

Présentation du contexte

Objectifs

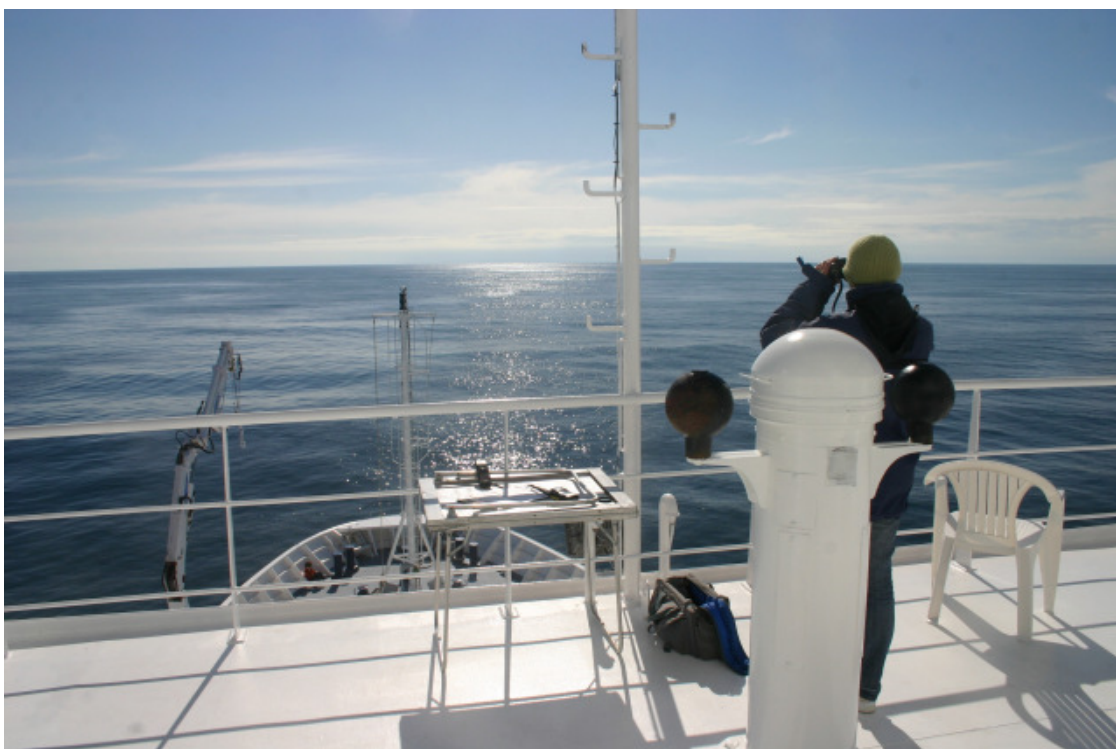
- Ajuster un modèle linéaire généralisé sur des données de comptages
- Vérifier l'ajustement du modèle grâce à la technique du "posterior predictive check"
- Changer la vraisemblance des données pour tenir compte de la surdispersion
- Ajouter des effets aléatoires
- Ajouter un effet spatial

Estimer la taille de groupe du dauphin commun dans le golfe de Gascogne

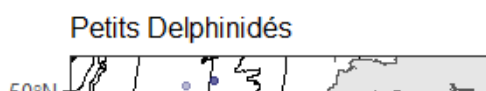
Le dauphin commun (*Delphinus delphis*) est l'espèce la plus abondante de mammifères marins dans le Golfe de Gascogne (Laran et al, 2017).



Depuis 2004, les campagnes océanographiques de l'Ifremer sur le navire Thalassa permettent de recenser et dénombrer les dauphins communs.

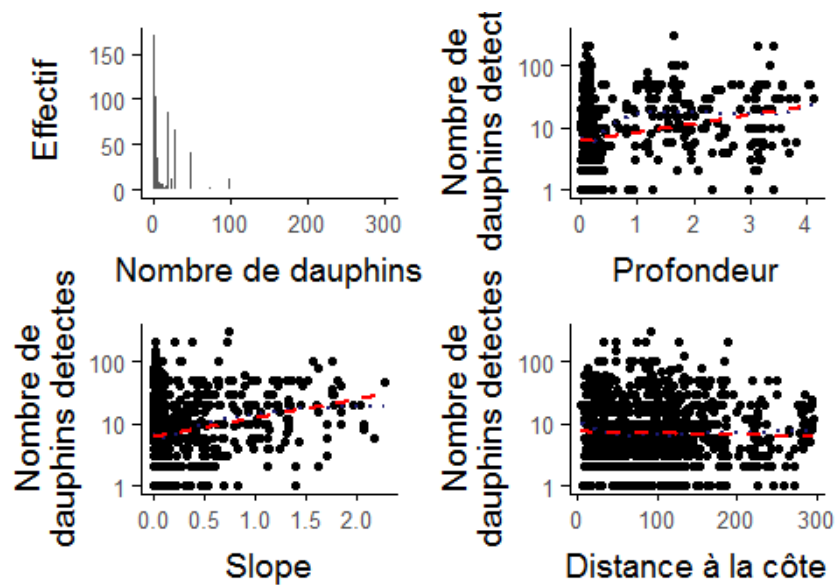


Les données collectées représentent 1269 observations de petits delphinidés (principalement des dauphins communs mais aussi quelques dauphins bleu et blanc *Stenella coeruleoalba*).



Question: Quels sont les facteurs déterminants de l'abondance de dauphins dans le Golfe de Gascogne?

Le but de l'analyse est ici de corrélérer le nombre de dauphins (la variable réponse) à un ensemble de covariables environnementales (par exemple la profondeur/bathymétrie, *etc.*) afin d'expliquer les variations observées.



Les données y_i sont des entiers naturels et la première approche est de considérer un modèle poissonien

$$y_i^{\text{obs}} \sim \mathcal{P}(\lambda_i)$$

où $\log(\lambda_i) = \beta_0 + \beta_1 \times \text{Bathy}_i + \dots$

Début de l'analyse

Charger les données et les bibliothèques de fonctions nécessaires à l'analyse

```
In [1]: lapply(c("ggplot2", "ggthemes", "sp", "sf", "loo"), library, character.only=TRUE)

### theme for graphics
theme_set(theme_bw(base_size = 20))

### clean up
rm(list = ls())
getwd()
WorkDir <- paste(getwd(), "dauphin_v2", sep = "/")
DataDir <- paste(WorkDir, "data", sep = "/")
OutDir <- paste(WorkDir, "output", sep = "/")

### load the data
load(paste(DataDir, "Filtered_data.RData", sep = "/"))
```

Linking to GEOS 3.6.1, GDAL 2.2.3, PROJ 4.9.3

This is loo version 2.0.0.

**NOTE: As of version 2.0.0 loo defaults to 1 core but we recommend using as many as possible. Use the 'cores' argument or set options(mc.cores = NUM_CORES) for an entire session. Visit mc-stan.org/loo/news for details on other changes.

1. 'ggplot2' 'stats' 'graphics' 'grDevices' 'utils' 'datasets' 'methods' 'base'
2. 'ggthemes' 'ggplot2' 'stats' 'graphics' 'grDevices' 'utils' 'datasets' 'methods' 'base'
3. 'sp' 'ggthemes' 'ggplot2' 'stats' 'graphics' 'grDevices' 'utils' 'datasets' 'methods' 'base'
4. 'sf' 'sp' 'ggthemes' 'ggplot2' 'stats' 'graphics' 'grDevices' 'utils' 'datasets' 'methods' 'base'
5. 'loo' 'sf' 'sp' 'ggthemes' 'ggplot2' 'stats' 'graphics' 'grDevices' 'utils' 'datasets' 'methods' 'base'

'C:/Users/biobayes/Desktop/EC_BIOBAYES/TD2_dauphin'

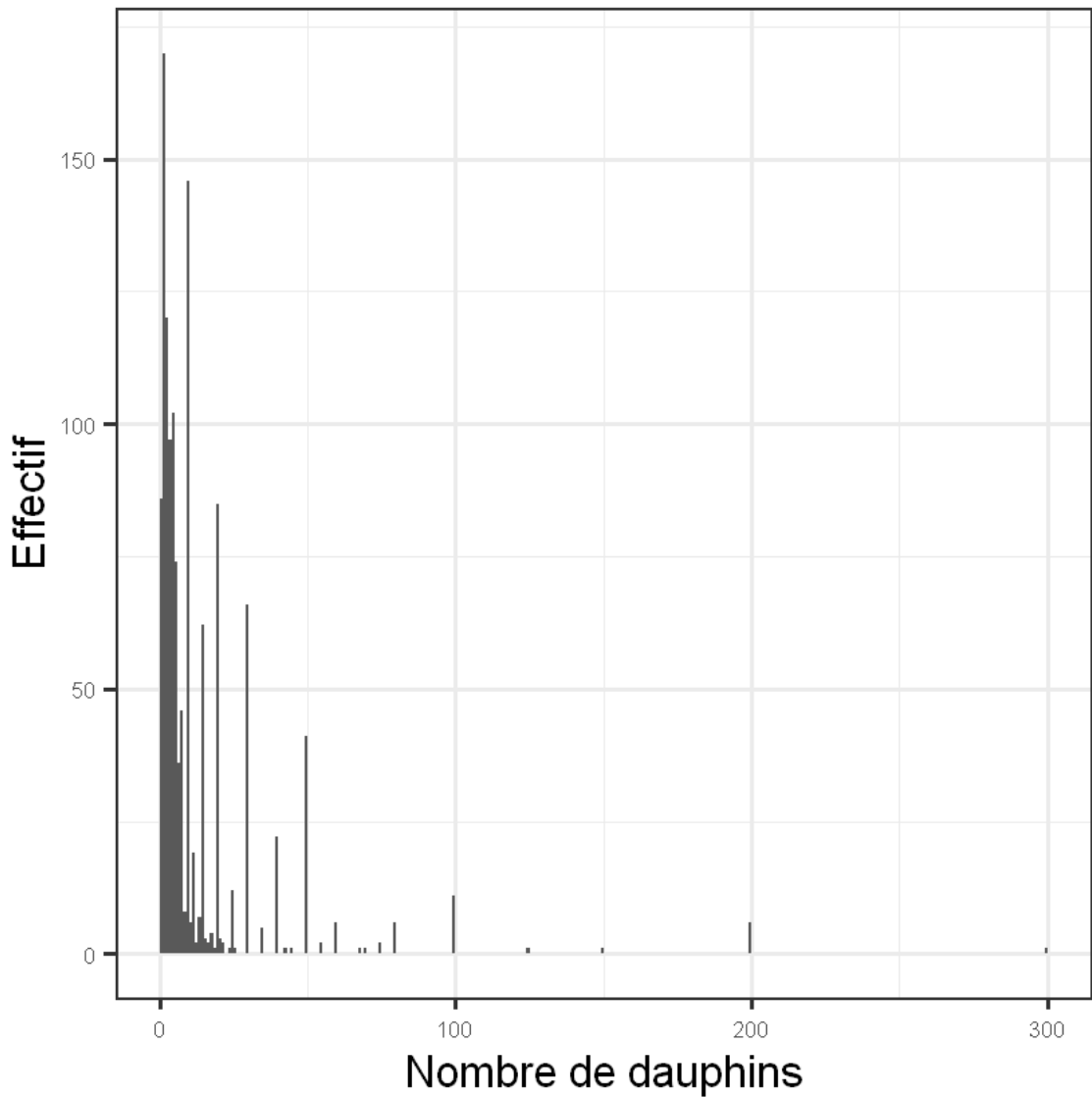
```
In [2]: ### view the different objects in your environment
ls()
```

'A' 'at_sea' 'DataDir' 'grid_at_sea' 'grid_sp' 'neighbors_mat' 'obs' 'OutDir' 'WorkDir'

```
In [3]: ### dataframe of observations
head(obs)
```

| Longitude | Latitude | X | Y | an | Famille | nombre | DepthGebco | Slope | DistCot | ce |
|-----------|----------|----------|---------|------|-------------|--------|------------|------------|---------|-----|
| -6.073679 | 47.91717 | 23721.71 | 6796318 | 2010 | Delphininae | 4 | 0.139192 | 0.02701630 | 92.2346 | 395 |
| -6.069807 | 47.91725 | 24010.08 | 6796294 | 2010 | Delphininae | 5 | 0.139192 | 0.02701630 | 92.2346 | 395 |
| -6.039052 | 47.91930 | 26318.33 | 6796258 | 2010 | Delphininae | 4 | 0.139827 | 0.01384940 | 90.5742 | 395 |
| -6.030971 | 47.91995 | 26926.16 | 6796260 | 2010 | Delphininae | 3 | 0.139827 | 0.01384940 | 90.5742 | 395 |
| -6.018800 | 47.92077 | 27839.73 | 6796247 | 2010 | Delphininae | 5 | 0.135000 | 0.00834959 | 88.8965 | 395 |
| -6.010793 | 47.92166 | 28445.18 | 6796278 | 2010 | Delphininae | 5 | 0.135000 | 0.00834959 | 88.8965 | 395 |

```
In [4]: ### Exploratory data analysis
ggplot(data = obs,
       aes(x = nombre)
       ) +
geom_histogram(breaks = 0:300) +
ylab("Effectif") + xlab("Nombre de dauphins") +
theme(plot.title = element_text(lineheight = 0.8, face = "bold"),
      axis.text = element_text(size = 10)
      )
```



Travaux préliminaires

L'expérience nous pousse à utiliser une loi de Poisson, commune pour représenter les comptages. Les nombres observés étant toujours positifs, on peut décaler d'une unité. Il est d'usage de centrer et de réduire les variables explicatives. L'approche fréquentiste consisterait à *caler* un glm, ce qu'on fait.

