# SNLP Assignment 3

**Clemens Damke - 7011488**

## 1. Naïve Bayes classification

Using $\alpha = 1$ and $c_s := \text{spam}, c_n := \text{normal}$.

$$P(c_s) = \frac{\#docs_{c_s}}{\#docs} = \frac{2}{3}, \; P(c_n) = \frac{\#docs_{c_n}}{\#docs} = \frac{1}{3}$$

### 1.1. Equal weights

$$\hat{P}(\text{free}|c_s) = \frac{2+\alpha}{11+\alpha|V|} = \frac{3}{19} \qquad \hat{P}(\text{free}|c_n) = \frac{0+\alpha}{4+\alpha|V|} = \frac{1}{12}$$

$$\hat{P}(\text{bitcoins}|c_s) = \frac{4}{19} \qquad \hat{P}(\text{bitcoins}|c_n) = \frac{1}{12}$$

$$\hat{P}(\text{bank}|c_s) = \frac{2}{19} \qquad \hat{P}(\text{bank}|c_n) = \frac{1}{12}$$

$$\hat{P}(\text{account}|c_s) = \frac{2}{19} \qquad \hat{P}(\text{account}|c_n) = \frac{1}{12}$$

$$\hat{P}(\text{credit}|c_s) = \frac{3}{19} \qquad \hat{P}(\text{credit}|c_n) = \frac{2}{12}$$

$$\hat{P}(\text{card}|c_s) = \frac{2}{19} \qquad \hat{P}(\text{card}|c_n) = \frac{2}{12}$$

$$\hat{P}(\text{wallet}|c_s) = \frac{2}{19} \qquad \hat{P}(\text{wallet}|c_n) = \frac{2}{12}$$

$$\hat{P}(\text{wood}|c_s) = \frac{1}{19} \qquad \hat{P}(\text{wood}|c_n) = \frac{2}{12}$$

$$\hat{P}(c_s|d_4) = \frac{2 \cdot 4 \cdot 2 \cdot 3 \cdot 1}{3 \cdot 19^4} \approx 0.000123 \qquad \hat{P}(c_n|d_4) = \frac{1 \cdot 1 \cdot 2 \cdot 1 \cdot 2}{3 \cdot 12^4} \approx 0.000064$$

$$\hat{P}(c_s|d_5) = \frac{2 \cdot 2 \cdot 3 \cdot 3 \cdot 2}{3 \cdot 19^4} \approx 0.000184 \qquad \hat{P}(c_n|d_5) = \frac{1 \cdot 2 \cdot 1 \cdot 2 \cdot 2}{3 \cdot 12^4} \approx 0.000129$$

$$\implies c_{d_4} = c_s = \text{spam}, \; c_{d_5} = c_s = \text{spam}$$

### 1.2. Different weights

It was unclear to me what exactly was meant by tripling the weight. I chose to triple the title weight by counting each title word thrice. A word appearing in the title of a training sample thus has the same higher weight in the title and the text of test samples:

$$\hat{P}(\text{free}|c_s) = \frac{\overbrace{3 \cdot 1}^{\text{title}} + \overbrace{1}^{\text{text}} + \alpha}{3 \cdot 4 + 7 + \alpha|V|} = \frac{5}{27} \qquad \hat{P}(\text{free}|c_n) = \frac{\overbrace{3 \cdot 0}^{\text{title}} + \overbrace{0}^{\text{text}} + \alpha}{3 \cdot 2 + 2 + \alpha|V|} = \frac{1}{16}$$

$$\hat{P}(\text{bitcoins}|c_s) = \frac{6}{27} \qquad \hat{P}(\text{bitcoins}|c_n) = \frac{1}{16}$$

$$\hat{P}(\text{bank}|c_s) = \frac{2}{27} \qquad \hat{P}(\text{bank}|c_n) = \frac{1}{16}$$

$$\hat{P}(\text{account}|c_s) = \frac{2}{27} \qquad \hat{P}(\text{account}|c_n) = \frac{1}{16}$$

$$\hat{P}(\text{credit}|c_s) = \frac{5}{27} \qquad \hat{P}(\text{credit}|c_n) = \frac{2}{16}$$

$$\hat{P}(\text{card}|c_s) = \frac{4}{27} \qquad \hat{P}(\text{card}|c_n) = \frac{2}{16}$$

$$\hat{P}(\text{wallet}|c_s) = \frac{2}{27} \qquad \hat{P}(\text{wallet}|c_n) = \frac{4}{16}$$

$$\hat{P}(\text{wood}|c_s) = \frac{1}{27} \qquad \hat{P}(\text{wood}|c_n) = \frac{4}{16}$$

$$\hat{P}(c_s|d_4) = \frac{2 \cdot 6 \cdot 2 \cdot 5 \cdot 1}{3 \cdot 27^4} \approx 0.000075 \qquad \hat{P}(c_n|d_4) = \frac{1 \cdot 1 \cdot 4 \cdot 1 \cdot 4}{3 \cdot 16^4} \approx 0.000081$$

$$\hat{P}(c_s|d_5) = \frac{2 \cdot 2 \cdot 5 \cdot 5 \cdot 4}{3 \cdot 27^4} \approx 0.000251 \qquad \hat{P}(c_n|d_5) = \frac{1 \cdot 4 \cdot 1 \cdot 2 \cdot 2}{3 \cdot 16^4} \approx 0.000081$$

$$\implies c_{d_4} = c_n = \text{normal}, \; c_{d_5} = c_s = \text{spam}$$