

CS5560 Knowledge Discovery and Management

Problem Set 6
July 10 (T), 2017

Name: Yalamanchili Sowmya
Class ID: 30

References

<https://www.analyticsvidhya.com/blog/2015/09/naive-bayes-explained/>
<https://nlp.stanford.edu/IR-book/html/htmledition/text-classification-and-naive-bayes-1.html>
<http://www.nltk.org/book/ch06.html>

- I. Consider the problem of classifying the origination point of passenger travel itineraries. Suppose we have the following training set of travel itineraries:

Itinerary	Document	Class
1	"smith: new york - chicago - san francisco - new york"	JFK
2	"chen: san francisco - london - paris - san francisco"	SFO
3	"chen: san francisco - tokyo - singapore- san francisco"	SFO
4	"o'brien: chicago - buenos aires - new york - chicago"	ORD

- a) Assume that we use a Bernoulli (i.e., binary) Naive Bayes model. Compute the following feature probabilities:
- $P(X_{\text{francisco}}=\text{true} \mid \text{Class}=\text{SFO})$
 - $P(X_{\text{london}}=\text{true} \mid \text{Class}=\text{SFO})$
 - $P(X_{\text{francisco}}=\text{true} \mid \text{Class}=\text{JFK})$
- b) Assume that we use a multinomial NB model instead. Compute the following probabilities:
- $P(X=\text{francisco} \mid \text{Class}=\text{SFO})$
 - $P(X=\text{london} \mid \text{Class}=\text{SFO})$
 - $P(X=\text{francisco} \mid \text{Class}=\text{JFK})$
- c) Consider a standard Naive Bayes classifier trained on the training set and applied to a similar test set. How accurate is this classifier for:
- the Bernoulli model, and
 - the multinomial model?
- d) Construct a non-standard feature representation that is 100% accurate for either model.

- II. This problem concerns smoothing Naïve Bayes classifiers. Consider the following formula for Laplace (add-1) smoothing for Naïve Bayes

$$\hat{P}(w_i | c) = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} (\text{count}(w, c) + 1)}$$

$$= \frac{\text{count}(w_i, c) + 1}{\left(\sum_{w \in V} \text{count}(w, c) \right) + |V|}$$

- a) Suppose we build a Naive Bayes classifier (multinomial or Bernoulli) with no smoothing of the respective $P(\text{word} | \text{class})$ probabilities. If a word was unseen in a class, it will thus have a probability of 0. Describe in words the decision procedure of this classifier (emphasizing the effect of the lack of smoothing, and how its decisions will differ from a smoothed Naive Bayes classifier).
- b) Suppose we take a smoothed multinomial classifier and double the amount of smoothing (e.g., for a variant of “add 1 smoothing”, add 2 to each count, and add to the denominator $2k$, where k is the number of samples). What qualitative effect will this have on decisions of the classifier?

- III. An IR system returns 3 relevant documents, and 2 irrelevant documents. There are a total of 8 relevant documents in the collection.

- a) What is the precision of the system on this search, and what is its recall?
- b) Instead of using recall/precision for evaluating IR systems, we could use accuracy of classification. Consider a classifier that classifies documents as being either relevant or non-relevant. The accuracy of a classifier that makes c correct decisions and i incorrect decisions is defined as: $c/(c+i)$.
 - (i) Why do the recall and precision measures reflect the utility (i.e., quality or usefulness) of an IR system better than accuracy does?
 - (ii) Suppose that we have a collection of 10 documents, and two different boolean retrieval systems A and B. Give an example of two result sets, A_q and B_q , assumed to have been returned by the system in response to a query q , constructed such that A_q has clearly higher utility and a better score for precision than B_q , but such that A_q and B_q have the same scores on accuracy.

① Document models:

Text classifiers often don't use any kind of deep representation about language. Often a document is represented as a bag of words. Consider a document D , whose class is given by c . In the case of email spam filtering there are 2 classes $c=S$ (spam) and $c=H$ (ham). We classify D as the class which has the highest posterior probability $P(c|D)$, which can be re-expressed using Bayes' Theorem:

$$P(c|D) = \frac{P(D|c)P(c)}{P(D)} \propto P(D|c)P(c)$$

There are 2 probabilistic models of documents, both of which represent documents as a bag of words, using the naive Bayes assumption. Both models represent documents using feature vectors whose components correspond to word types. If we have a vocabulary V , containing $|V|$ word types, then the feature vector dimension $d = |V|$.

