

Applied Bayesian Analysis : NCSU ST 540

Homework 6

Bruce Campbell

For this problem we use the 2016 election data described at <https://www4.stat.ncsu.edu/~reich/ABA/code/election2016data> with data provided in the R workspace https://www4.stat.ncsu.edu/~reich/ABA/code/election_2008_2016.RData

We restrict our analysis to the counties in North and South Carolina. In JAGS, we fit the logistic regression model

$$P(Z_i = 1) = \frac{1}{1 + e^{-\beta_0 - \sum_{j=1}^p X_{ij}\beta_j}}$$

where Z_i is the binary indicator that GOP support in county i increased by at least 5% from 2012 to 2016 $Z_i = 1 : Y_i > 5$ and $Z_i = 0$ otherwise, where Y is the change variable in the R workspace. X_{ij} are the covariates in the R workspace. We standardize each covariate to have mean zero and variance one before fitting the model. The priors are $\beta_j \sim N(0, \tau^2)$

(1) Fit the model with $\tau = 1$ and $\tau = 100$

First we make some notes on the data and the preprocessing steps.

`county_facts.csv`

<https://www.kaggle.com/benhamner/2016-us-election/data>

demographic data on counties from US census

3195 rows and 54 columns.

The metadata in the county facts table is located in a dictionary. We might need this for interpretation of the coefficients.

`county_facts_dictionary.csv`

description of the columns in `county_facts`

<https://www.kaggle.com/benhamner/2016-us-election/data>

The election data

`County_Election_08_16.csv`

<https://www4.stat.ncsu.edu/~reich/ABA/code/election2016data>

county-level voting patterns in the 2016 Presidential elections

3112 rows and 14 columns

The county data is joined to the election data via the key `fips_code`. We load the processed data below and start our analysis with the joined dataframe `all_dat`.

```

library(rjags)
library(coda)
library(choroplethr)
library(modeest)
load("election_2008_2016.RData")
carolinas <- all_dat[all_dat$state_abbreviation == "NC" | all_dat$state_abbreviation ==
  "SC", ]

gop.percent.increase <- 100 * (carolinas$gop_2016 - carolinas$gop_2012)/carolinas$gop_2012
gop.percent.increase.gt.5 <- gop.percent.increase > 5

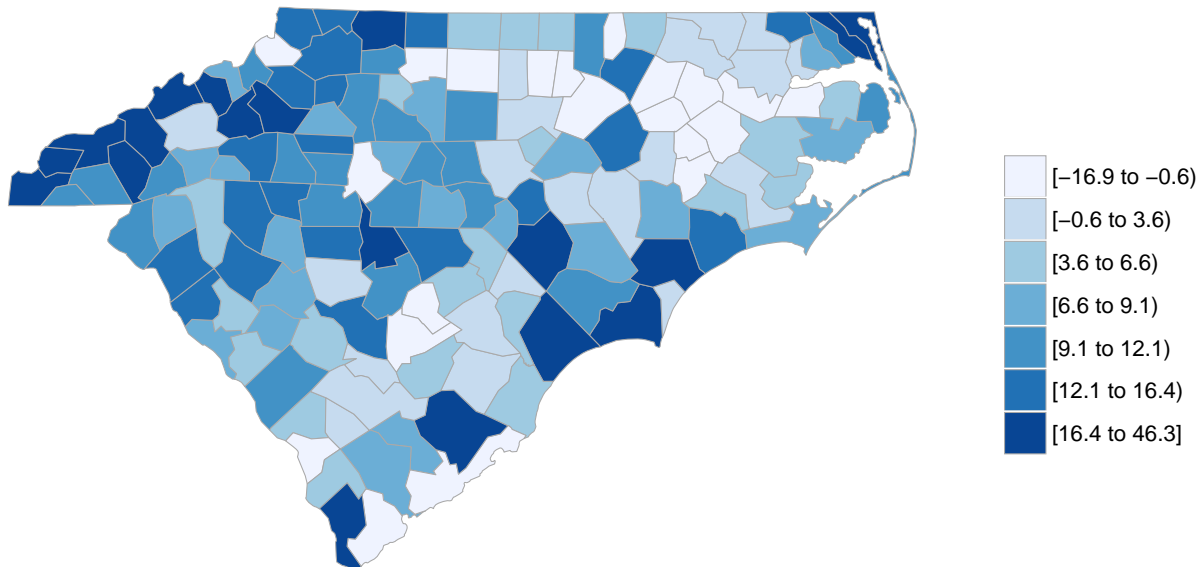
# Borrowed from instructor codebase to create our predictors
# and reposenses for the carolinas data set.
fips <- carolinas[, 1]
Y <- round(100 * (carolinas$gop_2016 - carolinas$gop_2012)/carolinas$gop_2012,
  1)
Z <- Y > 5
these <- c(3, 7, 10, 15, 20, 21, 25, 27, 31, 32, 47, 51, 22,
  44:45)
X <- as.matrix(carolinas[, these + 15])
X <- scale(X)
names <- dict[these, ]
colnames(X) <- names[, 1]

county_plot <- function(fips, Y, main = "", units = "") {
  temp <- as.data.frame(list(region = fips, value = Y))
  # county_choropleth(temp, title=main, legend=units)
  county_choropleth(temp, title = main, county_zoom = fips)
}

county_plot(fips, Y, "Percent change in GOP support from 2012 to 2016",
  unit = "Percent increase")

