

CS5560 Knowledge Discovery and Management

Problem Set 6

July 10 (T), 2017

Name: SUJITHA POTHANA

Class ID: 24

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
1 a) Bernoulli Model :- It is equivalent to the binary independence model, which generates an indicator for each term in the vocabulary either "1" indicating presence of the term in the document or "0" indicating absence.

→ $P(x_{\text{francisco}} = \text{true} \mid \text{class SFO})$

• Total no. of class of SFO = 2

• Presence of term "francisco" in no. of doc = 2

$$\Rightarrow 2/2 = \boxed{1}$$

→ $P(x_{\text{london}} = \text{true} \mid \text{class SFO})$

• Total SFO classes = 2

• Presence of "london" in doc = 1

$$\Rightarrow 1/2 = \boxed{0.5}$$

→ $P(x_{\text{francisco}} = \text{true} \mid \text{class JFK})$

• Total JFK classes = 1

• Presence of "francisco" in doc = 1

$$\Rightarrow 1/1 = \boxed{1}$$

b) Multinomial Model :- It generates one term from the vocabulary in each position of the document, where we assume a generative model

	Multinomial Model	Bernoulli model
Event model	Generation of token	Generation of document
Random Variable	$X = t$ iff t occurs at given pos	$U_i = 1$ iff t occurs in doc
document representation	$d = \{t_1, \dots, t_k, \dots, t_n\}, t_k \in V$	$d = \{c_1, \dots, c_i, \dots, c_n\}, c_i \in \{0,1\}$
Parameter estimation	$\hat{P}(X=t c)$	$\hat{P}(U_i=c c)$
Decision rule: maximize	$\hat{P}(c) \prod_{k=1}^n \hat{P}(X=t_k c)$	$\hat{P}(c) \prod_{i=1}^n \hat{P}(U_i=c_i c)$
Multiple Occurrences	taken into account	ignored
Length of docs	can handle longer docs	Works best for short docs
# feature	can handle more	Works best with fewer
estimate for term freq	$\hat{P}(X=t_k c) \approx 0.5$	$\hat{P}(U_{t_k}=1 c) \approx 1.0$

$$\rightarrow P(X = \text{francisco} | \text{Class} = \text{SFO}) =$$

- Occurrence of Francisco in class SFO = 4
- Word Count in class SFO = 14

$$\Rightarrow 4/14$$

$$\rightarrow P(X = \text{london} | \text{Class} = \text{SFO})$$

- Occurrence of Word "london" = 1
- Word Count in class SFO = 14

$$\Rightarrow 1/14$$

$$\rightarrow P(X = \text{francisco} | \text{Class} = \text{JFK})$$

- Occurrence of word francisco in class JFK = 1
- Word Count in class JFK = 8

$$\Rightarrow 1/8$$

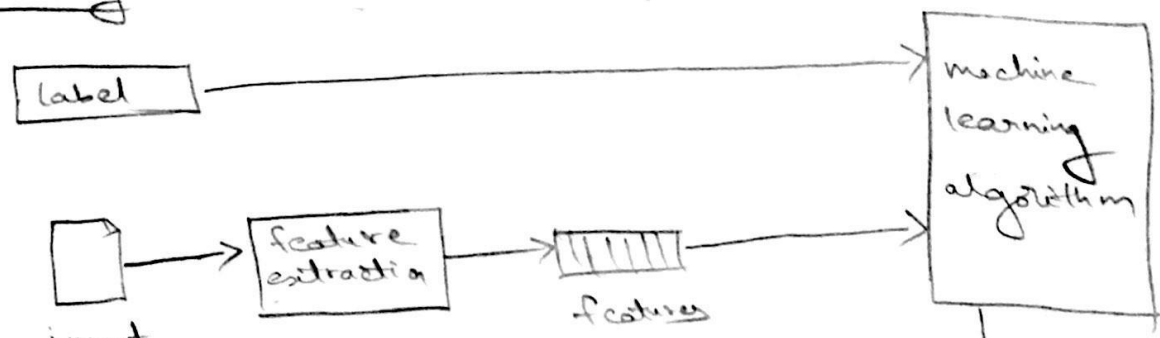
c) → Multinomial model takes into account multiple occurrence of the word and can handle more docs and features. So multinomial model is more accurate, because it uses frequency information.

→ Bernoulli Model in contrast does not consider the frequent occurrence of words, so it is not much accurate, as it models absence of terms explicitly.

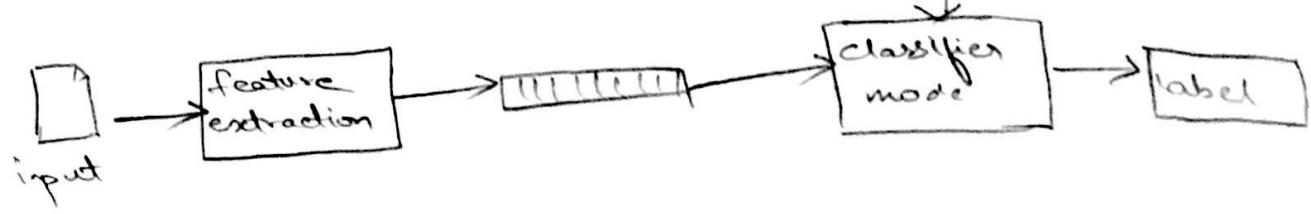
2) Classification :-

Classification is the task of choosing the correct class label for a given input.

