

## Bayesian Learning, 6 hp

### Computer lab 4

You can use any programming language for the labs, but my hints, help and solutions will be in R.

You are supposed to work and submit your labs in pairs, but do make sure that both of you are contributing. Submit your solutions by Lisam no later than **May 24 at midnight**. You should submit the report (pdf only) and an executable file containing your code.

1. *Poisson regression - the MCMC way.*

Consider the following Poisson regression model

$$y_i|\beta \sim \text{Poisson}[\exp(x_i'\beta)], \quad i = 1, \dots, n,$$

where  $y_i$  is the count for the  $i$ th observation in the sample and  $x_i$  is the  $p$ -dimensional vector with covariate observations for the  $i$ th observation. Use the data set `eBayNumberOfBidderData.dat`. This dataset contains observations from 1000 eBay auctions of coins. The response variable is **nBids** and records the number of bids in each auction. The remaining variables are covariates ( $x$ ):

- **Const** (for the intercept)
  - **PowerSeller** (is the seller selling large volumes on eBay?)
  - **VerifyID** (is the seller verified by eBay?)
  - **Sealed** (was the coin sold sealed in never opened envelope?)
  - **MinBlem** (did the coin have a minor defect?)
  - **MajBlem** (a major defect?)
  - **LargNeg** (did the seller get a lot of negative feedback from customers?)
  - **LogBook** (logarithm of the coins book value according to expert sellers. Standardized)
  - **MinBidShare** (a variable that measures ratio of the minimum selling price (starting price) to the book value. Standardized).
- (a) Obtain the maximum likelihood estimator of  $\beta$  in the Poisson regression model for the eBay data [Hint: `glm.R`, don't forget that `glm()` adds its own intercept so don't input the covariate `Const`]. Which covariates are significant?

- (b) Let's now do a Bayesian analysis of the Poisson regression. Let the prior be  $\beta \sim \mathcal{N}[\mathbf{0}, 100 \cdot (X'X)^{-1}]$  where  $X$  is the  $n \times p$  covariate matrix. This is a commonly used prior which is called Zellner's g-prior. Assume first that the posterior density is approximately multivariate normal:

$$\beta|y \sim \mathcal{N}(\tilde{\beta}, J_{y,\tilde{\beta}}^{-1}),$$

where  $\tilde{\beta}$  is the posterior mode and  $J_{\tilde{\beta}}$  is the negative Hessian at the posterior mode.  $\tilde{\beta}$  and  $J_{\tilde{\beta}}$  can be obtained by numerical optimization (`optim.R`) exactly like you already did for the probit regression on the spam data (but with the log posterior function replaced by the corresponding one for the Poisson model).

- (c) Now, let's simulate from the actual posterior of  $\beta$  using the Metropolis algorithm and compare with the approximate results in b). Program a general function that uses the Metropolis algorithm to generate random draws from an *arbitrary* posterior density. Let the proposal density be the multivariate normal density mentioned in Lecture 9 (random walk Metropolis):

$$\theta_p|\theta_c \sim N(\theta_c, \tilde{c} \cdot \Sigma),$$

where  $\Sigma = J_{\tilde{\theta},y}^{-1}$  and  $J_{\tilde{\theta},y}$  is the negative Hessian at the posterior mode  $\tilde{\theta}$ . These are already available from question b) above. The value  $\tilde{c}$  is a tuning parameter and should be an input to your Metropolis function (so that a user can change it). The user of your Metropolis function should be able to supply her own posterior density function and still be able to use your Metropolis function. This is not so straightforward, unless you have come across *function objects* in R and the triple dot (...) wildcard argument. Let me give you the roadmap. Ok, take a deep breath ... this is good stuff ...

First, one of the input arguments of your Metropolis function should be `logPostFunc` (or some other suitable name). `logPostFunc` is a *function object* that computes the log posterior density at any value of the parameter vector. This is needed when you compute the acceptance probability of the Metropolis algorithm. I suggest always to program the *log* posterior density, since logs are more stable and avoids problems with too small or large numbers (overflow). Note that the ratio of posterior densities in the Metropolis acceptance probability can be written

$$\frac{p(\theta_p|y)}{p(\theta_c|y)} = \exp[\log p(\theta_p|y) - \log p(\theta_c|y)]$$

This is smart since the large or small common factors in  $p(\theta_p|y)$  and  $p(\theta_c|y)$  cancel out before we evaluate the exponential function (which can otherwise overflow).

Second, the first argument of your (log) posterior function should be `theta`, the (vector) of parameters for which the posterior density is evaluated. You can of course use some other name for the variable, but it must be the *first* argument of the posterior density function.

Third, the user's posterior density is also a function of the data and prior hyper-parameters and those need to be supplied to the Metropolis function in some way. A nice way of doing that is to use the triple dot (...) argument which is like a wildcard for *any* parameters supplied by the user. This makes it possible to use the Metropolis function for any problem, even when you as a programmer don't know what the user's posterior density function looks

like or what kind of data and hyper-parameters being used in that particular problem. To make all of this very clear (I hope!) I give a simple code below where the log posterior density for the Bernoulli model with a Beta prior is coded. That log posterior density is then used in a useless, but illustrative, function `MultiplyByTwo` that returns 2 times the log posterior density evaluated at  $\theta = 0.3$ . Note how the `MultiplyByTwo` takes a function object as input and how it uses the triple dot (`...`) argument to supply the data  $s$  and  $f$ , and the prior hyper-parameters  $a$  and  $b$  without explicitly using these symbols inside the function. That makes the useless `MultiplyByTwo` function applicable for *any* function.

[MATLAB users: Here is one instance where MATLAB's lack of proper function objects becomes a pain. You can do the same thing with `eval()` and `varargin`, however, but it is not as pretty as in R. See the MATLAB help.]

Use your new Metropolis function to sample from the posterior of  $\beta$  in the Poisson regression for the eBay dataset. Assess MCMC convergence by graphical methods and compute the effective sample size (`CODA.R`). The parameters  $\phi_j = \exp(\beta_j)$  are usually considered more interpretable than the  $\beta_j$ . Compute the posterior distribution of  $\phi_j$  for all variables.

- (d) Use the MCMC draws from c) to simulate from the predictive distribution of the number of bidders in a new auction with the characteristics below. Plot the predictive distribution. What is the probability of no bidders in this new auction?

- **PowerSeller** = 1
- **VerifyID** = 1
- **Sealed** = 1
- **MinBlem** = 0
- **MajBlem** = 0
- **LargNeg** = 0
- **LogBook** = 1
- **MinBidShare** = 0.5

MAY BAYES BE WITH YOU!

```

# This is the log posterior density of the beta(s+a,f+b) density
LogPostBernBeta <- function(theta, s, f, a, b){
logPost <- (s+a-1)*log(theta) + (f+b-1)*log(1-theta)
return(logPost)
}

# Testing if the log posterior function works
s <- 8;f <- 2;a <- 1;b <- 1
logPost <- LogPostBernBeta(theta = 0.1, s, f, a, b)
print(logPost)

# This is a rather useless function that takes the function myFunction,
# evaluates it at x = 0.3, and then returns two times the function value.
MultiplyByTwo <- function(myFunction, ...){
x <- 0.3
y <- myFunction(x,...)
return(2*y)
}
#Let's try if the MultiplyByTwo function works:
MultiplyByTwo(LogPostBernBeta,s,f,a,b)

```