

7

Naive Bayes

Lab Objective: *Create a Naïve Bayes Classifier. Use this classifier, and Sklearn's premade classifier to make an SMS spam filter.*

About Naïve Bayes

Naïve Bayes classifiers are a family of machine learning classification methods that use Bayes' theorem to probabilistically categorize data. They are called naïve because they assume independence between the features. The main idea is to use Bayes' theorem to determine the probability that a certain data point belongs in a certain class, given the features of that data. Despite what the name may suggest, the naïve Bayes classifier is not a Bayesian method. This is because naïve Bayes is based on likelihood rather than Bayesian inference.

While naïve Bayes classifiers are most easily seen as applicable in cases where the features have, ostensibly, well defined probability distributions (such as classifying sex given physical characteristics), they are applicable in many other cases. While it is generally a bad idea to assume independence naïve Bayes classifiers are still very effective, even when we can be confident there is nonzero covariance between features.

The Classifier

You are likely already familiar with Bayes' Theorem, but we will review how we can use Bayes' Theorem to construct a robust machine learning model.

Given the feature vector of a piece of data we want to classify, we want to know which of all the classes is most likely. Essentially, we want to answer the following question

$$\operatorname{argmax}_{k \in K} P(C = k | \mathbf{x}), \quad (7.1)$$

where C is the random variable representing the class of the data. Using Bayes' Theorem, we can reformulate this problem into something that is actually computable. We find that for any $k \in K$ we have

$$P(C = k | \mathbf{x}) = \frac{P(C = k)P(\mathbf{x} | C = k)}{P(\mathbf{x})}.$$

Now we will examine each feature individually and use the chain rule to expand the following expression

$$\begin{aligned} P(C = k)P(\mathbf{x} | C = k) &= P(x_1, \dots, x_n, C = k) \\ &= P(x_1 | x_2, \dots, x_n, C = k)P(x_2, \dots, x_n, C = k) \\ &= \dots \\ &= P(x_1 | x_2, \dots, x_n, C = k)P(x_2 | x_3, \dots, x_n, C = k) \cdots P(x_n | C = k)P(C = k), \end{aligned}$$

and applying the assumption that each feature is independent we can drastically simplify this expression to the following

$$P(x_1 | x_2, \dots, x_n, C = k) \cdots P(x_n | C = k) = \prod_{i=1}^n P(x_i | C = k).$$

Therefore we have that

$$P(C = k | \mathbf{x}) = \frac{P(C = k)}{P(\mathbf{x})} \prod_{i=1}^n P(x_i | C = k),$$

which reforms Equation 7.1 as

$$\operatorname{argmax}_{k \in K} P(C = k | \mathbf{x}) = \operatorname{argmax}_{k \in K} P(C = k) \prod_{i=1}^n P(x_i | C = k). \quad (7.2)$$

We drop the $P(\mathbf{x})$ in the denominator since it is not dependent on k .

This problem is approximately computable, since we can use training data to attempt to find the parameters which describe $P(x_i | C = k)$ for each i, k combination, and $P(C = k)$ for each k . In reality, a naïve Bayes classifier won't often find the actual correct parameters for each distribution, but in practice the model does well enough to be robust. Something to note here is that we are actually computing $P(C = k | \mathbf{x})$ by finding $P(C = k, \mathbf{x})$. This means that naïve Bayes is a generative classifier, and not a discriminative classifier.

Spam Filters

A spam filter is just a special case of a classifier with two classes: spam and not spam (or ham). We can now describe in more detail how we are to calculate Equation 7.2 given that we know what the features are. We can use a labeled training set to determine $P(C = \text{spam})$ the probability of spam and $P(C = \text{ham})$ the probability of ham. To do this we will assume that the training set is a representative sample and define

$$P(C = \text{spam}) = \frac{N_{\text{spam}}}{N_{\text{samples}}}, \quad (7.3)$$

and

$$P(C = \text{ham}) = \frac{N_{\text{ham}}}{N_{\text{samples}}}. \quad (7.4)$$

Using a bag of words model, we can create a simple representation of $P(x_i | C = k)$ where x_i is the i^{th} word in a message, and therefore \mathbf{x} is the entire message. This results in the simple definition of

$$P(x_i | C = k) = \frac{N_{\text{occurrences of } x_i \text{ in class } k}}{N_{\text{words in class } k}}. \quad (7.5)$$

Note that the denominator in Equation 7.5 is not the number of unique words in class k , but the total number of occurrences of any word in class k . In the case we have some word x_u that is not found in the training set, we can may choose $P(x_u | C = k)$ so that the computation is not effected, i.e. letting $P(x_u | C = k) = 1$ for unique words.

