

10-701 Machine Learning, Spring 2011: Homework 2

Due: Friday Feb. 4 at 4pm in Sharon Cavlovich's office (GHC 8215)

Instructions There are 3 questions on this assignment. The last question involves coding. Please submit your writeup as 3 **separate** sets of pages according to TAs, with your name and userid on each set. If you choose to work with a partner for question 3, you (as a team) should submit one set of pages for that question, labeled with both names/userids. You should still submit the solutions to questions 1-2 separately.

1 Probability [Xi Chen, 30 points]

1. This problem studies the relationship between entropy, conditional entropy, mutual information, conditional independence, and expected values.

Consider random variables X and Y with joint probability density $p(X, Y)$. The expected value of any function $f(X, Y)$ of these variables is defined as $\mathbb{E}_p f(X, Y) = \int p(X, Y) f(X, Y) dx dy$ (i.e., it is the value of $f(X, Y)$ averaged over the different values X and Y can take on, weighted by their probabilities.) The expected value of a quantity is also sometimes called its "expectation".

The entropy, joint entropy and conditional entropy can be expressed as the following expectations:

- Entropy: $H(X) = -\mathbb{E}_p \ln p(X) = -\int p(X) \ln p(X) dx = -\int p(X, Y) \ln p(X) dx dy$
- Joint entropy: $H(X, Y) = -\mathbb{E}_p \ln p(X, Y) = -\int p(X, Y) \ln p(X, Y) dx dy$
- Conditional entropy: $H(X|Y) = -\mathbb{E}_p \ln p(X|Y) = -\int p(X, Y) \ln p(X|Y) dx dy$

- (a) In class we defined mutual information as

$$I(X, Y) \triangleq H(X) - H(X|Y).$$

Please use the linearity property of the expectation (i.e. $\mathbb{E}_p(X + Y) = \mathbb{E}_p(X) + \mathbb{E}_p(Y)$) to express $I(X, Y)$ as the expected value of a specific quantity. [2pt]

- (b) Use the linearity property of the expectation to prove the chain rule for entropy¹[3pt]:

$$H(X, Y) = H(Y) + H(X|Y).$$

Notice that given this chain rule for entropy, we can easily derive that $I(X, Y) = H(X) + H(Y) - H(X, Y)$.

- (c) Recall the *conditional mutual information* between random variables X and Y given Z is defined by:

$$I(X, Y|Z) \triangleq H(X|Z) - H(X|Y, Z).$$

Express $I(X, Y|Z)$ as the expected value of a specific quantity. Also, state a conditional independence assumption that will guarantee $I(X, Y|Z) = 0$ [5pt].

2. Given two random variables X and Y , let $\mathbb{E}X$ and $\mathbb{E}Y$ denote the means (i.e., expected values) of X and Y and let σ_X and σ_Y denote the standard deviations of X and Y . Here are three quantities that are commonly used to characterize the relationship between X and Y :

¹For your own interest, you may compare the chain rule for entropy to the chain rule in probability theory: $P(X, Y) = P(Y)P(X|Y)$

- Covariance: $\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}X)(Y - \mathbb{E}Y)] = \mathbb{E}(XY) - \mathbb{E}X\mathbb{E}Y$
- Correlation: $\rho_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$
- Mutual Information: $I(X, Y) = H(X) - H(X|Y) = KL(p(X, Y)||p(X)p(Y))$
where $KL(p||q) \equiv -\int p(x) \ln \frac{q(x)}{p(x)} dx$ is the Kullback-Leibler(KL) distance ².

(a) Recall the Cauchy-Schwarz inequality: for any two non-degenerate ³ random variables X and Y :

$$\{\mathbb{E}(XY)\}^2 \leq \mathbb{E}(X^2)\mathbb{E}(Y^2),$$

where the equality holds if and only if $P(X = aY) = 1$ for some non-zero constant a .