Il nous faut écrire le modèle suivant:

$$y_i^{\text{obs}} \sim \mathcal{P}(\lambda_i)$$

avec $\log(\lambda_i) = \beta_0 + \beta_1 \times \text{Bathy}_i + \beta_2 \times \text{Pente}_i + \beta_3 \times \text{DistanceCote}_i$

```
In [5]: obs$nombreShift <- obs$nombre - 1
dauphin2.glm = glm(nombreShift ~ scale(DepthGebco) + scale(Slope) + scale(DistCot),
                    family = "poisson", data = obs
                    )
summary(dauphin2.glm)
dauphin2.glm$coefficients
```

Call:

```
glm(formula = nombreShift ~ scale(DepthGebco) + scale(Slope) +
     scale(DistCot), family = "poisson", data = obs)
```

Deviance Residuals:

```
      Min       1Q   Median       3Q      Max
-7.384  -3.494  -2.029   0.398  32.342
```

Coefficients:

| | Estimate | Std. Error | z value | Pr(> z) |
|-------------------|-----------|------------|---------|------------|
| (Intercept) | 2.472917 | 0.008339 | 296.55 | <2e-16 *** |
| scale(DepthGebco) | 0.233272 | 0.007683 | 30.36 | <2e-16 *** |
| scale(Slope) | 0.128827 | 0.007698 | 16.74 | <2e-16 *** |
| scale(DistCot) | -0.184177 | 0.009476 | -19.44 | <2e-16 *** |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for poisson family taken to be 1)

```
Null deviance: 25774 on 1268 degrees of freedom
Residual deviance: 23626 on 1265 degrees of freedom
AIC: 28052
```

Number of Fisher Scoring iterations: 6

| | |
|--------------------------|--------------------|
| (Intercept) | 2.4729165567972 |
| scale(DepthGebco) | 0.233272395505298 |
| scale(Slope) | 0.128827100897259 |
| scale(DistCot) | -0.184177449274629 |

Codage du modèle en bayésien avec Jags

Écrivons dans un premier temps le modèle avec le langage spécifique à Jags. Ce langage est structuré et la spécification d'un modèle se fait grâce à un ensemble de blocs de type BUGS:

Priors

Dans le paradigme bayésien, les inconnues sont munies de lois a priori qu'il nous faut spécifier. On prendra

$$\beta_0 \sim \mathcal{N}(0.0, 1.0)$$

$$\beta_1 \sim \mathcal{N}(0.0, \frac{\log 2}{2})$$

$$\beta_2 \sim \mathcal{N}(0.0, \frac{\log 2}{2})$$

$$\beta_3 \sim \mathcal{N}(0.0, \frac{\log 2}{2})$$

```
In [6]: library(R2jags)
```

```
model.jags <- '  
model{  
  # Prior  
  intercept ~ dnorm(0.0, hyperparam_inter)  
  slope[1] ~ dnorm(0.0, hyperparam_slope[1])  
  slope[2] ~ dnorm(0.0, hyperparam_slope[2])  
  slope[3] ~ dnorm(0.0, hyperparam_slope[3])  
  
  # Likelihood  
  for (i in 1:n_obs) {  
    lambda[i] <- exp(intercept + slope[1] * X1[i] + slope[2] * X2[i] + slope[3]  
* X3[i])  
    Y[i] ~ dpois(lambda[i])  
  }  
}
```

Loading required package: rjags

Loading required package: coda

Linked to JAGS 4.3.0

Loaded modules: basemod,bugs

Attaching package: 'R2jags'

The following object is masked from 'package:coda':

traceplot

Préparer les données et initialiser les chaînes (en trichant un peu)

```
In [7]: data.jags <- list(n_obs = nrow(obs),  
  Y = obs$nombreShift,  
  X1 = as.vector(scale(obs$DepthGebco)),  
  X2 = as.vector(scale(obs$Slope)),  
  X3 = as.vector(scale(obs$DistCot)),  
  hyperparam_inter = 1.0,  
  hyperparam_slope = 1/rep(log(2)/2, 3)^2  
)  
  
# Inits function  
# On triche un peu: on va prendre les estimations de glm  
inits.jags <- function(idchain, glm) {  
  list(intercept = rnorm(1, glm$coefficients[1], sqrt(10 * diag(vcov(glm)))[1]),  
    slope = rnorm(3, glm$coefficients[2:4], sqrt(10 * diag(vcov(glm)))[2:4])  
  )  
}
```

On nomme les inconnues, on va spécifier le nombre d'itérations puis lancer la machinerie d'inférence. On met les sorties au format de sortie de BUGS avant de les regarder sommairement.


```
In [8]: # parameter of inferential interest
params.jags <- c("intercept", "slope")

# MCMC settings
ni <- 2000 # TOTAL Number of iterations including Burn in
nb <- 1000 #Number of iterations for Burn in
nt <- 1
nc <- 3

# Load/check model, load data and inits
# Start Gibbs sampler
temp <- jags(data = data.jags,
             inits = lapply(1:nc, inits.jags, glm = dauphin2.glm),
             param = params.jags,
             model.file = textConnection(model.jags),
             n.chains = nc,
             n.burnin = nb,
             n.iter = ni,
             n.thin = nt
            )

out = temp$BUGSoutput
# Inferences
# After Burn in period let's work inference assuming ergodic regime is reached
print(out$summary, dig = 2)
```

module glm loaded

Compiling model graph
 Resolving undeclared variables
 Allocating nodes
 Graph information:
 Observed stochastic nodes: 1269
 Unobserved stochastic nodes: 4
 Total graph size: 9591

Initializing model

| | mean | sd | 2.5% | 25% | 50% | 75% | 97.5% | Rhat |
|-----------|----------|---------|----------|----------|----------|----------|----------|------|
| deviance | 30147.60 | 1.2e+02 | 3.0e+04 | 3.0e+04 | 30149.23 | 30219.39 | 30359.12 | 1 |
| intercept | 2.79 | 9.8e-03 | 2.8e+00 | 2.8e+00 | 2.79 | 2.80 | 2.80 | 1 |
| slope[1] | 0.30 | 8.6e-03 | 2.8e-01 | 2.9e-01 | 0.30 | 0.30 | 0.31 | 1 |
| slope[2] | 0.10 | 7.7e-03 | 8.9e-02 | 9.8e-02 | 0.10 | 0.11 | 0.12 | 1 |
| slope[3] | -0.23 | 8.4e-03 | -2.4e-01 | -2.3e-01 | -0.23 | -0.22 | -0.21 | 1 |
| | n.eff | | | | | | | |
| deviance | 1100 | | | | | | | |
| intercept | 3000 | | | | | | | |
| slope[1] | 490 | | | | | | | |
| slope[2] | 670 | | | | | | | |
| slope[3] | 320 | | | | | | | |

Avant toute chose

On doit vérifier (rapidement) la convergence du modèle, par exemple en regardant le diagnostic de Gelman-Rubin.

Prendre un moment pour comprendre la construction de cette statistique d'analyse de variance pour les séries dans la partie théorie de l'aide sous R du programme: `?coda::gelman.diag`

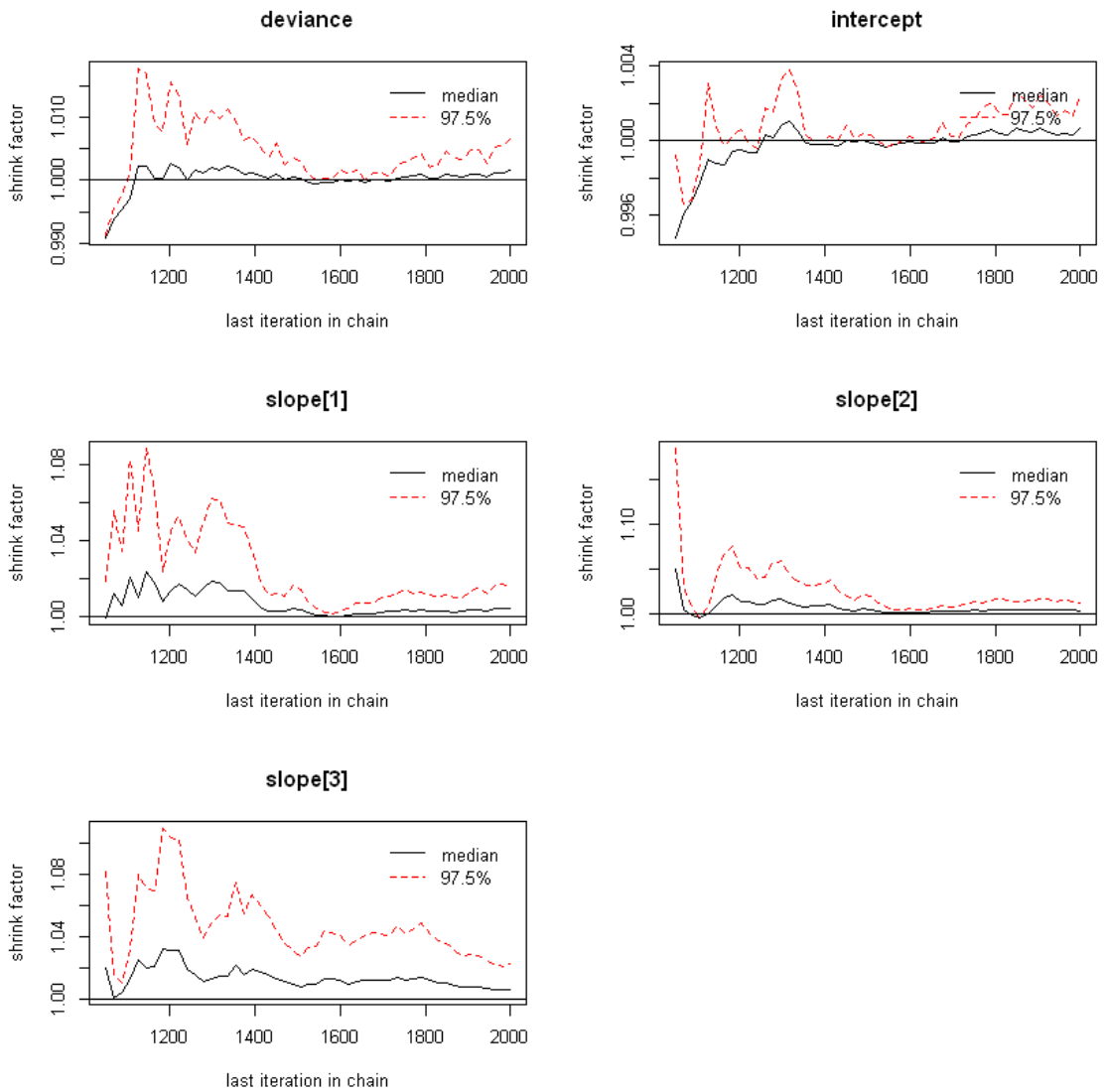
```
In [9]: library(coda)
        gelman.diag(x=out, confidence = 0.95, transform=FALSE, autoburnin=TRUE)
        gelman.plot(x=out)
```

Potential scale reduction factors:

| | Point est. | Upper C.I. |
|-----------|------------|------------|
| deviance | 1.00 | 1.01 |
| intercept | 1.00 | 1.00 |
| slope[1] | 1.00 | 1.02 |
| slope[2] | 1.00 | 1.01 |
| slope[3] | 1.01 | 1.02 |

