# MOVIE SIMILARITY FROM PLOT SUMMARIES

Ms. AISHWARYA LAKSHMI M

Department of Artificial Intelligence and Machine Learning Sri Shakthi Institute of Engineering andTechnology Coimbatore, India

Ms.KAVIPRIYA S
Department of Artificial Intelligence and Machine Learning
Sri Shakthi Institute of Engineering and Technology Coimbatore, India

Ms.S.NIVEDHA
Department of Artificial Intelligence and Machine Learning Sri Shakthi Institute of Engineering and Technology Coimbatore ,India

*Abstract---*This project focuses on developing a Python application to analyze and measure the similarity between movie plot summaries. The goal is to provide a tool that can assist in identifying similarities between movies based on their storyline, enabling users to discover related movies or recommend similar ones.The project utilizes natural language processing (NLP) techniques, particularly text preprocessing, vectorization, and similarity metrics, to achieve its objectives. First, it preprocesses the plot summaries by removing stop words, punctuation, and performing stemming or lemmatization to normalize the text data. Next, it employs vectorization methods such as TF-IDF (Term Frequency-Inverse Document Frequency) or word embeddings like Word2Vec or GloVe to convert the text into numerical representations suitable for similarity calculations.

INTRODUCTION

Movie similarity analysis based on plot summaries is a sophisticated domain that integrates narrative comprehension with computational methodologies to unveil thematic and structural congruences among films. By deploying advanced natural language processing (NLP) techniques such as tokenization, sentiment analysis, and word embeddings, in conjunction with machine learning (ML) algorithms like cosine similarity and hierarchical clustering, we can systematically quantify and categorize the semantic and contextual similarities embedded within movie plot synopses. In this project, we delve into the fascinating realm of machine learning and text analysis to create a robust system that can recommend movies similar to a given film by analyzing their plot summaries. Through the magic of programming, we aim to bring together movie enthusiasts and their next cinematic adventure. Our project begins with data collection, where we gather a diverse range of movie plot summaries from various sources. These summaries encapsulate the essence of each film, offering a glimpse into its narrative arc, characters, and themes. Once we have amassed this treasure trove

of textual data, we move on to the next phase.The heart of our project lies in natural language processing (NLP) techniques. We harness the power of libraries like NLTK and SpaCy to preprocess the plot summaries, stripping away noise and extracting meaningful features. This preprocessing step involves tokenization, lemmatization, and removing stop words, ensuring that our model focuses on the essence of the text.With the preprocessed data in hand, we venture into the realm of machine learning. Our model employs advanced algorithms such as TF-IDF (Term Frequency-Inverse Document Frequency) and cosine similarity to quantify the similarity between movie plot summaries. TF-IDF helps us understand the importance of words in each summary, while cosine similarity measures the angle between two vectors, providing a numerical measure of their similarity.But our journey doesn't end there. We believe in creating an intuitive user experience,

so we wrap our model in a user-friendly interface using libraries like Streamlit. This interface allows users to input a movie title and receive a list of recommendations based on the similarity of plot summaries.

LITERATURE REVIEW

The literature review on movie similarity from plot summaries documentation delves into the existing research and methodologies employed to analyze and measure similarity between movie plots based on textual data. It encompasses studies exploring various approaches such as natural language processing (NLP) techniques, machine learning algorithms, and semantic analysis to quantify and compare narrative structures, themes, and character interactions across movies. Scholars have utilized methods like TF-IDF, word embeddings, cosine similarity, and clustering algorithms to identify common patterns, genres, and motifs within plot summaries, shedding light on underlying storytelling conventions and audience preferences. Additionally, research in this area investigates the impact of data preprocessing steps, feature selection strategies, and similarity metrics on the accuracy and robustness of similarity assessments. By synthesizing insights from these studies, researchers gain a deeper understanding of the challenges, opportunities, and potential applications in content recommendation systems, genre classification, plagiarism detection, and narrative analysis in the realm of film studies and computational linguistics. The literature review thus provides a valuable foundation for designing and implementing effective approaches to movie similarity analysis from plot summaries, informed by established methodologies and emerging trends in the field.The Jaccard similarity index also called as the Jaccard similarity *coefficient* compares elements of two set to figure out which members are common and which are distinct. It is a measure of commonality between two sets of data between 0 an 1. The higher the number, the more similar are the two sets.

