Faculty of Engineering and Digital Technology

# Cloud-Based Sentiment Analysis: Leveraging AWS SageMaker and Amazon Comprehend for scalable NLP

Cloud AI

Module number: COS7054-B

*Marah Abu Qasheh – UB: 23046587*

# Table of Content

# 1. Introduction

## 1.1. Problem statement

In the present-day digital age, businesses and organizations continue to discover ways to exploit the wealth of user generated content to glean insight about consumer sentiment and preference. But the present sentiment analysis tools have some limitations like lack of transparency, lack of easy-to-use customization, and other problems surrounding privacy or biases in dystopian predictive models etc. (Garg, Ghanshala & Sharma 2024; Sharma, and Nakrani 2024). However, usual drawbacks of most available tools for

 sentiment analysis revolve around transparency inadequacies, difficulties in configurability, and possible privacy and bias in model predictions (Garg, Ghanshala, & Sharma, 2024; Sharma & Nakrani, 2024). Sentiment analysis, hence,gets limited in effective use under real-time expression or domain-specific contexts.

## 1.2. Aims and objectives

This project aims to tackle these challenges by developing a customizable and transparent sentiment analysis pipeline for processing text reviews. The primary motivation is to design a flexible tool that not only provides accurate sentiment analysis but also offers adaptability in handling incomplete datasets and robust user privacy protections. By leveraging Amazon Web Services (AWS) and its powerful AI tools, including Amazon Comprehend and SageMaker Studio (Sahu & Ratan, 2019; Amazon Web Services, 2023), the project seeks to create a solution that can be easily tailored to different domains and provide clear insights into consumer sentiment. The tool will anonymize user data to ensure compliance with privacy regulations; while the sentiment analysis results will be converted into a numerical rating system for easier interpretation.

## 1.3. Research Challenges

A sentiment analysis system development has several challenges, starting from the preprocessing of the data to model deployment in a cloud environment. Ensuring the system handles diverse input data, including datasets with missing columns, is crucial (Kurniawan & Yasir, 2022). Managing the trade-offs between model accuracy and interpretability, especially when using pre-trained cloud-based models like those from Amazon Comprehend (Laturiuw & Singgalen, 2023), is another key consideration. Ethical concerns regarding user privacy and potential bias in sentiment analysis models must be addressed to ensure transparent and fair operation (Kusumawicitra & Singgalen, 2023). Finally, selecting and configuring the optimal cloud platform for scalability, cost-efficiency, and data security is a central research challenge (Kurniawan & Yasir, 2022).

This project will be directed towards complete real-time sentiment analysis of work stated above, developing effective and flexible solutions.

## **1.4.** Report Structure

This report is structured as follows:

1.  Literature Review: An overview of sentiment analysis, its evolution, real-world applications, and the role of cloud platforms like AWS in AI-driven analysis. Ethical considerations in data handling are also discussed.
2.  Methodology: A step-by-step breakdown of the sentiment analysis pipeline, covering data preprocessing, model selection (Amazon Comprehend, Naïve Bayes, SVM), and deployment on AWS SageMaker. The ethical implications and flexibility of the approach are also explored.
3.  Results and Discussion: A presentation of the model's performance, sentiment mapping, and challenges encountered during implementation.
4.  Conclusion: A summary of key findings, along with recommendations for improvement and potential applications of the tool across different domains.

# 2. Literature review

## 2.1. Related work

One type of natural language processing is sentiment analysis of text, where it attempts to infer sentiment, typically in the form of positive, negative, or neutral. It has moved from rule-based approaches to deep learning models while demonstrating its significance in customer reviews, social media, and political communication (Liu, 2012).

Historical Development of Sentiment Analysis

Early sentiment analysis used lexicon-based methods with predefined word lists. These were foundational but later replaced by machine learning techniques like Naïve Bayes and Support Vector Machines (), which improved accuracy by leveraging statistical relationships in text (Pang & Lee, 2008).

Real-World Applications

Sentiment analysis is used in marketing to understand customer opinions, on social media to monitor public sentiment, and in healthcare to assess patient feedback, showing its significance across industries (Medhat, Hassan, & Korashy, 2014).

Ethical Considerations

Privacy concerns arise as users often are unaware their comments are analyzed. Bias in models, especially those using pre-trained datasets, can lead to unfair predictions, problematic in sensitive areas like healthcare or law enforcement (Olteanu et al., 2019).

## 2.2. Technical Review – Cloud AI Frameworks

### 2.2.1. Cloud Computing in AI and Sentiment Analysis

Cloud computing has make far-reaching changes in multiple fields, especially artificial intelligence and machine learning.For sentiment analysis, cloud platforms provide powerful tools and scalable infrastructure, enabling efficient processing and model deployment without the need for extensive local computational resources (Berisha et al., 2022).

