

CMPUT 466/551—Machine Learning

Assignment 4

Winter 2004

Department of Computing Science

University of Alberta

Due: in class, *Tuesday, March 23*

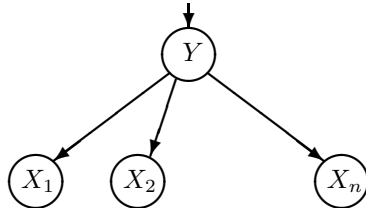
Worth: 15% of final grade

Instructor: Dale Schuurmans, Ath409, x2-4806, dale@cs.ualberta.ca

Note One question requires you to write a few small Matlab programs which are to be submitted by email. When finished, please send a *single* tar file containing all of your .m files **to the TA, Alexander Kovarsky, at kovarsky@cs.ualberta.ca** with a subject heading “CMPUT 466/551 A4 solutions”. The write-up should be printed out in hard copy and handed in at the start of class on March 23.

Question 1 (Probability models)

In class we discussed the “Naive Bayes” (NB) probability model, given by the generator



which specifies that the probability of a joint assignment to the variables is given by

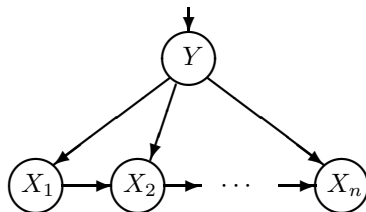
$$P(y, x_1, x_2, \dots, x_n) = P(Y = y) \prod_{j=1}^n P(X_j = x_j | Y = y).$$

Recall that for boolean attributes $X_j \in \{0, 1\}$, $Y \in \{0, 1\}$ this model is determined by $1 + 2n$ free parameters: $\theta = P(Y = 1)$, $\theta_{j1} = P(X_j = 1 | Y = 1)$, and $\theta_{j0} = P(X_j = 1 | Y = 0)$. Also remember that the minimum-expected-error classification rule for this model

$$\hat{y} = 1 \Leftrightarrow P(Y = 1 | x_1, \dots, x_n) \geq P(Y = 0 | x_1, \dots, x_n)$$

is equivalent to a linear discriminant $\hat{y} = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$ that uses weights $w_j = \log \frac{\theta_{j1} (1 - \theta_{j0})}{\theta_{j0} (1 - \theta_{j1})}$ and bias $b = \log \frac{1 - \theta}{\theta} + \sum_{j=1}^n \log \frac{1 - \theta_{j0}}{1 - \theta_{j1}}$. In this question we will explore a simple generalization of this model that has some interesting properties.

A “Chain Augmented Naive Bayes” model (CAN) is a Naive Bayes-like model where we include an additional chain of dependencies that runs through the variables as follows



This generative model specifies that we compute the probability of a joint assignment to the variables according to

$$P(y, x_1, x_2, \dots, x_n) = P(Y = y) P(X_1 = x_1 | Y = y) \prod_{j=2}^n P(X_j = x_j | Y = y, X_{j-1} = x_{j-1}).$$

Assume boolean attributes $X_j \in \{0, 1\}$, $Y \in \{0, 1\}$.

- (a) (1%) Specify the parameters of the CAN model in a manner similar to the NB model above. How many free parameters are there (as a function of n)? Also, given training data

$$\begin{bmatrix} x_{11} & \cdots & x_{1n} & y_1 \\ \vdots & & \vdots & \vdots \\ x_{t1} & \cdots & x_{tn} & y_t \end{bmatrix}$$

what are the maximum likelihood estimates for the CAN parameters?

- (b) (1%) Derive a simple form for the minimum-expected-error classification rule for this model in a similar manner to the NB model above. What form does this decision rule have? (Hint: it is not a linear discriminant, but it is close.)

For the remainder of this question we will investigate the consequences of learning probability models when incorrect assumptions are made about the structure of the distribution. The lesson is that even principled learning under bad modeling assumptions can lead to suboptimal results.

Assume that data is generated according to a CAN model where

$$\begin{aligned} P(Y = 1) &= \frac{2}{3}, \\ P(X_1 = 1|Y = 1) &= P(X_1 = 1|Y = 0) = \frac{1}{2}, \\ \text{and } X_j &= X_{j-1} \oplus Y \text{ for } j = 2, \dots, n, n \geq 2 \end{aligned}$$

where \oplus denotes “exclusive or”. (That is, X_j is a deterministic function of X_{j-1} and Y and therefore $P(X_j = x_{j-1} \oplus y | X_{j-1} = x_{j-1}, Y = y) = 1$.) Assume also that we are (incorrectly) attempting to learn from this data source by using a Naive Bayes model.

