

Part1

The useful point of Laplace smoothing is to change the zero-values of your data into small positive number to prevent fail your entire process. When you change the zero value, you should pay attention to reduce other values at the same time since the total sum of probability maintains 1.

Considering, $y_{pred} = \operatorname{argmax}_y P(y)P(x_1 | y)P(x_2 | y) \cdots P(x_n | y)$

for instance,

Step1

First, you had 30 spam emails and 70 non-spam emails, then you had

$$P(\text{spam}) = 0.3, P(\text{non-spam}) = 0.7 \cdots (*1)$$

Now, maintaining the total sum of probability is 1, changes to 31 spam emails and 69 non-spam emails.

$$P(\text{spam}) = 0.31, P(\text{non-spam}) = 0.69 \cdots (*2)$$

(*1) and (*2) belong to $P(y)$ and the latter value should be decreased considering the augmentation in the process of (*3) \rightarrow (*4)

Step2

Think about $P(x_1 | y)$

Assume you originally had 'buy' in 10 out of 40 spam emails and in no non-spam emails,

$$P(\text{buy} | \text{spam}) = 0.25, P(\text{buy} | \text{non-spam}) = 0 \cdots (*3)$$

" $P(\text{buy} | \text{non-spam}) = 0$ " leads the entire process to fail by being multiplied since the value is zero. That is,

$$y_{pred} = \operatorname{argmax}_y P(y)P(\text{buy} | \text{non-spam})P(x_2 | y) \cdots P(x_n | y) = 0$$

So, this phenomenon yields this model bad. To prevent this disadvantageous process, you should use Laplace smoothing.

Based on this concept, I try to change the zero value into the small positive number. See (*3 \rightarrow *4) To change, suppose you had 'buy' in 12 out of 41 spam emails and in 2 out of 200 non-spam emails.

$$P(\text{buy} | \text{spam}) = 0.29, P(\text{buy} | \text{non-spam}) = 0.01 \cdots (*4)$$

$$y_{pred} = \operatorname{argmax}_y P(y)P(\text{buy} | \text{non-spam})P(x_2 | y) \cdots P(x_n | y) \neq 0$$

You can hamper the entire-process from failing by the Laplace smoothing. In specifically,

Suppose $P_{\text{empirical}} = \frac{x_i}{N} = 0$, then you should employ the Laplace smoothing like this

$$P_{\alpha\text{-smoothed}} = \frac{x_i + \alpha}{N + \alpha d} \text{ where } d \text{ is constant.}$$

You could see that α varies and there is tradeoff between variance and bias:

If α is small, it leads to high variance.

If α is large, it leads to high bias.