

10-401 Machine Learning: Homework 2

Due 5:00 p.m. Wednesday, February 21, 2018

Instructions

- **Submit your homework on time electronically by submitting to Autolab by 5:00 pm, Wednesday, February 21, 2018.**

We recommend that you use L^AT_EX, but we will accept scanned solutions as well. On the Homework 2 Autolab page, you can click on the “download handout” link to download the tar archive containing Octave .m files for each programming question and some datasets you will need. Replace each of these files with your solutions for the corresponding problem, create a new tar archive of the top-level directory, and submit your archived solutions online by clicking the “Submit File” button. You should submit a file called hw2.tar. Inside that should be the directory hw2/. Inside of that should be all of your code files. Don’t submit the data files, or you’ll get an error because your submission is too big.

DO NOT change the name of any of the files or folders in the submission template. In other words, your submitted files should have exactly the same names as those in the submission template. Do not modify the directory structure.

- **Late homework policy:** Homework is worth full credit if submitted before the due date. Up to 50 % credit can be received if the submission is less than 48 hours late. The lowest homework grade at the end of the semester will be dropped. Please talk to the instructor in the case of extreme extenuating circumstances.
- **Collaboration policy:** You are welcome to collaborate on any of the questions with anybody you like. However, you must write up your own final solution, and you must list the names of anybody you collaborated with on this assignment.

Problem 1: Independent events and Bayes Theorem

1. [8 Points] For events A, B prove the following expression using only the Law of Total Probability and the definition of conditional probabilities:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\neg A)P(\neg A)}.$$

($\neg A$ denote the event that A does not occur.) You should work from the left hand side to get to the right hand side.

2. We define three random variables in $\{0, 1\}$ ($0 = \text{false}$, $1 = \text{true}$): person X likes tartan, person X is hard-working, and person X attends CMU. The probabilities are given below:

	\neg attends CMU		attends CMU	
	\neg likes tartan	likes tartan	\neg likes tartan	likes tartan
\neg hard-working	0.15	0.15	0.05	0.1
hard-working	0.125	0.125	0.1	0.2

For example, $P(\neg \text{tartan}, \text{hard-working}, \neg \text{CMU}) = 0.125$.

- (a) [5 Points] Is “Person X likes tartan” independent of “person X is hard-working”? Why or why not?
- (b) [5 Points] Is “person X likes tartan” conditionally independent of “person X is hard-working” given “person X attends CMU”? Why or why not?
- (c) [5 Points] Remember the variable are either 0 or 1. Calculate $\mathbb{E}(\text{CMU} | \text{hard-working} = 1)$

Problem 2: Maximum Likelihood Estimation

This problem explores maximum likelihood estimation (MLE), which is a technique for estimating an unknown parameter of a probability distribution based on observed samples. Suppose we observe the values of n iid¹ random variables X_1, \dots, X_n drawn from a single Bernoulli distribution with parameter θ . In other words, for each X_i , we know that

$$P(X_i = 1) = \theta \quad \text{and} \quad P(X_i = 0) = 1 - \theta.$$

Our goal is to estimate the value of θ from these observed values of X_1 through X_n .

For any hypothetical value $\hat{\theta}$, we can compute the probability of observing the outcome X_1, \dots, X_n if the true parameter value θ were equal to $\hat{\theta}$. This probability of the observed data is often called the *data likelihood*, and the function $L(\hat{\theta}) = P(X_1, \dots, X_n | \hat{\theta})$ that maps each $\hat{\theta}$ to the corresponding likelihood is called the *likelihood function*. A natural way to estimate the unknown parameter θ is to choose the $\hat{\theta}$ that maximizes the likelihood function. Formally,

$$\hat{\theta}^{\text{MLE}} = \underset{\hat{\theta}}{\operatorname{argmax}} L(\hat{\theta}).$$

¹iid means Independent, Identically Distributed.