Use the Cauchy-Schwarz inequality to prove that for any two non-degenerate random variables X and Y , the absolute value of the correlation between X and Y is less than or equal to 1, i.e. $|\rho_{XY}| \leq 1$. [3pt]

- (b) Show conditions under which we have $\rho_{XY} = 1$ and the conditions under which $\rho_{XY} = -1$. [2pt]
- (c) If $I(X, Y) = 0$, can we conclude that $\rho_{XY} = 0$ (or equivalently, $\text{cov}(X, Y) = 0$)? If so, please give the proof; if not, please find the two random variables X and Y such that $I(X, Y) = 0$ but $\rho_{XY} \neq 0$. [5pt]
- (d) If $\rho_{XY} = 0$, can we conclude that $I(X, Y) = 0$? If so, please give the proof; if not, please find the two random variables X and Y such that $\rho_{XY} = 0$ but $I(X, Y) \neq 0$. [10pt]

2 Generative and Discriminative Classifiers: Gaussian (Naive) Bayes and Logistic Regression [Yi Zhang, 30 points]

Recall that a generative classifier estimates $P(\mathbf{X}, Y) = P(Y)P(\mathbf{X}|Y)$, while a discriminative classifier directly estimates $P(Y|\mathbf{X})$ ⁴. For clarity, we highlight \mathbf{X} in bold to emphasize that it usually represents a vector of multiple attributes, i.e., $\mathbf{X} = \langle X_1, X_2, \dots, X_n \rangle$. However, this question does **not** require students to derive the answer in vector/matrix notation.

In class we have observed an interesting relationship between a discriminative classifier (logistic regression) and a generative classifier (Gaussian naive Bayes): the form of $P(Y|\mathbf{X})$ derived from the assumptions of **a specific class** of Gaussian naive Bayes classifiers is precisely the form used by logistic regression. The derivation can be found in the required reading, **Mitchell: Naive Bayes and Logistic Regression, Section 3.1 (page 8 - 10)**. In that reading, we made the following assumptions for Gaussian naive Bayes classifiers to model $P(\mathbf{X}, Y) = P(Y)P(\mathbf{X}|Y)$ (rephrased from the required reading, with a few comments):

1. Y is a boolean variable following a Bernoulli distribution, with parameter $\pi = P(Y = 1)$ and thus $P(Y = 0) = 1 - \pi$.
2. $\mathbf{X} = \langle X_1, X_2, \dots, X_n \rangle$, where each attribute X_i is a continuous random variable. For each X_i , $P(X_i|Y = k)$ is a Gaussian distribution $N(\mu_{ik}, \sigma_i)$. Note that σ_i is the standard deviation of the Gaussian distribution (and thus σ_i^2 is the variance), which does **not** depend on k .
3. For all $i \neq j$, X_i and X_j are conditionally independent given Y . This is why this type of classifier is called “naive”.

We say this is **a specific class** of Gaussian naive Bayes classifiers because we have made an assumption that the standard deviation σ_i of $P(X_i|Y = k)$ does not depend on the value k of Y . This is **not** a general assumption for Gaussian naive Bayes classifiers.

²As shown in homework 1, the definitions of mutual information based on KL divergence and entropy are equivalent.

³A random variable that can only take on one value (i.e., a constant) is called a degenerate random variable.

⁴Note that certain discriminative classifiers are non-probabilistic: they directly estimate a function $f: \mathbf{X} \rightarrow Y$ instead of $P(Y|\mathbf{X})$. We will see such classifiers later this semester, but it is beyond the scope of this question.

2.1 General Gaussian naive Bayes Classifiers and Logistic Regression [15 points]

Let's make our Gaussian naive Bayes classifiers a little more general by removing the assumption that the standard deviation σ_i of $P(X_i|Y = k)$ does not depend on k . As a result, for each X_i , $P(X_i|Y = k)$ is a Gaussian distribution $N(\mu_{ik}, \sigma_{ik})$, where $i = 1, 2, \dots, n$ and $k = 0, 1$. Note that now the standard deviation σ_{ik} of $P(X_i|Y = k)$ depends on both the attribute index i and the value k of Y .

Question: is the new form of $P(Y|\mathbf{X})$ implied by this more general Gaussian naive Bayes classifier still the form used by logistic regression? Derive the new form of $P(Y|\mathbf{X})$ to prove your answer.

2.2 Gaussian Bayes Classifiers and Logistic Regression [15 points]

Students in 10-701 are all smart, so clearly we will not be satisfied by only studying a “naive” classifier. In this part, we will turn our attention to a specific class of Gaussian Bayes classifiers (without “naive”, yeah!). We consider the following assumptions for our Gaussian Bayes classifiers:

