

Assignment 5.1

Anant Chaturvedi(IIT2016506)

Question Description

1. Using Naive Bayesian classifier predict where a given mail is spam or not. Use the data set provided for this purpose. (structured data set)

Introduction

We are given a dataset comprising of email body text followed by the label depicting whether the email was a spam mail or ham mail.

Concepts Used

Bayesian Formula

$$P(A | B) = P(B | A) * P(A) / P(B)$$

Formula Definition

$P(\text{ham} | \text{bodyText})$ = Probability that email is ham given that it contains document- bodyText (lets say bodyText = content of email)

$P(\text{spam} | \text{bodyText})$ = Probability that email is spam given that it contains document- bodyText

$$P(\text{ham} | \text{bodyText}) = (P(\text{ham}) * P(\text{bodyText} | \text{ham})) / P(\text{bodyText})$$

$$P(\text{spam} | \text{bodyText}) = (P(\text{spam}) * P(\text{bodyText} | \text{spam})) / P(\text{bodyText})$$

$$P(\text{ham} | \text{bodyText}) = (P(\text{ham}) * P(\text{bodyText} | \text{ham})) / P(\text{bodyText})$$

$$P(\text{spam} | \text{bodyText}) = (P(\text{spam}) * P(\text{bodyText} | \text{spam})) / P(\text{bodyText})$$

$P(\text{ham})$ = no of documents belonging to category ham / total no of documents

$P(\text{spam})$ = no of documents belonging to category spam / total no of documents

$P(\text{bodyText} \mid \text{spam}) = P(\text{word1} \mid \text{spam}) * P(\text{word2} \mid \text{spam}) * \dots$

$P(\text{bodyText} \mid \text{ham}) = P(\text{word1} \mid \text{ham}) * P(\text{word2} \mid \text{ham}) * \dots$

$P(\text{word1} \mid \text{spam}) = \text{count of word1 belonging to category spam} / \text{total count of words belonging to category spam.}$

$P(\text{word1} \mid \text{ham}) = \text{count of word1 belonging to category ham} / \text{total count of words belonging to category ham.}$

RESULTS

75% of the data was used for training while the remaining was used for testing the predictor.

The Accuracy on the Training data is around :- 98.82051282051282%

The Accuracy on the Testing data is around :- 97.188995215311%