```

Percent change in GOP support from 2012 to 2016



Fit with $\tau = 1$

```
n.chains <- 1
DEBUG <- FALSE
if(DEBUG)
{
  nSamples <- 1000
} else
{
  nSamples <- 1000
}

n <- nrow(X)
tau <- 1
p <- ncol(X)

logistic_model <- "model{

  # Likelihood

  for(i in 1:n){
```

```

    Z[i] ~ dbern(q[i])
    logit(q[i]) <- intercept + inprod(X[i,], beta[])
  }

  #Priors
  intercept ~ dnorm(0, tau)
  for(j in 1:p){
    beta[j] ~ dnorm(0, tau)
  }
}

}"

model.carolinas <- jags.model(textConnection(logistic_model), data = list(Z=Z,X=X,n=n,p=p,tau=tau),
  ## Compiling model graph
  ##   Resolving undeclared variables
  ##   Allocating nodes
  ## Graph information:
  ##   Observed stochastic nodes: 146
  ##   Unobserved stochastic nodes: 16
  ##   Total graph size: 2941
  ##
  ## Initializing model

update(model.carolinas, nSamples, progress.bar="none"); # Burnin
samp.coeff <- coda.samples(model.carolinas, variable.names=c("intercept", "beta"), n.iter=2*nSamples)

```

Fit with $\tau = 100$

```

tau <- 100
model.carolinas.uninformative <- jags.model(textConnection(logistic_model), data = list(Z=Z,X=X,n=n,p=p,tau=tau),
  ## Compiling model graph
  ##   Resolving undeclared variables
  ##   Allocating nodes
  ## Graph information:
  ##   Observed stochastic nodes: 146
  ##   Unobserved stochastic nodes: 16
  ##   Total graph size: 2941
  ##
  ## Initializing model

update(model.carolinas.uninformative, nSamples, progress.bar="none"); # Burnin
samp.coeff.uninformative <- coda.samples(model.carolinas.uninformative, variable.names=c("intercept", "beta"), n.iter=2*nSamples)

```

(2) Assess convergence of the sampler for both priors.

In this section we sample from our model after burn in. Although all of the plots are not presented 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
- we fit a frequentist model and evaluated the estimated coefficients against the posterior means and modes
- compared the 95% prediction intervals for the intercepts against the p-values from the logistic regression maximum likelihood model

Code for this is below, we run some of it conditionally through the `DEBUG` variable. We did run the model without standardizing the feature data and noted evidence that the chain might be experiencing convergence issues. There was significant autocorrelation of the chains.

