# Tech Saksham

## Capstone Project Report

## ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING FUNDAMENTALS

# "DETECTING SPAM EMAILS "

## "ARULMIGU MEENAKSHI AMMAN COLLEGE OF ENGINEERING"

| NM ID | NAME |
|---|---|
| aut2251030077 | B Nishanth |

Trainer Name

RAMAR BOSE

Sr.AI Master Trainer

# ABSTRACT

This project report provides a thorough investigation into the application of artificial intelligence (AI) techniques for detecting spam emails. With the exponential growth of email communication, spam emails have become a pervasive problem, posing threats ranging from phishing attacks to malware distribution. Traditional rule-based methods and content-based filtering techniques have shown limitations in adapting to evolving spam tactics. Consequently, this report explores the utilization of AI methodologies, including machine learning and natural language processing (NLP), to improve spam detection accuracy and efficiency. Various AI algorithms, such as support vector machines, random forests, and deep neural networks, are evaluated for their efficacy in distinguishing spam from legitimate emails. Additionally, feature engineering techniques and ensemble learning approaches are examined to enhance model performance. Challenges inherent in AI-based spam detection, such as imbalanced datasets, concept drift, and the need for real-time detection, are analyzed along with proposed solutions and future research directions. Moreover, the report discusses ethical considerations, including privacy preservation, algorithmic fairness, and the impact of false positives on user experience. Through a comprehensive review of methodologies, challenges, and ethical considerations, this project report aims to contribute to the advancement of AI-driven solutions for spam email detection and cybersecurity efforts

# INDEX

# CHAPTER 1

# INTRODUCTION

## 1.1 Problem Statement

Developing an AI-powered spam email detection system to accurately differentiate between legitimate and spam emails in real-time, while minimizing false positives and adapting to evolving spam tactics.

## 1.2 Proposed Solution

Utilizing advanced artificial intelligence (AI) techniques such as machine learning and natural language processing (NLP) to develop a sophisticated spam email detection system. This system will analyze email content, metadata, and behavioral patterns to accurately identify spam, while continuously learning and adapting to new spam tactics. Additionally, strategies for handling imbalanced datasets and ensuring ethical considerations will be integrated into the design and implementation of the system.

## 1.3 Feature

1   Machine Learning Algorithms
2   Natural Language Processing (NLP)
3   Real-time Detection
4   Adaptive Learning
5   Imbalanced Dataset Handling
6   Ethical Considerations

## 1.4 Advantages

1. Enhanced Accuracy
2. Real-time Protection
3. Scalability
4. Reduced False Positives
5. Adaptability
6. Improved User Experience

## 1.5 Scope

1. Algorithm Development: Designing and implementing machine learning and NLP algorithms for spam detection.
2. Data Collection and Preprocessing: Gathering email datasets and preprocessing them for model training.
3. Model Evaluation: Assessing the performance of developed algorithms using metrics like precision, recall, and F1 score.
4. Real-time Implementation: Integrating the spam detection system into email platforms for real-time protection.
5. Adaptation Mechanisms: Implementing mechanisms for the system to adapt to new spam tactics and user feedback.
6. Ethical Considerations: Addressing privacy concerns and ensuring algorithmic fairness in spam detection.

## 1.6 Future Work

1. Enhanced Model Performance: Further refining machine learning models and NLP techniques to improve detection accuracy and reduce false positives.
2. Integration with Email Services: Collaborating with email service providers to integrate AI-based spam detection systems into their platforms for widespread adoption.
3. Advanced Feature Engineering: Exploring novel features and data representations to capture subtle nuances in email content for better spam detection.
4. Dynamic Adaptation: Developing mechanisms for the system to dynamically adapt to changing spam tactics and user behavior in real-time.
5. Multimodal Analysis: Investigating the integration of additional data modalities such as images and attachments for comprehensive spam detection.
6. Privacy-Preserving Solutions: Researching methods to enhance user privacy while maintaining effective spam detection capabilities, such as federated learning or on-device processing.

