

problem 2.

1) 3-grams: $p_r(W | W_{i-1}, W_{i-2}) = \frac{\text{number of sequence } (W_{i-2}, W_{i-1}, W)}{\text{number of sequence } (W_{i-2}, W_{i-1})}$ (MLE).

clearly, it follows multinomial Distribution.

2) For a given data set, we first need to ~~split~~ divided it into training and testing data randomly.

2) For the training data set: calculate the ~~total~~ proportion of email and proportion of spam: $p(\text{email})$ and $p(\text{spam})$

$$\text{And } p(\text{spam}) + p(\text{email}) = 1$$

3) For the new data set we make the decision based on Lazy learning Approach and idea of bayes theorem: $\hat{y} = \underset{y}{\text{argmax}} p(y | \text{new data}) = \underset{y}{\text{argmax}} p(y) p(\text{new data} | y)$.

$$\begin{aligned} p(\text{spam} | \text{new data}) &\propto p(\text{spam}) \times p(\text{new data} | \text{spam}) \\ &= p(\text{spam}) \times \prod_{i=3}^m (\text{3-gram probability of } i) \\ &= p(\text{spam}) \times \frac{\prod_{i=3}^m \# \text{ of sequence } (W_{i-2}, W_{i-1}, W)}{\prod_{i=3}^m \# \text{ of sequence } (W_{i-2}, W_{i-1})} \end{aligned}$$

~~W~~ W is the words in the new data set.

in the same way: $p(\text{email} | \text{new data}) \propto p(\text{email}) \times p(\text{new data} | \text{email})$

$$= p(\text{email}) \times \frac{\prod_{i=3}^m \# \text{ of sequence } (W_{i-2}, W_{i-1}, W)}{\prod_{i=3}^m \# \text{ of sequence } (W_{i-2}, W_{i-1})}$$

Finally, $y = \begin{cases} \text{spam} & \text{if } p(\text{spam} | \text{new data}) \geq p(\text{email} | \text{new data}) \\ \text{email} & \text{otherwise.} \end{cases}$

3) As mentioned above, for a new message: $\hat{y} = \underset{y}{\text{argmax}} p(\hat{y} | \text{new data}) = \underset{y}{\text{argmax}} p(y) \cdot p(\text{new data} | y)$

$$y = \begin{cases} \text{spam} & \text{if } p(\text{spam} | \text{new data}) \geq p(\text{email} | \text{new data}) \\ \text{email} & \text{otherwise} \end{cases}$$

By the way, we can use this method for test data and calculate the Accuracy of the model.