

# Probit regression for ordinal data: MAR

---

Setting:

- For each observations, we have 3 associated variables

X1: Continuous variable which we will generate from  $Unif(0, 3)$  X2: Nomial variable with 3 possible outcomes  $\{1, 2, 3\}$  which we will generate from  $Mult(1, [0.2, 0.3, 0.5])$  Y: Ordinal variable with 5 levels,  $\{1, 2, 3, 4, 5\}$  generated using the following process

$$\epsilon_i \sim N(0, 1) \quad z_i = \beta^T X_i + \epsilon \quad \text{where } X = [X_1, X_2 == 1, X_2 == 2, X_2 == 3] \quad g(z_i) = Y_i.$$

Here  $\beta = [-3, -5, 5, 10]$  and function  $g$  will be the binning function which will bin data into 5 different bins corresponding to 5 possible ordinal outcome.

- There will be data missing at random (MAR) in Y. The probability of missingness is govern by  $\text{logit}(w_1^T x_{i1} + w_2^T x_{i2} - 1)$  where  $w_1 = -0.1$  and  $w_2 = [-0.1, -0.1, -0.1]$ . This results in 24% of missing data in the last feature.
- We will try to model Y conditioned on other variables (X1, and X2) using probit regression with latent variable Z with rank likelihood on parameter Z.
- We have two unknown parameters  $\beta$  and  $z_i$  which will be sampled from the full conditional posterior distribution using blocked gibbs sampling.

Summary on modelling process

$$\epsilon_i \sim N(0, 1) \quad z_i = \beta^T X_i + \epsilon \quad z_i \in R(Y) \quad \text{where } R(Y) = \{z_i: z_i > z_j \text{ if } Y_i > Y_j \text{ and } z_i < z_j \text{ if } Y_i < Y_j\}$$

```
# Data generating process
set.seed(0)
n = 100
beta = c(-3, -5, 5, 10)

# noise term
epsilon = rnorm(n, mean = 0, sd = 1)

# X1
X1 = runif(n, 0, 3)

# X2
X2 = t(rmultinom(n, size = 1, prob = c(0.2, 0.3, 0.5)))

# X
X = cbind(X1, X2)
colnames(X) <- c('X1', 'X2_cat1', 'X2_cat2', 'X2_cat3')

# Z
Z = X %*% beta + epsilon

# Cut-off points and Y
g = quantile(Z, probs = c(0.2, 0.4, 0.6, 0.8))
Y = rep(NA, n)
Y[Z < g[1]] = 1
Y[Z >= g[1] & Z < g[2]] = 2
Y[Z >= g[2] & Z < g[3]] = 3
```

```
Y[Z>=g[3] & Z<g[4]] = 4
Y[Z>=g[4]] = 5
Z_original = Z
Y_original = Y
```

Generate missingness in data

```
# Define parameter of logistic function
w1 = -0.1
w2 = c(-0.1, -0.1, -0.1)

# Calculate probability of missingness of features 3
prob = w1*X1 + apply(t(w2*t(X2)), MARGIN = 1, FUN = sum) - 1
prob = 1/(exp(-prob)+1)

# Indicator for X3miss
indicator = rbernoulli(n = n, p = prob)
Y[indicator] = NA
Z[indicator] = NA
```

Prior specifications:

$$\beta \sim \text{multiN}(0, n(X^T X)^{-1})$$

Blocked Gibbs Sampling here consists of two major steps

1. Sample new  $\beta$  from its full conditional

$$\beta \mid z, X, Y, z \in R(Y) \sim \text{multiN}\left(\frac{n}{n+1}(X^T X)^{-1} X^T z, \frac{n}{n+1}(X^T X)^{-1}\right)$$

2. for each  $i$ , using inverse cdf method, sample new  $z_i$  from its full conditional which is a truncated normal distribution:

$$z_i \mid \beta, X, Y, z_i \in R(Y) \sim N(\beta^T x_i, 1) * I\{z_i \in (a, b)\}$$

where  $a = \max(z_j \text{ for } Y_j < Y_i)$   $b = \min(z_j \text{ for } Y_j > Y_i)$

For observations where the target is missing, we need not condition on  $z_i \in R(Y)$  and the resulting full conditional distribution becomes unconstrained normal distribution.

