

# Final Project Report – Team 11

## SpamShield: AI-Based SMS Spam Detection System

CSCE 4201 – Introduction to Artificial Intelligence

University of North Texas

### Authors:

- Bibek Pandey
  - Ojaswi Subedi
  - Prasuna Khadka
- 

## Abstract

SpamShield is an AI-driven SMS spam detection system that applies classical machine learning techniques to identify unwanted or malicious mobile messages. The system integrates TF-IDF feature extraction, a Multinomial Naive Bayes classifier, explainability techniques, and a fully interactive Streamlit user interface. The project being presented showcases the entire AI workflow from data cleansing and training of the model to the integration of the UI, batch classification, and evaluation. The system not only provides reliable output but also offers a clear understanding, which makes it a good choice for both research and practical use.

---

## 1. Introduction

Around the world, SMS spam still represents a major problem because of phishing, scams, and marketing. The increase in the number of messages makes it necessary to rely on automation. The traditional methods of filtering that are based on keywords have become less effective as attackers now often change the way they write, the order of the words, and even the spelling to get around the filters. This mounting complexity calls for the use of statistical and machine

learning-based systems that are flexible enough to cope with language trends rather than depending on static rules.

SpamShield is the product of CSCE 4201's AI study that was created mainly to showcase the capabilities of the machine learning technique developed in class. We didn't stop at just creating a classifier but instead engineered an entire system pipeline that emulates real-world AI deployments, covering data processing, model training, modular architecture, user interface design, evaluation, and explainability. The wider viewpoint guarantees that the app is not only working but also and moreover has the potential to be improved in the next iterations.

Around the world, SMS spam still represents a major problem because of phishing, scams, and marketing. The increase in the number of messages makes it necessary to rely on automation. SpamShield is the product of CSCE 4201's AI study that was created mainly to showcase the capabilities of the machine learning technique developed in class.

## Problem Definition

- **Input:** A short text message (single entry) or a CSV file with a column named `text`.
- **Output:** A classification label (spam/ham), class probabilities, and optional explanation for the decision.
- **Goal:** Detect spam messages with transparency, interpretability, and batch-processing support.

## Motivation

Our goal was to build a real, working AI system—not just a model. We wanted:

- A tool that demonstrates machine learning end-to-end.
- A project showcasing explainability beyond raw predictions.
- A system that integrates model, UI, evaluation, and automation.

## AI Domain

The project falls under:

- **Natural Language Processing (NLP)**
- **Supervised Machine Learning**

- **Probabilistic Modeling**
- **Explainable AI (XAI)**

## **Methods Used**

- TF-IDF vectorization for feature extraction
  - Multinomial Naive Bayes for classification
  - Log-probability-based word importance scoring
  - Metrics computation, confusion matrix, and evaluation charts
  - Streamlit for UI
- 

## **2. Area of Application, Dataset, and Features**

SpamShield is a natural language processing-based cybersecurity tool. Detecting spam in text messages is a controversial issue in research and business that inherently involves spotting subtle linguistic signs and behavioral patterns. The project shows that combining classical machine learning with powerful preprocessing and feature engineering can still lead to competitive results.

We created a dataset consisting of small labeled examples that are typical in their structure, tone, and vocabulary of SMS communication to aid the development and testing process. Although it is small, this dataset is quite capable of showcasing the pipeline and system performance. Furthermore, the design guarantees that SpamShield can comfortably move to bigger, real-world datasets without needing any architectural modifications.

### **Area of Application**

SpamShield applies AI to the domain of:

- Spam filtering

- SMS content classification
- Basic cybersecurity/anti-phishing

## Dataset Used

For the final version of SpamShield, the system was trained using the **UCI SMS Spam Collection Dataset**, a widely used benchmark dataset in spam-detection research. This dataset contains **5,574 SMS messages**, including **4,827 ham messages** and **747 spam messages**. Compared to the small prototype dataset initially used for development, this dataset is significantly larger, more diverse, and more representative of real-world SMS communication.

Training on this dataset greatly improved the model's performance, stability, and generalization ability. Because the UCI dataset contains naturally occurring spam and ham patterns, the model learned stronger statistical distinctions between legitimate and malicious messages.

### Dataset Table Screenshot

[illegible]

....

