

BDA A3

Piero Birello

April 2025

Exercise A

Part 1

Let y_i be the number of tulips grown in row i (equivalently, from bag i). Let n_i be the number of bulbs contained in bag i and b_i the bag (bulb) type. The model can be written as:

$$y_i \sim \text{Binomial}(n_i, p_{b_i}) \quad (1)$$

$$n_i = \text{round}(\hat{n}_i) \quad (2)$$

$$\hat{n}_i \sim \text{Normal}(\mu = 30, \sigma = 1.5) \quad (3)$$

$$p_{b_i=1} \sim \text{Beta}(40, 10) \quad (4)$$

$$p_{b_i=2} \sim \text{Beta}(1, 1) \quad (5)$$

$$b_i \sim \text{Bernoulli}(1 - \theta) + 1 \quad (6)$$

$$\theta \sim \text{Beta}(1, 1) \quad (7)$$

In Jags, this translates to:

```
1 modelString = "  
2 # JAGS model specification  
3 model {  
4   # Likelihood  
5   for (i in 1:nBags) {  
6     y[i] ~ dbin( p[b[i]], n[i] )  
7     n[i] <- round(nContinuous[i])  
8     nContinuous[i] ~ dnorm(normMu, normTau)  
9     b[i] <- bBernoulli[i] + 1  
10    bBernoulli[i] ~ dbern(1-catTheta) # equivalent to dcat  
11  }  
12  # Priors  
13  p[1] ~ dbeta(betaA, betaB) # when bBernoulli=0  
14  p[2] ~ dbeta(1,1)         # when bBernoulli=1  
15  catTheta ~ dbeta(1,1)  
16  # Constants  
17  normMu <- 30  
18  normTau <- 1/(normSigma*normSigma) #Convert to precision  
19  normSigma <- 1.5  
20  betaA <- 40  
21  betaB <- 10  
22 }  
23 # ... JAGS model specification ends.  
24 " # close quote to end modelString
```

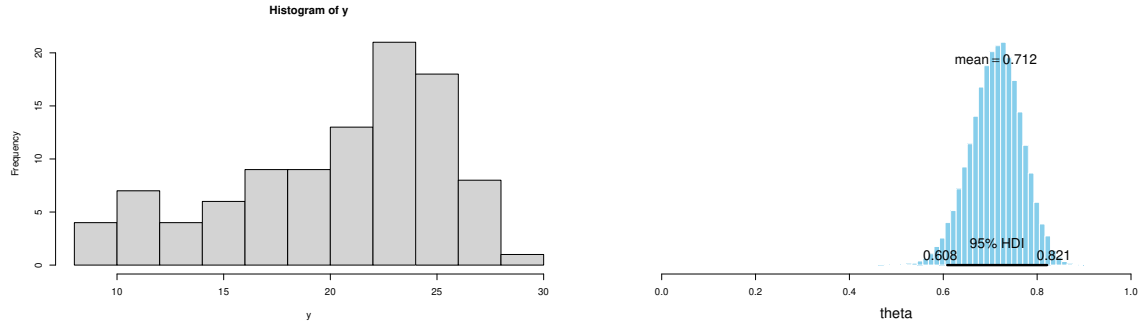


Figure 1: Histogram of y and posterior estimate of parameter θ .

A histogram for the occurrence of the growth of a number of tulips y is shown in Figure 1. Parameter estimates are shown in Figure 1 and Figure 2. Results indicate that around 70% of the bags are of type A (mean=0.712, 95% HDI=[0.608,0.821]). Type A bulbs actually have a fertility $p_{b_i=1}$ close to 0.8 (mean=0.785, 95% HDI=[0.761,0.809]), while type B bulbs have lower fertility (mean=0.468, 95% HDI=[0.417,0.514]). We may exclude the hypothesis that the two bulb types have the same fertility, since a ROPE centered in 0 of half-length 0.1 contains no samples of the difference $p_{b_i=1} - p_{b_i=2}$. If we assume prior knowledge that the value of θ is close to 0.5, i.e., $\theta \sim \text{Beta}(15,15)$, we get the results reported in Figure 3 and Figure 4. While estimates for $p_{b=1}$ and $p_{b=2}$ are similar to the previous ones, posterior estimates for θ are pulled towards 0.5, resulting in mean=0.652, 95% HDI=[0.556,0.743] (previous was mean=0.712, 95% HDI=[0.608,0.821]). We can perform a proper model comparison to select the most plausible prior distribution. We can implement this in JAGS by introducing a dummy variable for the model index, distributed according to a categorical distribution with prior probability of 0.5 on each of the indices. This reads:

```

1 # JAGS MODEL
2
3 # Specify the model in JAGS language, save it as a string in R:
4 modelString = "
5 # JAGS model specification
6 model {
7 # Likelihood
8 for (i in 1:nBags) {
9 y[i] ~ dbin( p[b[i]], n[i] )
10 n[i] <- round(nContinuous[i])
11 nContinuous[i] ~ dnorm(normMu,normTau)
12 b[i] <- bBernoulli[i] + 1
13 bBernoulli[i] ~ dbern(1-catTheta) # equivalent to dcat(pCat[]) // pCat[1] = catTheta
14 }
15 # Priors
16 p[1] ~ dbeta(p1BetaA,p1BetaB) # when b=0
17 p[2] ~ dbeta(1,1) # when b=1
18 catTheta ~ dbeta(thetaBetaA[m],thetaBetaB[m])
19 # Constants
20 normMu <- 30
21 normTau <- 1/(normSigma*normSigma) #Convert to precision
22 normSigma <- 1.5
23 p1BetaA <- 40
24 p1BetaB <- 10
25 thetaBetaA[1] <- 1
26 thetaBetaB[1] <- 1
27 thetaBetaA[2] <- 15
28 thetaBetaB[2] <- 15

```

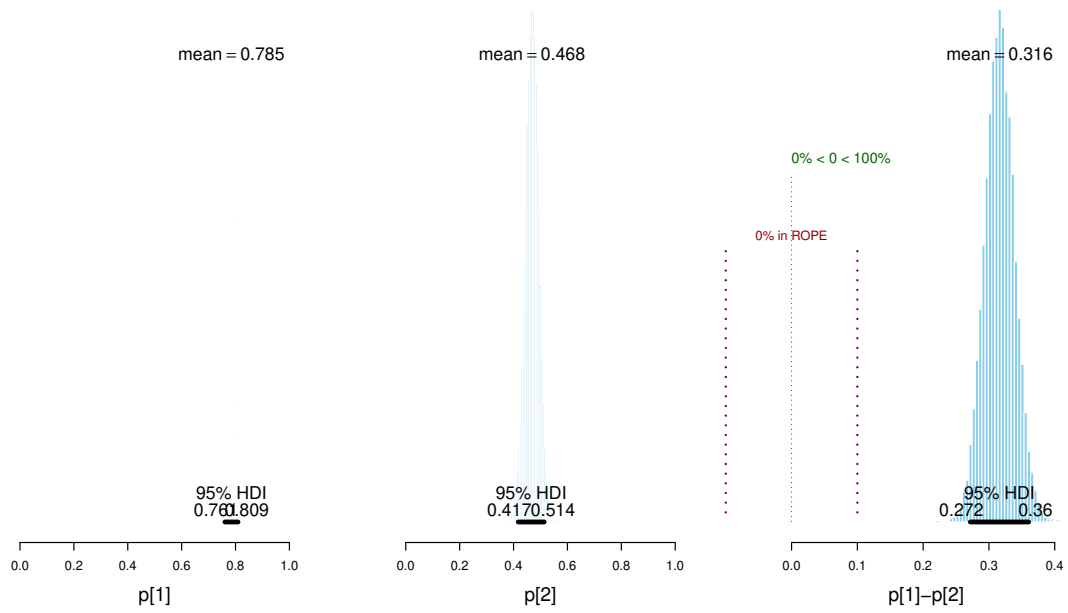


Figure 2: Posterior estimates for the fertility of the two bulb types $p_{b=1}$ and $p_{b=2}$.

