

STAT/BIOSTAT 534 Statistical Computing

Spring Quarter 2017

Homework 4

Adrian Dobra
adobra@uw.edu

This homework is due on Friday, April 28 at 11:00pm. You should use the dropbox to submit your code. (Last update: 4/22)

The background material you need to use is presented in Section 1. The full description of the homework is given in Section 2. Note: you do not need any knowledge of Bayesian statistics to complete this assignment.

1 Bayesian Inference in Univariate Logistic Regression

We assume to have observed the n independent samples

$$\mathcal{D} = \{(y_1, x_1), \dots, (y_n, x_n)\}.$$

The response variable Y is binary, i.e. $y_i \in \{0, 1\}$ for $i = 1, 2, \dots, n$. The explanatory variable X can be continuous or discrete. We consider the univariate logistic regression model

$$\log \frac{P(y = 1|x)}{P(y = 0|x)} = \beta_0 + \beta_1 x. \quad (1)$$

Our model assumptions say that each y_i follows a Bernoulli distribution with probability of success $\pi_i = P(y_i = 1|x_i)$:

$$y_i \sim \text{Ber}(\pi_i).$$

Since the samples are assumed to be independent, the likelihood is:

$$L(\beta_0, \beta_1 | \mathcal{D}) = \prod_{i=1}^n [P(y_i = 1|x_i)]^{y_i} [1 - P(y_i = 1|x_i)]^{1-y_i},$$

where

$$\pi_i = P(y_i = 1|x_i) = \text{logit}^{-1}(\beta_0 + \beta_1 x_i) \in (0, 1).$$

We define the *logit* function:

$$\text{logit} : (0, 1) \longrightarrow (-\infty, +\infty), \quad \text{logit}(p) = \log \frac{p}{1-p}.$$

Its inverse is:

$$\text{logit}^{-1} : (-\infty, +\infty) \longrightarrow (0, 1), \quad \text{logit}^{-1}(x) = \frac{\exp(x)}{1 + \exp(x)}.$$

The log-likelihood is

$$\begin{aligned} l(\beta_0, \beta_1 | \mathcal{D}) &= \log L(\beta_0, \beta_1 | \mathcal{D}), \\ &= \sum_{i=1}^n (y_i \log \pi_i + (1 - y_i) \log [1 - \pi_i]). \end{aligned}$$

Simple calculations show that

$$\begin{aligned} \frac{\partial l(\beta_0, \beta_1 | \mathcal{D})}{\partial \beta_0} &= \sum_{i=1}^n [y_i - \pi_i], \\ \frac{\partial l(\beta_0, \beta_1 | \mathcal{D})}{\partial \beta_1} &= \sum_{i=1}^n [y_i x_i - \pi_i x_i]. \end{aligned}$$

We assume that the logistic regression coefficients follow independent $N(0, 1)$ priors. The joint posterior distribution of β_0 and β_1 is therefore given by

$$P(\beta_0, \beta_1 | \mathcal{D}) = \frac{1}{P(\mathcal{D})} \exp(l^*(\beta_0, \beta_1)), \quad (2)$$

where

$$l^*(\beta_0, \beta_1) = -\log(2\pi) - \frac{1}{2} (\beta_0^2 + \beta_1^2) + l(\beta_0, \beta_1 | \mathcal{D}),$$

and

$$P(\mathcal{D}) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(l^*(\beta_0, \beta_1)) d\beta_0 d\beta_1. \quad (3)$$

We call $P(\mathcal{D})$ the marginal likelihood associated with the univariate logistic regression (1).

The gradient of $l^*(\beta_0, \beta_1)$ is

$$\nabla l^*(\beta_0, \beta_1) = \begin{pmatrix} \frac{\partial l^*(\beta_0, \beta_1)}{\partial \beta_0} \\ \frac{\partial l^*(\beta_0, \beta_1)}{\partial \beta_1} \end{pmatrix}.$$

The Hessian matrix associated with $l^*(\beta_0, \beta_1)$ is

$$D^2 l^*(\beta_0, \beta_1) = \begin{bmatrix} \frac{\partial^2 l^*(\beta_0, \beta_1)}{\partial \beta_0^2} & \frac{\partial^2 l^*(\beta_0, \beta_1)}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 l^*(\beta_0, \beta_1)}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 l^*(\beta_0, \beta_1)}{\partial \beta_1^2} \end{bmatrix}.$$

1.1 The Newton-Raphson Algorithm

We determine the mode of the posterior distribution (2), i.e.

