

ml-classifier

More Information on the Naive Bayes Classifier

This document provides a more complete walkthrough of the natural language processing (NLP) and machine learning (ML) techniques behind the classifier implemented in this project, in case you're interested. Nothing in this document should be required to implement the project - it's just extra information if you're interested. There is also considerable overlap between this document and what is described in the [main specification](#).

At a high level, the classifier we implement works by assuming a probabilistic model of how Piazza posts are composed, and then finding which **label** (e.g. our categories of "euchre", "exam", etc.) is the most probable source of a particular post.

There are many different kinds of classifiers that have subtle differences. The classifier we describe here is a version of a "Multi-Variate Bernoulli Naive Bayes Classifier".

Background: Conditional Probabilities and Bayes Theorem

We write $P(A)$ to denote the probability (a number between 0 and 1) that some event A will occur. $P(A|B)$ denotes the probability that event A will occur given that we already know event B has occurred. For example, if event A is "it will rain today", we might guess that $P(A)$ (assuming we hadn't looked outside or at the weather forecast). However, if we were to look outside and observe event B , "it is cloudy", then we might believe it is more likely that it will rain. So we could say something like $P(A|B) > P(A)$.

The Bag of Words Model

We will treat a Piazza post as a "bag of words" - each post is simply characterized by which words it includes. The ordering of words is ignored, as are multiple occurrences of the same word. These two posts would be considered equivalent:

- "the left bower took the trick"
- "took took trick the left bower bower"

Thus, we could imagine the post generation process as a person sitting down and going through every possible word and deciding which to toss into a bag.

Which Label is Most Likely?

Given a post p , we must determine the most probable label from which it could have been generated. This would be whatever label **maximizes the probability** of p . Using Bayes' theorem to rewrite the probability, we have:

$$P(L = l | p) = \frac{P(p | L = l) P(L = l)}{P(p)}$$

The denominator $P(p)$ does not depend on l , so it will be the same for all labels and can be ignored. Thus we want to find the label l that maximizes the **probability score**:

These two quantities have intuitive meanings for our application.

- $P(L = l)$ is the **prior probability** of label l . It is the probability a randomly selected post has label as l . (i.e. How common is each label?)
- $P(p | L = l)$ is the **likelihood** of post p given label l . If someone sat down and wrote a post for label l this is the probability that post p would result.

Let's consider the likelihood $P(p | L = l)$ in more detail. According to our generative process and bag of words model, we could think of the probability of the whole post p as the probability of the particular combination of words chosen to put into the bag:

This is the **joint probability** of several individual events for each word that was included in the post. Technically, we could also consider events for each word that was left out, but this turns out not to matter in many cases (and specifically for our Piazza classifier).

However, we have a problem - the joint distribution is difficult to learn from data. An informal way to think about this is that the compound event of several different words occurring is quite rare. Although we may have seen many posts about euchre, we may never have seen precisely the post "the left bower took the trick".

To solve this problem we introduce the **Naive Bayes assumption**, which posits that the occurrence of each word is **conditionally independent** from the occurrence of other words, given a particular

label. That is, for posts within a particular label, the presence of one word doesn't influence the chance of seeing another. We can rewrite the likelihood of the whole post simply as the product of the likelihoods of its words:

Of course, this is not completely true, and this is why the assumption is called "naive". For example, even given the class "euchre", occurrences of the words "left" and "bower" are going to be related. That said, the assumption is close to true for many words (e.g. "dealer" and "trick" given the class "euchre"), and Naive Bayes classifiers work well for many applications in practice. It is also much easier to estimate individual word likelihoods by learning from data. We'll address this in the next section.

But first, one crucial detail. Our goal is now to find the label that maximizes:

But computing this product can be problematic due to the limited precision of floating point numbers. The probability of a particular word occurring is generally pretty low (e.g.

), and when we multiply many of them together the result becomes very close to zero and problems with limited floating-point precision can occur.

We avoid this problem using a neat trick: **work with the natural logarithm of probabilities instead of the probabilities themselves**. (In C++, `std::log()` gives the natural logarithm.) Because the logarithm is a monotonically increasing function, our goal is still to find the label with the highest log-probability, but we can remove the multiplications that can cause underflow problems with this property of logarithms:

So we must find the label with the highest **log-probability** score given the post:

(Note: If multiple classes are tied, predict whichever comes first alphabetically.)

Important: Because we're using the bag-of-words model, the words are only the unique words in the post, not including duplicates!

Learning Classifier Parameters (Training the Classifier)

To compute the log-probability score, we need concrete values for $\theta_{w,l}$ and θ_l for each word w . These are called the **classifier parameters**. To find these, we train the classifier on a dataset of already labeled Piazza posts (this is called supervised learning). Based on observations from the training dataset, the classifier can estimate the parameters. If the dataset is large enough, these estimates should be quite good.

To estimate the **log-prior probability** of a label l , the classifier should observe the proportion of posts with that label in the training set:

$$\theta_l = \frac{\text{number of posts with label } l}{\text{total number of posts}}$$

To estimate the **log-likelihood** of a word w given a label l , the classifier should observe the proportion of posts with label l that contain the word w :

$$\theta_{w,l} = \frac{\text{number of posts with label } l \text{ that contain } w}{\text{number of posts with label } l}$$

If we need to compute $\theta_{w,l}$, but w was never seen in a post with label l in the training data, we get likelihood of 0 and a corresponding log-likelihood of $-\infty$. So instead, use this alternate formula for occurrences of the word through the entire training set:

$$\theta_{w,l} = \frac{\text{number of posts with label } l \text{ that contain } w + 1}{\text{number of posts with label } l + \text{number of words in training set}}$$

(Use when w does not occur in posts labeled l but does occur in the training data overall.)

If the word has never been seen anywhere in the entire training set we just pretend the word had been seen in one document:

$$\theta_{w,l} = \frac{1}{\text{number of words in training set} + 1}$$

(Use when w does not occur anywhere at all in the training set.)

[Return to the main specification.](#)