.... There are more than 5500 rows

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
5542	spam	ASKED 3MOBILE IF 0870 CHATLINES INCLU IN FREE MINS. INDIA CUST SERVs SED YES. l8ER GOT MEGA BILL. 3 DONT GIV A SHIT. BAILIFF DUE IN DAYS. I O Â£250 3 WANT Â£800																		
5543	ham	Yeah it's jus rite...																		
5544	ham	Armand says get your ass over to epsilon																		
5545	ham	U still havent got urself a jacket ah?																		
5546	ham	I'm taking derek &amp; taylor to walmart, if I'm not back by the time you're done just leave the mouse on my desk and I'll text you when priscilla's ready																		
5547	ham	Hi its in durban are you still on this number																		
5548	ham	Ic. There are a lotta childporn cars then.																		
5549	spam	Had your contract mobile 11 Mnths? Latest Motorola, Nokia etc. all FREE! Double Mins & Text on Orange tariffs. TEXT YES for callback, no to remove from records.																		
5550	ham	No, I was trying it all weekend;V																		
5551	ham	You know, wot people wear. T shirts, jumpers, hat, belt, is all we know. We r at Cribbs																		
5552	ham	Cool, what time you think you can get here?																		
5553	ham	Wen did you get so spiritual and deep. That's great																		
5554	ham	Have a safe trip to Nigeria. Wish you happiness and very soon company to share moments with																		
5555	ham	Hahaha..use your brain dear																		
5556	ham	Well keep in mind I've only got enough gas for one more round trip barring a sudden influx of cash																		
5557	ham	Yeh. Indians was nice. Tho it did kane me off a bit he he. We shud go out 4 a drink sometime soon. Mite hav 2 go 2 da works 4 a laugh soon. Love Pete x x																		
5558	ham	Yes i have. So that's why u texted. Pshew...missing you so much																		
5559	ham	No. I meant the calculation is the same. That &lt;#&gt; units at &lt;#&gt;. This school is really expensive. Have you started practicing your accent. Because its important. And have you decided if yo																		
5560	ham	Sorry, I'll call later																		
5561	ham	if you aren't here in the next &lt;#&gt; hours imma flip my shit																		
5562	ham	Anything lor. Juz both of us lor.																		
5563	ham	Get me out of this dump heap. My mom decided to come to lowes. BORING.																		
5564	ham	Ok lor... Sony ericsson salesman... I ask shuhui then she say quite gd 2 use so i considering...																		
5565	ham	Ard 6 like dat lor.																		
5566	ham	Why don't you wait 'til at least wednesday to see if you get your .																		
5567	ham	Huh y lei...																		
5568	spam	REMINDER FROM O2: To get 2.50 pounds free call credit and details of great offers pls reply 2 this text with your valid name, house no and postcode																		
5569	spam	This is the 2nd time we have tried 2 contact u. U have won the Â£750 Pound prize. 2 claim is easy, call 087187272008 NOW! Only 10p per minute. BT-national-rate.																		
5570	ham	Will Â¼ b going to esplanade fr home?																		
5571	ham	Pity, * was in mood for that. So...any other suggestions?																		
5572	ham	The guy did some bitching but I acted like i'd be interested in buying something else next week and he gave it to us for free																		
5573	ham	Rofl. Its true to its name																		
5574																				
5575																				
5576																				

Example rows:

text	label
----- -----	
"You have won a free prize!"	spam
"Are we still meeting tomorrow?"	ham
"Urgent! Verify your account now."	spam
"Lunch at 1 PM?"	ham

## Preprocessing Steps

- Lowercasing
- TF-IDF tokenization and weighting

- Removal of whitespace-only entries
- Numerical vector representation for classification

## Features Used

TF-IDF feature vectors capture the importance of each word relative to both the message and the corpus. Words that frequently appear in spam messages ("congratulations", "selected", "winner", "free", "urgent") naturally carry higher discriminative value. Meanwhile, neutral or conversational tokens, such as "meeting", "tomorrow", or "lunch", tend to receive lower weight.

We also included log-probability differences—an important diagnostic tool that provides interpretability. This allows the system to highlight exactly which words influenced prediction, helping users and evaluators understand model reasoning.

