# 732A91: Lab 4
# Bayesian Learning

Sarah Alsaadi, Carles Sans Fuentes

May 24, 2017

---

## Poisson regression-the MCMC way

Consider the following Poisson regression model

$$y_i|\beta \sim \text{Poisson}[\exp(\mathbf{x}_i^T \beta)], i = 1, \ldots, n,$$

where $y_i$ is the count for the $i$th observation in the sample and $\mathbf{x}_i$ is the p-dimensional vector with covariate observations for the $i$th observation. The data set **eBayNumberOfBidderData.dat** 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 features/covariates ($\mathbf{x}$):

- **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) Using **glm** in R, we obtain the following results:

```
 1 Call:
 2 glm(formula = nBids ~ 0 + ., family = poisson, data = ebay)
 3
 4 Deviance Residuals:
 5     Min       1Q    Median       3Q       Max
 6 -3.5800  -0.7222  -0.0441    0.5269    2.4605
 7
 8 Coefficients:
 9              Estimate Std. Error z value Pr(>|z|)
10 Const         1.07244    0.03077  34.848  < 2e-16 ***
11 PowerSeller  -0.02054    0.03678  -0.558   0.5765
12 VerifyID     -0.39452    0.09243  -4.268 1.97e-05 ***
13 Sealed        0.44384    0.05056   8.778  < 2e-16 ***
14 Minblem      -0.05220    0.06020  -0.867   0.3859
15 MajBlem      -0.22087    0.09144  -2.416   0.0157 *
16 LargNeg       0.07067    0.05633   1.255   0.2096
17 LogBook      -0.12068    0.02896  -4.166 3.09e-05 ***
18 MinBidShare  -1.89410    0.07124 -26.588  < 2e-16 ***
19 ---
20 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
21
22 (Dispersion parameter for poisson family taken to be 1)
23
24     Null deviance: 6264.01  on 1000  degrees of freedom
```

```
25  Residual deviance:  867.47  on  991  degrees of freedom
26  AIC: 3610.3
27
28  Number of Fisher Scoring iterations: 5
```

The output shows that the significant MLE of the $\beta$:s are those for the covariates Const, VeryfyID, Sealed, MayBlem, LogBook and MinBidShare.

(b) In this exercise we did a Bayesian analysis of the Poisson regression. We let $\beta \sim N[0, 100(\mathbf{X}^T\mathbf{X})^{-1}]$ apriori, where $\mathbf{X}$ is the $n \times p$ covariate matrix. Next we assumed that the posterior density is approximately multivariate normal:

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

where $\tilde{\beta}$ is the posterior mode and $J_{\mathbf{y}}(\tilde{\beta}) = -\frac{\partial^2 lnp(\beta|\mathbf{y})}{\partial\beta\partial\beta^T}|_{\beta=\tilde{\beta}}$ is the observed Hessian evaluated at the posterior mode. To obtain $\tilde{\beta}$ and $J_{\mathbf{y}}^{-1}(\tilde{\beta})$ we used numerical optimization (**optim.R**). The results of the estimated $\beta$:s are given in the table below. We also made histograms, one for each variable, of 10000 draws of $\beta$, the marginal distribution looks normal, see Figure 1.

| Variable | $\tilde{\beta}$ |
|---|---|
| **Const** | 1.06984118 |
| **PowerSeller** | $-0.02051246$ |
| **VerifyID** | $-0.39300599$ |
| **Sealed** | 0.44355549 |
| **MinBlem** | $-0.05246627$ |
| **MajBlem** | 0.22123840 |
| **LargNeg** | 0.07069683 |
| **LogBook** | $-0.12021767$ |
| **MinBidShare** | $-1.89198501$ |

