenreID'], test_size=.2, random_state=42)
 2 testIndices = y_test.index
 3 y_train = to_categorical(y_train)
 4 y_test = to_categorical(y_test)
 5
 6 y_train_small = y_train.copy()
 7 X_train_small = X_train.copy()
 8
 9 y_train_add = to_categorical(news['Comedy'])
10 X_train_add = news_data
11 X_train = np.concatenate([X_train,X_train_add],axis=0)
12 y_train = np.concatenate([y_train, y_train_add],axis=0)
```

Figure 10: Splitting dataset with and without extra news data.

Secondly, 3 different CNN models with 3 layers were trained for 10 epochs and 128 batch size.

1. Without GloVE encodings, news and films data.
2. With GloVE encodings, news and movies data.
3. With GloVE encodings, only movies data.

---

[3]The GloVE 6B 100d link: https://www.kaggle.com/datasets/rtatman/glove-global-vectors-for-word-representation

The example of a CNN model is presented on Fig. 11.

```python
def model_train(epochs, batch):
    sequence_input = Input(shape=(max_len,), dtype='int32')
    embedded_sequences = layer_emb(sequence_input)
    x = Conv1D(128, 9, activation='relu')(embedded_sequences)
    x = MaxPooling1D(9)(x)
    x = Conv1D(128, 9, activation='relu')(x)
    x = Dropout(.4)(x)
    x = MaxPooling1D(9)(x)
    x = Conv1D(128, 9, activation='relu')(x)
    x = Dropout(.4)(x)
    x = MaxPooling1D(9)(x)

    x = Flatten()(x)
    x = Dense(128, activation='relu')(x)
    x = Dense(len(genres), activation='softmax')(x)

    model = Model(sequence_input, x)
    model.compile(loss='categorical_crossentropy',
                  optimizer='adam',
                  metrics=['acc'])
    X_train.shape
    model.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=epochs, batch_size=batch)

    return model
```

Figure 11: CNN model.

Thirdly, 2 different LSTM-CNN models were trained for 10 epochs and 128 batch size:

1. With GloVE encodings, news and movies data.

2. With GloVE encodings, only movies data.

The example of a LSTM-CNN model is presented on Fig. 12.

```python
def lstm_model_train(epochs, batch):
    word_indices = Input(shape=(max_len,), dtype='int32')
    embeddingsLSTM = layer_emb(word_indices)
    x=Conv1D(128, 9, activation='relu')(embeddingsLSTM)
    x=MaxPooling1D(9)(x)
    x=Conv1D(128, 9, activation='relu')(x)
    x = Dropout(.4)(x)
    x=MaxPooling1D(9)(x)
    x=Conv1D(128, 9, activation='relu')(x)
    x = Dropout(.4)(x)
    x=MaxPooling1D(9)(x)
    X =  LSTM(128,return_sequences=False)(x)
    X = Dropout(.65)(X)
    X = Dense(len(genres),activation='softmax')(X)
    X = Activation('softmax')(X)
    LSTMmodel = Model(inputs = word_indices, outputs = X)
    LSTMmodel.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
    LSTMmodel.fit(X_train, y_train, validation_data=(X_test, y_test), epochs=epochs, batch_size=batch)

    return LSTMmodel
```

Figure 12: LSTM-CNN model.

Additionally, special function was created to transform any input plot into vector and send it into model. This function is presented on Fig. 13.

```
1 def pred_your_synop(test_plot):
2    lower = test_plot.lower()
3    removed = []
4    removed.append(' '.join([word for word in lower.split() if word not in nltk.corpus.stopwords.words('english')]))
5    test_token = tokenizer.texts_to_sequences(removed)
6    test_token = np.array(test_token, dtype='int32')
7    test_data = pad_sequences(test_token, maxlen=max_len)
8    test_preds = modelSmall.predict(test_data)
9    answ = 'comedy' if test_preds[0][0]<test_preds[0][1] else 'drama'
10   return answ
```

Figure 13: Function to vectorize input text and make predictions

Finally, the ensemble of all 5 models mentioned above was created. To determine the best combination of weights of each model special algorithm was prepared. It is shown on Fig. 14.