Multivariate psrf

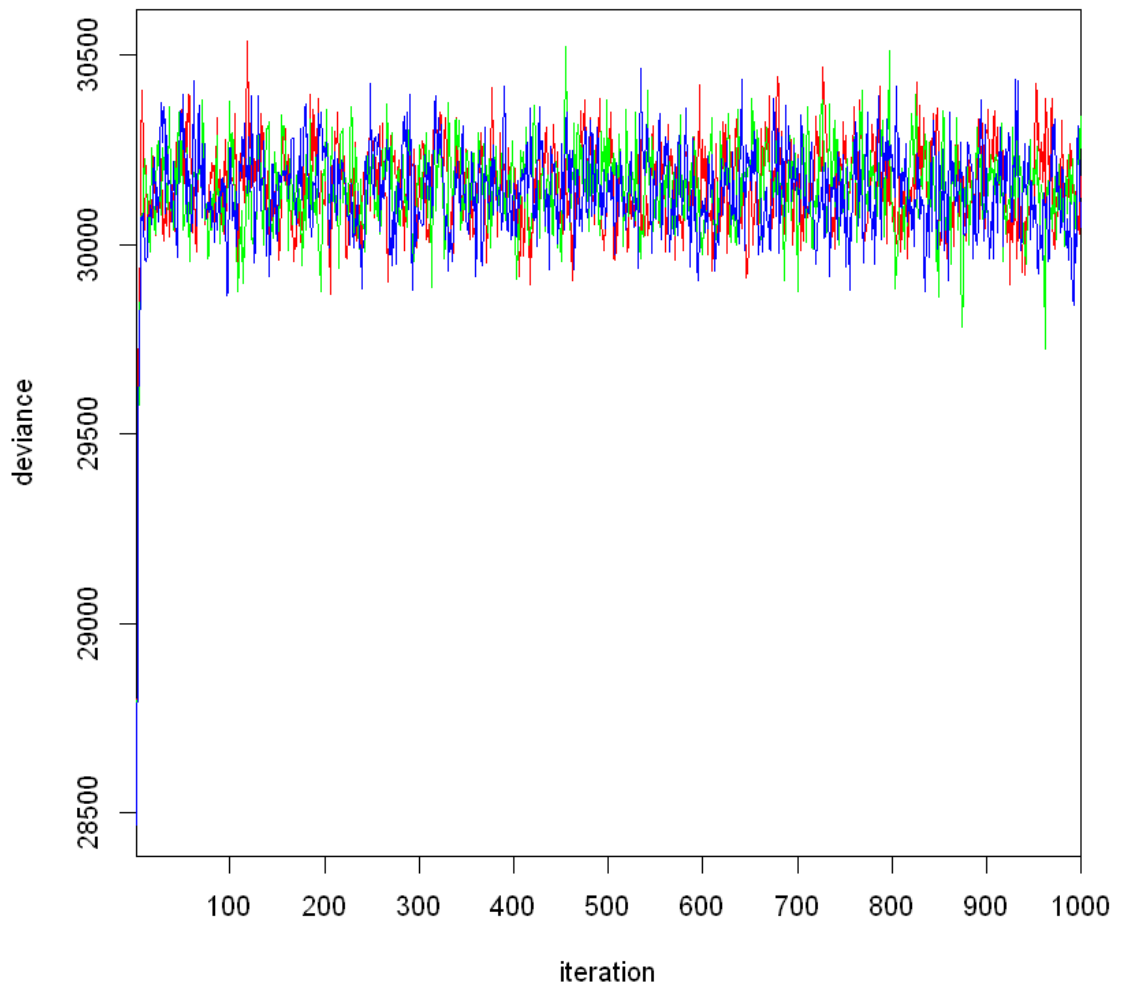
1.01



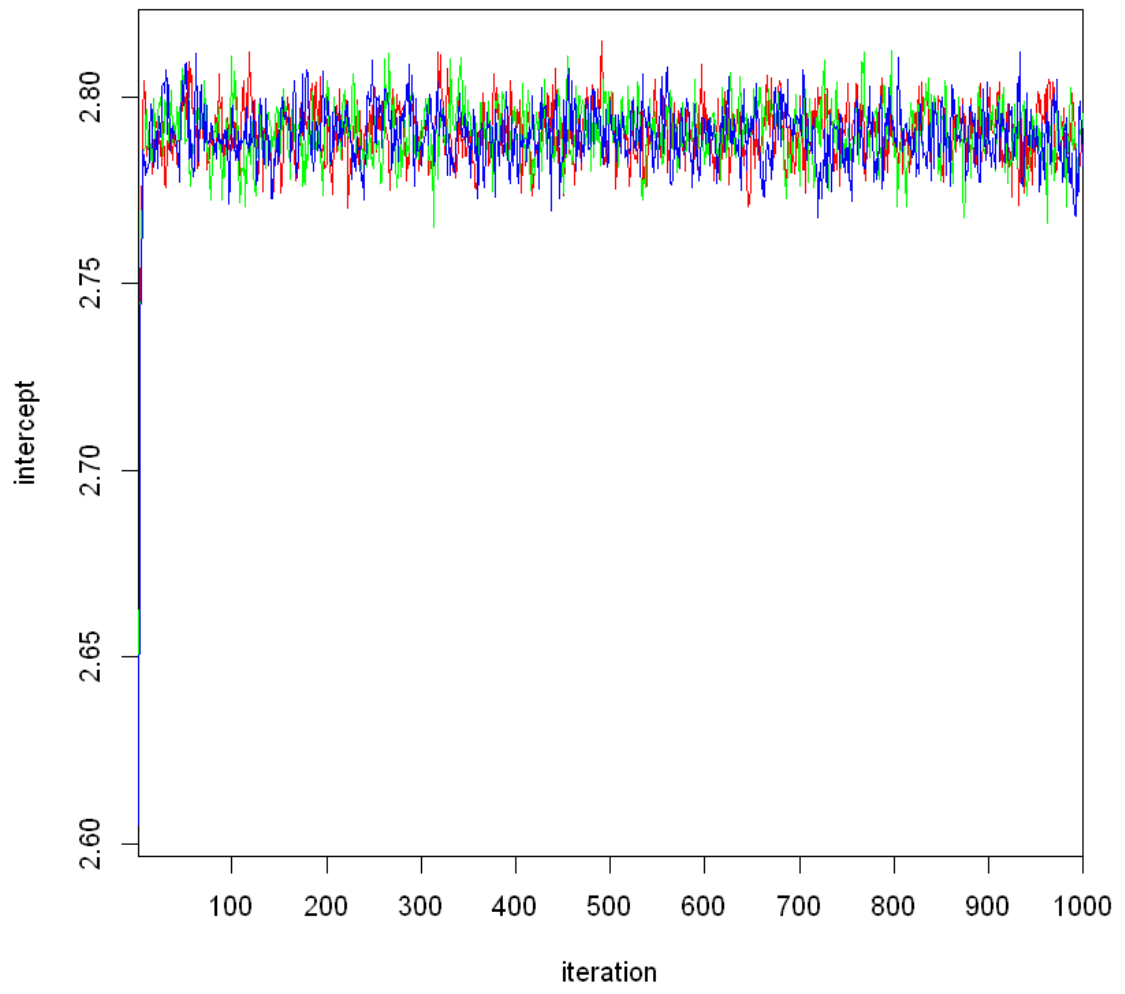
Visualisons la trace des paramètres : on peut voir que les paramètres convergent rapidement (au bout de seulement quelques itérations).

```
In [10]: traceplot(out)
```

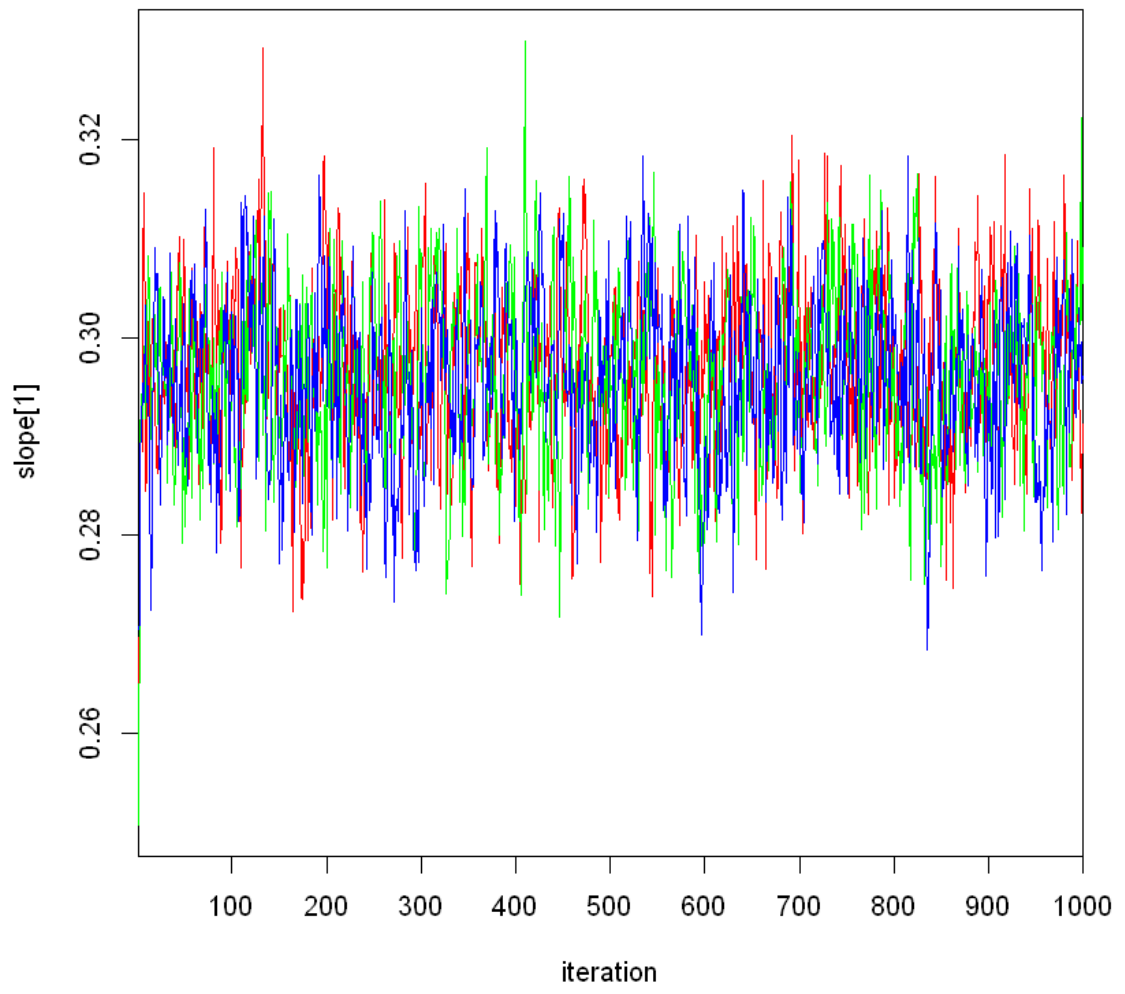
deviance



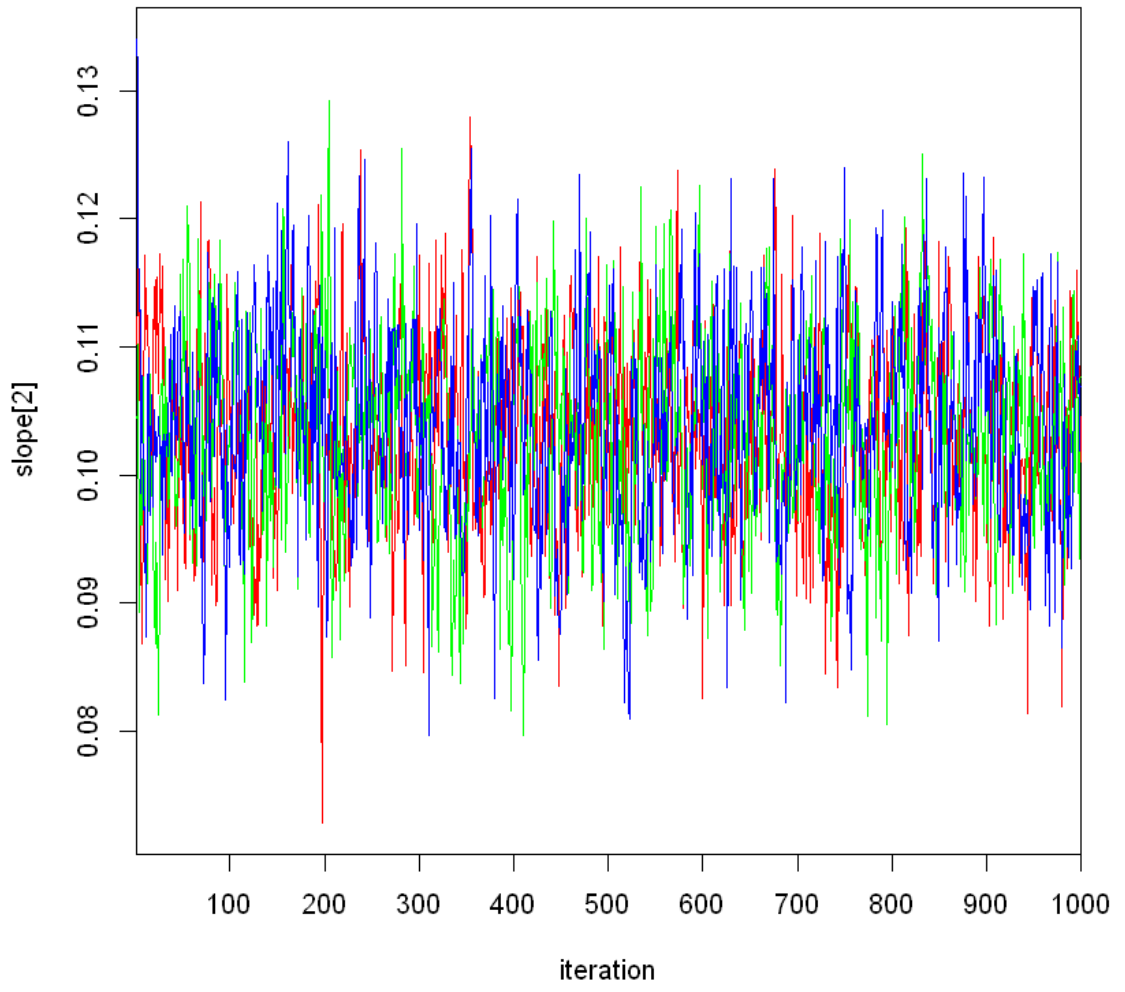
intercept

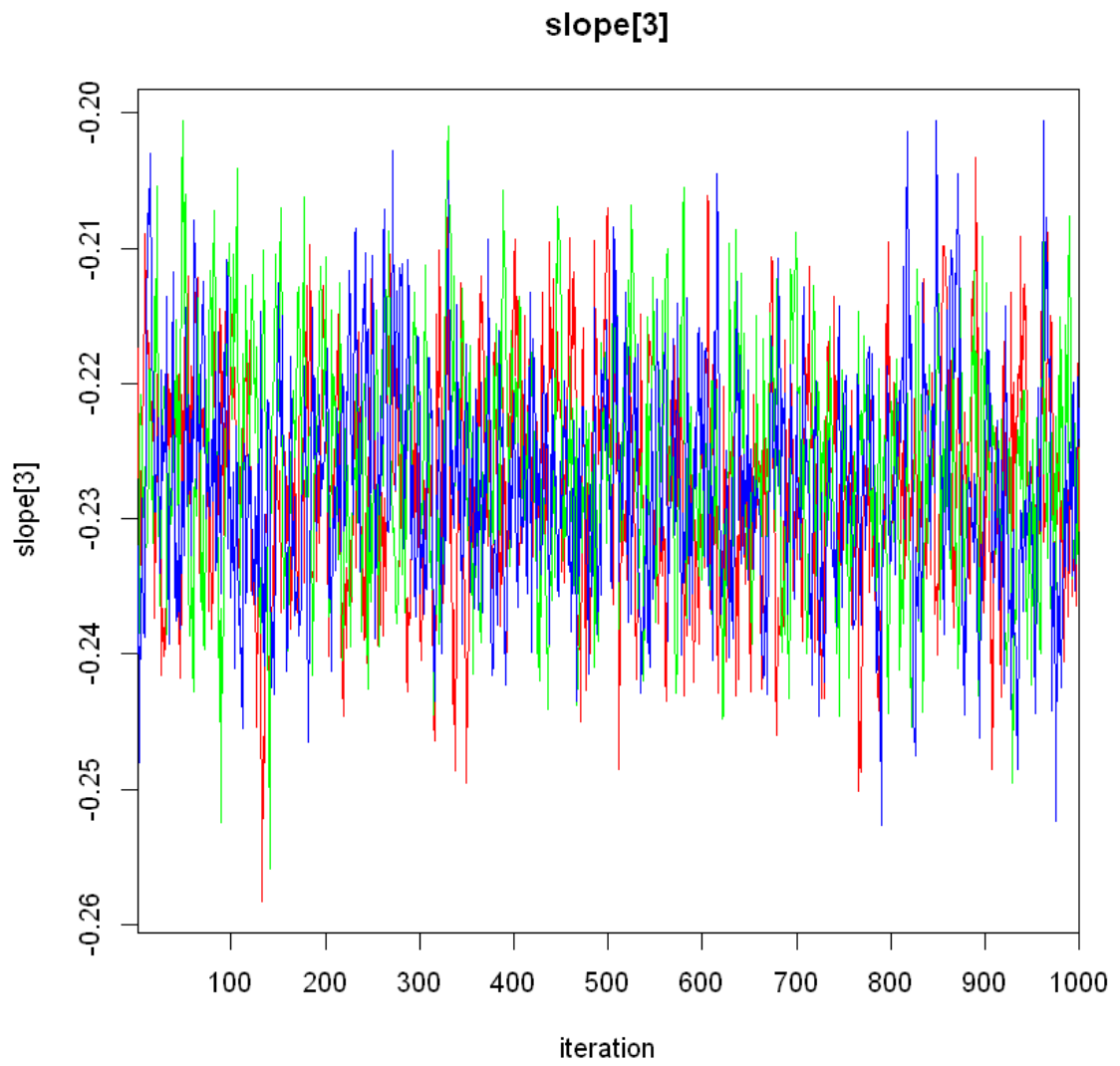


slope[1]



slope[2]