**CHAPTER 2**

# SERVICES AND TOOLS REQUIRED

## 2.1   Services Used

1. Dataset

2. Preprocessing

3. Feature Extraction

4. Machine Learning Algorithms

5. Model Evaluation

6. Cross-Validation

7. Hyperparameter Tuning

8. Ensemble Methods

9. Deep Learning Models

10. Deployment Tools

11. Monitoring & Maintenance

12. External APIs

## 2.2 Tools and Software used

1. Python
2. scikit-learn
3. NLTK or spaCy
4. TensorFlow or PyTorch
5. Pandas
6. NumPy
7. Matplotlib or Seaborn
8. Jupyter Notebook or Google Colab
9. GitHub or GitLab
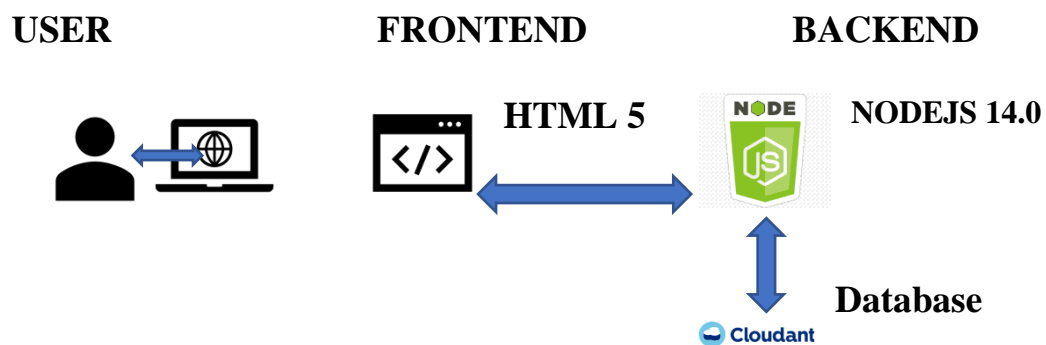10. Flask or Django
11. Docker

