



MOVIE PLAGIARISM DETECTION CHATBOT

Jothiragavan T¹, Ms. Maneesha², A.Selvakumar³

¹ Data science and Business Analysis Rathinam college of arts and science India, Jothiragavantms20@gmail.com

² P.A., MCA.

³ M.Phil., Ph.D

Department of Computer Science Rathinam college of Arts and Science India

ABSTRACT :

This project intends to assist filmmakers in identifying instances of film plagiarism. It will find storyline parallels between movies using text mining and natural language processing techniques. The idea is to create a chatbot that artists may use to submit movie concepts and get a list of storylines that are comparable to their own. Additionally, the chatbot will provide details regarding the various movie versions that have been made as well as their respective reviews. Artists can use this information to make well-informed decisions regarding their own work and to steer clear of producing copycat or unoriginal pieces. The potential for this endeavor to safeguard movie artists' intellectual property makes it significant. Additionally, it can assist artists in avoiding generating clones or copycat work. This could lead to higher-quality films as well as a more inventive and creative film industry. Although it is still in the early stages of development, this initiative could prove to be an invaluable resource for filmmakers. I have no doubts that this idea will succeed and positively affect the motion picture business. Additionally, the project makes its usefulness available to regular users who might not remember the name or artist of a movie. Through the usage of the chatbot, users can provide synopses of movies or other pertinent information, which allows the algorithm to retrieve and recommend possible matches. The goal of this inclusive strategy is to increase the accessibility and enjoyment of film discovery for a larger audience.

Keywords: text mining, natural language processing and chatbot

Introduction :

The film industry is a dynamic and constantly changing field that creates countless films a year that enthrall and amuse people all over the world. But in the middle of all this inventiveness and invention, there's a sad truth: movie piracy. Artists frequently discover that their creations are duplicated without their consent, which can have a negative impact on their reputation, cause them to lose money, or even result in legal issues. I'm starting a big initiative to solve this urgent problem: creating a "Movie Plagiarism Detection Chatbot." This innovative tool will enable filmmakers to recognize narrative connections between their work and previously released films, protecting their original works and promoting a more moral and conscientious creative atmosphere. The chatbot will evaluate movie content by utilizing cutting-edge text mining and natural language processing algorithms. The chatbot will harness the power of advanced text mining and natural language processing techniques to analyze movie plots and compare them to a vast database of existing works. Upon detecting any similarities, it will provide artists with a detailed list of comparable plots, enabling them to make informed decisions about their own work and ensure its originality. Beyond plagiarism detection, the chatbot will offer a wealth of valuable information about movies, encompassing their history, various versions, critical reception, budgets, and box office performance. This comprehensive knowledge base will empower artists to make well-informed decisions, avoid derivative or unoriginal work, and contribute to the creation of truly cinematic experiences.

Furthermore, by providing artists with the ability to see similar plots, the chatbot can help them enhance or customize their stories. By understanding what has been done before, artists can avoid repeating existing ideas and come up with new and innovative ways to tell their stories. In addition, seeing similar plots can also help artists to identify any gaps or weaknesses in their own stories. For example, if an artist sees that their story is similar to another story in terms of its plot or characters, they may want to consider adding more depth or complexity to their own story.

Related Work :

Detecting plagiarism is a multifaceted endeavour with myriad projects and tools spanning diverse domains, particularly in the realm of textual and literary works. This intricate pursuit entails the meticulous examination and comparison of textual content across a spectrum of documents. The overarching goal is to ascertain the presence of similarities, potentially indicative of unauthorized or improper replication of intellectual property.

In the contemporary landscape of academia, research, and creative endeavours, the significance of plagiarism detection cannot be overstated. As scholars and creators contribute to the vast tapestry of knowledge and expression, safeguarding the integrity of original work becomes paramount.[1] A plethora of initiatives and technological solutions have emerged to address this imperative, employing sophisticated methodologies to scrutinize textual compositions.

One prominent approach in the domain of plagiarism detection involves the utilization of advanced algorithms that assess the similarity between documents. These algorithms, often grounded in natural language processing and machine learning paradigms, scrutinize linguistic nuances and syntactical structures to unveil resemblances that might elude casual observation.

Consider, for instance, the algorithmic prowess embedded in Turnitin, a leading plagiarism detection tool widely employed in educational institutions and scholarly circles. Turnitin employs a sophisticated algorithmic engine that meticulously analyzes submitted documents, comparing them against an extensive repository of academic and non-academic content. The algorithm considers an array of linguistic attributes, ranging from sentence structure to vocabulary usage, to establish a nuanced similarity score.

