

CS5560 Knowledge Discovery and Management

Problem Set 6

July 10 (T), 2017

Name: *Revanth Chakram*

Class ID: *02*

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

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

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

- a) Suppose we build a Naïve 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 Naïve 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

a) Bernoulli - Naive Bayes Model

$$P(X_{\text{france}} = \text{true} \mid \text{class} = \text{SFO}) = 1.0$$

$$P(X_{\text{london}} = \text{true} \mid \text{class} = \text{SFO}) = 0.5$$

$$P(X_{\text{france}} = \text{true} \mid \text{class} = \text{JFK}) = 1.0$$

b) Multinomial - Naive Bayes Model

$$P(X = \text{france} \mid \text{class} = \text{SFO}) = 4/14 \quad (\text{assuming no tokenization of punctuation})$$

$$P(X = \text{london} \mid \text{class} = \text{SFO}) = 1/14$$

$$P(X = \text{france} \mid \text{class} = \text{JFK}) = 1/8$$

c) Considering a standard Naive Bayes classifier trained on the training set by applying it to similar set

(i) Bernoulli's model:-

'Not very accurate' model because it ignores the frequency information which is important in this domain.

(ii) Multinomial model:-

'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 item or from one occurring in the middle.

- d) Construct a non-standard feature representation that is 100% accurate for either model.

Sol: Use as a feature the term that occurs in the last position of each document.

II.

$$\begin{aligned}\hat{P}(w_i | c) &= \frac{\text{Count}(w_i, c) + 1}{\left(\sum_{w \in V} (\text{Count}(w, c)) + 1 \right)} \\ &= \frac{\text{Count}(w_i, c) + 1}{\left(\sum_{w \in V} \text{Count}(w, c) \right) + |V|}\end{aligned}$$

without smoothing:

- 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 b/w the classes for which all words were seen similarly to the smoothed classifier (but with possible differences due to the smoothing)

- b) Doubling the smoothing value:-

It will ~~not~~ be more likely to choose categories for which some / many of the words in the document were unseen.

III. Given Data :

- a) Relevant Documents # 3 } Retrieved
 Irrelevant Documents # 2 }

Relevant Documents # 8

	Retrieved	Not Retrieved	
Relevant	3 (TP)	5 (FN)	8
Non Relevant	2 (FP)	3 (TN)	5
	<u>5</u>	<u>8</u>	Total = 13

$$\begin{aligned} \text{Relevant not Retrieved} &= \text{Relevant} - \text{Retrieved Relevant doc} \\ &= 8 - 3 = 5 \text{ (FN)} \end{aligned}$$

$$\text{Precision} = \frac{tp}{tp+fp} = \frac{3}{5}$$

$$\text{Recall} = \frac{tp}{tp+fn} = \frac{3}{8}$$

- b) (i) An Info-Retrieval system always ~~result~~ 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 non-relevant and thus accuracy is close to 1. Recall and precision are two different measures that can jointly capture the tradeoff b/w returning more relevant results vs returning fewer irrelevant results.

(99) Assuming a collection of 10 documents.
Two Retrieval Systems A & B.
Let's assume, document 1 is the only relevant document

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

$B_q \rightarrow \{3\}$

Both A_q & B_q made 2 mistakes, so they have same accuracy = 80%

The precision of $A_q = 1/3$

" " $B_q = 0$.

Since B_q didn't return any relevant documents
by there is no use of it.