Visualisation des lois a posteriori

On peut faire appel aux potentialités graphiques de R pour représenter les lois a posteriori. Voici un exemple avec ggplot


```

In [11]: library(tidyverse)

other =
  tibble(
    param = paste0("slope[",1:3,"]"),
    class = "slope",
    value = dauphin2.glm$coefficients[2:4],
    type = "mle"
  ) %>%
  rbind(list("intercept", "intercept", dauphin2.glm$coefficients[1], "mle"))

out$sims.matrix[,-1] %>% as_tibble() %>%
  tidyr::gather(param, value) %>%
  mutate(class = stringr::str_replace(param, "\\[[\\d+\\]", "")) %>%
  ggplot(aes_string(x="value", fill="class")) +
  geom_density(alpha=0.5) +
  facet_grid(param ~ .) +
  geom_vline(data = other, aes_string(xintercept="value", color="type", linetype
="type"), size=0.8) +
  scale_colour_manual(values=c("red","purple")) +
  guides(fill=FALSE)+labs(x="valeur de l'inconnue", y="densité")

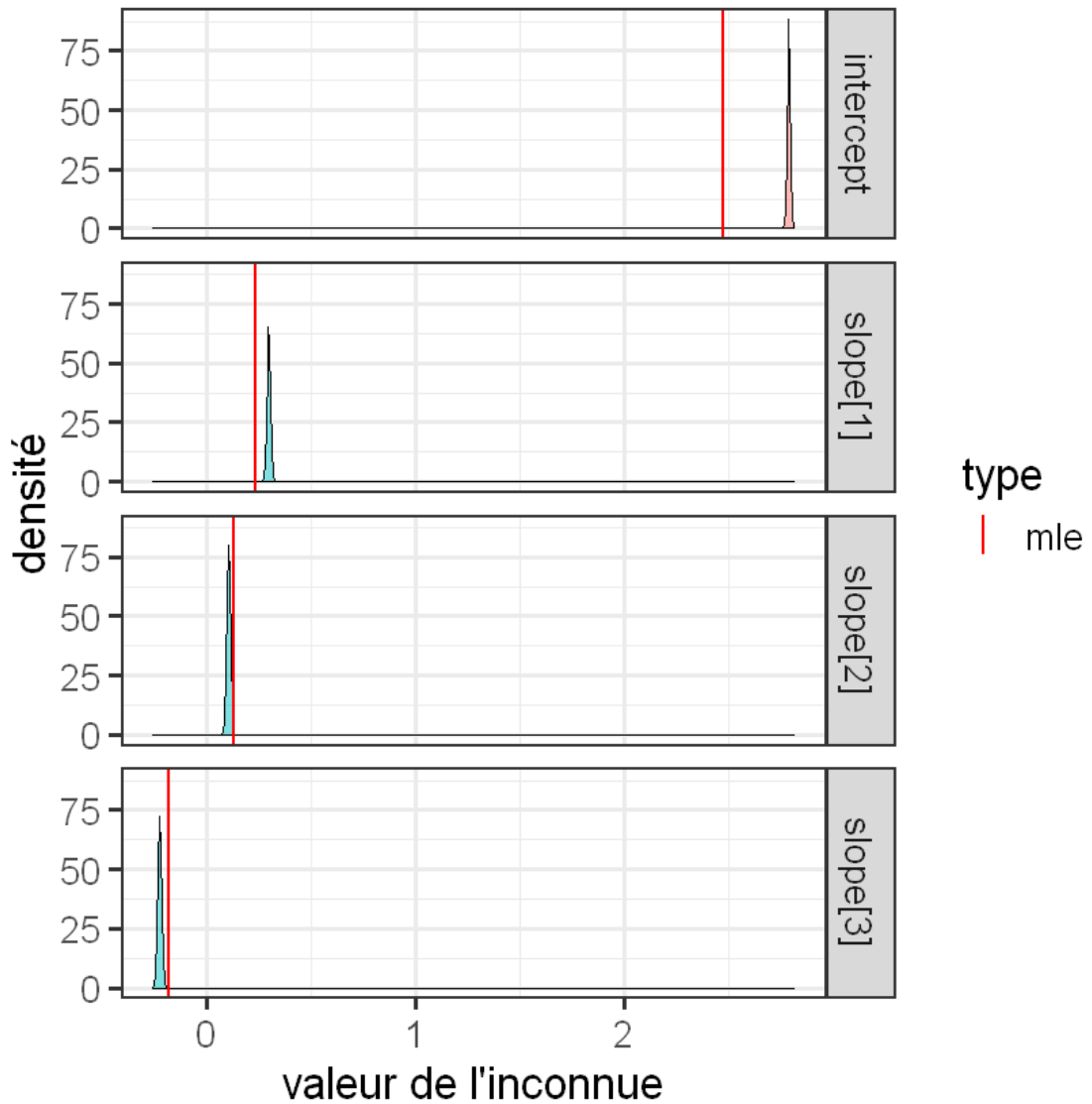
#essayer facet_grid(~param, scales = "free")

```

```

Warning message:
"package 'tidyverse' was built under R version 3.5.3"-- Attaching packages ---
----- tidyverse 1.2.1 --
v tibble  2.0.1      v purrr   0.3.0
v tidyr   0.8.2      v dplyr   0.8.0.1
v readr   1.3.1      v stringr 1.4.0
v tibble  2.0.1      v forcats 0.4.0
Warning message:
"package 'readr' was built under R version 3.5.3"Warning message:
"package 'forcats' was built under R version 3.5.3"-- Conflicts -----
----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()

```



Exercice

- dans la matrice des sorties jags , les 3 chaines sont empilées à la suite. Créer une nouvelle variable catégorielle donnant le numéro de la chaine
- refaire sous *ggplot2* l'équivalent de *traceplot*

Critère d'information

Pour calculer le critère d'information de Watanabee, on a besoin de calculer la log vraisemblance sous le modèle. On fait appel au package *loo*.

```
In [12]: ### calcul du WAIC
## besoin de calculer la log vraisemblance sous le modèle
library(loo)
y_rep_1 <- log_lik1 <- array(NA, dim = c(nc * (ni - nb)/nt, nrow(obs)))

for(i in 1:nrow(obs)) {
  linpred <- out$sims.list$intercept +
    out$sims.list$slope[, 1] * data.jags$X1[i] +
    out$sims.list$slope[, 2] * data.jags$X2[i] +
    out$sims.list$slope[, 3] * data.jags$X3[i]
  log_lik1[, i] <- dpois(data.jags$Y[i], lambda = exp(linpred), log = TRUE)
  y_rep_1[, i] <- rpois(nrow(y_rep_1), lambda = exp(linpred)) + 1
}; rm(i, linpred)

loo::waic(log_lik1)
loo::loo(log_lik1)
```

```
Warning message:
"86 (6.8%) p_waic estimates greater than 0.4. We recommend trying loo instead.
Warning message:
"86 (6.8%) p_waic estimates greater than 0.4. We recommend trying loo instead.
"
```

Computed from 3000 by 1269 log-likelihood matrix

| | Estimate | SE |
|-----------|----------|--------|
| elpd_waic | -15097.4 | 922.7 |
| p_waic | 183.5 | 33.1 |
| waic | 30194.9 | 1845.3 |

```
Warning message:
"Relative effective sample sizes ('r_eff' argument) not specified.
For models fit with MCMC, the reported PSIS effective sample sizes and
MCSE estimates will be over-optimistic."Warning message:
"Some Pareto k diagnostic values are too high. See help('pareto-k-diagnostic')
for details.
```

```
"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"Warning message in log(z):
"production de NaN"
```

Computed from 3000 by 1269 log-likelihood matrix

| | Estimate | SE |
|----------|----------|--------|
| elpd_loo | -15214.3 | 934.3 |
| p_loo | 300.4 | 45.4 |
| looic | 30428.7 | 1868.5 |

Monte Carlo SE of elpd_loo is NA.

```
Pareto k diagnostic values:
               Count Pct.    Min. n_eff
(-Inf, 0.5]  (good)  1028  81.0%   623
 (0.5, 0.7]  (ok)    127   10.0%  1389
  (0.7, 1]   (bad)    83    6.5%   57
   (1, Inf)  (very bad)  31    2.4%    1
See help('pareto-k-diagnostic') for details.
```

Posterior Predictive Checks

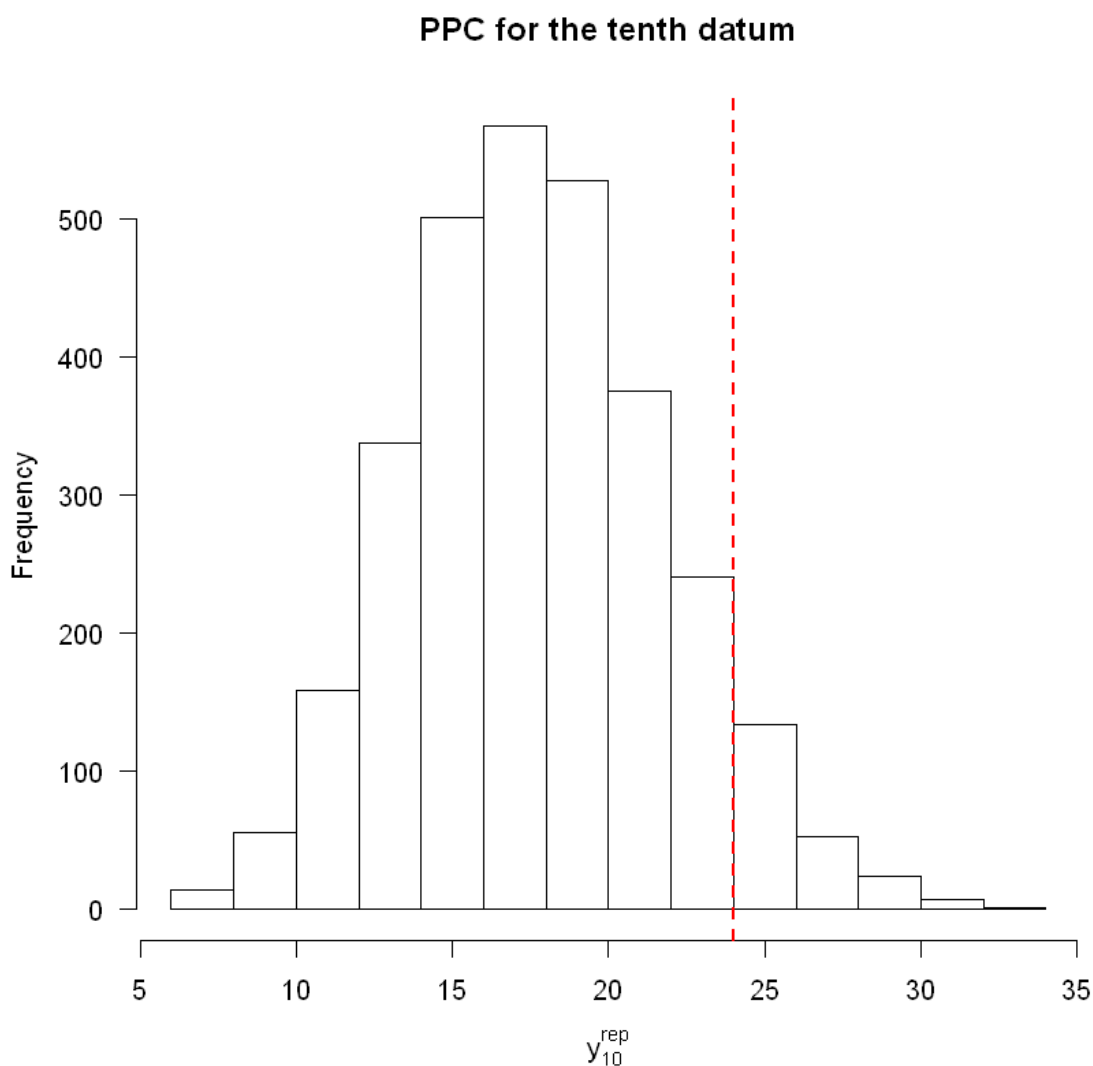
On va extraire directement les valeurs prédites y^{rep_i}

```
In [13]: y_rep_poisson_1 <- y_rep_1  
dim(y_rep_poisson_1)
```

```
3000 1269
```

Cet objet est une matrice de 3000 lignes et 1269 colonnes, soit une colonne pour chaque donnée et une ligne pour chaque itération.

```
In [14]: hist(y_rep_poisson_1[, 10], las = 1, xlab = quote(y[10]^rep), main = "PPC for the tenth datum")  
abline(v = obs$nombreShift[10], lty = 2, col = "red", lwd = 2)
```