$\tau = 1$ Posterior quantiles

```
summary(samp.coeff)
```

```
##
## Iterations = 2001:4000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 2000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean      SD Naive SE Time-series SE
## beta[1]    1.05507 0.4858 0.010863      0.031639
## beta[2]   -0.64952 0.3651 0.008163      0.018698
## beta[3]   -1.86430 0.3855 0.008619      0.021261
## beta[4]   -0.77923 0.2874 0.006427      0.016252
## beta[5]   -0.39266 0.5017 0.011219      0.032259
## beta[6]   -1.36573 0.5723 0.012796      0.044015
## beta[7]    0.84493 0.4037 0.009026      0.023454
## beta[8]    0.21034 0.5578 0.012473      0.036191
## beta[9]   -0.45400 0.5967 0.013343      0.047776
## beta[10]  -0.27003 0.5167 0.011553      0.036933
## beta[11]  -0.26092 0.3221 0.007201      0.014539
## beta[12]  -0.25652 0.5869 0.013124      0.032762
## beta[13]   0.20429 0.5622 0.012570      0.034273
## beta[14]   0.08753 0.4904 0.010966      0.024892
```

```
## beta[15] -0.61064 0.5835 0.013046      0.030398
## intercept 0.98707 0.2695 0.006027      0.009892
##
## 2. Quantiles for each variable:
##
##          2.5%    25%    50%    75%    97.5%
## beta[1]    0.11255 0.7346 1.04860 1.37912 2.01529
## beta[2]   -1.36469 -0.8958 -0.65090 -0.39549 0.06931
## beta[3]   -2.59542 -2.1296 -1.86370 -1.59808 -1.08560
## beta[4]   -1.35705 -0.9683 -0.77279 -0.58071 -0.22738
## beta[5]   -1.36020 -0.7287 -0.40385 -0.06353 0.63279
## beta[6]   -2.54044 -1.7312 -1.36088 -1.00485 -0.16842
## beta[7]    0.05588 0.5688 0.84813 1.12150 1.66187
## beta[8]   -0.89388 -0.1664 0.21158 0.59269 1.33688
## beta[9]   -1.63167 -0.8530 -0.44859 -0.04644 0.72251
## beta[10]  -1.27853 -0.6124 -0.26767 0.06614 0.73710
## beta[11]  -0.88426 -0.4802 -0.26490 -0.03382 0.36890
## beta[12]  -1.42041 -0.6604 -0.24348 0.16024 0.80235
## beta[13]  -0.97779 -0.1558 0.23416 0.58858 1.25625
## beta[14]  -0.87262 -0.2492 0.08282 0.43493 1.01978
## beta[15]  -1.82770 -0.9838 -0.59772 -0.23522 0.52522
## intercept 0.49016 0.8027 0.97695 1.15919 1.53502
```

$\tau = 1$ Sample again and estimate the mean and MAP mode of the posterior distributions.

```
samp.coeff.jags <- jags.samples(model.carolinas, variable.names = c("intercept",
  "beta"), n.iter = nSamples, progress.bar = "none")
posterior_means <- lapply(samp.coeff.jags, apply, 1, "mean")
pander(posterior_means, caption = "posterior means second sample")
```

- **beta:** *1.058, -0.6486, -1.875, -0.7711, -0.4092, -1.304, 0.8278, 0.1211, -0.4353, -0.3133, -0.2622, -0.2803, 0.2038, 0.1023 and -0.656*
- **intercept:** *0.9845*

```
posterior_modes <- lapply(samp.coeff.jags, apply, 1, "mlv")
posterior_modes
```

```
## $beta
## $beta[[1]]
## Mode (most likely value): 1.004637
## Bickel's modal skewness: 0.072
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[2]]
## Mode (most likely value): -0.619813
## Bickel's modal skewness: -0.072
```

```

## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[3]]
## Mode (most likely value): -1.891976
## Bickel's modal skewness: 0.014
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[4]]
## Mode (most likely value): -0.762805
## Bickel's modal skewness: -0.038
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[5]]
## Mode (most likely value): -0.3594715
## Bickel's modal skewness: -0.108
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[6]]
## Mode (most likely value): -1.30435
## Bickel's modal skewness: 0.014
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[7]]
## Mode (most likely value): 0.8531151
## Bickel's modal skewness: -0.066
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[8]]
## Mode (most likely value): 0.1820475
## Bickel's modal skewness: -0.098
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[9]]
## Mode (most likely value): -0.4119206
## Bickel's modal skewness: -0.034
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[10]]
## Mode (most likely value): -0.2987245
## Bickel's modal skewness: -0.042
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[11]]
## Mode (most likely value): -0.242295
## Bickel's modal skewness: -0.064
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[12]]

