# Sentiment Analysis on Financial News Using Machine Learning Techniques

Project Overview:

This project aims to perform sentiment analysis on financial news articles using multiple machine learning models. The objective is to classify news articles into sentiment categories such as positive, neutral, or negative, which can be valuable for financial decision-making. The project involves data preprocessing, sentiment analysis using VADER and TextBlob, and the implementation of machine learning models including Logistic Regression, Support Vector Classifier (SVC), and Random Forest.

1. Data Acquisition and Preparation:

The dataset consists of financial news articles, which were loaded and preprocessed to ensure the quality of the input data. This process included removing duplicates, handling missing values, and standardizing the data format. The preprocessing is crucial for the accurate performance of the sentiment analysis models.

2. Sentiment Analysis:

Two primary methods of sentiment analysis were employed:

- VADER (Valence Aware Dictionary and sEntiment Reasoner): A lexicon and rule-based sentiment analysis tool specifically attuned to sentiments expressed in social media. VADER calculates the sentiment polarity (positive, negative, neutral) of the text.

- TextBlob: A text processing library that provides a simple API for diving into common natural language processing (NLP) tasks, including sentiment analysis. TextBlob also computes the polarity of the text and assigns it to the sentiment categories.

3. Machine Learning Models:

Three machine learning models were implemented and compared:

- Logistic Regression: A statistical model that in this context is used for binary classification. It predicts the probability of

a categorical dependent variable.

- Support Vector Classifier (SVC): A type of Support Vector Machine that is used for classification tasks. It works by finding the hyperplane that best divides the dataset into classes.

- Random Forest: An ensemble learning method that operates by constructing multiple decision trees during training and outputting the mode of the classes for classification.

Each model was trained and tested on the preprocessed dataset, with hyperparameter tuning performed to optimize performance.

4. Model Evaluation:

The models were evaluated using standard classification metrics:

- Accuracy: The proportion of correct predictions.

- Precision, Recall, F1-Score: Metrics that provide deeper insights into the performance of the models, especially in handling imbalanced datasets.

- Confusion Matrix: A tool used to visualize the performance of an algorithm. It shows the true vs. predicted classifications.

The evaluation results for the models are as follows:

- Logistic Regression: Achieved an accuracy of 89.6%, with high precision, recall, and F1-scores.

- SVC: Delivered slightly better performance with an accuracy of 91.4%.

- Random Forest: Provided the best accuracy at 94.7%, demonstrating strong performance across all evaluation metrics.

5. Key Findings:

- The Random Forest model outperformed the other models, offering the highest accuracy and balanced performance across various sentiment classes.

- VADER and TextBlob provided complementary insights into the sentiment analysis, with VADER being particularly

effective for nuanced financial text.

- The models' performance metrics indicate their potential for real-world applications in sentiment analysis of financial news, which can be used to inform trading strategies or risk management.

Conclusion:

This project successfully implemented sentiment analysis on financial news using both lexicon-based methods (VADER, TextBlob) and machine learning models. The Random Forest model was identified as the most effective for this task, highlighting the potential of ensemble learning methods in sentiment classification tasks. The results demonstrate the value of sentiment analysis in financial contexts, where it can support decision-making by providing insights into market sentiment.