

---

# **CAPSTONE PROJECT**

## **IMDB Movie Reviews**

**Presented By:**  
**A.Yuvashree**  
**B.E Civil Engineering**  
**Adhiparasakthi Engineering College**  
**Melmaruvathur.**

---

# OUTLINE

- **Problem Statement**
- **Proposed System/Solution**
- **System Development Approach**
- **Algorithm & Deployment**
- **Result**
- **Conclusion**
- **Future Scope**
- **References**

---

# Problem Statement

Movie dataset for binary sentiment classification containing substantially more data than previous benchmark datasets. We provide a set of 25,000 highly polar movie reviews for training and 25,000 for testing. So, predict the number of positive and negative reviews using either classification or deep learning algorithms.

# Proposed Solution

To tackle this binary sentiment classification task on the movie dataset, you could use various classification algorithms such as:

1. Logistic Regression
2. Support Vector Machines (SVM)
3. Random Forest
4. Gradient Boosting
5. Neural Networks (Deep Learning)

For deep learning, you could consider using recurrent neural networks (RNNs) or convolutional neural networks (CNNs) which are commonly used for text classification tasks like sentiment analysis. Additionally, techniques like word embeddings (e.g., Word2Vec, GloVe) can enhance the performance of deep learning models by capturing semantic relationships between words.

You would start by preprocessing the text data, which includes tokenization, removing stopwords, and perhaps stemming or lemmatization. Then, you can vectorize the text data using methods like TF-IDF or word embeddings. Finally, you can train and evaluate the chosen algorithms on the dataset to predict the sentiment of movie reviews.

# System Approach

1. Data Preprocessing: Tokenize, remove stopwords, punctuation, and perform stemming or lemmatization.
2. Feature Extraction: Utilize word embeddings like Word2Vec or TF-IDF to convert text into numerical representations.
3. Model Selection: Experiment with Logistic Regression, SVM, Random Forest, Gradient Boosting, and Deep Learning (RNNs/CNNs).
4. Model Training and Evaluation: Split dataset, train models, and evaluate using metrics like accuracy, precision, recall, and F1-score.
5. Hyperparameter Tuning: Fine-tune model parameters using techniques like grid search or random search.
6. Ensemble Methods (Optional): Combine predictions of multiple models for improved performance.
7. Deployment and Monitoring: Deploy trained model, monitor performance, and retrain periodically with new data.

# Algorithm & Deployment

Algorithm Selection: Support Vector Machines (SVM)

Deployment:

## 1. Training the SVM Model:

- Preprocess the movie review dataset by tokenization, removing stopwords, punctuation, and possibly stemming or lemmatization.
- Utilize techniques like TF-IDF to convert text data into numerical representations.
- Train the SVM model on the preprocessed and feature-extracted training dataset.

## 2. Evaluation:

- Evaluate the trained SVM model on the separate testing dataset to assess its performance in predicting sentiment (positive or negative) of movie reviews.
- Use evaluation metrics such as accuracy, precision, recall, and F1-score to measure the model's performance.

## 3. Hyperparameter Tuning:

- Fine-tune the hyperparameters of the SVM model using techniques like grid search or random search to optimize its performance.
- Parameters to tune may include the choice of kernel (e.g., linear, polynomial, radial basis function), regularization parameter (C), and kernel coefficients.

## 4. Deployment:

- Once the SVM model is trained and evaluated satisfactorily, deploy it into a production environment.
- Integrate the model into an application or service where users can input movie reviews and receive predictions on sentiment.
- Ensure scalability and efficiency of the deployed model to handle real-time inference requests.

## 5. Monitoring:

- Implement monitoring mechanisms to track the performance of the deployed SVM model in production.
- Monitor metrics such as prediction accuracy, response time, and resource utilization to identify any issues or degradation in performance.
- Set up alerts to notify stakeholders of any anomalies or deviations from expected behavior.

