

Spam Email Classifier Using Machine Learning

1. Project Title

Spam Email Classifier using Natural Language Processing and Logistic Regression

2. Objective

To develop a machine learning-based spam detection system that can classify messages as spam or ham (not spam) using natural language processing (NLP) techniques and a logistic regression classifier.

3. Dataset

- Source: SMS Spam Collection Dataset
- Attributes:
 - label: Classification as 'spam' or 'ham'
 - text: The actual SMS message content
- Size: 5,574 records

4. Tools & Technologies Used

Language: Python

Libraries: pandas, nltk, scikit-learn

Techniques:

- Text Preprocessing
- TF-IDF Vectorization
- Logistic Regression

- Model Evaluation Metrics

5. Methodology

Step 1: Data Loading

- Read the CSV file and selected only the relevant columns (label, text).
- Renamed columns for clarity.

Step 2: Text Preprocessing

- Lowercasing the text
- Removing punctuation
- Removing stopwords
- Stemming using PorterStemmer

Step 3: Feature Extraction

- Applied TF-IDF vectorization to convert text data into numerical vectors.

Step 4: Train-Test Split

- Dataset split into 80% training and 20% testing sets.

Step 5: Model Training

- Trained a Logistic Regression classifier on the TF-IDF features.

Step 6: Evaluation

- Used accuracy, confusion matrix, and classification report to assess performance.

6. Results

- Accuracy: 95.07%
- Evaluation Metrics:
 - High precision and recall for both spam and ham categories
 - Balanced performance indicating good generalization

7. Conclusion

The spam classifier developed using logistic regression and NLP techniques achieved 95.07% accuracy on unseen data. This demonstrates the effectiveness of the model in detecting spam messages. The project serves as a strong foundation for deploying spam filters in communication systems.

8. Future Enhancements

Experiment with other classifiers like Naive Bayes and SVM

Explore deep learning models such as LSTM for improved context understanding

Integrate the classifier into real-time email or SMS platforms