

Applied Bayesian Analysis : NCSU ST 540

Homework 7

Bruce Campbell

In this assignment we perform Bayesian linear regression for the microbiome data on the course website

<https://www4.stat.ncsu.edu/~reich/ABA/assignments/homes.RData>

Let Y_i be the precipitation for observation i and X_{ij} equal one if OTU j is present in sample i .

First, extract the 50 OTU with the largest absolute correlation between X_{ij} and Y_i . Then fit a Bayesian linear regression model precipitation as the response and with these 50 covariates (and an intercept term) using two priors:

- (1) Uninformative normal priors: $\beta_j \sim \text{Normal}(0, 100^2)$
- (2) Hierarchical normal priors: $\beta_j | \tau \sim \text{Normal}(0, \tau^2)$ where $\tau^2 \sim \text{InvGamma}(0 : 01, 0 : 01)$
- (3) Bayesian LASSO: $\beta_j | \tau^2 \sim \text{DE}(0, \tau^2)$ where $\tau^2 \sim \text{InvGamma}(0 : 01, 0 : 01)$

Compare convergence and the posterior distribution of the regression coefficients under these three priors. In particular, are the same OTU's significant in all three fits?

Load data and select 50 most ocrrelated OUT variables.

```
library(rjags)
library(coda)
library(modeest)
load("homes.RData")

X <- OTU != 0
Y <- homes$MeanAnnualPrecipitation

C_xy <- cor(X, Y)

top <- function(x, n) {
  tail(order(x), n)
}

# One of the X is all 1's -
# resulting in an NA for the
# correlation.
indices <- top(C_xy, 51)
# Remove the NA - I'm sure there's
# a more elegant way...
indices <- indices[1:50]
```

```

X <- X[, indices]

predictor.names <- names(OTU)[indices]
predictor.names[51] <- "intercept"

top.corr <- C_xy[indices]

# Y <- scale(Y) X <- scale(X)

DEBUG <- FALSE
if (DEBUG) {
  nSamples <- 20000
  n.chains <- 1
} else {
  nSamples <- 20000
  n.chains <- 1
}

```

We sample from our model after burn in. Not all of the diagnostic plots are not presented. See the diagnostic plots in <https://github.com/brucebcampbell/bayesian-learning-with-R.git> we assessed convergence by; - viewing the time series for the intercept and each of the predictors. For this we utilized the coda package. - ran multiple chains and viewed evaluated the autocorrelation plots. - calculated the posterior means for the intercept and the β_j - utilized the mlv functions in the modeest to calculate the MAP estimated of the posterior modes - compared the 95% prediction intervals for the intercepts against the p-values from the logistic regression maximum likelihood model - Gelman plots are optionally produced when the number of MCMC chains is greater than one.

Some of the code is run conditionally through the DEBUG flag. We ran under debug mode and noted that all models converged. All but one of the predictors were the same for all the models. Depending on the run OTU_624 was swapped with another predictor in the uninformative model.

This was an interesting project. I iterated several versions and had to debug working jags models to get things running well. Also we encountered a NaN in the correlation due to one of the predictors being all 1's. If this was not accounted for the run times went up and convergence was bad. The width of the credible intervals could be investigated. We'll be running a longer simulation as a follow on task for fun. We set the precision of the model error to be 0.01 and 0.1 and got similar results.

Normal Uninformative

It's not specified what the prior variance is for $E[Y_j|X_j]$. We will assume $Y|\beta \sim N(y \cdot \beta, \sigma^2)$ where $\sigma^2 \sim \text{InvGamma}(0.1, 0.1)$

```

n <- nrow(X)

sigma.beta <- 100
inv.gamma.param <- 0.1
p <- ncol(X)