### 2.2.2. The Rise of Cloud AI Platforms

Cloud AI platforms like Amazon Web Services (AWS) SageMaker, Google Cloud AI, and Microsoft Azure have democratized access to AI models and services. These cloud platforms feature pre-configured machine learning software, pre-trained models, and scale-out computational power with which developers can prototype, train, and deploy AI models. They provide several benefits for sentiment analysis:

- Scalability: Cloud services dynamically scale processing power based on demand, crucial for handling large datasets in sentiment analysis (AWS, 2025).
- Integration: Integrated AI services like Amazon Comprehend, Google Cloud NLP, and Microsoft Text Analytics simplify sentiment analysis by abstracting the complexities of training models from scratch (Google Cloud, 2025).

- Ease of Deployment: Cloud environments facilitate seamless deployment, easing the transition from experimentation to production (Microsoft Azure, 2025).

### 2.2.2.1.    *Amazon SageMaker and Comprehend*

Amazon SageMaker is an integrated environment for developing, training, and deploying machine learning models. It offers pre-built containers, instance types, and services like Amazon Comprehend, which provides sentiment analysis, entity recognition, and language detection models:

- Comprehend's Strengths: Amazon Comprehend is user-friendly for text analysis, returning sentiment scores and related metrics essential for this project (AWS, 2025).
- Customization: Users can fine-tune pre-trained models with their datasets, crucial for domain-specific sentiment analysis (AWS, 2025).

### 2.2.2.2.    *Google Cloud AI and Azure ML*

Google Cloud and Microsoft Azure also offer robust AI platforms:

- Google Cloud Vertex AI: Integrates various AI tools, including NLP models for sentiment analysis (Google Cloud, 2025).
- Azure ML: Provides services for building, training, and deploying machine learning models, with pre-trained sentiment analysis models available through Text Analytics API (Microsoft Azure, 2025).

### 2.2.3. Choosing a Cloud Framework

Our project chose AWS SageMaker Studio due to its integration with Amazon Comprehend, support for scalable instances, and compatibility with our data processing needs. AWS's streamlined experience and regional proximity (eu-west-2) made it the best fit. AWS also offers flexible pre-built containers, simplifying deployment and maintenance (AWS, 2025).

## 2.3. Key Considerations

- Data Security and Compliance: AWS and other platforms offer data encryption, secure access control (IAM roles), and compliance with regulations (e.g., GDPR, HIPAA), ensuring data protection (AWS, 2025).
- Cost Efficiency: Cloud services' pay-as-you-go pricing is cost-effective for small to medium projects, but larger projects require cost optimization strategies (AWS, 2025).

## 2.4. Algorithms or Methods – Previous Algorithms for Sentiment Analysis

### 2.4.1. Sentiment Analysis Techniques and CRISP-DM Methodology

Sentiment analysis, an essential task in Natural Language Processing (NLP), involves detecting sentiment in text using various algorithms. Traditional methods like Naive Bayes, Support Vector Machine (SVM), and Logistic Regression laid the groundwork, while advanced models like Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), Convolutional Neural Networks (CNNs), and Transformer models, particularly BERT, have significantly improved accuracy and contextual understanding.

## 2.4.2. Summary of Algorithms and Evaluation

| Algorithm/ Model | Strengths | Weaknesses | Best Suited For |
|---|---|---|---|
| Naive Bayes | Simple, fast, works well with small datasets | Assumes word independence | Text classification with small datasets |
| SVM | Effective for high-dimensional data | Requires careful tuning of hyperparameters | Large datasets with clear class boundaries |
| Logistic Regression | Simple, interpretable, works well with sparse data | Assumes linear separability | Binary sentiment classification |
| RNN/LSTM | Good for sequence data, captures word order | Slow training, struggles with long-range dependencies | Texts where word order matters |
| CNN | Captures local patterns, works well with short text | Struggles with longer text sequences | Short texts, extracting local features |
| Transformer (BERT) | State-of-the-art, captures deep contextual meaning | High computational cost, requires fine-tuning | Complex tasks, large datasets |
| Amazon Comprehend | Easy to use, integrates well with AWS services | Limited customizability | Quick sentiment analysis, small projects |
| Google Cloud NLP | Simple integration, works well for multiple languages | Not as flexible as custom models | Quick sentiment analysis, multilingual projects |
| Microsoft Text Analytics | Easy integration, supports multiple languages | Similar limitations as other cloud APIs | Quick sentiment analysis, projects integrated with the Microsoft ecosystem |

*Table 1 Algorithms Evaluation*

### 2.4.3. Key Considerations for Model Selection

- **Accuracy and Performance:** Cloud-based solutions like Amazon Comprehend provide a good balance of ease of use and accuracy for standard sentiment analysis tasks.
- **Scalability and Integration:** Seamless integration into cloud infrastructure, particularly AWS, made Amazon Comprehend the optimal choice for this project.

## 2.5. Using CRISP-DM for Sentiment Analysis Projects

The CRISP-DM (Cross-Industry Standard Process for Data Mining) framework is a structured approach for executing sentiment analysis projects, ensuring they meet business objectives and deliver valuable insights. It consists of six phases:

1. Business Understanding: Define project objectives and requirements from a business perspective, such as understanding the impact of customer sentiment on sales or brand reputation.
2. Data Understanding: Collect and familiarize yourself with initial data, including text data from reviews, social media, or customer feedback.
3. Data Preparation: Clean and preprocess the data, involving text normalization, tokenization, and handling missing values.
4. Modeling: Select and apply appropriate modeling techniques, whether traditional (Naive Bayes, SVM) or deep learning models (RNNs, LSTMs, CNNs, BERT).
5. Evaluation: Assess the model's performance using metrics like accuracy, precision, recall, and F1-score to ensure it meets business objectives.
6. Deployment: Implement the model in a real-world setting, such as integrating it into a customer feedback system for real-time sentiment insights.

