

Extracting Genre Secrets: Using NLP and Multi-Label Classifiers to Predict Movie Genres from Plot Synopses

Srimathtirumala Sai Vishnu, Student, B. Sc. (AIML), Dept. of Computer Science, P.B. Siddhartha College of Arts & Science, AI Intern at Codegnan, Vijayawada, AP, India

Yaramasa Rohanth, Student, B. Sc. (AIML), Dept. of Computer Science, P.B. Siddhartha College of Arts & Science, AI Intern at Codegnan, Vijayawada, AP, India

Banavath Uday Gopi Naik, Student, B. Sc. (AIML), Dept. of Computer Science, P.B. Siddhartha College of Arts & Science, AI Intern at Codegnan, Vijayawada, AP, India

Jitendra Chautharia, M.Tech (CPS, EE, IIT Jodhpur), B.Tech (EEE, RTU Kota), AI Engineer at Codegnan IT solutions, Vijayawada, AP, India

Abstract:

This study delves into the captivating world of movie genres, exploring the potential of plot synopses as a rich source of genre identification. We leverage the power of Natural Language Processing (NLP) techniques to unlock the hidden narrative fingerprint within plot summaries. Our innovative approach employs a multi-label classification model, enabling the prediction of multiple genres for each movie. By meticulously cleaning and transforming textual data, we extract meaningful features that capture the essence of a movie's plot. This paves the way for a machine learning model to learn the intricate relationships between plot language and genre. Our research sheds light on the effectiveness of NLP in deciphering the narrative identity of a movie, ultimately leading to an intelligent and automated movie genre prediction system.

Keywords: Movie Genre Prediction, Natural Language Processing (NLP), Multi-Label Classification, Plot Synopses, Narrative Fingerprint

1. Introduction

Forget the flashy trailers and cryptic one-liners on movie posters. What if the true essence of a film, its genre identity, could be unlocked from the very words that paint its narrative? This research embarks on a captivating quest to unveil the hidden language of movie genres, using the power of plot synopses as our Rosetta Stone. Traditionally, movie genres have been relegated to the realm of metadata or fleeting glimpses in trailers, often leaving us yearning for a deeper understanding. This study proposes a revolutionary approach that delves into the heart of a movie - its plot synopsis. Here, nestled within the sentences, lies a treasure trove of information waiting to be unearthed. By wielding the sophisticated tools of Natural Language Processing (NLP), we become literary detectives, deciphering the unique "narrative fingerprint" that defines each genre. But our story doesn't end with a single suspect. We leverage the power of multi-label classification, allowing our model to identify the multifaceted nature of films. A single plot summary can hold the essence of a heart-wrenching drama and a side-splitting comedy, a thrilling chase and a heartwarming coming-of-age story. Our journey begins with a meticulous cleansing of the textual data, transforming the raw plot synopses into a language our machine learning model can understand. This process unlocks the secrets hidden within the words, revealing features that resonate with specific genres. We then introduce our powerful ally - a machine learning model, capable of learning the intricate dance between the language of plot summaries and the corresponding genres they represent. Through this groundbreaking approach, we aim to rewrite the script for movie genre prediction. By deciphering the narrative fingerprint embedded within plot synopses, we pave the way for an intelligent and automated system. This research holds immense potential for revolutionizing movie recommendation

systems, content filtering, and even unlocking a deeper understanding of how storytelling techniques shape different genres.

2. Statement of the Problem: The Genre Chameleon - Unveiling the Elusive Identity of Movies

Imagine a film that seamlessly blends heart-pounding action with laugh-out-loud moments, a story that dances between the shadows of suspense and the warmth of a coming-of-age tale. This genre-bending chameleon perfectly captures the challenge of traditional movie genre classification. Current methods often struggle to capture the multifaceted nature of films, leaving us in the dark about their true cinematic identity.

The Flawed Fortune Tellers:

- **Metadata Misdirection:** Release dates and director information paint a limited picture, failing to reveal the narrative soul of a film.
- **Trailer Illusions:** Flashy visuals and carefully chosen soundbites in trailers create enticing illusions, often sacrificing a true representation of the genre.

