

1. Calculate the prior probabilities, spam and not spam
2. Use the Naive Bayes classifier to predict whether a new email from gmail.com that contains "Free" but no "Discount" is spam or not

Drill Task:

1. You are given a dataset of emails labelled as spam or not spam. Each email is represented by a few key features. Your task is to manually build a Naive Bayes classifier to predict whether a new email is spam or not based on these features.

Email-ID	Sender Domain	Contains Offer	Contains Free	Contains "Discount"	Is Spam
1	Gmail.com	Yes	No	No	No
2	Yahoo.com	Yes	Yes	No	Yes
3	Gmail.com	No	No	Yes	Yes
4	Outlook.com	Yes	No	Yes	Yes
5	Yahoo.com	No	Yes	No	Yes
6	Gmail.com	Yes	Yes	Yes	Yes
7	Gmail.com	No	No	No	No
8	Outlook.com	Yes	No	No	No
9	Gmail.com	No	Yes	No	Yes
10	Yahoo.com	Yes	Yes	Yes	Yes

Ans:- Prior probability of spam ($p(\text{spam})$)

$$P(\text{spam}) = \frac{\text{Num of spam emails}}{\text{Total num of emails}}$$

$$= \frac{7}{10} = 0.7$$

$$P(\text{not spam}) = \frac{\text{Num of not spam emails}}{\text{Total num of emails}}$$

$$= \frac{3}{10} = \underline{\underline{0.3}}$$

2. Use the naive bayes classifier to predict whether a new email from gmail.com that contains free but no discount is spam or not.

but no discount is spam?

Mail: Contain Free, Discount = spam or not spam
Yes No

Step 1:- Sender Domain gmail.com.

Step 1:- Sender

$$P(\text{gmail} | \text{spam}) = \frac{\text{Num of spam emails from gmail.com}}{\text{Total number of spam emails.}}$$

$$= \frac{3}{7} = \underline{\underline{0.42}}$$

$$P(\text{gmail} | \text{Not spam}) = \frac{\text{Num of } \overset{\text{not spam}}{\text{gmail}} \text{ emails from gmail}}{\text{Total num of not spam email}}$$

$$= \frac{2}{3} = \underline{\underline{0.66}}$$

Step 2 :- find Contains Free yes | ~~gmail~~ spam & not spam.

$$\rightarrow P(\text{Contains Free} = \text{yes} | \text{spam}) = \frac{\text{Num of spam emails containing Free}}{\text{Total num of spam email}}$$

$$= \frac{5}{7} = \underline{\underline{0.71}}$$

$$\rightarrow P(\text{Contains Free} = \text{yes} | \text{not spam}) = \frac{\text{Num of } \overset{\text{not}}{\text{spam}} \text{ emails containing Free}}{\text{Total num of not spam email}}$$

$$= \frac{0}{3} = \underline{\underline{0}}$$

Step 3: Find contain discount Yes No | spam & not

$$\rightarrow P(\text{contain Discount} = \text{No} | \text{spam})$$

$$= \frac{\text{Num of spam emails Contain Discount No}}{\text{Total num of spam email}}$$

$$= \frac{3}{7} = \underline{0.42}$$

$$\rightarrow P(\text{contain Discount} = \text{No} | \text{not spam})$$

$$= \frac{\text{Num of not spam emails Contain Discount No}}{\text{Total num of not spam email}}$$

$$= \frac{3}{3} = \underline{1}$$

II Calculate posterior probabilities using Naive Bayes

$P(\text{features} | \text{spam})$

\rightarrow Probability of every feature with all spam

$$P(x | \text{spam}) = 0.42 \times 0.71 \times 0.42$$

$$= \underline{0.125}$$

→ Probability of every features with not spam

$$P(x|\text{not spam}) = 0.66 \times 0 \times 1 \\ = \underline{\underline{0}}$$

III §

→ probability of this spam into corresponding spam

$$P(x|\text{spam}) \times P(\text{spam}) \\ = 0.125 \times 0.7 \\ = \underline{\underline{0.08}}$$

→ probability of this not spam into corresponding not spam

$$P(x|\text{not spam}) \times P(\text{not spam}) \\ = 0 \times 0.3 \\ = \underline{\underline{0}}$$

$$\therefore P(x|\text{spam}) > P(x|\text{Not spam})$$

The new email from gmail-com that contains "free" but no "Discount" is predicted to be spam.