

How to use RSGHB

Jeffrey Dumont

21 May 2014

Abstract

This vignette describes the process for specifying, estimating and analyzing the output of a choice model with RSGHB. The document is structured into two case studies using different model structures. The synthetic choice data used in this document (as well as number of other examples using different model structures) can be downloaded from the RSGHB github page - <https://github.com/jeffdumont/RSGHB>

RSGHB modeling file structure

The typical RSGHB model file has the following 4 main sections. We will walk through each of these 4 sections in the examples that follow.

1. Setup and data preparation
2. Setting the controls for model estimation
3. Defining the likelihood function
4. Calling the doHB function to start the estimation process

EXAMPLE 1: MNL Model with Fixed Parameters

In this section, the code for estimating a multinomial logit model with fixed (non-random) parameters is explained. In this example, synthetic respondents were presented with a choice between two travel alternatives - one that is toll free but slower and one that was priced but faster. Each respondent was presented with a panel of 8 choice tasks.

Setup and Data Preparation

```
library(RSGHB)

# load example data
data(choicedata)

# We typically work with one row per choice observations. This isn't
# necessary however but it does lend it self to faster computation of the
# choice likelihoods
head(choicedata)

##      ID thecount tt1 tt2 toll2 asc1 Choice
## 1 8738         1  60  51  1.25     1      1
## 2 8738         2  60  51  0.75     1      1
## 3 8738         3  63  59  0.50     1      2
## 4 8738         4  60  54  0.75     1      1
## 5 8738         5  60  54  0.50     1      2
## 6 8738         6  63  54  0.75     1      1
```

```

# We can then specify any variables from the choicedata data.frame that
# you'd like to use in the utility equations in the likelihood function
# below. These can be any variables within the data or transformations of
# those variables. This example comes from transport so each alternative is
# defined by travel times and toll costs.
TT1 <- choicedata$tt1
TT2 <- choicedata$tt2
TOLL2 <- choicedata$toll2

# Here we specify the choice vectors. Note in this example there are only
# two alternatives. Also, dummyming coding the choice vector is not necessary
# but allows for easier coding of the likelihood.
choice1 <- (choicedata$Choice == 1)
choice2 <- (choicedata$Choice == 2)

# Frequency of choice for the first alternative
table(choice1)

## choice1
## FALSE TRUE
## 3560 6682

# Frequency of choice for the second alternative.
table(choice2)

## choice2
## FALSE TRUE
## 6682 3560

```

Controlling the Estimation Process

There are number of options for controlling the estimation process. Please see the help file for doHB or the final section of this document for more details. Note that a number of controls have default values and do not need to be directly specified if the default is acceptable.

```

# ----- ESTIMATION CONTROL -----

# Setting control list for estimation (see ?doHB for more estimation
# options)

# modelname is used for naming the output files
modelname <- "MNL"

# gVarNamesFixed contains the names for the fixed (non-random) variables in
# your model. This will be used in output and also when displaying iteration
# detail to the screen.
gVarNamesFixed <- c("ASC1", "BTime", "BCost")

# FC contains the starting values for the fixed coefficients.
FC <- c(0, 0, 0)

# ITERATION SETTINGS

# gNCREP contains the number of iterations to use prior to convergence
gNCREP <- 30000
# gNEREP contains the number of iterations to keep for averaging after
# convergence has been reached
gNEREP <- 20000
# gNSKIP contains the number of iterations to do in between retaining draws
# for averaging

```

```

gNSKIP <- 1
# gINFOSKIP controls how frequently to print info about the iteration
# process
gINFOSKIP <- 250

# To simplify the doHB functional call, we put all of the control parameters
# into a single list that can be passed directly to doHB.
control <- list(modelname = modelname, gVarNamesFixed = gVarNamesFixed, FC = FC,
  gNCREP = gNCREP, gNEREP = gNEREP, gNSKIP = gNSKIP, gINFOSKIP = gINFOSKIP)

```

Writing the likelihood function

RSGHB is expecting the user-specified likelihood function to take the parameters fc and b (even if they are not used within the function calculate the likelihood). The fc parameter is a vector of fixed coefficients (they do not vary across individuals in your data). The b parameter is a matrix of individual coefficients which are generated from the random coefficients in the model. In this example, we only focus on the fc vector.