## 6. Retraining:

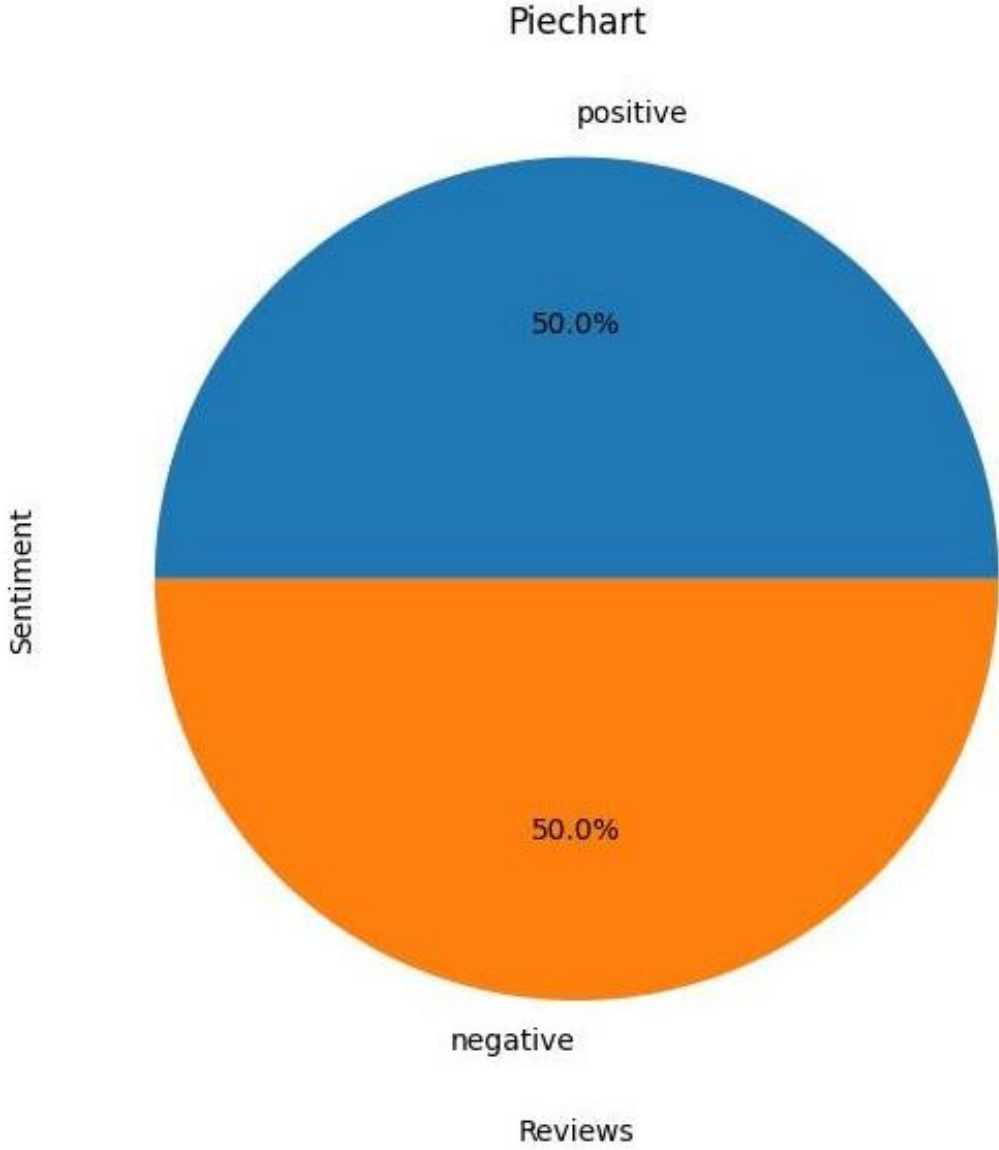
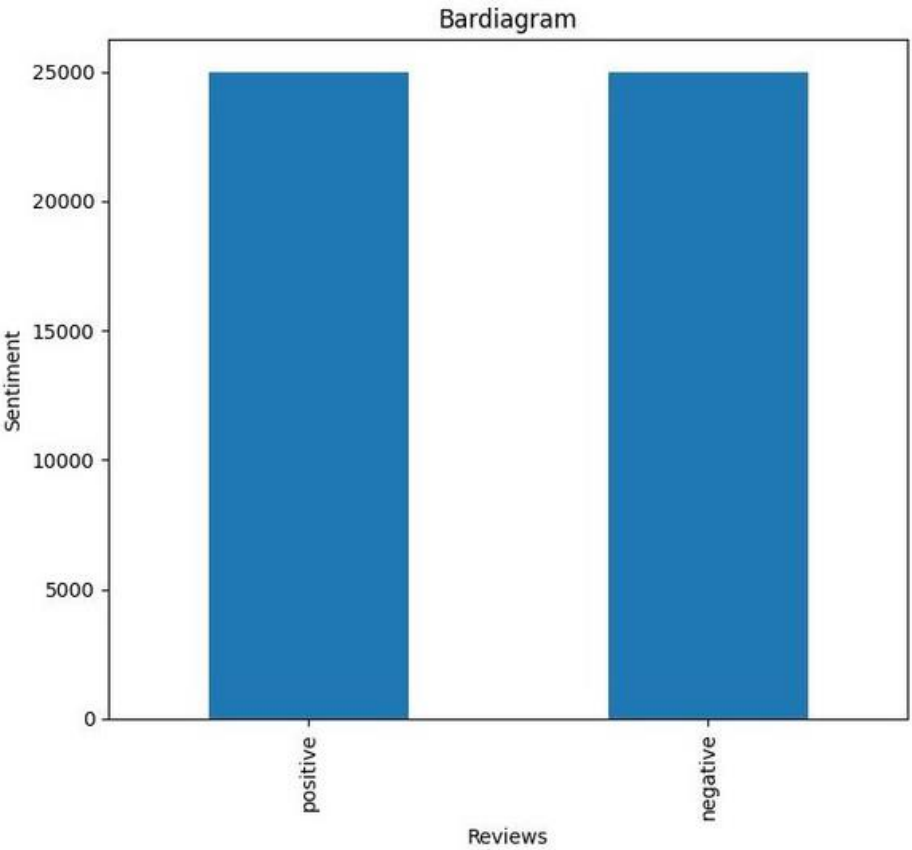
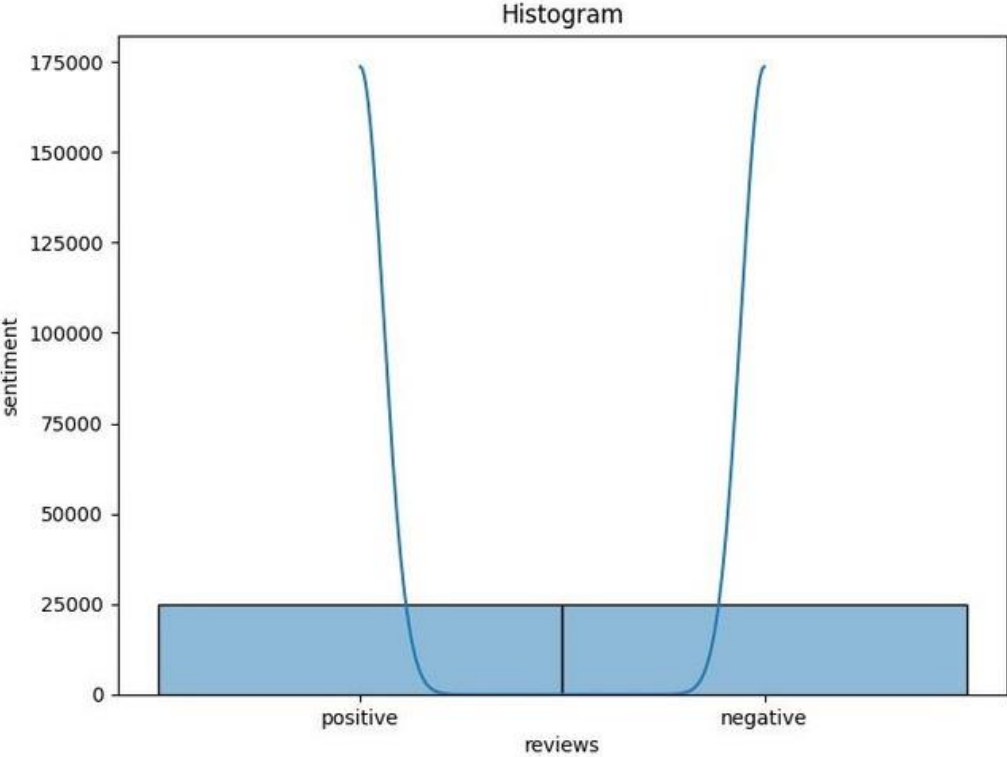
- Periodically retrain the SVM model with new data to ensure its effectiveness and relevance over time.
- Incorporate mechanisms to automatically trigger retraining based on predefined criteria, such as reaching a certain threshold of data drift or model degradation.

