

Project 1 Naïve Bayes

COMP4901K and MATH 4824B
Fall 2018

Notes:

- Due: October 18 at 23:59. **No late submission is accepted.**
- Answers Release: October 19.
- Submission in PDF (only) via canvas (emails with attachments will be ignored).
- Only typed PDF documents are accepted, no handwritten scanned documents.

1 Bayes' Rule

Q1 Probabilities Computation (4*5 points)

Suppose you are given a chance to win bonus grade points:

There are three boxes. Only one box contains a special prize that will grant you 1 bonus points. After you have chosen a box B_1 (B_1 is kept closed), one of the two remaining boxes will be opened (called B_2) such that it must not contain the prize (note that there is at least one such box).

- (i) What is the probability of B_1 contains prize, $P(B_1 = 1)$?
- (ii) What is the likelihood probability of B_2 does not contains prize if B_1 contains prize, $P(B_2 = 0|B_1 = 1)$?
- (iii) What is the probability of B_1 contains prize given B_2 does not contain prize, $P(B_1 = 1|B_2 = 0)$?
- (iv) Now you are are given a second chance to choose boxes. You can either stick to B_1 or choose the only left box B_3 . According to the Bayes decision rule, should you change your choice or not?

2 Naïve Bayes Classifier

Q2 Building a Naïve Bayes Classifier for Text (4*10 points)

Suppose you want to build a Naïve Bayes classifier for text. First of all, you represent a text document as if it were a bag-of-words. Naïve Bayes is a probabilistic classifier, meaning that for a document d , out of all classes $k \in \{1, 2, \dots, K\}$ the classifier returns the class \hat{y} which has the maximum probability given the document.

$$\hat{y} = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} P(y = k|d) \quad (1)$$

This idea of Bayesian inference has been known since the work of Bayes (1763), and was first applied to text classification by Mosteller and Wallace (1964). The intuition of Bayesian classification is to use Bayes' rule to transform Eq.(1) into other probabilities that have some useful properties. Bayes' rule is presented in Eq.(2); it gives us a way to break down any conditional probability $P(x|y)$ into three other probabilities:

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)} \quad (2)$$

We can then substitute Eq.(2) into Eq.(1) to get Eq.(3):

$$\hat{y} = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} P(y = k|d) = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} \frac{P(d|y = k)P(y = k)}{P(d)} \quad (3)$$

Because $P(d)$ doesn't change for each class, we can choose the class that maximizes this simpler formula:

$$\hat{y} = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} P(y = k|d) = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} P(d|y = k)P(y = k) \quad (4)$$

We thus compute the most probable class \hat{y} given some document d by choosing the class which has the highest product of two probabilities: the probability of the class $P(y = k)$ and the likelihood of the document $P(d|y = k)$:

$$\hat{y} = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} \underbrace{P(d|y = k)}_{\text{likelihood}} \overbrace{P(y = k)}^{\text{probability of class } k} \quad (5)$$

To apply the Naïve Bayes classifier to text, we need to consider how to represent documents. Bag of words is an easy method to model documents. And it assumes the position doesn't matter. Thus we can consider word positions, by simply walking an index through every word position in the document:

$$\hat{y} = \operatorname{argmax}_{k \in \{1, 2, \dots, K\}} P(y = k) \prod_{i \in \text{positions}} P(w_i|y = k) \quad (6)$$

Naïve Bayes calculations, like calculations for language modeling, are done in log space, to avoid underflow and increase speed. Thus Eq.(6) is generally instead expressed as

$$\hat{y} = \underset{k \in \{1, 2, \dots, K\}}{\operatorname{argmax}} \log P(y = k) + \sum_{i \in \text{positions}} \log P(w_i | y = k) \quad (7)$$

By considering features in log space Eq.(7) computes the predicted class as a linear function of input features (here are words in the document). Classifiers that use a linear combination of the inputs to make a classification decision are called linear classifiers.

- (i) Let N_k be the number of documents in our training data with class k and N_{doc} be the total number of documents. What is the maximum likelihood estimate of the probability of each class $P(y = k)$?
- (ii) We use bag-of-words representation. We first concatenate all documents with class k into one big text. Let V be the union of all the words in all classes, w_i be the i -th word in the document, and $\text{count}(w_i, k)$ be the frequency of w_i . Then what is the maximum likelihood estimate of the probability $P(w_i | y = k)$?
- (iii) But suppose there are no training documents that both contain the word "fantastic" and are classified as class k . In such a case the probability for this feature will be zero. But since Naive Bayes naively multiplies all the feature likelihoods together, zero probabilities in the likelihood term for any class will cause the probability of the class to be zero. What will you do to avoid this situation? Please write down your solution and improved formula.
- (iv) Because some words did not occur in any training document in any class, they are not in our vocabulary V . What should you do for such unknown words? And very frequent words like *the* and *a* and number are unhelpful for prediction, what should you do for such words?

Q3 Running the Classifier on an Example (4*10 points)

According to what you have done in Q2, we will use a sentiment analysis domain with the two classes positive (+) and negative (-), and take the following miniature training and test documents. You need to compute:

- (i) the probability of each class, $P(y = +)$ and $P(y = -)$
- (ii) the probabilities for every word in the test document, $P(w_i | y = +)$ and $P(w_i | y = -)$
- (iii) apply the strategy in Q2 (iii) to Q3 (i) and Q3 (ii)
- (iv) the probabilities for each class, $P(y = + | d)$ and $P(y = - | d)$, and your final decision, positive or negative

Table 1: training and test documents

	Class	Documents
Training	+	great
	−	very annoying
	+	the best one
	−	this version is rubbish
	+	awesome
	+	this version is easy to use
	+	good
	−	so terrible
Test	?	very easy and great

You are required to solve these questions with the help of code in Lab 4. Then you can compute by hand to check the results.

You don't need to handle very frequent words!