It is important to note that the computation of the likelihood is the most computational taxing part of the estimation process. So coding the likelihood efficiently is essential to reduce run time of the model.

```

# ----- likelihood -----

likelihood <- function(fc, b) {

  # defining the parameters

  # using cc var to index the fc vector simplifies the addition/subtraction of
  # new parameters
  cc <- 1
  ASC1 <- fc[cc]
  cc <- cc + 1
  Btime <- fc[cc]
  cc <- cc + 1
  Btoll <- fc[cc]
  cc <- cc + 1

  # utility functions
  v1 <- ASC1 + Btime * TT1
  v2 <- Btime * TT2 + Btoll * TOLL2

  # mnl probability statement
  p <- (exp(v1) * choice1 + exp(v2) * choice2)/(exp(v1) + exp(v2))

  return(p)
}

```

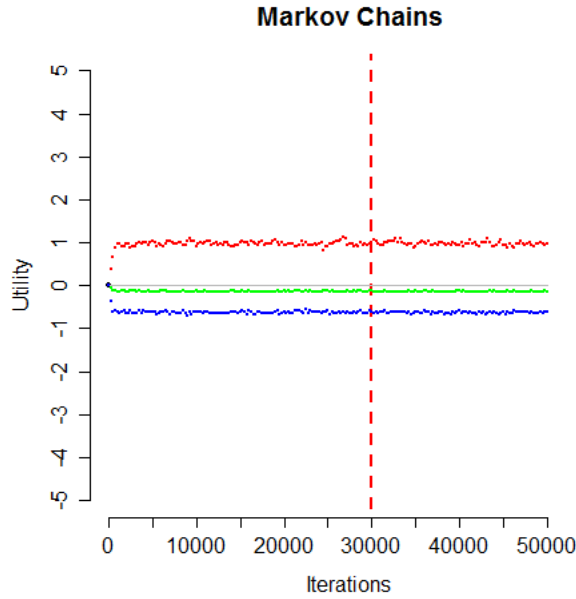
Estimating the model

To start the model estimation process, the analyst needs to call the *doHB* function passing in the *likelihood* function, the *choicedata* data.frame and the *control* list.

```
doHB(likelihood, choicedata, control)
```

RSGHB will first perform a series of diagnostics on your model to catch common errors in model setup. In addition, it will provide you with some basic summary statistics on your choice data and model. Before estimation begins, RSGHB will present you with a confirmation prompt allowing you to cancel the model estimation.

Figure 1: Plotting of the Markov Chains during estimation



During the estimation, current estimates of the markov chains will be plotted to the screen. This plot is updated based on the control parameter *gINFOSKIP* (see figure 1). In addition, it will provide numerical iteration details in the R Console.

Evaluating the output

There are two main output files for this particular model - *log* file and the *_F* file. RSGHB comes with some basic tools for plotting the contents of these files.

The *.log* file contains some statistics that can be used to understand if model convergence has been reached. Because this model contains only fixed coefficients, the log file contains just the root likelihood (RLH) and log-likelihood at each iteration defined by *gINFOSKIP* - see figure 2.

The *_F* file contains the set of fixed (non-random) coefficients for each iteration after convergence of the markov chain.

```
data(Example1_F)

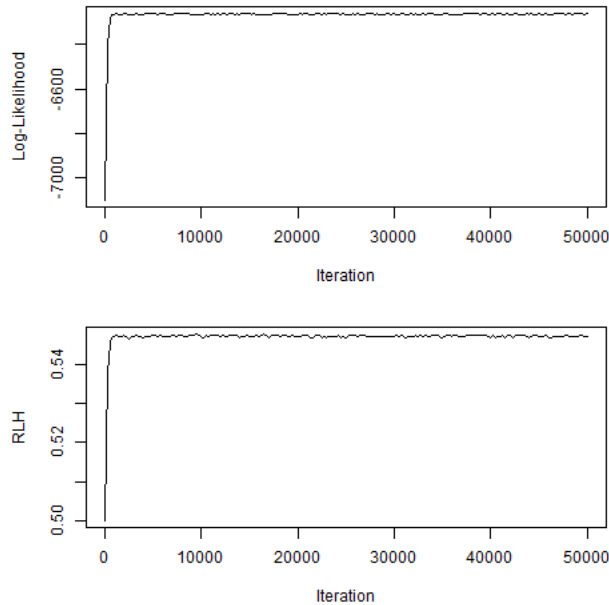
head(Example1_F)

##   iteration  ASC1   BTime  BCost
## 1         1  1.030 -0.1460 -0.6586
## 2         2  1.030 -0.1460 -0.6586
## 3         3  1.030 -0.1460 -0.6586
## 4         4  1.030 -0.1460 -0.6586
## 5         5  1.030 -0.1460 -0.6586
## 6         6  1.027 -0.1451 -0.6579
```

EXAMPLE 2: MNL Model with Random Coefficients

In this section, we expand on the model estimated in EXAMPLE 1 by allowing the coefficients to vary across the individuals in our dataset. This type of model is referred to by many names - Random Coefficients Logit, Random Parameters Logit or Mixed Logit.

