

# BDA A2

Piero Birello

March 2025

## Exercise A

### Task 1

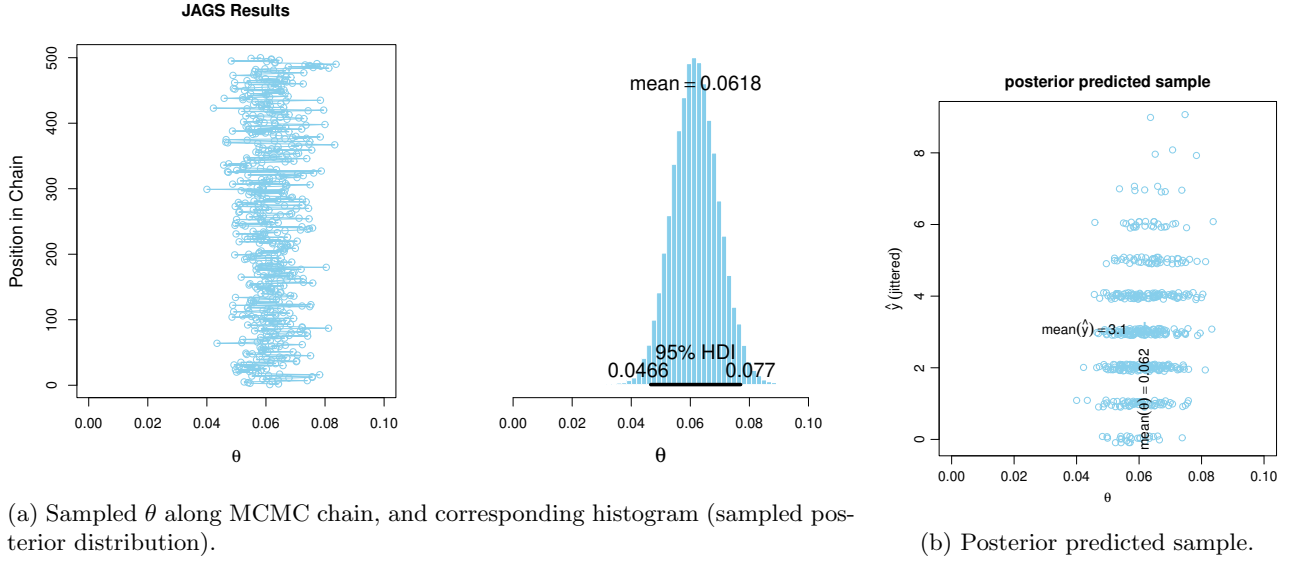
The modified code is as follows:

```
1 rm(list = ls())
2 graphics.off()
3 source("openGraphSaveGraph.R")
4 save_path = "Figures/"
5
6 require(rjags)          # Kruschke, J. K. (2011). Doing Bayesian Data Analysis:
7 # A Tutorial with R and BUGS. Academic Press / Elsevier.
8 #-----
9 # THE MODEL.
10
11 # Specify the model in JAGS language, but save it as a string in R:
12 modelString = "
13 # JAGS model specification begins ...
14 model {
15 # Likelihood:
16 for ( i in 1:nStores ) {
17 y[i] ~ dbin( theta,50 )
18 }
19 # Prior distribution:
20 theta ~ dnorm( 0.05,1/(0.01)^2 )T( 0.0,0.1 )
21 # theta ~ dnorm( 0.1,1/(0.02)^2 )T( 0.05,0.15 )
22 }
23 # ... JAGS model specification ends.
24 " # close quote to end modelString
25
26 # Write the modelString to a file, using R commands:
27 writelines(modelString,con="model.txt")
28
29 #-----
30 # THE DATA.
31
32 # Specify the data in R, using a list format compatible with JAGS:
33 dataList = list(
34   nStores = 5 ,
35   y = c( 7,5,4,4,2 )
36 )
37
38 #-----
39 # INITIALIZE THE CHAIN.
40
41 # Can be done automatically in jags.model() by commenting out inits argument.
```

```

42 # Otherwise could be established as:
43 # initsList = list( theta = sum(dataList$y)/length(dataList$y) )
44
45 #-----
46 # RUN THE CHAINS.
47
48 parameters = c( "theta" )      # The parameter(s) to be monitored.
49 adaptSteps = 500                # Number of steps to "tune" the samplers.
50 burnInSteps = 1000             # Number of steps to "burn-in" the samplers.
51 nChains = 3                    # Number of chains to run.
52 numSavedSteps=50000            # Total number of steps in chains to save.
53 thinSteps=1                    # Number of steps to "thin" (1=keep every step).
54 nIter = ceiling( ( numSavedSteps * thinSteps ) / nChains ) # Steps per chain.
55 # Create, initialize, and adapt the model:
56 jagsModel = jags.model( "model.txt" , data=dataList , # inits=initsList ,
57                        n.chains=nChains , n.adapt=adaptSteps )
58 # Burn-in:
59 cat( "Burning in the MCMC chain...\n" )
60 update( jagsModel , n.iter=burnInSteps )
61 # The saved MCMC chain:
62 cat( "Sampling final MCMC chain...\n" )
63 codaSamples = coda.samples( jagsModel , variable.names=parameters ,
64                             n.iter=nIter , thin=thinSteps )
65 # resulting codaSamples object has these indices:
66 #   codaSamples[[ chainIdx ]][ stepIdx , paramIdx ]
67
68
69 #-----
70 # EXAMINE THE RESULTS.
71
72 # Convert coda-object codaSamples to matrix object for easier handling.
73 # But note that this concatenates the different chains into one long chain.
74 # Result is mcmcChain[ stepIdx , paramIdx ]
75 mcmcChain = as.matrix( codaSamples )
76
77 thetaSample = mcmcChain
78
79 # Make a graph using R commands:
80 trunc_low = 0.0
81 trunc_high = 0.1
82 openGraph(width=10,height=6)
83 layout( matrix( c(1,2) , nrow=1 ) )
84 plot( thetaSample[1:500] , 1:length(thetaSample[1:500]) , type="o" ,
85       xlim=c(trunc_low,trunc_high) , xlab=bquote(theta) , ylab="Position in Chain" ,
86       cex.lab=1.25 , main="JAGS Results" , col="skyblue" )
87 source("plotPost.R")
88 histInfo = plotPost( thetaSample , xlim=c(trunc_low,trunc_high) , xlab=bquote(theta) )
89 saveGraph( file=paste0(save_path,"BernBetaJagsFullPost") , type="eps" )
90
91 # Posterior prediction:
92 # For each step in the chain, use posterior theta to flip a coin:
93 chainLength = length( thetaSample )
94 yPred = rep( NULL , chainLength ) # define placeholder for flip results
95 for ( stepIdx in 1:chainLength ) {
96   pHead = thetaSample[stepIdx]
97   #yPred[stepIdx] = sample( x=c(0,1), prob=c(1-pHead,pHead), size=1 )
98   yPred[stepIdx] = rbinom( n=1,size=50,prob=pHead )
99 }
100 # Jitter the y values for plotting purposes:
101 yPredJittered = yPred + runif( length(yPred) , -.1 , +.1 )
102 # Now plot the jittered values:
103 openGraph(width=5,height=5.5)