```

```
## Mode (most likely value): -0.312061
## Bickel's modal skewness: 0.084
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[13]]
## Mode (most likely value): 0.2714668
## Bickel's modal skewness: -0.078
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[14]]
## Mode (most likely value): 0.1705194
## Bickel's modal skewness: -0.096
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[15]]
## Mode (most likely value): -0.4825561
## Bickel's modal skewness: -0.198
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
##
## $intercept
## $intercept[[1]]
## Mode (most likely value): 0.9417978
## Bickel's modal skewness: 0.098
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
```

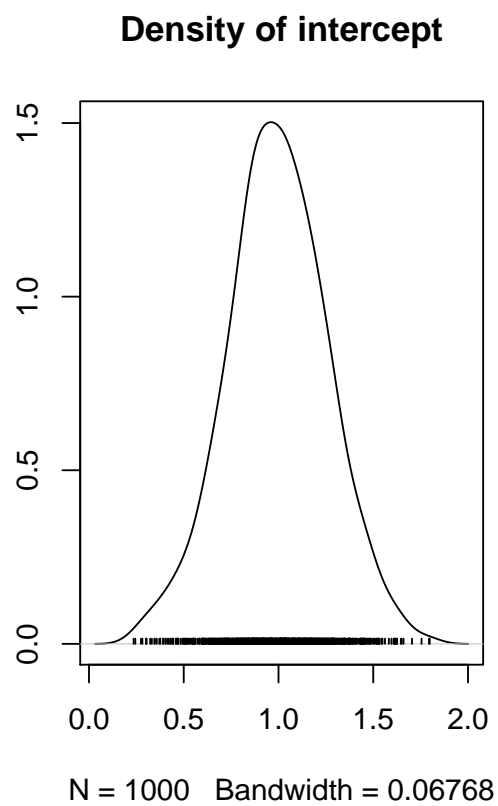
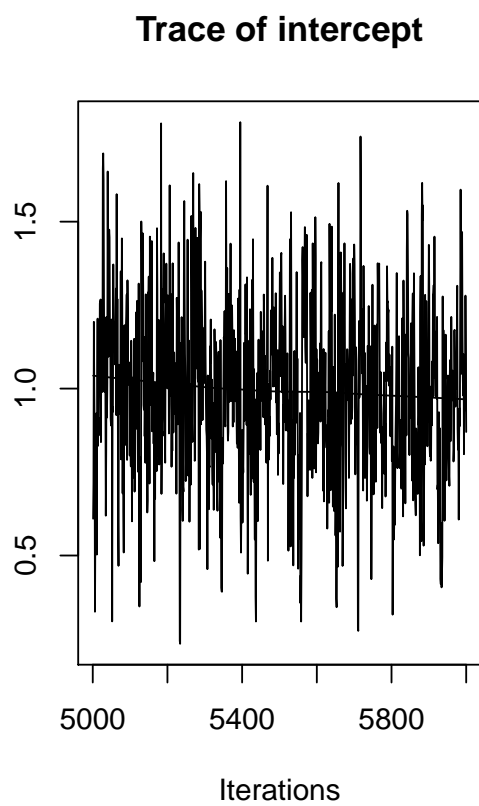
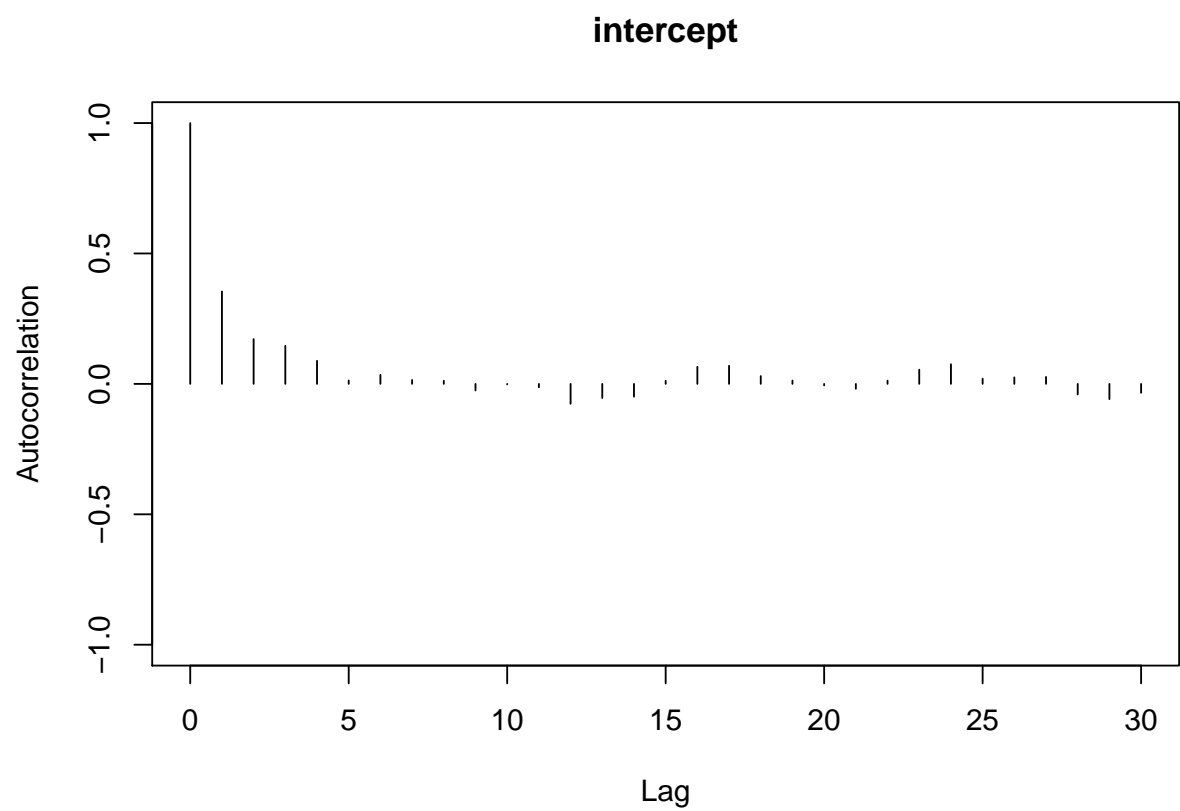
$\tau = 1$ Plot the time series, empirical posterior distribution, and the autocorrelation function for the coefficients