Figure 2: Plotting of the Log file



Setup and Data Preparation

The setup and data preparation are very similar to the first model.

```
library(RSGHB)

# load example data
data(choicedata)

# We typically work with one row per choice observations. This isn't
# necessary however but it does lend it self to faster computation of the
# choice likelihoods
head(choicedata)

##      ID thecount tt1 tt2 toll2 asc1 Choice
## 1 8738         1  60  51  1.25    1      1
## 2 8738         2  60  51  0.75    1      1
## 3 8738         3  63  59  0.50    1      2
## 4 8738         4  60  54  0.75    1      1
## 5 8738         5  60  54  0.50    1      2
## 6 8738         6  63  54  0.75    1      1

# We can then specify any variables from the choicedata data.frame that
# you'd like to use in the utility equations in the likelihood function
# below. These can be any variables within the data or transformations of
# those variables. This example comes from transport so each alternative is
# defined by travel times and toll costs.
TT1 <- choicedata$tt1
TT2 <- choicedata$tt2
TOLL2 <- choicedata$toll2

# Here we specify the choice vectors. Note in this example there are only
# two examples. Also, dummyming coding the choice vector is not necessary
# but allows for easier coding of the likelihood.
choice1 <- (choicedata$Choice == 1)
```

```

choice2 <- (choicedata$Choice == 2)

# Frequency of choice for the first alternative
table(choice1)

## choice1
## FALSE TRUE
## 3560 6682

# Frequency of choice for the second alternative.
table(choice2)

## choice2
## FALSE TRUE
## 6682 3560

```

Controlling the Estimation Process

To allow for mixing of the parameters, we need to specify a few more controls to pass into the *doHB* function.

```

# ----- ESTIMATION CONTROL -----

# Setting control list for estimation (see ?doHB for more estimation
# options)

# modelname is used for naming the output files
modelname <- "MNL"

# gVarNamesNormal provides names for the random parameters
gVarNamesNormal <- c("ASC1", "BTime", "BCost")

# gDIST specifies the type of continuous distribution to use for the random
# parameters. gDIST must have an entry for each value in gVarNamesNormal
# The options are: 1. normal 2. log-normal 3. negative log-normal 4. normal
# with all values below zero massed at zero 5. normal with all values
# greater than zero massed at zero 6. Johnson SB with a specified min and
# max

# In this example, we use normal distributions for all 3 of the parameters.
gDIST <- c(1, 1, 1)

# svN contains the starting values for the means of the normal distributions
# for each of the random parameters
svN <- c(0, 0, 0)

# ITERATION SETTINGS

# gNCREP contains the number of iterations to use prior to convergence
gNCREP <- 30000
# gNEREP contains the number of iterations to keep for averaging after
# convergence has been reached
gNEREP <- 20000
# gNSKIP contains the number of iterations to do in between retaining draws
# for averaging
gNSKIP <- 1
# gINFOSKIP controls how frequently to print info about the iteration
# process

```

```
gINFOSKIP <- 250

# To simplify the doHB functional call, we put all of the control parameters
# into a single list that can be passed directly to doHB.
control <- list(modelname = modelname, gVarNamesNormal = gVarNamesNormal, gDIST = gDIST,
  svN = svN, gNCREP = gNCREP, gNEREP = gNEREP, gNSKIP = gNSKIP, gINFOSKIP = gINFOSKIP)
```

Writing the likelihood function

To introduce mixing into the model, we switch from using fc vector to using the b matrix in the coding of the likelihood. The b matrix contains the individual conditionals for the sample-level random coefficients. The matrix b has one row per observation (an individual's coefficients are repeated across their choice observations automatically by RSGHB) and one column for each of the random parameters.

```
likelihood <- function(fc, b) {

  # the change from using fc to b is the only change in the likelihood
  # function required to allow for mixing.
  cc <- 1
  ASC1 <- b[, cc]
  cc <- cc + 1
  Btime <- b[, cc]
  cc <- cc + 1
  Btoll <- b[, cc]
  cc <- cc + 1

  v1 <- ASC1 + Btime * TT1
  v2 <- Btime * TT2 + Btoll * TOLL2

  p <- (exp(v1) * choice1 + exp(v2) * choice2)/(exp(v1) + exp(v2))

  return(p)
}
```