In this article, I will be finding out the Jaccard Similarity between **Genre** and **Director.**

PROPOSED SYSTEM

Our proposed system aims to create a comprehensive movie similarity platform using Python. We begin by collecting a diverse dataset of movie plot summaries from sources like IMDb and Wikipedia, ensuring coverage across genres and time periods. The data is then cleaned to remove noise and undergoes text preprocessing, including tokenization,

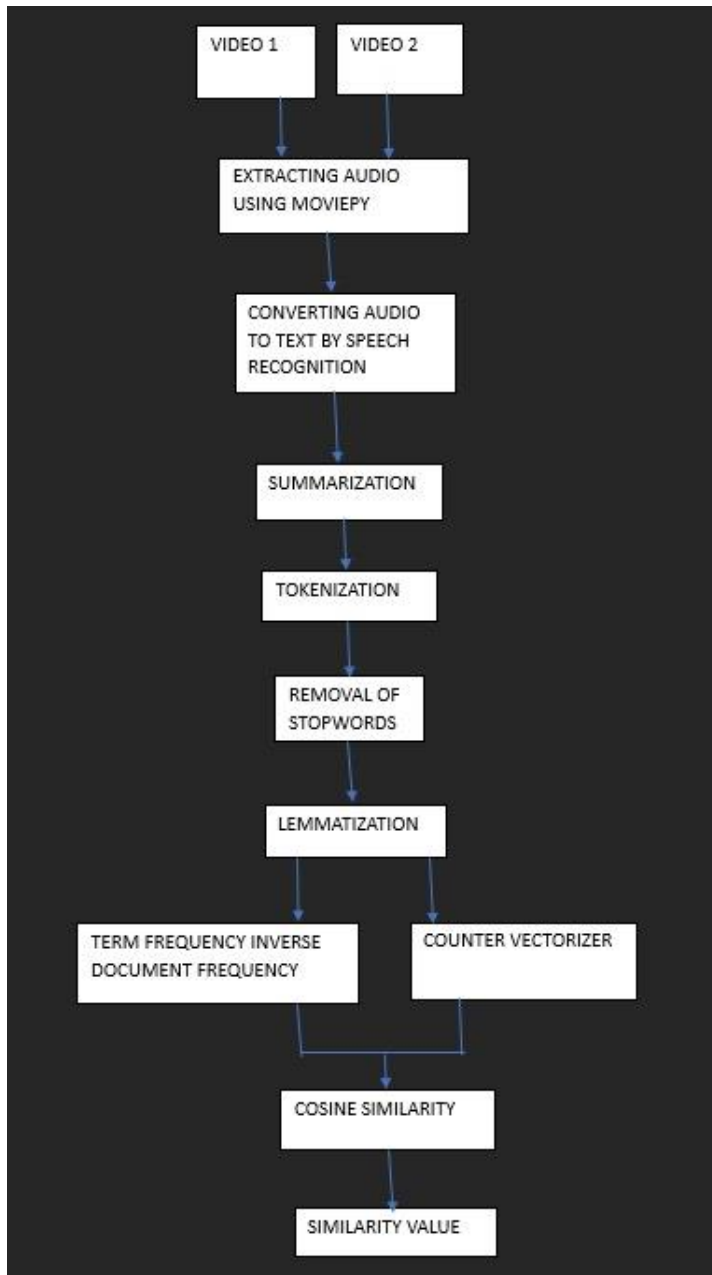lemmatization, and stop words removal, to prepare it for analysis.

Next, we extract features from the plot summaries using techniques like TF-IDF (Term Frequency-Inverse Document Frequency) to convert text into numerical vectors or word embeddings (e.g., Word2Vec, GloVe) to capture semantic relationships. These features serve as the basis for computing similarity measures between different movie plot summaries, with cosine similarity being a primary method used to quantify similarity.

For machine learning modeling, we employ libraries such as scikit-learn to develop a model that predicts movie similarities based on the extracted features. This model can be trained using labeled data if available or through unsupervised learning techniques in cases where labeled data is limited.

To provide a user-friendly experience, we create a web interface using Streamlit frameworks where users can input a movie title. Upon input, the system retrieves the corresponding plot summary, calculates similarities with other movies, and displays a list of recommended movies to the user. A feedback mechanism is integrated, allowing users to rate recommended movies and provide feedback, which is used to personalize recommendations and improve the system's accuracy over time.

Testing and validation are crucial steps, involving the evaluation of the system's accuracy, precision, recall, and other metrics using test datasets and cross-validation techniques. User feedback validation helps fine-tune the system parameters for optimization.