- (c) (1%) What are the parameters of the optimal NB model for this CAN data source? (Note: we haven’t formally discussed optimal approximations for probability distributions in this course, but for this question just assume it means computing the corresponding conditional probabilities under the true model.) What is the minimum-expected-error classification rule for the optimal NB model?
- (d) (1%) What is the misclassification probability of this NB classifier when test examples are drawn from the CAN model? What is the misclassification probability of the optimal CAN classifier? (Note: just compute the probability—you do not have to show a derivation of the optimal CAN classifier.)
- (e) (1%) For this part, assume that $n = 2$. (That is, we are only considering the variables Y , X_1 and X_2 .) Write down the coefficients w_1 , w_2 and b of a linear discriminant that achieves better misclassification probability than the solution of Part (c). What is the misclassification probability of your alternative linear discriminant?

Question 2 (EM—application to Naive Bayes)

In this question you will derive an EM training algorithm for the Naive Bayes probability model introduced in Question 1. The outcome will be a learning method that can produce a predictor $h : X^n \rightarrow Y$ based on exploiting both labeled $\langle \mathbf{x}_i, y_i \rangle$ and unlabeled $\langle \mathbf{x}_i, - \rangle$ training examples. Recall that the Naive Bayes model assumes there is a single random variable $Y \in \{0, 1\}$ whose value determines the conditional distribution over the remaining variables X_1, \dots, X_n ($X_j \in \{0, 1\}^n$), where, in particular, the conditional distributions of X_1, \dots, X_n are independent given $Y = y$. This model is defined by $2n + 1$ parameters $\boldsymbol{\theta} = \langle \theta, \{\theta_{j0}\}_{j=1}^n, \{\theta_{j1}\}_{j=1}^n \rangle$.

Consider a scenario where all of the variables X_j are always observed but the Y labels are sometimes missing. That is, consider training data of the form

$$\begin{bmatrix} x_{11} & \dots & x_{1n} & z_1 \\ \vdots & & \vdots & \vdots \\ x_{t1} & \dots & x_{tn} & z_t \end{bmatrix}$$

where the z_i labels are chosen from the set $\{0, 1, -1\}$ such that $z_i = -1$ indicates that the original y_i value was *unobserved*. Recall that EM works as follows:

EM

- Start with an initial parameter vector $\boldsymbol{\theta}^{(0)}$.
- Repeatedly update the parameter vector $\boldsymbol{\theta}^{(k)} \rightarrow \boldsymbol{\theta}^{(k+1)}$ in two sub-stages:
 - **E step:** For each training example $x_{i1}, \dots, x_{in}, z_i$, compute the conditional probability of each possible completion $y_i = 0$ or $y_i = 1$ given the observed x -values x_{i1}, \dots, x_{in} :

$$Q_i^{(k)}(y) = \begin{cases} 0 & \text{if } z_i \in \{0, 1\} \text{ and } y \neq z_i \\ 1 & \text{if } z_i \in \{0, 1\} \text{ and } y = z_i \\ P(Y = y | X_1 = x_{i1}, \dots, X_n = x_{in}, \boldsymbol{\theta}^{(k)}) & \text{if } z_i = -1 \end{cases} \quad (1)$$

- **M step:** Compute:

$$\boldsymbol{\theta}^{(k+1)} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^t \sum_{y=0}^1 Q_i^{(k)}(y) \log P(X_1 = x_{i1}, \dots, X_n = x_{in}, Y = y | \boldsymbol{\theta}) \quad (2)$$

In this question you will derive a simplified form of the EM update for the Naive Bayes model, which will allow you to implement it efficiently in Question 3 below.

- (a) (1%) Express the summation in (2) as a function of $\theta, \theta_{j0}, \theta_{j1}, x_{ij}, y$ and q_{yi} , where $q_{yi} = Q_i^{(k)}(y)$. (Hint: you can express $P(Y = y)$ as $\theta^y(1 - \theta)^{1-y}$.)
- (b) (1%) Compute the partial derivatives of the summation in (2) with respect to $\theta, \theta_{j0}, \theta_{j1}$.
- (c) (1%) Solve for the values of $\theta, \theta_{j0}, \theta_{j1}$ that maximize (2). (Hint: The final solutions will have a reasonably simple form. To simplify the expressions you will need to make use of the identity $\sum_{y=0}^1 q_{yi} = \sum_{y=0}^1 Q_i^{(k)}(y) = 1$.)
- (d) (1%) Express $Q_i^{(k)}(y)$ as a function of $\theta, \theta_{j0}, \theta_{j1}, x_{ij}$ and y .
- (e) (1%) Combine the E and the M steps into one step by writing the updated parameter values $\theta^{(k+1)}, \theta_{j0}^{(k+1)}, \theta_{j1}^{(k+1)}$ as a function of the previous parameter values $\theta^{(k)}, \theta_{j0}^{(k)}, \theta_{j1}^{(k)}$ and x_{ij} and y .