```

```

model_string.normal_uniformative <- "model{
  # Likelihood
  for(i in 1:n){
    Y[i] ~ dnorm(mu[i],inv.var)
    mu[i] <- intercept +inprod(X[i,],beta[])
  }

  # Prior for beta
  for(j in 1:p){
    beta[j] ~ dnorm(0,1/sigma.beta^2)
  }
  intercept ~ dnorm(0,1/sigma.beta^2)

  # Prior for the inverse variance
  inv.var ~ dgamma(inv.gamma.param, inv.gamma.param)
  sigma <- 1/sqrt(inv.var)
}"

model.normal_uniformative <- jags.model(textConnection(model_string.normal_uniformative), data

## Compiling model graph
##   Resolving undeclared variables
##   Allocating nodes
## Graph information:
##   Observed stochastic nodes: 1133
##   Unobserved stochastic nodes: 52
##   Total graph size: 61100
##
## Initializing model

update(model.normal_uniformative, nSamples, progress.bar="none"); # Burnin
samp.coeff.normal_uniformative <- coda.samples(model.normal_uniformative, variable.names=c("in

sum.normal_uniformative <- summary(samp.coeff.normal_uniformative)
quantiles<-sum.normal_uniformative$quantiles
left.05.quantile.sign <- sign(quantiles[,1])==-1
right.95.quantile.sign <- sign(quantiles[,5])==1
significant <- xor(left.05.quantile.sign ,right.95.quantile.sign)
beta.significant <- quantiles[significant,]

pander(data.frame(beta.significant), caption = "significant normal uninformativ ")

```

Table 1: significant normal uninformativ

	X2.5.	X25.	X50.	X75.	X97.5.
beta[7]	-8.013	-5.475	-4.233	-3.019	-0.7797

	X2.5.	X25.	X50.	X75.	X97.5.
beta[20]	-8.109	-5.743	-4.412	-3.207	-1.274
beta[25]	0.0936	2.131	3.152	4.213	6.159
beta[29]	1.708	4.029	5.23	6.456	8.821
beta[36]	0.2507	2.511	3.788	5.215	7.726
beta[37]	1.45	3.234	4.219	5.247	7.082
beta[38]	0.4348	2.667	3.974	5.283	7.807
beta[41]	2.064	4.31	5.491	6.676	8.913
beta[45]	1.977	4.3	5.535	6.753	9.057
beta[46]	2.951	5.029	6.14	7.178	9.132
beta[49]	1.82	3.893	4.914	6.02	8.139
beta[50]	3.591	5.838	7.06	8.167	10.33
intercept	51.44	53.22	54.13	55.05	56.62

```
credible.widths <- beta.significant[,5]-beta.significant[,1]

predictor.names.significant <- predictor.names[significant]

pander(data.frame(predictor.names.significant,credible.widths), caption = "credible widths normal uninformativ
```

Table 2: credible widths normal uninformativ

	predictor.names.significant	credible.widths
beta[7]	OTU_54646	7.234
beta[20]	OTU_9405	6.835
beta[25]	OTU_624	6.066
beta[29]	OTU_999	7.113
beta[36]	OTU_43955	7.476
beta[37]	OTU_66	5.632
beta[38]	OTU_51578	7.372
beta[41]	OTU_8086	6.848
beta[45]	OTU_72918	7.08
beta[46]	OTU_97	6.181
beta[49]	OTU_277	6.318
beta[50]	OTU_18758	6.735
intercept	intercept	5.18

```
if (DEBUG)
{
  autocorr.plot(samp.coeff.normal_uninformative)

  plot(samp.coeff.normal_uninformative)

  #Sample again and estimate posterior means and MAP posterior modes.
  samp.coeff.normal_uninformative.jags <- jags.samples(model.normal_uninformative, variable.names
```

```

posterior_means.normal_uniformative <- lapply(samp.coeff.normal_uniformative.jags, apply, 1,
pander(posterior_means.normal_uniformative, caption = "posterior means second sample")

posterior_modes.normal_uniformative <- lapply(samp.coeff.normal_uniformative.jags, apply, 1,
posterior_modes.normal_uniformative

if(n.chains>1)
{
gelman.plot(samp.coeff)
}
}

```

Hierarchical Normal Priors

$\beta_j | \tau \sim \text{Normal}(0, \tau^2)$ where $\tau^2 \sim \text{InvGamma}(0 : 01, 0 : 01)$

```

beta.inv.gamma.param <- 0.01
variance.inv.gamma.param <- 0.1
p <- ncol(X)

model_string.normal_hierarchical <- "model{
  # Likelihood
  for(i in 1:n){
    Y[i] ~ dnorm(mu[i],inv.var)
    mu[i] <- intercept +inprod(X[i,],beta[])
  }

  # Prior for beta
  for(j in 1:p){
    beta[j] ~ dnorm(0,beta.inv.gamma.param)
  }
  intercept ~ dnorm(0,beta.inv.gamma.param)

  # Prior for the inverse variance
  inv.var ~ dgamma(variance.inv.gamma.param, variance.inv.gamma.param)
  sigma <- 1/sqrt(inv.var)

  #Beta Prior for the inverse variance
  inv.var.beta ~ dgamma(beta.inv.gamma.param, beta.inv.gamma.param)
}"

model.normal_hierarchical <- jags.model(textConnection(model_string.normal_hierarchical), data

## Compiling model graph
## Resolving undeclared variables
## Allocating nodes
## Graph information:

```

```

## Observed stochastic nodes: 1133
## Unobserved stochastic nodes: 53
## Total graph size: 61098
##
## Initializing model
update(model.normal_hierarchical, nSamples, progress.bar="none"); # Burnin

samp.coeff.normal_hierarchical <- coda.samples(model.normal_hierarchical, variable.names=c("intercept", "beta[7]", "beta[20]", "beta[25]", "beta[29]", "beta[37]", "beta[41]", "beta[45]", "beta[46]", "beta[49]", "beta[50]"), n = 10000)

sum.normal_hierarchical <- summary(samp.coeff.normal_hierarchical)
quantiles<-sum.normal_hierarchical$quantiles
left.05.quantile.sign <- sign(quantiles[,1])==-1
right.95.quantile.sign <- sign(quantiles[,5])==1
significant <- xor(left.05.quantile.sign ,right.95.quantile.sign)
beta.significant <- quantiles[significant,]

pander(data.frame(beta.significant), caption = "significant normal hierarchical ")

```

Table 3: significant normal hierarchical

	X2.5.	X25.	X50.	X75.	X97.5.
beta[7]	-7.659	-5.371	-4.17	-2.928	-0.473
beta[20]	-7.474	-5.162	-3.851	-2.603	-0.1802
beta[25]	0.4377	2.17	3.105	4.16	6.244
beta[29]	1.564	3.927	5.116	6.326	8.52
beta[37]	1.235	3.096	4.176	5.158	7.152
beta[41]	1.683	4.026	5.346	6.575	9.098
beta[45]	1.873	4.328	5.439	6.624	8.853
beta[46]	2.913	5.116	6.207	7.365	9.338
beta[49]	1.689	3.74	4.817	5.892	7.653
beta[50]	3.114	5.389	6.628	7.865	10.14
intercept	51.09	52.64	53.46	54.22	56.09

```

credible.widths <- beta.significant[,5]-beta.significant[,1]

predictor.names.significant <- predictor.names[significant]

pander(data.frame(predictor.names.significant,credible.widths), caption = "credible widths normal hierarchical ")

```

Table 4: credible widths normal hierarchical

	predictor.names.significant	credible.widths
beta[7]	OTU_54646	7.186
beta[20]	OTU_9405	7.294
beta[25]	OTU_624	5.806
beta[29]	OTU_999	6.956

	predictor.names.significant	credible.widths
beta[37]	OTU_66	5.917
beta[41]	OTU_8086	7.415
beta[45]	OTU_72918	6.98
beta[46]	OTU_97	6.425
beta[49]	OTU_277	5.964
beta[50]	OTU_18758	7.023
intercept	intercept	5.005

```

if (DEBUG)
{
  autocorr.plot(samp.coeff.normal_hierarchical)

  plot(samp.coeff.normal_hierarchical)

  #Sample again and estimate posterior means and MAP posterior modes.
  samp.coeff.normal_hierarchical.jags <- jags.samples(model.normal_hierarchical, variable.names, 1, 10000)
  posterior_means.normal_hierarchical <- lapply(samp.coeff.normal_hierarchical.jags, apply, 1, FUN = function(x) {
    pandor(posterior_means.normal_hierarchical, caption = "posterior means second sample")
  })
  posterior_modes.normal_hierarchical <- lapply(samp.coeff.normal_hierarchical.jags, apply, 1, FUN = function(x) {
    posterior_modes.normal_hierarchical
  })
  if(n.chains>1)
  {
    gelman.plot(samp.coeff)
  }
}