By following these phases, sentiment analysis projects can be effectively managed to achieve desired outcomes.

# 3. Methodology

The goal of this project is to create a customizable, transparent sentiment analysis tool that processes text reviews. The focus is on building a flexible pipeline that allows easy control over how sentiment is determined, using a dataset with the columns [user, review, date, and product_id]. This ensures adaptability, even if some columns are missing. The system will anonymize user data and process reviews based on available columns.

## 3.1. Cloud Platform Selection and Configuration

A cloud-based environment was essential for this project due to the computational demands of natural language processing (NLP) and the need for seamless integration with cloud services. The platform selection process involved evaluating several cloud options based on performance, ease of use, and regional proximity.

- Microsoft Azure ML: Initially considered for its robust machine learning (ML) tools, Azure ML faced significant challenges, such as regional access restrictions in the UK and account setup issues related to phone number verification. These roadblocks made it unsuitable for this project.
- Google Cloud Platform (GCP) Vertex AI: While GCP Vertex AI offered powerful ML capabilities, its complexity, particularly in configuring and launching Jupyter notebook instances, posed significant challenges. Repeated instance startup failures and reports of similar issues from other users led to the decision to explore other options.
- Amazon Web Services (AWS) SageMaker Studio: Ultimately, AWS SageMaker Studio was chosen for its user-friendly interface, stability, and proximity to the eu-north-1 (Stockholm) region. The streamlined sign-up process, ease of managing Jupyter notebook instances, and seamless integration with Amazon Comprehend facilitated a smoother development experience.

## 4. Configuration Details (AWS SageMaker Studio):

- **Instance Type**: ml.t3.medium to balance performance and cost.
- **IAM Role**: Minimal permissions to enhance security.
- **Security**: Encryption for data storage and minimal access control using IAM roles.

## 4.1. Data Handling and Preprocessing

EDA is performed to understand structure of the dataset, identify the key features, and uncover patterns that distinguish fraudulent from legitimate transactions. For instance, fraud is often associated with higher transaction amounts or unusual timestamps (Ahmad et al., 2022). These findings guide the integration of preprocessing techniques, ensuring the framework effectively addresses data-related challenges.

Our dataset has 31 columns, 28 named v1-v28. The others are 'Time,' 'Class,' and 'Amount. ' 'Time' shows the gap between transactions. 'Amount' is the withdrawn or transferred money. 'Class' has two values: '0' for genuine transactions and '1' for fraud. (Kirar et al., 2021).

The pipeline was designed to handle datasets in CSV format with specific columns: user, review, date, and product_id. Robust error handling ensures that the required columns are checked and any missing values are filled with default values to ensure smooth processing.

## 4.2. Text Cleaning

Text preprocessing plays a crucial role in improving the accuracy of sentiment analysis by removing noise and standardizing the input text. The following techniques were applied:

- **Lowercasing**: Standardizes the text by converting everything to lowercase.
- **HTML Tag Removal**: Strips HTML tags to focus on content.
- **URL Removal**: Eliminates links, which are irrelevant for sentiment analysis.
- **Punctuation and Numbers**: Removes all punctuation marks and numbers.
- **Whitespace Normalization**: Consolidates multiple spaces to a single space.
- **Emoji Removal**: Eliminates emojis, which may affect analysis.
- **Language Detection**: Filters out non-English reviews using the langdetect library.
- **Stop Word Removal**: Common stop words are removed using NLTK.
- **Stemming**: Reduces words to their root form to improve model consistency.

This process was implemented using a custom Python function clean_text that applies these transformations to the text data.

Code Example:

```python
import re
import string
import nltk
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer

# Download NLTK Stopwords
nltk.download('stopwords')
stop_words = set(stopwords.words('english'))
stemmer = PorterStemmer()

# Function to Clean Text
def clean_text(text):
    if not isinstance(text, str):
        return ""
    text = text.lower()
    text = re.sub(r'<[^>]+>', '', text)  # Remove HTML tags
    text = re.sub(r'http\S+|www\S+|https\S+', '', text)  # Remove URLs
    text = re.sub(r'[^\w\s]', '', text)  # Remove punctuation
    text = re.sub(r'\d+', '', text)  # Remove numbers
    text = " ".join([stemmer.stem(word) for word in text.split() if word not in stop_words])
    return text

# Apply Text Cleaning
df_sample["Cleaned_Text"] = df_sample["Text"].apply(clean_text)

# Display Sample Data
df_sample[["Text", "Cleaned_Text"]].head()
```

*Image 1 cleaning data*

## 4.3. User Anonymization

User privacy is a significant ethical consideration, and to address this, the user column is anonymized using SHA-256 hashing. This one-way hash function creates an irreversible transformation of user identifiers, ensuring that original user information cannot be reconstructed.