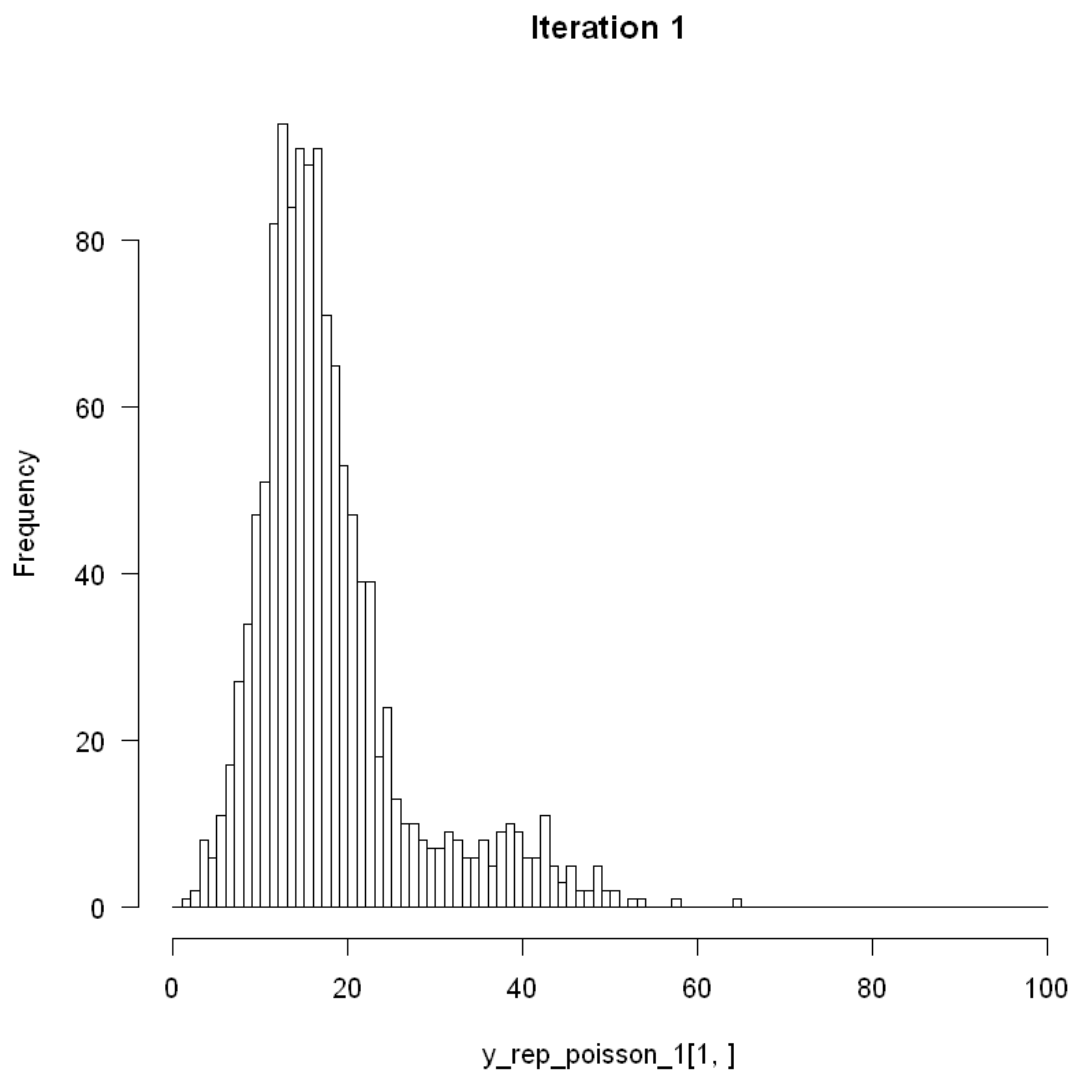
On peut par exemple calculer la probabilité $Pr(y_{10}^{\text{rep}} > y_{10}^{\text{obs}})$

```
In [15]: round(mean(ifelse(y_rep_poisson_1[, 10] > (obs$nombreShift[10]), 1, 0)), dig=3)
```

```
0.073
```

On peut aussi comparer les histogrammes de l'ensemble des données (observées et prédites) au lieu de regarder donnée par donnée.

```
In [16]: hist(y_rep_poisson_1[1, ], las = 1, main = "Iteration 1", breaks = c(0:100))
```



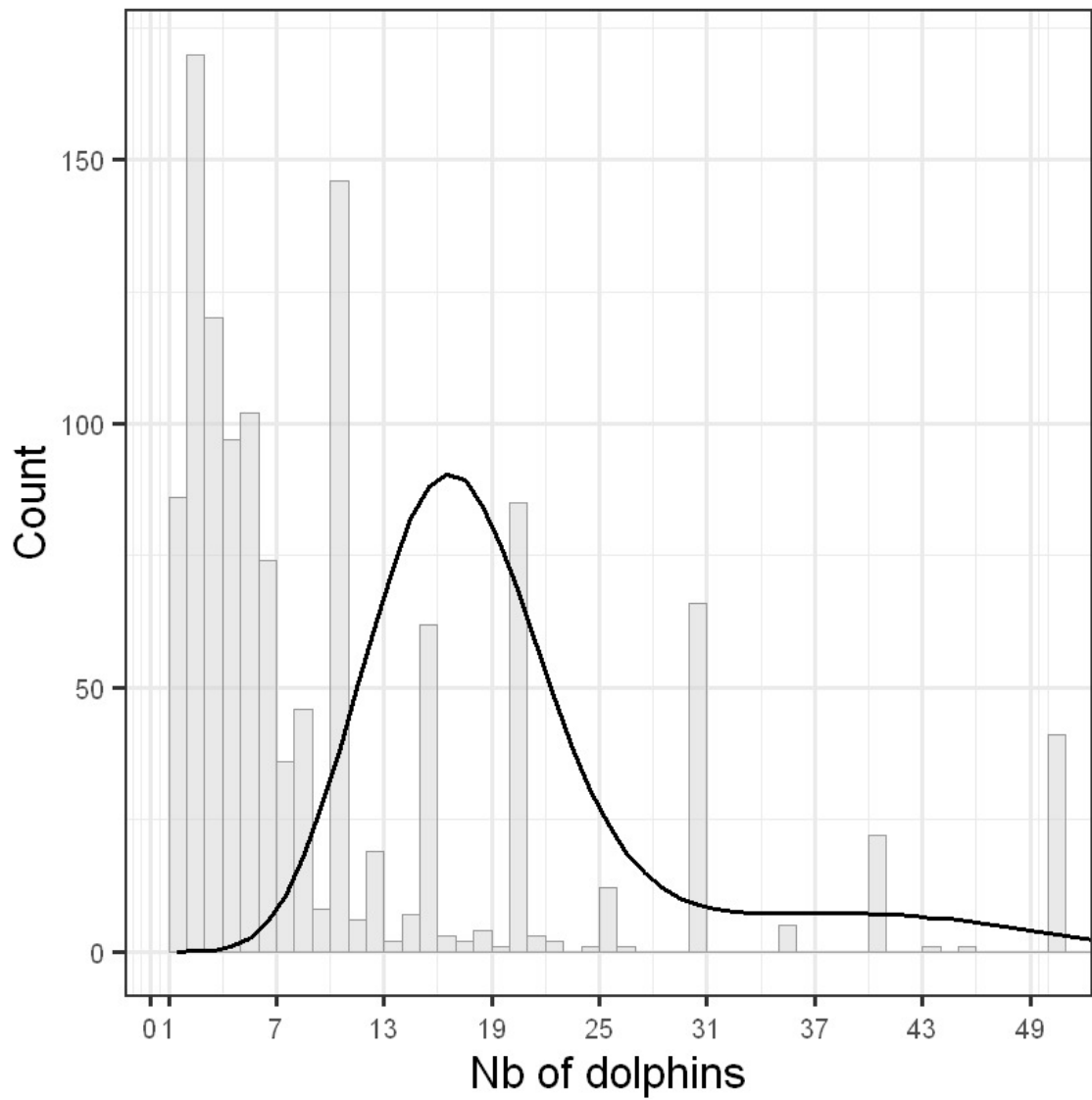
On va moyenner l'histogramme sur l'ensemble des itérations en faisant appel à une fonction

```

In [17]: # rootograms
rootogram <- function(countdata, y_rep, min_obs = 1) {
  f_histogram <- function(x, max_obs) { table(factor(x, levels = min_obs:max_obs
)) }
  max_obs <- max(countdata, as.numeric(y_rep))
  dd <- ggplot(data.frame(mids = (min_obs:max_obs + (min_obs + 1):(max_obs + 1))
/2,
                        y_obs = as.numeric(f_histogram(x = countdata, max_obs
= max_obs)),
                        y_rep = apply(apply(y_rep, 1, f_histogram, max_obs = m
ax_obs), 1, mean)
                        ),
            aes(x = mids, y = y_rep)
            ) +
  geom_rect(aes(xmin = mids - 0.5, xmax = mids + 0.5,
               ymax = y_obs, ymin = 0),
            fill = "lightgrey", color = grey(0.6), alpha = 0.5) +
  geom_line(aes(x = mids, y = y_rep), color = "black", size = 1.0) +
  scale_y_continuous(name = "Count") +
  scale_x_continuous(name = "Nb of dolphins", breaks = c(0, seq(1, max_obs, ma
x_obs/50))) +
  guides(size = "none") +
  theme(plot.title = element_text(lineheight = 0.8, face = "bold"),
        axis.text = element_text(size = 12),
        strip.text = element_text(size = 12),
        strip.background = element_rect(fill = grey(0.95))
        )
  return(dd)
}

rootogram_model_1 <- rootogram(countdata = obs$nombre, y_rep = y_rep_1+1)
rootogram_model_1 +
  coord_cartesian(xlim = c(1, 50))

```



Un ajustement très (trop?) optimiste quant aux incertitudes

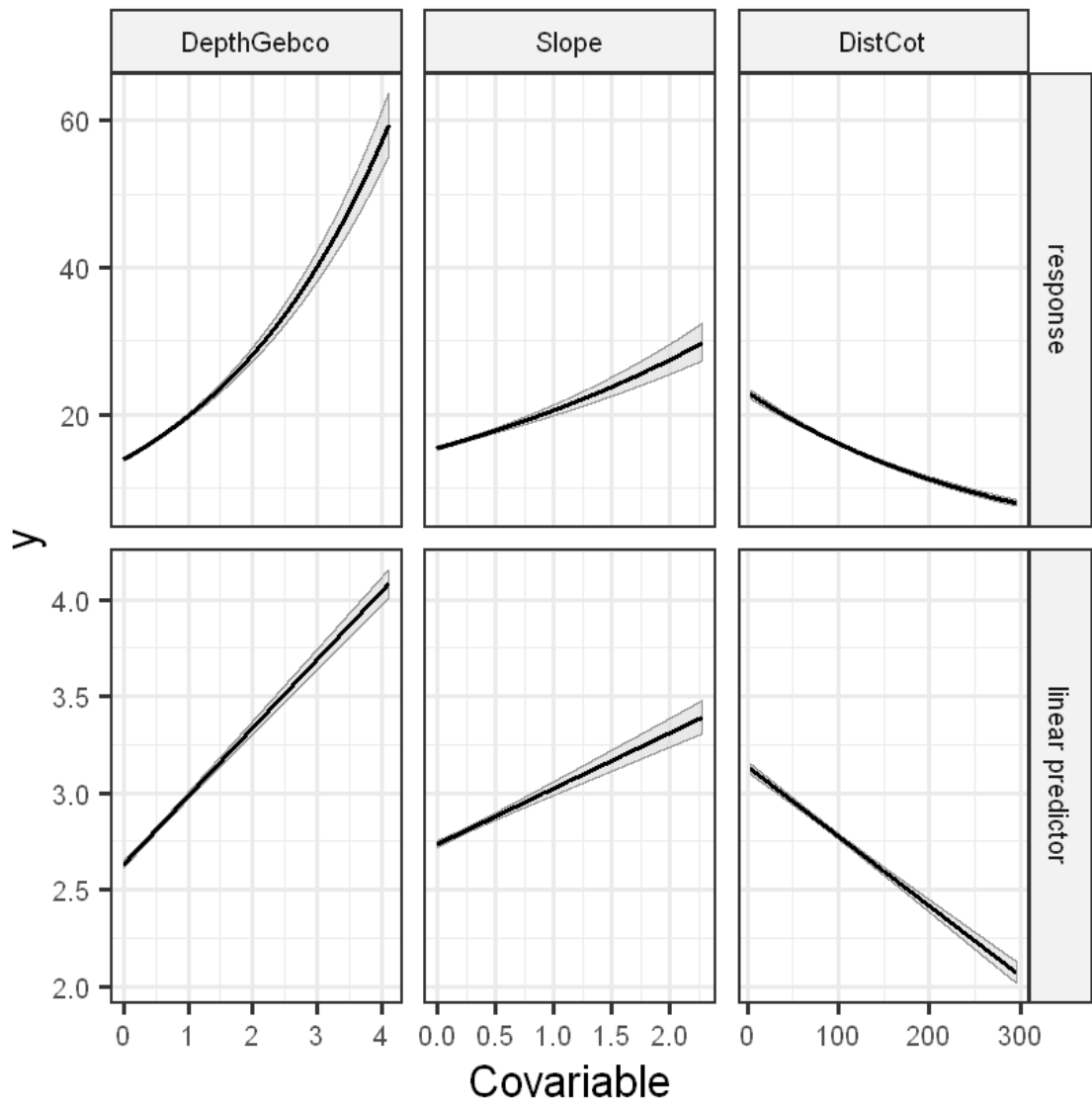
Nous allons reconstruire les effets des prédicteurs avec les gammes d'incertitude issues de la connaissance a posteriori des paramètres. *Trop beau pour être honnête?*

```

In [18]: plot_relationships <- function(data_df, jagsfit, cov_name) {
  X <- data_df[, cov_name]
  x_scale <- apply(X, 2, function(x) {c(mean(x), sd(x))})
  x_range <- apply(apply(X, 2, scale), 2, range)
  stdX <- cbind(rep(1, 1e4),
                apply(x_range, 2, function(vec) {seq(vec[1], vec[2], length.out
= 1e4)}))
  )
  beta <- cbind(jagsfit$sims.list$intercept,
                jagsfit$sims.list$slope
                )
  if(ncol(stdX) != ncol(beta)) { stop("Please check model and covariable names:\n\tdimension mismatch between slope parameters and covariables") }
  else {
    dd <- data.frame(x = numeric(0),
                     y = numeric(0),
                     lower = numeric(0),
                     upper = numeric(0),
                     what = character(0),
                     covariable = character(0)
                     )
    for(j in 1:length(cov_name)) {
      Xmat <- cbind(rep(1, 1e4),
                    matrix(0, nrow = 1e4, ncol = length(cov_name))
                    )
      Xmat[, j+1] <- stdX[, j+1]
      linpred <- beta %*% t(Xmat)
      dd <- rbind(dd,
                  data.frame(x = Xmat[, j+1] * x_scale[2, cov_name[j]] + x_scale
[1, cov_name[j]],
                             y = apply(linpred, 2, mean),
                             lower = apply(linpred, 2, stats::quantile, prob = 0
.025),
                             upper = apply(linpred, 2, stats::quantile, prob = 0
.975),
                             what = rep("linear predictor", 1e4),
                             covariable = rep(cov_name[j], 1e4)
                             ),
                  data.frame(x = Xmat[, j+1] * x_scale[2, cov_name[j]] + x_scale
[1, cov_name[j]],
                             y = apply(exp(linpred), 2, mean),
                             lower = apply(exp(linpred), 2, stats::quantile, pro
b = 0.025),
                             upper = apply(exp(linpred), 2, stats::quantile, pro
b = 0.975),
                             what = rep("response", 1e4),
                             covariable = rep(cov_name[j], 1e4)
                             )
                  )
    }
    dd$what <- factor(as.character(dd$what), levels = c("response", "linear pred
ictor"))
    dd$covariable <- factor(as.character(dd$covariable), levels = cov_name)

    the_plot <- ggplot(data = dd,
                       aes(x = x, y = y)
                       ) +
      geom_ribbon(aes(x = x, ymin = lower, ymax = upper),
                  fill = "lightgrey", color = grey(0.6), alpha = 0.5
                  ) +
      geom_line(aes(x = x, y = y), color = "black", size = 1.0) +
      scale_y_continuous(name = "y") +
      scale_x_continuous(name = "Covariable") +
      facet_grid(what~covariable, scales = "free") +
      theme(plot.title = element_text(lineheight = 0.8, face = "bold"),
            axis.text = element_text(size = 12),
            strip.text = element_text(size = 12),
            strip.background = element_rect(fill = grey(0.95))
            )
  }
}