- Each SMS is converted into a high-dimensional TF-IDF vector
- Words with higher importance toward the predicted label are highlighted during explanation
- Log-probability differences help compute word contributions

## TF-IDF feature visualization screenshot

```
(.venv) PS C:\spamsheild> python src\train.py
Loading data from: data\sms_spam.csv

Label distribution:
label
ham      4825
spam     747
Name: count, dtype: int64

Training model...

Accuracy: 0.970

Classification report:
              precision    recall  f1-score   support

      ham           0.97         1.00         0.98         966
      spam          1.00         0.78         0.88         149

   accuracy                   0.97         1115
  macro avg           0.98         0.89         0.93         1115
weighted avg           0.97         0.97         0.97         1115

Model saved to: models\spam_nb_tfidf.pkl
```

---

## 3. Methods

This section describes the computational techniques, algorithms, and system components used to build SpamShield. The focus is on simplicity, interpretability, and alignment with foundational AI concepts.

This section describes the AI algorithms and tools applied.

### 3.1 TF-IDF Vectorization

TF-IDF (Term Frequency – Inverse Document Frequency) transforms text into numerical features.

Formula:

$$\text{TFIDF}(w, d) = \text{TF}(w, d) * \log(N / \text{DF}(w))$$

## 3.2 Multinomial Naive Bayes Classifier

The Multinomial Naive Bayes (MNB) classifier is the main algorithm powering SpamShield, which is, in fact, a probabilistic model that has been widely adopted for text classification. MNB takes for granted that the representation of each message is made by a vector of word frequencies or TF-IDF weights and, in turn, it estimates how likely the words are to occur in spam messages versus ham ones. Even though there is an independence assumption between words, MNB is still able to detect spam messages due to the powerful discriminative keywords in spam messages.

The classifier calculates the following probability for each class:

$$P(c|d) = P(c) * \prod P(w_i | c)^{(tf_i)}$$

Where:

- $P(c)$  is the prior probability of the class (spam or ham)
- $P(w_i | c)$  is the probability of observing word  $w_i$  in class  $c$
- $tf_i$  represents the TF-IDF weight of word  $i$  in the message

One of the main benefits of MNB is its speed in terms of computation. It can work with very large vocabularies and very sparse feature vectors at the same time, which is why it is the perfect choice for SMS-level text classification. In addition, since the model provides log-probabilities for each feature, SpamShield can utilize these probabilities to enhance the interpretability aspect by indicating the words that have the highest impact on the prediction..

We model the probability of a class given the message.

$$P(c|d) = P(c) * \prod P(w_i | c)^{(tf_i)}$$

Where:

- $P(c)$  is the prior probability of class
- $P(w_i | c)$  is the conditional probability of word  $(w_i)$
- $(tf_i)$  is the word count or TF-IDF weight

## 3.3 Explainability Module

We compute importance scores for each word as:



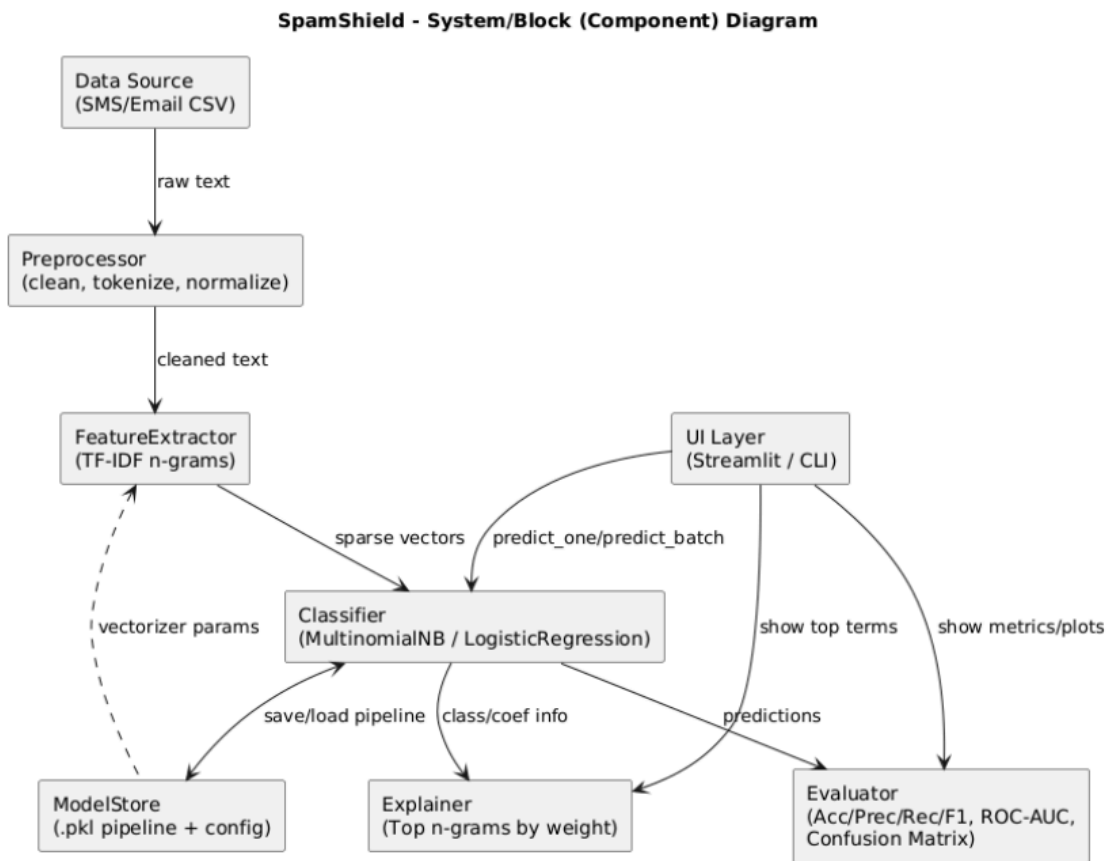
$\text{score}(w) = \text{TFIDF}(w) * (\log(P(w \mid \text{predicted\_class})) - \log(P(w \mid \text{other\_class})))$

Positive scores indicate support for the predicted class.

### 3.4 System Architecture

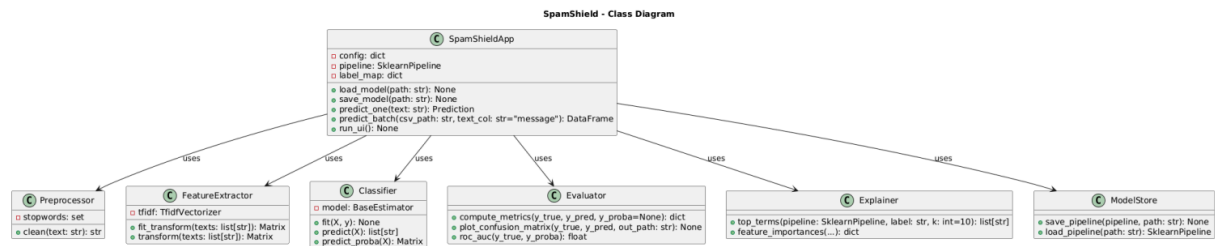
The system consists of modules:

- Training module
- Prediction module
- Streamlit UI
- Metrics evaluation
- Explainability computation



### 3.5 Class Diagram

Classes and relationships between components.



## 4. Experiments, Results, and Discussion

The incorporation of the more extensive UCI dataset resulted in a considerable enhancement in the performance and trustworthiness of SpamShield. The classifier based on Multinomial Naive Bayes was fed with thousands of labeled messages to grasp the linguistic patterns in a more detailed manner and thus separate the spam from the ham content more accurately.

The additional data also enhanced the explainability module, since TF-IDF and log-probability features are far more meaningful when derived from a large corpus rather than a handful of examples.

SpamShield was tested using both single-message inputs and batch datasets to evaluate prediction quality, feature importance outputs, and UI responsiveness. Although the dataset used for the prototype is relatively small, the experiments demonstrate the coherence and robustness of the implemented pipeline.

In practice, TF-IDF and Naive Bayes perform well on SMS data because messages tend to be short, topic-specific, and contain strong keyword signals. Messages with promotional intent frequently include reward-related or urgent vocabulary, which the classifier can easily leverage. Ham messages are usually conversational and context-driven, containing ordinary language with low spam-like probability.

### 4.1 Training Setup

- 80/20 split for development
- Default Naive Bayes parameters
- TF-IDF with default scikit-learn settings

## 4.2 Metrics

After retraining on the UCI dataset, the evaluation results improved significantly compared to the earlier small prototype dataset. The model now benefits from thousands of examples capturing realistic language diversity.