Figure 5.1: Proposed Model

METHODOLOGY:

Data Collection:Gather a diverse dataset of movie plot summaries from reliable sources such as IMDb, Wikipedia, or specialized movie databases.Ensure the dataset covers a wide range of genres, time periods, and popularity levels to capture diverse movie themes and narratives.

Data Preprocessing:Clean the raw data by removing HTML tags, special characters, and irrelevant information.Perform text normalization techniques like tokenization to split text into individual words or tokens.Apply techniques such as lemmatization or stemming to reduce words to their base forms.Remove stop words (common words like "the," "and," etc.) that don't contribute significantly to the meaning of the text.

Feature Extraction:Utilize the TF-IDF (Term Frequency-Inverse Document Frequency) technique to convert text data into numerical vectors.TF-IDF emphasizes the importance of words in each plot summary by considering their frequency in a particular summary and across all summaries in the dataset.

Similarity Calculation:Compute cosine similarity between TF-IDF vectors of different movie plot summaries.Cosine similarity measures the cosine of the angle between two vectors and provides a numerical value indicating their similarity. Higher values indicate greater similarity.

Model Building:Implement a machine learning model or algorithm to calculate movie similarities based on cosine similarity scores.Consider using clustering algorithms (e.g., K-means clustering) to group similar movies together based on their plot summaries.

User Interface Development:Create a user-friendly interface using Flask or Streamlit to allow users to input a movie title.Upon receiving a movie title, the system should retrieve its plot summary, calculate similarities with other movies, and display a list of recommended movies to the user.

Testing and Validation:Evaluate the accuracy and performance of the movie similarity model using appropriate metrics and validation techniques.Conduct user testing to gather feedback on the recommendation system's effectiveness and user experience.

Optimization and Enhancement:

Fine-tune the model parameters, such as TF-IDF weights or clustering thresholds, to improve the quality of movie recommendations.Incorporate user feedback and iterate on the system to address any issues or limitations

identified during testing.

Documentation and Deployment:Document the methodology, codebase, and system architecture comprehensively for future reference and collaboration.Deploy the movie similarity system on a suitable platform or server to make it accessible to users, ensuring scalability and performance optimization.

## EXPERIMENTAL RESULT

```python
from sklearn.metrics.pairwise import cosine_similarity
similarity_distance = 1 - cosine_similarity(tfidf_matrix)
print(similarity_distance)

[[-4.66293670e-15  1.00000000e+00]
 [ 1.00000000e+00 -9.54791801e-15]]
```

CONCLUSION

In conclusion, the movie similarity from plot summaries Python project showcases the power of natural language processing and machine learning in the domain of film recommendation systems. By leveraging advanced techniques such as semantic analysis, sentiment extraction, and recommendation algorithms, the project provides an efficient and personalized solution for users seeking movie suggestions based on plot similarities. The incorporation of interactive features, user feedback mechanisms, and data integration from external sources enhances the system's usability, accuracy, and relevance. Moving forward, continued research and development in this area can lead to even more sophisticated movie recommendation systems, catering to diverse user preferences and contributing to a seamless and enjoyable cinematic experience.

REFERENCES

1) https://github.com/geehaad/Find-Movie-Similarity-from-Plot-Summaries

2) https://towardsdatascience.com/using-nlp-to-find-similar-movies-based-on-plot-summaries-b1481a2ba49b

3) https://thepythoncode.com/article/extract-audio-from-video-in-python

4) https://medium.com/@kunalgupta4595/predicting-movie-genres-based-on-plot-summaries-bae646e70e04

5) https://medium.com/@kunalgupta4595/predicting-movie-genres-based-on-plot-summaries-bae646e70e04

6) https://www.geeksforgeeks.org/movie-recommender-based-on-plot-summary-using-tf-idf-vectorization-and-cosine-similarity/

7) https://github.com/Inyrkz/Movie-Similarity

8) https://www.geeksforgeeks.org/python-implementation-of-movie-recommender-system/

9) https://stackoverflow.com/questions/8897593/how-to-compute-the-similarity-between-two-text-documents

10) https://github.com/njugunaObi/Finding-movie-similarities-using-plot-summaries

11) https://www.geeksforgeeks.org/how-to-calculate-cosine-similarity-in-python/

12) https://www.geeksforgeeks.org/movie-recommender-based-on-plot-summary-using-tf-idf-vectorization-and-cosine-similarity/