# 10-701/15-781 Machine Learning, Fall 2005

#### Assignment 2 SOLUTIONS

Out: 9/27/05 Due: beginning of class 10/06/05

If you have questions, please contact Mike Stilman <robot+ta@cmu.edu>.

## Linear Regression

1. (Noise in Linear Regression)[25 pts] Linear regression is applied when we assume that there is an underlying linear function  $\mathbf{f}(x) = \mathbf{w}x + E$  that is generating data. Here, E is a random variable representing noise.

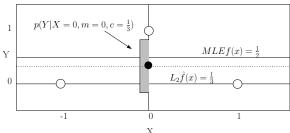
In class you saw that selecting  $\mathbf{w}$  to minimize the sum of squares of residual error yields the Maximum Likelihood  $\mathbf{w}$ . To prove this, we assumed that E is independent, normally distributed and of constant variance.

Your classmate has some trouble believing that noise is normally distributed. He claims that the noise distribution should be bounded by some region [-a, a] where p(E) is 0 outside these bounds and strictly greater than 0 within the bounds.

(a) Lets look at a simple bounded distribution of noise:

$$p(E) = \begin{cases} 1 & -\frac{1}{2} \le E \le \frac{1}{2} \\ 0 & \end{cases}$$

In this case, least squares  $(L_2)$  regression will not necessarily maximize the likelihood of the data. Demonstrate an example where  $L_2$  regression finds a  $\mathbf{w}$  that is not MLE. (i.e. Present a univariate linear function and data points that do not violate the uniform error distribution. Then show that  $L_2$  regression results in a suboptimal model.)



Suppose your underlying function is  $\mathbf{f}(x) = \frac{1}{2} + E$ . We place three data points as illustrated at (-1,0),(0,1),(1,0). These points are valid for the error distribution. Suppose our  $L_2$  estimate is of the form  $\mathbf{f}'(x) = mx + c$ . We wish to minimize the sum of squared residual error:

Residual Error 
$$\mathbf{e}(\mathbf{m}, \mathbf{c}) = (-1 \times m + c - 0)^2 + (0 \times m + c - 1)^2 + (1 \times m + c - 0)^2$$
  

$$= m^2 - 2mc + c^2 + c^2 - 2c + 1 + m^2 + 2mc + c^2$$

$$= 2m^2 + 3c^2 - 2c + 1$$

1

Setting the partial derivatives to 0 trivially yields:

$$\frac{\delta \mathbf{e}}{\delta m} = 0$$

$$\frac{\delta \mathbf{e}}{\delta c} = 0$$

$$4m = 0$$

$$6c - 2 = 0$$

$$c = \frac{1}{3}$$

Now, our  $L_2$  estimate  $\mathbf{f}'(x_1) = \mathbf{f}'(0) = \frac{1}{3}$ . Yet the data point  $y_1 = 1$ . Hence the probability of this data point is 0 given the  $L_2$  parameters. The likelihood of our three data points is:

$$\prod_{i=1}^{n} p(y_i|w, x_i) = 1 \times 0 \times 1 = 0$$

Notice that if we were to select  $b = \frac{1}{2}$  then we would get a probability of 1 for all the data points, yielding a cumulative likelihood of 1. That is the MLE estimate. Clearly, in this case  $L_2$  does not find a parameter that maximizes the likelihood of the data.

(b) Your classmate accepts that uniform distributions are not necessarily a good model for noise. The bigger question, however, remains unresolved: Is there a bounded distribution for noise such that least squares regression always finds the MLE  $\mathbf{w}$ ? Either define such a distribution and prove that  $L_2$  regression must find the MLE  $\mathbf{w}$  (for any valid linear function/data points), or prove that no such distribution exists. (Visual arguments are acceptable, but they must be supported by clear and correct logical reasoning.)

Your classmate will be disappointed: there is NO error distribution that follows his rules and guarantees that least squares will find the MLE solution.