A First Model

When building a naïve Bayes classifier we need to choose what probability distribution we believe our features to have. For this first model, we will assume that the words are a categorically distributed random variable. This means the random variable may take on say N different values, each value has a certain probability of occurring. This distribution can be thought of as a Bernoulli trial with N outcomes instead of 2.

In our situation we may have N different words which we expect may occur in a spam or ham message, so we need to use the training data to find each word and its associated probability. In order to do this we will make a DataFrame that will allow us to calculate the probability of the occurrence of a certain word x_i based on what percentage of words in the training set were that word x_i . This DataFrame that will allow us to more easily compute Equation 7.5, assuming the words are categorically distributed. While we are creating this DataFrame, it will also be a good opportunity to compute Equations 7.3 and 7.4.

Throughout the lab we will use an SMS spam dataset contained in `sms_spam_collection.csv`. The following code makes full test and train sets, but we will also provide you with code to check against specific subsets.

```
>>> import pandas as pd
>>> from sklearn.model_selection import train_test_split

>>> # load in the sms dataset
>>> df = pd.read_csv('sms_spam_collection.csv')

>>> # separate the data into the messages and labels
>>> X = df.Message
>>> y = df.Label

>>> # split the data into test and train sets
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=.7)
```

Training The Model

Problem 1. Create a class `NaiveBayesFilter`, with an `__init__()` method that is empty. Add a `fit()` method which takes as arguments `X`, the training data, and `y` the training labels. In this case `X` is a `pandas.Series` containing strings that are SMS messages. Create a new `DataFrame` with two rows and a column for each vocabulary word with `'spam'` and `'ham'` being the index. Each entry will be the number of times a word appears in spam or ham messages.

For example, `self.data.loc['ham', 'in']` is the number of times the word "in" appears in ham messages. Save this DataFrame as `self.data`.

Hint: be sure you count the number of occurrences of a word and not a string. For example, when searching the string `'find it in there'` for the word `'in'`, make sure you get 1 and not 2 (because of the `'in'` in `'find'`). The methods `pd.Series.str.split()` and `count()` may be helpful.

```
>>> # checkout what the DataFrame looks like
```

```

>>> NB = NaiveBayesFilter()
>>> NB.fit(X[:300], y[:300])
>>> NB.data.loc['ham', 'in']
47
>>> NB.data.loc['spam', 'in']
4

```

Predictions

Now that we have implemented the `fit()` method, we can begin to classify new data. We will do this with two methods, the first will be a method that calculates $P(S | \mathbf{x})$ and $P(H | \mathbf{x})$, and the other will determine the more likely of the two and assign a label. While it may seem like we should have $P(C = S | \mathbf{x}) = 1 - P(C = H | \mathbf{x})$, we do not. This would only be true if we assume the S and H are independent of \mathbf{x} . It is clear that we shouldn't make this assumption, because we are trying to determine the likelihood of S and H based on what \mathbf{x} tells us. Therefore we must compute both $P(C = S | \mathbf{x})$ and $P(C = H | \mathbf{x})$.

Problem 2. Implement the `predict_proba()` method in your naïve Bayes classifier. It should take as an argument \mathbf{X} , the data that needs to be classified, and it will compute the product portion of equation 7.2.

Notice that $P(x_i | C)$ is the same for every repeated instance of word x_i in message \mathbf{x} . To save time, we only want to calculate this probability once. To do this, find

$$\prod_{i=1}^l P(x_i | C)^{n_i}$$

for each message \mathbf{x} in \mathbf{X} where l is the number of unique words in the message and n_i is the number of times the i^{th} unique word (x_i) occurs.

The method should return an $(N, 2)$ array, where N is the length of \mathbf{X} , whose entries are the probabilities of each message \mathbf{x} in \mathbf{X} belonging to each category. The first column corresponds to $P(C = H | \mathbf{x})$, and the second to $P(C = S | \mathbf{x})$.