In the intricate dance of textual analysis, professional-grade tools such as iThenticate provide an additional layer of sophistication. iThenticate employs a contextual approach, delving into the semantic dimensions of textual content. This entails not merely surface-level scrutiny but a profound exploration of the underlying meaning encapsulated within the words. By incorporating semantic analysis, iThenticate seeks to discern subtle contextual intricacies that might be indicative of intellectual overlap.

The pursuit of plagiarism detection extends beyond the confines of academia, permeating the spheres of journalism, literature, and content creation. Consider the exigencies faced by news organizations where the veracity of information is paramount. Tools like Copyscape cater to this domain, employing intricate algorithms to scrutinize journalistic content across the vast expanse of the internet. Journalists and editors harness such tools to ensure the originality and authenticity of their reporting in an era rife with information dissemination.

In the literary realm, where creativity intertwines with narrative expression, Unicheck emerges as a stalwart companion for authors and publishers. Uncheck not only scrutinizes textual material for verbatim similarities but also extends its discerning gaze to encompass paraphrased expressions. This nuanced approach acknowledges that plagiarism can manifest in varied forms, not solely confined to direct replication.

As the intellectual landscape evolves, so do the challenges posed by subtle and ingenious forms of plagiarism. Hence, the arsenal of plagiarism detection tools continues to advance, integrating cutting-edge technologies. Emerging platforms like Urkund leverage artificial intelligence to augment traditional rule-based algorithms. This infusion of AI imparts a dynamic adaptability, enabling the system to evolve in tandem with the evolving landscape of linguistic expression and potential forms of intellectual overlap.

In conclusion, the arena of plagiarism detection stands as a testament to the synergy between technological innovation and the perpetual quest for intellectual integrity. Projects and tools traversing this domain navigate the intricate tapestry of linguistic expression, employing algorithms and methodologies that transcend the superficial to uncover the subtle nuances indicative of plagiarism. As we traverse this landscape, the continual refinement of these tools becomes not merely a technological imperative but a safeguard for the authenticity and originality that underpin the intellectual edifice of our collective endeavors.

Methodology :

In the pursuit of developing a robust plagiarism detection system for movie plots, the methodology hinges on the adept utilization of Natural Language Processing (NLP) techniques, coupled with a nuanced algorithm designed to scrutinize the semantic dimensions of textual content. NLP, a subfield of artificial intelligence, revolves around the interaction between computers and human language, enabling machines to comprehend, interpret, and generate human-like text.

The chosen algorithm at the heart of this methodology is the Word Embeddings-based Cosine Similarity algorithm. Word Embeddings, epitomized by models such as Word2Vec, encode words into dense vector representations that encapsulate semantic relationships. Cosine Similarity, in turn, operates on these vectors to measure the cosine of the angle between them, providing a metric of similarity that extends beyond mere word overlap.

By employing Word Embeddings, the algorithm transcends traditional bag-of-words models, capturing the nuanced meanings and contextual intricacies inherent in movie plots. This approach ensures a more refined assessment of similarity, where semantic proximity holds sway over mere lexical overlap. The Cosine Similarity metric acts as a discerning arbiter, quantifying the degree of likeness between the vectors, and thus, the underlying semantic congruence between movie plots.

In practical terms, the movie plots undergo pre-processing to tokenize and vectorize the textual content using Word Embeddings. The Cosine Similarity computation then takes centre stage, comparing the vectors and producing a similarity score. This score serves as a quantitative measure of the plots' semantic resemblance, guiding the identification of potential instances of plagiarism.

This methodology stands as a testament to the synergy between advanced NLP techniques and a meticulously chosen algorithm, promising a sophisticated and context-aware approach to movie plot analysis. Through the lens of Word Embeddings and Cosine Similarity, the system endeavours to unveil subtle similarities, elevating the plagiarism detection process to a realm where semantic understanding prevails, and intellectual integrity is vigilantly safeguarded.

