



Table of contents



(1) Introduction

Background, problems statements, objectives, literature review



Naive Bayes, SVM, Random Forest, and Decision Tree

1) 2 Data Understanding

Data collection, exploration, quality, web-scraping

1) Evaluation

Visualize model performance, comparison and sentiment

1) 3 Data Preparation

Data cleaning and sentiment extraction

1) Conclusion

Summarize: Findings and implications







Barbie (2023)

Director: Greta Gerwig

Writers: Greta Gerwig & Noah Baumbach

• Main Casts: Margot Robbie, Ryan Gosling, Issa Rae

• Genres: Adventure, Comedy, Fantasy

Overall Review: 6.9 / 10 (IMDb)

Awards: 106 wins & 325 nominations

'She's everything. He's just Ken'



Sentiment Analysis

23

Also known as opinion mining, a natural language processing techniques that involves determine the sentiment or emotional tone expressed in a piece of text

Deminon Application feed under on d

Widely used in various industries and applications such as business for feedback, social media to understand public opinion on different topics

Complexity of language such as sarcasm, irony and ambiguity. Sentiment can vary based on individual perspectives.

Thallenges & Model

Rule-based systems to machine learning and deep learning approaches. Rule-based rely on predefined rules and dictionaries while ML trained on labeled dataset.

Problem Statement

Ambiguity and Contextual

- Face difficulties with uncertain language, sarcasm and contextual complexities particularly in movie reviews.
- Need of model that effectively understand and interprets the intricate aspects of language in movie critiques.

Performance of the models

- Some of the models having difficulties to classify the result accurately as it may struggle to understand the language nuance.
- Performance evaluation of the models must be done to determine the best model to use for sentiment analysis for movie critiques.

Objectives



To develop and train a model for movie reviews that uses advanced techniques to analyze sentiment and comprehend complex language.

Sentiment Classification

To determine whether the majority of reviews express positive or negative sentiments toward Barbie movies.

Evaluate and Compare

To evaluate and compare different models of sentiment analysis for movie reviews using accuracy, recall and precision.



Literature Review

Literature Review	Topic	Citation
Research has determined that reviews that present both sides of an argument are more beneficial to users than reviews that contain a bias.	Domain	(Kong et al., 2023)
Data obtained from a representative set of online reviews penned by IMDb users residing in the United States, as well as Douban users based in China, were consolidated.	Dataset	(Lou et al., 2023)
Users utilize platforms such as social networks, review websites, forums, and blogs to express their views on products and services. Analyzing user expressions can provide valuable patterns for decision making.	Dataset	(Saraswathi et al., 2023)
The Sentiment Analysis technique extracts relevant information from User Reviews dataset and classifies them into positive and negative comments for decision-making.	Technique	(Rooba et al., 2023)

Literature Review

Literature Review	Topic	Citation
Understanding the author's opinion and the user	Technique	(Pavitha et al., 2023)
experience is valuable. Opinion mining involves		
extracting and categorizing opinions from online forums		
or platforms. This allows for a deeper understanding of		
user sentiment towards a specific topic.		
The linguistic barrier hinders sentimental analysis. Only	Techniques (Limitation)	(Raut et al., 2023)
English reviews can be analyzed. Sarcastic or ironic reviews lead to incorrect classification.		
Research on sentiment analysis of Indian regional	Dataset (Limitation)	(Mohan et al., 2023)
language texts is impeded by the absence of regional		05
language datasets.		



Data Collection



Source of Data: IMDb



Dataset:



Link: https://www.imdb.com/title/tt1517268/reviews



No. of reviews: 1,471 Attributes:

- Rating
- User
- Date
- Headline
- Description



Has **special characters**. Has **missing values**.

Web-Scraping

```
1 from urllib.request import urlopen as uReq
2 from bs4 import BeautifulSoup as soup
4 import requests
6 def get reviews(url):
     headers = {
          'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML,
     reviews = []
         response = requests.get(url, headers=headers)
         page html = response.text
         page_soup = soup(page_html, "html.parser")
         for item in page soup.find all("div", class ="lister-item-content"):
             rating = item.find("span", class ="rating-other-user-rating").text
              user = item.find("span", class = "display-name-link").text
              date = item.find("span", class ="review-date").text
              title = item.find("a", class ="title").text.strip()
             description = item.find("div", class = "text show-more_control").text.strip()
              reviews.append([rating, user, date, title, description])
         load more button = page soup.find("div", class ="load-more-data")
         if load more button:
             data_key = load_more_button['data-key']
     return reviews
```





