

CS-534 Machine Learning

Implementation Assignment 2

Mohammand Velayati (33.3%) Lucas Wells (33.3%) Chirag Shah (33.3%)

Introduction

In this assignment we are given a set of tweets from Hillary Clinton's and Donald Trump's Twitter accounts. We will train a Naïve Bayes classifier to predict the account that a new tweet comes from. Throughout this document we use the variable C to represent the account that a given tweet comes from, Hillary's or Trump's, and \mathbf{x} represents a tweet, such that each x_i is a single word in the tweet \mathbf{x} of class C . Also, let V be the vocabulary of unique words found in all tweets of the training set.

The Naïve Bayes classifier

The Naïve Bayes classifier is based on Bayes' Theorem,

$$p(C|\mathbf{x}) = \frac{p(C)p(\mathbf{x}|C)}{p(\mathbf{x})} \propto p(C)p(\mathbf{x}|C), \quad (1)$$

that states that the probability of observing a class C , given a set of features \mathbf{x} , is equal to the prior probability of observing class C times the likelihood of the feature begin in that class divided by the evidence $p(\mathbf{x})$. As the denominator does not depend on the class we can state that the probability of observing a class is proportional to the numerator. This is equivalent to the joint probability $p(C, x_1, \dots, x_n)$. In practice this computation on high dimensional features is prohibitive since we must estimate $k \times (2^d - 1)$ parameters. Naïve Bayes relaxes the conditional probability model and assumes that each feature is conditionally independent from all other features. Thus, we can write the conditional distribution over C as

$$p(C|x_1, \dots, x_n) \propto p(C) \prod_{i=1}^n p(x_i|C). \quad (2)$$

Combining the Naïve Bayes probability model with the Maximum A Posteriori (MAP) decision rule we get

$$\hat{y} = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i|C_k). \quad (3)$$

We will compare the performance of two Naïve Bayes event models, Bernoulli and Multinomial, on our Twitter classification problem.

Bernoulli event model

In the Bernoulli event model we represent tweets as binary vectors, so each word $x_i \in \{0, 1\}$ and $|\mathbf{x}| = |V|$. Thus, $x_i = 1$ means that the word V_i is contained in the tweet at least once, and $x_i = 0$ means that word V_i is not in the tweet \mathbf{x} . In the Bernoulli model we use

$$p(\mathbf{x}|C) = \prod_{i=1}^{|V|} p_{x_i|C}^{x_i} (1 - p_{x_i|C})^{(1-x_i)} \quad (4)$$

to compute the likelihood that a given tweet \mathbf{x} is from account C . And the likelihood that word x_i is found in class C_k is estimated by

$$p_{x_i|C_k} = \frac{n_k(x_i)}{N_k}, \quad (5)$$

where $n_k(x_i)$ denotes the number of occurrences of word x_i in class k and N_k is the number of tweets from class k . Since the product of many probabilities could result in underflow, we compute the log of the likelihoods by

$$\log p(\mathbf{x}|C) = \log \left(\prod_{i=1}^{|V|} p_{x_i|C}^{x_i} (1 - p_{x_i|C})^{(1-x_i)} \right) \quad (6)$$

$$= \sum_{i=1}^{|V|} [x_i \cdot \log p_{x_i|C} + (1 - x_i) \log(1 - p_{x_i|C})]. \quad (7)$$

Combining equations (1) and (7) we perform MAP estimation in log space with

$$\hat{y} = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \log p(C_k) \sum_{i=1}^{|V|} [x_i \cdot \log p_{x_i|C_k} + (1 - x_i) \log(1 - p_{x_i|C_k})]. \quad (8)$$

The priors for each class are calculated using Maximum Likelihood Estimation

$$p(C_k) = \frac{N_k}{N}, \quad (9)$$

where N is the total number of tweets in both classes.