```

Travaux dirigés R

- Retour à la définition: Construire une fonction qui calcule le WAIC à la main à partir d'une matrice comme `log_lik1`. Soit θ le vecteur des paramètres inconnus d'un modèle d'observations $f(\cdot | \theta)$. Soit (y_1, \dots, y_n) un échantillon de réalisations de $f(\cdot | \theta)$. Soit $(\theta^1, \dots, \theta^S)$ un échantillon de valeurs issues de la loi a posteriori jointe de θ et générées à l'aide d'un algorithme MCMC. Le critère d'information WAIC (Watanabe-Akaike ou Widely Applicable Information Criterion) du modèle peut alors être approximé (voir Bayesian Data Analysis, Gelman et al., 2016) en remplaçant les espérances *a posteriori* par des sommes sur les tirages dans le posterior:

$$W\hat{AIC} = -2\hat{llpd} + 2p_{W\hat{AIC}}$$

avec

$$\hat{llpd} = \sum_{i=1}^n \log \left(\frac{1}{S} \sum_{s=1}^S f(y_i | \theta^s) \right)$$

et

$$p_{W\hat{AIC}} = \sum_{i=1}^n V_{s=1}^S (\log f(y_i | \theta^s))$$

$$V_{s=1}^S(a_s) = \frac{1}{S-1} \sum_{s=1}^S (a_s - \bar{a})^2, \bar{a} = \frac{1}{S} \sum_{s=1}^S a_s.$$

- On pourra également y ajouter le DIC comme sortie (qu'on comparera aux résultats de la fonction `dic.samples` de `rjags`).
- Faire un modèle à 2 covariables sous Jags (par exemple, supprimer la distance à la côte).
- Est-ce que ce nouveau modèle est meilleur que le précédent? On pourra utiliser le critère WAIC pour faire une comparaison de modèle.
- Les critères DIC et WAIC mènent-ils aux mêmes conclusions? Quel(s) avantage(s) voyez-vous à l'utilisation du critère WAIC par rapport au critère DIC?

Pour aller plus loin: Les modèles fonctionnent mais l'ajustement est atroce, que faire?

Ici, nous avons un problème d'ajustement sérieux des modèles envisagés aux données. Le problème provient en fait d'une hypothèse implicite lorsque l'on utilise un modèle poissonien: l'équidispersion. Cela signifie que moyenne et variance sont égales. Cependant, le cas général est plutôt d'avoir des données surdispersées, c'est-à-dire des données dont la variance est supérieure à la moyenne. Cette surdispersion peut être le reflet de différents processus, par exemple une structure spatiale non prise en compte par les covariables, etc.

Ici, nous allons modifier le modèle pour ne plus supposer une vraisemblance poissonienne pour les données, mais une vraisemblance Negative Binomiale:

$$y_i^{\text{obs}} \sim \mathcal{NB}(\omega, \lambda_i)$$

où $\omega > 1$ est le paramètre de surdispersion. Il permet de prendre en compte une variance supérieure à la moyenne dans des données de comptages:

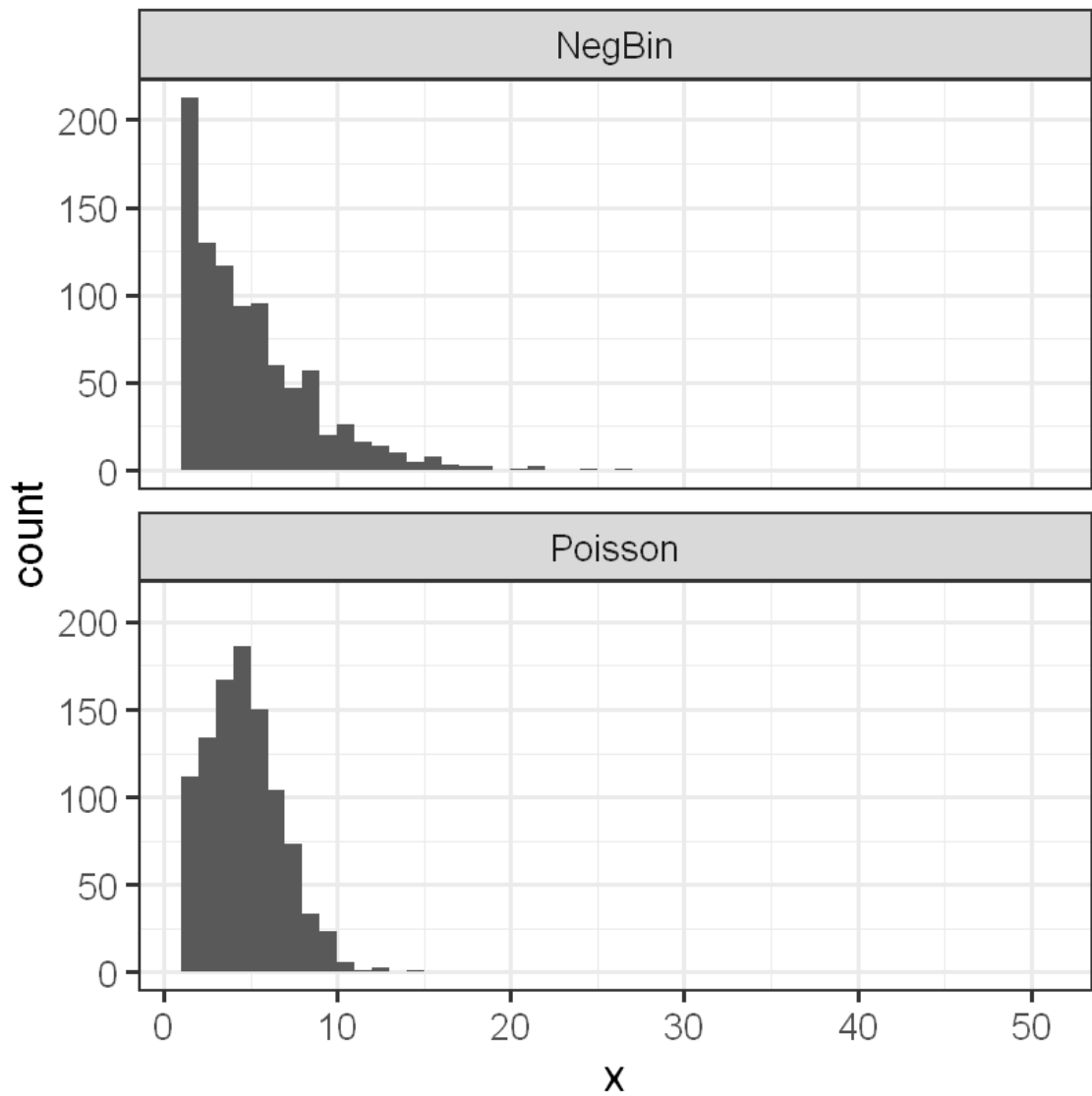
$$V(y) = \omega E(y)$$

Nous allons simuler quelques données avec R pour regarder la différence entre une distribution Poisson et une distribution Negative Binomiale.

```

In [19]: n <- 1e3 # nombre de données à simuler
lambda <- 5 # moyenne du processus
omega <- 3 # paramètre de surdispersion
ggplot(data = data.frame(x = c(rpois(n, lambda),
                             MASS::rnegin(n, mu = lambda, theta = lambda/(ome
ga - 1))
                             ),
                             distribution = rep(c("Poisson", "NegBin"), each = n)
                             ),
        aes(x = x, group = distribution)
        ) +
geom_histogram(breaks = seq(0:50)) +
facet_wrap(~distribution, ncol = 1)

```



Negative Binomiale et Poisson

La fonction de densité de probabilité d'une variable aléatoire y qui suit une loi de Poisson d'intensité $\lambda (> 0)$ est donnée par :

$$p(y|\lambda) = \frac{e^{-\lambda} \lambda^y}{y!}$$

La moyenne de y est égale à la variance : $\mathbb{V}(y) = \mathbb{E}(y) = \lambda$

La fonction de densité de probabilité d'une variable aléatoire y qui suit une loi Negative Binomiale d'intensité $\lambda (> 0)$ et de paramètre de surdispersion $\omega (> 0)$ est donnée par :

$$p(y|\lambda, \omega) = \binom{y + \frac{\lambda}{\omega} - 1}{y} \left(\frac{\lambda}{\omega}\right)^{\frac{\lambda}{\omega}} \left(1 - \frac{\lambda}{\omega}\right)^y$$

La distribution Negative binomiale peut aussi être vue comme une généralisation de la Poisson avec une intensité qui est elle-même une variable aléatoire suivant une loi Gamma de paramètre de forme $\frac{\lambda}{\omega-1}$ et de paramètre de taux $\frac{1}{\omega-1}$. Cette caractérisation permet d'écrire facilement un modèle sous Jags ou Bugs.

Voir aussi https://www.johndcook.com/negative_binomial.pdf (https://www.johndcook.com/negative_binomial.pdf) pour les différentes paramétrisations possible.

Pour modifier la distribution dans le code Jags il va falloir faire quelques changements. En particulier, il faut spécifier un prior pour ω .

La paramètre de surdispersion est compris entre 1 et $+\infty$. On va ici mettre un prior uniforme entre 0 et 1 sur l'inverse du paramètre de surdispersion :

$$p\left(\frac{1}{\omega}\right) \propto 1$$

Il faut également modifier le calcul de la vraisemblance. La paramétrisation la plus prisée par les écologistes utilise un paramètre $\phi = \lambda\omega - 1$.

Nota: Ce modèle, plus complexe, nécessite plus de temps pour tourner

```

In [20]: model.jags <- '
model{
  # Prior
  intercept ~ dnorm(0.0, hyperparam_inter)
  slope[1] ~ dnorm(0.0, hyperparam_slope[1])
  slope[2] ~ dnorm(0.0, hyperparam_slope[2])
  slope[3] ~ dnorm(0.0, hyperparam_slope[3])
  inv_overdispersion ~ dunif(0, 1)
  overdispersion <- 1/inv_overdispersion

  # Likelihood
  for (i in 1:n_obs) {
    lambda[i] <- exp(intercept + slope[1] * X1[i] + slope[2] * X2[i] + slope[3]
* X3[i])
    r[i] <- lambda[i] / (overdispersion - 1)
    Y[i] ~ dnegbin(inv_overdispersion, r[i])
  }
}
'

data.jags <- list(n_obs = nrow(obs),
                 Y = obs$nombreShift,
                 X1 = as.vector(scale(obs$DepthGebco)),
                 X2 = as.vector(scale(obs$Slope)),
                 X3 = as.vector(scale(obs$DistCot)),
                 hyperparam_inter = 1.0,
                 hyperparam_slope = 1/rep(log(2)/2, 3)^2
                 )

# Inits function
inits.jags <- function(idchain, glm) {
  list(intercept = rnorm(1, glm$coefficients[1], sqrt(10 * diag(vcov(glm))[1])),
        slope = rnorm(3, glm$coefficients[2:4], sqrt(10 * diag(vcov(glm))[2:4])),
        inv_overdispersion = runif(1, 0, 1)
  )
}

params.jags <- c("intercept", "slope", "overdispersion")

# Start Gibbs sampler
temp <- jags(data = data.jags,
             inits = lapply(1:nc, inits.jags, glm = dauphin2.glm),
             param = params.jags,
             model.file = textConnection(model.jags),
             n.chains = nc,
             n.burnin = nb,
             n.iter = ni,
             n.thin = nt
             )

out3 = temp$BUGSoutput
# Inferences
# After Burn in period let's work inference assuming ergodic regime is reached
print(out3$summary, dig = 2)

```

```

Compiling model graph
  Resolving undeclared variables
  Allocating nodes
Graph information:
  Observed stochastic nodes: 1269
  Unobserved stochastic nodes: 5
  Total graph size: 10501

```

```

Initializing model

```

| | mean | sd | 2.5% | 25% | 50% | 75% | 97.5% | Rhat |
|----------------|----------|-------|----------|---------|----------|----------|----------|------|
| deviance | 8877.892 | 3.023 | 8873.925 | 8875.64 | 8877.262 | 8879.594 | 8885.081 | 1 |
| intercept | 2.498 | 0.034 | 2.431 | 2.47 | 2.498 | 2.521 | 2.563 | 1 |
| overdispersion | 18.069 | 0.872 | 16.396 | 17.48 | 18.023 | 18.634 | 19.819 | 1 |
| slope[1] | 0.143 | 0.028 | 0.088 | 0.12 | 0.144 | 0.162 | 0.196 | 1 |
| slope[2] | 0.124 | 0.028 | 0.070 | 0.11 | 0.125 | 0.143 | 0.177 | 1 |
| slope[3] | -0.098 | 0.028 | -0.154 | -0.12 | -0.098 | -0.079 | -0.044 | 1 |
| | n.eff | | | | | | | |
| deviance | 1400 | | | | | | | |
| intercept | 1600 | | | | | | | |
| overdispersion | 1700 | | | | | | | |
| slope[1] | 500 | | | | | | | |
| slope[2] | 530 | | | | | | | |
| slope[3] | 500 | | | | | | | |

On peut vérifier graphiquement la convergence du modèle.