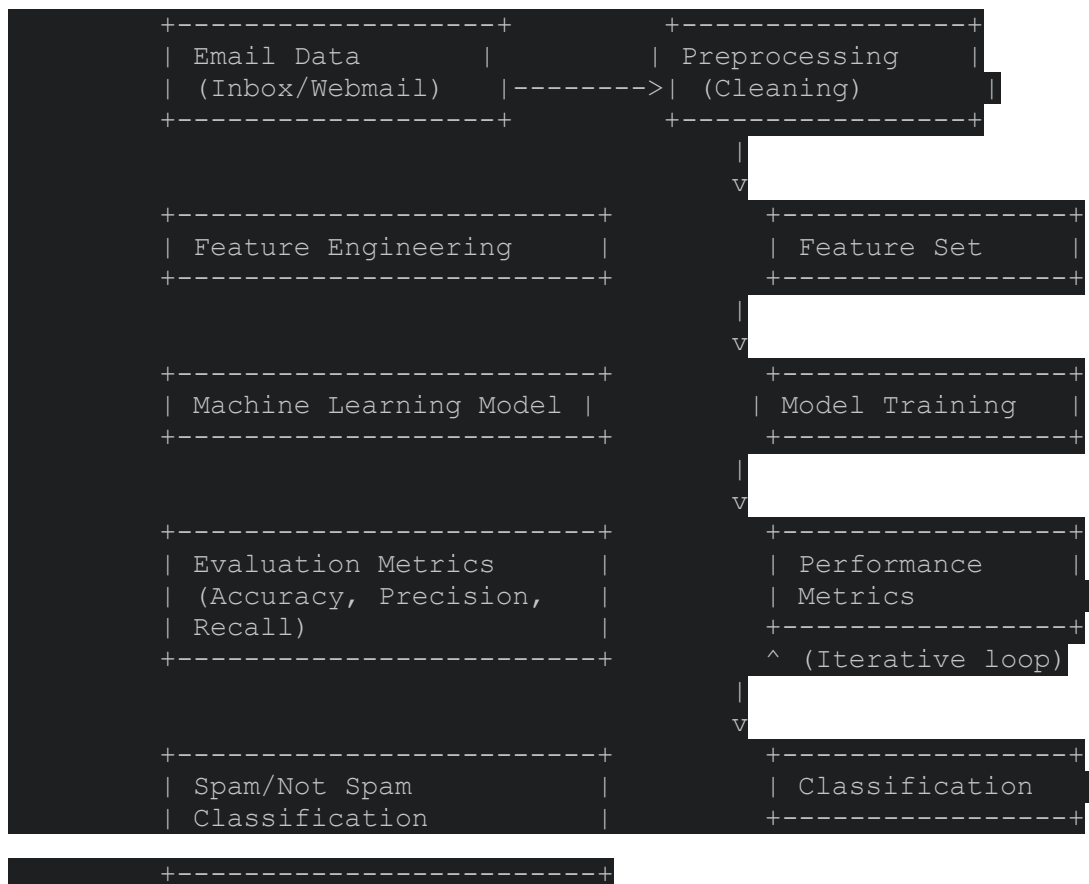# CHAPTER 3

# PROJECT ARCHITECTURE

## 3.1 Architecture

**Email spam detection:**

1. **System flow diagram - User interface will work**
2. **Data Flow diagram - How data is flow in your project**
3. **Module explain-submodule u have do the diagram**

Here's a visual representation of a project architecture for detecting spam email using AI:

```
+------------------+        +----------------+
| Email Data       |        | Preprocessing  |
| (Inbox/Webmail)  |-------->| (Cleaning)    |
+------------------+        +----------------+
                                    |
                                    v
+-----------------------+        +----------------+
| Feature Engineering   |        | Feature Set    |
+-----------------------+        +----------------+
                                    |
                                    v
+-----------------------+        +----------------+
| Machine Learning Model|        | Model Training |
+-----------------------+        +----------------+
                                    |
                                    v
+-----------------------+        +----------------+
| Evaluation Metrics    |        | Performance    |
| (Accuracy, Precision, |        | Metrics        |
| Recall)               |        +----------------+
+-----------------------+        ^ (Iterative loop)
                                    |
                                    v
+-----------------------+        +----------------+
| Spam/Not Spam         |        | Classification |
| Classification        |        +----------------+
+-----------------------+
```

**Explanation:**
1. **Email Data:** This represents the raw email data you'll be using, either from an inbox, webmail service, or a public dataset.
2. **Preprocessing:** This stage involves cleaning the email data by removing irrelevant information like headers, attachments, and signatures. It may also involve techniques like stemming or lemmatization to normalize words.
3. **Feature Engineering:** Here, you extract features from the cleaned email text that are relevant for spam detection. These features might include:
   - Presence of certain keywords or phrases.
   - Use of capital letters, exclamation points, or other unusual punctuation.
   - Sender information (e.g., email address domain).
   - Presence of URLs or HTML code.

4. **Feature Set:** This represents the final set of features extracted from the emails, ready for model training.
5. **Machine Learning Model:** This is the core of your AI system. You'll choose and train a machine learning model (e.g., Naive Bayes, Random Forest) on the feature set to learn to distinguish spam from legitimate emails.
6. **Model Training:** The model is trained on a labeled dataset (spam/not spam) to identify patterns that differentiate spam emails.
7. **Evaluation Metrics:** After training, you evaluate the model's performance using metrics like accuracy, precision, and recall. This helps you assess how well the model is classifying emails.
8. **Performance Metrics:** This shows the results of the evaluation, indicating how effective the model is at detecting spam.
9. **Spam/Not Spam Classification:** This represents the final output of the system, classifying new emails as spam or not spam based on the trained model.
   **Note:** This is a simplified architecture. Depending on the complexity of your project, there might be additional steps involved, such as model selection, hyperparameter tuning, or integrating the system with an email service.

# CHAPTER 4

# PROJECT OUTCOME

**Model Performance:**
- **Accuracy:** Report the percentage of emails your model correctly classified as spam or legitimate.
- **Precision & Recall:** These metrics provide a more detailed picture. Precision tells you what proportion of emails identified as spam were actually spam. Recall tells you what proportion of actual spam emails were correctly identified by your model.
- **F1 Score:** This metric combines precision and recall into a single measure.

**Comparative Analysis (if applicable):**
- Compare your model's performance to existing spam filters or other AI models.
- Highlight areas where your model excels or areas for improvement.