Question 3 (Maximum likelihood learning—Naive Bayes)

In this exercise you will write simple Matlab functions for learning maximum likelihood Naive Bayes models from training data, and classifying text vectors with these models. In particular, you will implement the EM algorithm derived in Question 2. For this question you will consider training and test data generated according to the following generative model.

Data generation:

```
n = 5                % dimension
t = 100             % training size
p = 0.5            % prob y=1
p0 = 0.33*ones(1,n) % prob x_j=1 given y=0
p1 = 0.67*ones(1,n) % prob x_j=1 given y=1
pm = 0.8           % prob y missing
y = rand(t,1) > p  % true labels
pos = find(y == 1)
neg = find(y == 0)
X = rand(t,n)
X(neg,:) = X(neg,:) <= repmat(p0,length(neg),1) % positive examples
X(pos,:) = X(pos,:) <= repmat(p1,length(pos),1) % negative examples
b = rand(t,1) <= pm
z = (~b&y)-b       % blocked labels (-1 means unobserved)
```

- (a) (1%) Write a Matlab function `classify(Xtest,p,p0,p1)` which takes a $te \times n$ matrix of test vectors `Xtest` and the $2n + 1$ parameters of a Naive Bayes model—one scalar parameter `p` and two $1 \times n$ vectors of parameters `p0` and `p1`—and returns a $te \times 1$ vector of classifications `ytest` on the test vectors. Classifications are to be determined according to the optimal decision rule

$$\hat{y} = 1 \Leftrightarrow P(Y = 1|x_1, \dots, x_n) \geq P(Y = 0|x_1, \dots, x_n)$$

where `p` specifies θ in the Naive Bayes model, `p0` specifies $\{\theta_{j0}\}_{j=1}^n$ and `p1` specifies $\{\theta_{j1}\}_{j=1}^n$ (see Question 1). Your function must handle `Xtest`, `p`, `p0` and `p1` with arbitrary dimensions te and n .

- (b) (1%) Write a Matlab function `[p,p0,p1] = maxlike(X,z)` which takes a $t \times n$ matrix of training vectors `X` and a $t \times 1$ vector of partially blocked training labels `z`, and returns the parameters of a Naive Bayes model that maximizes the likelihood *only* of the training vectors with *observed* labels (i.e. ignoring any training vectors that have a label $z_i = -1$). The parameters returned are one scalar `p` and two $1 \times n$ vectors `p0` and `p1`. (Hint: You do not use EM for this part—only simple plug-in formulas.) Your function must handle `X` and `z` with arbitrary dimensions t and n .
- (c) (2%) Write a Matlab function `[p,p0,p1] = em(X,z,tol)` which takes a $t \times n$ matrix of training vectors `X`, a $t \times 1$ vector of partially blocked training labels `z` and scalar tolerance `tol`, and returns the parameters of a Naive Bayes model that maximizes

the likelihood of *all* the training vectors; i.e. including those vectors with *unobserved* labels ($z_i = -1$). This function is to be implemented using the EM algorithm derived in Question 2. EM should be run until the maximum change in any parameter is less than `tol`. (You can start EM with a random parameter vector, although it might be interesting to start it with the parameter vector returned by `maxlike`.) Your function must handle \mathbf{X} and \mathbf{z} with arbitrary dimensions t and n .

(d) (1%) Compare the quality of these two maximum likelihood training algorithms by testing the accuracy of the classifiers they produce as follows.

A: Generate a random training set \mathbf{X}, \mathbf{z} using the generative model outlined above, and solve for optimal Naive Bayes parameters using the two procedures:

$$\begin{aligned} [\mathbf{q}, \mathbf{q}_0, \mathbf{q}_1] &= \text{maxlike}(\mathbf{X}, \mathbf{z}) \\ [\mathbf{r}, \mathbf{r}_0, \mathbf{r}_1] &= \text{em}(\mathbf{X}, \mathbf{z}, 0.01) \end{aligned}$$

B: Use the `classify` function to label the training data using each learned model $[\mathbf{q}, \mathbf{q}_0, \mathbf{q}_1]$ and $[\mathbf{r}, \mathbf{r}_0, \mathbf{r}_1]$. Report the *mean* misclassification error that each of these models obtain on the training data *with respect to the original (hidden) y -labels*.

C: Generate `te = 1000` test examples from the same generative model and report the mean misclassification error of both models on *test* data.

Repeat parts A, B and C 100 times and accumulate the mean misclassification errors for both models in two matrices: one for training error and the other for test error. Report the *averages* for each kind of mean error at each degree in two tables (one training error and the other testing error).