

SMS Spam Classification

Presented by Alfiyyah Ajeng Nurardita - SI602003



Table of contents

01

Introduction

02

Methodology

03

Model Evaluation

04

Conclusion

Introduction

Problem Identification

E-commerce platforms often face spam in product reviews, user messages, and customer support tickets.

This leads to poor **user experience**, reduced **customer trust**, and potential security risks such as **phishing** and **scams**.

Solution

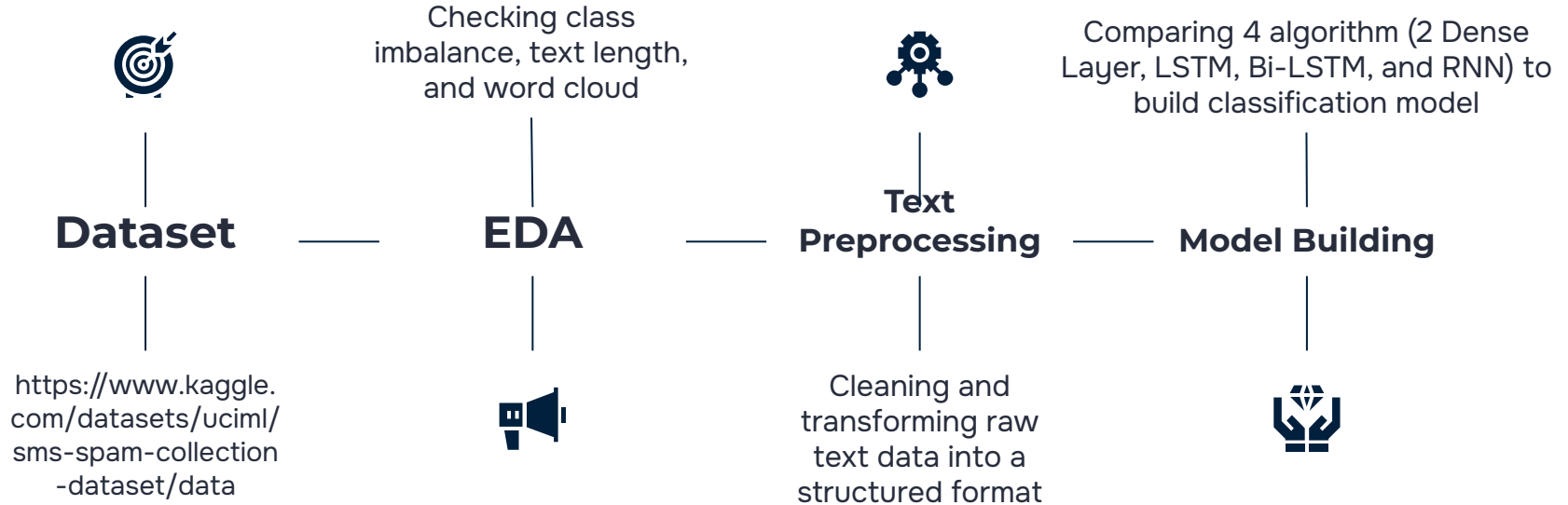
A machine learning-based **spam classification system** can automatically distinguish between spam and legitimate content, ensuring only relevant and safe information reaches users.

Objective

The project aims to:

1. Develop a reliable spam vs. ham classification model.
2. Apply effective text preprocessing to improve accuracy.
3. Evaluate the model using precision, recall, and F1-score.

Methodology





Dataset

No	v1	v2
1	ham	Go until jurong point, crazy.. Available only ...
2	ham	Ok lar... Joking wif u oni...
3	spam	Free entry in 2 a wkly comp to win FA Cup fina...
4	ham	U dun say so early hor... U c already then say...
5	ham	Nah I don't think he goes to usf, he lives aro...



Data Cleaning

1. Rename columns name

v1 -> Target

v2 -> Text

```
1 df.rename(columns={'v1': 'Target', 'v2': 'Text'}, inplace=True)
```

2. Checking missing values

```
1 #checking missing values
2 df.isna().sum()
```

```
Target    0
Text      0
dtype: int64
```

3. Checking and drop duplicates values

```
1 #checking duplicate values
2 df.duplicated().sum()
```

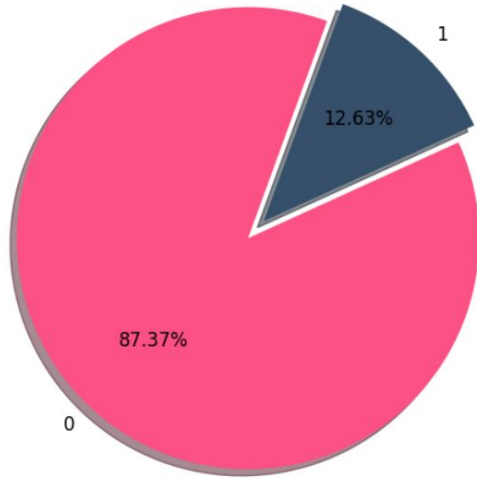
```
403
```

```
1 #drop duplicate values
2 df = df.drop_duplicates()
```