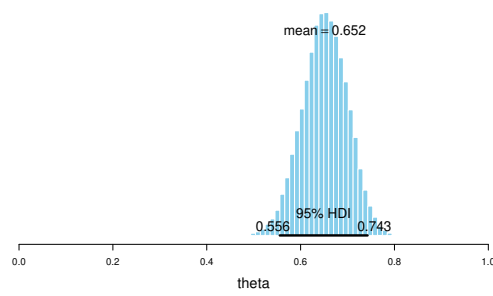


Figure 3: Posterior estimate for θ when the prior is $Beta(15, 15)$.

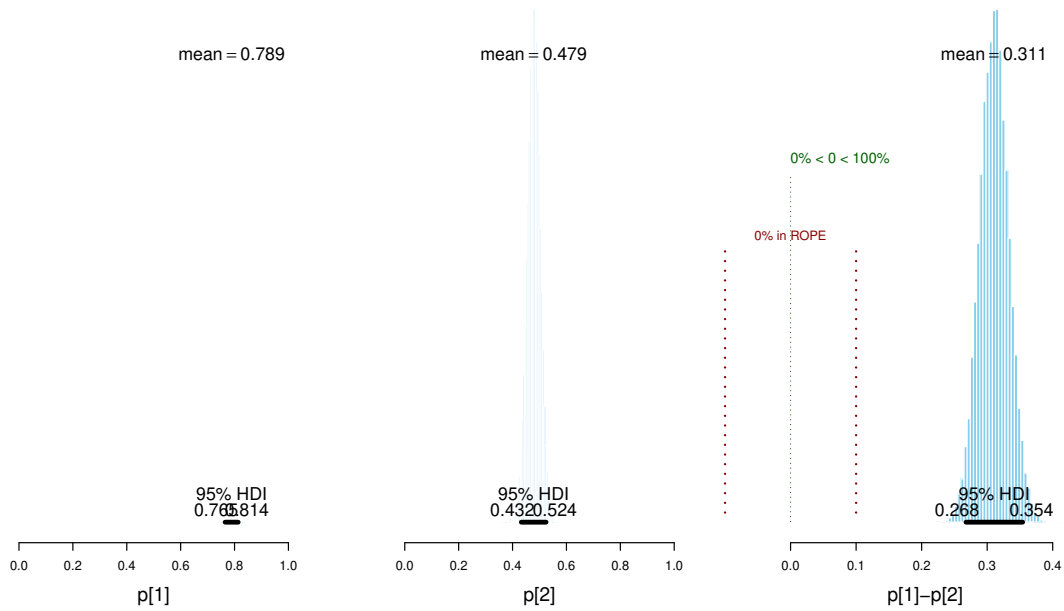


Figure 4: Posterior estimates for $p_{b=1}$ and $p_{b=2}$ when the prior for θ is $Beta(15,15)$.

```

29 # Model Dummy Variable
30 m ~ dcat(mPriorProb)
31 mPriorProb[1] <- 0.5
32 mPriorProb[2] <- 0.5
33 }
34 # ... JAGS model specification ends.
35 " # close quote to end modelString
36
37 # Write the modelString to a file, using R commands:
38 writelines(modelString,con="modelExercise1c.txt")
39
40 #-----
41 # INITIALIZE THE CHAIN
42 initsList = list( p=c(0.8,0.8), catTheta=0.5 , m=1 )
43
44 #-----
45 # RUN THE CHAINS
46
47 # Assign parameters
48 parameters = c( "catTheta","p[1]","p[2]","m" ) # The parameter(s) to be monitored.
49 adaptSteps = 1000 # Number of steps to "tune" the samplers.
50 burnInSteps = 5000 # 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
56 # Create, initialize, and adapt the model:
57 jagsModel = jags.model( "modelExercise1c.txt" , data=dataList , inits=initsList ,
58 n.chains=nChains , n.adapt=adaptSteps )

```

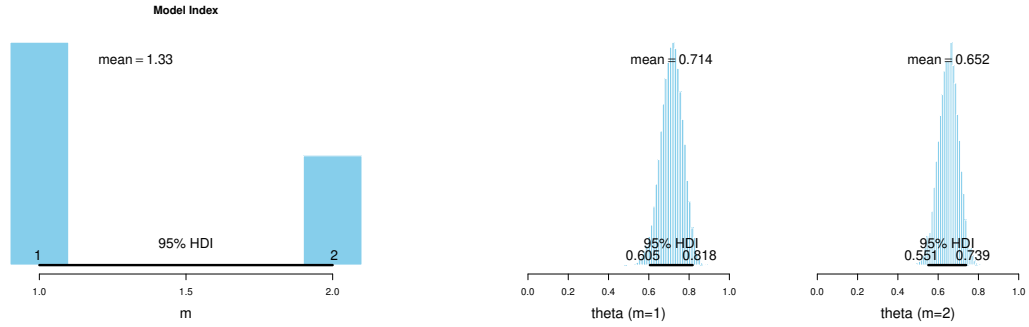


Figure 5: Posterior estimates for the model index m and parameter θ under the two models.

```

59 # Burn-in:
60 cat( "Burning in the MCMC chain...\n" )
61 update( jagsModel , n.iter=burnInSteps )
62 # The saved MCMC chain:
63 cat( "Sampling final MCMC chain...\n" )
64 codaSamples = coda.samples( jagsModel , variable.names=parameters ,
65                             n.iter=nIter , thin=thinSteps )
66
67 #-----
68 # EXAMINE THE RESULTS
69 mcmcChain = as.matrix( codaSamples )
70
71 # Extract the posterior sample from JAGS:
72 catThetaSample = mcmcChain[, "catTheta"]
73 p1Sample = mcmcChain[, "p[1]" ]
74 p2Sample = mcmcChain[, "p[2]" ]
75 pDiffSample = p1Sample - p2Sample
76 mSample = mcmcChain[, "m"]
77
78 # Compute the proportion of m at each index value:
79 pM1 = sum( mSample == 1 ) / length( mSample )
80 pM2 = 1 - pM1
81 print(paste0('model 1 prob:', pM1))
82
83 # Extract param values for each model index:
84 catThetaSampleM1 = catThetaSample[ mSample == 1 ]
85 catThetaSampleM2 = catThetaSample[ mSample == 2 ]
86 p1SampleM1 = p1Sample[ mSample == 1 ]
87 p1SampleM2 = p1Sample[ mSample == 2 ]
88 p2SampleM1 = p2Sample[ mSample == 1 ]
89 p2SampleM2 = p2Sample[ mSample == 2 ]
90 pDiffSampleM1 = p1SampleM1 - p2SampleM1
91 pDiffSampleM2 = p1SampleM2 - p2SampleM2

```

Results for posterior estimates of the model index m , parameter θ and the fertilities $p_{b=1}$ and $p_{b=2}$ under the two models are shown in Figure 5, Figure 6 and Figure 7. Again, we see that the two models result in equivalent estimates for the fertilities. However, model 1 estimates θ to be larger than model 2, whose prior pulls the estimate towards the value $\theta = 0.5$. Model 1 is also the most likely. Precisely, the posterior probability for model 1 is 0.6708. Equivalently, posterior odds are $p(m = 1|D)/p(m = 2|D) = 2.0377$, i.e., model 1 is two times as likely as model 2.