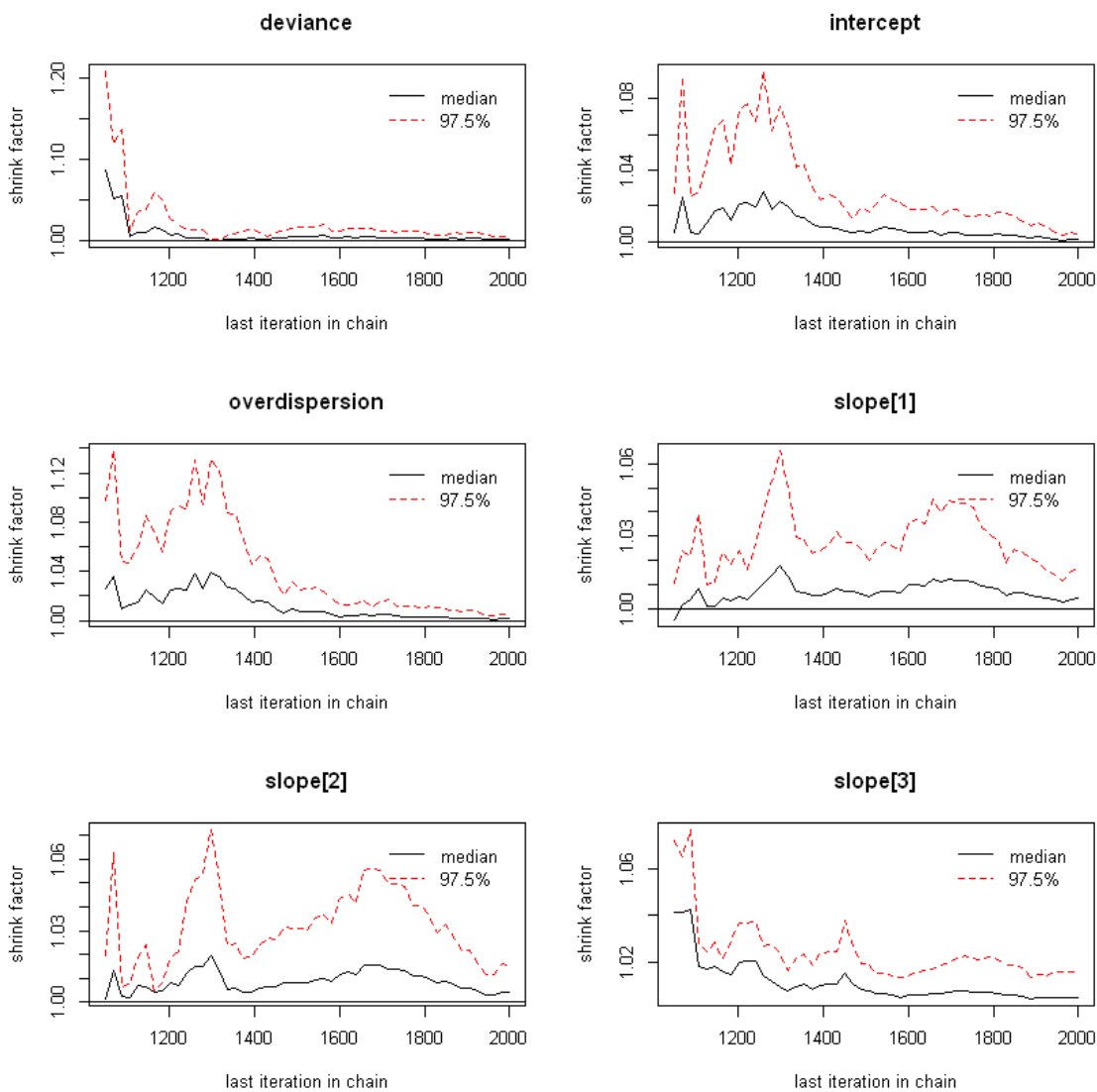
```
In [21]: library(coda)
gelman.diag(x=out3, confidence = 0.95, transform=FALSE, autoburnin=TRUE)
gelman.plot(x=out3)
```

Potential scale reduction factors:

| | Point est. | Upper C.I. |
|----------------|------------|------------|
| deviance | 1 | 1.00 |
| intercept | 1 | 1.00 |
| overdispersion | 1 | 1.00 |
| slope[1] | 1 | 1.02 |
| slope[2] | 1 | 1.02 |
| slope[3] | 1 | 1.02 |

Multivariate psrf

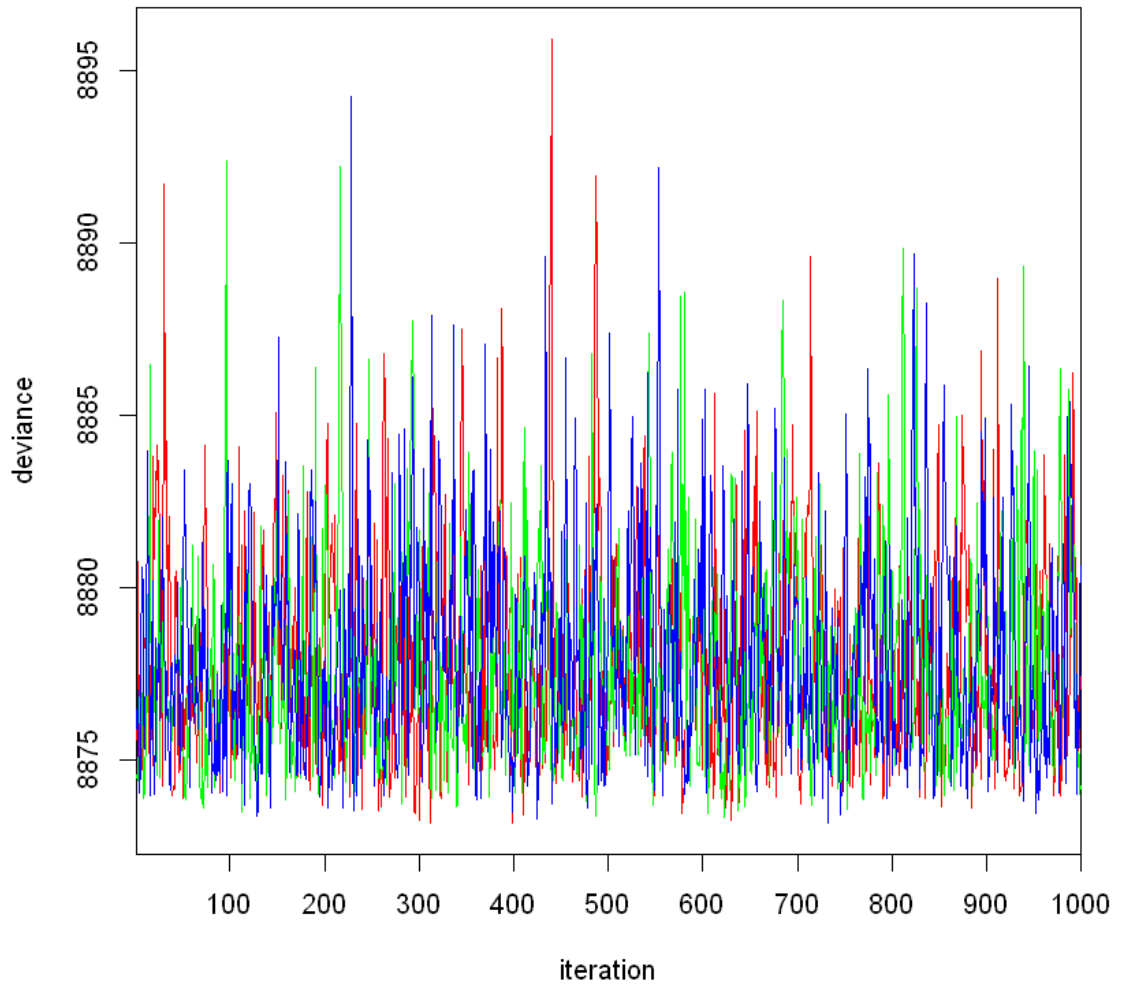
1.01



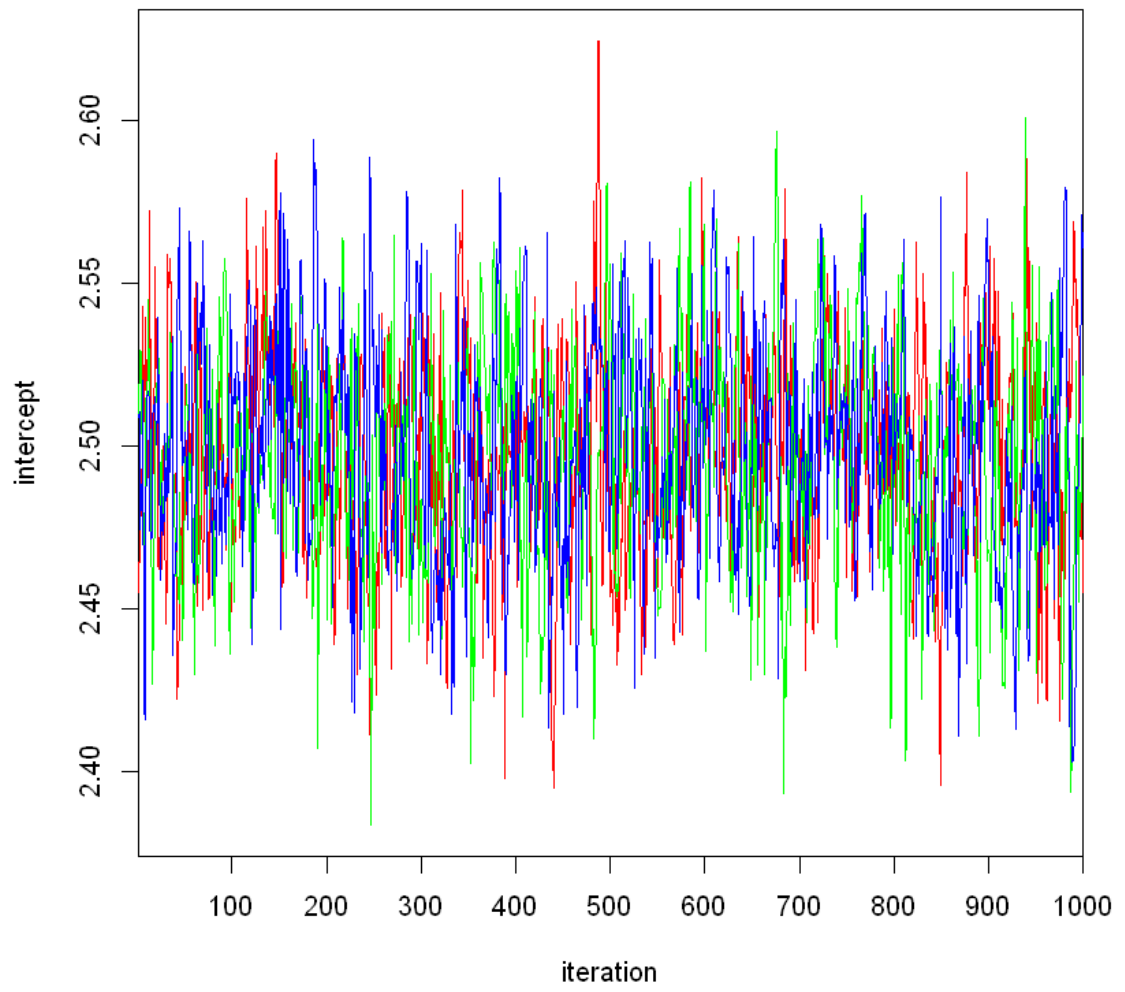
Et la trace des paramètres du modèle :

```
In [22]: traceplot(out3)
```