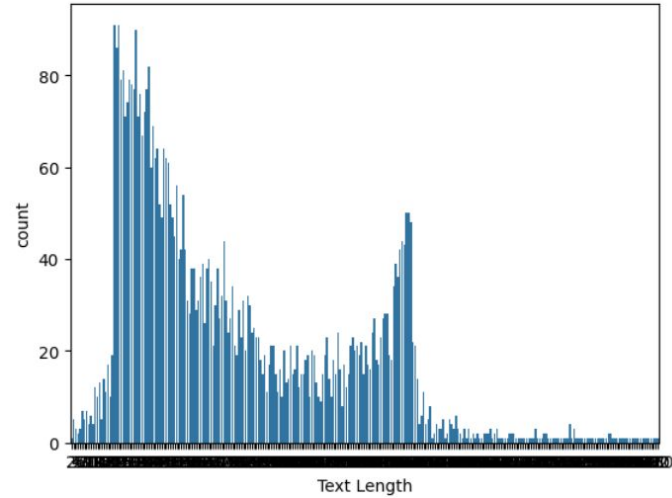
4. Encoding (ham -> 0, spam -> 1)

```
1 #encoding target column
2 from sklearn.preprocessing import LabelEncoder
3 encoder = LabelEncoder()
4 df['Target'] = encoder.fit_transform(df['Target'])
```

Exploratory Data Analysis (EDA)

