# Naïve Bayes Classifier

Course: Data Mining

Professor: Dr. Tahaei

Author: Sina Asghari

**Subject: Classification Problem and Solution**

December 2024

# Question

A company wants to classify incoming emails as **Spam** or **Not Spam** based on the occurrence of the words *offer* and *discount.* You are given the following information:

## Data:

- **Training Emails:**

  - 6 emails classified as **Spam**.
  - 4 emails classified as **Not Spam**.

- **Occurrences of words:**

  - In **Spam emails**:

    * *offer* appears in 4 emails.
    * *discount* appears in 3 emails.

  - In **Not Spam emails**:

    * *offer* appears in 1 email.
    * *discount* appears in 2 emails.

**Task:** Using Naïve Bayes, calculate whether an email containing the words *offer* and *discount* is classified as **Spam** or **Not Spam**.

# Solution

## Step 1: Calculate Prior Probabilities

$$P(\text{Spam}) = \frac{\text{Number of Spam emails}}{\text{Total emails}} = \frac{6}{10} = 0.6$$

$$P(\text{Not Spam}) = \frac{\text{Number of Not Spam emails}}{\text{Total emails}} = \frac{4}{10} = 0.4$$

## Step 2: Calculate Likelihoods

**For Spam:**

$$P(\text{"offer"}|\text{Spam}) = \frac{\text{Number of Spam emails with "offer"}}{\text{Total Spam emails}} = \frac{4}{6} = 0.67$$

$$P(\text{"discount"}|\text{Spam}) = \frac{\text{Number of Spam emails with "discount"}}{\text{Total Spam emails}} = \frac{3}{6} = 0.5$$

**For Not Spam:**

$$P(\text{"offer"}|\text{Not Spam}) = \frac{\text{Number of Not Spam emails with "offer"}}{\text{Total Not Spam emails}} = \frac{1}{4} = 0.25$$

$$P(\text{"discount"}|\text{Not Spam}) = \frac{\text{Number of Not Spam emails with "discount"}}{\text{Total Not Spam emails}} = \frac{2}{4} = 0.5$$

## Step 3: Combine Probabilities Using Independence Assumption

**For Spam:**

$$P(\text{"offer and discount"}|\text{Spam}) = P(\text{"offer"}|\text{Spam}) \times P(\text{"discount"}|\text{Spam}) = 0.67 \times 0.5 = 0.335$$

**For Not Spam:**

$$P(\text{"offer and discount"}|\text{Not Spam}) = P(\text{"offer"}|\text{Not Spam}) \times P(\text{"discount"}|\text{Not Spam}) = 0.25 \times 0.5 = 0.1$$

## Step 4: Calculate Posterior Probabilities

**For Spam:**

$$P(\text{Spam}|\text{"offer and discount"}) \propto P(\text{Spam}) \times P(\text{"offer and discount"}|\text{Spam}) = 0.6 \times 0.335 = 0.201$$

**For Not Spam:**

$$P(\text{Not Spam}|\text{"offer and discount"}) \propto P(\text{Not Spam}) \times P(\text{"offer and discount"}|\text{Not Spam}) = 0.4 \times 0.125 = 0.05$$

## Step 5: Compare Probabilities

$$P(\text{Spam}|\text{"offer and discount"}) > P(\text{Not Spam}|\text{"offer and discount"})$$

**Conclusion:** The email is classified as **Spam**.