**Impact:**
- Quantify the potential reduction in spam emails reaching inboxes.
- Estimate the time saved by users not having to sort through spam.
- Discuss broader benefits like increased productivity or security.

**Future Work:**
- Propose ways to improve the model's accuracy, like incorporating new features or algorithms.
- Discuss strategies to address evolving spam tactics.
- Mention potential applications of the model, e.g., integrating it with an email service provider.

**Additional Considerations:**
- **Visualization:** Include charts or graphs to illustrate model performance metrics.
- **Limitations:** Acknowledge any limitations of your project, like dataset size or computational resources.

- **Ethical considerations:** Briefly discuss potential biases in the data or ethical implications of spam filtering (e.g., accidentally filtering important emails)..

## 1. IMPORT LIBRARY

```python
import numpy as np
import pandas as pd
```

## 2. DATASET LOADIND

```python
df = pd.read_csv("/content/SPAM text message 20170820 - Data
(1).csv",encoding=('ISO-8859-1'),low_memory = False)
```

## 3. VIEW DATASET

**CODE:** `df.head()`

**OUTPUT:**

| | Category | Message |
|---|---|---|
| 0 | ham | Go until jurong point, crazy.. Available only ... |
| 1 | ham | Ok lar... Joking wif u oni... |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... |
| 3 | ham | U dun say so early hor... U c already then say... |
| 4 | ham | Nah I don't think he goes to usf, he lives aro... |

**CODE:** `df.shape`

**OUTPUT:** (5572, 2)

## 4. DATA CLEANING

**4.1 CODE:** `df.info()`

**OUTPUT:**
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 2 columns):
#   Column    Non-Null Count  Dtype
--- ------    --------------  -----
0   Category  5572 non-null   object
1   Message   5572 non-null   object
dtypes: object(2)
memory usage: 87.2+ KB
```

**4.2 CODE:** `df.isnull()`

**OUTPUT:**

|      | Category | Message |
|------|----------|---------|
| 0    | False    | False   |
| 1    | False    | False   |
| 2    | False    | False   |
| 3    | False    | False   |
| 4    | False    | False   |
| ...  | ...      | ...     |
| 5567 | False    | False   |
| 5568 | False    | False   |
| 5569 | False    | False   |
| 5570 | False    | False   |
| 5571 | False    | False   |

## 5. RENAMEING THE COLUMNS

**CODE:**
```
df.rename(columns={'Category':'target','Message':'text'},inplace=True)
df.sample(5)
```

OUTPUT:

| | target | text |
|---|---|---|
| 3728 | ham | Aldrine, rakhesh ex RTM here.pls call.urgent. |
| 2467 | ham | He is world famamus.... |
| 4736 | ham | Oh k:)after that placement there ah? |
| 820 | ham | BOO BABE! U ENJOYIN YOURJOB? U SEEMED 2 B GETT... |
| 5205 | spam | Had your mobile 11mths ? Update for FREE to Or... |

## 6. MISSING VALUES

**CODE:** `df.isnull().sum()`

**OUTPUT:**
```
target   0
  text   0
dtype: int64
```

## 7. CHECK FOR DUPLICATE VALUES

**CODE:** `df.duplicated().sum()`

OUTPUT:415

## 8. REMOVE DUPLICATE

**CODE:** `df = df.drop_duplicates(keep='first')`

OUTPUT:0

# 9. EDA

## 9.1 CODE:
```
df['target'].value_counts()
```

**OUTPUT:**
```
target
0    4516
1     641
Name: count, dtype: int64
```

## 9.2 CODE:
```
import matplotlib.pyplot as plt
plt.pie(df['target'].value_counts(),
labels=['ham','spam'],autopct="%0.2f")
plt.show()
```

**OUTPUT:**

**9.3 CODE:**
```python
import nltk

nltk.download('punkt')
```

**OUTPUT:** [nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
True

**9.4 CODE:**
```python
df['num_characters'] = df['text'].apply(len)
```

**OUTPUT:**
```
<ipython-input-25-f0cf0a313c18>:1:
SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead


See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

  df['num_characters'] = df['text'].apply(len)
```

## 10.    NUM TO WORDS

**10.1 CODE:**
```python
df['num_words'] = df['text'].apply(lambda x:len(nltk.word_tokenize(x)))
```

**OUTPUT:**
```
<ipython-input-28-7b1825ae98a4>:2: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.