```



```

104 par( mar=c(3.5,3.5,2.5,1) , mgp=c(2,0.7,0) )
105 plot( thetaSample[1:500] , yPredJittered[1:500] , xlim=c(trunc_low,trunc_high) ,
106       main="posterior predicted sample" ,
107       xlab=expression(theta) , ylab=expression(hat(y)*" *(jittered))" ,
108       col="skyblue" )
109 points( mean(thetaSample) , mean(yPred) , pch="+" , cex=2 , col="skyblue" )
110 text( mean(thetaSample) , mean(yPred) ,
111       bquote( mean(hat(y)) == .(signif(mean(yPred),2)) ) ,
112       adj=c(1.2,.5) )
113 text( mean(thetaSample) , mean(yPred) , srt=90 ,
114       bquote( mean(theta) == .(signif(mean(thetaSample),2)) ) ,
115       adj=c(1.2,.5) )
116 #abline( trunc_low , trunc_high , lty="dashed" , lwd=2 )
117 saveGraph( file=paste0(save_path,"BernBetaJagsFullPred") , type="eps" )

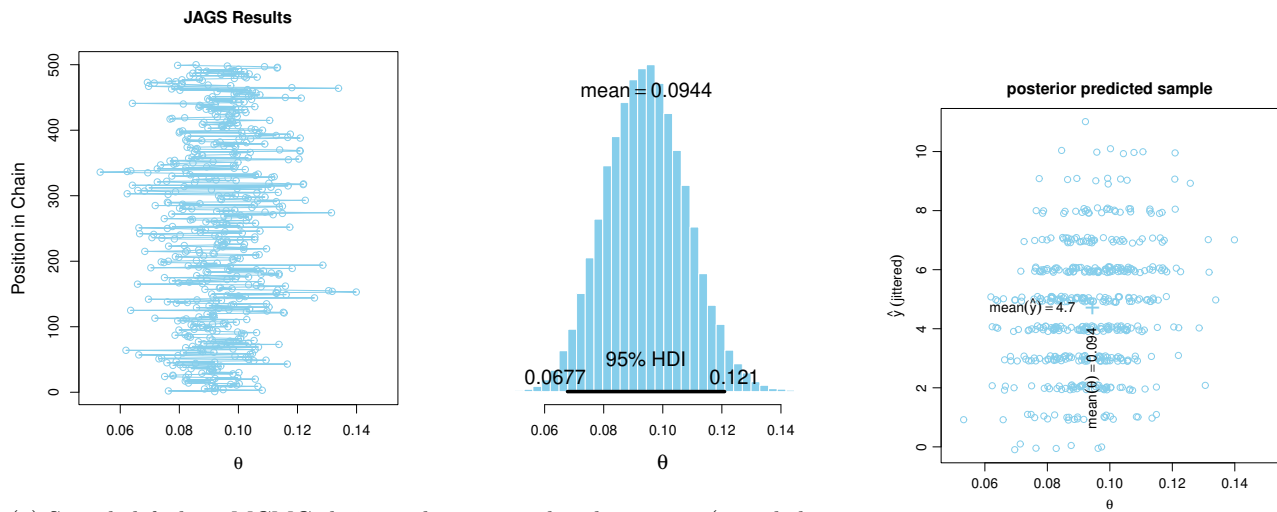
```