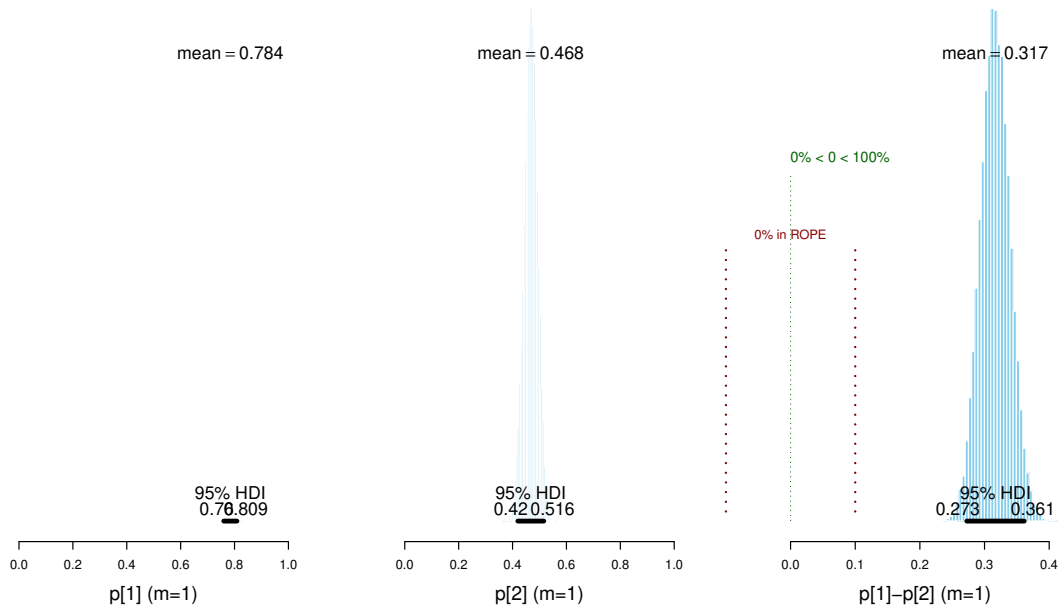


Figure 6: Posterior estimates for $p_{b=1}$ and $p_{b=2}$ under model 1.

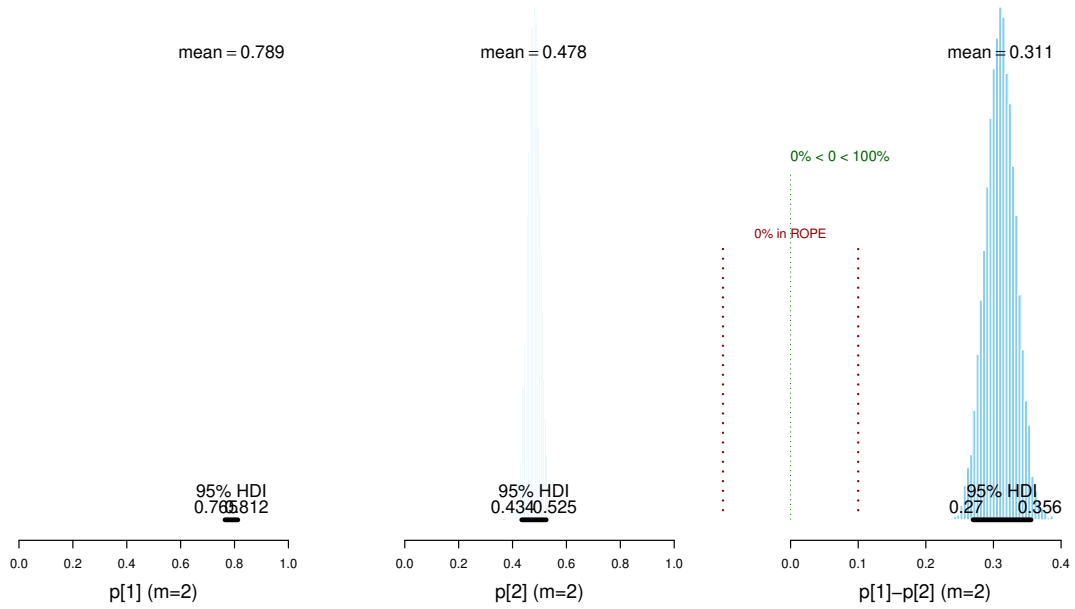


Figure 7: Posterior estimates for $p_{b=1}$ and $p_{b=2}$ under model 2.

Part 2

We now perform a power analysis to estimate the number of bags we will need to test in order for HDI's of the two fertility parameters not to overlap, given hypothesized hyper-parameters $p_{b_i=1} = 0.8, p_{b_i=2} = 0.5, \theta = 0.3$. We generate representative samples of n_i and b_i given the hypothesized hyper-parameters. From the representative parameters, we generate a random sample of the observable data y_i . We run Bayesian inference according to the model used so far (using $\theta \sim \text{Beta}(1,1)$). We then check whether the goal of non overlapping HDI between $p_{b_i=1}$ and $p_{b_i=2}$ is achieved. We repeat for various sample sizes, i.e., number of bags, starting from 4 and increasing until the power crosses 0.8. For each sample size, we repeat the sampling procedure 10 times. The developed code is the following:

```
1 # Jags-Ybin-Nnorm-Bbern-Pbeta-Power.R
2 graphics.off() # This closes all of R's graphics windows.
3 rm(list=ls()) # Careful! This clears all of R's memory!
4 fileNameRoot = "Jags-Ydich-Xnom1subj-MbernBeta-Power-" # for future use
5
6 # Load the functions genMCMC, smryMCMC, and plotMCMC:
7 # (This also sources DBDA2E-utilities.R)
8 source("Jags-Ybin-Nnorm-Bbern-Pbeta.R")
9
10 #-----
11 # Define function that assesses goal achievement for a single set of data:
12 goalAchievedForSample = function( data ) {
13   # Generate the MCMC chain:
14   mcmcCoda = genMCMC( data=data , numSavedSteps=10000 , saveName=NULL )
15   # Check goal achievement. First, compute the HDI:
16   p1HDI = HDIoMCMC( as.matrix(mcmcCoda[, "p[1]"] ) )
17   p2HDI = HDIoMCMC( as.matrix(mcmcCoda[, "p[2]"] ) )
18   # Define list for recording results:
19   goalAchieved = list()
20   # Goal: No overlap between HDIs
21   goalAchieved = c(goalAchieved,
22     "noOverlap"=(( p1HDI[2] < p2HDI[1]) || (p2HDI[2] < p1HDI[1])) )
23   # More goals can be inserted here if wanted...
24   # Return list of goal results:
25   return(goalAchieved)
26 }
27
28 #-----
29 # Specify hypothetical parameters:
30 p = c(0.8,0.5)
31 catTheta=0.3
32 normMu = 30
33 normSigma = 1.5
34
35 # Define proportion success list
36 proportionSuccessList = c()
37
38 # Iterate over sample size (number of bags):
39 for (sampleN in 4:100) {
40   # Run a bunch of simulated experiments:
41   nSimulatedDataSets = 10
42   # Track number of times goal is achieved
43   goalCount=0
44
45   for ( simIdx in 1:nSimulatedDataSets ) {
46     # Generate random value from hypothesized parameter distribution:
47     nContinuous = rnorm(n=sampleN, mean=normMu, sd=normSigma)
48     n = round(nContinuous)
49     bBernoulli = rbinom(n=sampleN, size=1, prob=1-catTheta)
```

```

50     b = bBernoulli + 1
51     # Generate random data based on parameter value:
52     sampleY = sapply(1:length(n), function(i) {
53         rbinom(n = 1, size = n[i], prob = p[b[i]])
54     })
55     # Do Bayesian analysis on simulated data:
56     goalAchieved = goalAchievedForSample( sampleY )
57     # Tally the results:
58     if (!exists("goalTally")) { # if goalTally does not exist, create it
59         goalTally=matrix( nrow=0 , ncol=length(goalAchieved) )
60     }
61     goalTally = rbind( goalTally , goalAchieved )
62     # save( goalTally ,
63     #       file="Jags-Ydich-Xnom1subj-MbernBeta-Power-goalTally.Rdata" )
64     if (goalAchieved[["noOverlap"]] == TRUE) {
65         goalCount = goalCount + 1
66     }
67 }
68
69 # Calculate proportion of successes:
70 proportionSuccess = goalCount / nSimulatedDataSets
71 cat("SampleN =", sampleN, " | Success rate:", round(proportionSuccess * 100, 2), "%\n")
72 proportionSuccessList[as.character(sampleN)] = proportionSuccess
73 # Check stopping condition:
74 if (proportionSuccess > 0.8) {
75     print(proportionSuccessList)
76     cat("Stopping early: goal achieved in over 80% of simulations for sample size", sampleN,
77         "\n")
78     break
79 }
80 }

```

Note that the sourced script is the one developed for the previous part of the exercise, with the only difference that the relevant section has been wrapped in the function `genMCMC`. We obtain that a power of 0.8 or more can be obtained already by testing 7 bags. More precisely, obtained power estimates are 0.6, 0.7, 0.6, 0.9 for 4, 5, 6, 7 bags respectively. Note that a larger number of trials per each sample size would guarantee a more accurate estimate.