Note that in the gibbs sampling process, we force the first threshold ( $g_1$ ) to be at the true value in the data generating process otherwise the parameters will be unidentifiable (we can definitely find infinite combinations of weights parameter to order the outcome according to their target  $Y$ : scaling, shifting, etc.)

```
# Blocked Gibbs Sampling
set.seed(2)

# Prior Parameter
g1 = g[1] #Fix g1 the first cutoff

# Imputation for first trial
df <- data.frame(cbind(X1, apply(t(t(X2)*c(1,2,3)),1,sum), Y))
colnames(df) <- c('X1', 'X2', 'Y')
df$Y <- as.factor(df$Y)
df$X2 <- as.factor(df$X2)
mod <- multinom(Y~., data=df[!is.na(df$Y),])

## # weights: 25 (16 variable)
## initial value 122.317281
## iter 10 value 54.940347
```

```

## iter 20 value 34.715949
## iter 30 value 33.134798
## iter 40 value 32.304953
## iter 50 value 32.144763
## iter 60 value 31.984214
## iter 70 value 31.930541
## iter 80 value 31.912062
## iter 90 value 31.888047
## iter 100 value 31.768083
## final value 31.768083
## stopped after 100 iterations

Y_i = Y
Y_i[indicator] = predict(mod, newdata = df[is.na(df$Y),], 'class')
Z_i = g1 + Y_i - 1.5
variance_beta = (n/(n+1))*solve(t(X)%*%X)
mean_beta_hat = (n/(n+1))*solve(t(X)%*%X)%*%t(X)

# Initialize the sampling matrix
S = 30000
SAMPLED_Z = matrix(nrow=S,ncol=n)
SAMPLED_Y = matrix(nrow = S, ncol = n)
BETA = matrix(nrow=S,ncol=4)

for (round in 1:S) {
  # Step 1: Sample Beta
  mean_beta = mean_beta_hat%*%Z_i
  beta_sampled = rmvnorm(mean = mean_beta,
                        V = variance_beta, method = 'choleski')
  BETA[round,] <- beta_sampled

  # Step 2: Sample Z using inverse cdf approach
  for (i in 1:n) {
    # Get the lower and upper bound (a, b) of truncated normal
    a = max(-Inf, Z_i[Y_i<Y[i]], na.rm = TRUE)
    b = min(Z_i[Y_i>Y[i]], Inf, na.rm = TRUE)

    # Force the lowest cutoff to be at g1
    if(indicator[i] == FALSE){
      if (Y_i[i] == 1) {
        b = g1
      }else if (Y_i[i] == 2) {
        a = g1
      }
    }
    # Sample using inverse cdf
    ez = t(beta_sampled)%*%X[i,]
    u = runif(1, pnorm(a - ez), pnorm(b-ez))
    Z_i[i] = ez + qnorm(u)
    if (indicator[i]==TRUE) {
      # Impute Y for missing values
      Y_i[i] = -1
      if (Z_i[i] < g1) {
        Y_i[i] = 1
      }
    }
  }
}