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmax}_{(\beta_0, \beta_1) \in \mathbb{R}^2} l^*(\beta_0, \beta_1),$$

by employing the Newton-Raphson algorithm. The procedure starts with the initial values $(\beta_0^{(0)}, \beta_1^{(0)}) = (0, 0)$. At iteration k , we update our current estimate $(\beta_0^{(k-1)}, \beta_1^{(k-1)})$ of the mode $(\hat{\beta}_0, \hat{\beta}_1)$ to a new estimate $(\beta_0^{(k)}, \beta_1^{(k)})$ as follows:

$$\begin{pmatrix} \beta_0^{(k)} \\ \beta_1^{(k)} \end{pmatrix} = \begin{pmatrix} \beta_0^{(k-1)} \\ \beta_1^{(k-1)} \end{pmatrix} - [D^2 l^*(\beta_0^{(k-1)}, \beta_1^{(k-1)})]^{-1} \nabla l^*(\beta_0^{(k-1)}, \beta_1^{(k-1)}).$$

The procedure stops when the estimates of the mode do not change after performing a new update, i.e. $|\beta_0^{(k)} - \beta_0^{(k-1)}| < \epsilon$ and $|\beta_1^{(k)} - \beta_1^{(k-1)}| < \epsilon$. Here ϵ is some small positive number, e.g. 0.0001. *Please note that we have already discussed and implemented the Newton-Raphson algorithm for univariate logistic regression.*

1.2 The Laplace Approximation

Since the integral (3) cannot be explicitly calculated, we need to approximate it numerically. We calculate the marginal likelihood $P(\mathcal{D})$ using the Laplace approximation, i.e.

$$P(\widehat{\mathcal{D}}) = 2\pi \exp(l^*(\hat{\beta}_0, \hat{\beta}_1)) \left\{ \det [-D^2 l^*(\hat{\beta}_0, \hat{\beta}_1)] \right\}^{-1/2}, \quad (4)$$

where $(\hat{\beta}_0, \hat{\beta}_1)$ is the mode of the posterior distribution (2). We note that you should not actually calculate $P(\widehat{\mathcal{D}})$. Instead, you should calculate the logarithm of the marginal likelihood $\log P(\widehat{\mathcal{D}})$.

1.3 The Metropolis-Hastings Algorithm

Sampling from the posterior distribution (2) can be done using the Metropolis-Hastings algorithm. The procedure starts with the initial values $(\beta_0^{(0)}, \beta_1^{(0)}) = (\hat{\beta}_0, \hat{\beta}_1)$, i.e. we start right at the mode of the distribution (2). We update the current state $(\beta_0^{(k-1)}, \beta_1^{(k-1)})$ of the Markov chain to its next state $(\beta_0^{(k)}, \beta_1^{(k)})$ as follows.

We generate a candidate state $(\tilde{\beta}_0, \tilde{\beta}_1)$ by sampling from the bivariate normal distribution

$$N_2 \left(\begin{pmatrix} \beta_0^{(k-1)} \\ \beta_1^{(k-1)} \end{pmatrix}, -[D^2 l^*(\hat{\beta}_0, \hat{\beta}_1)]^{-1} \right). \quad (5)$$

Note that the covariance matrix of the proposal (5) is the negative of the inverse of the Hessian matrix evaluated at the mode of (2). The R function `mvrnorm` might be useful here.

We accept the move to the proposed state, i.e. we set $(\beta_0^{(k)}, \beta_1^{(k)}) = (\tilde{\beta}_0, \tilde{\beta}_1)$ with probability

$$\min \left\{ 1, \exp \left[l^* (\tilde{\beta}_0, \tilde{\beta}_1) - l^* (\beta_0^{(k-1)}, \beta_1^{(k-1)}) \right] \right\}. \quad (6)$$

Otherwise the Markov chain stays at its current state, i.e. we set $(\beta_0^{(k)}, \beta_1^{(k)}) = (\beta_0^{(k-1)}, \beta_1^{(k-1)})$. We see that the proposal distribution (5) is symmetric, i.e. the probability of proposing $(\tilde{\beta}_0, \tilde{\beta}_1)$ if the chain is currently in $(\beta_0^{(k-1)}, \beta_1^{(k-1)})$ is equal with the probability of proposing $(\beta_0^{(k-1)}, \beta_1^{(k-1)})$ if the chain is currently in $(\tilde{\beta}_0, \tilde{\beta}_1)$. As such, the proposal distribution (5) cancels when we calculate the acceptance probability (6).