We only plot the intercept for the final report. Set the DEBUG flag to TRUE in order to include all of the coefficients.

```
if (DEBUG) {
  for (i in 1:p) {
    samp.coeff <- coda.samples(model.carolinas, variable.names = c(paste("beta[",
      i, "]", sep = "")), n.iter = nSamples, progress.bar = "none")
    autocorr.plot(samp.coeff)
    plot(samp.coeff)
  }
  samp.coeff <- coda.samples(model.carolinas, variable.names = "intercept",
    n.iter = nSamples, progress.bar = "none")
  autocorr.plot(samp.coeff)
  plot(samp.coeff)
} else {
  samp.coeff <- coda.samples(model.carolinas, variable.names = "intercept",
    n.iter = nSamples, progress.bar = "none")
```



```
    autocorr.plot(samp.coeff)  
    plot(samp.coeff)  
}
```



$\tau = 100$ Posterior quantiles

```
summary(samp.coeff.uninformative)
```

```
##
## Iterations = 2001:4000
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 2000
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##              Mean      SD Naive SE Time-series SE
## beta[1]    0.061827 0.09343 0.002089      0.003155
## beta[2]    0.008002 0.08793 0.001966      0.002499
## beta[3]   -0.225427 0.09297 0.002079      0.002818
## beta[4]   -0.026195 0.08581 0.001919      0.002468
## beta[5]   -0.017534 0.09352 0.002091      0.003146
## beta[6]   -0.099206 0.09338 0.002088      0.003167
## beta[7]    0.175354 0.09130 0.002042      0.002654
## beta[8]   -0.013477 0.09435 0.002110      0.002833
## beta[9]    0.005270 0.09618 0.002151      0.002937
## beta[10]  -0.084892 0.08955 0.002002      0.002668
## beta[11]  -0.061597 0.08979 0.002008      0.002682
## beta[12]  -0.079881 0.09176 0.002052      0.002664
## beta[13]  -0.050911 0.09330 0.002086      0.003101
## beta[14]  -0.057537 0.09177 0.002052      0.002706
## beta[15]  -0.079683 0.09525 0.002130      0.002818
## intercept 0.153714 0.08847 0.001978      0.002627
##
## 2. Quantiles for each variable:
##
##              2.5%      25%      50%      75%      97.5%
## beta[1]   -0.118361 0.001564 0.059059 0.122377 0.24538
## beta[2]   -0.160020 -0.052301 0.006501 0.067780 0.17801
## beta[3]   -0.408350 -0.288623 -0.225182 -0.160749 -0.04535
## beta[4]   -0.191232 -0.083567 -0.028848 0.031315 0.14231
## beta[5]   -0.195967 -0.079941 -0.020729 0.047207 0.16540
## beta[6]   -0.272659 -0.162113 -0.100087 -0.033563 0.07969
## beta[7]    0.005217 0.111713 0.174538 0.238808 0.35726
## beta[8]   -0.200976 -0.073993 -0.013263 0.046368 0.17550
## beta[9]   -0.182505 -0.060202 0.004981 0.069165 0.19881
## beta[10]  -0.256676 -0.143351 -0.084511 -0.027995 0.09666
## beta[11]  -0.234842 -0.121763 -0.060659 -0.001192 0.11889
## beta[12]  -0.258739 -0.145157 -0.078109 -0.015326 0.09561
## beta[13]  -0.237169 -0.112525 -0.051503 0.011133 0.13422
```

```
## beta[14] -0.240131 -0.118330 -0.055760 0.002551 0.12974
## beta[15] -0.271583 -0.140421 -0.079406 -0.013979 0.10178
## intercept -0.020493 0.092947 0.151799 0.214466 0.32459
```

$\tau = 100$ Sample again and estimate the mean and MAP mode of the posterior distributions.

```
samp.coeff.jags.uninformative <- jags.samples(model.carolinas.uninformative,
  variable.names = c("intercept", "beta"), n.iter = 2 * nSamples,
  progress.bar = "none")
posterior_means.uninformative <- lapply(samp.coeff.jags.uninformative,
  apply, 1, "mean")
pander(posterior_means.uninformative, caption = "posterior_means.uninformative")
```