1. Y is a boolean variable following a Bernoulli distribution, with parameter $\pi = P(Y = 1)$ and thus $P(Y = 0) = 1 - \pi$.
2. $\mathbf{X} = \langle X_1, X_2 \rangle$, i.e., we only consider **two** attributes, where each attribute X_i is a continuous random variable. X_1 and X_2 are **not** conditionally independent given Y . We assume $P(X_1, X_2|Y = k)$ is a **bivariate Gaussian distribution** $N(\mu_{1k}, \mu_{2k}, \sigma_1, \sigma_2, \rho)$, where μ_{1k} and μ_{2k} are means of X_1 and X_2 , σ_1 and σ_2 are standard deviations of X_1 and X_2 , and ρ is the **correlation** between X_1 and X_2 . Note that μ_{1k} and μ_{2k} depend on the value k of Y , but σ_1 , σ_2 , and ρ do **not** depend on Y . Also recall that the density of a bivariate Gaussian distribution, given $(\mu_{1k}, \mu_{2k}, \sigma_1, \sigma_2, \rho)$, is:

$$P(X_1, X_2|Y = k) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{\sigma_2^2(X_1 - \mu_{1k})^2 + \sigma_1^2(X_2 - \mu_{2k})^2 - 2\rho\sigma_1\sigma_2(X_1 - \mu_{1k})(X_2 - \mu_{2k})}{2(1-\rho^2)\sigma_1^2\sigma_2^2}\right]$$

Question: is the form of $P(Y|\mathbf{X})$ implied by such not-so-naive Gaussian Bayes classifiers still the form used by logistic regression? Derive the form of $P(Y|\mathbf{X})$ to prove your answer.

SPECIAL NOTES: we intentionally choose only two attributes $\mathbf{X} = \langle X_1, X_2 \rangle$ because the density of *bivariate* Gaussian distribution can be expressed using scalars instead of vectors and matrices. This way, students do not need to derive the answer using vector/matrix algebra. However, if you are comfortable with vector and matrix algebra, please feel free to use following **alternative assumptions** and derive your answer in vector/matrix notation (which may turn out to be less time-consuming):

1. Y is a boolean variable following a Bernoulli distribution, with parameter $\pi = P(Y = 1)$ and thus $P(Y = 0) = 1 - \pi$.
2. $\mathbf{X} = \langle X_1, X_2, \dots, X_n \rangle$ are **not** conditionally independent given Y , and $P(\mathbf{X}|Y = k)$ follows a **multivariate normal distribution** $N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$. Note that $\boldsymbol{\mu}_k$ is the $n \times 1$ mean vector depending on the value of Y , and $\boldsymbol{\Sigma}$ is the $n \times n$ **covariance** matrix, which does **not** depend on Y . Also, you should be familiar with the density of multivariate normal distribution in vector/matrix notation.

You may choose to derive your answer for **either** the bivariate normal version using scalars **or** the multivariate normal version using vector/matrix notation. Derive both versions will **not** gain extra credits.

3 Naive Bayes Document Classifier [Carl Doersch, 40 points]

In this question, you will implement the Naive Bayes document classifier and apply it to the classic 20 newsgroups dataset⁵. In this dataset, each document is a posting that was made to one of 20 different usenet newsgroups. Our goal is to write a program which can predict which newsgroup a given document was posted to.

⁵<http://people.csail.mit.edu/jrennie/20Newsgroups/>

For this question, you may write your code and solution in teams of at most 2. If you decide to do this, you should submit one copy of your solutions to question 3 (both code and answers to questions) per team. This copy should be clearly marked with the names of both team members.

3.1 Model

Say we have a document D containing n words; call the words $\{X_1, \dots, X_n\}$. The value of random variable X_i is the word found in position i in the document. We wish to predict the label Y of the document, which can be one of m categories. We could use the model:

$$P(Y|X_1\dots X_n) \propto P(X_1\dots X_n|Y)P(Y) = P(Y) \prod_i P(X_i|Y)$$

That is, each X_i is sampled from some distribution that depends on its position X_i and the document category Y . As usual with discrete data, we assume that $P(X_i|Y)$ is a multinomial distribution over some vocabulary V ; that is, each X_i can take one of $|V|$ possible values corresponding to the words in the vocabulary. Therefore, in this model, we are assuming (roughly) that for any pair of document positions i and j , $P(X_i|Y)$ may be completely different from $P(X_j|Y)$.

Question 3.1: In your answer sheet, explain in a sentence or two why it would be difficult to accurately estimate the parameters of this model on a reasonable set of documents (e.g. 1000 documents, each 1000 words long, where each word comes from a 50,000 word vocabulary). **[3 points]**

To improve the model, we will make the additional assumption that:

$$\forall i, j \quad P(X_i|Y) = p(X_j|Y)$$

Thus, in addition to estimating $P(Y)$, you must estimate the parameters for the single distribution $P(X|Y)$, which we define to be equal to $P(X_i|Y)$ for all X_i . Each word in a document is assumed to be an *iid* draw from this distribution.

