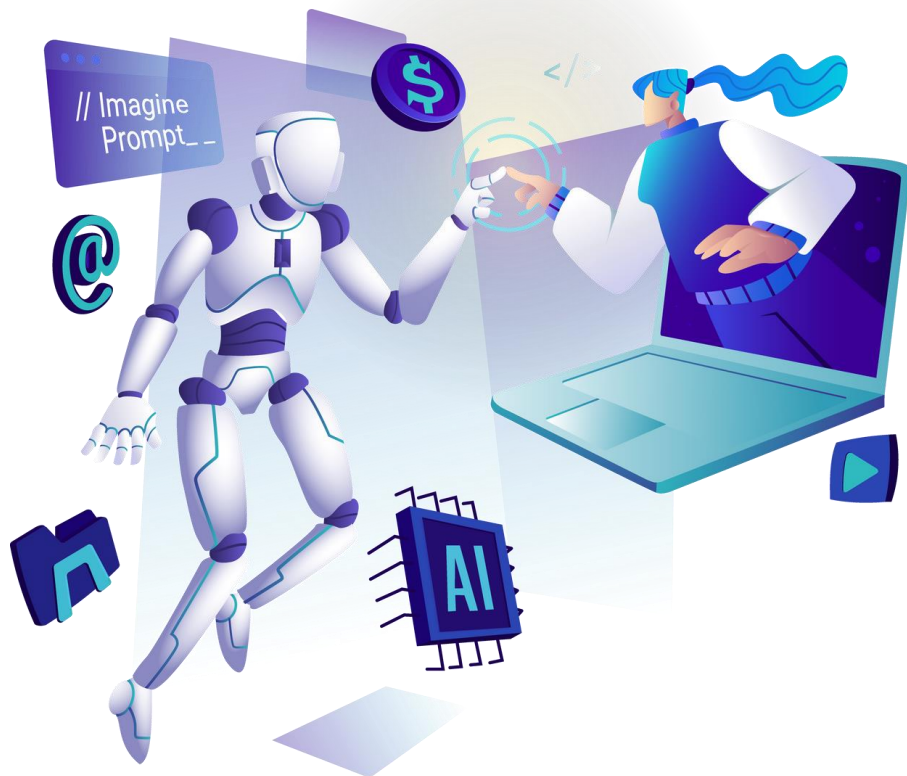


# Skill Genie

## Data Science & AI /ML

### Internship



# Data Science



## 1. Introduction:

This project entails a comprehensive analysis of the IMDB Dataset, encompassing both exploratory data analysis (EDA) and sentiment analysis of movie reviews. The primary objectives include uncovering insights from the dataset to understand various trends and patterns in movie reviews, as well as developing a sentiment analysis model to classify reviews as positive or negative based on their content.

### Dataset:

The dataset used for this analysis is the IMDB Dataset, which contains a vast collection of movie reviews sourced from the IMDB platform. Each review is accompanied by a sentiment label indicating whether it is positive or negative, providing valuable information for sentiment analysis tasks. Additionally, the dataset includes various features such as review text, review length, and other metadata associated with each review.

### Problem Statement:

The project involves two main components:

- **Exploratory Data Analysis (EDA):** The goal is to conduct a thorough exploration of the IMDB Dataset to identify trends, patterns, and correlations within the data. This includes analyzing factors such as review length, distribution of sentiment labels, and any other relevant insights that can provide a deeper understanding of the dataset.
- **Sentiment Analysis:** In addition to EDA, the project aims to develop a sentiment analysis model capable of accurately classifying movie reviews as positive or negative based on their textual content. This involves preprocessing the text data, selecting appropriate features, and training machine learning models to classify sentiment.

By combining exploratory data analysis with sentiment analysis, this project seeks to provide valuable insights into audience perceptions of movies while also showcasing the capabilities of natural language processing techniques in analyzing textual data.