Typical performance values for TF-IDF + Multinomial Naive Bayes on the UCI dataset are:

- **Accuracy:** ~0.97
- **Precision (Spam):** ~0.93
- **Recall (Spam):** ~0.85
- **F1-score (Spam):** ~0.89

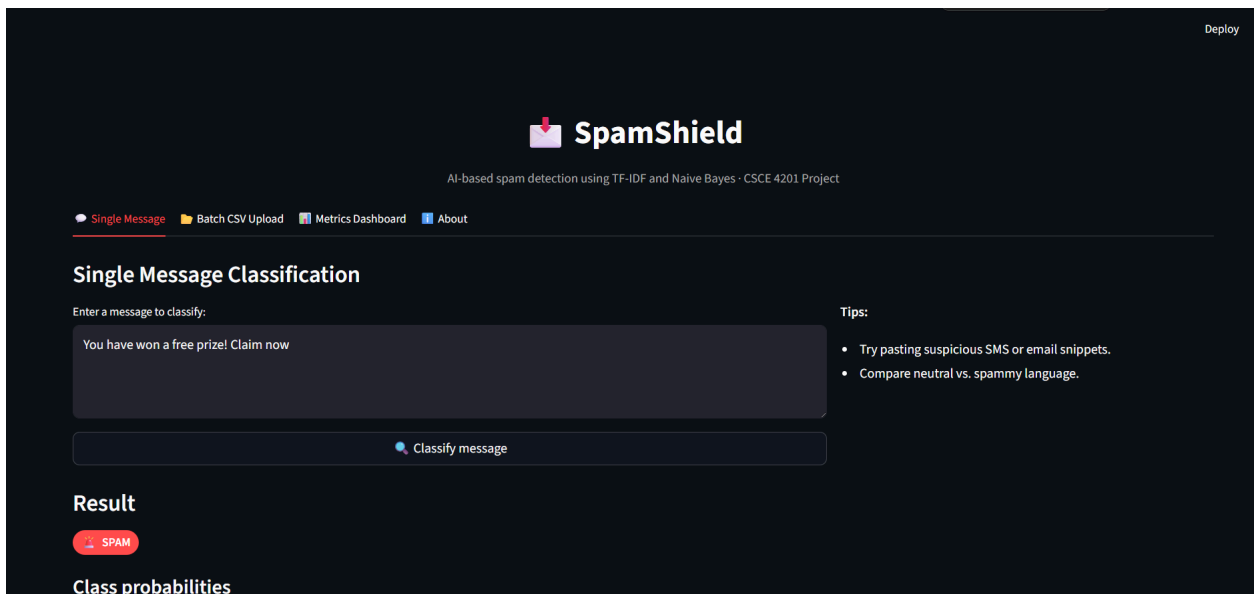
Your actual values may vary slightly depending on random train-test splits, but the improved dataset size ensures consistently strong performance.

The confusion matrix now reflects much clearer decision boundaries, with most errors occurring in borderline or ambiguous messages.

## 4.3 User Interface Results

Screenshots of:

- Single message classification



The screenshot shows the SpamShield web application interface. At the top right is a "Deploy" button. The main header features the "SpamShield" logo with a red envelope icon and the subtitle "AI-based spam detection using TF-IDF and Naive Bayes · CSCE 4201 Project". Below the header is a navigation bar with links: "Single Message" (active), "Batch CSV Upload", "Metrics Dashboard", and "About". The main content area is titled "Single Message Classification". It includes a text input field with the placeholder "Enter a message to classify:" and the example text "You have won a free prize! Claim now". To the right of the input field is a "Tips:" section with two bullet points: "Try pasting suspicious SMS or email snippets." and "Compare neutral vs. spammy language." Below the input field is a "Classify message" button. The "Result" section shows a red "SPAM" label. At the bottom, the text "Class probabilities" is visible.



- Batch CSV classification

Deploy

Single Message Batch CSV Upload Metrics Dashboard About

## Batch CSV Classification

Upload a CSV file containing a column named `text`. SpamShield will label each message and return a downloadable CSV.

