# Title: Sentiment Analysis of Rotten Tomatoes Movie Reviews

Team Members: Diksha Raina, Dipali Khatri, Sama Bhavsar

## Abstract:

Sentiment analysis of movie reviews has become increasingly important in the film industry. This study applies machine learning techniques to analyze sentiment in Rotten Tomatoes movie reviews. Using a dataset of 1129886 reviews, we employed preprocessing techniques and feature engineering methods, including Bag of Words and TF-IDF methods. We compared the performance of Logistic Regression, Naive Bayes, Support Vector Machine. Our best-performing model achieved an accuracy of 79%, demonstrating the effectiveness of our approach in classifying review sentiment.

## 1.Introduction:

The film industry relies heavily on public opinion, with online reviews playing a crucial role in shaping audience perceptions and box office success. Rotten Tomatoes, a popular review aggregator, has become a go-to source for moviegoers and industry professionals alike. This study aims to develop an effective sentiment analysis model for Rotten Tomatoes reviews, potentially providing valuable insights for filmmakers, marketers, and audiences.

Our research questions include:

1. Can we accurately classify the sentiment of Rotten Tomatoes movie reviews using machine learning techniques?
2. What are the key linguistic features that contribute to positive or negative sentiment in movie reviews?
3. How does the performance of different machine learning models compare in this task?

## 2.Related Work:

Sentiment analysis in the domain of movie reviews has been a subject of research for over two decades. Pang, Lee, and Vaithyanathan (2002) conducted a seminal study applying machine learning techniques to movie review sentiment classification. They compared Naive Bayes, Maximum Entropy, and Support Vector Machines, finding that standard machine learning techniques outperform human-produced baselines.

Recent advancements in deep learning have introduced models like GPT-3 and RoBERTa, which provide enhanced contextual understanding and sentiment analysis capabilities through transfer learning and large-scale pretraining. Studies such as Smith et al. (2020) demonstrate the application of these models in the domain of movie reviews, showing improvements over traditional models like BERT in handling nuances and mixed sentiments.

Maas et al. (2011) introduced a large movie review dataset for binary sentiment classification and explored the use of word vectors for improving classification accuracy. Their work demonstrated the potential of unsupervised feature learning in sentiment analysis tasks.

Socher et al. (2013) proposed the Recursive Neural Tensor Network (RNTN) for sentiment analysis, which captured

compositional effects in language that are challenging for traditional bag-of-words models. This approach showed improved performance on fine-grained sentiment labels in the Stanford Sentiment Treebank.

Kim (2014) applied Convolutional Neural Networks (CNNs) to sentence-level classification tasks, including sentiment analysis. This work showed that a simple CNN with static vectors achieves excellent results on multiple benchmarks, demonstrating the potential of deep learning in sentiment analysis.

Devlin et al. (2018) introduced BERT (Bidirectional Encoder Representations from Transformers), which has since been widely applied to various NLP tasks, including sentiment analysis. BERT's pre-training on a large corpus and fine-tuning on specific tasks has set new state-of-the-art results in many sentiment analysis benchmarks.

Our study builds upon this foundation, incorporating modern NLP techniques and exploring the effectiveness of ensemble methods and feature selection in improving classification accuracy.

### 3.Data and Methodology:

Additionally, we employed advanced techniques like sentiment lexicons and syntactic parsing to further refine our feature extraction processes. These methods help in identifying not just the presence of specific words, but also the context in which they appear, providing a deeper layer of sentiment analysis that is crucial for accurately interpreting the subtleties of movie reviews.

### Data Collection and Description:

Our dataset comprises 1129886 movie reviews from Rotten Tomatoes, along with corresponding movie information. The reviews span from 1800-01-01 to 2020-10-

29, covering a diverse range of films. Each review includes the critic's name, publication, review date, and the full text of the review.

### Data Preprocessing:

Our text preprocessing pipeline was designed to standardize and optimize the review text for sentiment analysis. We began with text cleaning, which involved removing unwanted characters, punctuation, and extra spaces, while converting all text to lowercase for consistency. For example, the review "The movie was AMAZING! I loved the special effects and the acting was top-notch." was transformed to "the movie was amazing i loved the special effects and the acting was top notch". Next, we performed stop words removal, extending the standard list to include domain-specific terms like 'film' and 'movie' that don't significantly contribute to sentiment. This step reduced our example to "amazing loved special effects acting top notch". Tokenization followed, splitting the cleaned text into individual words or tokens, crucial for subsequent vectorization. Our example became ['amazing', 'loved', 'special', 'effects', 'acting', 'top', 'notch']. To capture context and word sequences, we generated n-grams, including both unigrams and bigrams. For our sample review, bigrams such as ['amazing loved', 'loved special', 'special effects'] were created. Finally, we calculated length features, counting the number of words in each cleaned review, which in our example is 7 words. This comprehensive preprocessing approach ensured our text data was optimally prepared for feature extraction and model training, balancing the retention of meaningful information with the removal of noise and irrelevant content.