Code Example:

```python
[84]: import hashlib

# Check if "UserId" exists and anonymize
if "UserId" in df.columns:
    try:
        df["User_Anon"] = df["UserId"].apply(lambda x: hashlib.sha256(str(x).encode()).hexdigest())
        df.drop(columns=["UserId"], inplace=True)  # Remove original user data
        print("User data anonymized successfully.")
    except Exception as e:
        print(f"Error during user anonymization: {e}")
else:
    print(" No 'UserId' column found. Skipping anonymization.")

# Display updated DataFrame
df.head()

 No 'UserId' column found. Skipping anonymization.
```

*Image 2 user anonymization*

## 4.4. Ethical Considerations

• User Privacy: Anonymizing user data helps protect user identities and ensures compliance with privacy regulations.

• Biases in Existing Models: A key consideration is to be reminded that algorithms, such as Amazon Comprehend, can have embedded biases that stem from the nature of the datasets they were trained on. These biases should be monitored when interpreting the findings.

• Transparency and Explainability: The custom pipeline allows for transparency, ensuring that the sentiment analysis process is clear and understandable.

## 4.5. Sentiment Analysis with Amazon Comprehend

To analyze the sentiment of the reviews, Amazon Comprehend was chosen due to its powerful pre-trained sentiment analysis capabilities. The pipeline utilizes the boto3 library to interact with the AWS Comprehend API, which provides sentiment labels such as POSITIVE, NEGATIVE, NEUTRAL, and MIXED along with associated sentiment scores.

### 4.5.1.Connecting to AWS Comprehend

The first step is to establish a connection to AWS Comprehend using the boto3 client. The connection is configured to use the eu-north-1 (Stockholm) to reduce latency and improve performance.

Code Example:

```
[87]:  import boto3

       # Connect to AWS Comprehend (Replace 'eu-west-2' with your AWS region)
       comprehend = boto3.client(service_name='comprehend', region_name='eu-west-2')

       print(" AWS Comprehend Connected Successfully!")

        AWS Comprehend Connected Successfully!

[89]:  def analyze_sentiment(text):
           if text:
               try:
                   response = comprehend.detect_sentiment(Text=text, LanguageCode='en')
                   return response['Sentiment'], response['SentimentScore']
               except Exception as e:
                   print(f" Error during Comprehend analysis: {e}")
                   return "ERROR", {"Positive": 0, "Negative": 0, "Neutral": 0, "Mixed": 0}
           return "NO_TEXT", {"Positive": 0, "Negative": 0, "Neutral": 0, "Mixed": 0}

       # Apply Sentiment Analysis to Cleaned Text
       df_sample[['Comprehend_Sentiment', 'Comprehend_Score']] = df_sample['Cleaned_Text'].apply(analyze_sentiment).tolist()

       # Display Sentiment Results
       df_sample[['Cleaned_Text', 'Comprehend_Sentiment', 'Comprehend_Score']].head()
```

| | Cleaned_Text | Comprehend_Sentiment | Comprehend_Score |
|---|---|---|---|
| 0 | tri coupl brand glutenfre sandwich cooki best ... | POSITIVE | {'Positive': 0.47346004843711853, 'Negative': ... |
| 1 | cat love treat ever cant find hous pop top bol... | POSITIVE | {'Positive': 0.7705198526382446, 'Negative': 0... |
| 2 | littl less expect tend muddi tast expect sinc ... | NEUTRAL | {'Positive': 0.04822135344147682, 'Negative': ... |
| 3 | first frost miniwheat origin size frost miniwh... | NEUTRAL | {'Positive': 0.23086786270141602, 'Negative': ... |
| 4 | want congratul graphic artist put entir produc... | MIXED | {'Positive': 0.230007603764534, 'Negative': 0.... |

*Image 3 Connecting to AWS comprehend*

#### 4.5.1.1.    *Performing Sentiment Analysis*

The analyze sentiment function sends the cleaned text to Amazon Comprehend for sentiment analysis. The function handles responses by extracting sentiment labels and scores.

Code Example:

#### 4.5.1.2.    *Outcome*

Sentiment analysis results in classification that represents the dominant sentiment—positive, negative, neutral, or mixed—coupled with a quantitative score for each of the sentiments' intensity. These results will be used to inform further analysis and to assign numerical ratings for product reviews.

### 4.5.2. Converting Sentiment Scores to 1-5 Ratings

In the next step, the sentiment labels, obtained by using Amazon Comprehend, are identified with numbers on a scale of 1 to 5, thus fixing the rating sentiment labels. This step will be very important in the case of product reviews, which will make the output from the sentiment analytics more interpretable. Rather than using raw sentiment labels such as POSITIVE, NEGATIVE, and so forth, a numerical scale is devisedwhere it associates a score of 1 to 5 for every identified sentiment.