The implementation of an iteration of the Metropolis-Hastings algorithm proceeds as follows. If the proposed state leads to an increase of l^* , i.e.

$$l^* (\tilde{\beta}_0, \tilde{\beta}_1) \geq l^* (\beta_0^{(k-1)}, \beta_1^{(k-1)}),$$

we accept the move to the proposed state and set $(\beta_0^{(k)}, \beta_1^{(k)}) = (\tilde{\beta}_0, \tilde{\beta}_1)$. Otherwise, if the proposed state leads to a decrease in l^* , i.e.

$$l^* (\tilde{\beta}_0, \tilde{\beta}_1) < l^* (\beta_0^{(k-1)}, \beta_1^{(k-1)}),$$

we sample u from a Uniform(0, 1) distribution. If

$$\log(u) \leq l^* (\tilde{\beta}_0, \tilde{\beta}_1) - l^* (\beta_0^{(k-1)}, \beta_1^{(k-1)}),$$

we accept the move to the proposed state and set $(\beta_0^{(k)}, \beta_1^{(k)}) = (\tilde{\beta}_0, \tilde{\beta}_1)$. If

$$\log(u) > l^* (\tilde{\beta}_0, \tilde{\beta}_1) - l^* (\beta_0^{(k-1)}, \beta_1^{(k-1)})$$

we reject the move and the chain stays at the current state, i.e. we set $(\beta_0^{(k)}, \beta_1^{(k)}) = (\beta_0^{(k-1)}, \beta_1^{(k-1)})$.

2 Homework Problems

2.1 Problem 1 (30 points)

Write a function in R that computes the Laplace approximation (4) of the log of the marginal likelihood (3) associated with a univariate logistic regression. It should be named `laplaceLogLik(explanatory, response, data)`, where `explanatory` is the column number of the desired predictor in the dataset, `response` is the column number of the response, and `data` is a data frame.

2.2 Problem 2 (30 points)

Write a function in R that estimates the coefficients β_0 and β_1 of the univariate logistic regression (1) as follows:

1. Use the Metropolis-Hastings algorithm to simulate 10000 samples

$$\left\{ \left(\beta_0^{(k)}, \beta_1^{(k)} \right) : k = 1, \dots, 10000 \right\},$$

from the posterior distribution (2).

2. The estimates for β_0 and β_1 are given by the sample means:

$$\bar{\beta}_0 = \frac{1}{10000} \sum_{k=1}^{10000} \beta_0^{(k)}, \quad \bar{\beta}_1 = \frac{1}{10000} \sum_{k=1}^{10000} \beta_1^{(k)}.$$

It should be named `metropolisHastings(explanatory, response, data, iterations)`, where `iterations` denotes the number of iterations to run the algorithm.

2.3 Problem 3 (40 points)

The dataset you need to use is contained in the file “534binarydata.txt”. This file has 148 rows (samples) and 61 columns (variables). The first 60 columns are associated with 60 explanatory variables X , while column 61 (the last column) corresponds with the response binary variable Y .

You need to write a parallel program using the R package `snow` that computes *in parallel* the Laplace approximation (4) and the estimates of β_0 and β_1 associated with each univariate logistic regression of the response binary variable Y on each of the 60 explanatory variables. You can use the MLEs of β_0 and β_1 , given by `mleLogistic`, to give yourself a way to check your estimates from Problem 2.

Please start from the code below. Your task is to write the function `bayesLogistic` as well as all the other functions that are called in `bayesLogistic`.

```
bayesLogistic = function(apredictor,response,data,NumberOfIterations)
{
  ....
}

#PARALLEL VERSION
#datafile = the name of the file with the data
#NumberOfIterations = number of iterations of the Metropolis-Hastings algorithm
#clusterSize = number of separate processes; each process performs one or more
#univariate regressions
main <- function(datafile,NumberOfIterations,clusterSize)
{
  #read the data
  data = read.table(datafile,header=FALSE);
```

```

#the sample size is 148 (number of rows)
#the explanatory variables are the first 60 columns for '534binarydata.txt'
#the last column is the binary response
response = ncol(data);
lastPredictor = ncol(data)-1;

#initialize a cluster for parallel computing
cluster <- makeCluster(clusterSize, type = "SOCK")

#run the MC3 algorithm from several times
results = clusterApply(cluster, 1:lastPredictor, bayesLogistic,
                        response,data,NumberOfIterations);

#print out the results
for(i in 1:lastPredictor)
{
  cat('Regression of Y on explanatory variable ',results[[i]]$apredictor,
      ' has log marginal likelihood ',results[[i]]$logmarglik,
      ' with beta0 = ',results[[i]]$beta0bayes,' (',results[[i]]$beta0mle,')',
      ' and beta1 = ',results[[i]]$beta1bayes,' (',results[[i]]$beta1mle,')',
      '\n');
}

#destroy the cluster
stopCluster(cluster);
}

#NOTE: YOU NEED THE PACKAGE 'SNOW' FOR PARALLEL COMPUTING
require(snow);

#this is where the program starts
main('534binarydata.txt',10000,10);

```