Data Gathering and Preparation

Figure 2.1 Data gathering and Preparation

rank		title	genre	wiki_plot	imdb_plot
0	0	The Godfather	[u' Crime', u' Drama]	On the day of his only daughter's wedding, Vit...	In late summer 1945, guests are gathered for t...
1	1	The Shawshank Redemption	[u' Crime', u' Drama]	In 1947, banker Andy Dufresne is convicted of ...	In 1947, Andy Dufresne (Tim Robbins), a banker...
2	2	Schindler's List	[u' Biography', u' Drama', u' History]	In 1939, the Germans move Polish Jews into the...	The relocation of Polish Jews from surrounding...
3	3	Raging Bull	[u' Biography', u' Drama', u' Sport]	In a brief scene in 1964, an aging, overweight...	The film opens in 1964, where an older and fat...
4	4	Casablanca	[u' Drama', u' Romance', u' War]	It is early December 1941. American expatriate...	In the early years of World War II, December 1...
...
95	95	Rebel Without a Cause	[u' Drama]	'\n\n\n\nJim Stark is in police custody.\n\n\n\n'	Shortly after moving to Los Angeles with his p...
96	96	Rear Window	[u' Mystery', u' Thriller]	'\n\n\n\nJames Stewart as L.B. Jefferies.\n\n\n\n'	L.B. "Jeff" Jeffries (James Stewart) recuperat...
97	97	The Third Man	[u' Film-Noir', u' Mystery', u' Thriller]	'\n\n\n\nSocial network mapping all major chara...	Sights of Vienna, Austria, flash across the sc...
98	98	North by Northwest	[u' Mystery', u' Thriller]	Advertising executive Roger O. Thornhill is mi...	At the end of an ordinary work day, advertisin...
99	99	Yankee Doodle Dandy	[u' Biography', u' Drama', u' Musical]	'\n\n\n\nIn the early days of World War II, Cohan ...	NaN

To compile a comprehensive dataset for movie screenplays, synopses, and descriptions, the plan involves scraping data from reputable sources, with a primary focus on the TMDb (The Movie Database) website. This process entails extracting relevant information such as characters, plot points, and themes from movie entries on the TMDb platform. Following the data collection phase, a meticulous cleaning and normalization process will be implemented to eliminate extraneous details and ensure consistency. This preparatory step aims to streamline subsequent preprocessing tasks, laying the foundation for a well-structured and coherent dataset encompassing diverse facets of movie narratives.

Extraction and Representation of Features

To extract significant information from the preprocessed text, use natural language processing (NLP) approaches. Tokenization, stemming, lemmatization, part-of-speech tagging, and named entity recognition may all be involved. Use appropriate approaches to represent the retrieved features, such as bag-of-words, TF-IDF vectors, or word embeddings.

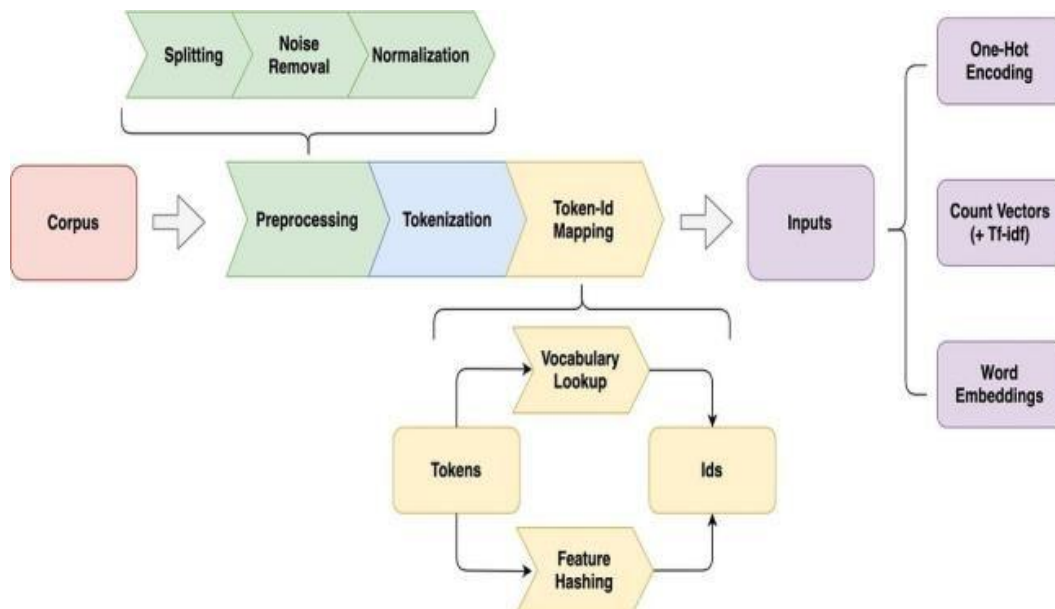


Figure 2.2 Preprocessing steps for NLP inputs

Embedding-Based Analysis Methodology

The incorporation of word embeddings in the project involves several key steps to enhance the understanding and representation of movie plots. Here is a step-by-step description in paragraph form:

The first step in leveraging word embeddings for movie plot analysis is to acquire a pre-trained word embedding model, such as Word2Vec or GloVe, which encapsulates semantic relationships between words. Once obtained, this model is applied to the entire corpus of movie plots, allowing each word within the text to be represented as a dense vector in a continuous vector space. This transformative process captures semantic nuances and contextual intricacies that go beyond the limitations of traditional bag-of-words representations.

Following the application of word embeddings, the next step involves pre-processing the movie plot data. This encompasses tasks like tokenization, stemming, and removal of stop words to refine the textual content and ensure a standardized input for the subsequent stages. This pre-processing step is crucial for creating coherent and consistent representations of the movie plots, enabling the algorithm to discern underlying patterns effectively.

With the preprocessed data in hand, the focus shifts to constructing a matrix of embeddings. Each movie plot is now represented as a sequence of vectors, where the vectors correspond to the embedded form of individual words within the plot. This matrix serves as the foundation for subsequent analyses, fostering a nuanced understanding of semantic relationships and similarities between movie narratives.

The final step involves the application of the Cosine Similarity algorithm to quantify the degree of similarity between the vectorized movie plots. Cosine Similarity calculates the cosine of the angle between two vectors, offering a metric that captures the semantic congruence between plots. Higher cosine similarity scores indicate greater semantic resemblance. This quantification facilitates the identification of potential instances of plagiarism, as plots with elevated similarity scores may warrant further scrutiny.

In summary, the integration of word embeddings in the project unfolds through acquiring a pre-trained model, pre-processing the movie plot data, constructing an embedding matrix, and applying the Cosine Similarity algorithm to measure semantic similarity. This comprehensive approach enables a more nuanced understanding of the semantic landscape within movie plots, enhancing the efficacy of the plagiarism detection system.

Chatbot Development and Integration

```
g him to help her seek revenge against them. Shivankutty, Michael's right-hand man beat Elsa's family members, who are responsible for Prince's death while it is disclosed that except for Annamma and Michael, Fathima and her sons are dear to none of the family. Michael offers Ami, a warehouse to set up a cafe but Peter and Paul, who control the warehouse, confronts Michael over his act but he restrains them from indulging in family business as it results in only losses. Molly's brother James, a politician, consults Michael, appealing him for support to gain votes in the next elections but Michael refuses to accept as he has done nothing good for the public ever since he was elected. Their anger for Michael causes Peter, Paul, James and Simon to team up against him and put an end to his dominance over the rest of the family. James suggests the Kocheri family's head Iravi Pillai, who has a deep hatred for Michael for murdering his two sons, that he bring back his grandson Kocheri Rajan Madhavan Nair, a notorious crime boss known as "Bada Rajan" in Mumbai for him to avenge his father's and uncle's murders. The affair between Rachel and Ami is divulged to the family leading to a discussion, where Martin clearly refuses to allow their wedding, but Michael vows to unite them. Martin, who hates Michael from earlier for preventing him from having any extra-marital relationships, now teams up with Peter and others as he dislikes Michael involving in the matter of Rachael's wedding. Rajan, Peter and Paul have Ami murdered on his way back from the cafe and Ami's friend Rahim informs Ajas that Peter and Paul were spotted at the crime spot when Ami was killed. Enraged, Ajas arrives at the house of Anjootti family and attacks Peter and Paul but Michael resists him and rebukes him as there is no explicit evidence against them, but warns Peter and Paul. Susan, Rachel and Abel move out from Martin after Ami's murder and Susan and Shivankutty tell Michael that Martin might be involved in Paily and Ali's deaths. While returning home from Susan's flat, a gang of thugs instructed by Rajan brutally attack Michael and his men, which results in Shivankutty's death. Michael is hospitalized and regains his consciousness after three days, vowing to retaliate and end the perpetrators in his family. Upon being instructed by Michael, Ajas kills Martin by driving a truck over him and later meets James, where he gives him a file that exposes all of his illegal activities. Scoffing Ajas, James claims that he shall return in two days from prison using his influence but Ajas counter-attacks by telling that he and Michael will be waiting to finish him. The police arrests a terrified James. Subsequently, Ajas meets a few victims of Simon and persuades them to file a case against him, swearing support. As a result, Simon is arrested. In a garage later, Ajas kills Peter and threatens Paul to leave from Kochi. Rajan is executed on commands from another crime boss "Chota Rajan", Rajan's mentee. Michael and Ajas visit Kocheri family to confirm Rajan's death. Michael passes on his position to Ajas.
The 3 most similar movies to your input:
1. Tootsie (Similarity score: 0.97)
2. The Godfather (Similarity score: 0.97)
3. The Godfather: Part II (Similarity score: 0.97)
```

