

Lab4 Report

Ludvig Noring, Michael Sörsäter

May 17, 2017

Poisson regression

a)

Maximum likelihood estimation is performed with the generalized model using the family poisson. The result can be shown in the following table.

	Const	PowerSeller	VerifyID	Sealed	Minblem	MajBlem	LargNeg	LogBook	MinBidShare
1	1.0724	-0.0205	-0.3945	0.4438	-0.0522	-0.2209	0.0707	-0.1207	-1.8941

Table 1: Coefficients for glm

The covarities that significantly influence the response are VerifyID, Sealed and MinBidShare.

b)

We implemented the log posterior method for the Poisson model. This produced similar results to what we got with glm. The result can be shown in the following table.

	Const	PowerSeller	VerifyID	Sealed	Minblem	MajBlem	LargNeg	LogBook	MinBidShare
1	1.0698	-0.0205	-0.3930	0.4436	-0.0525	-0.2212	0.0707	-0.1202	-1.8920

Table 2: Coefficients for optim

c)

We implemented a general random walk-metropolis algorithm and used it with our log-Poisson-Posterior function. With 1000 draws the theta converges as shown in Figure 1.

	Const	PowerSeller	VerifyID	Sealed	Minblem	MajBlem	LargNeg	LogBook	MinBidShare
1	1.0671	-0.0154	-0.4038	0.4414	-0.0506	-0.2268	0.0700	-0.1172	-1.8949

Table 3: Coefficients for optim

d)

Using the MCMC from above the probability for k number of bids are presented in Figure 2. As shown in the Figure, the probability of 0 bids are 35.6 %.