Table 1: The obtained $\beta$:s through maximization of the posterior distirbution, one for each variable.
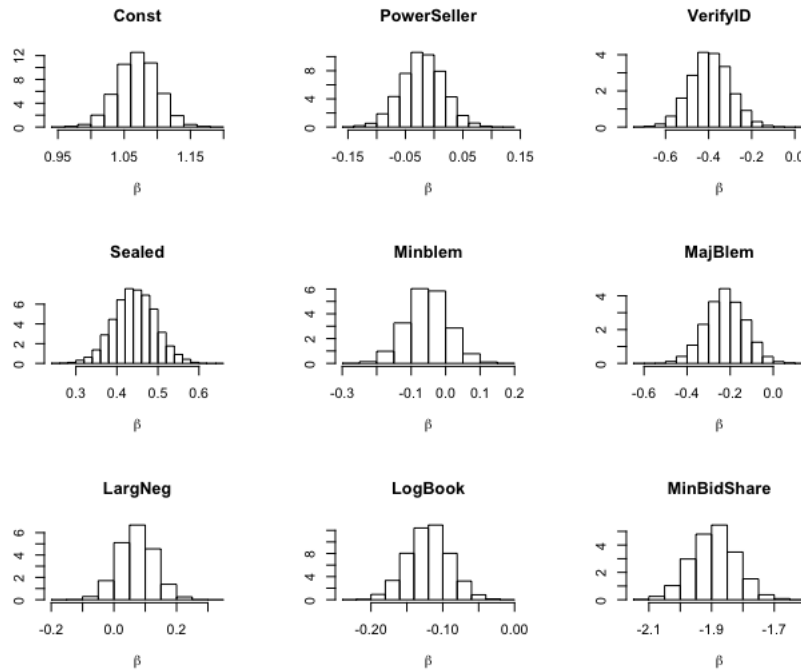


Figure 1: Histogram of the 10000 $\beta$:s drawn from the marginal posterior distribution

(c) In this exercise we simulate from the actual posterior of $\beta$ using the Metropolis algorithm and compare with the approximate results in b). We program a general function that uses

2

the Metropolis algorithm to generate random draws from an arbitrary posterior density. We use the multivariate normal density as proposal density:

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

where $\Sigma = J_{\mathbf{y}}^{-1}(\tilde{\beta})$ obtained in the previous exercise and $\theta_c$ is the current draw (hence the subscript c). The value $\tilde{c}$ is a tuning parameter and is an input to our Metropolis function (so that a user can change it). The user of our Metropolis function is able to supply her own posterior density function, not necessarily for the Poisson regression, and still be able to use our Metropolis function.

First, one of the input arguments of our Metropolis function is called logPostFunc. logPost-Func is a function object that computes the log posterior density at any value of the parameter vector. This is needed when we compute the acceptance probability of the Metropolis algorithm. We 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

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

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

Second, the first argument of our (log) posterior function is theta, the (vector) of parameters for which the posterior density is evaluated.

Third, the user's posterior density is also a function of the data and prior hyperparameters and those are supplied to the Metropolis function where we 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 a programmer don't know what the user's posterior density function looks like or what kind of data and hyperparameters being used in that particular problem.

Now, we use Metropolis function to sample from the posterior of $\beta$ in the Poisson regression for the eBay dataset. We assess MCMC convergence by graphical methods. The parameters $\phi_j = \exp(\beta_j)$ are usually considered more interpretable than the $\beta_j$. We compute the posterior distribution of $\phi_j$ for all variables.

Figure 2 shows that the Metropolis algorithm seem to converge, with some burn in period for all the $\beta$:s.

In Table 2 we have the mean posterior of the simulated $\beta$:s which are very similar to the ones obtained through maximization.