3.2 Data

The data file (available on the website) contains six files:

1. **vocabulary.txt** is a list of the words that may appear in documents. The line number is word's id in other files. That is, the first word ('archive') has wordId 1, the second ('name') has wordId 2, etc.
2. **newsgrouplabels.txt** is a list of newsgroups from which a document may have come. Again, the line number corresponds to the label's id, which is used in the .label files. The first line ('alt.atheism') has id 1, etc.
3. **train.label** Each line corresponds to the label for one document from the training set. Again, the document's id (docId) is the line number.
4. **test.label** The same as train.label, except that the labels are for the test documents.
5. **train.data** Specifies the counts for each of the words used in each of the documents. Each line is of the form "docId wordId count", where count specifies the number of times the word with id wordId in the training document with id docId. All word/document pairs that do not appear in the file have count 0.
6. **test.data** Same as train.data, except that it specified counts for test documents.

If you are using matlab, the functions `textread` and `sparse` will be useful in reading these files.

3.3 Implementation

Your first task is to implement the Naive Bayes classifier specified above. You should estimate $P(Y)$ using the MLE, and estimate $P(X|Y)$ using a MAP estimate with the prior distribution $Dirichlet(1 + \alpha, \dots, 1 + \alpha)$, where $\alpha = 1/|V|$ and V is vocabulary.

Question 3.2: In your answer sheet, report your overall testing accuracy (Number of correctly classified documents in the test set over the total number of test documents), and print out the confusion matrix (the matrix C , where c_{ij} is the number of times a document with ground truth category j was classified as category i). [7 points]

Question 3.3: Are there any newsgroups that the algorithm confuses more often than others? Why do you think this is? [2 points]

3.4 Priors and Overfitting

In your initial implementation, you used a prior $Dirichlet(1 + \alpha, \dots, 1 + \alpha)$ to estimate $P(X|Y)$, and I told you set $\alpha = 1/|V|$. Hopefully you wondered where this value came from. In practice, the choice of prior is a difficult question in Bayesian learning: either we must use domain knowledge, or we must look at the performance of different values on some validation set. Here we will use the performance on the testing set to gauge the effect of α ⁶.

Question 3.4: Re-train your Naive Bayes classifier for values of α between .00001 and 1 and report the accuracy over the test set for each value of α . Create a plot with values of α on the x -axis and accuracy on the y -axis. Use a logarithmic scale for the x -axis (in Matlab, the `semilogx` command). Explain in a few sentences why accuracy drops for both small and large values of α [5 points]

3.5 Identifying Important Features

One useful property of Naive Bayes is that its simplicity makes it easy to understand why the classifier behaves the way it does. This can be useful both while debugging your algorithm and for understanding your dataset in general. For example, it is possible to identify which words are strong indicators of the category labels we're interested in.

Question 3.5: Propose a method for ranking the words in the dataset based on how much the classifier 'relies on' them when performing its classification (hint: information theory will help). Your metric should use only the classifier's estimates of $P(Y)$ and $P(X|Y)$. It should give high scores to those words that appear frequently in one or a few of the newsgroups but not in other ones. Words that are used frequently in general English ('the', 'of', etc.) should have lower scores, as well as words that only appear extremely rarely throughout the whole dataset. Finally, your method this should be an overall ranking for the words, not a per-category ranking. [3 points]

Question 3.6: Implement your method, set α back to $1/|V|$, and print out the 100 words with the highest measure. [2 points]

Question 3.7: If the points in the training dataset were not sampled independently at random from the same distribution of data we plan to classify in the future, we might call that training set *biased*. Dataset bias is a problem because the performance of a classifier on a biased dataset will not accurately reflect its future performance in the real world. Look again at the words your classifier is 'relying on'. Do you see any signs of dataset bias? [3 points]

3.6 Hand-in for Question 3

Print out the code that you used to solve problem 3. This should include the code for the Naive Bayes classifier, the code for modifying α , and the code for identifying important features. Submit this code in class on the due date, along with your answers to questions 3.1-3.7. [15 points]

⁶It is tempting to choose α to be the one with the best performance on the testing set. However, if we do this, then we can no longer assume that the classifier's performance on the test set is an unbiased estimate of the classifier's performance in general. The act of choosing α based on the test set is equivalent to training on the test set; like any training procedure, this choice is subject to overfitting.