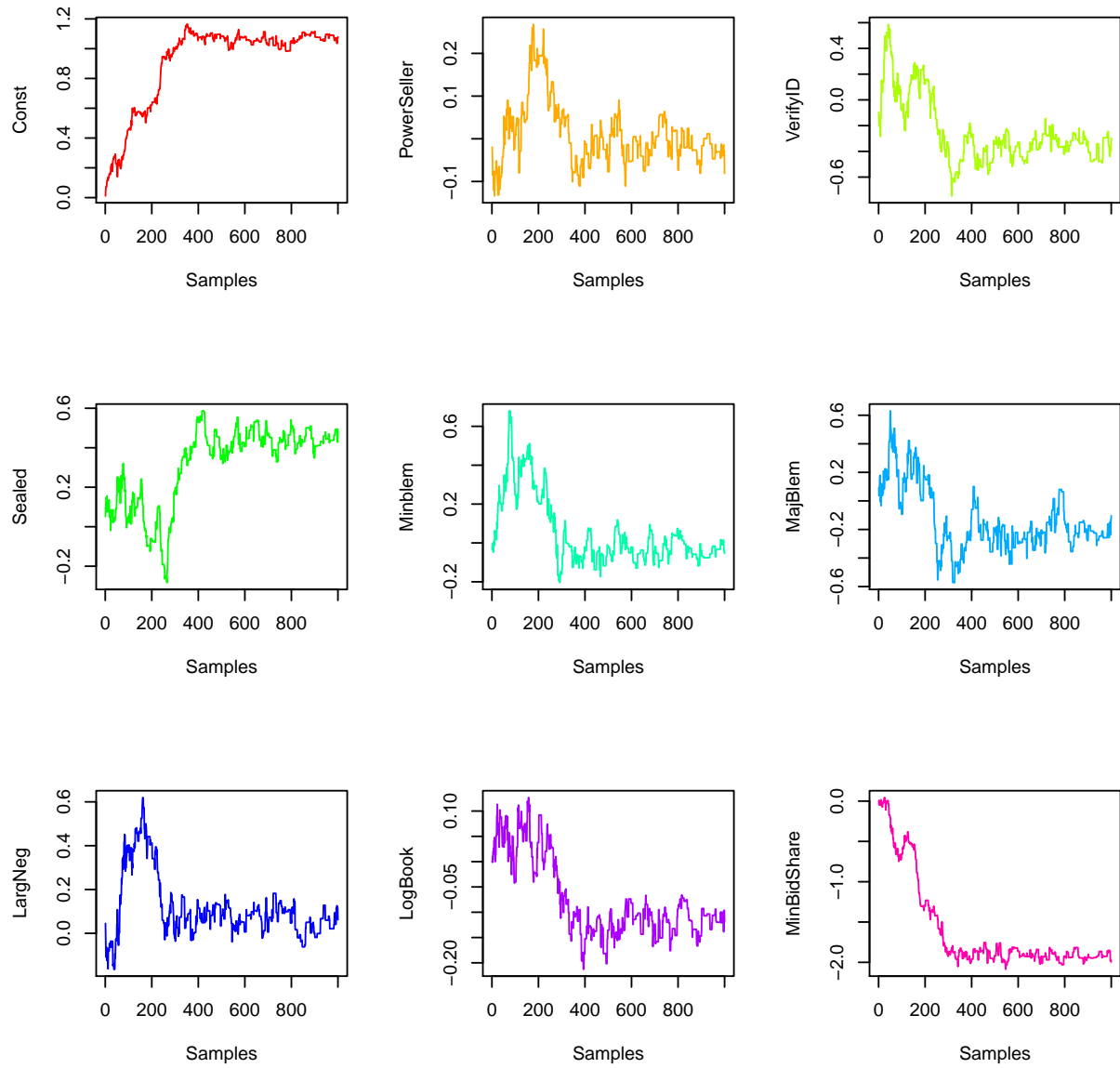


Figure 1: Convergence of θ

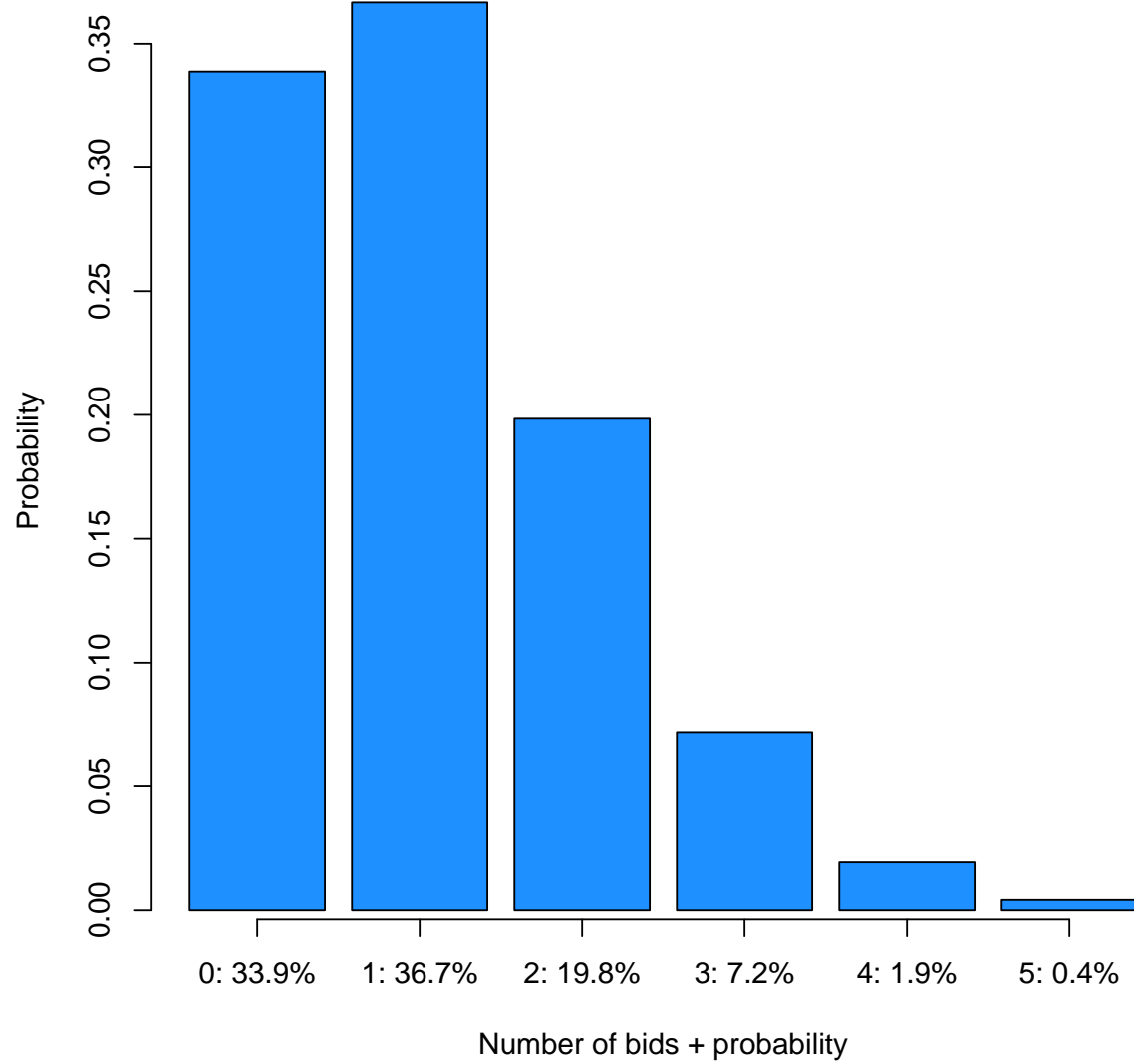


Figure 2: Bid Predictions

Appendix - Code for the assignment

```
require('msm')
require('mvtnorm')

data = read.table('WomenWork.dat', header=TRUE)
Y = data$Work
X = as.matrix(data[, -data$Work])
nObs = nrow(X)
nFeats = ncol(X)

# Calculate the accuracy of the beta values
performance = function(betas){
  y_hat = as.vector(X %*% betas)
  y_hat = sign(y_hat)
  y_hat[y_hat==--1] = 0

  print(table(y_hat, Y))
  print(mean(y_hat != Y))
}

uGenerator = function(betas) {
  curMean = X %*% t(betas)
  u = rep(0, nObs)
  for(i in 1:nObs){
    # Define bounds = (lower, upper)
    if(Y[i] == 0){
      bounds = c(-Inf, 0)
    }else{
      bounds = c(0, Inf)
    }
    u[i] = rtnorm(n=1, mean=curMean[i], sd=1, lower=bounds[1], upper=bounds[2])
  }
  return(u)
}

# From the given u, draw new beta-values
betaGenerator = function(u, sigma2=1) {
  # From lecture 5, slide 8
  XtX = t(X) %*% X
  beta_hat = solve(XtX) %*% t(X) %*% as.matrix(u)
  mu_n = solve(XtX + omega_0) %*% (XtX %*% beta_hat + omega_0 %*% mu_0)
  omega_n = XtX + omega_0

  beta = rmvnorm(n=1, mean=mu_n, sigma=(sigma2*solve(omega_n)))
  return (beta)
}

# Initial values
mu_0 = 0
tau = 10
mu_0 = rep(0, nFeats)
omega_0 = tau^2 * diag(nFeats)
```

```

beta_prior = as.matrix(rnorm(n=nFeats, mean=mu_0, sd=sqrt(diag(omega_0))))
u = uGenerator(t(beta_prior))

# Necessary with 1000 iterations, draws more to get smoother histograms
draws = 10000
# One column for each beta parameter
result_beta = matrix(0, draws, nFeats)
for(i in 1:draws) {
  print(i)
  beta = betaGenerator(u)
  u = uGenerator(beta)