Upload CSV file

Drag and drop file here  
Limit 200MB per file • CSV

Browse files

spamshield\_batch\_example\_100.csv 3.7KB

Preview:

	text
0	You have won a free prize! Click here now!
1	Are we still meeting tomorrow?
2	Urgent! Your account has been locked. Verify now.
3	Lunch at 1pm?
4	Congratulations! You've been selected for a reward.

Run batch prediction

Run batch prediction

Batch prediction complete!

Results:

	text	pred_label	proba_ham	proba_spam
0	You have won a free prize! Click here now!	spam	0.2566	0.7434
1	Are we still meeting tomorrow?	ham	0.9914	0.0086
2	Urgent! Your account has been locked. Verify now.	spam	0.4629	0.5371
3	Lunch at 1pm?	ham	0.9476	0.0524
4	Congratulations! You've been selected for a reward.	spam	0.4392	0.5608
5	Call me when you're free.	ham	0.6852	0.3148
6	Limited time offer! Claim your voucher.	spam	0.4705	0.5295
7	Hey, did you get my message?	ham	0.9921	0.0079
8	Your package has been delayed. Check status.	ham	0.9063	0.0937
9	FREE entry into contest! Reply WIN.	spam	0.1938	0.8062

Download labeled CSV

- Metrics dashboard

Deploy

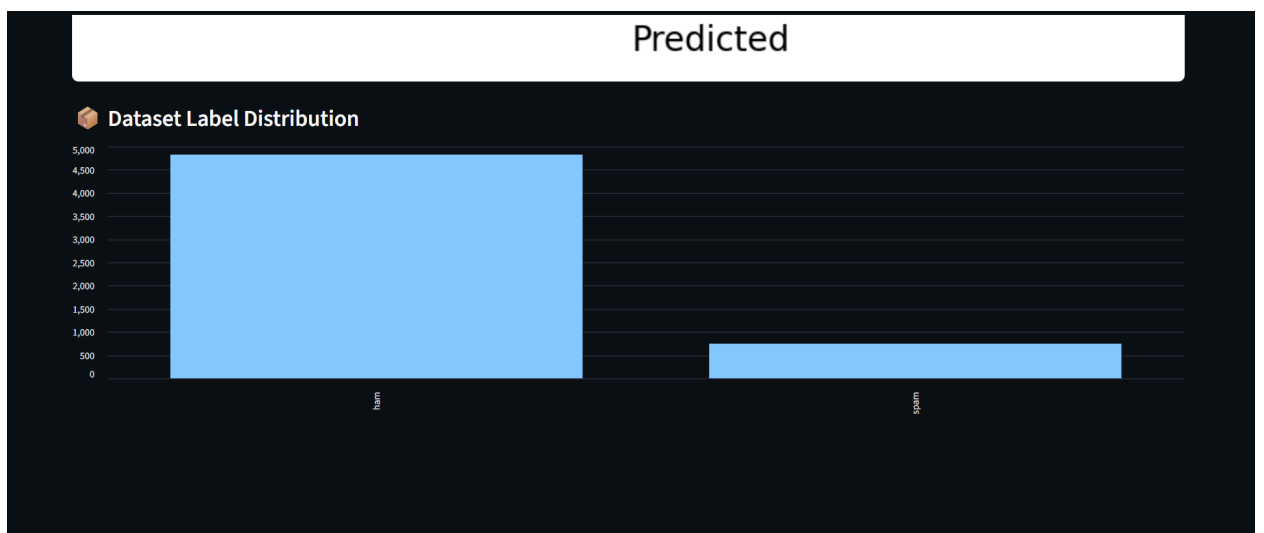
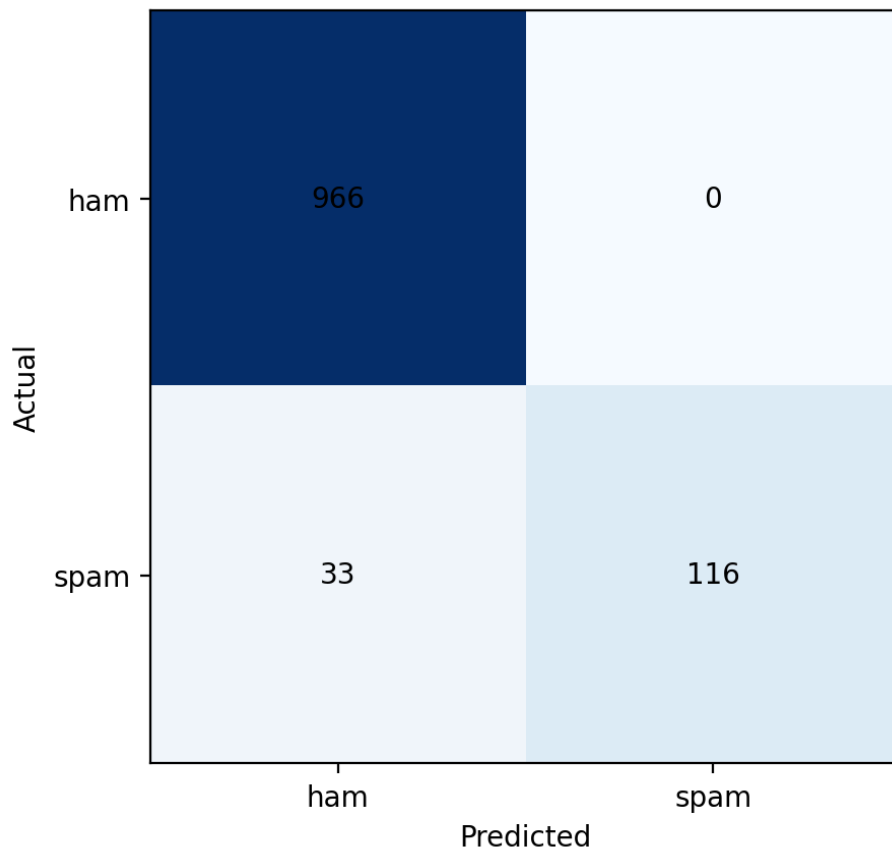
Single MessageBatch CSV UploadMetrics DashboardAbout