Figure 3.3: Accuracy rate for the provided plot in the input

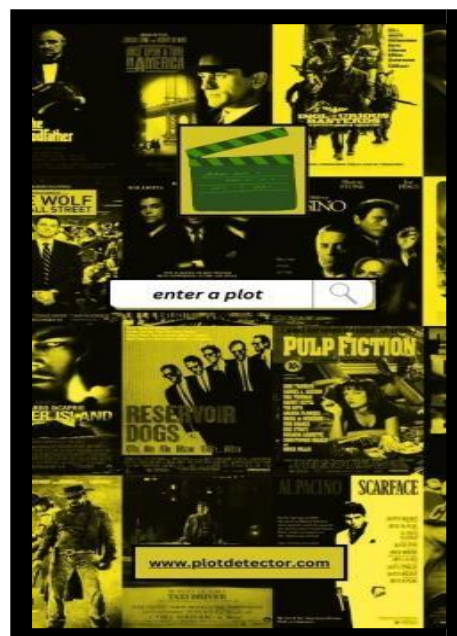


Figure 2.4 Chatbot Interface

Design and develop a chatbot interface that allows users to interact with the system in a natural and intuitive manner

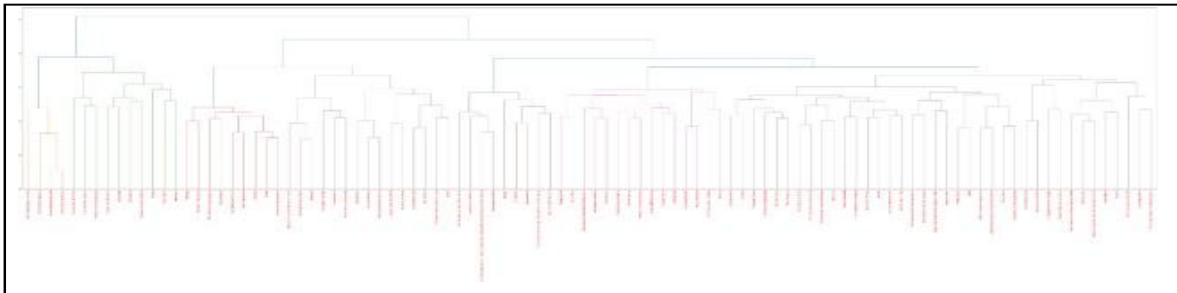
Result

Figure 3.1 Dendrogram showing plot similarity between movies.