```

BLASSO

```

beta.inv.gamma.param <- 0.01
variance.inv.gamma.param <- 0.1
p <- ncol(X)

model_string.normal_blasso <- "model{
  # Likelihood
  for(i in 1:n){
    Y[i] ~ dnorm(mu[i],inv.var)
    mu[i] <- intercept +inprod(X[i,],beta[])
  }

  # Prior for beta
  for(j in 1:p){
    beta[j] ~ ddexp(0,beta.inv.gamma.param)
  }
  intercept ~ ddexp(0,beta.inv.gamma.param)
}

```

```

# Prior for the inverse variance
inv.var ~ dgamma(variance.inv.gamma.param, variance.inv.gamma.param)
sigma   <- 1/sqrt(inv.var)

#Beta Prior for the inverse variance
inv.var.beta ~ dgamma(beta.inv.gamma.param, beta.inv.gamma.param)
}"

model.normal_blasso <- jags.model(textConnection(model_string.normal_blasso), data = list(Y=Y,X=X))

## Compiling model graph
##   Resolving undeclared variables
##   Allocating nodes
## Graph information:
##   Observed stochastic nodes: 1133
##   Unobserved stochastic nodes: 53
##   Total graph size: 61098
##
## Initializing model

update(model.normal_blasso, nSamples, progress.bar="none"); # Burnin

samp.coeff.normal_blasso <- coda.samples(model.normal_blasso,variable.names=c("intercept","beta"),n=10000)

sum.normal_blasso <- summary(samp.coeff.normal_blasso)
quantiles<-sum.normal_blasso$quantiles
left.05.quantile.sign <- sign(quantiles[,1])==-1
right.95.quantile.sign <- sign(quantiles[,5])==1
significant <- xor(left.05.quantile.sign ,right.95.quantile.sign)
beta.significant <- quantiles[significant,]

pander(data.frame(beta.significant), caption = "significant normal BLASSO ")

```

Table 5: significant normal BLASSO

	X2.5.	X25.	X50.	X75.	X97.5.
beta[7]	-7.754	-5.281	-3.955	-2.781	-0.4584
beta[20]	-8.102	-5.585	-4.27	-3.079	-0.8157
beta[29]	1.608	3.847	5.127	6.473	9.154
beta[37]	1.172	3.168	4.127	5.085	7.077
beta[38]	0.4863	2.8	4.144	5.373	7.789
beta[41]	1.743	4.223	5.442	6.752	9.431
beta[45]	1.387	4.077	5.44	6.604	8.941
beta[46]	3.01	5.172	6.225	7.383	9.476
beta[49]	1.543	3.749	4.833	5.875	7.894
beta[50]	3.429	5.862	6.964	8.139	10.36
intercept	51.89	53.34	54.18	54.98	56.58


```
credible.widths <- beta.significant[,5]-beta.significant[,1]

predictor.names.significant <- predictor.names[significant]

pander(data.frame(predictor.names.significant,credible.widths), caption = "credible widths normal BLASSO")
```

Table 6: credible widths normal BLASSO

	predictor.names.significant	credible.widths
beta[7]	OTU_54646	7.295
beta[20]	OTU_9405	7.287
beta[29]	OTU_999	7.546
beta[37]	OTU_66	5.905
beta[38]	OTU_51578	7.303
beta[41]	OTU_8086	7.688
beta[45]	OTU_72918	7.554
beta[46]	OTU_97	6.465
beta[49]	OTU_277	6.351
beta[50]	OTU_18758	6.936
intercept	intercept	4.697

```
if (DEBUG)
{
  autocorr.plot(samp.coeff.normal_blasso)

  plot(samp.coeff.normal_blasso)

  #Sample again and estimate posterior means and MAP posterior modes.
  samp.coeff.normal_blasso.jags <- jags.samples(model.normal_blasso, variable.names = c("intercept", "beta[7]", "beta[20]", "beta[29]", "beta[37]", "beta[38]", "beta[41]", "beta[45]", "beta[46]", "beta[49]", "beta[50]"), n.iter = 10000, thin = 10)
  posterior_means.normal_blasso <- lapply(samp.coeff.normal_blasso.jags, apply, 1, "mean")
  pander(posterior_means.normal_blasso, caption = "posterior means second sample")
  posterior_modes.normal_blasso <- lapply(samp.coeff.normal_blasso.jags, apply, 1, "mlv")
  posterior_modes.normal_blasso
  if(n.chains>1)
  {
    gelman.plot(samp.coeff)
  }
}
```