For any distribution that is positive iff it is in the range [-a,a]: Let h be the width of the distribution, h = 2a. Let your underlying function be horizontal with m = 0 and c = a. Thus, for any x, the error distribution always lies between 0 and h.

$$y(x) = 0 \times x + \frac{h}{2}$$

• Select three data points (-1,0),(0,h),(1,0).

Observe that there exists a linear function  $(y(x) = 0 \times x + \frac{h}{2})$  such that  $p(y_i|w, x_i) > 0$  for all three points. Therefore, the cumulative likelihood of the data given this function is greater than 0.

• The function that maximizes the likelihood of the data must also yield a likelihood that is greater than 0.

The least squares solution is found in the same way as (a):

$$L_2 y'(x) = 0 \times x + \frac{h}{3}$$

Clearly  $\frac{h}{3} + \frac{h}{2} < h$  and hence the data point at x = 0 is outside the noise bounds for the  $L_2$  solution (e.g.  $p(Y = h|L_2, X = 0) = 0$ ).

• Regardless of  $p(y_i|L_2, x_i)$  for the remaining two data points, the likelihood of the data (product of probabilities) for the  $L_2$  solution will be 0.

Thus, for any distribution that is positive iff it is in the range [-a,a], we see an example where  $L_2$  does not find the MLE solution.

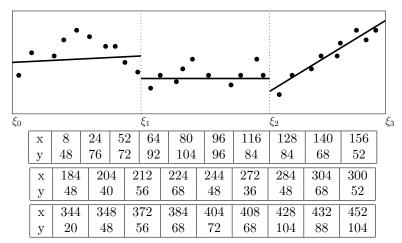
2. (Interpolating with Regression) [20 pts] One fun application of regression is interpolation. Basic linear regression fits a line to points. By introducing basis functions we were able to model polynomials.

Some functions, however, are quite complex and could require very high order polynomials to achieve accurate modeling. Often it is better to split these functions into segments at some number of points that we call knots.

In the following exercises you may use these functions to simplify your notation:

$$I(a,b,x) = \begin{cases} 1 & a \le x < b \\ 0 & otherwise \end{cases} \qquad G(x) = \begin{cases} 1 & 0 < x \\ 0 & otherwise \end{cases}$$

(a) We have split the following data-set into three evenly spaced regions at knots  $\xi_1$  and  $\xi_2$  (where  $\xi_{i+1} - \xi_i = 160$ ). The data set represents some function  $y = \mathbf{f}(x) + E$ . Suppose we applied linear regression to each segment separately. Draw a rough sketch of the resulting model function.



(b) Your model should be discontinuous. It is possible to achieve the same model with a single linear regression of the form:

$$y' = \sum_{i=1}^{m} \beta_i h_i(x)$$

Find a set of six (6) basis functions,  $h_i(x)$ , such that their linear combination can represent any piecewise linear function with discontinuities at the knots  $\xi$ . Find the parameters  $\beta_i$  that satisfy the maximum likelihood model for the given data. Bases:

$$h_1(x) = I(\xi_0, \xi_1, x), \quad h_2(x) = I(\xi_1, \xi_2, x), \quad h_3(x) = I(\xi_2, \xi_3, x),$$
  
 $h_4(x) = I(\xi_0, \xi_1, x)x, \quad h_5(x) = I(\xi_1, \xi_2, x)x, \quad h_6(x) = I(\xi_2, \xi_3, x)x$ 

Parameters:

$$(\beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6) = (75.8177, 41.4472, -196.9013, 0.0206, 0.0408, 0.6719)$$

(c) Generally, we do want our model y' to be continuous. Find a set of four (4) basis functions,  $h_i(x)$ , such that  $y' = \sum_{i=1}^m \beta_i h_i(x)$  can represent any continuous piecewise linear function that has discontinuous derivatives only at the knots. Estimate  $\beta_i$  and formulate a new y' that minimizes the squared residual error. Bases:

$$h_1(x) = 1$$
,  $h_2(x) = x$ ,  $h_3(x) = G(x - \xi_1)(x - \xi_1)$ ,  $h_4(x) = G(x - \xi_2)(x - \xi_2)$ 

