

# **DESIGNTHINKING FOR AI SPAM DETECTOR**

## **PROBLEM DEFINITION:**

Spam detection is a supervised machine learning problem. This means you must provide your machine learning model with a set of examples of spam and ham messages and let it find the relevant patterns that separate the two different categories.

## **DEFINITION:**

In order to more effectively analyze the content and not trash a real message, sophisticated spam filters use artificial intelligence (AI) techniques that look for key words and attempt to decipher their meaning in sentences (see Bayesian filtering). See spam trap, spam relay and spamdexing.

## **BACKGROUND:**

Spam emails, messages, and content pose a significant nuisance and security threat to individuals and organizations. Traditional rule-based methods for spam detection often fall short in identifying sophisticated and evolving spam patterns. To address this challenge, the objective is to develop a robust and adaptive spam detector using Artificial Intelligence (AI) techniques.

## **DATACOLLECTION:**

### **Data Sources:**

Identify the sources from which you will collect spam

data. These sources may include email inboxes, SMS messages, website comments, or any other platform where spam is prevalent.

### **Data Variety:**

Spam can take various forms and may contain different types of content (e.g., text, images, links). Ensure your dataset reflects this diversity so that your AI model can learn to detect a wide range of spam

### **Negative Samples:**

In addition to spam samples, it's important to have a representative set of non-spam (ham) samples. This helps the model distinguish between spam and legitimate content effectively.

## **DATA PROCESSING:**

### **Data Collection:**

Gather a large dataset of emails, messages, or content that includes both spam and non-spam (ham) examples. These should be representative of the types of messages the detector will encounter in real-world applications.

### **Data Labeling:**

Annotate each message as either spam or ham. This can be done manually or using pre-existing labels if available.

Ensure that the labeling is accurate, as it forms the basis for training and evaluating your model.

### **Data Normalization/Scaling:**

Scale the numerical features to have a consistent range, often between 0 and 1 or using standardization.

This step is particularly important if you're using machine learning algorithms that are sensitive to feature scales.

## **FUTURE EXTRATION:**

### **Feature Extraction:**

Convert the text messages into numerical features that can be fed into the AI model. Common techniques include:

Bag of Words (BoW): Represents each message as a vector of word frequencies.

Term Frequency-Inverse Document Frequency (TF-IDF): Measures the importance of words in a document relative to a corpus.

Word Embeddings: Pre-trained word embeddings like Word2Vec or GloVe can be used to capture semantic relationships between words.

N-grams: Capture sequences of words, which can help detect patterns in text

## **MODEL SELECTION:**

### **Naive Bayes:**

Pros: Simple and computationally efficient, making it a good choice for small to medium-sized datasets. It often performs well when the assumption of feature independence holds reasonably.

Cons: May not capture complex relationships in data, as it assumes that features are independent.

### **Support Vector Machines (SVM):**

Pros: Effective for high-dimensional data and can capture complex decision boundaries. SVMs can perform well with proper tuning.

Cons: Training SVMs can be computationally expensive, and the choice of kernel function and hyperparameters can impact performance.

### **Random Forest:**

Pros: Ensemble models like Random Forest can handle complex relationships in data, are robust to overfitting, and can provide feature importance scores.

Cons: Requires more data than simpler models to perform well, and hyperparameter tuning is necessary.

### **EVALUATION:**

#### **F1-Score:**

The F1-Score is the harmonic mean of precision and recall and provides a balance between the two metrics. It is calculated as  $2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall})$ .

The F1-Score is especially useful when you want to balance precision and recall in your spam detector.

#### **Confusion Matrix:**

A confusion matrix is a fundamental tool for evaluating classification models, including spam detectors. It provides a breakdown of the model's predictions into four categories: true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN).

**From the confusion matrix, you can calculate other evaluation metrics.**