

## 1. Problem Definition

- **Goal:** Classify text messages (emails, SMS, etc.) as **spam** or **ham** (not spam).
- 

## 2. Data Collection

- Use datasets like the **SMS Spam Collection** dataset from UCI Machine Learning Repository or any relevant labeled dataset.
  - **Format:** Typically has two columns: **label** (spam or ham) and **message** (the actual text).
- 

## 3. Data Preprocessing

- Convert labels to binary (e.g., **spam** → 1, **ham** → 0).
- **Clean the text:**
  - Lowercasing
  - Remove punctuation, numbers, and special characters
  - Remove stopwords (like "the", "is", "in")
  - Tokenization (splitting text into words)
  - (Optional) Lemmatization or stemming
  - **Explanation**

### 1. Tokenization

**Meaning:** Breaking a sentence or paragraph into smaller parts like words or punctuation.

**Example:**

**Text:** "I love playing football."

**Tokens:** ['I', 'love', 'playing', 'football', '.']

**This helps machines understand and process each word separately.**

---

## 2. Stemming

**Meaning:** Cutting the word to its root form by removing prefixes or suffixes, even if the result is not a valid word.

**Example:**

- "playing", "played", "plays" → "play"
- "running" → "runn" (not a real word, but stemmed)

*It's fast but sometimes gives rough or inaccurate root words.*

---

## 3. Lemmatization

**Meaning:** Converting a word to its proper dictionary root form, considering grammar and context.

**Example:**

- "playing" → "play"
- "better" → "good"
- "studies" → "study"

*More accurate than stemming, but a bit slower.*

○

---

## 4. Feature Extraction (Text Vectorization)

- Convert text into numerical format using:
  - **Bag of Words (CountVectorizer)**

- **TF-IDF** (*TfidfVectorizer*)
- (Optional) Word embeddings like Word2Vec, BERT for advanced models

## 5 Split Your Data

Divide your dataset into **training** and **testing** sets:

---

## 6 Model Selection

- Choose a classification model:
  - **Naive Bayes (MultinomialNB)** — best for text classification
  - Logistic Regression
  - Support Vector Machine (SVM)
  - Random Forest / XGBoost
  - (Advanced) Deep Learning: LSTM, BERT

---

## 7. Model Evaluation

- Use metrics:
  - Accuracy
  - Precision, Recall, F1-Score (important for imbalanced datasets)

- Confusion Matrix
  - Validate with cross-validation if needed
- 

## 8. Prediction and Testing

- Use the trained model to classify new messages.
  - Evaluate its real-world performance.
- 

## 9. Save the Model

- Use `joblib` or `pickle` to save the trained model and vectorizer:

```
python
CopyEdit
import joblib
joblib.dump(model, 'spam_classifier.pkl')
joblib.dump(vectorizer, 'vectorizer.pkl')
```

---

## 10. Deployment(local in vs\_code)

- Build a web or mobile app using Flask, Django, or Streamlit to input new messages and show classification results.