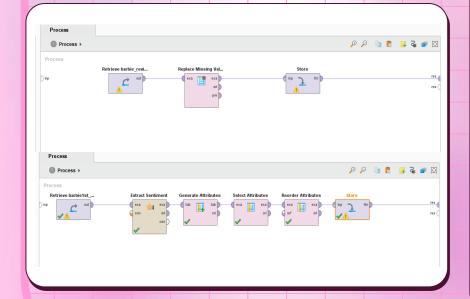








Data Preparation: Cleaning and Extract Sentiment





The Operator

- Replace Missing Value: Remove missing value
- **Extract Sentiment**: Using Vader technique to analyze the Description
- Generate Attributes: Use IF-Else to create "Sentiment" column for "Positive" and "Negative" attributes.
- **Select Attributes**: Select only the important attributes. (Sentiment and Description)
- **Reorder Attributes**: Reorder the attributes position
- **Store**: Store the data into the repository.



Before

After





Row No.	Rating	User	Headline	Description
1	1/10	BA_Harrison	Not for me.	Barbie (Marg
2	1/10	kobalplatt	Did we all se	I was literally
3	1/10	nialImcenteg	Marketing co	A word of advi
4	1/10	STARONTHE	Overrated	This movie is
5	1/10	thedarkknight	Atrocious	Don't know w
6	1/10	yurixander10	Bad movie.	I bought it on
7	1/10	protsenkodv	Barbie Movie	Oh boy, wher
8	1/10	evisadh	Rubbish	Watched it at
9	1/10	d-s-allen82	Terrible	This movie w
10	1/10	ericvankeimp	Worst movie	We saw the
11	1/10	owenjmiddlet	This years fil	When I went t
12	1/10	chriswalls-08	Don't believe	Rarely have I
13	1/10	mariekeherwi	Very weak.	I would not re
14	1/10	pensacolaco	Really really	I am at a com
15	1/10	eva-rebac	So bad	I really wasn't
16	1/10	li0904426	Total Disaste	The movie "B
17	1/10	bdewiyah	BAD	I am just ama
18	1/10	amandakelly	Rubbish mov	The only goo
19	1/10	hamedkhaza	Failed to be f	It was one of t
20	1/10	mr_asdfg	Terribble	The Barbie m