Figure 3 shows the marginal posterior for each $e^\beta$ obtained by the Metropolis algorit. The distribution looks normal.
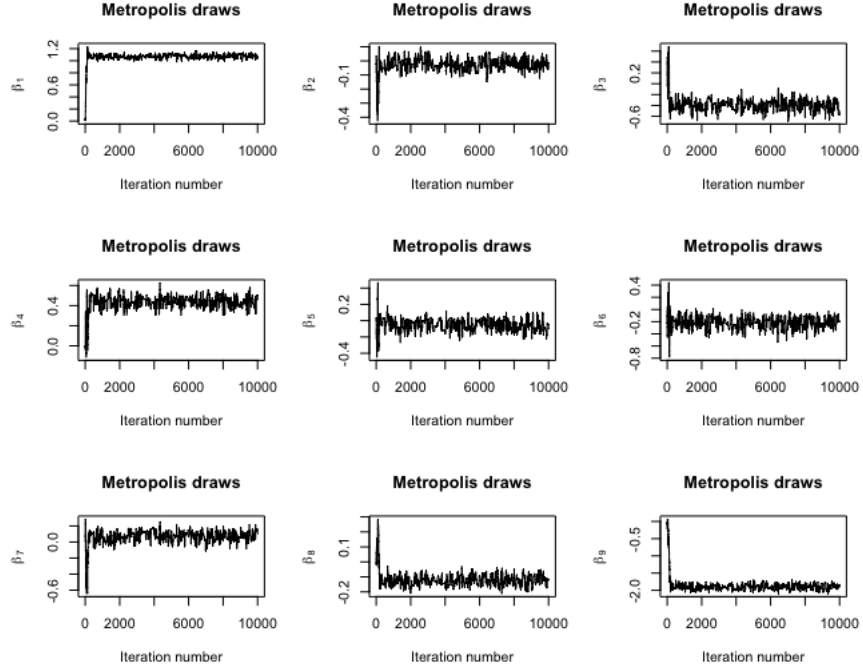
Figure 2: Plot of the 10000 $\beta$:s drawn from the marginal posterior distribution obtained by the Metropolis algorithm.

| Variable | $\tilde{\beta}$ |
|---|---|
| **Const** | 1.06538266 |
| **PowerSeller** | $-0.02529706$ |
| **VerifyID** | $-0.38552635$ |
| **Sealed** | 0.43553156 |
| **MinBlem** | $-0.05410433$ |
| **MajBlem** | $-0.22019352$ |
| **LargNeg** | 0.06493763 |
| **LogBook** | $-0.12021880$ |
| **MinBidShare** | $-1.88132053$ |

Table 2: The posterior mean of the simulated $\beta$:s, simulated using the Metropolis algorithm, the values are almost identical as the ones obtained through maximization
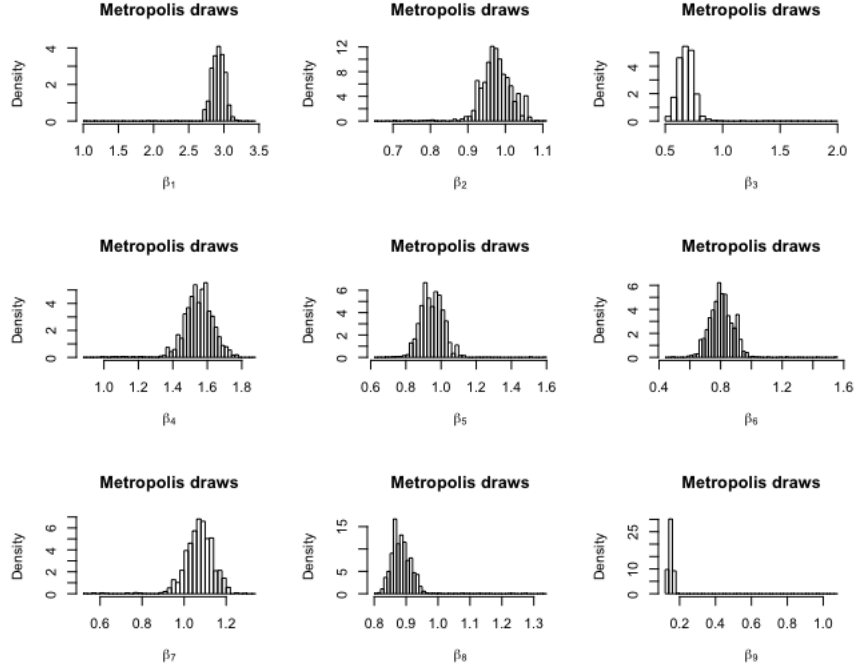
Figure 3: Histogram of the 10000 $e^{\beta}$ drawn from the marginal posterior distribution obtained by the Metropolis algorithm.