Bernoulli document model: a document is represented by a feature vector with binary elements taking value 1 if the corresponding word is present in the document and 0 if the word is not present.

multinomial document model: a document is represented by a feature vector with integer elements whose value is the frequency of that word in the document.

a) Bernoulli Naive Bayes model

The likelihood of a document given a class c_k is given by

$$P(x|c_k) = \prod_{i=1}^n p_{ki}^{x_i} (1-p_{ki})^{(1-x_i)}$$

$$p(x_{\text{francisco}} = \text{true} \mid \text{class} = \text{sfo})$$

which indicates whether francisco appears in the document with class sfo. If it appears probability = 1 else 0.

Ans: $P(x_{\text{francisco}} = \text{true} \mid \text{class} = \text{sfo}) = 1.0$

$$p(x_{\text{london}} = \text{true} \mid \text{class} = \text{sfo})$$

Ans: 0.5

$$p(x_{\text{francisco}} = \text{true} \mid \text{class} = \text{jfk})$$

Ans: 1.0

(1b) Multinomial NB model.

$$P(x = \text{francisco} \mid \text{class} = \text{sfo}) = 4/14$$

$$P(x = \text{london} \mid \text{class} = \text{sfo}) = 1/14$$

$$P(x = \text{francisco} \mid \text{class} = \text{jfk}) = 1/8$$

(1c)

(i) When the Bernoulli Naive Bayes model is applied to the test set after trained on the training set it is not very accurate, because it ignores frequency information, which is important in this domain.

(ii) The multinomial model is more accurate, because it uses frequency information. However, it ignores position information, so doesn't distinguish between a city name occurring at the beginning/end of the itinerary from the one which is occurring in the middle of the

(1d) we can use as a feature the term that occurs in the last position of each document. Non-standard feature represented with using non-standard words. The non-standard words are classified to 6 categories using skip2 collection to official, literature, informative, popular, educational and scientific.

$$P(x_{\text{newyork}} = \text{true} | \text{class} = \text{JFK}) = 1.0$$

$$P(x_{\text{sanfrancisco}} = \text{true} | \text{class} = \text{SFO}) = 1.0$$

$$P(x_{\text{chicago}} = \text{true} | \text{class} = \text{ORD}) = 1.0$$

2
(a) It will never choose a category unless all words in a document were seen for that category for the training set (unless there is no category for which all words were seen, and then all categories are tied for the classifier). It will rank between classes for which all words were seen similarly to the smoothed classifier (but with possible differences due to the smoothing).

(b) Here it is given that they have doubled the amount of smoothing.

Formula for Laplace (add-1) smoothing for naive Bayes

$$\hat{P}(w_i | c) = \frac{\text{count}(w_i, c) + 1}{\sum (\text{count}(w, c) + 1)}$$

$$= \frac{\text{count}(w, c) + 1}{\left(\sum_{w \in V} \text{count}(w, c) \right) + |V|}$$

It will be more likely to choose categories for which some many of the words in the document were unseen.

III

Given that

system returns 3 relevant documents

2 irrelevant documents.

Total 8 relevant documents in the collection.

a) $\text{precision} = \frac{TP}{TP + FP}$

Here $TP = 3$

$FP = 2$

$$= \frac{3}{3+2}$$

$$= \frac{3}{5}$$

$\text{Recall} = \frac{TP}{TP + FN}$

Here $TP = 3$

$FN = 5$

$$= \frac{3}{3+5}$$

$$= \frac{3}{8}$$

Therefore $\text{precision} = \frac{3}{5}$

$\text{recall} = \frac{3}{8}$

(b) (i) An IR system which always returns no results will have high accuracy for most queries, since the corpus usually contains only a few relevant documents. Documents that are truly relevant are the only ones that will be mistakenly classified as nonrelevant, and thus the accuracy is close to 1. Recall and precision are two different measures that can jointly capture the tradeoff between returning more relevant results and returning fewer irrelevant results:

(ii) There are of course many correct answers. one simple correct answer is

Assume document 1 is the only relevant document.

$$A_q = \{1, 2, 3\}$$

$$B_q = \{3\}$$

Both A_q and B_q made 2 mistakes, so they have the same accuracy: 80%.

The precision of A_q is $1/3$, the precision for B_q is 0, since B_q didn't return any relevant documents, it is of no utility.