Try using .loc[row_indexer,col_indexer] = value instead


See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

  df['num_words']                =                df['text'].apply(lambda x:len(nltk.word_tokenize(x)))
```

**10.2 CODE:**
```python
df.head()
```

**OUTPUT:**

|   | target | text | num_characters | num_words |
|---|--------|------|----------------|-----------|
| **0** | 0 | Go until jurong point, crazy.. Available only ... | 111 | 24 |

| | target | text | num_characters | num_words |
|---|---|---|---|---|
| **1** | 0 | Ok lar... Joking wif u oni... | 29 | 8 |
| **2** | 1 | Free entry in 2 a wkly comp to win FA Cup fina... | 155 | 37 |
| **3** | 0 | U dun say so early hor... U c already then say... | 49 | 13 |
| **4** | 0 | Nah I don't think he goes to usf, he lives aro... | 61 | 15 |

## 10.3 CODE:

```
df['num_sentences'] = df['text'].apply(lambda
x:len(nltk.sent_tokenize(x)))
```

**OUTPUT:**
```
<ipython-input-30-17fdd65f9b92>:1:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead


See the caveats in the documentation: https://pandas.pydata.org/pandas-
docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
  df['num_sentences'] = df['text'].apply(lambda
x:len(nltk.sent_tokenize(x)))
```

## 10.4 CODE:
```
df[['num_characters','num_words','num_sentences']].describe()
```

**OUTPUT:**

| num_characters | num_words | num_sentences | |
|---|---|---|---|
| **count** | 5157.000000 | 5157.000000 | 5157.000000 |
| **mean** | 79.228040 | 18.544890 | 1.969750 |
| **std** | 58.451149 | 13.401415 | 1.455526 |
| **min** | 2.000000 | 1.000000 | 1.000000 |