Often it is more convenient to work with the log likelihood function $\ell(\hat{\theta}) = \log L(\hat{\theta})$. Since the log function is increasing, we also have

$$\hat{\theta}^{\text{MLE}} = \underset{\hat{\theta}}{\operatorname{argmax}} \ell(\hat{\theta}).$$

1. **[6 Points]** Write a formula for the log likelihood function, $\ell(\hat{\theta})$. Your function should depend on the random variables X_1, \dots, X_N , the hypothetical parameter $\hat{\theta}$, and should be simplified as far as possible (i.e., don't just write the definition of the log likelihood function).
2. **[6 Points]** Compute a closed form expression for the maximum likelihood estimate (hint: recall that x^* is a critical point of $f(x)$ if $f'(x^*) = 0$. If you use part (1), remember to argue that finding the maximizer of $\log f(x)$ yields the maximizer of $f(x)$).

Use your formula to compute $\hat{\theta}^{\text{MLE}}$ for the following sequence of 10 samples:

$$X = (0, 1, 1, 1, 1, 0, 1, 0, 1, 1).$$

3. **(Bonus) [5 Points]** Now we will consider a related distribution. Suppose we observe the values of m iid random variables Y_1, \dots, Y_m drawn from a single Binomial distribution $B(n, \theta)$. A Binomial distribution models the number of 1's from a sequence of n independent Bernoulli variables with parameter θ . In other words,

$$P(Y_i = k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} = \frac{n!}{k!(n-k)!} \cdot \theta^k (1 - \theta)^{n-k}.$$

Write a single formula for the log likelihood function, $\ell(\hat{\theta})$ for Y_1, \dots, Y_m whose values are k_1, \dots, k_m , respectively. Your function should depend on m, n, k_1, \dots, k_m , and $\hat{\theta}$.

4. **(Bonus) [5 Points]** Compute a closed form expression for the maximum likelihood estimate. Then use your formula to compute $\hat{\theta}^{\text{MLE}}$ for the two Binomial random variables Y_1 and Y_2 which both have parameters $n = 5$ and θ . The Bernoulli variables for Y_1 and Y_2 resulted in $(0, 1, 1, 1, 1)$ and $(0, 1, 0, 1, 1)$, respectively, so $Y_1 = 4$ and $Y_2 = 3$.
5. **(Bonus) [5 Points]** How do your answers for parts 2 and 4 compare? If you got the same or different answers, why was that the case?

Problem 3: Implementing Naive Bayes

In this question you will implement a Naive Bayes classifier for a text classification problem. You will be given a collection of text articles, each coming from either the serious European magazine *The Economist*, or from the not-so-serious American magazine *The Onion*. The goal is to learn a classifier that can distinguish between articles from each magazine.

We have pre-processed the articles so that they are easier to use in your experiments. We extracted the set of all words that occur in any of the articles. This set is called the *vocabulary* and we let V be the number of words in the vocabulary. For each article, we produced a feature vector $X = \langle X_1, \dots, X_V \rangle$, where X_i is equal to 1 if the i^{th} word appears in the article and 0 otherwise.²

²Note the standard approach is to set the feature X_i equal to the number of times the i^{th} word appears in the article. Then the feature can take on any integer value from 0 to the maximum length of an article. This is the standard “bag of words” approach. However, the MLE computation and the programming become more involved, so in this homework, we simplify to boolean features 0,1.

Each article is also accompanied by a class label of either 1 for The Economist or 2 for The Onion. Later in the question we give instructions for loading this data into Octave.

When we apply the Naive Bayes classification algorithm, we make two assumptions about the data: first, we assume that our data is drawn iid from a joint probability distribution over the possible feature vectors X and the corresponding class labels Y ; second, we assume for each pair of features X_i and X_j with $i \neq j$ that X_i is conditionally independent of X_j given the class label Y (this is the Naive Bayes assumption). Under these assumptions, a natural classification rule is as follows: Given a new input X , predict the most probable class label \hat{Y} given X . Formally,

$$\hat{Y} = \underset{y}{\operatorname{argmax}} P(Y = y|X).$$

1. **[5 points]** Prove the classification rule can be rewritten as

$$\hat{Y} = \underset{y}{\operatorname{argmax}} \left(\prod_{w=1}^V P(X_w|Y = y) \right) P(Y = y).$$

2. **[5 points]** How many parameters are needed to represent the distribution $P(X|Y = y)$ when using the Naive Bayes assumption? How many are needed if we do not use the Naive Bayes assumption? Based on this difference, in which cases is there a big gain from making this assumption?