- **beta:** 0.05786, 0.007732, -0.2241, -0.03007, -0.02445, -0.1013, 0.1788, 0.0004141, 0.003337, -0.08039, -0.06151, -0.07349, -0.04971, -0.05554 and -0.08523
- **intercept:** 0.1495

```
posterior_means.uninformative <- lapply(samp.coeff.jags.uninformative,
  apply, 1, "mlv")
posterior_means.uninformative
```

```
## $beta
## $beta[[1]]
## Mode (most likely value): 0.06074651
## Bickel's modal skewness: -0.019
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[2]]
## Mode (most likely value): 0.009226344
## Bickel's modal skewness: 0.001
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[3]]
## Mode (most likely value): -0.2244602
## Bickel's modal skewness: 0.016
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[4]]
## Mode (most likely value): -0.02350581
## Bickel's modal skewness: -0.049
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[5]]
## Mode (most likely value): -0.02151864
## Bickel's modal skewness: -0.041
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
```

```

##
## $beta[[6]]
## Mode (most likely value): -0.1114117
## Bickel's modal skewness: 0.086
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[7]]
## Mode (most likely value): 0.1851249
## Bickel's modal skewness: -0.027
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[8]]
## Mode (most likely value): -0.005341771
## Bickel's modal skewness: 0.039
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[9]]
## Mode (most likely value): -0.002769133
## Bickel's modal skewness: 0.052
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[10]]
## Mode (most likely value): -0.08876738
## Bickel's modal skewness: 0.056
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[11]]
## Mode (most likely value): -0.0672086
## Bickel's modal skewness: 0.062
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[12]]
## Mode (most likely value): -0.07706274
## Bickel's modal skewness: 0.039
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[13]]
## Mode (most likely value): -0.04893353
## Bickel's modal skewness: -0.02
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[14]]
## Mode (most likely value): -0.05963779
## Bickel's modal skewness: 0.03
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
## $beta[[15]]
## Mode (most likely value): -0.08004244

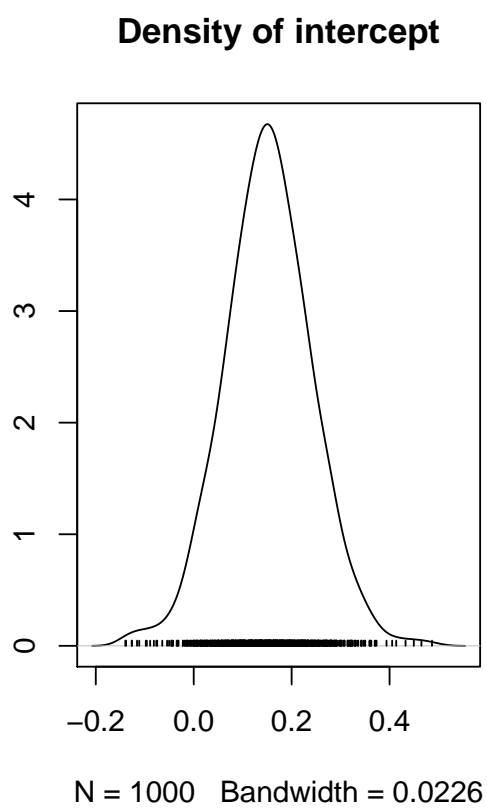
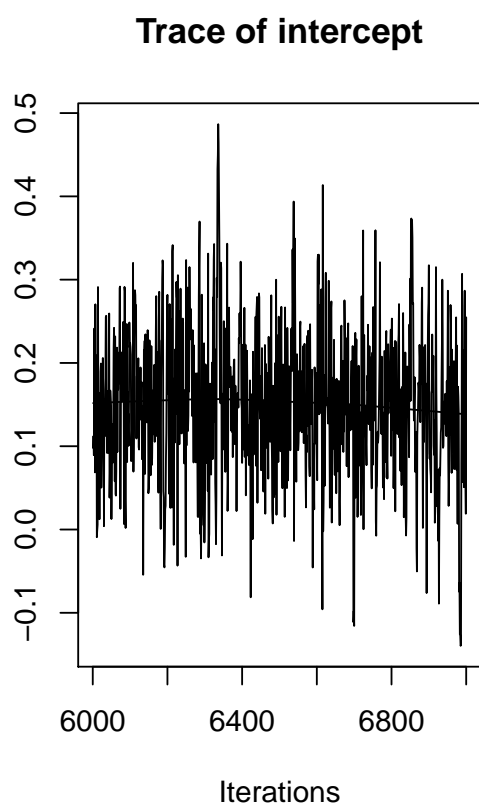
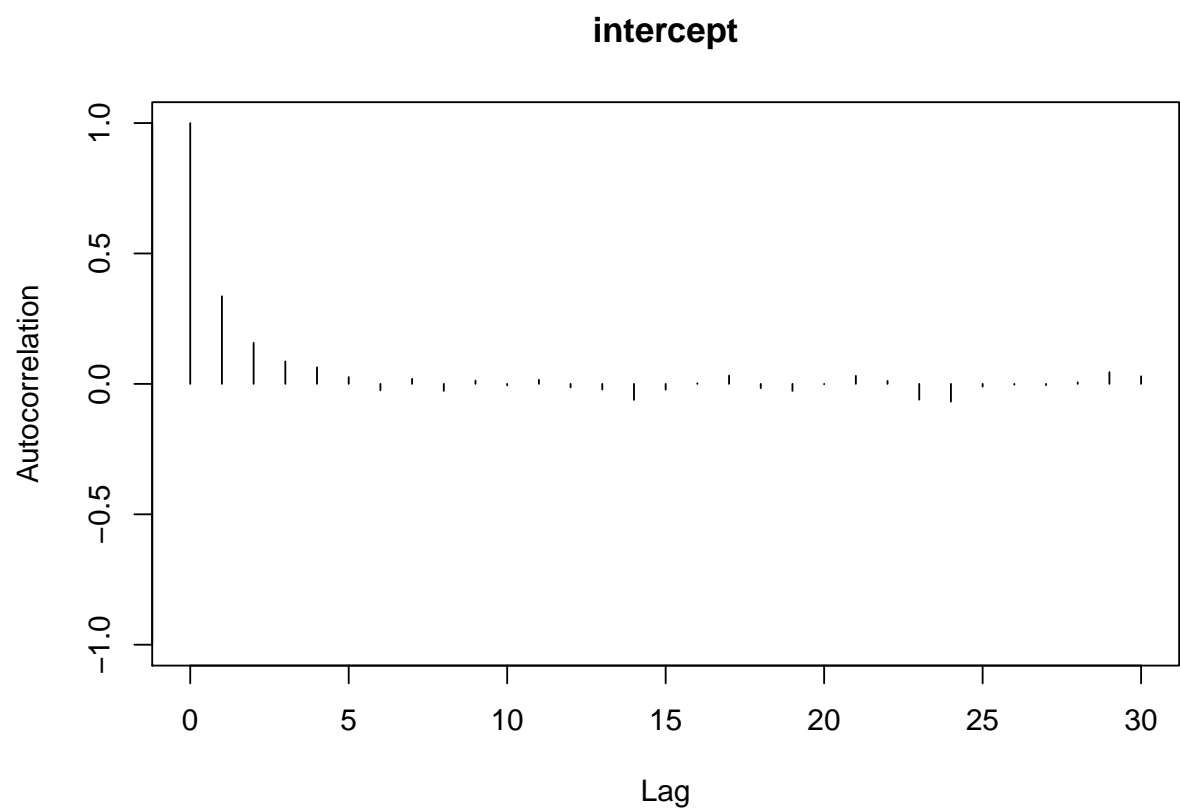
```

```
## Bickel's modal skewness: -0.026
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
##
##
## $intercept
## $intercept[[1]]
## Mode (most likely value): 0.1523908
## Bickel's modal skewness: -0.019
## Call: mlv.default(x = array(newX[, i], d.call, dn.call))
```

