

# SMS Spam Classification using Markov Chains and Bayes Theorem

## Task Overview

Implement a spam classifier that combines Markov chains with Bayes' theorem to identify spam SMS messages using spam.csv dataset.

## Dataset

- File name: spam.csv
- Contains two columns:
  - 'label': marks if message is 'spam' or 'ham'
  - 'message': contains the SMS text

## Bayes' Theorem Application

For a message M, we calculate  $P(\text{spam}|M)$  using Bayes' theorem:

$$P(\text{spam}|M) = P(M|\text{spam})P(\text{spam}) / P(M)$$

Breaking this down:

1. Prior Probability (in log space):

$$\log P(\text{spam}) = \log(\text{number of spam messages} / \text{total messages})$$

$$\log P(\text{ham}) = \log(\text{number of ham messages} / \text{total messages})$$

2. Likelihood  $P(M|\text{spam})$  for message  $M = [\text{word}_1, \text{word}_2, \dots, \text{word}_T]$ :

$$\log P(M|\text{spam}) = \log P(\text{word}_1|\text{spam}) + \# \text{ Initial probability}$$

$$\log P(\text{word}_2|\text{word}_1) + \# \text{ First transition}$$

$$\log P(\text{word}_3|\text{word}_2) + \# \text{ Second transition}$$

$$\dots \# \text{ Remaining transitions}$$

$$\log p(s_{1...T}) = \log \pi_{s_1} + \sum_{t=2}^T \log A_{s_{t-1}, s_t}$$

3. Final Classification:

- Calculate  $\log P(\text{spam}|M) = \log P(M|\text{spam}) + \log P(\text{spam})$
- Calculate  $\log P(\text{ham}|M) = \log P(M|\text{ham}) + \log P(\text{ham})$
- Compare: if  $\log P(\text{spam}|M) > \log P(\text{ham}|M)$ , classify as spam

## Required Formulas

1. Initial State Probabilities:

$$\hat{\pi}_i = \frac{\text{count}(s_1 = i) + \varepsilon}{N + \varepsilon M}$$

2. Transition Probabilities:

$$\hat{A}_{ij} = \frac{\text{count}(i \rightarrow j) + \varepsilon}{\text{count}(i) + \varepsilon M}$$

## Implementation Steps

1. Data Processing:

- Load spam.csv
- Split into train (80%) and test (20%)
- Preprocess messages (lowercase, basic cleaning)

2. Training:

- Calculate log prior probabilities  $P(\text{spam})$  and  $P(\text{ham})$
- Build vocabulary for spam and ham
- Compute initial state probabilities
- Compute transition probabilities

3. Prediction: For each message:

- Calculate  $\log P(M|\text{spam})$ :
  - Add log of initial probability for first word
  - Add log of transition probabilities for each word pair
  - Add log prior  $P(\text{spam})$
- Calculate  $\log P(M|\text{ham})$  similarly
- Compare and classify

4. Evaluation:

- Calculate accuracy on test set