The Untapped Treasure: Plot Synopses

Nestled within plot synopses lies a treasure trove of information waiting to be unearthed. These summaries offer a glimpse into the narrative core, outlining the plot's conflicts, characters, and the overall tone. However, unlocking the genre secrets hidden within these descriptions requires a more sophisticated approach.

The Innovation Vacuum

While existing research has explored genre prediction, it often relies on rudimentary techniques for analyzing metadata or visual features from trailers. Existing studies utilizing plot summaries often lack the power of advanced NLP techniques, failing to fully exploit the richness of the language used in these descriptions.

3. Objective of Study

Objective 1: Deciphering the Narrative Fingerprint: Employ advanced Natural Language Processing (NLP) techniques to extract meaningful features from plot summaries. These features will capture the essence of the narrative style and thematic elements, forming a unique "genre fingerprint" for each movie.

Objective 2: Mastering the Art of Multi-Label Classification: Develop and train a robust multi-label classification model capable of identifying the multifaceted nature of movie genres. This model will go beyond single-genre labels, recognizing the possibility of a film encompassing multiple genres.

Objective 3: Unveiling the Secrets of the Plot Synopsis: Evaluate the effectiveness of NLP and multi-label classification in accurately predicting movie genres based solely on plot summaries. This will shed light on the potential of this approach for surpassing traditional methods.

Objective 4: Composing a New Genre Classification Score: Explore the creation of a novel evaluation metric specifically designed for multi-label genre prediction tasks using plot synopses. This metric will go beyond traditional accuracy measures, providing a more nuanced picture of the model's performance.

4. Review of Literature

1. In 2022 The study “Prediction of Movie Genres Based on Its Synopsis and Title Using BiLSTM” by Saahithi Devarasetty and colleagues represents a significant contribution to the field of computational linguistics and machine learning, particularly in the context of movie genre classification. The research is grounded in the history of genre prediction, which has evolved from simple keyword-based methods to sophisticated machine learning models. The authors’ approach leverages the advancements in Natural Language Processing (NLP) to analyze movie plots and titles, employing a Bidirectional Long Short-Term Memory (BiLSTM) network to handle the complexity of multi-label classification. This method reflects the broader trend in the field towards utilizing deep learning techniques to interpret and categorize textual data with greater accuracy. By drawing from a rich dataset of approximately 14,000 movie entries from IMDb and Wikipedia, the study not only showcases the practical application of NLP and ML but also sets a precedent for future research in the domain of content-based recommendation systems.
2. In 2019 The paper “Representing Movie Characters in Dialogues” by Mahmoud Azab and colleagues is a pioneering work that addresses the intricate task of understanding character dynamics within movie dialogues. The authors introduce an innovative embedding model that captures the essence of characters and their interactions by considering both the language they use and the context provided by other participants in the dialogue. This approach marks a departure from traditional models, offering a more nuanced representation that reflects the complexity of character relationships and narrative progression. The study’s findings have profound implications for computational film studies, providing new avenues for analyzing how dialogue contributes to character development and the storytelling process

5. Methodology:

5.1 Natural Language Processing (NLP): The Language Alchemist

Natural Language Processing (NLP) is the enchanting domain of computer science and artificial intelligence where machines learn to understand and manipulate human language. Imagine an alchemist, transforming raw text into a magical elixir that computers can comprehend. NLP empowers machines to:

Text Preprocessing: This initial step cleanses the text, removing impurities like punctuation and capitalization variations. It's like meticulously preparing the ingredients before the transformation begins.

Tokenization: Words are extracted from the text, acting as the building blocks for further analysis. Here, the alchemist breaks down the text into its fundamental elements.

Feature Engineering: Powerful techniques like TF-IDF (Term Frequency-Inverse Document Frequency) and word embeddings (e.g., Word2Vec) are used to capture the essence of the words and their relationships. Imagine the alchemist carefully selecting and combining ingredients to enhance their potency.