```

```

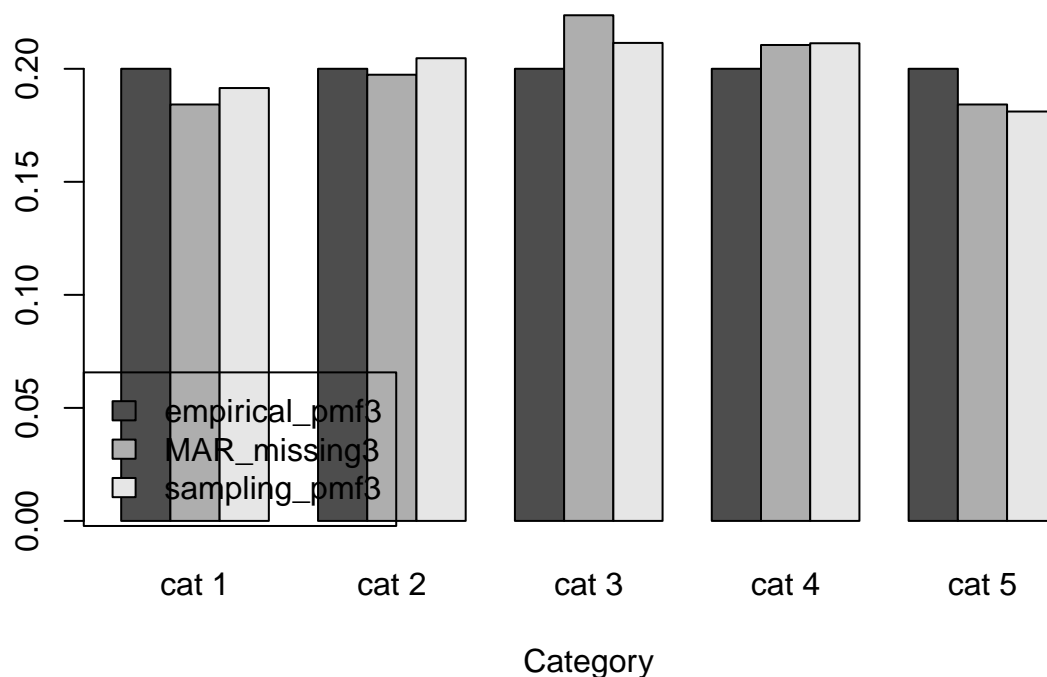
    }else if (Z_i[i] < min(Z_i[Y_i==3],na.rm = TRUE)) {
      Y_i[i] = 2
    }else if (Z_i[i] < min(Z_i[Y_i==4],na.rm = TRUE)){
      Y_i[i] = 3
    }else if (Z_i[i] < min(Z_i[Y_i==5],na.rm = TRUE)){
      Y_i[i] = 4
    }else{
      Y_i[i] = 5
    }
  }
}
SAMPLED_Z[round,] <- Z_i
SAMPLED_Y[round,] <- Y_i
}

burnin = 10000
thining = 100
# Imputation accuracy
empirical_pmf3 = table(Y_original)/n
MAR_missing3 = table(Y[!indicator])/(n-sum(indicator))
sampling_pmf3 = table(SAMPLED_Y[seq(burnin, dim(BETA)[1], thining),])/
  sum(table(SAMPLED_Y[seq(burnin, dim(BETA)[1], thining),]))

df3 = rbind(empirical_pmf3, MAR_missing3, sampling_pmf3)
colnames(df3)<- c('cat 1', 'cat 2', 'cat 3', 'cat 4', 'cat 5')
barplot(df3, xlab = 'Category', beside = TRUE,
  legend = TRUE, args.legend=list(x='bottomleft'),
  main = 'Blocked Gibbs Sampling Assessment: Marginal Y pmf')