Parameters:

$$(\beta_1, \beta_2, \beta_3, \beta_4) = (80.3877, -0.0606, -0.1613, 0.7108)$$

Other solutions were possible. Credit was deducted for not finding  $\beta$ s or constructing incomplete sets of bases. In particular, discontinuous bases can lead to discontinuous functions.

#### Theoretical Analysis of Logistic Regression and Naive Bayes

[20 pts] In class and in Tom's draft chapter handout we showed that when Y is Boolean and  $X = (X_1 ... X_n)$  is a vector of continuous variables, then the assumptions of the Gaussian Naive Bayes classifier imply that P(Y|X) is given by the logistic function with appropriate parameters W. In particular:

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^{n} w_i X_i)}$$

and

$$P(Y = 0|X) = \frac{\exp(w_0 + \sum_{i=1}^n w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}$$

Consider instead the case where Y is Boolean and  $X = (X_1 ... X_n)$  is a vector of Boolean variables. Prove for this case also that P(Y|X) follows this same form (and hence that logistic regression is also the discriminative counterpart to a Naive Bayes generative classifier over Boolean features).

Hints:

- 1. Simple notation will help. Since the  $X_i$  are Boolean variables, you need only one parameter to define  $P(X_i|Y=y_k)$ . Define  $\theta_{i1} \equiv P(X_i=1|Y=1)$ , in which case  $P(X_i=0|Y=1) = (1-\theta_{i1})$ . Similarly, use  $\theta_{i0}$  to denote  $P(X_i=1|Y=0)$ .
- 2. Notice with the above notation you can represent  $P(X_i|Y=1)$  as follows

$$P(X_i|Y=1) = \theta_{i1}^{X_i} (1-\theta_{i1})^{(1-X_i)}$$

Note when  $X_i = 1$  the second term is equal to 1 because its exponent is zero. Similarly, when  $X_i = 0$  the first term is equal to 1 because its exponent is zero.

$$P(Y=1|X) = \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)}$$

$$P(Y = 1|X) = \frac{1}{1 + \exp(\ln\frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)})}$$
$$= \frac{1}{1 + \exp(\ln\frac{1-\pi}{\pi} + \sum_{i} \ln\frac{P(X_{i}|Y=0)}{P(X_{i}|Y=1)})}$$

Until now we are simply replicating the notation in Tom's chapter. We now expand the denominator term using the notation from Hint 2:

$$\begin{split} \sum_{i} \ln \frac{P(X_{i}|Y=0)}{P(X_{i}|Y=1)} &= \sum_{i} \ln \frac{\theta_{i0}^{X_{i}}(1-\theta_{i0})^{(1-X_{i})}}{\theta_{i1}^{X_{i}}(1-\theta_{i1})^{(1-X_{i})}} \\ &= \sum_{i} \ln \theta_{i0}^{X_{i}} + \ln(1-\theta_{i0})^{(1-X_{i})} - \ln \theta_{i1}^{X_{i}} - \ln(1-\theta_{i1})^{(1-X_{i})} \\ &= \sum_{i} X_{i} \ln \theta_{i0} + (1-X_{i}) \ln(1-\theta_{i0}) - X_{i} \ln \theta_{i1} - (1-X_{i}) \ln(1-\theta_{i1}) \\ &= \sum_{i} X_{i} (\ln \theta_{i0} - \ln \theta_{i1}) + (1-X_{i}) (\ln(1-\theta_{i0}) - \ln(1-\theta_{i1})) \\ &= \sum_{i} X_{i} \ln \frac{\theta_{i0}}{\theta_{i1}} + (1-X_{i}) \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \\ &= \sum_{i} X_{i} \ln \frac{\theta_{i0}(1-\theta_{i1})}{\theta_{i1}(1-\theta_{i0})} + \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \end{split}$$