Obtained results are shown in Figure 1a, Figure 1b. Based on the histogram in Figure 1a, we may say that, qualitatively, we suspend belief about the statement that the panels are produced by company A. The posterior distribution for  $\theta$  (the probability of failure) qualitatively matches a truncated normal distribution with mean 0.05 and standard deviation 0.01, and the expected value of the mean of  $\theta$  is contained within the HDI. However, the HDI would not be contained in a ROPE around the expected value of  $\theta$ .

In fact, a more precise assessment would require to evaluate the posterior distributions for the parameters of the *dnorm* themselves, and check whether the corresponding HDI intervals contain a ROPE defined around the declared values. Otherwise, we could adopt a model comparison approach and compare the results under the current model with those under an alternative one.

## Task 2

We now run the model under a different prior, namely  $\theta \sim \text{dnorm}(0.1, 1/(0.02)^2)T(0.05, 0.15)$ . Results are shown in Figure 2a, Figure 2b. Again, we may say that the obtained posterior distribution matches qualitatively the prior. Qualitatively, we suspect that this model is more accurate than the previous one. Figure 2b also shows a posterior predicted sample that is more similar to the one at disposal.



(a) Sampled  $\theta$  along MCMC chain, and corresponding histogram (sampled posterior distribution).

(b) Posterior predicted sample.

### Task 3

Considering a ten times larger sample reporting the same fractions of failures, we have more evidence in support of factory B than factory A, see Figure 3a, Figure 3b, Figure 4a, Figure 4b. The HDI for  $\theta$  given prior A does not contain the expected value for  $\theta$ , while the HDI given prior B does -it ends right there. Still, a more accurate analysis would require to evaluate the mean and standard deviation parameters of the *dnorm*, or perform model comparison.