Of course, since we don't know the true joint distribution over feature vectors X and class labels Y , we need to estimate the probabilities $P(X|Y = y)$ and $P(Y = y)$ from the training data. For each word index $w \in \{1, \dots, V\}$ and class label $y \in \{1, 2\}$, the distribution of X_w given $Y = y$ is a Bernoulli distribution with parameter θ_{yw} . In other words, there is some unknown number θ_{yw} such that

$$P(X_w = 1|Y = y) = \theta_{yw} \quad \text{and} \quad P(X_w = 0|Y = y) = 1 - \theta_{yw}.$$

For both The Economist and The Onion, we believe that each word w has a non-zero chance of appearing, but it is more likely that w will not occur in any particular document. We incorporate this belief by computing a MAP estimate using a Beta(1.001, 1.9) prior on θ_{wy} . This has the added benefit of ensuring that none of our estimates of θ_{wy} are equal to 0 or 1 (which can cause problems for Naive Bayes).

Similarly, the distribution of Y (when we consider it alone) is a Bernoulli distribution (except taking values 1 and 2 instead of 0 and 1) with parameter ρ . In other words, there is some unknown number ρ such that

$$P(Y = 1) = \rho \quad \text{and} \quad P(Y = 2) = 1 - \rho.$$

In this case, since we have many examples of articles from both The Economist and The Onion, there is no risk of having zero-probability estimates, so we will instead use the MLE.

Programming Instructions

Question 3 through 5 of this question each ask you to implement one function related to the Naive Bayes classifier. You will submit your code online through the Autolab system, which will execute it remotely against a suite of tests. Your grade will be automatically determined from the testing results. Since you get immediate feedback after submitting your code and you are allowed to submit as many different versions as you like (without any penalty), it is easy for you to check your code as you go.

Download the template for Homework 2 and extract the contents (i.e., by executing `tar xvf hw2.tar` at the command line). In the archive you will find one `.m` file for each of the functions that you are asked to implement and a file that contains the data for this problem, `HW2Data.mat`. To finish each programming part of this problem, open the corresponding `.m` file and complete the function defined in that file. When you are ready to submit your solutions, you will create a new tar archive of the top-level directory (i.e., by executing `tar cvf hw2.tar hw2`) and upload that through the Autolab website.

The file `HW2Data.mat` contains the data that you will use in this problem. You can load it from Octave by executing `load("HW2Data.mat")` in the interpreter. After loading the data, you will see that there are 5 variables: `Vocabulary`, `XTrain`, `yTrain`, `XTest`, and `yTest`.

- `Vocabulary` is a $V \times 1$ dimensional cell array that contains every word appearing in the documents. When we refer to the j^{th} word, we mean `Vocabulary(j,1)`.
- `XTrain` is a $n \times V$ dimensional matrix describing the n documents used for training your Naive Bayes classifier. The entry `XTrain(i,j)` is 1 if word j appears in the i^{th} training document and 0 otherwise.
- `yTrain` is a $n \times 1$ dimensional matrix containing the class labels for the training documents. `yTrain(i,1)` is 1 if the i^{th} document belongs to The Economist and 2 if it belongs to The Onion.
- Finally, `XTest` and `yTest` are the same as `XTrain` and `yTrain`, except instead of having n rows, they have m rows. This is the data you will test your classifier on and it should not be used for training.

Logspace Arithmetic

When working with very large or very small numbers (such as probabilities), it is useful to work in *logspace* to avoid numerical precision issues. In logspace, we keep track of the logs of numbers, instead of the numbers themselves. For example, if $p(x)$ and $p(y)$ are probability values, instead of storing $p(x)$ and $p(y)$ and computing $p(x) * p(y)$, we work in log space by storing $\log(p(x))$, $\log(p(y))$, and we can compute the log of the product, $\log(p(x) * p(y))$ by taking the sum: $\log(p(x) * p(y)) = \log(p(x)) + \log(p(y))$.