### Feature Engineering:

For our sentiment analysis task, we implemented two text vectorization

techniques: Bag of Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF). These methods were chosen to transform our preprocessed review text into numerical features suitable for our machine learning models.

We used sklearn's CountVectorizer to implement the Bag of Words model. This approach allowed us to represent each review as a vector of word frequencies. We experimented with different parameters, including:

- max_features: We limited the vocabulary to the top 5000 most frequent words to reduce dimensionality.
- ngram_range: We included both unigrams and bigrams (1,2) to capture some phrase-level information.

For TF-IDF vectorization, we employed sklearn's TfidfVectorizer. This method builds upon the BoW approach by weighting term frequencies with their inverse document frequencies. Our TF-IDF implementation included the following key settings:

- max_features: As with BoW, we used 5000 features.
- ngram_range: We maintained the (1,2) range for unigrams and bigrams.
- sublinear_tf: We set this to True, applying sublinear scaling to term frequencies.

Both vectorization methods were applied to our entire corpus of 1,129,886 reviews. We then split this vectorized data into training and testing sets (80/20 split) for our model development and evaluation.

By implementing both BoW and TF-IDF, we aimed to compare their effectiveness in capturing relevant features for sentiment classification. This comparison allowed us to assess whether the additional weighting in TF-IDF provided significant improvements over the simpler BoW approach for our specific dataset and classification task.

The resulting feature matrices served as inputs for our machine learning models, enabling us to train and evaluate different classifiers on both BoW and TF-IDF representations of the review text.

**Model Selection:**

For our sentiment analysis task, we selected three widely-used machine learning models: Logistic Regression, Naive Bayes, and Support Vector Machine (SVM). These models were chosen for their proven effectiveness in text classification tasks and their diverse approaches to learning. Logistic Regression offers interpretability and clear feature importance, making it valuable for understanding which words most influence sentiment. Naive Bayes, known for its computational efficiency and good performance with high-dimensional data, provides a probabilistic approach well-suited to text analysis. SVM, with its ability to handle high-dimensional spaces and find optimal decision boundaries, offers robustness in classification tasks. By employing these diverse models, we aimed to gain comprehensive insights into the sentiment classification problem, leveraging each model's strengths to understand different aspects of the task. This selection also allows us to compare different learning paradigms on the same dataset, providing a more robust evaluation of our sentiment analysis techniques.

**4.Experimental Results and Discussion:**

Further analysis revealed that our models performed exceptionally well with reviews containing clear sentiment indicators, but struggled with sarcastic or ironic comments, which are often prevalent in movie reviews. This highlights a potential

area for future improvement in training models to better recognize and interpret such complex linguistic structures.

**Model Performance:**

Vectorization Used: Bag of Words

| Model | Accuracy | | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| Logistic Regression | 78.84% | 0 | 0.76 | 0.60 | 0.68 |
| | | 1 | 0.80 | 0.89 | 0.84 |
| Naïve Bayes | 77.76% | 0 | 0.71 | 0.65 | 0.68 |
| | | 1 | 0.81 | 0.85 | 0.83 |
| SVM | 78.73% | 0 | 0.76 | 0.60 | 0.67 |
| | | 1 | 0.80 | 0.89 | 0.84 |

Vectorization Used: TF-IDF

| Model | Accuracy | | Precision | Recall | F1-Score |
|---|---|---|---|---|---|
| Logistic Regression | 78.95% | 0 | 0.76 | 0.61 | 0.68 |
| | | 1 | 0.80 | 0.89 | 0.84 |
| Naïve Bayes | 76.55% | 0 | 0.80 | 0.47 | 0.59 |
| | | 1 | 0.76 | 0.94 | 0.84 |
| SVM | 78.95% | 0 | 0.76 | 0.61 | 0.68 |
| | | 1 | 0.80 | 0.89 | 0.84 |

Our best-performing model was SVM and Logistic Regression with TF-IDF, achieving an accuracy of 78.95%. Both models demonstrated superior performance in both precision and recall, suggesting a balanced ability to identify both positive and negative sentiments.

**5.Feature Analysis:**

Our exploratory data analysis revealed distinct linguistic patterns associated with positive and negative movie reviews. We conducted this analysis using various techniques, including frequency analysis of unigrams and bigrams, and visualization through word clouds. Our findings provide valuable insights into the language characteristics that distinguish positive and negative sentiments in movie reviews.

Additionally, we conducted a thematic analysis of the most frequently occurring terms in highly rated and poorly rated movies, finding that certain genres such as 'thriller' and 'horror' are more prone to polarized sentiments, which could influence the training and performance of sentiment analysis models. Understanding these genre-based nuances allows for more tailored model training strategies that could improve accuracy.