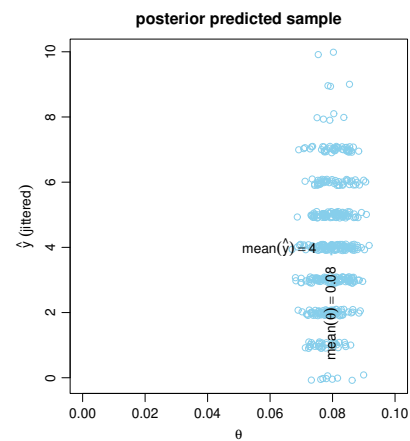
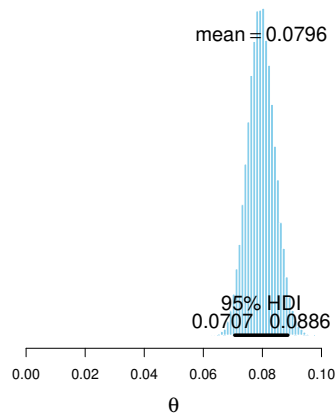
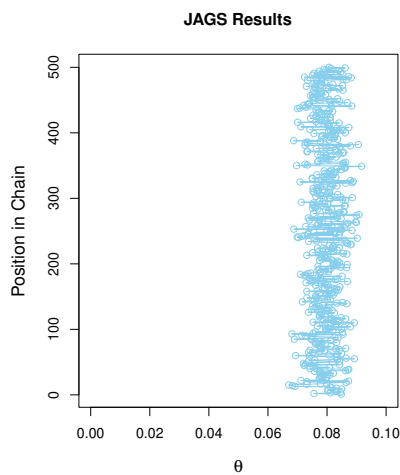
### Task 4

We add the following piece of code:

```

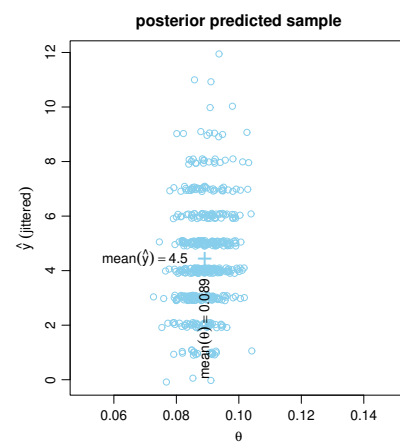
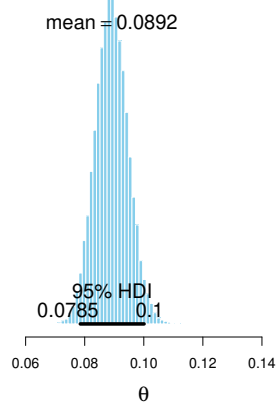
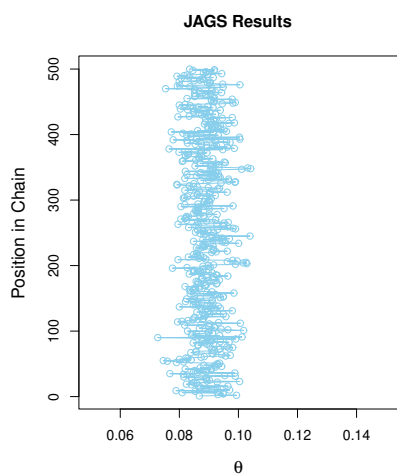
1 trunc_low = 0.05
2 trunc_high = 0.15
3 dataN = 50*10
4 mu = 0.1
5 sigma = 0.02
6
7 # COMPUTE PROBABILITY OF DATA
8
9 # Define truncated normal
10 truncNorm = function(x,mean,sd,a,b) {
11   pdfValues = dnorm(x, mean = mean, sd = sd) /
12   (pnorm(b, mean = mean, sd = sd) - pnorm(a, mean = mean, sd = sd))
13   return(pdfValues)
14 }
15
16 # Compute mean and standard deviation of MCMC values:
17 meanTheta = mean(thetaSample)
18 sdTheta = sd(thetaSample)
19
20 # Compute 1/p(D):
21 logOneOverPD = mean( log( truncNorm( thetaSample,meanTheta,sdTheta,trunc_low,trunc_high ) ) -
   # h(theta_i)

```



(a) Sampled  $\theta$  along MCMC chain, and corresponding histogram (sampled posterior distribution).

(b) Posterior predicted sample.



(a) Sampled  $\theta$  along MCMC chain, and corresponding histogram (sampled posterior distribution).

(b) Posterior predicted sample.

```

22      ( rowSums( dbinom( dataList$y,dataN,thetaSample, log=TRUE ) ) +      #
23              p(D|theta_i)
24      log( truncNorm( thetaSample,mu,sigma,trunc_low,trunc_high ) ) ) )      # p(
25      theta_i)
logPD = -(logOneOverPD)
show(logPD)

```

We run the script under the two prior distributions for  $\theta$ , then compute the log Bayes Factor as:

$$\log BF = \log p(D|m = A) - \log p(D|m = B) = -11.20936 - (-7.179099) = -4.030261 \quad (1)$$

Hence, as expected, model comparison lets us select the hypothesis that panels are produced by company B.