We provide the function `logProd(x)` so you can use it in your implementation. This function takes as input a vector of numbers in logspace (i.e., $x_i = \log p_i$) and returns the product of those numbers in logspace—i.e., $\text{logProd}(\mathbf{x}) = \log(\prod_i p_i)$.

Training Naive Bayes

3. **[8 Points]** Complete the function `[D] = NB_XGivenY(XTrain, yTrain)`. The output `D` is a $2 \times V$ matrix, where for any word index $w \in \{1, \dots, V\}$ and class index $y \in \{1, 2\}$, the entry `D(y,w)` is the MAP estimate of $\theta_{yw} = P(X_w = 1 | Y = y)$ with a `Beta(1.001,1.9)` prior distribution.
4. **[8 Points]** Complete the function `[p] = NB_YPrior(yTrain)`. The output `p` is the MLE for $\rho = P(Y = 1)$.
5. **[8 Points]** Complete the function `[yHat] = NB_Classify(D, p, X)`. The input `X` is an $m \times V$ matrix containing m feature vectors (stored as its rows). Inputs `D` and `p` come from questions 3

and 4 respectively. The output `yHat` is a $m \times 1$ vector of predicted class labels, where `yHat(i)` is the predicted label for the i^{th} row of `X`. [Hint: In this function, you will want to use the `logProd` function to avoid numerical problems.]

Evaluating Naive Bayes

To help you evaluate your results, we provide the function `[error] = ClassificationError(yHat, yTruth)`, which takes two vectors of equal length and returns the proportion of entries that they disagree on.

Questions

6. **[5 Points]** Train your classifier on the data contained in `XTrain` and `yTrain` by running

```
D = NB_XGivenY(XTrain, yTrain);  
p = NB_YPrior(yTrain);
```

Use the learned classifier to predict the labels for the article feature vectors in `XTrain` and `XTest` by running

```
yHatTrain = NB_Classify(D, p, XTrain);  
yHatTest = NB_Classify(D, p, XTest);
```

Use the function `ClassificationError` to measure and report the training and testing error by running

```
trainError = ClassificationError(yHatTrain, yTrain);  
testError = ClassificationError(yHatTest, yTest);
```

How do the train and test errors compare? Which is more representative of the error we would expect to have on a new collection of articles? Does Naive Bayes attempt to minimize the training error?

7. **[8 Points]** In this question we explore how the size of the training data set affects the test and train error. For each value of m in $\{100, 130, 160, \dots, 580\}$, train your Naive Bayes classifier on the first m training examples (that is, use the data given by `XTrain(1:m,:)` and `yTrain(1:m)`). Plot the training and testing error for each such value of m . The x -axis of your plot should be m , the y -axis should be error, and there should be one curve for training error and one curve for testing error. Explain the general trend of both the training and testing error curves (Note: to see the trends better, you can plot the training and testing error individually).
8. **[8 Points]** Finally, we will try to interpret the learned parameters. Train your classifier on the data contained in `XTrain` and `yTrain`. For each class label $y \in \{1, 2\}$, create two lists according to the following criteria (Note that some of the words may look a little strange because we have run them through a stemming algorithm that tries to make words with common roots look the same. For example, “stemming” and “stemmed” would both become “stem”):

- Top five words that the model says are most likely to occur in a document from class y (i.e when $y = 1$ it refers to the Economist articles and when $y = 2$ it refers to Onion articles). That is, the top five words according to this metric:

$$P(X_w = 1|Y = y)$$

- Top five words w according to this metric:

$$\frac{P(X_w = 1|Y = y)}{P(X_w = 1|Y \neq y)}.$$

Which list of words is more informative about the class y ? Briefly explain your reasoning.

1 Submission Instructions

Below are the files you need to submit

1. logProd.m
2. NB_XGivenY.m
3. NB_YPrior.m
4. NB_Classify.m
5. HW2Data.mat

Please put these files, together with your writeup, into a folder called hw2, and run the following command

```
$ tar -cvf hw2.tar hw2
```

Then submit your tarfile to Autolab.