```

### Blocked Gibbs Sampling Assessment: Marginal Y pmf



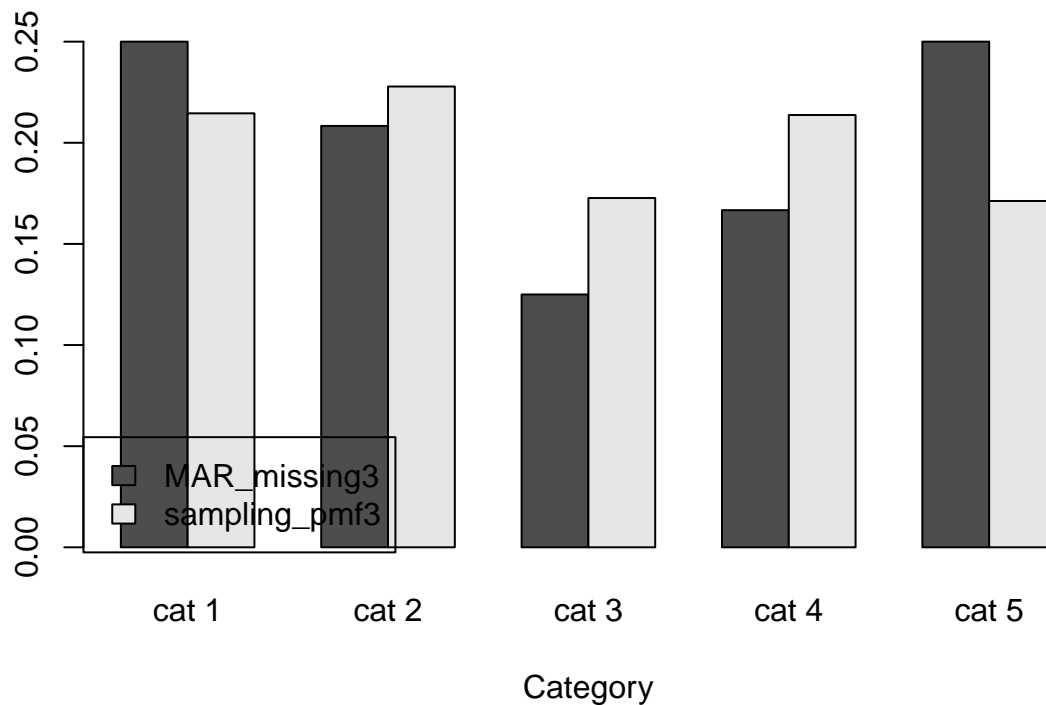
```

# Imputation accuracy
MAR_missing3 = table(Y_original[indicator])/sum(indicator)
sampling_pmf3 = table(SAMPLED_Y[seq(burnin, dim(BETA)[1], thinning),indicator])/
  sum(table(SAMPLED_Y[seq(burnin, dim(BETA)[1], thinning), indicator]))

df3 = rbind(MAR_missing3, sampling_pmf3)
colnames(df3)<- c('cat 1', 'cat 2', 'cat 3', 'cat 4', 'cat 5')
barplot(df3, xlab = 'Category', beside = TRUE,
  legend = TRUE, args.legend=list(x='bottomleft'),
  main = 'Blocked Gibbs Sampling Assessment: Missing Y pmf')