The final expression:

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^{n} w_i X_i)}$$

can be obtained obtained by setting the weights to:

$$w_i = \ln \frac{\theta_{i0}(1 - \theta_{i1})}{\theta_{i1}(1 - \theta_{i0})}$$

$$w_0 = \ln \frac{1-\pi}{\pi} + \sum_{i} \ln \frac{1-\theta_{i0}}{1-\theta_{i1}}$$

## Programming Generative and Discriminative Classifiers

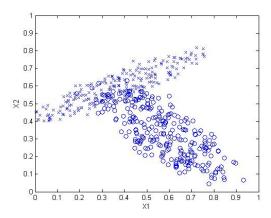
Discriminative classifiers learn the parameters of P(Y|X) directly, whereas generative classifiers instead learn the parameters of P(X|Y) and P(Y).

In this exercise you are asked to implement and compare both types of classifiers. You may use any programming language you like (Matlab, C++, C, Java...). All programming must be done from first principles. You are only permitted to use existing tools for simple linear algebra such as matrix multiplication/inversion. Do NOT use any toolkit that performs machine learning functions.

For this assignment, please submit all answers and any plots that are requested in the following questions. Also, *print out* and clearly label any code you wrote for this assignment and append it to the back of your submission. We do not require comments, however the clarity of your code and explanations will affect how much partial credit we can give. We encourage you to discuss the questions, but you must write/submit your own code and your own answers.

The provided data has two real variables  $X_1$ ,  $X_2$  and the boolean variable Y representing a class. Each line in the data files represents a data point  $(X_1, X_2, Y)$ .

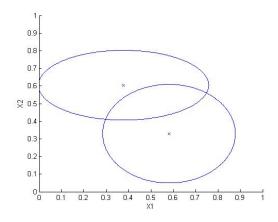
Make a 2D plot  $(X_1, X_2)$  of "test-1", using different symbols for data points that belong to different classes.



- 1. (Naive Bayes) [9 pts] You may notice that given the class (Y), the data follows a bi-variate normal distribution. For Naive Bayes, we make the assumption that  $X_1$  is conditionally independent of  $X_2$  given the class.
  - (a) Write the equation for the probability of some data point coming from class Y = 1 i.e.  $P(Y = 1|X_1 = x_1, X_2 = x_2)$ .

$$\begin{split} P(Y=1|X_1,X_2) &= \frac{P(X_1,X_2|Y=1)P(Y=1)}{P(X_1,X_2|Y=1)P(Y=1) + P(X_1,X_2|Y=0)P(Y=0)} \\ &= \frac{N(\mu_{11},\mu_{12},\sigma_{11}^2,\sigma_{12}^2)P(Y=1)}{N(\mu_{11},\mu_{12},\sigma_{12}^2,P(Y=1) + N(\mu_{01},\mu_{02},\sigma_{01}^2,\sigma_{02}^2)P(Y=0)} \end{split}$$

- (b) Write a program that estimates the means and variances of the Gaussians, as well as P(Y) for each class. This program should yield all the information necessary to complete your equation in (a). Run this program on "train-1." (Use the maximum likelihood estimates for  $\mu$  and  $\sigma^2$ ).
- (c) Part of your program should estimate the two normal distributions that maximize the probability of the data. Submit one plot that shows the mean and the two standard deviation iso-contour for each Gaussian.



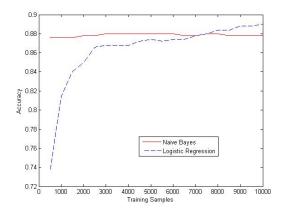
(d) Write a program that uses the parameters learned from "train-1" to classify the points in "test-1". Report your test error.