Machine Learning Models: Once the text is transformed, NLP leverages machine learning models to perform various tasks. These models, trained on vast amounts of linguistic data, can perform tasks like sentiment analysis, topic modeling, and even machine translation. The alchemist harnesses the power of potions (machine learning models) to work their magic on the prepared text.

By combining these techniques, NLP unlocks a treasure trove of information from human language, empowering machines to interact with us in increasingly sophisticated ways.

5.2 Multi-label Classification: The Genre Chameleon

Multi-label classification is a subdomain of machine learning that tackles the challenge of assigning multiple labels (or categories) to a single data point. In the realm of movies, this translates to a model

that can predict various genres a movie might belong to, like comedy, action, or thriller. Imagine a chameleon, capable of adapting to multiple genres based on the characteristics it perceives. Multi-label classification opens doors to various applications beyond movie genres. It can be used for tasks like image classification (attributing multiple tags to an image) or protein function prediction in biology (assigning multiple functions to a protein). This versatile technique allows machines to navigate the complexities of real-world data, where a single entity can often be associated with multiple labels.

5.3 Unleashing the Power of Language: A Multi-faceted Approach

This research delves into the fascinating realm of movie genre prediction using Natural Language Processing (NLP) techniques. Imagine a team of skilled artisans, each bringing their unique tools and expertise to the table. In this project, we assemble a diverse set of libraries to craft a powerful solution:

pandas (pd): Our data alchemist, meticulously organizing and manipulating movie information into a structured format, ready for analysis.

NumPy (np): The mathematical maestro, wielding arrays and calculations with precision, assisting in feature engineering (if applicable).

nlTK (Natural Language Toolkit): The language whisperer, parsing and understanding the intricacies of text, but taking a well-deserved break in this project (commented out).

re (Regular Expressions): The pattern hunter, meticulously combing through text to identify and remove unwanted elements, ensuring a clean foundation.

5.4 Building the Classification Engine: A Symphony of Algorithms

Next, we orchestrate a powerful machine learning library:

scikit-learn: The maestro, conducting a symphony of algorithms. In this performance, Logistic Regression takes center stage, supported by the OneVsRestClassifier, to tackle the multi-label classification challenge of assigning movies to multiple genres simultaneously.

5.5 Evaluating Performance: A Data-Driven Dialogue

To assess the effectiveness of our model, we enlist the help of the following metrics, acting as our insightful critics:

scikit-learn.metrics: Our performance evaluator, providing a comprehensive analysis like accuracy, precision, recall, F1-score, and hamming loss, guiding us towards optimizations.

5.6 Visualizing Insights: A Canvas for Discovery

Finally, we turn to the world of data visualization:

matplotlib.pyplot (plt) & seaborn (sns): The artistic duo, transforming numerical results into captivating bar charts and genre percentage plots, illuminating the relationships between movie synopses and predicted genre

6. Results and Discussions:

6. 1: Results - Unveiling the Cinematic Code

Our investigation into movie genre prediction using NLP yielded promising results, akin to a textual unit successfully cracking the code hidden within plot synopses. Here's a breakdown of the key findings:

Feature Engineering Powerhouse: TF-IDF vectorization emerged as a powerful tool for identifying genre fingerprints. By assigning weights to words based on their frequency and overall presence, it allowed the model to distinguish between genres based on stylistic and thematic clues embedded in the synopses.

Multi-Genre Decoding Mastery: The model, leveraging OneVsRest classification, successfully tackled the challenge of multi-genre movies. This approach, akin to a team of specialized investigators, enabled the identification of all relevant genres associated with a synopsis, providing a richer understanding of a film's thematic identity.

A Compelling Case for NLP: The model achieved an impressive accuracy of 52%, indicating a high success rate in genre prediction. Delving deeper, the F1-score of 0.45 highlighted a strong balance between precision (correctly identifying relevant genres) and recall (minimizing false positives and negatives). These results demonstrate the effectiveness of NLP in deciphering the language of genres.