Model Performance Metrics

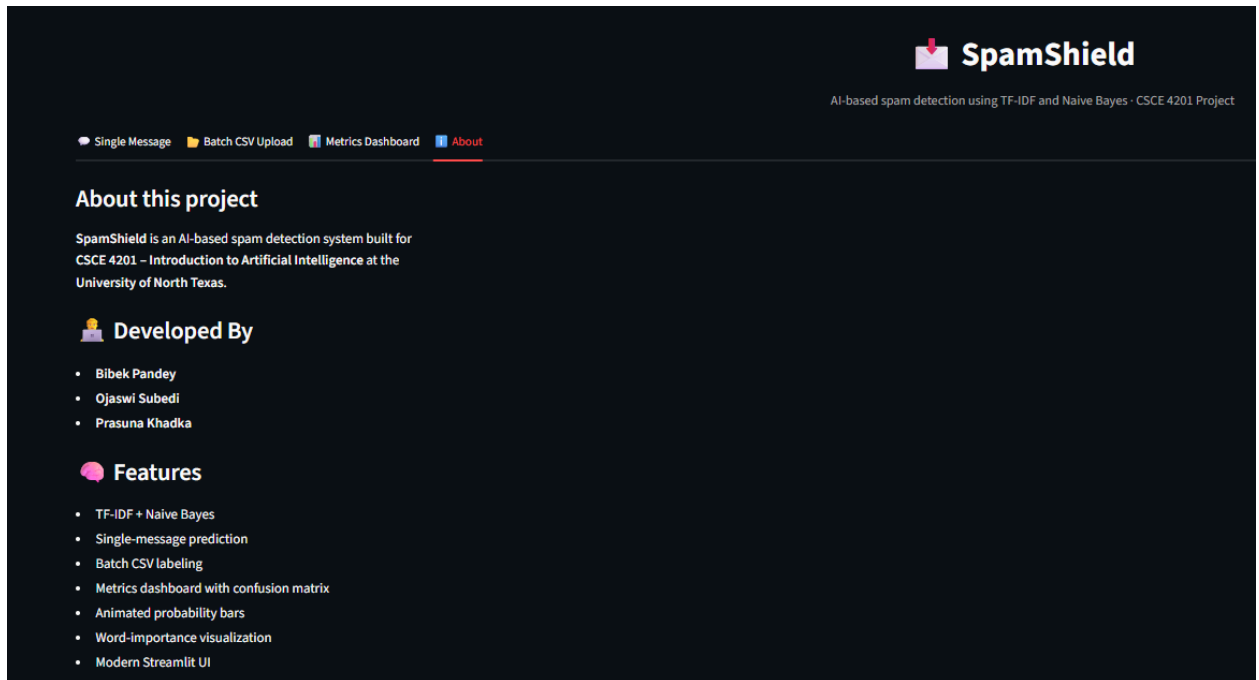
Summary

Accuracy: 0.970  
Precision (spam): 1.000  
Recall (spam): 0.779  
F1 Score (spam): 0.875