```
1 weights = np.arange(.0, 2.1, 0.1)
2
3 weights_combs = list(combinations(weights, 5))
4 weights_combs_list = []
5 for tpl in weights_combs:
6    temp_arr = []
7    for val in tpl:
8        temp_arr.append(val)
9    weights_combs_list.append(temp_arr)
```

```
1 from itertools import combinations
2
3 models_preds = [CNNpreds,
4                 LSTMpreds,
5                 noGloveCNNpreds,
6                 LSTMpreds_small,
7                 CNNpreds_small]
8 max_score = 0
9 res_dict = {}
10 for comb in range(len(weights_combs_list)):
11    avgPreds = np.average(models_preds, weights=weights_combs_list[comb], axis=0)
12    preds_df = pd.concat([pd.DataFrame(avgPreds),df.loc[testIndices,['OriginalPlot','Genre','Title']].reset_index()],axis=1)
13    preds_df['Predicted Genre'] = (preds_df[0]>preds_df[1]).replace(True,'drama').replace(False,'comedy')
14    accuracy = (preds_df['Predicted Genre'] == preds_df['Genre']).sum()/len(preds_df)
15    if accuracy > max_score:
16        max_score = accuracy
17        res_dict['best_comb'] = (weights_combs_list[comb], accuracy)
```

```
1 res_dict
```

```
{'best_comb': ([1.2000000000000002,
  1.3,
  1.7000000000000002,
  1.8,
  1.9000000000000001],
 0.6916384726921218)}
```

Figure 14: Algorithm finding the best weights for ensemble.

# 6 Results

The final metrics of CNN models with training set are presented on the Tab. 3.

| Model | Accuracy |
|---|---|
| CNN, no GloVE encodings, news and films data | 0.66 |
| CNN, GloVE encodings, news and movies data | 0.65 |
| CNN, GloVE encodings, only movies data | 0.86 |

Table 3: Accuracy scores of CNN models.

The final metrics of LSTM-CNN models with training set are presented on the Tab. 4.

| Model | Accuracy |
|---|---|
| LSTM, GloVE encodings, news and films data | 0.62 |
| LSTM, GloVE encodings, only movies data | 0.8 |

Table 4: Accuracy scores of LSTM-CNN models.

The accuracy scores of all 5 models with validation data are shown on the Tab. 5 below.

| Model | Accuracy |
|---|---|
| CNN, GloVE encodings, news and movies data | 0.623 |
| LSTM, GloVE encodings, news and films data | 0.657 |
| CNN, no GloVE encodings, news and films data | 0.681 |
| LSTM, GloVE encodings, only movies data | 0.656 |
| CNN, GloVE encodings, only movies data | 0.666 |

Table 5: Accuracy scores of models with validation set.

As shown on Fig. 14, the highest accuracy score achieved with ensemble is 1.2 percentage point higher than the best individual model.

## 7   Conclusion

Eventually, the preparation and processing of datasets with news and short film plots during the execution of this course project. The final result of this work can be considered several models, which are capable of determine whether the input plot is comedy or drama. The example of using your own plot as input is shown on. Fig. 15.

```
1 ### Мальчишник в Вегасе ###
2
3 test_plot = "They dreamed of having an unforgettable bachelor party in Vegas. But now they need to remember exactly what
4 answ = pred_your_synop(test_plot)
5 answ

1/1 [==============================] - 0s 44ms/step
'drama'
```

```
1 ### Мусорная охота ###
2
3 test_plot = "Milton Parker (Vincent Price), an eccentric game inventor, dies after losing a video game with his nurse. P
4 answ = pred_your_synop(test_plot)
5 answ

1/1 [==============================] - 0s 46ms/step
'comedy'
```

```
1 ### 1+1 ###
2
3 test_plot = 'Having suffered an accident, the wealthy aristocrat Philip hires the person who is least suitable for this
4 answ = pred_your_synop(test_plot)
5 answ

1/1 [==============================] - 0s 40ms/step
'drama'
```

Figure 15: Predictions with your own plots.

## References

[1] Quan Hoang. Predicting Movie Genres Based on Plot Summaries, 2018.

[2] Alex Blackstock and Matt Spitz. Classifying movie scripts by genre with a memm using nlp-based features, 2008.

[3] Ka wing Ho. Movies genres classification by synopsis, 2011.