

# HW6

Andrew Liang

11/4/2020

```
library(MASS)
library(matlib)
library(tidyverse)
```

```
## -- Attaching packages -----

## v ggplot2 3.3.2      v purrr  0.3.4
## v tibble  3.0.3      v dplyr  1.0.2
## v tidyr   1.1.2      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0

## -- Conflicts -----
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()     masks stats::lag()
## x dplyr::select() masks MASS::select()
```

## 7.4

**a**

Intuitively, most married couples are probably between the ages of 25 to 80. I will set  $\theta_0 = (\frac{25+80}{2}, \frac{25+80}{2}) = (52.5, 52.5)^T$ .

Consequently, I will pick a semiconjugate prior distribution. It seems more likely that there will be more married couples around age 50, justifying a bell curve centered around 52.5 with variance  $13.75^2 = 189$  such that 95% of my prior is within (25, 80)

Additionally, I believe that ages of couples are tightly correlated, and thus aim for a prior correlation of around 0.75. So the covariance between the two variables would be:

$$0.75 = \frac{\sigma_{1,2}}{189}$$
$$\sigma_{1,2} = 141.75$$

Thus the covariance matrix will be:

$$\Sigma_0 = \begin{bmatrix} 189 & 141.75 \\ 141.75 & 189 \end{bmatrix}$$

```

agehw <- as.data.frame(read.csv('agehw.dat', sep = ""))

Y <- agehw
n <- nrow(agehw)
p <- ncol(agehw)

mu0 <- rep(52.5, p)
lambda0 <- s0 <- rbind(c(189, 141.75), c(141.75, 189))

#nu0 = p+2
nu0 <- p + 2 + 10

```

**b**

```

N <- 100
S <- 12

set.seed(1651)

Y_preds <- lapply(1:S, function(s) {

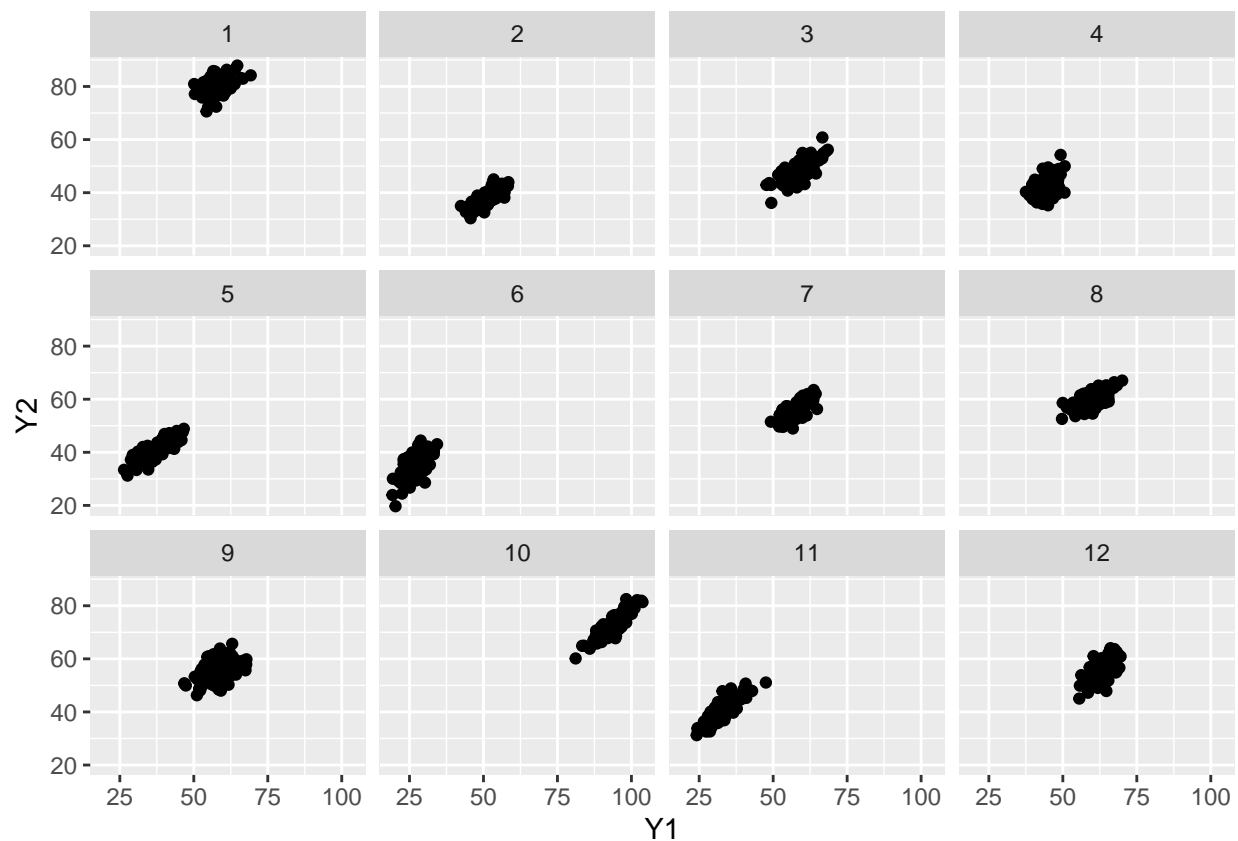
  # Sample THETA according using prior
  theta <- mvrnorm(n = 1, mu0, lambda0)
  sigma <- inv(rWishart(1, nu0, inv(s0))[, , 1])
  Y_s <- mvrnorm(n = 100, theta, sigma)
  data.frame(Y1 = Y_s[, 1], Y2 = Y_s[, 2], dataset = s)

})

Y_comb <- do.call(rbind, Y_preds)

ggplot(Y_comb, aes(x = Y1, y = Y2)) +
  geom_point() +
  facet_wrap(~ dataset)