We repeat the experiment considering the case in which $p_{b_i=2} = 0.5$. Unfortunately here the sampler repeatedly fails to initialize when testing large sample sizes, despite efforts to improve initialization and catchment of errors. The code developed to -unsuccessfully- try solve the issue is reported in section .

Exercise B

We perform a linear regression of the form:

$$y_i \sim \text{Normal}(\mu_i, \sigma) \quad (8)$$

$$\mu_i = \beta_0 + \beta_1 x_i \quad (9)$$

$$\beta_0 \sim \text{Normal}(0, 10) \quad (10)$$

$$\beta_1 \sim \text{Normal}(0, 10) \quad (11)$$

$$\sigma = 1/\sqrt{\tau} \quad (12)$$

$$\tau \sim \text{Gamma}(0.01, 0.01) \quad (13)$$

where variables x and y have been standardized. Estimates for the coefficients are shown in Table 1. Posterior estimates are shown in Figure 8 insetad. The interesting parameter to observe in order to understand whether

Parameter	Mean	Median	Mode	HDI Low	HDI High
β_0	684.00	684.14	685.15	660.50	707.83
β_1	-1.55	-1.56	-1.58	-2.75	-0.34
σ	18.36	18.32	18.32	16.74	19.99

Table 1: Posterior summary statistics with 95% highest density intervals (HDI)

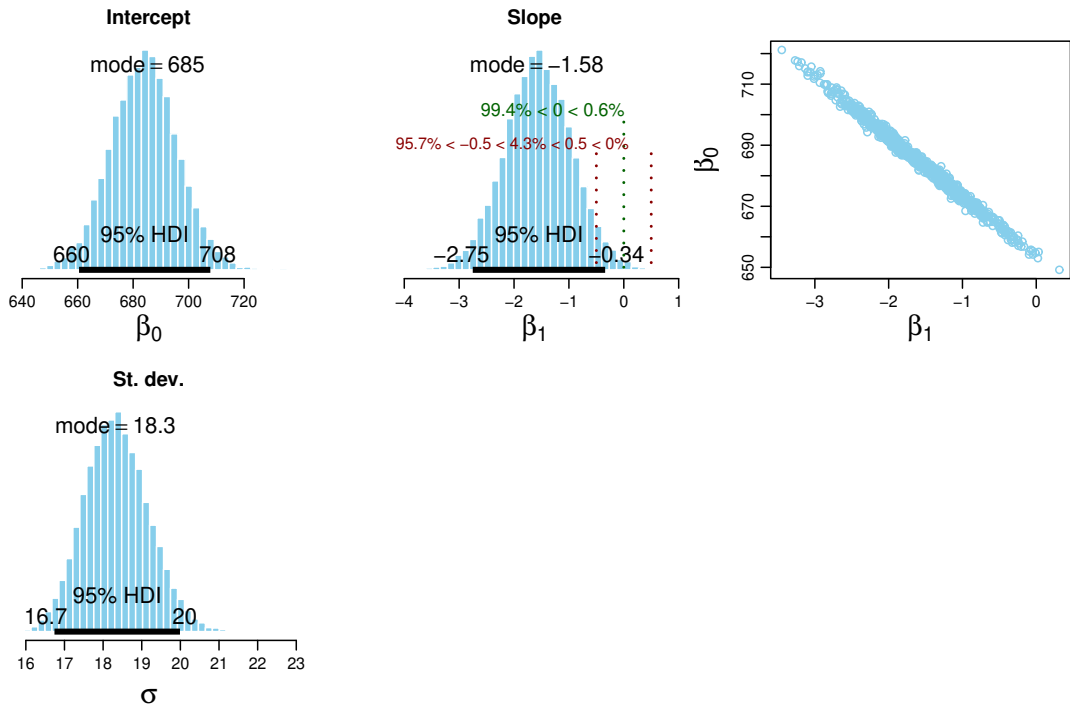


Figure 8: Posterior estimates of linear regression parameters.

the students to teacher ratio influences test scores is the slope β_1 . Posterior estimates suggests that it is very

likely that higher student to teacher ratios correspond to lower test scores. The HDI does not stand completely out of the ROPE centered in 0 with half-length 0.5, however only the 4.255% of the HDI falls inside the ROPE, and only the 0.6% of sampled values is higher than zero. In Figure 9 we show a posterior predictive check of the linear regression model, with plausible regression lines and sampled noise distributions. The model seems to correctly fit the data, with few outliers with respect to the noise distribution. A fit using a t-distribution could improve robustness and include these outliers better. Hence, although model improvement would be recommended (also testing a hierarchical model including group variables), we may preliminary conclude that we would suggest investing in assuming more teachers. As an example, we show the posterior predicted distribution for test scores when the students to teacher ratio is 20 in Figure 10. The code used in this Exercise B is composed by a source script and an example script. The source script:

```

1 # Jags-Ymet-Xmet-Mrobust.R
2 source("DBDA2E-utilities.R")
3
4 #=====
5
6 genMCMC = function( data , xName="x" , yName="y" ,
7                     numSavedSteps=50000 , saveName=NULL ) {
8   require(rjags)
9   #-----
10  # THE DATA.
11  y = data[,yName]
12  x = data[,xName]
13  # Do some checking that data make sense:
14  if ( any( !is.finite(y) ) ) { stop("All y values must be finite.") }
15  if ( any( !is.finite(x) ) ) { stop("All x values must be finite.") }
16  #Ntotal = length(y)
17  # Specify the data in a list, for later shipment to JAGS:
18  dataList = list(
19    x = x ,
20    y = y
21  )
22  #-----
23  # THE MODEL.
24  modelString = "
25  # Standardize the data:
26  data {
27    Ntotal <- length(y)
28    xm <- mean(x)
29    ym <- mean(y)
30    xsd <- sd(x)
31    ysd <- sd(y)
32    for ( i in 1:length(y) ) {
33      zx[i] <- ( x[i] - xm ) / xsd
34      zy[i] <- ( y[i] - ym ) / ysd
35    }
36  }
37  # Specify the model for standardized data:
38  model {
39    for ( i in 1:Ntotal ) {
40      zy[i] ~ dnorm( zmu[i] , ztau )
41      zmu[i] <- zbeta0 + zbeta1 * zx[i]
42    }
43    # Priors vague on standardized scale:
44    zbeta0 ~ dnorm( 0 , 1/(10)^2 )
45    zbeta1 ~ dnorm( 0 , 1/(10)^2 )
46    ztau ~ dgamma( 1.0E-2 , 1.0E-2 ) # corresponds to mean=1, var=100
47    # Transform to original scale:
48    beta1 <- zbeta1 * ysd / xsd
49    beta0 <- zbeta0 * ysd + ym - zbeta1 * xm * ysd / xsd