(d) We 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. The histogram of the predictive distribution is shown in Figure 4. The probability of no bidders in this new auction is 0.3512.

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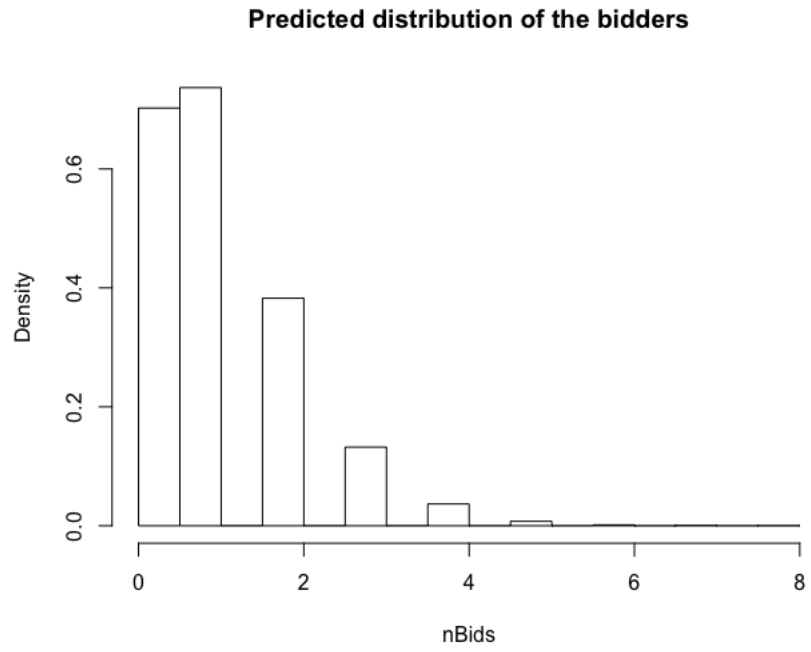- **LogBook**= 1
- **MinBidShare**= 0.5

Figure 4: Histogram of 10000 simulated predicted values of the variable nBid with given characteristics

# Contributions

All results and comments presented have been developed and discussed together by the members of the group.