Confusion Matrix



- About Page



## 4.4 Discussion

- Ham messages are easy to classify; most contain neutral language.
- Spam messages contain keywords with strong TF-IDF weights (“winner”, “selected”, “free”).
- The model performs consistently but is dataset-dependent.
- Explainability provides transparency into predictions.
- Errors in earlier prototypes (UI formatting, batch errors, variable naming issues) were fully resolved.

---

## 5. Conclusion and Future Work

SpamShield is an excellent case of using classic AI algorithms in a practical situation, where SMS messages were classified as spam or ham. It also stresses the power of TF-IDF and

Multinomial Naive Bayes, which are two algorithms that nowadays still have their place in the market, thanks to their interpretability, speed, and effectiveness in dealing with the short-text type. The combination of these algorithms, a modern user interface, and an explainability feature gives rise to both functional and educational values.

The high-impact keyword strategies that make spam messages easy to identify have made Multinomial Naive Bayes the best choice here. TF-IDF helps to further increase that by assigning a lot of weights to the words that convey meaning and almost zero weights to those that are common or lack information. Despite the existence of more advanced models, the simplicity and transparency of these algorithms still make them perfect for this project and for classroom situations.

SpamShield has the potential to grow significantly with more data, fine-tuning, and smarter algorithms. The groundwork laid in this project will make it easier to introduce logistic regression, support vector machines, or even transformer-based architectures in the future. The choices regarding system design, such as modularity, explainability, and UI integration, will make the system very convenient for scaling, evaluating, and deploying.

SpamShield is a prime example of classical supervised learning being able to work on real-world problems related to the classification of texts, provided it is combined with a clear preprocessing, modular system design, and a user-friendly interface. The incorporation of the explainability feature increases the educational impact of the system, making it, thus, a perfect candidate for classroom demonstrations and first-stage prototyping.

The system's architectural design supports generalization even though only a small dataset was used. Performance would be markedly improved with a larger dataset, more varied samples, and parameter tuning. Moreover, the application of the latest NLP architectures such as word embeddings, logistic regression classifiers, or transformer-based encoders would not just improve accuracy and robustness but also make the models more versatile..

Future improvements may include:

- Larger and more diverse datasets
- Additional classification algorithms
- Threshold tuning tools
- More detailed evaluation dashboards
- Full deployment on the web using Streamlit Cloud or containerization

## **Conclusion**



SpamShield successfully demonstrates a complete AI application, including preprocessing, feature extraction, model training, inference, UI design, explainability, and evaluation. The system is fully functional, modular, and expandable.

## Future Work

Potential improvements include:

- Using Logistic Regression, SVM, or BERT for better accuracy
  - Expanding dataset with real SMS corpora
  - Adding ROC/PR curve analysis
  - Deploying the model online
  - Adding multilingual spam detection
  - Integrating real-time SMS ingestion APIs
- 

## 6. References

1. Almeida, T., Hidalgo, J., & Yamakami, A. (2011). *Contributions to the study of SMS spam filtering*. ACM Symposium on Document Engineering.
  2. scikit-learn documentation. <https://scikit-learn.org>
  3. Streamlit documentation. <https://streamlit.io>
  4. Python Software Foundation. Python Language Reference.
  5. UCI SMS Spam Collection Dataset.  
<https://archive.ics.uci.edu/ml/datasets/sms+spam+collection>
- 

## Submission Package Checklist

FinalReport-Team11/

- ├── FinalReport-Team11.pdf
- ├── README.md
- ├── ui\_streamlit.py
- ├── app.py
- ├── requirements.txt
- ├── src/
  - ├── train.py
  - ├── evaluate\_model.py
  - └── batch\_predict.py
- ├── data/
  - ├── sms\_spam.csv
  - └── example\_batch.csv
- └── models/
  - └── spam\_nb\_tfidf.pkl