Section 4: Estimating the model

Again, to start the model estimation process, the analyst needs to call the *doHB* function passing in the *likelihood* function, the *choicedata* data.frame and the *control* list.

```
doHB(likelihood, choicedata, control)
```

Evaluating the output

As in the first example, current estimates of the markov chains will be plotted to the screen - see figure 3. In this model, these represent the means of the underlying normals for the random parameters.

There are more output files for this model. RSGHB comes with some basic tools for plotting the contents of these files.

The *.log* file contains some statistics that can be used to understand if model convergence has been reached. Because this model includes random coefficients, the log file now contains the average variance and parameter root mean square (RMS) at each iteration - see figure 4.

The *_A* file contain the sample-level means of the underlying normal at each iteration.

```
data(Example2_A)

head(Example2_A)
```

Figure 3: Plotting of the Markov Chains during estimation

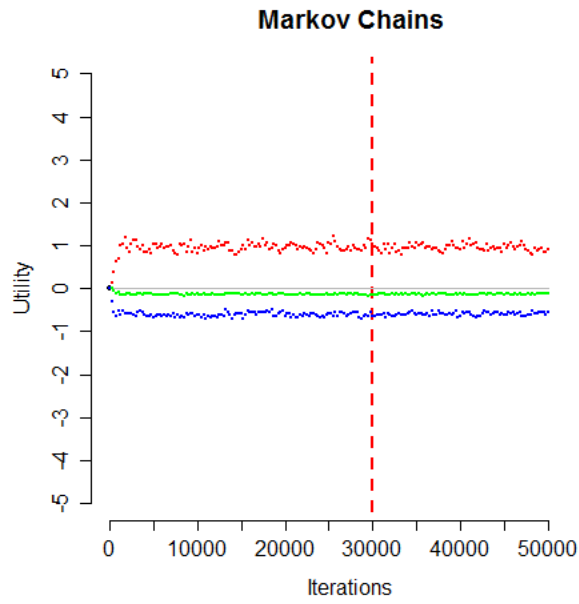
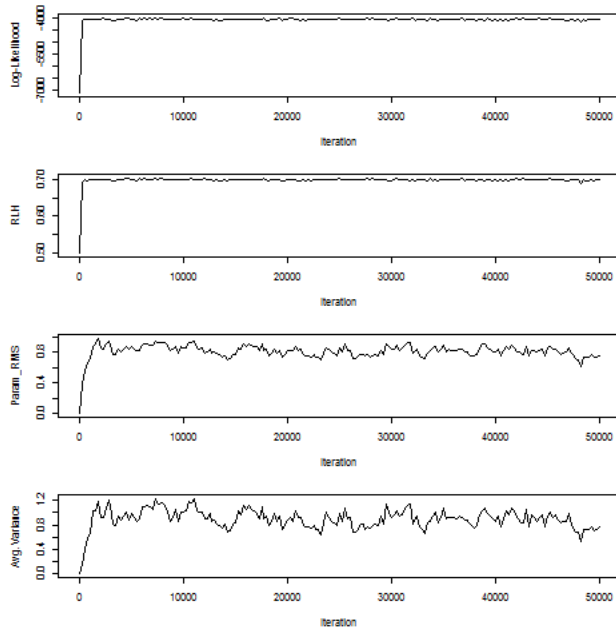


Figure 4: Plotting of the Log file




```
## iteration ASC1 BTime BCost
## 1 1 0.9502 -0.1249 -0.5975
## 2 2 0.9192 -0.1345 -0.6017
## 3 3 0.9334 -0.1276 -0.5895
## 4 4 0.9543 -0.1386 -0.5932
## 5 5 0.9488 -0.1187 -0.5895
## 6 6 0.9399 -0.1178 -0.5912
```