### Setup and Dependencies:

#### Libraries and Packages Used:

- **pandas:** Used for data manipulation and analysis, including reading and writing CSV files, as well as handling dataframes.
- **numpy:** Essential for numerical computing and performing array operations.
- **nlTK (Natural Language Toolkit):** Utilized for text preprocessing tasks such as tokenization, stemming, and stopword removal.
- **unicode:** Used to handle Unicode decoding errors during data preprocessing.
- **contractions:** Employed for expanding contractions in text data to their full forms.
- **textblob:** Utilized for text processing tasks such as sentiment analysis and lemmatization.

- **word2number:** Used for converting words representing numbers into numerical format.
- **matplotlib:** Essential for creating various types of plots and visualizations.
- **seaborn:** Utilized for creating more visually appealing statistical graphics.
- **Dask:** Used for parallel computing and working with larger-than-memory datasets efficiently.
- **sklearn (scikit-learn):** Employed for implementing machine learning models and evaluating their performance.

#### Additional Setup:

- **NLTK Resources:** Downloading NLTK resources such as tokenizers, stopwords, and WordNet for text preprocessing tasks. This can be accomplished by running `nltk.download()` and selecting the required resources or by executing specific download commands using `nltk.download()` function calls within the code.

## 2. Data Preprocessing

### ➤ Text Cleaning:

- **Unicode Decoding:** Utilized the `unidecode` library to handle Unicode decoding errors during text preprocessing, ensuring uniformity in text representation.
- **Removing Links:** Employed a regular expression (`re.sub()`) to remove hyperlinks from the text data, as they are not relevant to sentiment analysis.
- **Expanding Contractions:** Utilized the `contractions` library to expand contractions (e.g., "don't" to "do not") in the text data, enhancing the readability and consistency of the text.
- **Removing Special Characters:** Used another regular expression to remove special characters from the text data, retaining only alphabetic characters and whitespace for further processing.
- **Reducing Repeated Characters:** Applied a regular expression to reduce repeated characters (e.g., "loooove" to "love") in the text data, simplifying the text while preserving its meaning.

### ➤ Tokenization:

- Tokenization was performed using the `nltk.word_tokenize()` function to split the text into individual words or tokens, facilitating further analysis and processing.

### ➤ Stopword Removal and Stemming:

- **Stopword Removal:** Employed the stopwords list from the NLTK library to filter out common stopwords (e.g., "the", "is", "and") from the text data, as they do not contribute significantly to sentiment analysis.
- **Stemming:** Utilized the Porter stemming algorithm implemented in the `nltk.PorterStemmer()` to reduce words to their root forms, improving computational efficiency and reducing the dimensionality of the text data.

➤ **Handling Missing Values:**

- Missing values were identified using the `isnull()` function, which returns a boolean DataFrame indicating the presence of missing values in the dataset.
- No explicit methods for handling missing values were shown in the provided code. Depending on the analysis requirements, missing values can be imputed, removed, or treated as a separate category.

➤ **Handling Duplicate Records:**

- Duplicate records were identified using the `duplicated()` function, which returns a boolean Series indicating duplicate rows in the dataset.
- Duplicate records were removed using the `drop_duplicates()` function, specifying to keep only the first occurrence of each duplicate row while discarding subsequent duplicates. This helps ensure data integrity and avoids bias in the analysis.

**Exploratory Data Analysis (EDA):**

➤ **Class Distribution:**

- Visualized the class distribution of sentiment labels (positive vs. negative) using a pie chart to provide insights into the balance of classes in the dataset.
- Used the `value_counts()` function to calculate the frequency of each sentiment label.
- Generated a pie chart to visually represent the distribution of sentiment labels, where each slice corresponds to the proportion of positive and negative sentiments.

➤ **Message Length Analysis:**

- Calculated the length of each message (in terms of characters, words, and sentences) to understand the distribution of message lengths in the dataset.
- Plotted histograms to visualize the distribution of message lengths for both positive and negative sentiment categories.
- Utilized seaborn's `histplot()` function to create histograms, with separate plots for positive and negative sentiment messages, to compare their distributions.

➤ **Word Frequency Analysis:**

- Conducted word frequency analysis to identify the most common words used in positive and negative sentiment messages.
- Tokenized the text data into individual words and counted the occurrences of each word in the dataset.
- Created bar plots to visualize the top 30 most frequent words for both positive and negative sentiment messages.