```

### Blocked Gibbs Sampling Assessment: Missing Y pmf



```

# Imputation accuracy
true_label = Y_original[indicator]
sampled_label = SAMPLED_Y[seq(burnin, dim(BETA)[1], thinning),indicator]
mean(t(sampled_label) == true_label)

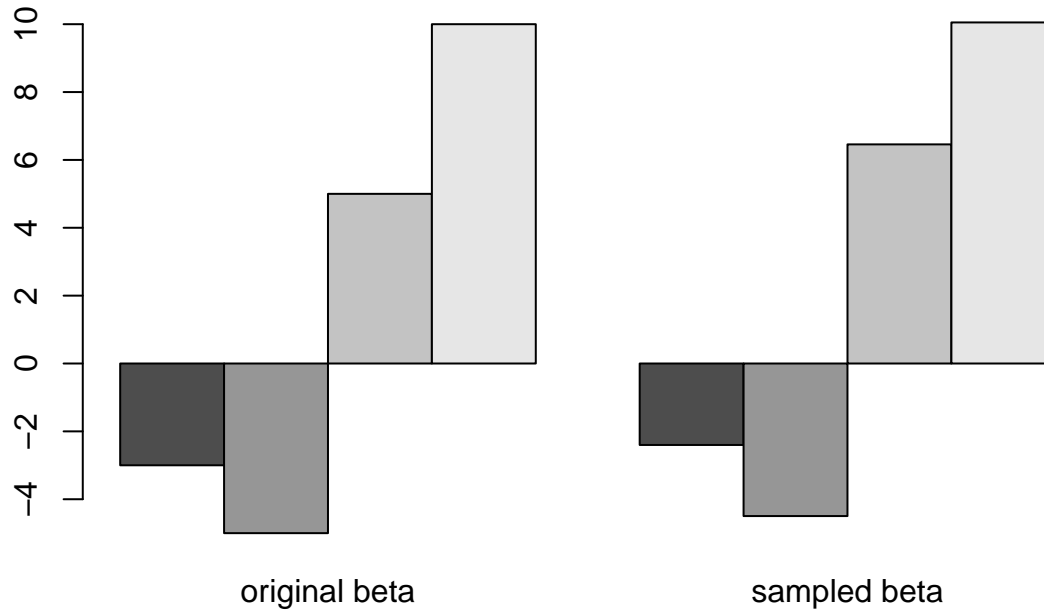
## [1] 0.6735075

# Check posterior expectation of beta
sampling_beta = apply(BETA[seq(burnin, dim(BETA)[1], thinning),], MARGIN = 2, mean)
df_beta= cbind(beta, sampling_beta)
colnames(df_beta)<- c('original beta', 'sampled beta')
barplot(df_beta, beside = TRUE,
  legend = TRUE, main = 'Blocked Gibbs Sampling Assessment: Beta')

# Check cut off points
g1 = apply(SAMPLED_Z[, Y == 1], MARGIN = 1, max, na.rm = TRUE)
g2 = apply(SAMPLED_Z[, Y == 2], MARGIN = 1, max, na.rm = TRUE)
g3 = apply(SAMPLED_Z[, Y == 3], MARGIN = 1, max, na.rm = TRUE)
g4 = apply(SAMPLED_Z[, Y == 4], MARGIN = 1, max, na.rm = TRUE)

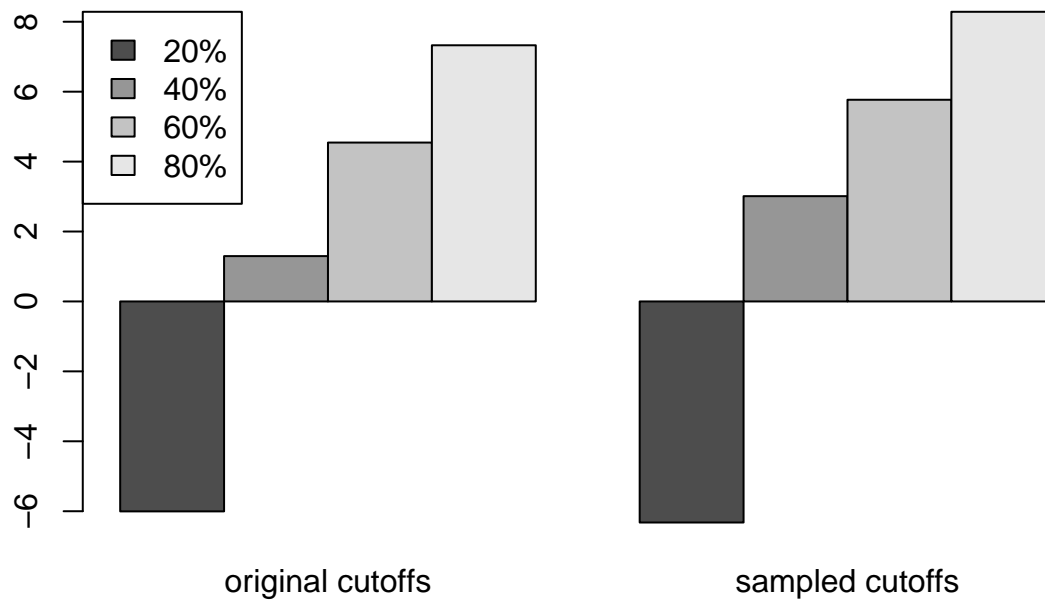
```