The `.B` file contains the average across iterations of the individual level draws for the underlying normals for the random parameters. The `.Bsd` file provides the standard deviations of those individual draws.

```
data(Example2_B)

head(Example2_B)

## Respondent ASC1 BTime BCost
## 1 8738 1.2066 -0.01700 -0.6707
## 2 8740 1.1857 -0.01575 -0.5232
## 3 8741 0.7790 -0.57719 -0.4940
## 4 8742 1.1593 -0.27220 -0.6376
## 5 8744 0.5109 -0.02025 -0.5083
## 6 8745 1.0068 -0.29565 -0.5575

data(Example2_Bsd)

head(Example2_Bsd)

## Respondent ASC1 BTime BCost
## 1 8738 0.6501 0.1616 0.4068
## 2 8740 0.6821 0.1693 0.4069
## 3 8741 0.6917 0.2080 0.4230
## 4 8742 0.6381 0.1567 0.4002
## 5 8744 0.7009 0.1650 0.4296
## 6 8745 0.6260 0.1676 0.4018
```

The `.C` file contains the average across iterations of the individual level draws for the random parameters including the appropriate transformations. The `.C` file also contains individual-specific measures of model fit (Root Likelihoods). The `.Csd` file provides the standard deviations of those individual draws. These two files are equivalent to the conditional distributions from models estimated using Maximum Simulated Likelihood methods.

```
data(Example2_C)

head(Example2_C)

## Respondent RLH ASC1 BTime BCost
## 1 8738 0.6762 1.2066 -0.01700 -0.6707
## 2 8740 0.6326 1.1857 -0.01575 -0.5232
## 3 8741 0.8420 0.7790 -0.57719 -0.4940
## 4 8742 0.5145 1.1593 -0.27220 -0.6376
## 5 8744 0.4781 0.5109 -0.02025 -0.5083
## 6 8745 0.4920 1.0068 -0.29565 -0.5575

data(Example2_Csd)

head(Example2_Csd)

## Respondent ASC1 BTime BCost
## 1 8738 0.6501 0.1616 0.4068
## 2 8740 0.6821 0.1693 0.4069
## 3 8741 0.6917 0.2080 0.4230
## 4 8742 0.6381 0.1567 0.4002
```

```
## 5      8744 0.7009 0.1650 0.4296
## 6      8745 0.6260 0.1676 0.4018
```

The `_D` file contains a row-based representation of the sample covariance for each iteration. Note the use below of the `xpnd` function to convert to a matrix representation of the sample covariance.

```
data(Example2_D)

head(Example2_D)

##   iteration   ASC1 ASC1xBTime ASC1xBCost      BTime BTimexBCost BCost
## 1         1 0.1996  -0.02085   0.06611 -0.0122684   0.03471 0.1830
## 2         2 0.2068  -0.02099   0.06575 -0.0053138   0.03063 0.1854
## 3         3 0.2206  -0.01673   0.06778 -0.0008272   0.03580 0.1873
## 4         4 0.2136  -0.02086   0.06905  0.0055351   0.03120 0.1801
## 5         5 0.2376  -0.01937   0.06593 -0.0054135   0.02932 0.1752
## 6         6 0.2234  -0.02707   0.07341 -0.0128438   0.03041 0.1791

# building the covariance matrix

covMat <- xpnd(colMeans(Example2_D[-1]))

rownames(covMat) <- c("ASC1", "BTime", "BCost")
colnames(covMat) <- c("ASC1", "BTime", "BCost")

covMat

##           ASC1      BTime      BCost
## ASC1    0.44636 -0.06659 -0.05475
## BTime -0.06659  0.07915  0.03123
## BCost -0.05475  0.03123  0.17248
```

RSGHB Control Parameters

Here is a list of the user-specified control parameters.

gVarNamesNormal - A vector of character-based names for the random parameters.

Default: NULL

gVarNamesFixed - A vector of character-based names for the fixed parameters.

Default: NULL

gDIST - A vector of integers (1-6) which indicate which type of distribution should be applied to the random parameters - 1 = Normal, 2 = Postive Log-Normal, 3 = Negative Log-Normal, 4 = Positive Censored Normal, 5 = Negative Censored Normal, 6 = Johnson SB. There should be an element for each name in `gVarNamesNormal`.

Default: NULL

FC - A vector of starting values for the fixed coefficients. There should be an element for each name in `gVarNamesFixed`.

Default: NULL

svN - A vector of starting values for the means of the underlying normals for the random parameters. There should be an element for each name in `gVarNamesNormal`.

Default: NULL

gNCREP - Number of burn-in iterations to use prior to convergence.

Default: 100000

gNEREP - Number of iterations to keep for averaging after convergence has been reached.

Default: 100000

gNSKIP - Number of iterations in between retaining draws for averaging.