Problem 3. Implement the `predict()` method in your naïve Bayes classifier. This should take as an argument \mathbf{X} , the data that needs to be classified. Using `predict_proba()`, finish implementing equation 7.2 and return an array that classifies each message \mathbf{x} in \mathbf{X} .

```

>>> # create the filter
>>> NB = NaiveBayesFilter()

>>> # fit the filter to the first 300 data points
>>> NB.fit(X[:300], y[:300])

>>> # test the predict function

```

```
>>> NB.predict(X[530:535])
array(['ham', 'spam', 'ham', 'ham', 'ham'], dtype=object)
```

Underflow

There are some issues that we encounter given this implementation. Notice that in the following example, the likelihoods for both spam and ham are 0 for each message.

```
>>> # find the likelihoods for messages 1085 and 2010
>>> NB.predict_proba(X[[1085,2010]])
array([[0., 0.],
       [0., 0.]])
```

This is because the messages are long, and thus involve the product of many numbers that are between 0 and 1. Because of this, we have encountered what is called underflow, where a number becomes so small it is not machine representable. Therefore, we should work in logspace, as to avoid inevitable underflow caused by long messages. If we take the log of equation 7.2 have

$$\operatorname{argmax}_{k \in K} \ln(P(C = k)) + \sum_{i=1}^n \ln(P(x_i | C = k)), \quad (7.6)$$

and this problem is still valid since logarithms are monotonically increasing. However, if any of the $P(x_i | C = k)$ are close to 0, we risk getting an overall value of $-\infty$. To prevent this from happening, we can perform *Laplace add-one smoothing* by adding 1 to the numerator of $P(x_i | C = k)$ and 2 to its denominator. This method is equivalent to using a Bayesian method for training. Thus, equation 7.5 becomes

$$P(x_i | C = k) = \frac{N_{\text{occurrences of } x_i \text{ in class } k} + 1}{N_{\text{words in class } k} + 2}. \quad (7.7)$$

Problem 4. Implement `predict_log_proba()` and `predict_log()` using equations 7.6 and 7.7. These methods will take the same arguments and return the same object types as the methods `predict_proba()` and `predict()`, respectively.

Notice how `X[[1085,2010]]` is now classifiable.

The Poisson Model

Now that we've examine one way to constructing a naïve Bayes classifier, let us look at one more method. In the Poisson model we assume that each word is Poisson random variable, occurring with potentially different frequencies among spam and ham messages. Because each of the messages is a different length, we can reparameterize the Poisson PMF to the following

$$P(n_i = x) = \frac{(rn)^x e^{-rn}}{x!} \quad (7.8)$$

where n_i is the number of times word i occurs in a message, n is the length of the message, and $\lambda = rn$ is the classical Poisson rate. In this case r represents the number of events per unit time/space/etc.

We could easily refactor this model to use Bayesian inference to determine r , which would allow greater control over the model. This would also create a condition where the training data doesn't completely determine the model's viability. However, in this lab we will use maximum likelihood estimation to determine r .

Training the Model

Similar to the other classifier that we made, training the model amounts to using the training data to determine how $P(x_i | C = k)$ is computed, as well as computing $P(C = k)$. As stated earlier, we will attempt to find the most likely value of r for each word that appears in the training set. To do this we will use maximum likelihood estimation. The parameter we choose is the one that maximizes the likelihood function

$$\hat{r} = \operatorname{argmax}_r L(r | \mathbf{x}) = \operatorname{argmax}_r P(\mathbf{x} | r).$$

In this case, since we are using a Poisson distribution (7.8) for each word, we will solve the following problem for both the spam class and the ham classes

$$r_{i,k} = \operatorname{argmax}_{r \in [0,1]} \frac{(rN_k)^{n_i} e^{-rN_k}}{n_i!}, \quad (7.9)$$

where $r_{i,k}$ is the Poisson rate for the i^{th} word in class k (either spam or ham), N_k is the total number of words in class k , and n_i is the number of times the i^{th} word occurs in class k . We have $r \in [0, 1]$ because a word cannot occur more than once per word in the message. If we take the derivative of the right side of equation 7.9 with respect to r , set it equal to 0, and solve for the maximizing r , we find that $r_{i,k} = n_i/N_k$.

Predictions

Making predictions with this model is exactly the same as we did earlier. To clarify the calculation, lets reformulate 7.6 to fit the Poisson case better. This gives

$$\operatorname{argmax}_{k \in K} \ln(P(C = k)) + \sum_{i=1}^l \ln \left(\frac{(r_{i,k}n)^{n_i} e^{-r_{i,k}n}}{n_i!} \right), \quad (7.10)$$

with l being the number of unique words in the message, n_i the number of times the i^{th} word occurs in the message, n the total number of words in the message, and $r_{i,k}$ the Poisson rate of the i^{th} word in class k . Notice, if $r_{i,k}$ is close to 0, we'll risk getting a total value of $-\infty$. We can fix this by using the *Laplace add-one smoothing* method as we did before, but this time on $r_{i,k}$. Thus, our new Poisson rate for the i^{th} word in class k becomes

$$r_{i,k} = \frac{n_i + 1}{N_k + 2}, \quad (7.11)$$

which has a Bayesian interpretation, as it did before.