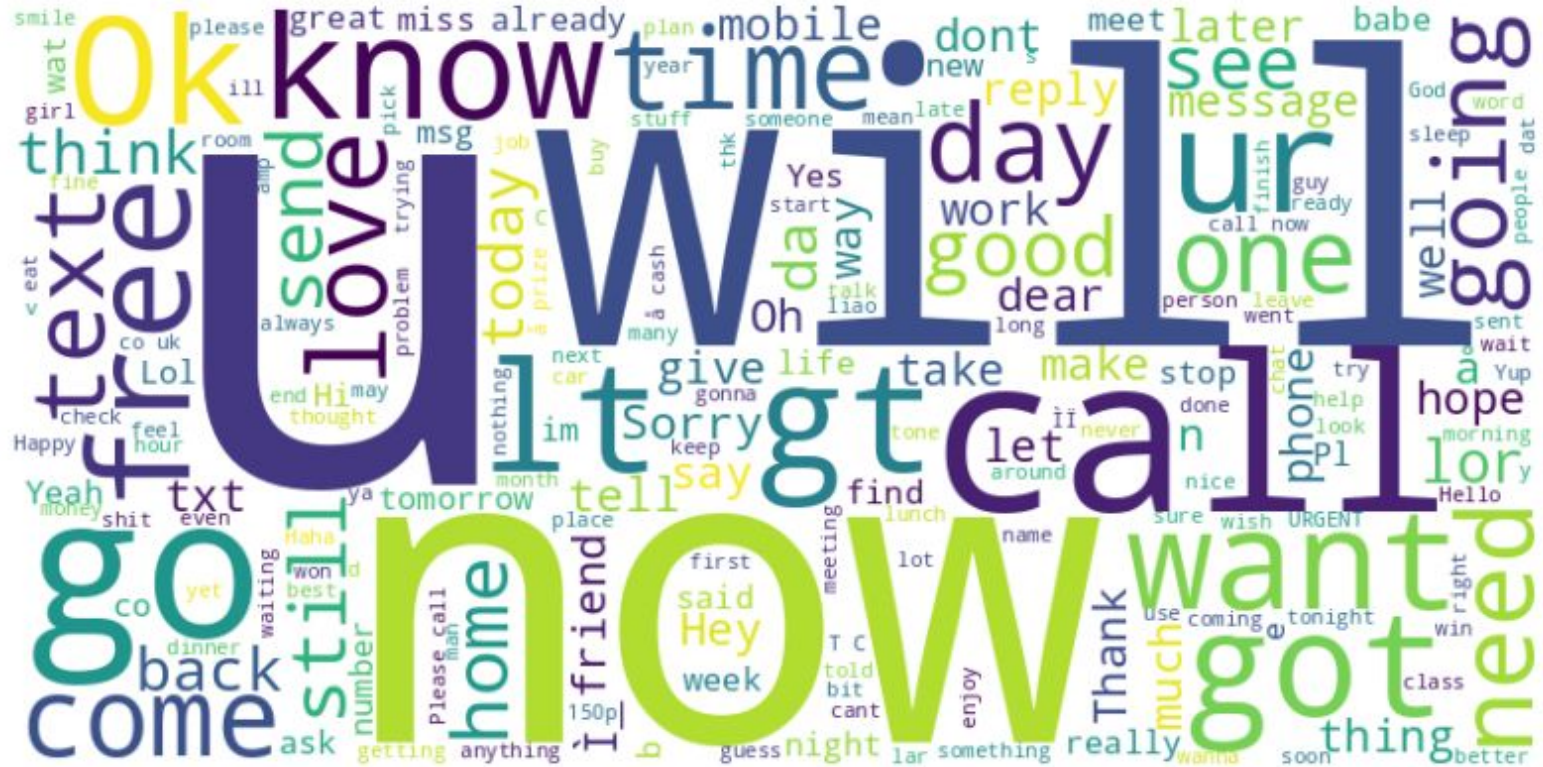


Data points 0 (ham) are 6.9 times more than points with 1 (spam) labels, which leads to **Imbalanced dataset**



We have 408237 words in our Dataframe the average word count in every sentence is 78

Word Cloud



Text Preprocessing

```
4 # Case folding
5 text = text.lower()
6
7 # Mention removal
8 text = re.sub("@[A-Za-z0-9_]+", " ", text)
9
10 # Hashtags removal
11 text = re.sub("#[A-Za-z0-9_]+", " ", text)
12
13 # Newline removal (\n)
14 text = re.sub(r"\\n", " ", text)
15
16 # Whitespace removal
17 text = text.strip()
18
19 # URL removal
20 text = re.sub(r"http\S+", " ", text)
21 text = re.sub(r"www.\S+", " ", text)
22
23 # Non-letter removal (such as emoticon, symbol (like μ, $, π), etc
24 text = re.sub("[^A-Za-z\s]", " ", text)
25
26 # Tokenization
27 tokens = word_tokenize(text)
28
29 # Stopwords removal
30 tokens = [word for word in tokens if word not in stopwords_id]
31
32 # Stemming
33 tokens = [stemmer.stem(word) for word in tokens]
34
35 # Combining Tokens
36 text = ' '.join(tokens)
```

1. **text.lower** → Converts all text to lowercase to maintain consistency.
2. **Removal (mention, hashtag, newline, whitespace, URL, non-letter, stopwords)** → Cleans the text by removing symbols, links, common words, and irrelevant characters.
3. **Tokenization** → Splits text into smaller units (tokens/words) for analysis.
4. **Stemming** → Reduces words to their root form by removing affixes.

Model Building

2 Dense Layers

The model uses a TextVectorization + Embedding layer, followed by GlobalAveragePooling and Flatten. It has two Dense layers: one hidden layer with 32 ReLU units and an output layer with 1 sigmoid unit for binary classification. Compiled with Adam optimizer and Binary Crossentropy loss.

LSTM

The model uses an Embedding layer (100 dimensions, input length 100), followed by an LSTM layer with 64 units, and a Dense output layer with 1 sigmoid unit for binary classification. It is compiled with Adam optimizer and Binary Crossentropy loss, trained for 5 epochs with batch size 32.

Bi-LSTM

The model uses a TextVectorization + Embedding layer, followed by two Bidirectional LSTM layers (first with return_sequences=True, both with 64 units, tanh activation). The output is passed through Flatten → Dropout (0.1) → Dense(32, ReLU) → Dense(1, Sigmoid) for binary classification. It is compiled with Adam optimizer and Binary Crossentropy loss, trained for 5 epochs.

LSTM

The model consists of an Embedding layer (64 dimensions, input length 100), followed by a SimpleRNN layer with 64 units, and a Dense output layer with 1 sigmoid unit for binary classification. It is compiled with Adam optimizer and Binary Crossentropy loss, trained for 5 epochs with batch size 32.

Model Evaluation

Model	Loss	Accuracy
2 Dense Layer	0.571600	0.979691
LSTM	0.347461	0.890716
Bi-LSTM	0.067345	0.983559
Simple RNN	0.354138	0.890716

Conclusion

- **metrics** : from all models, Bi-LSTM have a greatest accuracy which 98%.
- **problem** : the models use imbalanced dataset, most of data points have the label of ham.
For this case accuracy can't be a good metrics

Thank You !

Full Code : [colab](#)

Lets Connect



<https://www.linkedin.com/in/ajeng-nurardita/>



<https://github.com/alfiyyahajeng>