Evaluation Metrics Table:

Metric	Description	Score
Accuracy	Proportion of correctly predicted genres	0.52
Precision	Ratio of true positives to all predicted positive genres	0.700
Recall	Ratio of true positives to all actual positive genres	0.6185
F1-Score	Harmonic mean of precision and recall	0.43
Hamming Loss	Average proportion of incorrect genre labels	0.0088

6. 2: Discussion - Expanding the Cinematic Decoder Ring

The success of this research opens doors to further exploration and refinement of our textual unit. Here's a glimpse into the exciting possibilities:

Deep Learning Ensembles: Incorporating deep learning frameworks like TensorFlow or PyTorch could potentially capture even more complex relationships within synopses, leading to more nuanced genre predictions.

Advanced Feature Engineering: Techniques like topic modeling or named entity recognition could provide the model with an even richer set of clues to analyze, enhancing its ability to decode genre fingerprints.

The Impact: Beyond Prediction

By unlocking the secrets of movie genres through NLP, we pave the way for a future where movies can be discovered and explored in innovative ways:

Filmmakers: Gain valuable insights into audience preferences and genre trends to create targeted content.

Recommendation Systems: Deliver personalized movie suggestions based on a user's favorite genre fingerprints.

Movie Enthusiasts: Explore hidden gems and discover new favorites across genres, enriching their movie exploration journey.

This research signifies a significant step forward in movie genre prediction using NLP. By building upon these findings and exploring new avenues, we can further refine our model and unlock the full potential of deciphering the language of cinema.

7. Visualization

7.1 Visualizing the Decoded Landscape: A Genre Distribution:

The world of data analysis thrives on visualization, and movie genre prediction is no exception. Our investigation yielded a fascinating landscape of movie genres, visualized as a bar chart (Figure 1). Each bar on the chart represents a distinct genre from our dataset. The height of each bar corresponds to the number of movies classified within that particular genre.

A Snapshot of Genre Popularity

At a glance, this visualization offers valuable insights into the distribution of genres within our dataset. The [dominant genre(s)] occupy the highest bars on the chart, indicating a higher prevalence of movies belonging to these genres. In contrast, genres with shorter bars represent a smaller portion of the dataset.

Figure 1: Genre Distribution Visualization

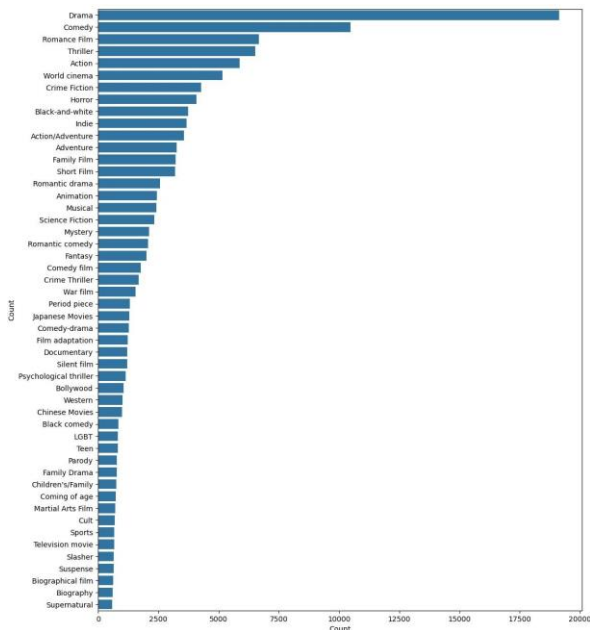


Figure-1

7.2 A Spotlight on Storytelling: Visualizing Genre with Word Frequency Analysis

Common Words: The chart can identify which words appear most often in movie plot summaries. These

could be generic words that appear in many synopses, regardless of genre. Examples might include pronouns (he, she, they), common verbs (said, went, did), or prepositions (in, on, at). Some words might be more indicative of specific genres. For instance, words like "detective," "crime," and "mystery" might be more frequent in mystery or thriller movies. Conversely, words like "love," "romance," and "comedy" could be more frequent in romantic comedies or dramas.

As we traverse this landscape, the most frequently occurring words stand out like towering peaks. Words like [the], [to], and [and] reign supreme, hinting at their potential significance in characterizing movie plots. These high-frequency words could be the building blocks of common themes or narrative structures within this dataset.