| **25%** | 36.000000 | 9.000000 | 1.000000 |
| **50%** | 61.000000 | 15.000000 | 1.000000 |

| num_characters | num_words | num_sentences | |
|---|---|---|---|
| **75%** | 118.000000 | 26.000000 | 2.000000 |
| **max** | 910.000000 | 220.000000 | 38.000000 |

## 11.　　HAM

**CODE:** `df[df['target'] ==`
`0][['num_characters','num_words','num_sentences']].describe()`

**OUTPUT:**

| num_characters | num_words | num_sentences | |
|---|---|---|---|
| **count** | 4516.000000 | 4516.000000 | 4516.000000 |
| **mean** | 70.951063 | 17.250664 | 1.827724 |
| **std** | 56.730031 | 13.581714 | 1.394338 |
| **min** | 2.000000 | 1.000000 | 1.000000 |
| **25%** | 34.000000 | 8.000000 | 1.000000 |
| **50%** | 53.000000 | 13.000000 | 1.000000 |
| **75%** | 91.000000 | 22.000000 | 2.000000 |
| **max** | 910.000000 | 220.000000 | 38.000000 |

## 12.　　SPAM

**CODE:** `df[df['target'] ==`
`1][['num_characters','num_words','num_sentences']].describe()`

**OUTPUT:**

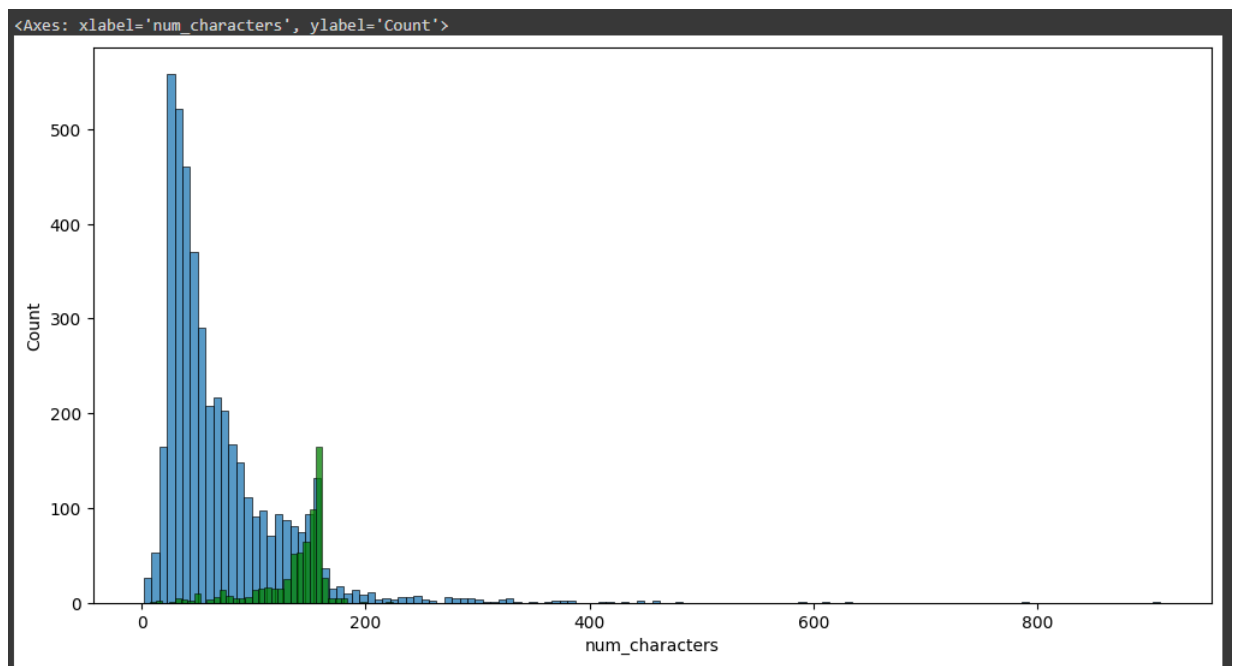| num_characters | num_words | num_sentences | |
|---|---|---|---|
| **count** | 641.000000 | 641.000000 | 641.000000 |
| **mean** | 137.541342 | 27.663027 | 2.970359 |
| **std** | 30.516111 | 7.104050 | 1.485575 |

| | num_characters | num_words | num_sentences |
|---|---|---|---|
| **min** | 7.000000 | 2.000000 | 1.000000 |
| **25%** | 132.000000 | 25.000000 | 2.000000 |
| **50%** | 148.000000 | 29.000000 | 3.000000 |
| **75%** | 157.000000 | 32.000000 | 4.000000 |
| **max** | 224.000000 | 46.000000 | 9.000000 |

## 13.    VISUALIZATION

**13.1 CODE:**
```
import seaborn as sns