## Exercise B

The model can be written as:

$$y_j \sim \text{Binomial}(p_c^{machines_j}) \quad (2)$$

$$p_c^i \sim \text{trunc.Normal}(\mu, \sigma, a, b) \quad (3)$$

$$\mu \sim \text{Beta}(1.1, 9.9) \quad (4)$$

$$\sigma \sim \text{Gamma}(2, 28) \quad (5)$$

$$a = 0.0 \quad (6)$$

$$b = 0.5 \quad (7)$$

It can be implemented in JAGS as:

```

1  # THE MODEL.
2
3  # Specify the model in JAGS language, but save it as a string in R:
4  modelString = "
5  # JAGS model specification
6  model {
7  # Likelihood
8  for (j in 1:nItemsTotal) {
9  y[j] ~ dbin( theta[machines[j]], 40 )
10 }
11 # Prior
12 for (i in 1:nMachines){
13 theta[i] ~ dnorm( normMu, 1/normSigma^2 )T( 0.0,0.5 )
14 }
15 # Hyperpriors
16 normMu ~ dbeta( 1.1, 9.9 )
17 normSigma ~ dgamma( 2, 28 )
18 }
19 # ... JAGS model specification ends.
20 "
21 # close quote to end modelString
22
23 # Write the modelString to a file, using R commands:
24 writeLines(modelString,con="model.txt")
25
26
27 # THE DATA.
28
29
30 dataList = list(

```

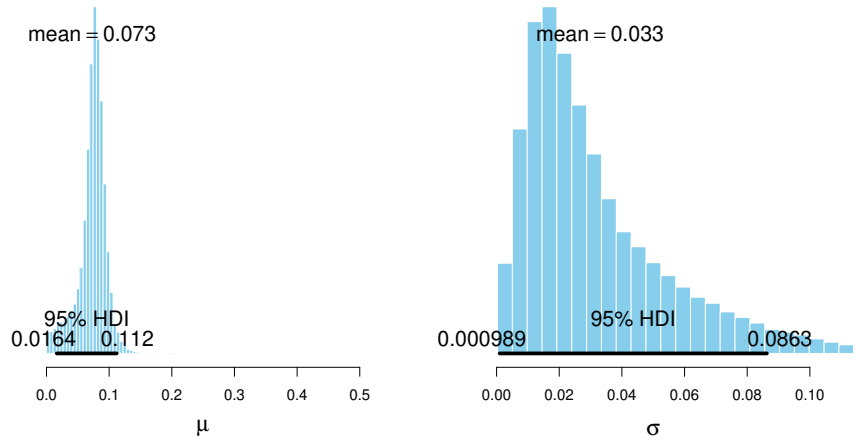


Figure 5: Sampled posterior distributions for hyperparameters of the truncated normal distribution.

```

31 y=c(1, 5, 5, 5, 2, 4, 2, 1, 3, 3, 3, 5, 3, 4, 6, 4, 2, 4, 5, 1, 1, 3, 2, 4, 1),
32 machines = rep(c(1,2,3,4,5),each=5),
33 nMachines = 5,
34 #nItems=rep(40, 25),
35 nItemsTotal = 25
36 )

```

The sampled posterior distributions for hyperparameters are reported in Figure 5.

## Exercise C

The problem can be formulated as follows. The number of successes under the three treatments are  $z = (8, 11, 14)$  with  $N = (30, 30, 30)$  trials. The model can be written as:

$$z_c \sim \text{Binomial}(\theta_c, N_c) \quad (8)$$

$$\theta_c \sim \text{Beta}(\omega(\kappa - 2) + 1, (1 - \omega)(\kappa - 2) + 1) \quad (9)$$

$$\omega \sim \text{Beta}(\omega_0(\kappa_0 - 2) + 1, (1 - \omega_0)(\kappa_0 - 2) + 1) \quad (10)$$

$$\kappa \sim \text{Gamma}(0.01, 0.01) \quad (11)$$

$$\omega_0 = 0.3 \quad (12)$$

$$\kappa_0 = 10 \quad (13)$$

Values for hyperpriors of  $\omega$  where set based on previous knowledge on recovery rate under no treatment. Values for hyperpriors of  $\kappa$  are set such that  $\kappa$  is vague.