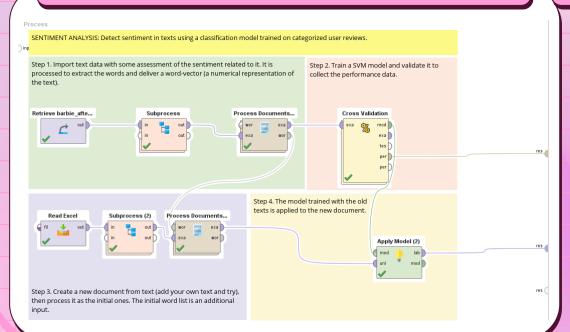






Data modelling

Main Process Structure

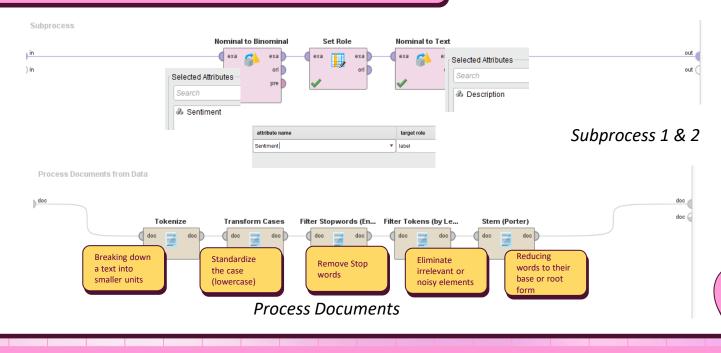


Model

- SVM
- Naïve Bayes
- **Random Forest**
- **Decision Tree**

Data modelling

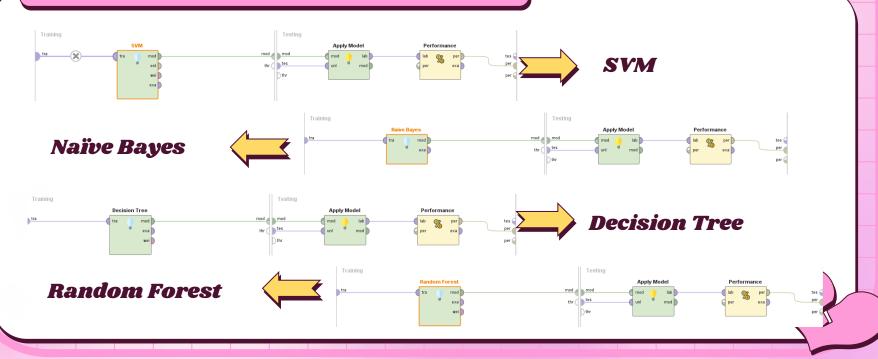
Sub Process Structure



Data modelling



Cross Validation Structure



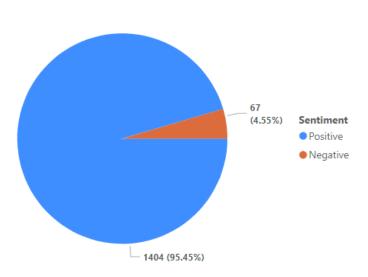


Visualize model performance, comparison and sentiment



Results







SVM

accuracy: 95.45% +/- 0.33% (micro average: 95.45%)

	true Positive	true Negative	class precision
pred. Positive	1404	67	95.45%
pred. Negative	0	0	0.00%
class recall	100.00%	0.00%	

Naïve Bayes

accuracy: 95.17% +/- 0.60% (micro average: 95.17%)

	true Positive	true Negative	class precision
pred. Positive	1400	67	95.43%
pred. Negative	4	0	0.00%
class recall	99.72%	0.00%	



Random Forest

accuracy: 95.45% +/- 0.33% (micro average: 95.45%)

	true Positive	true Negative	class precision
pred. Positive	1404	67	95.45%
pred. Negative	0	0	0.00%
class recall	100.00%	0.00%	

Decision Tree

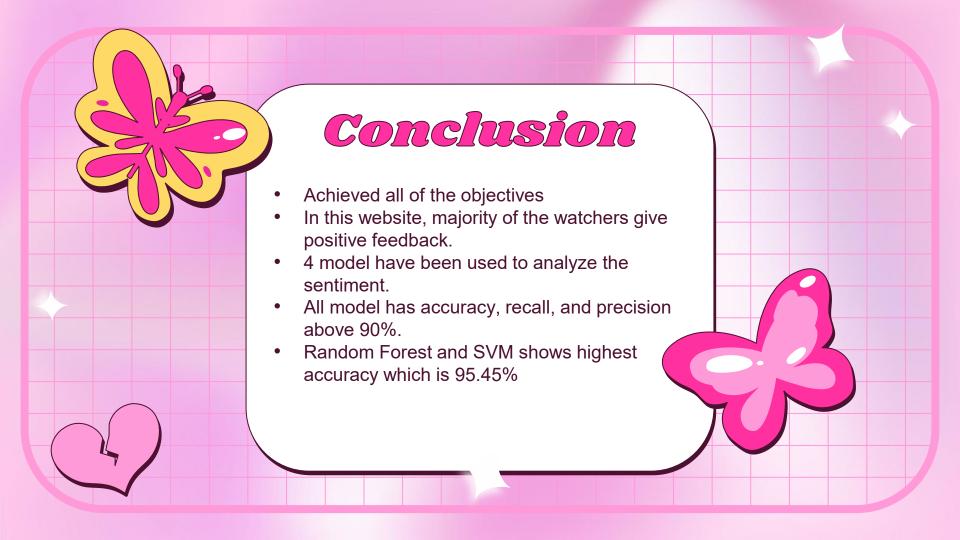
accuracy: 94.36% +/- 0.91% (micro average: 94.36%)

	true Positive	true Negative	class precision
pred. Positive	1383	62	95.71%
pred. Negative	21	5	19.23%
class recall	98.50%	7.46%	

Results







References

N Pavithaa, Vithika Pungliyaa, Ankur Raut, Roshita Bhonsle, Atharva Purohit, Aayushi Patel, R Shashidhar. (2023). *Movie recommendation and sentiment analysis using machine learning*. https://doi.org/10.1016/j.gltp.2022.03.012

N. Saraswathi, T. Sasi Rooba, S. Chakaravarthi. (2023). *Improving the accuracy of sentiment analysis using a linguistic rule-based feature selection method in tourism reviews*. https://doi.org/10.1016/j.measen.2023.100888

Juan Kong, Chen Lou. (2023). *Do cultural orientations moderate the effect of online review features on review helpfulness? A case study of online movie reviews.* https://doi.org/10.1016/j.jretconser.2023.103374

Syam Mohan E , R. Sunitha. (2023). MABSA: A curated Malayalam aspect based sentiment analysis dataset on movie reviews.

https://doi.org/10.1016/j.dib.2023.109452







Thanks!

Do you have any questions?