## Blocked Gibbs Sampling Assessment: Beta

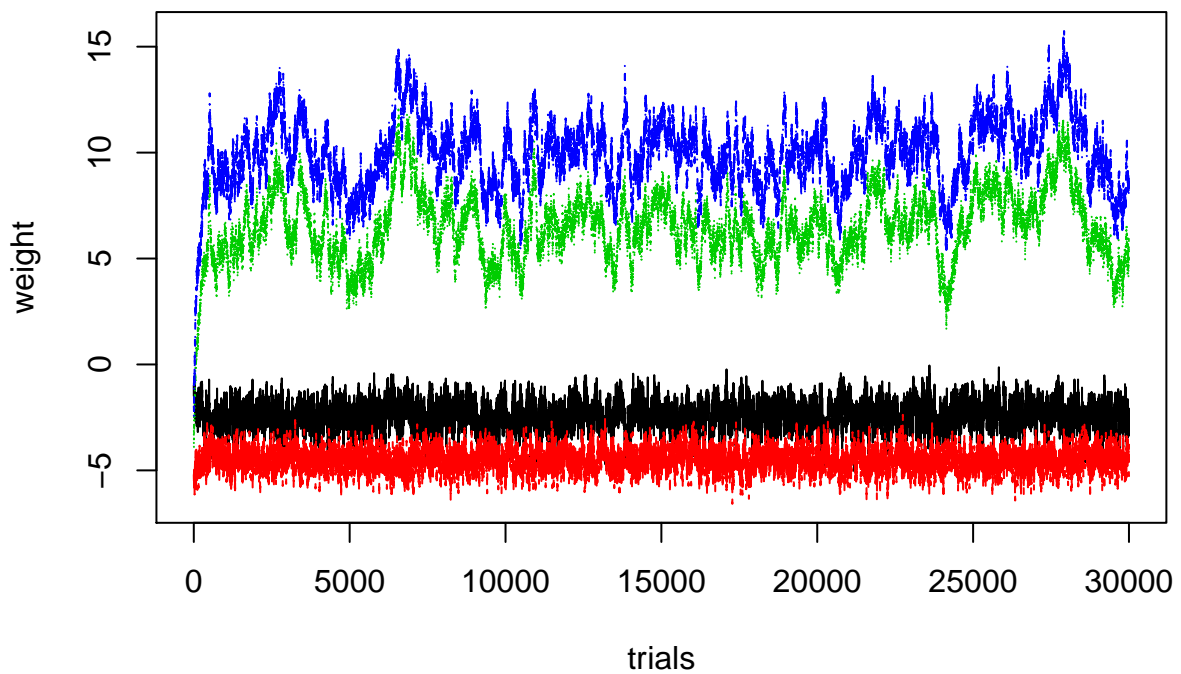


```
df_g= cbind(g, c(mean(g1[seq(burnin, dim(BETA)[1], thinning)]),
                 mean(g2[seq(burnin, dim(BETA)[1], thinning)]),
                 mean(g3[seq(burnin, dim(BETA)[1], thinning)]),
                 mean(g4[seq(burnin, dim(BETA)[1], thinning)])))
colnames(df_g)<- c('original cutoffs', 'sampled cutoffs')
barplot(df_g, beside = TRUE,
        legend = TRUE, main = 'Blocked Gibbs Sampling Assessment: Cutoffs',
        args.legend=list(x='topleft'))
```

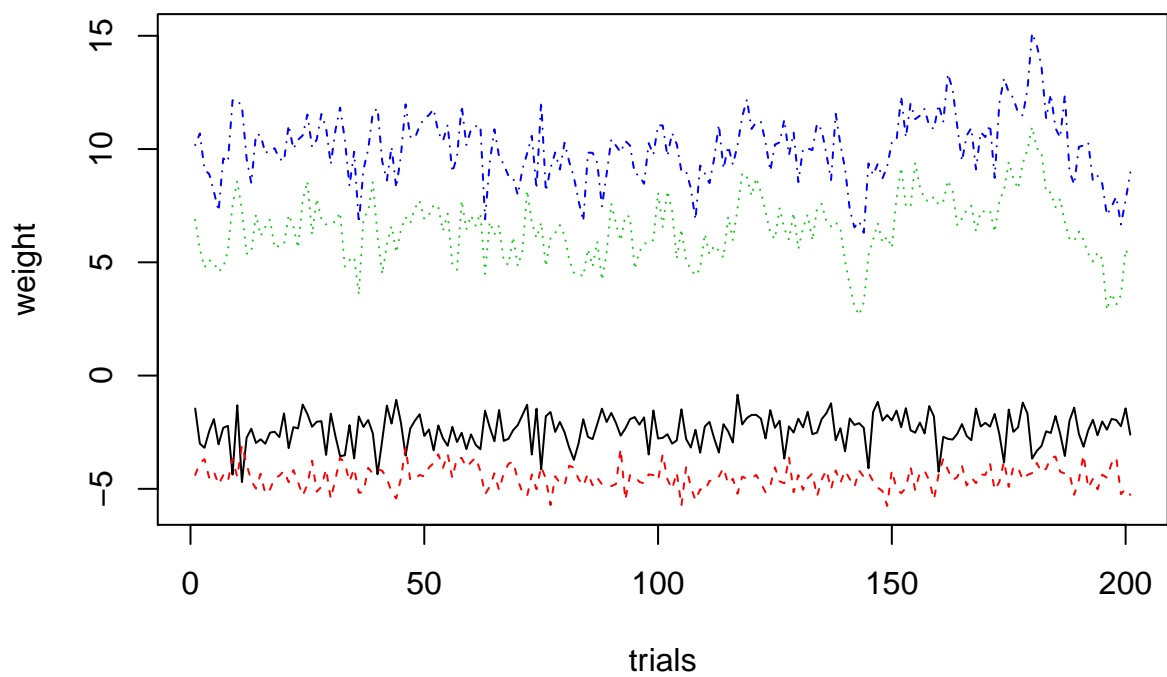
### Blocked Gibbs Sampling Assessment: Cutoffs