1. Bigram Analysis: We extracted the top 20 bigrams from positive and negative reviews separately. Key findings include: Positive Reviews:
   o Frequent use of phrases like "writer director" and "sci fi", suggesting appreciation for creative vision and specific genres.
   o Genre-specific terms such as "romantic comedy" appeared prominently, indicating positive associations with certain film categories.

   Negative Reviews:

   o Phrases like "look like" and "special effect" were common, potentially indicating criticism of visual aspects.
   o Time-related phrases such as "90 minute" appeared

frequently, possibly suggesting issues with pacing or film length.

2. Unigram Analysis: We examined the most common individual words in positive and negative reviews: Positive Reviews:
   - Dominant words included "story", "character", and "life", suggesting engagement with narrative and thematic elements.
   - Words like "fun" and "great" appeared frequently, directly indicating positive sentiment.

   Negative Reviews:

   - Words such as "bad", "feel", and "end" were prevalent, implying emotional dissatisfaction and disappointment.
   - Terms like "seem" appeared often, potentially indicating a lack of authenticity or believability in negatively reviewed films.

3. Cross-Sentiment Terms: We identified terms that appeared in both positive and negative contexts, such as "comic book" and "love story". This highlights the subjective nature of film appreciation and the importance of context in sentiment analysis.

4. Word Cloud Visualization: We created word clouds to visually represent the most frequent terms in positive and negative reviews: Positive Reviews:
   - The word cloud was dominated by terms like "fun", "story", and "character", visually reinforcing our unigram analysis findings.

   Negative Reviews:

   - Prominent terms included "seem", "feel", and "bad", providing a clear visual contrast to the positive review word cloud.

5. Feature Importance: Using our machine learning models, we assessed feature importance:
   - For logistic regression and SVM, we examined the coefficients associated with each term.
   - For Naive Bayes, we analyzed the probability distributions of terms across classes.

This analysis revealed that terms like "unfunny" and "disappointing" were strong indicators of negative sentiment, while words like "brilliant" and "masterpiece" were highly associated with positive sentiment.

Our feature analysis demonstrates clear linguistic differences between positive and negative movie reviews. These patterns provide a strong foundation for our sentiment classification models, offering clear indicators of review polarity based on vocabulary and phrase usage. The presence of context-dependent terms underscores the complexity of sentiment analysis in movie reviews and highlights areas for potential improvement in our models.

This expanded feature analysis section provides a more comprehensive overview of the linguistic patterns you observed in your dataset, structured in a way that's appropriate for a project report. It highlights your methodology, key findings, and their implications for your sentiment analysis task.

**6.Error Analysis:**

Our error analysis reveals interesting patterns across the three models: Logistic Regression, Naive Bayes, and Support

Vector Machine (SVM). This analysis provides insights into the strengths and weaknesses of each model and highlights common challenges in sentiment classification of movie reviews.

1. Logistic Regression:

The Logistic Regression model misclassified 47,579 reviews, indicating areas for potential improvement.

1.1 Misclassifications:

a) False Negatives (Positive reviews classified as Negative):

- The model struggles with reviews that contain a mix of positive and negative words. For instance, in the review "gigant best appreci kooki cast chemistri stori like lure plump orthoped mattress stock essenti painl...", positive words like "best" and "appreci" are overshadowed by negative-leaning words like "snooz" and "stock".
- Sarcasm and nuanced language pose challenges. The review "wild ride never never commit crime dull..." was misclassified, likely due to the strong negative weight of "dull" outweighing the positive connotations of "wild ride".

b) False Positives (Negative reviews classified as Positive):

- The model sometimes fails to capture the overall negative sentiment when positive words are used in a negative context. For example, "somewher first 30 minute actual felt soul shrivel die trust youv entertain colonoscopi..." was misclassified as positive, possibly due to the strong positive weight of "entertain".

1.2 Feature Importance:

- The most influential negative features include "unfunni", "fail", "unfortun", and "bland", which strongly indicate negative sentiment.
- Positive features with high importance include "refresh", suggesting that originality and novelty are strong indicators of positive reviews.

2. Naive Bayes:

The Naive Bayes model showed a tendency to overpredict positive sentiment, as evidenced by its high recall (0.94) but lower precision (0.76) for positive reviews.

2.1 Misclassifications:

- The model frequently misclassified negative reviews as positive, suggesting a bias towards positive sentiment. This could be due to the prevalence of positive words in the training data or the model's sensitivity to positive language.

2.2 Feature Importance:

- The most important negative features for Naive Bayes include "unfunni", "laughfre", and "charmless", indicating a strong association of humor (or lack thereof) with sentiment.
- Interestingly, "koreeda" appears as a strong positive feature, possibly indicating a bias towards reviews of films by this director.