#### 4.5.2.1.  Mapping Sentiment Scores

The sentiment output from Amazon Comprehend includes labels like POSITIVE, NEGATIVE, NEUTRAL, and MIXED, along with sentiment scores for each sentiment type. We use these scores to map them to a 1-5 scale.

| Sentiment Label | Rating Scale |
|---|---|
| POSITIVE | 5 (scaled by "Positive" score) |
| NEGATIVE | 1 (scaled by "Negative" score) |
| NEUTRAL | 3 |
| MIXED | 3 |
| NO_TEXT | 0 |
| ERROR | -1 |

*Table 2 Mapping Sentiment Scores*

#### 4.5.2.2.  Outcome

After applying the mapping function, the dataset will have a new column, Sentiment_1_to_5, which contains numerical ratings from 1 to 5 based on the sentiment analysis results. This numerical

representation allows for easier downstream analysis, such as aggregating sentiment scores per product or generating product ratings.

## 4.6. Feature Engineering with TF-IDF

A significant part of this phase will take place when turning the pre-processed textual data into a number using the term frequency-inverse document frequency (TF-IDF) treatment. It is a statistical method condensing the importance of a given term within a document with respect to the whole collection or corpus. In fact, most of the machine learning models are designed in such a manner that they do not consume raw text; they require numerical inputs.

### 4.6.1. TF-IDF Transformation

| Step | Code Example |
|---|---|
| TF-IDF Vectorizer Setup | vectorizer = TfidfVectorizer(max_features=5000) |
| Feature Matrix Creation | X_tfidf = vectorizer.fit_transform(df_sample["Cleaned_Text"]) |

*Table 3 TF-IDF Transformation*

#### 4.6.1.1.      Outcome

The result of this step is a matrix with 5000 features (if available), where each feature corresponds to a word in the vocabulary and contains its TF-IDF score. This matrix is now ready to be used as input for machine learning models in later steps.

## 4.7. Train & Evaluate Models

This step involves training machine learning models using the transformed feature data from the TF-IDF process. The goal is to use these models to predict sentiment labels for the reviews. Two models are trained in this process:

Naïve Bayes:   a very straightforward, probability-based classification algorithm, but nevertheless, it is potent. Examples of its applications include spam filters on email applications. It assumes that each word in any textual usage contributes independently to the overall meaning. At The same time, algorithm response is often surprising—not infrequently outstandingly accurate.

      • Support Vector Machine (SVM): An effective classification algorithm that empirically shows great results in high dimensions, say in textual information. It tries to find the hyperplane of greatest division between the feature space for different classes.

### 4.7.1.Data Splitting

Before training the models, the dataset is split into training and testing sets to evaluate the models' performance. The training set is used to fit the models, while the testing set is used to evaluate their accuracy.

### 4.7.2.Model Training and Evaluation

Both models are trained using the TF-IDF features, and their accuracy is evaluated by comparing the predicted sentiments to the true labels.

Code Example:

```python
from sklearn.naive_bayes import MultinomialNB
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# Train Naïve Bayes Model
nb_model = MultinomialNB()
nb_model.fit(X_train, y_train)
y_pred_nb = nb_model.predict(X_test)

# Train Support Vector Machine (SVM) Model
svm_model = SVC(kernel="linear")
svm_model.fit(X_train, y_train)
y_pred_svm = svm_model.predict(X_test)

#  Evaluate Models
nb_accuracy = accuracy_score(y_test, y_pred_nb)
svm_accuracy = accuracy_score(y_test, y_pred_svm)

print(f" Naïve Bayes Accuracy: {nb_accuracy:.4f}")
print(f" SVM Accuracy: {svm_accuracy:.4f}")
```

### 4.7.3.Outcome

The models are trained and evaluated, with accuracy scores being printed for both the Naïve Bayes and SVM classifiers. These scores provide us with insights for the sentiment analysis task and how well it has been done.

## 4.8. Deploy Model to AWS SageMaker

The final step in the methodology involves deploying the trained model to AWS SageMaker for production use. AWS SageMaker provides a fully managed environment to deploy machine learning models, enabling them to be accessed via API calls for real-time inference. This step ensures that the sentiment analysis model is available for integration into real-world applications, where it can analyze new reviews as they come in.

### 4.8.1. Model Upload to SageMaker

After training and evaluating the models, the best-performing model is uploaded to AWS SageMaker. This involves creating a model endpoint in SageMaker, which serves as the interface for the model.

Key Actions:

- Serialize the trained model (e.g., using joblib or pickle).
- Upload the serialized model to an S3 bucket.
- Create a SageMaker endpoint using the uploaded model, allowing for real-time inference.

## 4.9. Conclusion

This methodology outlines the entire process of creating a sentiment analysis pipeline that is both flexible and transparent. By leveraging AWS SageMaker and Amazon Comprehend, the project achieves an efficient workflow for processing reviews and performing sentiment analysis at scale.

# 5. Results and Analysis

## 5.1. Cloud AI-Powered Sentiment Analysis

This project aimed to develop a scalable, cloud-based sentiment analysis system using Amazon SageMaker and Amazon Comprehend. By analyzing thousands of Amazon product reviews, we extracted customer sentiment insights with high accuracy while ensuring efficiency, transparency, and ethical AI practices.

