

2. You are classifying emails as legitimate or spam and in doing so you would like to estimate the probability that a new email e containing the keywords (w_1, w_2, \dots, w_n) is spam by taking all the emails in the training set with those keywords. In other words the probability is calculated as:

$$P(\text{spam}|e) = \frac{\text{Number of spam emails with keywords } (w_1, w_2, \dots, w_n)}{\text{Number of total emails with keywords } (w_1, w_2, \dots, w_n)} \quad (1)$$

- Explain why the plan may not work.
- Describe the data sets for which this plan might work.
- Using Naive Bayes assumption, explain how to get the probability of an email being spam? Show it mathematically.
- How does Naive Bayes assumption caters the problem with your plan?

[a] The keywords are assumed independent of each other, which is unlikely as it might be possible that more than one keywords are dependent on each other.

- Also some keywords might be more important than another.
- Curse of dimensionality is also a issue. As for decently sized dataset, no. of spam emails with keywords w_1, w_2, \dots, w_n might be '0' or very small, as it is unlikely for an email to contain 'n' keywords for large n.
- This means we might estimate probability '0' for most new emails or undefined if denominator is '0'.

[b] The idea will work for large datasets. We want emails with each possible set of keywords.

E.g. for 'n' key words, we need 2^n emails.

S = set of 'n' keywords.

$PS = 2^n \rightarrow$ powerset of S .

[c] Let $y \rightarrow 1$ be spam
 $y \rightarrow 0$ be not-spam

Bayes Theorem:-

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad - (a)$$

$$P(y=1) = 1 - P(y=0)$$

$$\therefore P(\text{spam} | w_1, w_2, \dots, w_n) =$$

$$\frac{P(w_1, w_2, \dots, w_n | \text{spam}) \cdot P(\text{spam})}{P(w_1, w_2, \dots, w_n)} \quad - (b)$$

$$P(w_1, w_2, \dots, w_n)$$

Now, for independent events A & B ,

$$P(A, B) = P(A) \cdot P(B) \quad - (c)$$

By a, b, c

$$\therefore P(\text{spam} | w_1, \dots, w_n) =$$

$$\frac{P(\text{spam}) \cdot P(w_1, w_2, \dots, w_n | \text{spam})}{P(w_1) \cdot P(w_2) \dots P(w_n)}$$

$$= P(\text{spam}) \prod_{i=1}^n P(w_i | \text{spam})$$

$$\frac{P(w_1 | \text{spam}) \dots P(w_n | \text{spam}) + P(w_1 | \text{not spam}) \dots P(w_n | \text{not spam})}{P(w_1 | \text{spam}) \dots P(w_n | \text{spam}) + P(w_1 | \text{not spam}) \dots P(w_n | \text{not spam})}$$

$$= P(\text{spam}) \prod_{i=1}^n P(w_i | \text{spam})$$

$$\frac{P(\text{spam}) \prod_{i=1}^n P(w_i | \text{spam}) + P(\text{not spam}) \prod_{i=1}^n P(w_i | \text{not spam})}{P(\text{spam}) \prod_{i=1}^n P(w_i | \text{spam}) + P(\text{not spam}) \prod_{i=1}^n P(w_i | \text{not spam})}$$

[d] In Naive Bayes solution, we are calculating $P(w_i | \text{spam})$ for each word w_i .

Thus instead of set of keywords we are just looking for a single w_i , it decreases the chance of getting probability zero.

Also, the probability of each word is calculated instead of joint probability.