Here is the code we implemented:

```

1 rm(list = ls())
2 graphics.off()
3 source("openGraphSaveGraph.R")
4 require(rjags)

```

```

5 | savePath = "Figures/ExerciseC/"
6 |
7 | #-----
8 | # THE MODEL.
9 | modelString = "
10 |   model {
11 |     for ( cIdx in 1:Ncat ) {
12 |       z[cIdx] ~ dbin( theta[cIdx] , N[cIdx] )
13 |       theta[cIdx] ~ dbeta( omega*(kappa-2)+1 ,
14 |                           (1-omega)*(kappa-2)+1 )
15 |     }
16 |     omega ~ dbeta( omega0*(kappa0-2)+1, omega0*(kappa0-2)+1 )
17 |     kappa <- kappaMinusTwo + 2
18 |     kappaMinusTwo ~ dgamma( 0.01 , 0.01 ) # mean=1 , sd=10 (generic vague)
19 |     omega0 <- 0.3
20 |     kappa0 <- 10
21 |   }
22 |   " # close quote for modelString
23 | writelines( modelString , con="Cmodel.txt" )
24 |
25 | #-----
26 | # THE DATA.
27 |
28 |
29 | z = c(8,11,14)
30 | N = c(30,30,30)
31 | Ncat = length(unique(z))
32 | # Specify the data in a list, for later shipment to JAGS:
33 | dataList = list(
34 |   z = z ,
35 |   N = N ,
36 |   Ncat = Ncat
37 | )
38 |
39 | #-----
40 | # INITIALIZE THE CHAIN.
41 |
42 | # Initial values of MCMC chains based on data:
43 | initsList = function() {
44 |   thetaInit = rep(NA,Ncat)
45 |   for ( cIdx in 1:Ncat ) { # for each subject
46 |     resampledZ = rbinom(1, size=N[cIdx] , prob=z[cIdx]/N[cIdx] )
47 |     thetaInit[cIdx] = resampledZ/N[cIdx]
48 |   }
49 |   thetaInit = 0.001+0.998*thetaInit # keep away from 0,1
50 |   kappaInit = 100 # lazy, start high and let burn-in find better value
51 |   return( list( theta=thetaInit ,
52 |                 omega0=mean(thetaInit) ,
53 |                 kappaMinusTwo0=kappaInit-2 ) )
54 | }
55 |
56 | #-----
57 | # RUN THE CHAINS.
58 |
59 | parameters = c( "theta" ) # The parameter(s) to be monitored.
60 | adaptSteps = 1000          # Number of steps to "tune" the samplers.
61 | burnInSteps = 5000         # Number of steps to "burn-in" the samplers.
62 | nChains = 3                # Number of chains to run.
63 | numSavedSteps=50000        # Total number of steps in chains to save.
64 | thinSteps=1                # Number of steps to "thin" (1=keep every step).
65 | nIter = ceiling( ( numSavedSteps * thinSteps ) / nChains ) # Steps per chain.
66 | # Create, initialize, and adapt the model:

```



```

67 jagsModel = jags.model( "Cmodel.txt" , data=dataList , inits=initsList ,
68                          n.chains=nChains , n.adapt=adaptSteps )
69 # Burn-in:
70 cat( "Burning in the MCMC chain...\n" )
71 update( jagsModel , n.iter=burnInSteps )
72 # The saved MCMC chain:
73 cat( "Sampling final MCMC chain...\n" )
74 codaSamples = coda.samples( jagsModel , variable.names=parameters ,
75                             n.iter=nIter , thin=thinSteps )
76
77
78 #-----
79 # EXAMINE THE RESULTS.
80 mcmcChain = as.matrix( codaSamples )
81
82 # Extract the posterior samples
83 theta1Sample = mcmcChain[, "theta[1]"]
84 theta2Sample = mcmcChain[, "theta[2]"]
85 theta3Sample = mcmcChain[, "theta[3]"]
86
87 # Compute differences:
88 thetaDiff12 = theta1Sample - theta2Sample
89 thetaDiff13 = theta1Sample - theta3Sample
90 thetaDiff23 = theta2Sample - theta3Sample
91
92 # Load plotting function
93 source("plotPost.R")
94
95 # Open a plotting window
96 openGraph(width=12, height=8)
97 par(mfrow=c(2,3)) # 2 rows, 3 columns for all 6 plots
98
99 # Plot posteriors for theta1, theta2, theta3
100 plotPost(theta1Sample, xlab=expression(theta[1]), main="Posterior of theta1")
101 plotPost(theta2Sample, xlab=expression(theta[2]), main="Posterior of theta2")
102 plotPost(theta3Sample, xlab=expression(theta[3]), main="Posterior of theta3")
103
104 # Plot posteriors for differences
105 plotPost(thetaDiff12, xlab=expression(theta[1] - theta[2]), main="theta1 - theta2",
106          compVal=0.0, ROPE=c(-0.05,0.05))
107 plotPost(thetaDiff13, xlab=expression(theta[1] - theta[3]), main="theta1 - theta3",
108          compVal=0.0, ROPE=c(-0.05,0.05))
109 plotPost(thetaDiff23, xlab=expression(theta[2] - theta[3]), main="theta2 - theta3",
110          compVal=0.0, ROPE=c(-0.05,0.05))
111
112 # Save the graph
113 saveGraph(file = paste0(savePath, "ThetaPosteriors"), type="eps")