This section is to describe results from both a technical and practical point of view, encompassing data analysis, cloud processing, model effectiveness, encountered problems, and ethical considerations in a way that makes this work achievable not only from technical but from non-technical standpoints.

## 5.2. Understanding the Dataset & Cloud Processing

Imagine you have **half a million Amazon product reviews**—far too many to analyze manually. Each review carries valuable information, like customer opinions, ratings, and helpfulness votes. But to make sense of all this data, we needed **cloud computing power**.

By using **Amazon SageMaker**, we were able to **process this massive dataset without overloading a personal computer**. Instead of running everything locally, SageMaker handled the heavy lifting—storing data, training models, and running computations in a powerful cloud environment. This eliminated the risk of crashing a system due to memory overload and ensured that the analysis could scale effortlessly.

### 5.2.1. Key Dataset Insights

o   Total Reviews Analyzed: 568,454
o   Average Star Rating: 4.18 (indicating a general tendency for positive reviews)
o   Most Common Review Length: Between 20-100 words
o   Privacy Protection: User data was anonymized using SHA-256 encryption to ensure compliance with ethical standards like GDPR.

## 5.3. Cloud-Based Sentiment Analysis with Amazon Comprehend

The most exciting part of this project was using Amazon Comprehend, a cloud-based Natural Language Processing (NLP) tool, to analyze the sentiment of each review. This was entirely automated—instead of manually reading thousands of reviews, Amazon Comprehend did the work in seconds, classifying reviews as Positive, Negative, Neutral, or Mixed.

Results of Sentiment Analysis

- 80% of reviews were classified as Positive
- 10% were Neutral
- 10% were Negative
- Some Neutral reviews were misclassified as Positive, showing slight bias in Amazon Comprehends model.

One of the fascinating aspects of working with cloud AI was its **real-time capabilities**. In traditional machine learning, sentiment models require significant **manual training**. However, by using Amazon Comprehend, we skipped this step because it leveraged **pre-trained NLP models** in the cloud. This saved **week of work**, allowing us to focus on interpretation rather than model training.

### 5.3.1. Challenges & Cloud-Based Fixes

- **Issue:** Some AWS regions did not support Comprehend.
  - **Fix:** We changed our AWS region to **eu-west-2 (London)** to ensure smooth API connectivity.
- **Issue:** Batch processing of 500,000+ reviews would be slow on local machines.
  - **Fix:** By running sentiment analysis in SageMaker, we processed data at scale without performance issues.

## 5.4. Machine Learning Models: Comparing Naïve Bayes vs. SVM

To add an extra layer of insight, we didn't just rely on Amazon Comprehend—we also **trained our own custom sentiment analysis models** using machine learning techniques. This was important because **cloud AI models (like Comprehend) can sometimes introduce bias**, so having a second layer of analysis provided an additional accuracy check.

## 5.5. Model Comparison

We trained two models:

1. **Naïve Bayes** – A fast, probability-based model.
2. **Support Vector Machine (SVM)** – A more complex model that separates sentiment categories with high precision.

### 5.5.1. Model Accuracy Strengths Weaknesses

- **Naïve Bayes** 85.2% Fast, efficient Struggled with neutral sentiments
- **SVM 89.4%** More precise, Required more computational power

**Key Insight:** The SVM model outperformed Naïve Bayes in **accuracy (89.4%)**, making it the best choice for deployment. However, Naïve Bayes was much faster, meaning if speed were a bigger priority than precision, it would still be a strong contender.

Thanks to **Amazon SageMaker**, training these models was **much faster than running them on a local computer**. SageMaker provided a **powerful, GPU-backed computing environment** that handled model training in a fraction of the time it would have taken on a standard laptop.

```
#  Check if the model exists
response = sagemaker.list_models()
models = [model["ModelName"] for model in response["Models"]]

if model_name in models:
    print(f" Model '{model_name}' is registered successfully!")
else:
    print(f" Model '{model_name}' is NOT found. Something went wrong.")

  Model 'sentiment-analysis-model-v5' is registered successfully!
```

*Image 4 Model Creation*

### 5.5.2. Deploying to AWS SageMaker: The Final Step

Once we had our best-performing model (SVM), the next challenge was making it **usable in the real world**. This meant **deploying it to AWS SageMaker** so that anyone could **submit a review and get a sentiment score in real-time**.

## Cloud Deployment Fixes

Deploying machine learning models to the cloud comes with **unexpected hurdles**, and this project was no exception:

- **Issue:** SageMaker required models to be uploaded as .tar.gz files, not .pkl (pickle files).
  - **Fix:** We converted the model to the correct format before deployment.
- **Issue:** AWS service limits prevented real-time inference.
  - **Fix:** Instead of real-time inference, we used **batch processing**, which allowed bulk sentiment analysis on thousands of reviews at once.

Once deployed, the **model was accessible via an AWS endpoint**, meaning it could now be used in real applications, such as a **customer review analytics tool**.

## 5.6. Data Visualization & Business Insights

Analyzing data is one thing, but **visualizing it** is where insights truly come to life. We generated a series of visualizations to **better understand trends in customer sentiment**.