  result_beta[i,] = beta
}

performance(colMeans(result_beta))

# c

logPostProbit = function(betas, y, X) {
  yPred = as.matrix(X) %*% betas

  logLike = sum(y*pnorm(yPred, log.p = TRUE) + (1-y)*pnorm(yPred, log.p = TRUE, lower.tail = FALSE))
  logPrior = dmvnorm(betas, mu_0, omega_0, log=TRUE);

  # add the log prior and log-likelihood together to get log posterior
  return(logLike + logPrior)
}

initValues = rep(0, nFeats)
optimResults = optim(initValues,
                     logPostProbit, y=Y, X=X,
                     method=c('BFGS'), control=list(fnscale=-1), hessian=TRUE)

postMode = optimResults$par
postCov = -solve(optimResults$hessian)
approxPostStd = sqrt(diag(postCov)/nFeats)

performance(postMode)

pdf('plots/betas.pdf')
par(mfrow=c(4,2))
for(i in 1:nFeats) {
  mean = postMode[i]
  sd = approxPostStd[i]

  # Remove outliers. (2 < x < 98) %
  draws = result_beta[,i]
  threshold = 0.02
  bounds = quantile(draws, probs=c(threshold, 1 - threshold))
  draws = draws[draws > bounds[1]]
  trimmed = draws[draws < bounds[2]]
}

```

```

# Toggle to include/exclude outliers
values = result_beta[,i]
values = trimmed

xRange = c(values, mean + sd * 4, mean - sd * 4)
betaGrid = seq(min(xRange), max(xRange), length=1000)

gibbs_density = hist(values, breaks=50, plot=FALSE)
approx_density = dnorm(betaGrid, mean=mean, sd=sd)

yMax = max(c(approx_density, gibbs_density$density))

plot(gibbs_density, freq=FALSE, col='dodgerblue', border='dodgerblue',
      xlim=c(min(betaGrid), max(betaGrid)),
      ylim=c(0, yMax),
      xlab=substitute(beta[idx], list(idx=i)),
      main=colnames(X)[i])
lines(betaGrid, approx_density, col='tomato', lwd=2)
}
dev.off()

```