$\tau = 1$ Plot the time series, empirical posterior distribution, and the autocorrelation function for the coefficients

We only plot the intercept for the final report. Set the DEBUG flag to TRUE in order to include all of the coefficients.

```
if (DEBUG) {
  for (i in 1:p) {
    samp.coeff.uninformative <- coda.samples(model.carolinas.uninformative,
      variable.names = c(paste("beta[", i, "]", sep = "")),
      n.iter = nSamples, progress.bar = "none")
    autocorr.plot(samp.coeff.uninformative)
    plot(samp.coeff.uninformative)
  }
  samp.coeff.uninformative <- coda.samples(model.carolinas.uninformative,
    variable.names = "intercept", n.iter = nSamples, progress.bar = "none")
  autocorr.plot(samp.coeff.uninformative)
  plot(samp.coeff.uninformative)
} else {
  samp.coeff.uninformative <- coda.samples(model.carolinas.uninformative,
    variable.names = "intercept", n.iter = nSamples, progress.bar = "none")
  autocorr.plot(samp.coeff.uninformative)
  plot(samp.coeff.uninformative)
}
```



Fit frequentist logistic model for reference.

```
df <- data.frame(cbind(Z, X))
lm.logistic <- glm(Z ~ ., family = binomial, df)
summary(lm.logistic)

##
## Call:
## glm(formula = Z ~ ., family = binomial, data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.1986  -0.4521   0.1467   0.4048   3.0622
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  1.22344    0.33351   3.668 0.000244 ***
## PST120214    1.44242    0.68408   2.109 0.034984 *
## AGE775214   -0.81589    0.46258  -1.764 0.077771 .
## RHI225214   -2.05358    0.48616  -4.224 2.4e-05 ***
## RHI725214   -1.00849    0.37201  -2.711 0.006709 **
## EDU635213   -0.55048    0.64123  -0.858 0.390634
## EDU685213   -1.99593    0.88355  -2.259 0.023883 *
## HSG445213    0.96082    0.50725   1.894 0.058199 .
## HSG495213    0.87243    0.81258   1.074 0.282975
## INC110213   -0.97476    0.91641  -1.064 0.287475
## PVY020213   -0.36784    0.64185  -0.573 0.566582
## RTN131207   -0.44836    0.41945  -1.069 0.285097
## POP060210   -0.03261    0.76264  -0.043 0.965888
## VET605213    0.46278    0.81324   0.569 0.569321
## MAN450207    0.46521    0.69222   0.672 0.501550
## WTN220207   -1.01084    0.93321  -1.083 0.278726
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 190.144  on 145  degrees of freedom
## Residual deviance:  93.976  on 130  degrees of freedom
## AIC: 125.98
##
## Number of Fisher Scoring iterations: 6
```