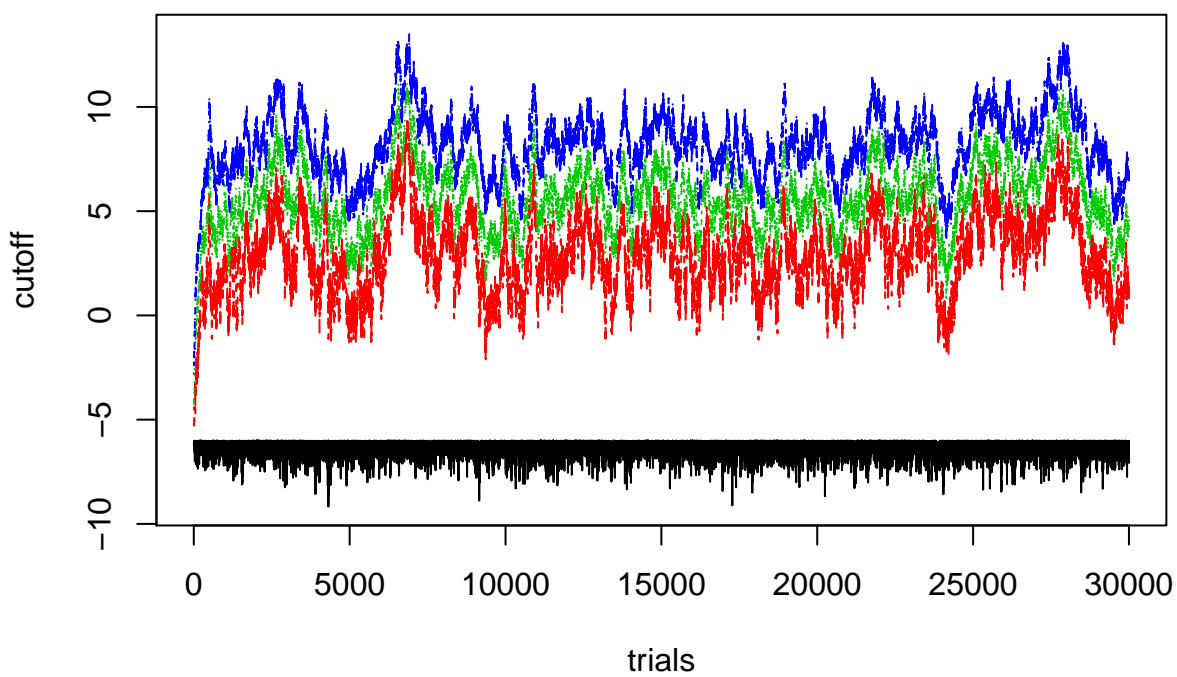
### Checking stability of sampled Beta



### Checking stability of sampled Beta: thinning



### Checking stability of sampled Cutoffs





### Checking stability of sampled Cutoffs: thinning