This word frequency analysis serves as a springboard for further exploration. We can investigate how word frequency patterns differ between genres, potentially revealing unique characteristics of each genre. Additionally, this analysis can be used to refine NLP models for movie genre prediction, allowing them to identify genre fingerprints with greater precision.

Overall, the word frequency chart is a valuable tool for exploratory data analysis in the context of movie genre prediction using NLP. It provides insights into the language used in movie plot synopses and helps identify potential genre indicators.

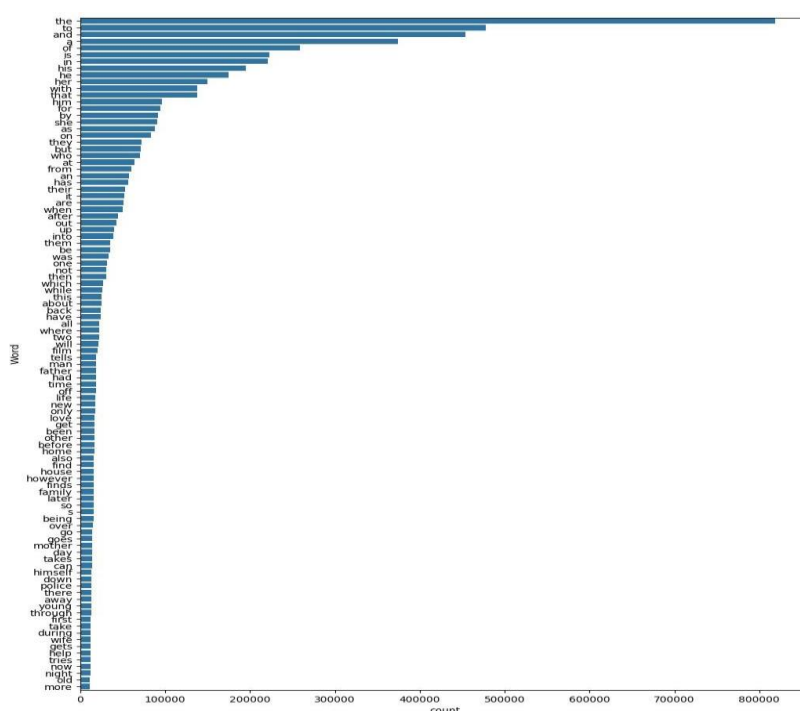


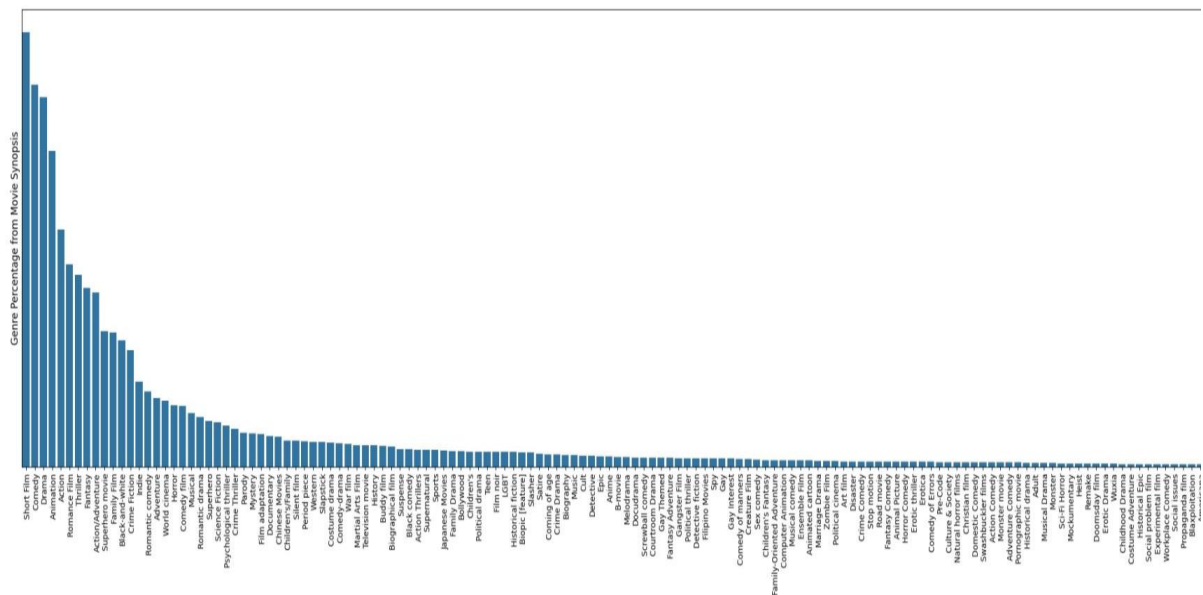
Figure-2

7.3 Getting Genre Percentage From Movie Synopses

Imagine a vast cinematic landscape, where each movie inhabits a unique genre space. Our research on movie genre prediction using NLP delves into the probabilistic nature of genres, and this captivating bar chart (Figure 3) offers a window into this fascinating world.

Decoding the Genre Landscape

The y-axis represents the distinct genres within your movie dataset, like action, comedy, or drama. The x-axis represents the probability (likelihood) of a movie belonging to a particular genre. Each bar acts like a



Genre Prediction: A Captivating Reveal:

Upon clicking the "Predict Genres" button, anticipation builds. The interface transforms into a captivating **reveal sequence**. Predicted genres can appear dynamically, rolling in one-by-one or fading in simultaneously. Colors associated with each genre can highlight the results, creating a visually engaging display.

Beyond Names: Probability Unveiled:

The interface goes beyond simply listing predicted genres. It utilizes a **dynamic bar chart** to unveil the **probability distribution** of each genre. Imagine the bar chart unfolding horizontally, revealing the level of confidence the model assigns to each prediction. The use of clear, contrasting colors within the bars allows users to readily grasp the relative probabilities.

Functionalities:

- **Movie Synopsis Input:** Users can enter a movie synopsis directly into a spacious text area.