Accuracy = 87.8% Test Error = 12.2%.

- 2. (Logistic Regression) [9 pts] Now, let us not make the conditional independence assumption and classify with Logistic Regression. For this exercise, use gradient descent as presented in (Ch. 3.2) of Tom's handout.
  - (a) Write a program that optimizes the weights w0, w1 and w2 to construct a logistic regression model of the data. In your program, set the step size  $\eta = 1 \times 10^{-6}$  and fix the number of iterations to 20000. (Perform exactly this number of iterations whether or not gradient descent converges).
  - (b) Train your program on "train-1" and report the learned weights.
  - (c) Write a program that uses the linear regression model to predict the class of the data based on observed  $(X_1, X_2)$ . Test it on "test-1". Report your test error.

Accuracy = 89% Test Error = 11%.

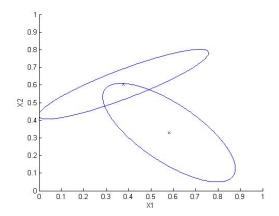
- 3. (Comparisons) [15 pts] Let us call Logistic Regression LR and Naive Bayes NB. So far, you have used the entire training set to train these classifiers. Suppose that less data was available. Limit your classifier to training on subsets of the provided training data. (e.g. 500, 1000, 1500 ... 10000.) Retrain LR and NB on subsets of "train-1" data and observe their performance on "test-1". To reduce needless variance in your experiment, be sure to use the same training data subsets (not just random subsets of the same size) to train both LR and FB.
  - (a) Construct a single plot showing the accuracy of each classifier as it depends on the number of training examples.



(b) What do you notice about the relationship between classifier accuracy and the number of available data points?

The Naive Bayes classifier performs significantly better than logistic regression with a small set of data. While the performance of NB does not improve much, Logistic Regression steadily improves in accuracy with added data. At the limit, LR converges to a greater accuracy than Naive Bayes.

- (c) Make a short list of observations about LR and NB that you have made throughout this exercise. Using knowledge from class/experience in programming and understanding of the assumptions briefly describe the causes for the differences you observe. (1-2 sentences per observation).
  - LR does not assume conditional independence of features and hence has the potential of yielding a better model (when the features are not conditionally independent).
  - LR performs gradient descent to get a solution and therefore for sparse data, many solutions can be optimal. (Not multiple minima, just one large minimal space). Only with sufficient data does it converge to a good predictive model.
  - Naive Bayes makes the Gaussian assumption. Therefore, for normally distributed data, it can quickly calculate a good predictive model.
  - Gradient descent (LR) requires many iterations to converge. It is significantly slower than the single linear calculation performed for Naive Bayes. How fast do you think conjugate gradient would be?
- 4. (Full Bayes) [2 pts Only if you have free time] Without making the conditional independence assumption, it is still possible to train a Bayesian classifier. Suppose that  $X_1$  is not conditionally independent of  $X_2$  given the class (i.e. you can't assume  $p(X_1|X_2,Y) = p(X_1|Y)$ ).
  - (a) Write the equation for the probability of some data point coming from class Y=1 i.e.  $P(Y=1|X_1=x_1,X_2=x_2)$ .
  - (b) Duplicate and modify your Naive Bayes classifier to reflect this changed formula. Train your classifier on "train-1." It should learn P(Y), the means and the full covariance matrices for the Gaussians.
  - (c) Plot the mean and 2 SD iso-contour for each Gaussian. How does your plot compare to Naive Bayes?



The covariance eigenvectors (principal axes of the Gaussians) are not axis aligned. This permits the Gaussians to fit the data more accurately than Naive Bayes.

(d) Test your classifier on "test-1." What do you notice?

Accuracy = 94.6% Test Error = 5.4%.

The accuracy improves substantially. This will not always be the case. However, when the underlying distributions have significant correlation between  $X_1$  and  $X_2$  the Full Bayes model is at an advantage.