Default: 1

gINFOSKIP - Number of iterations in between printing/saving information about the iteration process.
Default: 250

modelname - The model name which is used for creating output files.
Default: `paste("HBModel",round(runif(1)*10000000,0),sep="")`

gSIGDIG - The number of significant digits for reporting purposes.
Default: 10

priorVariance - The amount of prior variance assumed.
Default: 2.0

pvMatrix - A custom prior covariance matrix can be used in estimation. If specified in the control list, the custom matrix will override the default prior covariance matrix used by RSGHB. The prior covariance matrix needs to be a matrix object and of the correct size -

`length(gVarNamesNormal) x length(gVarNamesNormal)`

degreesOfFreedom - Additional degrees of freedom for the prior covariance matrix (not including the number of parameters.
Default: 5

rho - The initial proportionality fraction for the jumping distribution for the Metropolis-Hastings algorithm for the random parameters. This fraction is adjusted by the program after each iteration to attain an acceptance rate of about 0.3.
Default: 0.1

rhoF - The proportionality fraction for the jumping distribution for the Metropolis-Hastings algorithm for the fixed parameters. Unlike rho, this value is not adjusted as the markov chain proceeds.
Default: 0.0001

targetAcceptanceNormal - The target acceptance rate in the Metropolis-Hastings algorithm for the random parameters.
Defaults: 0.3

targetAcceptanceFixed - The target acceptance rate in the Metropolis-Hastings algorithm for the fixed parameters.
Defaults: 0.3

gFULLCV - A number that indicates if a full or independent covariance structure should be used for the random parameters. A value of 1 indicated full and 0 for an independent structure.
Default: 1

gMINCOEF - A vector of minimums for the Johnson SB distributions. If Johnson SB is used, each random parameter needs an element but only the elements that correspond to a JSB in gDIST are used.
Default: 0

gMAXCOEF - Like gMINCOEF but for the maximum of the Johnson SB distribution.
Default: 0

gStoreDraws - A boolean value to store the draws for the individual level coefficients.
Default: F

gSeed - The random seed.
Default: 0

constraintsNorm - This is a list of monotonic constraints to be applied during estimation. The structure of the constraints is `c(param1number - inequality - param2number)`. For constraints relative to 0, use 0 instead of the param2number. For the inequality, use 1 for < and 2 for >. Example

`constraintsNorm <- list(c(5,1,0),c(6,1,5),c(7,1,6),c(8,1,7))`

would constrain the 5th parameter \leq 0, the 6th parameter \leq 5th parameter, the 7th parameter \leq the 6th parameter, etc.
Default: NULL

nodiagnosics - If set to TRUE, the diagnostic report will not be reported to the screen with a prompt to continue. This makes batch processing easier to implement.
Default: FALSE

fixedA - This allows the analyst to fix means of the underlying normal distribution of random variables to certain values as opposed to estimating them. This would be important for example in an error components logit model or an integrated choice and latent variable model. The format for this input is a vector of length equal to the number of random parameters. Use NA for variables that should be estimated, example:

```
fixedA = c(NA, NA, NA, NA, NA, NA, NA, 0)
```

In this case, the mean of the underlying normal for the 8th random variable would be fixed to 0.

fixedD - This allows the analyst to fix the variance of the underlying normal distribution of the random variables to certain values as opposed to estimating them. This would be important for example in an integrated choice and latent variable model. The format for this input is a vector of length equal to the number of random parameters. Use NA for variables that should be estimated, example:

```
fixedD = c(NA, NA, NA, NA, NA, NA, NA, 1)
```

In this case, the variance of the underlying normal for the 8th random variable would be fixed to 1.

RSGHB Output Files

RSGHB generates a number of output files which can be used to evaluate the model.

A file - The A file contain the sample-level means of the underlying normal at each iteration.

B file, Bsd file - The B file contains the average across iterations of the individual level draws for the underlying normals for the random parameters. The Bsd file provides the standard deviations of those individual draws.

C file, Csd file - The C file contains the average across iterations of the individual level draws for the random parameters including the appropriate transformations. The Csd file provides the standard deviations of those individual draws. These two files are equivalent to the conditional distributions from models estimated using Maximum Simulated Likelihood methods.

D file - This file contains a row-based representation of the sample covariance for each iteration.

F file - This file contains the set of fixed (non-random) parameters for each iteration after convergence.

Log file - This contains some statistics that can be used to understand if model convergence has been reached.

PV Matrix - This file contains the prior covariance matrix that was assumed during the estimation of the model.