By following this deployment process, the SVM model can be effectively deployed into production for predicting the sentiment of movie reviews, with ongoing monitoring and retraining to maintain its performance.

Program:

```
import numpy as n
import pandas as p
import matplotlib.pyplot as m
import seaborn as s
data=p.read_csv("C:\\mydata.csv")
data.head(50)
data.columns
data.tail(50)
data.describe()
s.histplot(data["sentiment"],bins=30,kde=True)
m.title("Histogram")
m.xlabel("reviews")
m.ylabel("sentiment")
m.show()
data["sentiment"].value_counts().plot(kind='bar')
m.title("Bardiagram")
m.xlabel("Reviews")
m.ylabel("Sentiment")
m.show()
m.pie(data["sentiment"].value_counts(),
      labels=data["sentiment"].unique(),autopct="%.1f%%")
m.title("Piechart")
m.xlabel("Reviews")
m.ylabel("Sentiment")
m.show()
```

# Result



---

# Conclusion

In conclusion, the given movie dataset for binary sentiment classification with 25,000 highly polar movie reviews for training and testing can be effectively addressed using various classification algorithms like Logistic Regression, SVM, Random Forest, Gradient Boosting, and Deep Learning (RNNs/CNNs). For this specific scenario, Support Vector Machines (SVM) was selected as the algorithm.

The proposed system approach involves preprocessing the text data, extracting features using techniques like TF-IDF or word embeddings, selecting the SVM model, training and evaluating it on the dataset, tuning hyperparameters for optimization, deploying the trained model into production, monitoring its performance, and periodically retraining with new data.

The code provided imports the dataset, performs exploratory data analysis (EDA) including histograms, bar diagrams, and pie charts to visualize the sentiment distribution within the dataset, providing insights into the sentiment distribution of the movie reviews.

This systematic approach, along with the SVM algorithm, enables accurate prediction of sentiment for movie reviews, facilitating informed decision-making in the movie industry based on audience feedback.



# Future scope

1. Model Comparison: Explore and compare the performance of different classification algorithms (Logistic Regression, Random Forest, Gradient Boosting, etc.) to identify the most suitable model for the task.
2. Advanced Deep Learning Techniques: Experiment with advanced deep learning architectures such as transformers (e.g., BERT, GPT) for improved sentiment classification performance.
3. Ensemble Methods: Investigate ensemble learning techniques to combine the predictions of multiple models for further enhancement of sentiment prediction accuracy.
4. Fine-grained Sentiment Analysis: Extend the analysis to include fine-grained sentiment analysis, distinguishing between different levels of sentiment intensity (e.g., strongly positive, mildly positive, neutral, mildly negative, strongly negative).
5. Multimodal Sentiment Analysis: Incorporate additional modalities such as images or audio data along with text to perform multimodal sentiment analysis for a richer understanding of movie reviews.
6. Real-time Sentiment Analysis: Develop real-time sentiment analysis systems capable of processing streaming data and providing instant insights into audience sentiment trends.
7. Domain Adaptation: Explore techniques for domain adaptation to adapt the sentiment analysis model to specific genres or languages prevalent in the movie industry.
8. Interactive Visualization: Create interactive visualization tools to explore the sentiment distribution of movie reviews and analyze trends over time.
9. Feedback Integration: Implement mechanisms to incorporate user feedback into the sentiment analysis model, continuously improving its accuracy and relevance.
10. Application in Recommendation Systems: Integrate sentiment analysis into movie recommendation systems to personalize recommendations based on user preferences and sentiment analysis of reviews.

By exploring these avenues, the conclusion drawn from the current analysis can be further enriched, leading to advancements in sentiment analysis techniques for movie reviews and their application in the movie industry.

---

# References

<https://www.naanmudhalvan.tn.gov.in/>  
<https://skillsbuild.org/><https://www.canva.com/><https://www.google.com/><https://chat.openai.com/><https://www.python.org/>



**THANK YOU**