```

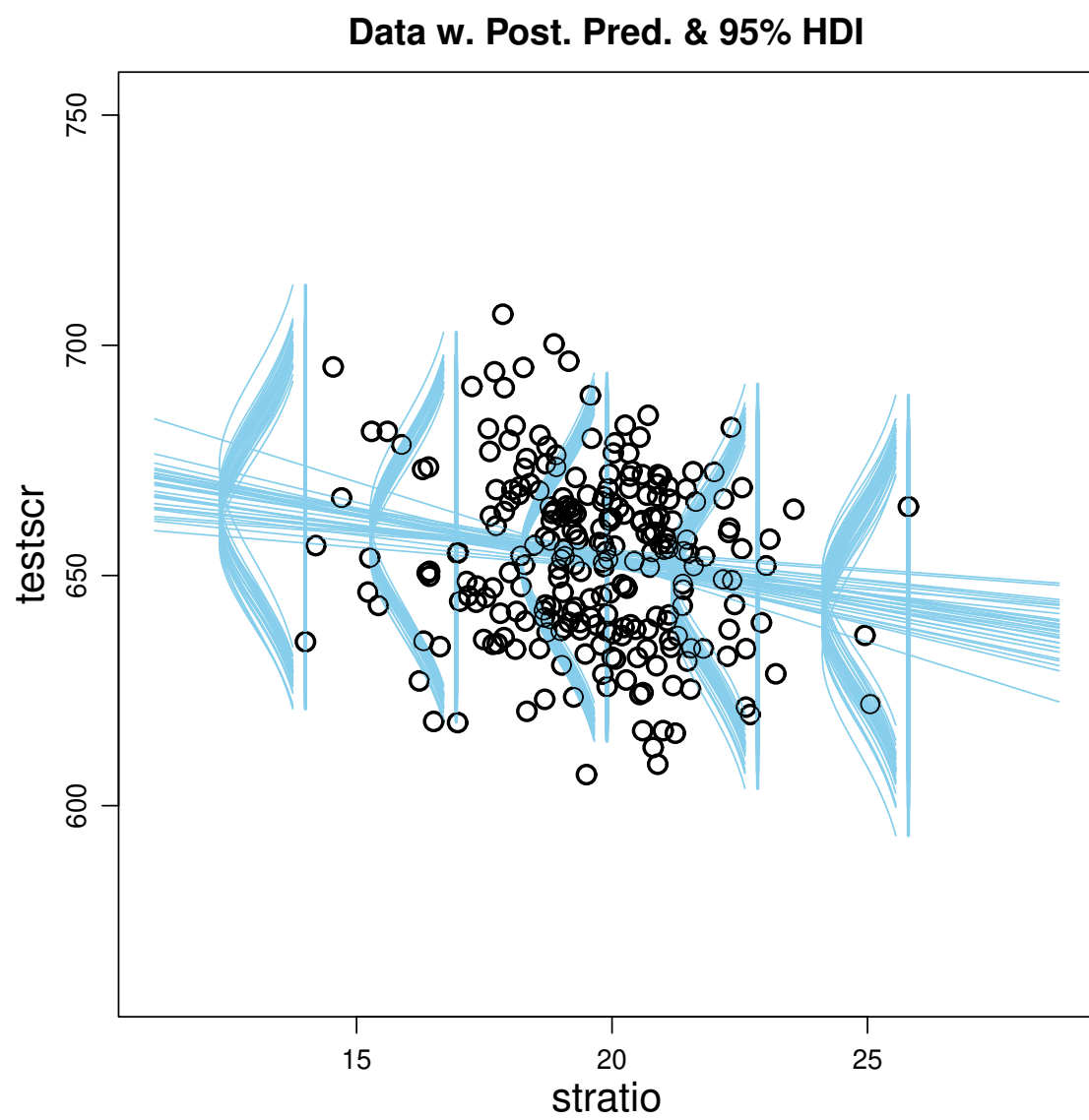


Figure 9: Posterior predictive check

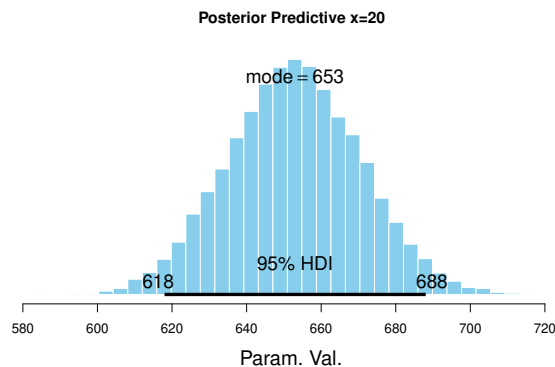


Figure 10: Posterior predicted distribution when $x=20$.

```

50   tau <- ztau / (ysd^2)
51   sigma <- 1/sqrt(tau)
52 }
53 " # close quote for modelString
54 # Write out modelString to a text file
55 writeLines( modelString , con="TEMPmodel.txt" )
56 #-----
57 # INITIALIZE THE CHAINS.
58 # Let JAGS do it...
59 #-----
60 # RUN THE CHAINS
61 parameters = c( "beta0" , "beta1" , "tau", "sigma",
62                "zbeta0" , "zbeta1" , "ztau" )
63 adaptSteps = 500 # Number of steps to "tune" the samplers
64 burnInSteps = 1000
65 nChains = 4
66 thinSteps = 1
67 nIter = ceiling( ( numSavedSteps * thinSteps ) / nChains )
68 # Create, initialize, and adapt the model:
69 jagsModel = jags.model( "TEMPmodel.txt" , data=dataList , #inits=initsList ,
70                        n.chains=nChains , n.adapt=adaptSteps )
71 # Burn-in:
72 cat( "Burning in the MCMC chain...\n" )
73 update( jagsModel , n.iter=burnInSteps )
74 # The saved MCMC chain:
75 cat( "Sampling final MCMC chain...\n" )
76 codaSamples = coda.samples( jagsModel , variable.names=parameters ,
77                             n.iter=nIter , thin=thinSteps )
78 # resulting codaSamples object has these indices:
79 #   codaSamples[[ chainIdx ]][ stepIdx , paramIdx ]
80 if ( !is.null(saveName) ) {
81   save( codaSamples , file=paste(saveName,"Mcmc.Rdata",sep="") )
82 }
83 return( codaSamples )
84 } # end function
85
86 #=====
87
88 smryMCMC = function( codaSamples ,
89                      compValBeta0=NULL , ropeBeta0=NULL ,
90                      compValBeta1=NULL , ropeBeta1=NULL ,

```

```

91         compValSigma=NULL , ropeSigma=NULL ,
92         saveName=NULL ) {
93     summaryInfo = NULL
94     mcmcMat = as.matrix(codaSamples,chains=TRUE)
95     summaryInfo = rbind( summaryInfo ,
96         "beta0" = summarizePost( mcmcMat[, "beta0"] ,
97                                 compVal=compValBeta0 ,
98                                 ROPE=ropeBeta0 ) )
99     summaryInfo = rbind( summaryInfo ,
100        "beta1" = summarizePost( mcmcMat[, "beta1"] ,
101                                compVal=compValBeta1 ,
102                                ROPE=ropeBeta1 ) )
103     summaryInfo = rbind( summaryInfo ,
104        "sigma" = summarizePost( mcmcMat[, "sigma"] ,
105                                compVal=compValSigma ,
106                                ROPE=ropeSigma ) )
107
108     if ( !is.null(saveName) ) {
109         write.csv( summaryInfo , file=paste(saveName,"SummaryInfo.csv",sep="") )
110     }
111     return( summaryInfo )
112 }
113
114 #=====
115
116 plotMCMC = function( codaSamples , data , xName="x" , yName="y" ,
117                     compValBeta0=NULL , ropeBeta0=NULL ,
118                     compValBeta1=NULL , ropeBeta1=NULL ,
119                     compValSigma=NULL , ropeSigma=NULL ,
120                     showCurve=FALSE , pairsPlot=FALSE ,
121                     saveName=NULL , saveType="jpg" ) {
122     # showCurve is TRUE or FALSE and indicates whether the posterior should
123     # be displayed as a histogram (by default) or by an approximate curve.
124     # pairsPlot is TRUE or FALSE and indicates whether scatterplots of pairs
125     # of parameters should be displayed.
126     #-----
127     y = data[,yName]
128     x = data[,xName]
129     mcmcMat = as.matrix(codaSamples,chains=TRUE)
130     chainLength = NROW( mcmcMat )
131     zbeta0 = mcmcMat[, "zbeta0"]
132     zbeta1 = mcmcMat[, "zbeta1"]
133     ztau = mcmcMat[, "ztau"]
134     beta0 = mcmcMat[, "beta0"]
135     beta1 = mcmcMat[, "beta1"]
136     tau = mcmcMat[, "tau"]
137     sigma = mcmcMat[, "sigma"]
138     #-----
139     if ( pairsPlot ) {
140         # Plot the parameters pairwise, to see correlations:
141         openGraph()
142         nPtToPlot = 1000
143         plotIdx = floor(seq(1,chainLength,by=chainLength/nPtToPlot))
144         panel.cor = function(x, y, digits=2, prefix="", cex.cor, ...) {
145             usr = par("usr"); on.exit(par(usr))
146             par(usr = c(0, 1, 0, 1))
147             r = (cor(x, y))
148             txt = format(c(r, 0.123456789), digits=digits)[1]
149             txt = paste(prefix, txt, sep="")
150             if(missing(cex.cor)) cex.cor <- 0.8/strwidth(txt)
151             text(0.5, 0.5, txt, cex=1.25 ) # was cex=cex.cor*r
152             on.exit(par(usr=usr))