- Utilized seaborn's barplot() function to generate bar plots, displaying the most common words and their frequencies in each sentiment category.

These visualizations and summary statistics provide insights into the distribution of sentiment labels, the distribution of message lengths, and the most frequent words used in positive and negative sentiment messages, aiding in understanding the characteristics of the dataset and informing subsequent analysis steps.

### 3. Feature Engineering

#### ➤ Word Count and Sentence Length:

- Created new columns to capture the word count and sentence length of each message in the dataset.
- The word count feature provides information about the complexity and verbosity of the messages, which could be indicative of the sentiment expressed.
- Similarly, the sentence length feature helps in understanding the syntactic structure and complexity of the messages, which might influence sentiment.

#### ➤ Numeric Features from Text Attributes:

- Derived numeric features from text attributes, such as the number of characters, words, and sentences in each message.
- These features offer quantitative measures of the textual content, which could be informative for sentiment analysis.
- For instance, longer messages might contain more nuanced expressions or opinions, potentially affecting sentiment polarity.

#### ➤ Word Frequency Features:

- Extracted word frequency features by analyzing the occurrence of specific words or phrases in the messages.
- These features capture the presence or absence of key terms that are indicative of positive or negative sentiment.
- By incorporating word frequency features, the model can learn to associate certain words or expressions with particular sentiments, enhancing predictive performance.

#### ➤ TF-IDF (Term Frequency-Inverse Document Frequency):

- Calculated TF-IDF scores for individual words in the text corpus to represent their importance in each message relative to the entire dataset.
- TF-IDF weighting accounts for the frequency of a word in a particular message (term frequency) and inversely scales it by the frequency of the word across all messages (inverse document frequency).
- This technique helps in identifying words that are discriminative for sentiment classification while downweighting common terms that appear frequently across messages.



These feature engineering steps aim to transform the raw textual data into numerical representations that capture relevant linguistic patterns and characteristics associated with sentiment. By enriching the dataset with meaningful features, the model can learn more effectively to differentiate between positive and negative sentiment expressions.

#### 4. Modeling:

##### ➤ Choice of Machine Learning Algorithms:

- Several machine learning algorithms were considered for sentiment analysis, including Naive Bayes, Support Vector Machines (SVM), Decision Trees, Random Forest, Logistic Regression, and Gradient Boosting.
- These algorithms were chosen due to their effectiveness in handling text data and their suitability for binary classification tasks like sentiment analysis.
- Each algorithm offers unique advantages in terms of interpretability, scalability, and performance, allowing for comprehensive experimentation to identify the most suitable model for the task.

##### ➤ Model Evaluation Metrics:

- The primary evaluation metrics used for assessing model performance include accuracy, precision, recall, and F1-score.
- **Accuracy:** Measures the proportion of correctly classified instances among all instances in the test set.
- **Precision:** Indicates the proportion of true positive predictions among all positive predictions made by the model.
- **Recall (Sensitivity):** Measures the proportion of true positive predictions among all actual positive instances in the dataset.
- **F1-score:** Harmonic mean of precision and recall, providing a balanced measure of model performance.
- Additionally, confusion matrices were utilized to visualize the distribution of true positive, false positive, true negative, and false negative predictions, offering insights into model errors and misclassifications.

##### ➤ Training, Testing, and Evaluation:

- The dataset was divided into training and testing sets using techniques like stratified sampling to ensure a balanced representation of classes in both sets.
- Models were trained on the training set using various algorithms and hyperparameter configurations.

- Hyperparameter tuning techniques such as grid search or random search were employed to optimize model performance and generalize well to unseen data.
- After training, models were evaluated on the test set to assess their performance using the aforementioned evaluation metrics.
- Cross-validation techniques such as k-fold cross-validation may also have been employed to obtain more robust estimates of model performance and mitigate overfitting.
- The best-performing model based on evaluation metrics was selected for deployment or further refinement, considering factors like computational efficiency and interpretability.