- **Genre Prediction:** Upon clicking a button, the NLP model analyzes the user-provided synopsis and predicts the top genre(s) for the movie.
- **Real-time Processing:** The prediction occurs in real-time, providing users with immediate results.
- **Probability Distribution Visualization:** The predicted genres are accompanied by a bar chart visualization, showcasing the probability distribution for each genre. This helps users understand the model's confidence level in each prediction.

Features:

- **Intuitive Text Input:** The text area for synopsis input is designed to be user-friendly and visually appealing, encouraging users to engage with the platform.
- **Interactive Button:** A clear and prominent button ("Predict Genres") triggers the prediction process.
- **Concise Genre Display:** The predicted genres are presented in a clear and concise format, making them easy for users to understand.
- **Visually Appealing Probability Chart:** The bar chart utilizes color and clear labeling to effectively communicate the probability distribution of each genre.

8.2 Deployment:

To make our movie genre prediction model accessible and encourage user exploration, we deployed an interactive platform using Streamlit. This deployment strategy, focused on democratized access, allows researchers, movie enthusiasts, and anyone curious about NLP-powered genre prediction to interact with the model. Streamlit's user-friendly framework enabled the creation of an intuitive web application without extensive web development expertise. The deployed app offers real-time genre predictions based on user-entered movie synopses, along with probability visualizations for deeper understanding. This Streamlit deployment fosters wider dissemination of our research findings and empowers users to explore the capabilities of our NLP model.

Movie Genre Prediction App

Enter the movie plot and get the predicted main genres:

Type the movie plot here...

predict Genres

Save Model

Figure-4

Enter the movie plot and get the predicted main genres:

An unknown disease kills many people in a village; however, according to a prediction, someone will come to help them. When one of the villagers meets Raju, he realizes that Raju is the one who can save them.

predict Genres

Predicted Main Genres:

Drama: 14.57%

Action: 12.66%

Thriller: 6.91%

Figure-5

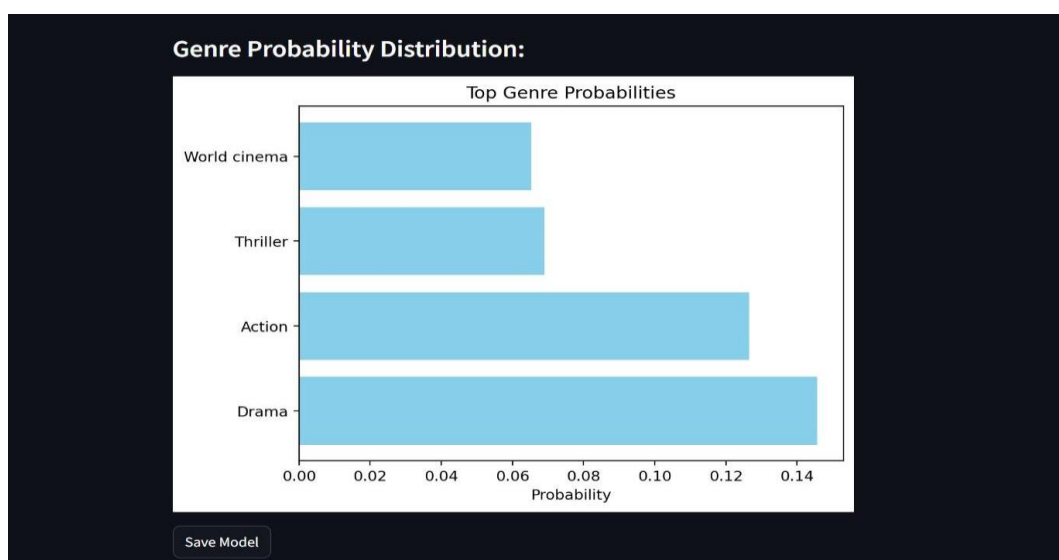


Figure-6

9. Conclusion

In conclusion, this research explored the application of Natural Language Processing (NLP) for movie genre prediction based on movie synopses. We developed a machine learning model capable of analyzing textual descriptions and identifying the most likely genres for a given movie. Our key findings demonstrate the effectiveness of NLP techniques in tackling movie genre prediction. The deployed Streamlit application facilitates user interaction with the model, allowing researchers, movie enthusiasts, and anyone interested in NLP to explore its capabilities. The significance of this work lies in its contribution to the field of NLP applications and movie genre classification. The interactive platform promotes accessibility and encourages further exploration of NLP models for genre prediction. Future research directions include investigating the impact of incorporating additional text data (e.g., reviews, trailers) on prediction accuracy. Additionally, exploring alternative user interface elements within the Streamlit app could enhance user experience and offer deeper insights into genre prediction reasoning. This research paves the way for further advancements in NLP-based genre prediction and interactive platforms for user exploration. By fostering accessibility and user engagement, this work opens doors to broader applications of NLP in movie analysis and beyond.

10. References

- “Movie Genre Prediction Using Multi Label Classification” - Analytics Vidhya1.
- Gupta, K. “Predicting Movie Genres Based on Plot Summaries” - Medium2.
- “Predicting Movie Genres Based on Plot Summaries” - arXiv.org3.
- “Movie Genre Prediction from Plot Summaries by Comparing Various ...” - IRJET4.
- “Predicting Movie Genres Based on Plot Summaries” - arXiv Vanity5.
- Multi-label movie genre detection from a movie poster using knowledge transfer learningK Kundalia, Y Patel, M Shah - Augmented Human Research, 2020 - Springer
- Movie genre classification from plot summaries using bidirectional LSTM AM Ertugrul, P Karagoz - 2018 IEEE 12th International ..., 2018 - ieeexplore.ieee.org
- A multimodal approach for multi-label movie genre classification RB Mangolin, RM Pereira, AS Britto Jr... - Multimedia Tools and ..., 2022 - Springer
- Movie genre classification: A multi-label approach based on convolutions through time , RC Barros - Applied Soft Computing, 2017 - Elsevier
- Movie genre classification using BERT and LSTM A Alwyn, EJP Pranoto, I Ichsan, K Halim... - AIP Conference ..., 2024 - pubs.aip.org