```

```

153 }
154 pairs( cbind( beta0 , beta1 , tau ) [plotIdx,] ,
155       labels=c( expression(beta[0]) , expression(beta[1]) ,
156               expression(tau) ) ,
157       lower.panel=panel.cor , col="skyblue" )
158 if ( !is.null(saveName) ) {
159   saveGraph( file=paste(saveName,"PostPairs",sep=""), type=saveType)
160 }
161 }
162 #-----
163 # Marginal histograms:
164 # Set up window and layout:
165 nPtToPlot = 1000
166 plotIdx = floor(seq(1,chainLength,by=chainLength/nPtToPlot))
167 openGraph(width=8,height=5)
168 layout( matrix( 1:6 , nrow=2, byrow=TRUE ) )
169 par( mar=c(4,4,2.5,0.5) , mgp=c(2.5,0.7,0) )
170 histInfo = plotPost( beta0 , cex.lab = 1.75 , showCurve=showCurve ,
171                     compVal=compValBeta0 , ROPE=ropeBeta0 ,
172                     xlab=bquote(beta[0]) , main=paste("Intercept") )
173 histInfo = plotPost( beta1 , cex.lab = 1.75 , showCurve=showCurve ,
174                     compVal=compValBeta1 , ROPE=ropeBeta1 ,
175                     xlab=bquote(beta[1]) , main=paste("Slope") )
176 plot( beta1[plotIdx] , beta0[plotIdx] ,
177       xlab=bquote(beta[1]) , ylab=bquote(beta[0]) ,
178       col="skyblue" , cex.lab = 1.75 )
179 histInfo = plotPost( sigma , cex.lab = 1.75 , showCurve=showCurve ,
180                     compVal=compValSigma , ROPE=ropeSigma ,
181                     xlab=bquote(sigma) , main=paste("St. dev.") )
182 if ( !is.null(saveName) ) {
183   saveGraph( file=paste(saveName,"PostMarg",sep=""), type=saveType)
184 }
185 #-----
186 # Data with superimposed regression lines and noise distributions:
187 openGraph()
188 par( mar=c(3,3,2,1)+0.5 , mgp=c(2.1,0.8,0) )
189 # Plot data values:
190 postPredHDI mass = 0.95
191 xRang = max(x)-min(x)
192 yRang = max(y)-min(y)
193 xLimMult = 0.25
194 yLimMult = 0.45
195 xLim= c( min(x)-xLimMult*xRang , max(x)+xLimMult*xRang )
196 yLim= c( min(y)-yLimMult*yRang , max(y)+yLimMult*yRang )
197 plot( x , y , cex=1.5 , lwd=2 , col="black" , xlim=xLim , ylim=yLim ,
198       xlab=xName , ylab=yName , cex.lab=1.5 ,
199       main=paste( "Data w. Post. Pred. & ",postPredHDI mass*100,"% HDI" , sep="" ) ,
200       cex.main=1.33 )
201 # Superimpose a smattering of believable regression lines:
202 nPredCurves=30
203 xComb = seq(xLim[1],xLim[2],length=501)
204 for ( i in floor(seq(from=1,to=chainLength,length=nPredCurves)) ) {
205   lines( xComb , beta0[i] + beta1[i]*xComb , col="skyblue" )
206 }
207 # Superimpose some vertical distributions to indicate spread:
208 #source("HDIofICDF.R")
209 nSlice = 5
210 curveXpos = seq(min(x),max(x),length=nSlice)
211 curveWidth = (max(x)-min(x))/(nSlice+2)
212 for ( i in floor(seq(from=1,to=chainLength,length=nPredCurves)) ) {
213   for ( j in 1:length(curveXpos) ) {
214     sigma = 1/sqrt(tau[i]) # Convert precision to standard deviation

```

```

215     yHDI = HDIoFICDF( qnorm , credMass=postPredHDI mass ) # Use normal distribution
216     yComb = seq(yHDI[1],yHDI[2],length=75)
217     xVals = dnorm( yComb ) # Use normal density
218     xVals = curveWidth * xVals / dnorm(0)
219     yPred = beta0[i] + beta1[i]*curveXpos[j]
220     yComb = yComb*sigma + yPred
221     lines( curveXpos[j] - xVals , yComb , col="skyblue" )
222     lines( curveXpos[j] - 0*xVals , yComb , col="skyblue" , lwd=2 )
223 }
224 }
225 # replot the data, in case they are obscured by lines:
226 points( x , y , cex=1.5 )
227 if ( !is.null(saveName) ) {
228     saveGraph( file=paste(saveName,"PostPred",sep=""), type=saveType)
229 }
230 # if you want to show the y intercept, set this to TRUE:
231 showIntercept=TRUE
232 if ( showIntercept ) {
233     openGraph()
234     par( mar=c(3,3,2,1)+0.5 , mgp=c(2.1,0.8,0) )
235     # Plot data values:
236     xRang = max(x)-min(x)
237     yRang = max(y)-min(y)
238     xLimMult = 0.25
239     yLimMult = 0.45
240     xLim= c( min(x)-xLimMult*xRang , max(x)+xLimMult*xRang )
241     xLim = c(0,max(xLim))
242     yLim= c( min(y)-yLimMult*yRang , max(y)+yLimMult*yRang )
243     nPredCurves=30
244     pltIdx = floor(seq(from=1,to=chainLength,length=nPredCurves))
245     intRange = range( beta0[pltIdx] )
246     yLim = range( c(yLim,intRange) )
247     postPredHDI mass = 0.95
248     plot( x , y , cex=1.5 , lwd=2 , col="black" , xlim=xLim , ylim=yLim ,
249           xlab=xName , ylab=yName , cex.lab=1.5 ,
250           main=paste( "Data w. Post. Pred. & ",postPredHDI mass*100,"% HDI" ,sep="") ,
251           cex.main=1.33 )
252     abline(v=0,lty="dashed")
253     # Superimpose a smattering of believable regression lines:
254     xComb = seq(xLim[1],xLim[2],length=501)
255     for ( i in pltIdx ) {
256         lines( xComb , beta0[i] + beta1[i]*xComb , col="skyblue" )
257     }
258     # Superimpose some vertical distributions to indicate spread:
259     #source("HDIofICDF.R")
260     nSlice = 5
261     curveXpos = seq(min(x),max(x),length=nSlice)
262     curveWidth = (max(x)-min(x))/(nSlice+2)
263     for ( i in floor(seq(from=1,to=chainLength,length=nPredCurves)) ) {
264         for ( j in 1:length(curveXpos)) {
265             sigma = 1/sqrt(tau[i]) # Convert precision to standard deviation
266             yHDI = HDIoFICDF( qnorm , credMass=postPredHDI mass ) # Use normal distribution
267             yComb = seq(yHDI[1],yHDI[2],length=75)
268             xVals = dnorm( yComb ) # Use normal density
269             xVals = curveWidth * xVals / dnorm(0)
270             yPred = beta0[i] + beta1[i]*curveXpos[j]
271             yComb = yComb*sigma + yPred
272             lines( curveXpos[j] - xVals , yComb , col="skyblue" )
273             lines( curveXpos[j] - 0*xVals , yComb , col="skyblue" , lwd=2 )
274         }
275     }
276     # replot the data, in case they are obscured by lines:

```

```

277     points( x , y , cex=1.5 )
278     if ( !is.null(saveName) ) {
279         saveGraph( file=paste(saveName,"PostPredYint",sep=""), type=saveType)
280     }
281 }
282 showPosteriorPredictiveSample = TRUE
283 if (showPosteriorPredictiveSample) {
284     # Function to sample from posterior predictive distribution
285     sample_posterior_predictive <- function(beta0,beta1,sigma, xNew) {
286
287         # Calculate mean for each posterior sample
288         muNew <- beta0 + beta1 * xNew
289
290         # Generate predicted y values (one for each posterior sample)
291         y_predicted <- rnorm(length(muNew), mean = muNew, sd = sigma)
292
293         return(y_predicted)
294     }
295     # Sample from posterior predictive at x=20
296     y_pred_x20 <- sample_posterior_predictive(beta0,beta1,sigma,xNew = 20)
297     # Plot the posterior predictive distribution
298     openGraph(width=8,height=5)
299     histInfo = plotPost(y_pred_x20,main=paste("Posterior Predictive x=20"))
300     if ( !is.null(saveName) ) {
301         saveGraph( file=paste(saveName,"PostPredx20",sep=""), type=saveType)
302     }
303 }
304 }
305
306 #=====

```

The example script:

```

1  # Bayesian Data Analysis 2025, Assignment 3, Exercise B
2  #-----
3  # Optional generic preliminaries:
4  graphics.off() # This closes all of R's graphics windows.
5  rm(list=ls()) # Careful! This clears all of R's memory!
6  #-----
7  # Preliminaries
8  source("openGraphSaveGraph.R")
9  source("plotPost.R")
10 require(rjags)
11 fileNameRoot = "Figures/Mnormal-"
12 graphFileType = "eps"
13 # Load the relevant model into R's working memory:
14 source("Jags-Ymet-Xmet-Mnormal.R")
15
16 #-----
17 # LOAD DATA
18 myData = read.table('../caschools.csv', header = TRUE, sep=',')
19 xName = "stratio" ; yName = "testscr"
20
21 # Summary
22 str(myData)
23 summary(myData)
24 # Hist and scatterplot
25 openGraph(width=12,height=6)
26 layout( matrix( c(1,2,3) , nrow=1 ) )
27 hist(myData[,xName])
28 hist(myData[,yName])
29 plot(myData[,xName], myData[,yName],

```



```

30     xlab = "Student-Teacher Ratio",
31     ylab = "Test Score"
32 )
33 saveGraph(file=paste0(fileNameRoot,"preliminary"),type="eps")
34
35 #-----
36 # Generate the MCMC chain:
37 #startTime = proc.time()
38 mcmcCoda = genMCMC( data=myData , xName=xName , yName=yName ,
39                   numSavedSteps=20000 , saveName=fileNameRoot )
40 #stopTime = proc.time()
41 #duration = stopTime - startTime
42 #show(duration)
43
44 #-----
45 # Display diagnostics of chain, for specified parameters:
46 parameterNames = varnames(mcmcCoda) # get all parameter names
47 for ( parName in parameterNames ) {
48     diagMCMC( codaObject=mcmcCoda , parName=parName ,
49             saveName=fileNameRoot , saveType=graphFileType )
50 }
51
52 #-----
53 # Get summary statistics of chain:
54 summaryInfo = smryMCMC( mcmcCoda ,
55                       compValBeta1=0.0 , ropeBeta1=c(-0.5,0.5) ,
56                       saveName=fileNameRoot )
57 show(summaryInfo)
58 # Display posterior information:
59 plotMCMC( mcmcCoda , data=myData , xName=xName , yName=yName ,
60         compValBeta1=0.0 , ropeBeta1=c(-0.5,0.5) ,
61         pairsPlot=TRUE , showCurve=FALSE ,
62         saveName=fileNameRoot , saveType=graphFileType )
63 #-----

```

appendix

The source code:

```
1 # Bayesian Data Analysis 2025, Assignment 3, Exercise 1
2 #-----
3 source("DBDA2E-utilities.R")
4 #=====
5
6 # Improved initialization function
7 genMCMC = function(data, numSavedSteps=20000, saveName=NULL) {
8   require(rjags)
9   savePath = "Figures/Exercise1/2-"
10
11   # LOAD DATA
12   if (class(data)=="data.frame") {
13     y = myData$y
14   } else {
15     y = data
16   }
17
18   # Do some checking that data make sense:
19   if (any(y != floor(y))) { stop("All y values must be integers.") }
20
21   # Compute length
22   nBags = length(y)
23
24   # Specify the data in a list
25   dataList = list(y = y, nBags = nBags)
26
27   # JAGS MODEL
28   modelString = "
29   model {
30     # Likelihood
31     for (i in 1:nBags) {
32       y[i] ~ dbin(p[b[i]], n[i])
33       n[i] <- round(nContinuous[i])
34       nContinuous[i] ~ dnorm(normMu, normTau)
35       b[i] <- bBernoulli[i] + 1
36       bBernoulli[i] ~ dbern(1-catTheta)
37     }
38     # Priors
39     p[1] ~ dbeta(betaA, betaB) # when b=0
40     p[2] ~ dbeta(1, 1)         # when b=1
41     catTheta ~ dbeta(1, 1)
42     # Constants
43     normMu <- 30
44     normTau <- 1/(normSigma*normSigma)
45     normSigma <- 1.5
46     betaA <- 40
47     betaB <- 10
48   }
49   "
50
51   writeLines(modelString, con="modelExercise1.txt")
52
53   # IMPROVED INITIALIZATION
54   # Create multiple initialization lists to increase chances of convergence
55   initsList = list()
56   for (i in 1:3) { # for each chain
57     initsList[[i]] = list(
58       p = c(runif(1, 0.7, 0.9), runif(1, 0.6, 0.8)), # randomize within reasonable range
```

```

59     catTheta = runif(1, 0.2, 0.4),                # randomize within reasonable range
60     bBernoulli = rbinom(prob=0.7, size=1, n=nBags),
61     nContinuous = pmax(10, rnorm(nBags, mean=30, sd=1.5)) # ensure positive values
62 )
63 }
64
65 # RUN THE CHAINS
66 parameters = c("catTheta", "p[1]", "p[2]")
67 adaptSteps = 2000                                # Increased for better adaptation
68 burnInSteps = 10000                              # Increased for better convergence
69 nChains = 3
70 thinSteps = 1
71 nIter = ceiling((numSavedSteps * thinSteps) / nChains)
72
73 # Create, initialize, and adapt the model with robust error handling
74 tryCatch({
75     jagsModel = jags.model("modelExercise1.txt", data=dataList, inits=initsList,
76                           n.chains=nChains, n.adapt=adaptSteps)
77
78     cat("Burning in the MCMC chain...\n")
79     update(jagsModel, n.iter=burnInSteps)
80
81     cat("Sampling final MCMC chain...\n")
82     codaSamples = coda.samples(jagsModel, variable.names=parameters,
83                               n.iter=nIter, thin=thinSteps)
84
85     if (!is.null(saveName)) {
86         save(codaSamples, file=paste(saveName, "Mcmc.Rdata", sep=""))
87     }
88     return(codaSamples)
89 }, error = function(e) {
90     cat("JAGS Error:", e$message, "\n")
91     cat("Attempting to recover with alternative initialization...\n")
92
93     # Alternative initialization as fallback
94     altInitsList = list()
95     for (i in 1:3) {
96         altInitsList[[i]] = list(
97             p = c(0.75, 0.65),
98             catTheta = 0.3,
99             bBernoulli = rep(0, nBags), # simplified initialization
100             nContinuous = rep(30, nBags) # simplified initialization
101         )
102     }
103
104     # Try again with simpler initialization
105     jagsModel = jags.model("modelExercise1.txt", data=dataList, inits=altInitsList,
106                           n.chains=nChains, n.adapt=adaptSteps*2)
107     update(jagsModel, n.iter=burnInSteps*2)
108     codaSamples = coda.samples(jagsModel, variable.names=parameters,
109                               n.iter=nIter, thin=thinSteps)
110
111     if (!is.null(saveName)) {
112         save(codaSamples, file=paste(saveName, "Mcmc.Rdata", sep=""))
113     }
114     return(codaSamples)
115 })
116 }
117
118 #=====
119 smryMCMC = function( codaSamples , compVal=NULL , rope=NULL , saveName=NULL ) {
120     summaryInfo = NULL