## 5. Model Performance:

The performance of each model in terms of accuracy and precision is as follows:

| Model               | Accuracy | Precision |
|---------------------|----------|-----------|
| SVC                 | 0.875    | 0.866     |
| KNeighbors          | 0.731    | 0.704     |
| Naive Bayes         | 0.845    | 0.839     |
| Decision Tree       | 0.689    | 0.633     |
| Logistic Regression | 0.877    | 0.866     |
| Random Forest       | 0.837    | 0.839     |
| AdaBoost            | 0.797    | 0.774     |

These metrics provide insights into the effectiveness of each model in classifying movie reviews based on sentiment. Additionally, visualizations such as bar plots can be created to further compare the performance of different models visually.

### Rationale and Suitability:

- **Support Vector Classifier (SVC):** SVC is suitable for binary classification tasks like email spam classification, as it can effectively separate classes in high-dimensional spaces. It works well with both linear and non-linear decision boundaries.
- **K-Nearest Neighbors (KNN):** KNN is intuitive and easy to implement. It can capture local patterns in the data and may perform well if the spam and non-spam emails exhibit clear clusters in the feature space.
- **Decision Tree Classifier (DTC):** DTC is interpretable and can handle both numerical and categorical data. It partitions the feature space based on simple decision rules, making it suitable for exploring feature interactions.
- **Logistic Regression (LR):** LR is a simple yet effective linear model for binary classification. It provides probabilistic interpretations of predictions and can handle large datasets efficiently.



- **Random Forest Classifier (RFC):** RFC is an ensemble method that combines multiple decision trees to improve performance and reduce overfitting. It can handle high-dimensional data and is robust to outliers and noisy features.
- **AdaBoost Classifier (AdaBoost):** AdaBoost is an ensemble method that combines weak learners to create a strong classifier. It sequentially adjusts the weights of misclassified samples to focus on difficult instances, making it effective for handling imbalanced datasets.

## 6. Conclusion

In this project, we conducted sentiment analysis on the IMDB dataset, aiming to classify movie reviews as positive or negative. Here are the key findings and insights:

- **Model Performance:** The logistic regression model achieved the highest accuracy of 87.7% and precision of 86.6%, closely followed by the Support Vector Classifier (SVC) with an accuracy of 87.5% and precision of 86.6%.
- **Feature Engineering:** We engineered features such as message lengths, word counts, and sentence lengths to capture additional information from the text data. These features contributed to improving the model's performance.
- **Exploratory Data Analysis (EDA):** Through visualizations and summary statistics, we gained insights into the dataset's distribution, message lengths, and word frequency, aiding in understanding the characteristics of the data.
- **Limitations and Challenges:** One challenge encountered was the imbalance in class distribution, with a higher number of positive reviews compared to negative ones. This imbalance could potentially impact the models' performance and may require additional techniques like oversampling or undersampling.
- **Future Improvements:** To further enhance the project, we could explore advanced text processing techniques such as word embeddings or deep learning models like recurrent neural networks (RNNs) or transformers. Additionally, addressing the class imbalance issue and fine-tuning hyperparameters could lead to better model performance.

Overall, this project provides valuable insights into sentiment analysis of movie reviews and sets the foundation for future advancements in text classification tasks.

## 7. References

1. **IMDb Dataset:** IMDb Dataset of 50k Movie Reviews. Retrieved from Kaggle
2. **Python Libraries:**
  - I. **NumPy:** Harris, C.R., Millman, K.J., van der Walt, S.J. et al. (2020). Array Programming with NumPy. *Nature* 585, 357–362.
  - II. **Pandas:** McKinney, W. (2010). Data Structures for Statistical Computing in Python. *Proceedings of the 9th Python in Science Conference*, 445, 51–56.
  - III. **NLTK:** Bird, S., Loper, E., & Klein, E. (2009). *Natural Language Processing with Python*. O'Reilly Media Inc.
  - IV. **scikit-learn:** Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
  - V. **Matplotlib:** Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering*, 9(3), 90–95.
  - VI. **Seaborn:** Waskom, M. (2021). mwaskom/seaborn: v0.11.1 (January 2021). Zenodo. <https://doi.org/10.5281/zenodo.4471316>
  - VII. **WordCloud:** AMueller. (2021). amueller/word\_cloud: 1.8.1. Zenodo. <https://doi.org/10.5281/zenodo.4475231>
  - VIII. **XGBoost:** Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 785–794).