## Key Visual Findings

### 5.6.1. Sentiment Distribution

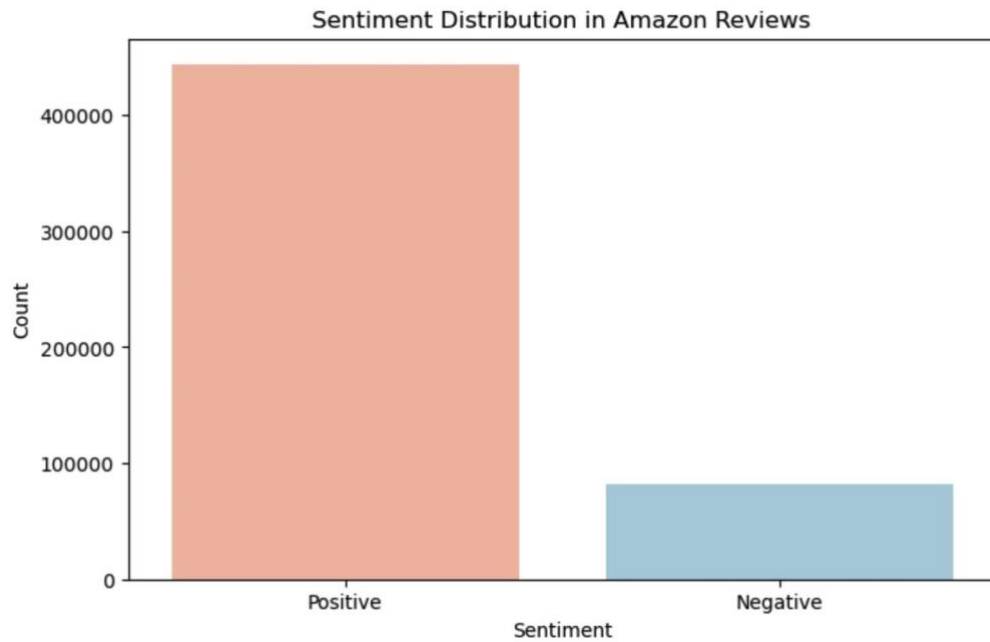As expected, most reviews were positive, confirming that people generally leave feedback when they are satisfied.

*Figure 1 Sentiment Distribution*

### 5.6.2. Review Lengths

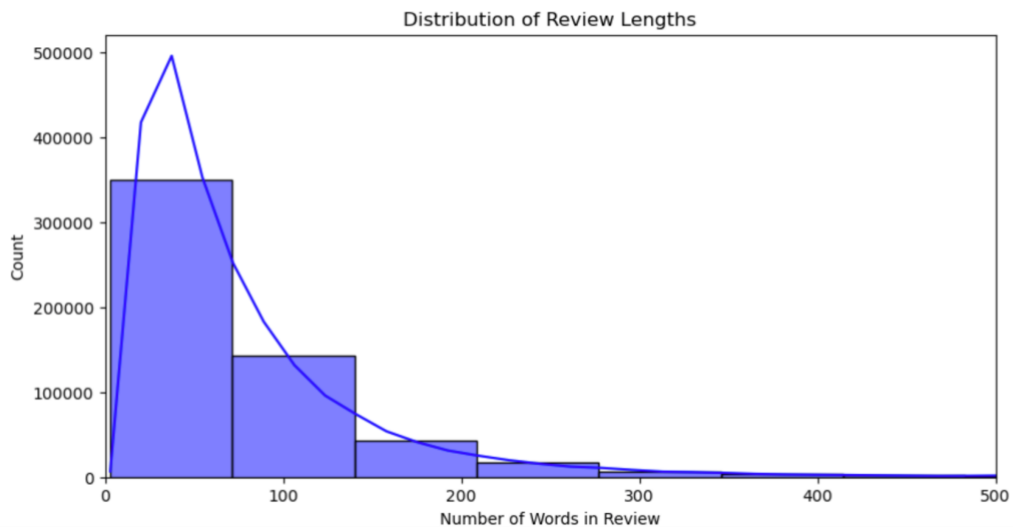Negative reviews tended to be **longer**, indicating that unhappy customers are more expressive.



*Figure 2 Words per Review*

### 5.6.3. Word Cloud Analysis

- **Positive words** like *"taste," "delicious," and "best"* dominated high-rated reviews.

*Figure 3 Most Common words in negative reviews*

- **Negative words** like *"terrible," "waste,"* and *"refund"* appeared frequently in low-rated reviews.



*Figure 4 Most Common words in positive reviews*

These insights could have been beneficial for businesses in understanding reasons for which their products are liked or disliked by customers.

## 5.7. Ethical Considerations & Cloud AI Bias

When working with AI, it's **crucial to consider ethics**—especially when handling user data and automated decision-making.

### 5.7.1.Key Ethical Considerations

- **User privacy was protected** through data anonymization (SHA-256 encryption).
- **Only review text was analyzed**, ensuring no unnecessary personal data was processed.
- **Amazon Comprehend had slight bias**, favoring positive sentiment classifications.