Problem 5. Create a new class called `PoissonBayesFilter` with an `__init__()` method that may be empty. Add a `fit()` method which takes as arguments training data `X` and training labels `y`.

Implement `fit()` by finding the MLE found in equation 7.11 to predict r for each word in both the spam and ham classes, thereby training the model. Store these computed rates in dictionaries called `self.spam_rates` and `self.ham_rates`, where the key is the word and the value is the associated r .

For example, `self.ham_rates['in']` will give the computed r value for the word "in" found in ham messages.

```
>>> #create a poisson bayes object to examine it
>>> PB = PoissonBayesFilter()
>>> PB.fit(X[:300], y[:300])

>>> # check spam and ham rate of 'in'
>>> PB.ham_rates['in']
0.012588512981904013
>>> PB.spam_rates['in']
0.004166666666666667
```

Problem 6. Implement the `predict_log_proba()` and `predict()` methods using equation 7.10. These methods will take the same arguments and return the same object types as the methods `predict_proba()` and `predict()` in the `NaiveBayesFilter` class, respectively. You may use `scipy.stats.poisson.pmf` if you wish.

Naive Bayes with Sklearn

Now that we've explored a few ways to implement our own naïve Bayes classifier, we can examine some robust tools from the `sklearn` library that will accomplish all the things we've coded up in a very simple manner.

The first thing we need to do is create a dictionary and transform the training data, which is what our first `fit()` method did. We instantiate a `CountVectorizer` model from `sklearn.feature_extraction.text`, and then use the `fit_transform()` method to create the dictionary and transform the training data.

```
>>> from sklearn.feature_extraction.text import CountVectorizer

>>> vectorizer = CountVectorizer()
>>> train_counts = vectorizer.fit_transform(X_train)
```

Now we can use the transformed training data to fit a `MultinomialNB` model from `sklearn.naive_bayes`.

```
>>> from sklearn.naive_bayes import MultinomialNB

>>> clf = MultinomialNB()
>>> clf = clf.fit(train_counts, y_train)
```

Testing data we want to classify must first be transformed by our vectorizer with the `transform()` method (not the `fit_transform()` method). We can then classify the data using the `predict()` method of the `MultinomialNB` model.

```
>>> test_counts = vectorizer.transform(X_test)
>>> labels = clf.predict(test_counts)
```

This naïve Bayes model uses the multinomial distribution where we have used the categorical and poisson distributions. Multinomial is very similar to the categorical implementation, as the multinomial distribution models the outcome of n categorical trials (in the same way that the binomial distribution models n Bernoulli trials).

Problem 7. Write a function that will classify messages. It will take as arguments training data `X_train` and `y_train`, and test data `X_test`. In this function use the `CountVectorizer` and `MultinomialNB` from `sklearn` and return the predicted classification of the model.

The results of Problem 7 can help you test the two Bayes Filters you created in this lab. Using the `accuracy_score` method of `sklearn.metrics`, you can compare your predicted labels with the ones from Problem 7. You should have very high accuracy, as demonstrated below.

```
>>> from sklearn.metrics import accuracy_score

>>> # labels returned by Problem 7
>>> actual_labels = sklearn_method(X_train, y_train, X_test)

>>> # test against NaiveBayesFilter
>>> NB = NaiveBayesFilter()
>>> NB.fit(X_train, y_train)
>>> NB_labels = NB.predict_log(X_test)
>>> accuracy_score(actual_labels, NB_labels)
0.9769289925660087

>>> # test against PoissonBayesFilter
>>> PB = PoissonBayesFilter()
>>> PB.fit(X_train, y_train)
>>> PB_labels = PB.predict(X_test)
>>> accuracy_score(actual_labels, PB_labels)
0.9782107152012305
```

References

Rish, Irina. (2001). An Empirical Study of the Naïve Bayes Classifier. IJCAI 2001 Work Empir Methods Artif Intell. 3.

Data from: <http://www.dt.fee.unicamp.br/~tiago/smsspamcollection/>