# Appendix

## Poisson regression-the MCMC way

```r
1
2  ### LAB 4
3
4  #1
5  Data<- read.csv("C:/Users/Carles/Desktop/Bayesian learning/Part4/eBayNumberOfBidderData.dat.txt
       ", sep = "")
6
7  ##a
8  #### Variables
9
10 PoisModel<-glm(nBids~.-Const, family="poisson", data = Data)
11 logLik(PoisModel)
12
13
14 mysum<-summary(PoisModel)[["coefficients"]]
15
16
17
18 ##Significants at a 95% level
19 significant <- mysum[which(mysum[,4]<0.05),]
20 non_significant <- mysum[which(mysum[,4]>=0.05),]
21 list(significant = significant, nonsignificant = non_significant)
22
23
24 ##b
25 chooseCov <- 1:(length(names(Data))-1);  # Here we choose which covariates to include in the
       model
26
27
28 y <- as.vector(Data[,1]); # Data from the read.table function is a data frame. Let's convert y
        and X to vector and matrix.
29 X <- as.matrix(Data[,2:length(names(Data))]);
30 covNames <- names(Data)[2:length(names(Data))];
31 X <- X[,chooseCov]; # Here we pick out the chosen covariates.
32 covNames <- covNames[chooseCov];
33 nPara <- dim(X)[2]
34
35
36 ##prior Beta
37 library(geoR)
38 library(mvtnorm)
39
40
41 mu_0    <- matrix(0, nPara,1)
42 Sigma_0 <- 100*solve(crossprod(X,X))
43 Beta_0  <- rmvnorm(1, mean = mu_0, sigma = Sigma_0)
44 InitVal <- matrix(0, ncol=nPara, nrow =1)
45
46
47 LogPostPois <- function(betaVect,y = y,X = X, mu =mu_0,Sigma= Sigma_0){
48   nPara <- length(betaVect);
49   linPred <- X%*%betaVect;
50
51   # The following is a more numerically stable evaluation of the log-likelihood in my slides:
52   # logLik <- sum(y*log(pnorm(linPred)) + (1-y)*log(1-pnorm(linPred)) )
53   logLik <- sum(linPred*y-exp(linPred))
54
55   # evaluating the prior
56
57   logPrior <- dmvnorm(betaVect, mu, Sigma, log=TRUE);
58
59   # add the log prior and log-likelihood together to get log posterior
60   return(logLik + logPrior)
61
62 }
63
64 OptimResults<-optim(InitVal,LogPostPois,gr=NULL,y = y,X = X, mu =mu_0,Sigma= Sigma_0 ,method=c(
       "BFGS"),control=list(fnscale=-1),hessian=TRUE)
65
66 BetaCoef<- OptimResults$par
67 colnames(BetaCoef)<- covNames
68 J<--solve(OptimResults$hessian)
69
70 mycoefvar<- data.frame(Coefficients= as.vector(BetaCoef), variance = diag(J))
71 rownames(mycoefvar)<- covNames
72 mycoefvar
73 ######C
74
75 set.seed(12345)
76
77 LogPostPoisson <- function(theta, priormu, priorsigma, X, Y, ...) {
```

```r
78     require(mvtnorm)
79     likelihood <- dpois(Y, lambda = as.vector(exp((X) %*% t(theta))), log = TRUE)
80     prior <- dmvnorm(theta, mean = priormu, sigma = priorsigma, log=TRUE)
81     return(sum(likelihood) + prior)
82  }
83
84
85  metropl<-function(logPostFunc, theta_0, constant, sigma, nIter,...){
86    #initialize chain
87    require(mvtnorm)
88    theta<- matrix(NA, nrow = nIter+1, ncol = dim(sigma)[1])
89    theta[1,]<- theta_0
90    rej_rate<-0
91    for(i in 2:(nIter+1)){
92      new_theta<- rmvnorm(1, mean = theta[i-1,], sigma = constant*sigma)
93      cur_theta<-t(as.matrix(theta[i-1,]))
94      U<-runif(1,0,1)
95      num<-logPostFunc(new_theta,...)+dmvnorm(cur_theta, mean =new_theta, sigma = sigma, log =
              TRUE)
96      den<-logPostFunc(cur_theta,...)+dmvnorm(new_theta, mean =cur_theta, sigma = sigma, log =
              TRUE)
97
98      if(U<min(1,exp(num-den))){
99        theta[i,] <-new_theta
100     }else{
101       theta[i,] <-cur_theta
102       rej_rate  <-rej_rate+1
103     }}
104
105
106   myres<- list(theta= theta, rej_rate= rej_rate/nIter)
107   return(myres)
108 }
109
110
111 res<-metropl(logPostFunc=LogPostPoisson,
112                    theta_0= rep(0,nPara),
113                    constant= 0.6,
114                    sigma = J ,
115                    priormu= mu_0,
116                    priorsigma= Sigma_0,
117                    X= X, Y=y,
118                    nIter= 10000)
119
120
121 res[[2]]
122
123 compbetas<- data.frame(OptimCoefficients= as.vector(BetaCoef), MCMCCoef =colMeans(res[[1]]))
124 rownames(compbetas)<- covNames
125 compbetas
126 ##d
127 Xpred <- matrix(c(1, 1, 1, 1, 0, 0, 0, 1, 0.5), nrow = 1)
128 predsamples <- rpois(10000, lambda = exp(Xpred %*% t(res[[1]])))
129 hist(predsamples, freq = FALSE)
130
131 ##probability of 0
132 length(which(predsamples==0))/length(predsamples)
```