```

```

121 mcmcMat = as.matrix(codaSamples,chains=TRUE)
122 summaryInfo = rbind( summaryInfo ,
123                       "catTheta" = summarizePost( mcmcMat[, "catTheta"] ,
124                                                    compVal=compVal , ROPE=rope ) )
125 if ( !is.null(saveName) ) {
126   write.csv( summaryInfo , file=paste(saveName,"SummaryInfo.csv",sep="") )
127 }
128 show( summaryInfo )
129 return( summaryInfo )
130 }
131
132 #=====
133 plotMCMC = function(codaSamples ,
134                     savePath = "Figures/Exercise1/c-",
135                     graphFileType = "eps" ) {
136
137   source("openGraphSaveGraph.R")
138   source("plotPost.R")
139
140   # EXAMINE THE RESULTS.
141   mcmcChain = as.matrix( codaSamples )
142
143   # Extract the posterior sample from JAGS:
144   catThetaSample = mcmcChain[, "catTheta"]
145   p1Sample = mcmcChain[, "p[1]"]
146   p2Sample = mcmcChain[, "p[2]"]
147   pDiffSample = p1Sample - p2Sample
148
149   # Make a graph using R commands:
150   openGraph(width=10,height=6)
151   #layout( matrix( c(1,3) , nrow=1 ) )
152   histInfo = plotPost( catThetaSample , xlim=c(0,1) , xlab=bquote("theta") )
153   saveGraph(file = paste0(savePath,"catTheta"), type="eps")
154
155   # Make a graph using R commands:
156   openGraph(width=10,height=6)
157   layout( matrix( c(1,2,3) , nrow=1 ) )
158   histInfo = plotPost( p1Sample , xlim=c(0,1) , xlab=bquote("p[1]") )
159   histInfo = plotPost( p2Sample , xlim=c(0,1) , xlab=bquote("p[2]") )
160   histInfo = plotPost( pDiffSample , xlab=bquote("p[1]-p[2]") , compVal=0. , ROPE=c(-0.1,0.1)
161                       )
162   saveGraph(file = paste0(savePath,"p[1]_p[2]"), type="eps")
163 } # Close plotMCMC

```

The power estimate experiment:

```

1 # Jags-Ybin-Nnorm-Bbern-Pbeta-Power.R
2 graphics.off() # This closes all of R's graphics windows.
3 rm(list=ls()) # Careful! This clears all of R's memory!
4 fileNameRoot = "Jags-Ydich-Xnom1subj-MbernBeta-Power-" # for future use
5
6 # Load the functions genMCMC, smryMCMC, and plotMCMC:
7 # (This also sources DBDA2E-utilities.R)
8 source("Jags-Ybin-Nnorm-Bbern-Pbeta.R")
9
10 #-----
11 # Improved power analysis
12 goalAchievedForSample = function(data) {
13   # Try to generate the MCMC chain with error handling
14   tryCatch({
15     mcmcCoda = genMCMC(data=data, numSavedSteps=10000, saveName=NULL)

```

```

16
17 # Check for convergence issues
18 gd = gelman.diag(mcmcCoda)
19 if (any(gd$psrf > 1.1)) {
20   warning("Potential convergence issues detected")
21 }
22
23 # Compute the HDI
24 p1HDI = HDIofMCMC(as.matrix(mcmcCoda[, "p[1]"]))
25 p2HDI = HDIofMCMC(as.matrix(mcmcCoda[, "p[2]"]))
26
27 # Define list for recording results
28 goalAchieved = list()
29
30 # Goal: No overlap between HDIs
31 goalAchieved = c(goalAchieved,
32   "noOverlap"=((p1HDI[2] < p2HDI[1]) || (p2HDI[2] < p1HDI[1])))
33
34 return(goalAchieved)
35 }, error = function(e) {
36   cat("Error in MCMC analysis:", e$message, "\n")
37   # Return NA for this simulation
38   return(list("noOverlap"=NA))
39 }
40 }
41
42 # Main power analysis loop with better data generation
43 powerAnalysis = function() {
44   # Specify hypothetical parameters
45   p = c(0.8, 0.7)
46   catTheta = 0.3
47   normMu = 30
48   normSigma = 1.5
49
50   proportionSuccessList = c()
51
52   # Iterate over sample size (number of bags)
53   for (sampleN in seq(200, 500, by=100)) {
54     nSimulatedDataSets = 10
55     goalCount = 0
56     validSimulations = 0
57
58     for (simIdx in 1:nSimulatedDataSets) {
59       cat("Running simulation", simIdx, "of", nSimulatedDataSets, "for sample size", sampleN,
60         "\n")
61
62       # Generate better data with appropriate checks
63       tryCatch({
64         # Generate random value from hypothesized parameter distribution
65         nContinuous = rnorm(n=sampleN, mean=normMu, sd=normSigma)
66         n = round(pmax(5, nContinuous)) # Ensure n is at least 5
67         bBernoulli = rbinom(n=sampleN, size=1, prob=1-catTheta)
68         b = bBernoulli + 1
69
70         # Generate random data based on parameter value
71         sampleY = sapply(1:length(n), function(i) {
72           rbinom(n=1, size=n[i], prob=p[b[i]])
73         })
74
75         # Verify data validity
76         if (all(is.finite(sampleY)) && all(sampleY >= 0) && all(sampleY <= n)) {
77           # Do Bayesian analysis on simulated data

```

```

77     goalAchieved = goalAchievedForSample(sampleY)
78
79     # Only count valid results
80     if (!is.na(goalAchieved[["noOverlap"]])) {
81         validSimulations = validSimulations + 1
82         if (goalAchieved[["noOverlap"]] == TRUE) {
83             goalCount = goalCount + 1
84         }
85     }
86     } else {
87         cat("Generated invalid data in simulation", simIdx, "\n")
88     }
89     }, error = function(e) {
90         cat("Error in simulation", simIdx, "for sample size", sampleN, ":", e$message, "\n")
91     })
92 }
93
94 # Calculate proportion of successes
95 if (validSimulations > 0) {
96     proportionSuccess = goalCount / validSimulations
97     cat("\nSampleN =", sampleN, "| Success rate:", round(proportionSuccess * 100, 2),
98         "% (", goalCount, "out of", validSimulations, "valid simulations)\n\n")
99
100     proportionSuccessList[as.character(sampleN)] = proportionSuccess
101
102     # Check stopping condition
103     if (proportionSuccess > 0.8 && validSimulations >= nSimulatedDataSets*0.8) {
104         print(proportionSuccessList)
105         cat("Stopping early: goal achieved in over 80% of valid simulations for sample size",
106             sampleN, "\n")
107         break
108     }
109     } else {
110         cat("No valid simulations completed for sample size", sampleN, "\n")
111     }
112
113     return(proportionSuccessList)
114 }
115
116 # Run the power analysis
117 results = powerAnalysis()
118 print(results)

```