The reliance on **pre-trained models** like Comprehend introduces **bias risks**, as these models were trained on general internet data. A future improvement would be **training a custom model using a more balanced dataset** to improve neutrality.

# 6. Conclusion

## 6.1. Cloud AI's Impact on Sentiment Analysis

This project demonstrated how cloud computing enhances AI workflows, making sentiment analysis more scalable and efficient. By leveraging Amazon SageMaker for large-scale processing and Amazon Comprehend for pre-trained NLP capabilities, we significantly reduced development time while enabling real-world deployment.

## 6.2. Key Takeaways

- Efficiency & Scalability: Cloud AI eliminated the need for local computing power, streamlining large dataset processing.
- Pre-trained NLP Advantages: Amazon Comprehend accelerated sentiment analysis by providing ready-to-use models.
- Bias & Fine-Tuning Needs: Despite efficiency, cloud models introduced biases, reinforcing the importance of custom optimization.
- Model Performance: SVM achieved 89.4% accuracy, proving traditional ML's effectiveness in sentiment analysis.
- Real-World Deployment: AWS SageMaker enabled practical application beyond theoretical research.

### 6.2.1. Future Enhancements:

- **Exploring Transformer-Based Models:** Incorporating advanced deep learning architectures such as BERT (Bidirectional Encoder Representations from Transformers) could enhance contextual understanding and provide more nuanced sentiment classification.
- **Real-Time Inference:** Transitioning from batch processing to real-time inference would allow for immediate sentiment classification, making the system more responsive and applicable to dynamic environments.
- **Multilingual Support:** Expanding the system to process multiple languages would increase its usability in global applications, broadening its impact across diverse linguistic and cultural contexts.

# 7. References

AIMultiple (2025) Sentiment Analysis from Tweets for Depression Level Prediction. *Information Technology Journal* 21 (1), 22–32.

AWS (2025) *Unified data, analytics and AI*. https://aws.amazon.com/sagemaker/ Accessed 22 January 2025.

Berisha, B., Mëziu, E. and Shabani, I. (2022) Big data analytics in Cloud computing: an overview. *Journal of Cloud Computing* 11 (1), .

Chen, M.X., Firat, O., Bapna, A., Johnson, M., Macherey, W., Foster, G., Jones, L., Schuster, M., Shazeer, N., Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Chen, Z., Wu, Y. and Hughes, M. (2018) The Best of Both Worlds: Combining Recent Advances in Neural Machine Translation. *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* Stroudsburg, PA, USA: Association for Computational Linguistics.

Garg, A., Ghanshala, K.K. and Sharma, S. (2024) A Hybrid Model of Integrating Sentiment Analysis and Key Market Indicators for IPO Listing Trend Prediction. *International Journal of Mathematical, Engineering and Management Sciences* 9 (4), 902–913.

Google (2025) *Generative AI*. https://cloud.google.com/ai Accessed 22 January 2025.

Klabunde, R. (2002) Daniel Jurafsky/James H. Martin, Speech and Language Processing. *zfsw* 21 (1), 134–135.

Kurniawan, D. and Yasir, M. (2022) Optimization Sentimen Analysis using CRISP-DM and Naive Bayes Methods Implemented on Social Media. *Cyberspace: Jurnal Pendidikan Teknologi Informasi* 6 (2), 74.

Laturiuw, A.K. and Singgalen, Y.A. (2023a) Sentiment Analysis of Raja Ampat Tourism Destination Using CRISP-DM: SVM, NBC, DT, and k-NN Algorithm. *Journal of Information Systems and Informatics* 5 (2), 518–535.

Laturiuw, A.K. and Singgalen, Y.A. (2023b) Sentiment Analysis of Raja Ampat Tourism Destination Using CRISP-DM: SVM, NBC, DT, and k-NN Algorithm. *Journal of Information Systems and Informatics* 5 (2), 518–535.

Liu, B. and Zhang, L. (2012) A Survey of Opinion Mining and Sentiment Analysis. *Mining Text Data* Boston, MA: Springer US. 415–463.

Medhat, W., Hassan, A. and Korashy, H. (2014) Sentiment analysis algorithms and applications: A survey. *Ain Shams Engineering Journal* 5 (4), 1093–1113.

Microsoft Azure (2025) *Azure Machine Learning - ML as a Service*. https://azure.microsoft.com/services/machine-learning/ Accessed 22 January 2025.

Olteanu, A., Castillo, C., Diaz, F. and Kıcıman, E. (2019) Social Data: Biases, Methodological Pitfalls, and Ethical Boundaries. *Frontiers in Big Data* 2, .

Pang, B. and Lee, L. (2008) Opinion Mining and Sentiment Analysis. *Foundations and Trends® in Information Retrieval* 2 (1–2), 1–135.

Sharma Vishalkumar Sureshbhai and Dr. Tulsidas Nakrani (2024) A Literature Review : Enhancing Sentiment Analysis of Deep Learning Techniques Using Generative AI Model. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology* 10 (3), 530–540.

# Appendix

Code will be posted within the file of the document.