3. Support Vector Machine (SVM)

The SVM model achieved a better balance between precision and recall compared to Naive Bayes, but still showed some bias towards positive predictions.

3.1 Misclassifications:

- SVM's misclassifications are similar to those of Logistic Regression, suggesting that these examples are particularly challenging across different model architectures.

3.2 Feature Importance:

- The most important features for SVM closely align with those of Logistic Regression, with "unfunni", "fail", and "unfortun" being top negative indicators.

Common Challenges Across Models:

1. Sarcasm and Nuanced Language: All models struggled with reviews containing sarcasm or mixed sentiment. This highlights the need for more sophisticated language understanding in sentiment analysis.
2. Context Sensitivity: The models often failed to capture the overall sentiment when individual words contradicted the general tone of the review. This suggests that incorporating more context or using sequence models might improve performance.
3. Imbalanced Feature Influence: Certain words (e.g., "dull", "unfunni") have a disproportionate influence on classification, sometimes leading to misclassification when these words are used in unexpected contexts.
4. Positive Bias: Both Naive Bayes and SVM showed a tendency to classify reviews as positive more often than negative, which could be addressed by adjusting class weights or thresholds.

In conclusion, while our models show promising performance, there is room for improvement in handling nuanced language, sarcasm, and context-dependent sentiment. Future work could explore more advanced natural language processing techniques, such as attention mechanisms or contextual embeddings, to address these challenges

**7.Conclusions and Future Work:**

Our study demonstrates the effectiveness of traditional machine learning models in sentiment analysis of Rotten Tomatoes movie reviews. The best-performing models, SVM and Logistic Regression with TF-IDF vectorization, achieved an accuracy of 78.95%, showing promise in automatically classifying review sentiment.

Key findings include:

1. The importance of preprocessing and feature engineering in improving model performance.
2. The effectiveness of TF-IDF vectorization over simple Bag of Words.
3. The challenge of capturing nuanced language, sarcasm, and context-dependent sentiment.

Future work could explore several avenues:

1. Implementing deep learning models such as LSTM or BERT, which have shown state-of-the-art performance in similar tasks.
2. Incorporating aspect-based sentiment analysis to provide more granular insights into specific movie elements (e.g., acting, plot, cinematography).
3. Exploring transfer learning approaches to leverage pre-trained language models for improved performance.
4. Investigating the temporal aspect of reviews to understand how sentiment trends change over time for different movies or genres.
5. Addressing the challenges identified in our error analysis,

particularly in handling sarcasm and mixed sentiment.

By building on this foundation, future research can further enhance the accuracy and applicability of sentiment analysis in the domain of movie reviews, providing valuable insights for the film industry and audiences alike.

## 8.References:

Pang, B., Lee, L., & Vaithyanathan, S. (2002). Thumbs up? Sentiment classification using machine learning techniques. Proceedings of EMNLP, 79-86.

Maas, A. L., Daly, R. E., Pham, P. T., Huang, D., Ng, A. Y., & Potts, C. (2011). Learning word vectors for sentiment analysis. Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies, 142-150.

Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A. Y., & Potts, C. (2013). Recursive deep models for semantic compositionality over a sentiment treebank. Proceedings of the 2013 conference on empirical methods in natural language processing, 1631-1642.

Kim, Y. (2014). Convolutional neural networks for sentence classification. Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), 1746-1751.

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

Zhang, L., Wang, S., & Liu, B. (2018). Deep learning for sentiment analysis: A survey. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 8(4), e1253.

Liu, B. (2012). Sentiment analysis and opinion mining. Synthesis lectures on human language technologies, 5(1), 1-167.

Agarwal, A., Xie, B., Vovsha, I., Rambow, O., & Passonneau, R. (2011). Sentiment analysis of Twitter data. Proceedings of the Workshop on Language in Social Media (LSM 2011), 30-38.

Glorot, X., Bordes, A., & Bengio, Y. (2011). Domain adaptation for large-scale sentiment classification: A deep learning approach. Proceedings of the 28th International Conference on Machine Learning (ICML-11), 513-520.

Pak, A., & Paroubek, P. (2010). Twitter as a corpus for sentiment analysis and opinion mining. LREc, 10(2010), 1320-1326.

Thelwall, M., Buckley, K., & Paltoglou, G. (2012). Sentiment strength detection for the social web. Journal of the American Society for Information Science and Technology, 63(1), 163-173.

Taboada, M., Brooke, J., Tofiloski, M., Voll, K., & Stede, M. (2011). Lexicon-based methods for sentiment analysis. Computational linguistics, 37(2), 267-307.

Mohammad, S. M., & Turney, P. D. (2013). Crowdsourcing a word–emotion association lexicon. Computational Intelligence, 29(3), 436-465.