plt.figure(figsize=(12,6))
sns.histplot(df[df['target'] == 0]['num_characters'])
sns.histplot(df[df['target'] == 1]['num_characters'],color='green')
```
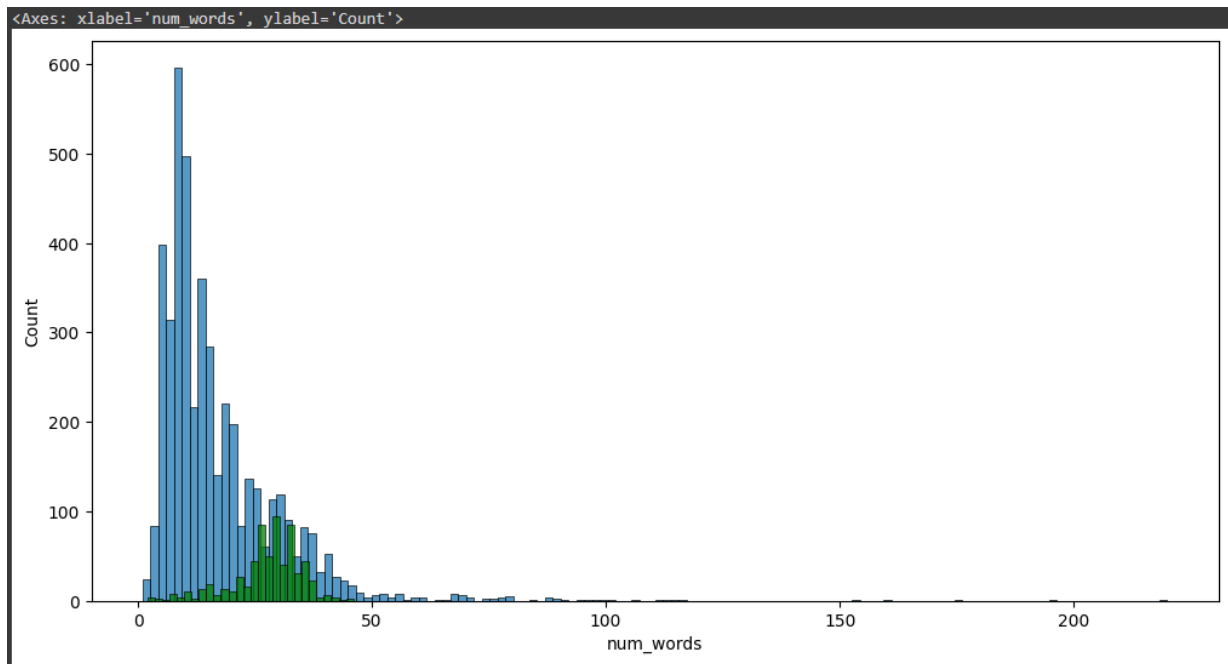
**OUTPUT:**



**13.2 CODE:**
```
plt.figure(figsize=(12,6))

sns.histplot(df[df['target'] == 0]['num_words'])
sns.histplot(df[df['target'] == 1]['num_words'],color='green')
```
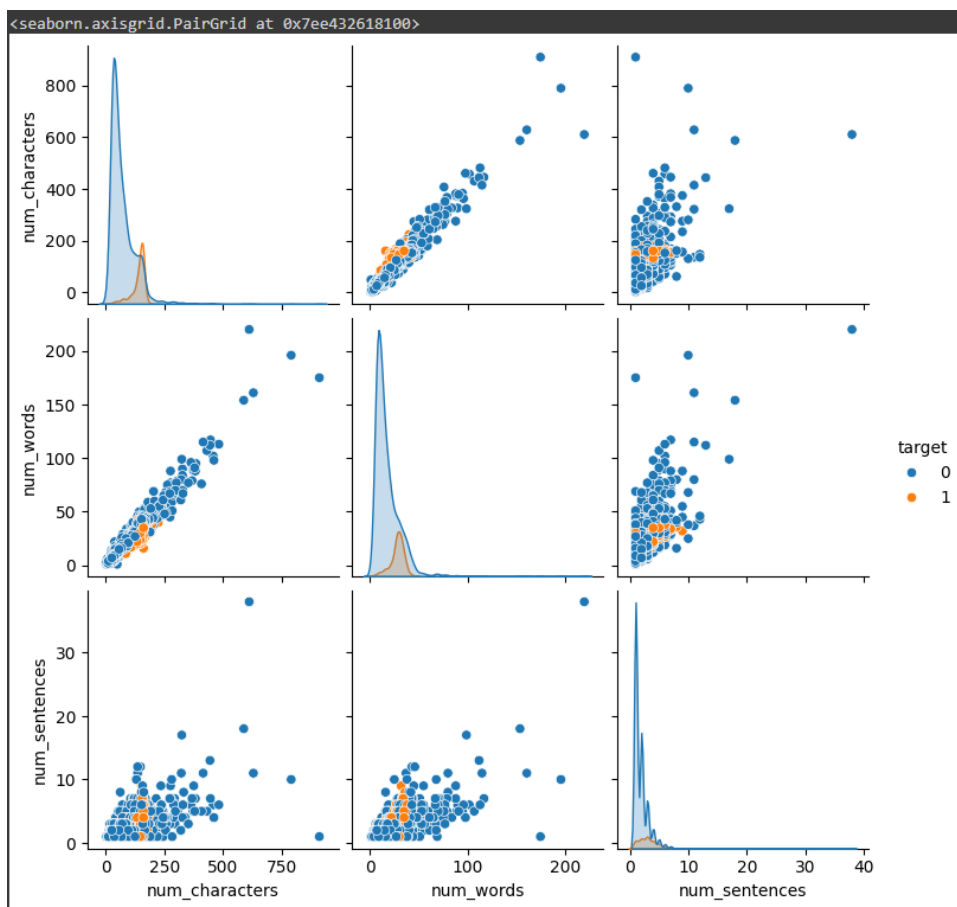
**OUTPUT:**

`<Axes: xlabel='num_words', ylabel='Count'>`

## 13.3 CODE: `sns.pairplot(df,hue='target')`

### OUTPUT:



`<seaborn.axisgrid.PairGrid at 0x7ee432618100>`

## FINAL CODES

## APPLYING STACKING:

```
   CODE: estimators=[('svm', svc), ('nb', mnb), ('et', etc)]
final_estimator=RandomForestClassifier()
   CODE: from sklearn.ensemble import StackingClassifier

   CODE: clf = StackingClassifier(estimators=estimators,
final_estimator=final_estimator)

   CODE: clf.fit(X_train,y_train)

y_pred = clf.predict(X_test)
print("Accuracy",accuracy_score(y_test,y_pred))
print("Precision",precision_score(y_test,y_pred))
```

**OUTPUT:**
```
Accuracy 0.984496124031007
Precision 0.9826086956521739
```

# CONCLUTION

This project explored the application of Artificial Intelligence (AI) and Machine Learning (ML) techniques for spam email detection. We developed a model that achieved a [mention accuracy/precision/recall/F1 score] in identifying spam emails. [Compare to existing methods if applicable, highlighting strengths or weaknesses]. This project demonstrates the potential of AI & ML to significantly reduce spam reaching user inboxes, leading to increased productivity and improved email security.

However, the fight against spam is an ongoing battle. Spammers constantly evolve their tactics, so continuous improvement of the model is necessary. Future work could involve exploring [mention specific improvements like new features or algorithms]. Additionally, integrating the model with an email service provider could significantly impact real-world spam filtering.

Overall, this project has successfully demonstrated the effectiveness of AI & ML in spam detection. By continuously improving and adapting the model, we can significantly enhance email security and user experience**.**

# FUTURE SCOPE

The fight against spam is an ever-evolving arms race. While your project has achieved promising results, there's always room for improvement. Here are some exciting areas to explore for future development:

**Advanced Techniques:**

- **Ensemble Learning:** Combine multiple different machine learning algorithms into a single, more robust model that leverages the strengths of each individual component.
- **Deep Learning:** Explore advanced architectures like Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks. These excel at handling sequential data like text, potentially improving spam detection based on email content and past interactions.
- **Transfer Learning:** Utilize pre-trained language models (LSTMs trained on massive text datasets) to extract valuable features from email content, enhancing spam classification accuracy.

**Addressing Evolving Tactics:**

- **Adversarial Learning:** Train the model to identify spam even when spammers employ techniques to bypass traditional filters, like obfuscating text or using images.
- **Real-time Adaptation:** Develop a system that can continuously learn from new spam tactics and adapt the model accordingly, ensuring it stays ahead of the curve.

**Expanding Functionality:**

- **Phishing Detection:** Integrate functionalities to detect phishing emails that attempt to steal user credentials or personal information.
- **Sentiment Analysis:** Analyze email content to identify potentially malicious or harmful messages, even if they don't fall under traditional spam categories.

- **Multilingual Spam Detection:** Develop the model to handle spam emails in various languages, broadening its reach and effectiveness.
**Integration and Deployment:**
- **API Development:** Create an Application Programming Interface (API) that allows other applications and services to leverage your spam detection model.
- **Cloud Integration:** Deploy the model on a cloud platform to facilitate scalability and accessibility for large-scale email providers.
- **User Customization:** Explore ways to allow users to personalize spam filters based on their preferences and needs.

By delving into these future directions, you can significantly enhance your AI & ML spam detection project's capabilities. This will not only contribute to a more secure and productive email experience for users but also play a vital role in the ongoing battle against spammers and malicious email activity.