Training



Prediction



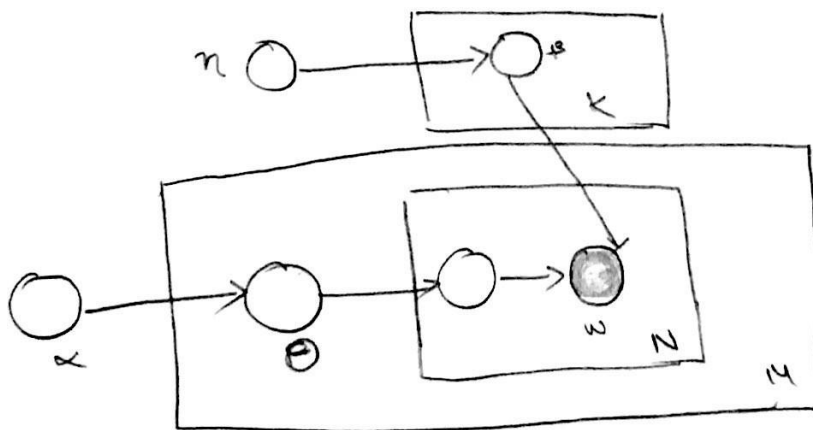
→ A classifier is called supervised if it is built based on training corpus containing the correct label for each input.

a) → In the naive bayes classifier when occurrence of word is "0" ⇒ probability of occurrence is 0
 → If $P(\text{word}|\text{class}) = 0$, then we can never choose a category.

→ Classification are based on training set, so we can rank for classes for which all words were seen similar to the smoothed classifier.

b) Smoothing Multinomial Model:-

- Large vocabulary size is characteristic of document corpora often problems with sparsity.
- Maximum likelihood estimates of the multinomial parameters θ assign zero probability to new words, and thus zero probability to new document.



Retrieval formula using general smoothing scheme $\log P(q|d) = \sum_{w \in V, c(w|q) > 0} c(w|q) \log P(w|d)$

$$P(w|d) = \begin{cases} P_{\text{seen}}(w|d) & \text{if } w \text{ is seen in } d \\ \alpha_d P(w|c) & \text{otherwise} \end{cases}$$

→ $P(\text{word}|\text{count})$

Let initial word = w
count = c

$$\Rightarrow \frac{1}{k}$$

→ Increase in smoothing

$$\Rightarrow \frac{1}{2k}$$

→ It'll be more likely to choose categories for which many of the words in the document

III) Total positive (tp) = 5

Total (tp + fp) = 8

	↓ Relevant	↓ Irrelevant
Retrieved →	True positive ③ relevant, retrieved	False positive ② irrelevant, retrieved
Not retrieved →	False negative ⑤ relevant, not retrieved	True Negative irrelevant, not retrieved

$$\text{Precision} = \frac{tp}{(tp + fp)} = \frac{3}{3+2} = \frac{3}{5}$$

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

b) (i)

- Precision, indicates how many of the items that we identified were relevant
- Recall, indicates how many of the relevant items that we identified.

	Relevant	Irrelevant
Retrieved	TP	FP
Not Retrieved	FN	TN

- For the information retrieval system, returns no results will have high accuracy for most queries, since the corpus usually contains only a few relevant documents.
- Recall & precision are two different measures that can jointly capture the tradeoff between retrieving more relevant results and retrieving fewer irrelevant results

(ii) Precision (A_q) > Precision (B_q)

For A_q

	Relevant	Irrelevant
Retrieved	1	2
Not retrieved	1	

$$\Rightarrow \text{Precision}(A_q) = \frac{1}{3}$$

For B_q

	Relevant	Irrelevant
Retrieved	0	2
Not retrieved	2	

$$\Rightarrow \text{Precision}(B_q) = 0$$

- In the above case both Aq & Bq made 2 mistakes \Rightarrow accuracy = 80% (7)
- \rightarrow Since Bq didn't return any relevant documents, it is of novelty.

(1d) 100% accurate Naive Bayes Model

- \rightarrow This can be achieved by using the term that occurs in the last position of each document.
- \rightarrow Non standard feature represented with using non-standard words. The non-standard words are classified to 6 categories using SKIPE \geq collection to official, literature, information, popular, educative and scientific.

$$P(X_{\text{New York}} = \text{true} | \text{class} = \text{JFK}) = 1.0$$

$$P(X_{\text{San Francisco}} = \text{true} | \text{class} = \text{SFO}) = 1.0$$

$$P(X_{\text{Chicago}} = \text{true} | \text{class} = \text{ORD}) = 1.0.$$