```



## c

S <- 10000

set.seed(1651)

*#monte carlo priors*

do\_mcmc <- function(Y, mu0, lambda0, s0, nu0) {

  ybar <- colMeans(Y)

  p <- ncol(Y)

  n <- nrow(Y)

  THETA <- matrix(nrow = S, ncol = p)

  SIGMA <- array(dim = c(p, p, S))

*#start with sigma sample*

  sigma <- cov(Y)

*#gibbs sampling*

*# more readable*

  inv <- solve

  for (s in 1:S) {

*#update theta*

    lambdan <- inv(inv(lambda0) + n \* inv(sigma))

    mun <- lambdan %\*% (inv(lambda0) %\*% mu0 + n \* inv(sigma) %\*% ybar)

    theta <- mvrnorm(n = 1, mun, lambdan)

```

    #update sigma
    resid <- t(Y) - c(theta)
    stheta <- resid %*% t(resid)
    sn <- s0 + stheta
    sigma <- inv(rWishart(1, nu0 + n, inv(sn))[, , 1])

    THETA[s, ] <- theta
    SIGMA[, , s] <- sigma
  }

  list(theta = THETA, sigma = SIGMA)
}

my_prior_mcmc <- do_mcmc(agehw, mu0, lambda0, s0, nu0)
THETA <- my_prior_mcmc$theta
SIGMA <- my_prior_mcmc$sigma

#for reuse later
print_quantiles <- function(THETA, SIGMA) {

  #husband
  print("Husband")
  print(quantile(THETA[, 1], probs = c(0.025, 0.5, 0.975)))

  #wife
  print("Wife")
  print(quantile(THETA[, 2], probs = c(0.025, 0.5, 0.975)))

  cors <- apply(SIGMA, MARGIN = 3, FUN = function(covmat) {
    covmat[1, 2] / (sqrt(covmat[1, 1] * covmat[2, 2]))
  })
  print("Correlation")
  print(quantile(cors, probs = c(0.025, 0.5, 0.975)))
}

print_quantiles(THETA, SIGMA)

```

```

## [1] "Husband"
##      2.5%      50%      97.5%
## 41.95711 44.52705 47.10299
## [1] "Wife"
##      2.5%      50%      97.5%
## 38.60235 40.99225 43.39706
## [1] "Correlation"
##      2.5%      50%      97.5%
## 0.8617723 0.9024292 0.9313893

```

d part iii

```
set.seed(1651)
```

```

mu0 <- rep(0, p)
lambda0 <- 10^5 * diag(p)
s0 <- 1000 * diag(p)
nu0 <- 3
diffuse_mcmc <- do_mcmc(agehw, mu0, lambda0, s0, nu0)
print_quantiles(diffuse_mcmc$theta, diffuse_mcmc$sigma)

```

```

## [1] "Husband"
##      2.5%      50%      97.5%
## 41.67009 44.45815 47.24395
## [1] "Wife"
##      2.5%      50%      97.5%
## 38.31861 40.91042 43.51476
## [1] "Correlation"
##      2.5%      50%      97.5%
## 0.7927651 0.8548633 0.8995693

```

e

Comparing the confidence intervals, it doesn't seem that prior info matters since the sample size is large. The diffuse prior is slightly different, but not significantly at all. Regardless of the prior info, the quantiles and correlations being quite similar. A smaller sample size may lead to different results.

```

set.seed(1651)
#USING MY PRIOR

mu0 <- rep(52.5, p)
lambda0 <- s0 <- rbind(c(189, 141.75), c(141.75, 189))

nu0 <- p + 2 + 10
my_prior_mcmc_short <- do_mcmc(agehw[1:25, ], mu0, lambda0, s0, nu0)
print_quantiles(my_prior_mcmc_short$theta, my_prior_mcmc_short$sigma)

```

```

## [1] "Husband"
##      2.5%      50%      97.5%
## 41.01928 45.45946 49.99407
## [1] "Wife"
##      2.5%      50%      97.5%
## 38.42809 43.12462 47.75715
## [1] "Correlation"
##      2.5%      50%      97.5%
## 0.8388558 0.9136948 0.9545245

```

```

#DIFFUSE PRIOR
mu0 <- rep(0, p)
lambda0 <- 10^5 * diag(p)
s0 <- 1000 * diag(p)
nu0 <- 3
diffuse_mcmc_short <- do_mcmc(agehw[1:25, ], mu0, lambda0, s0, nu0)
print_quantiles(diffuse_mcmc_short$theta, diffuse_mcmc_short$sigma)

```

```
## [1] "Husband"
##      2.5%      50%     97.5%
## 39.28346 45.19486 51.26876
## [1] "Wife"
##      2.5%      50%     97.5%
## 36.69468 42.82984 49.01497
## [1] "Correlation"
##      2.5%      50%     97.5%
## 0.5453413 0.7618878 0.8808638
```

By reducing the sample size, we can see a noticeable effect on the intervals.

## 7.6

```
diab <- read.csv('azdiabetes.dat', sep = ",")
diaby <- filter(diab, diabetes == "Yes"); diaby <- diaby[,-8] #yes diabetes
diabn <- filter(diab, diabetes == "No"); diabn <- diabn[,-8] # no diabetes

#prior info
mu0y <- colMeans(diaby)
mu0n <- colMeans(diabn)
p <- ncol(diaby)
v0 <- p + 2
covary <- var(diaby)
covarn <- var(diabn)
```

**a**