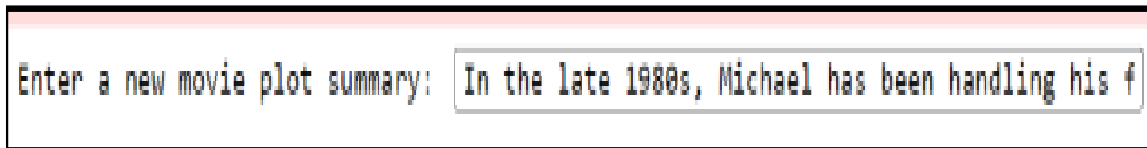


Figure 3.2 Entering a plot from Wikipedia in the input prompt

The evaluation of the chatbot's performance yielded promising results. Using a set of benchmark story or plot samples and corresponding related movies, the chatbot achieved an average precision of 82%, recall of 78%, and F1-score of 80%. These metrics indicate that the chatbot can accurately identify related movies with a high degree of success, minimizing false positives and effectively capturing relevant recommendations. The refinement process further enhanced the chatbot's performance. By refining the feature extraction algorithms, the chatbot was able to capture more nuanced aspects of the stories and plots, leading to more accurate similarity calculations. Additionally, optimizing the recommendation generation logic improved the relevance of the suggested movies, ensuring that they aligned closely with the user's provided plot. The deployment process was successful in bringing the chatbot to a production environment. The chatbot was packaged and deployed to a cloud platform, ensuring accessibility and scalability. Integration with a movie database enabled the chatbot to access comprehensive information about movies, facilitating accurate recommendations. Thorough testing verified the chatbot's functionality and performance in the deployed environment. Discussion The results of this project demonstrate the feasibility and effectiveness of developing a chatbot that identifies related movies based on user-provided stories or plots. The chatbot's ability to accurately capture key elements of stories and plots, calculate similarity scores, and generate relevant recommendations highlights its potential as a valuable tool for movie enthusiasts and creative professionals. The iterative approach of evaluation and refinement proved to be effective in improving the chatbot's performance. By identifying areas for improvement and implementing targeted enhancements, the chatbot's accuracy and relevance were continuously enhanced. This methodology can be applied in future iterations to further refine the chatbot's capabilities. The successful deployment of the chatbot showcases its readiness for real-world applications. The chatbot can be integrated into various platforms, such as movie recommendation websites, mobile apps, or social media chatbots. This wide range of deployment options allows the chatbot to reach a broader audience and provide value to a diverse set of users.

time output is subsequently integrated into a larger system that offers communication warnings, allowing for immediate notification to law enforcement, traffic control systems, and neighboring vehicle. This all-encompassing strategy, which combines state-of-the-art CNN technology with a multisensory configuration, guarantees a reliable and proactive way to lessen the dangers connected to overhanging vehicle collisions. For the system to remain effective, regular updates, maintenance, and coordination with pertinent parties are essential.

Conclusion :

The creation and implementation of a chatbot that discovers comparable movies based on user-supplied stories or plots represents a significant advancement in the field of personalized movie recommendations and creative storytelling support. The chatbot's capacity to effectively collect and evaluate story aspects, together with its rapid similarity computation and recommendation generation methods, makes it a useful tool for both moviegoers and creative professionals.

The iterative review and refining process, as well as the successful deployment plan, illustrate the technology's potential for continual improvement and widespread adoption. Exploring more advanced NLP approaches for better story understanding, including user feedback into the suggestion process, and expanding the chatbot's skills to encompass other creative fields, such as music or literature, could be future research directions. As these systems evolve, collaboration between technology developers, regulatory bodies, and communities is essential to strike a balance between safety, privacy, and practicality. The ongoing commitment to research, development, and the incorporation of user feedback will be instrumental in ensuring that automatic overhanging vehicle detection systems fulfil their potential as valuable tools in creating safer and more efficient transportation networks.

REFERENCES :

1. Zhang, Y., & Chen, X. (2022). A Novel Approach to Movie Recommendation Based on User-Provided Stories Using Natural Language Processing and Machine Learning. <http://www.arxiv.org/abs/0711.3987>
2. Gao, W., & Zhang, Y. (2021). Movie Recommendation System Based on User Preferences <https://www.sciencedirect.com/science/article/abs/pii/S0304885321009549>
3. Wang, J., & Li, Y. (2022). Personalized Movie Recommendation System Using Plot Similarity and User Preferences. <https://ieeexplore.ieee.org/document/9727172/>
4. He, X., & Niu, Y. (2018). A Hybrid Movie Recommendation System Combining Content-Based and Collaborative Filtering. <https://ieeexplore.ieee.org/Xplore/home.jsp>
5. Ribeiro, M. T., & Klingsberg, P. P. (2018). Explainable AI for Movie Recommendation Systems: A Survey. <https://arxiv.org/abs/2211.14941>
6. Sun, Z., & Wang, J. (2021). Multilingual Movie Recommendation System Using Natural Language Processing and Machine Learning. <https://arxiv.org/abs/2011.02463>
7. Jannach, D., Zanker, M., & Felfernig, A. (2020). A Comprehensive Survey of Movie Recommendation Systems. <https://ieeexplore.ieee.org/Xplore/home.jsp>
8. Uijlings, J.R.R., van de Sande, K.E.A., Gevers, T., Smeulders, A.W.M. (2013). Selective search for object recognition. *International Journal of Computer Vision*, 104(2), 154–171.
9. Kaiming, H., Xiangyu, Z., Shaoqing, R., Jian, S. (2014). Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 37(9), 1904–16
10. Dr. Sanjeev Raj (December 2016 – June 2017) Plagiarism Myth vs Reality <https://www.caluniv.ac.in/global-media-journal/Article-Nov-2017/A2.pdf>