Multinomial event model

In the Multinomial case we represent tweets as integer vectors, so each word $x_i \in \mathbb{Z}$, where each element x_i denotes the frequency of that word in the tweet. The likelihood is given by

$$p(\mathbf{x}|C) = \prod_{i=1}^{|V|} p_{x_i|C}^{x_i}, \quad (10)$$

and the probability $p_{x_i|C_k}$ is the relative frequency of word x_i in tweets of class k divided by the total number of words in the tweets of that class. Mathematically, this can be expressed as

$$p_{x_i|C_k} = \frac{\sum_{j=1}^N x_{ji} b_{jk}}{\sum_{m=1}^{|V|} \sum_{j=1}^N x_{jm} b_{jk}}, \quad (11)$$

where b_{jk} is an indicator variable that equal 1 when x_i is of class $C = k$, and 0 otherwise. The estimate for all word likelihood in log space is then given by

$$\log p(\mathbf{x}|C) = \log \left(\prod_{i=1}^{|V|} p_{x_i|C}^{x_i} \right) \quad (12)$$

$$= \sum_{i=1}^{|V|} [x_i \cdot \log p_{x_i|C}]. \quad (13)$$

Finally, combining equations (1) and (13) we perform MAP estimation with

$$\hat{y} = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \log p(C_k) \sum_{i=1}^{|V|} [x_i \cdot \log p_{x_i|C_k}], \quad (14)$$

where the priors are calculated using equation (9).

Implementation

1. (5 pts) Please explain how you use the log of probability to perform classification
2. (5 pts) Report the overall testing accuracy (number of correctly classified documents over the total number of documents) for both (Bernoulli and Multinomial) models.

Table 1: Overall testing accuracy for each model

Model	testing accuracy
Bernoulli	91.61%
Multinomial	92.08%

Table 2: Confusion matrices. H denotes Hillary Clinton’s account and T for Trump’s account.

(a) Bernoulli model			(b) Multinomial model		
class	$y = H$	$y = T$	class	$y = H$	$y = T$
$\hat{y} = H$	46.58	5.28	$\hat{y} = H$	47.36	5.59
$\hat{y} = T$	3.11	45.03	$\hat{y} = T$	2.33	44.72

- (5 pts) Whose tweets were confused more often than the other? Why do you think this is?
- (5 pts) Identify, for each class, the top ten words that have the highest probability.

Table 3: Top ten most likely words from each account

Rank	H	T
1	is	on
2	Hillary	you
3	in	of
4	for	is
5	Trump	a
6	of	in
7	and	I
8	a	and
9	the	to
10	to	the

Priors and overfitting (20 pts)

Identifying important features (20 pts)

We explored methods for reducing the vocabulary size to reduce the dimensionality of features and improve performance of classification. Intuitively we can think of the discriminative power of a word as the difference in the likelihood of that word appearing in one class versus the other. That is, for a given word x_i , the absolute difference of the likelihood $p_{x_i|C_k}$ where $k = 0$ and $k = 1$. Thus, as an example, if a word is very likely to appear in both classes, then the difference is small and it is said to have little discriminative power. This can be written as

$$\Delta\ell(\mathbf{x}_{\ell|k}) = |\mathbf{x}_{\ell|k=0} - \mathbf{x}_{\ell|k=1}|, \quad (15)$$

where $\mathbf{x}_{\ell|k} = \{p_{x_i|C_k}; \forall i, k\}$ given by equation (5) for the Bernoulli case and (11) in the Multinomial case. A new vocabulary set is generated by

$$V' = \{x_{m,\sigma} | x_{m,\sigma} = \operatorname{argmax}_{\mathbf{x}_{\ell|k}}^{\sigma} \{\Delta\ell(\mathbf{x}_{\ell|k})\}, \sigma \in \{1, \dots, \Sigma\}\}, \quad (16)$$

where the notation $\operatorname{argmax}^{\sigma}$ denotes a set of Σ maximum word likelihood differences. The number of words in the new vocabulary set is controlled by parameter λ in

$$\Sigma = \lfloor |V| \cdot (1 - \lambda) \rfloor. \quad (17)$$

We tested the effect of λ values in the range $(0.1, \dots, 0.9)$ on the classification accuracy of each model and the size of the vocabulary (Table 4).

Table 4: Change in model accuracy resulting from vocabulary reduction

λ	Bernoulli accuracy (%)	Multinomial accuracy (%)	$ V $
0.1	92.24	92.24	12,220
0.2	92.55	92.24	10,862
0.3	92.24	91.93	9,504
0.4	91.77	91.62	8,146
0.5	92.23	92.24	6,789
0.6	92.71	92.24	5,431
0.7	92.24	91.93	4,073
0.8	91.31	91.31	2,715
0.9	90.53	90.21	1,357

In both models, the accuracy is slightly improved by retaining only words with relatively large discriminative power in the dictionary. For the Bernoulli case, this improvement persists up to a 70% reduction in vocabulary size. And for the Multinomial case, up to a 60% reduction. In both cases, a comparable classification accuracy can be achieved while reducing the vocabulary size to 90% of its original dimension.

Bonus