```

Results are shown in Figure 6. The three posteriors exhibit different means. In particular,  $\theta_3$  has higher mean than  $\theta_2$  and  $\theta_2$  has higher mean than  $\theta_1$ . However, there is not sufficient evidence to defend the claim that any of the recovery rates under placebo is higher than the one under no treatment, nor that the expensive placebo is more effective than the cheap one. Actually, the ROPE  $[-0.05, 0.05]$  is always included in the HDIs of  $\theta$  differences.

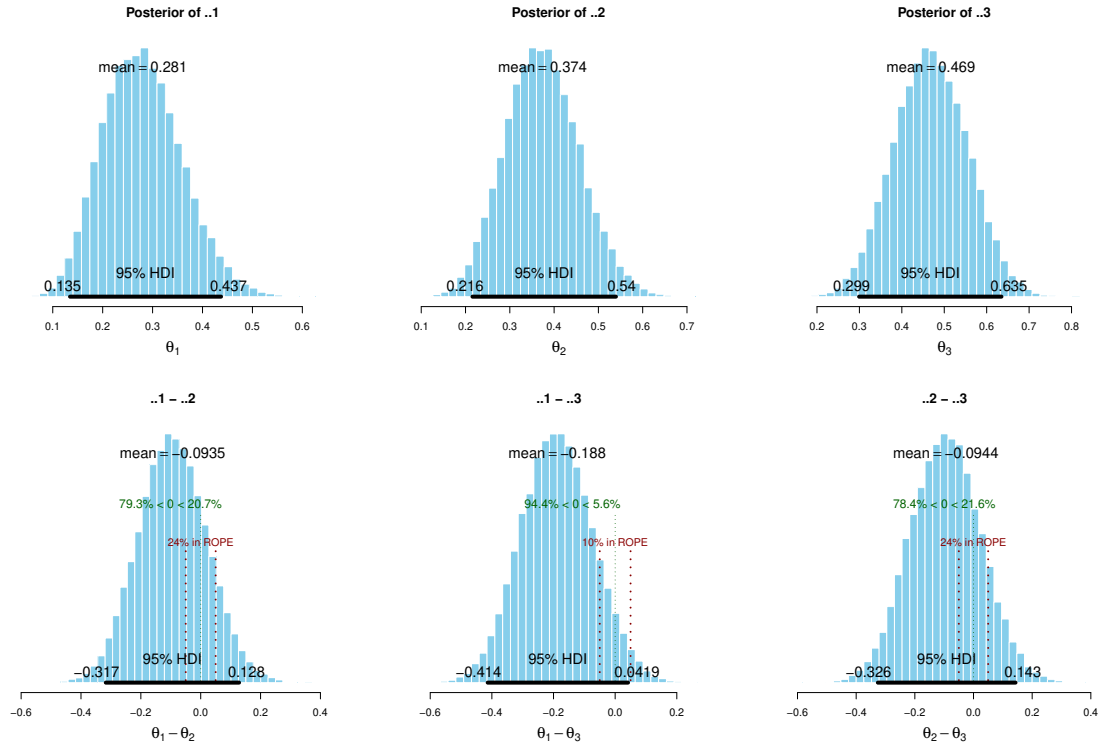


Figure 6: Posteriors and posterior differences for recovery rate  $\theta$  under the three different treatments of exercise C.