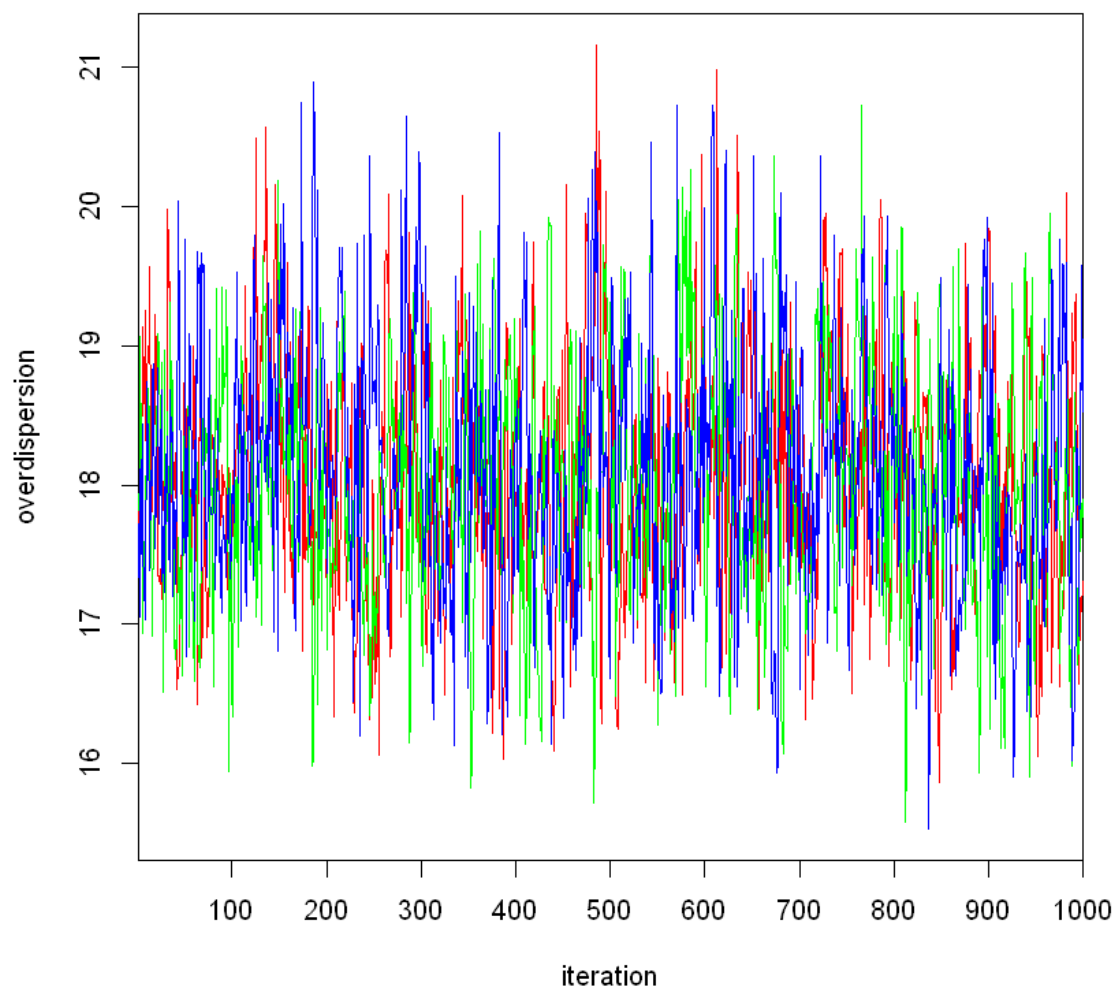

deviance



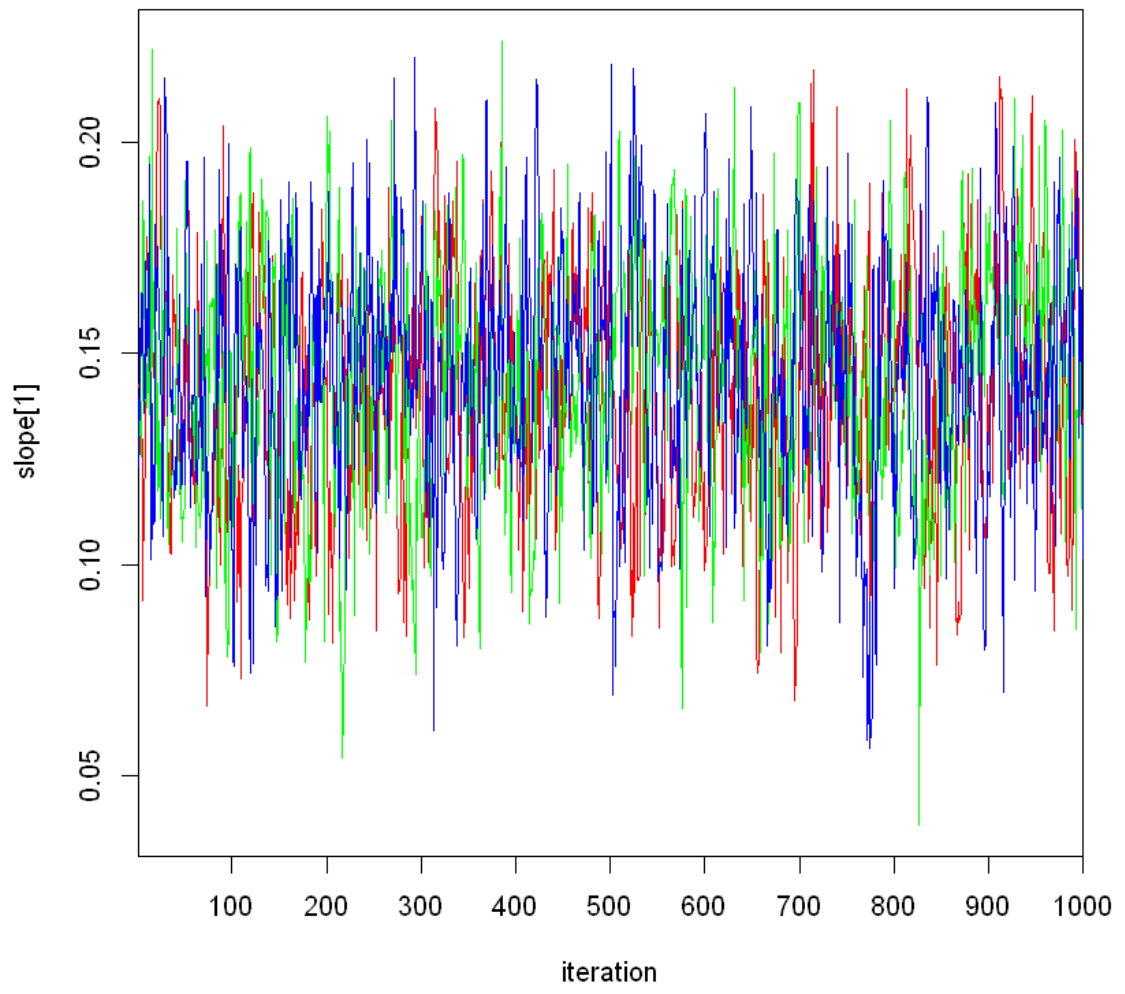
intercept



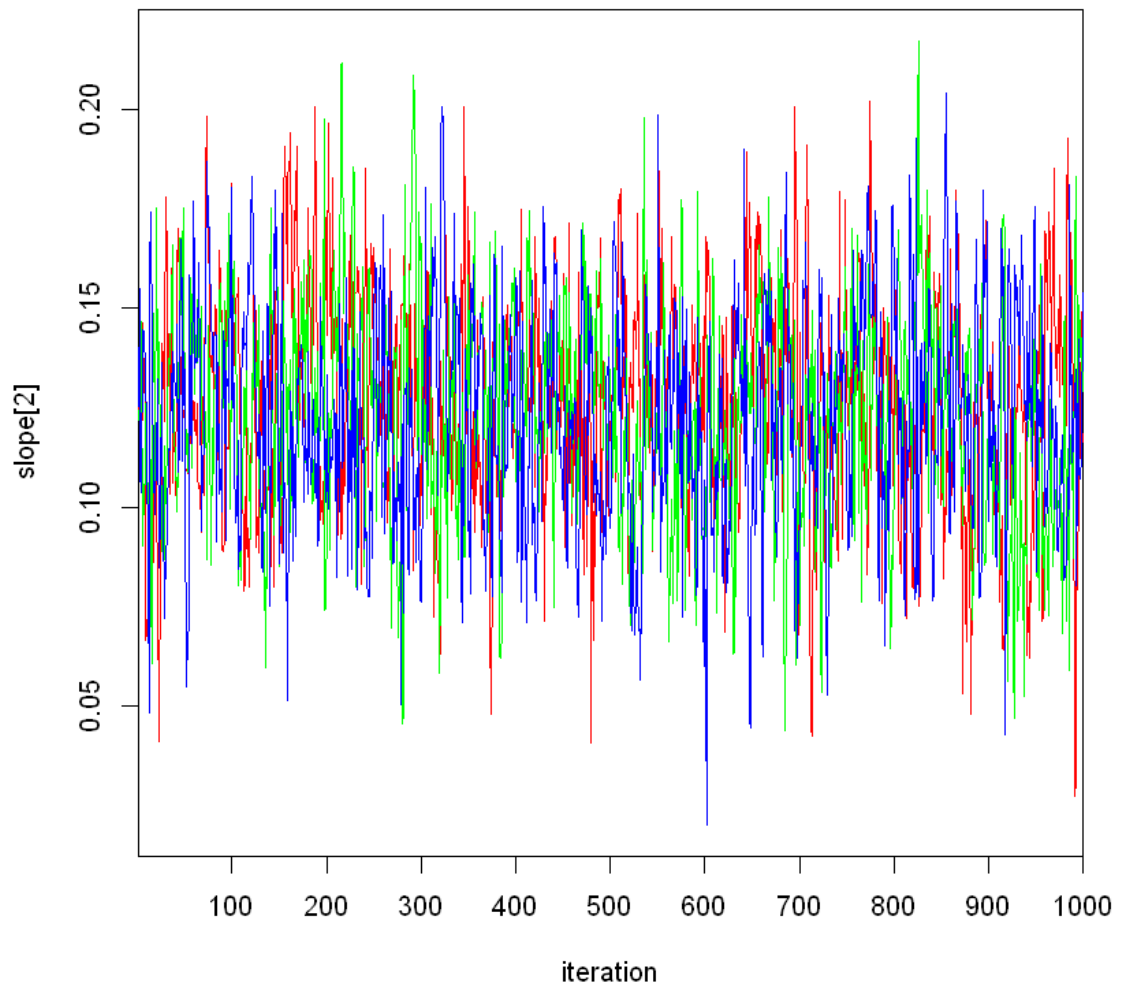
overdispersion

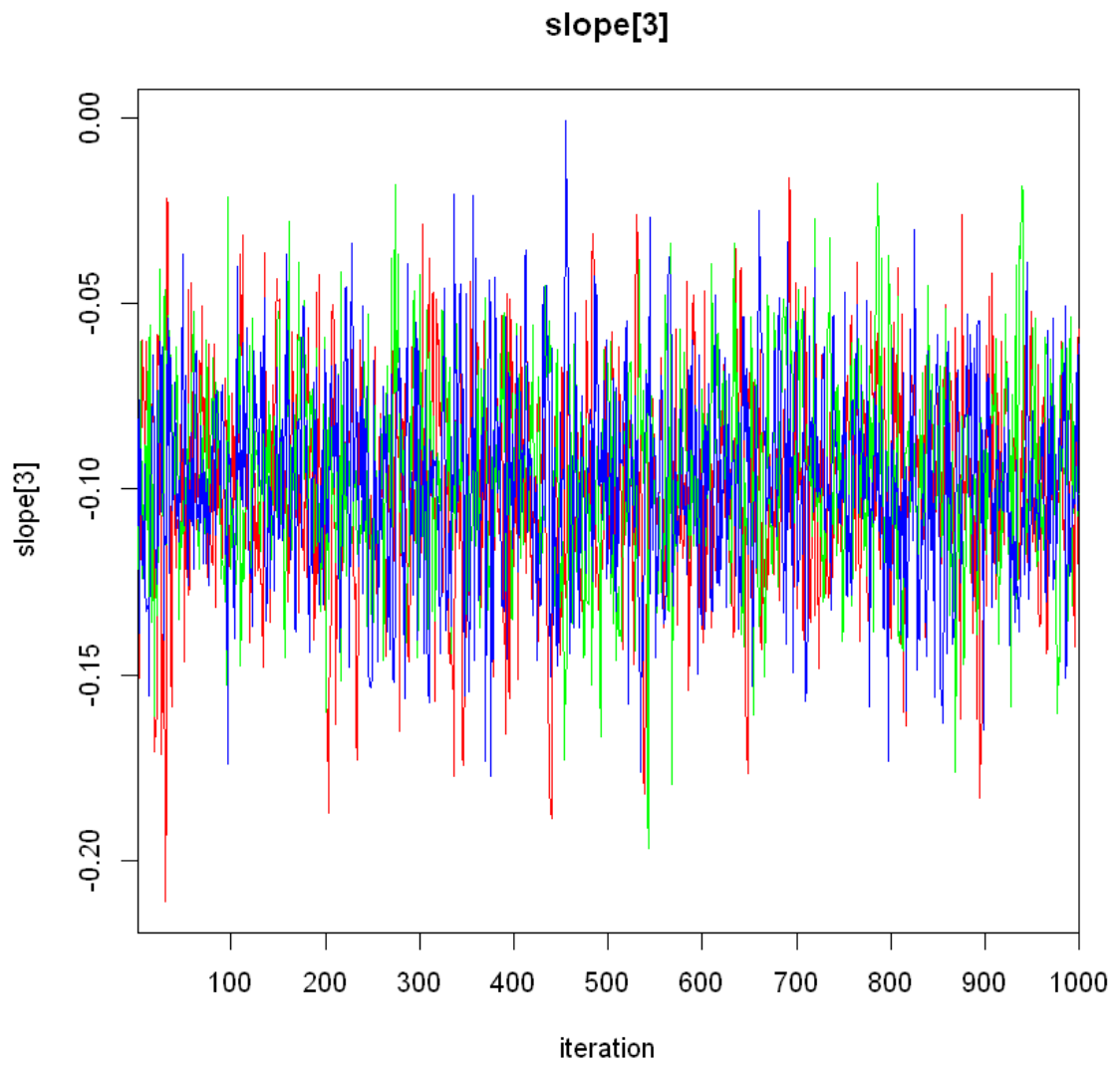


slope[1]



slope[2]





Calcul du critère d'information

Ce résultat est à comparer au modèle poisson sans surdispersion. Se souvenir que pour ces critères : *the lower, the better*

```
In [23]: y_rep_3 <- log_lik3 <- array(NA, dim = c(nc * (ni - nb)/nt, nrow(obs)))

for(i in 1:nrow(obs)) {
  overdispersion <- drop(out3$sims.list$overdispersion)
  linpred <- out3$sims.list$intercept +
    out3$sims.list$slope[, 1] * data.jags$X1[i] +
    out3$sims.list$slope[, 2] * data.jags$X2[i] +
    out3$sims.list$slope[, 3] * data.jags$X3[i]
  log_lik3[, i] <- dnbinom(data.jags$Y[i],
    prob = 1 / overdispersion,
    size = exp(linpred) / (overdispersion - 1),
    log = TRUE)
  y_rep_3[, i] <- rnbinom(nrow(y_rep_3),
    prob = 1 / overdispersion,
    size = exp(linpred) / (overdispersion - 1)
  ) + 1
}; rm(i, linpred, overdispersion)

loo::waic(log_lik3)
loo::loo(log_lik3)
```

Warning message:

"1 (0.1%) p_waic estimates greater than 0.4. We recommend trying loo instead."

Warning message:

"1 (0.1%) p_waic estimates greater than 0.4. We recommend trying loo instead."

Computed from 3000 by 1269 log-likelihood matrix

| | Estimate | SE |
|-----------|----------|-------|
| elpd_waic | -4442.1 | 54.2 |
| p_waic | 6.2 | 0.8 |
| waic | 8884.2 | 108.4 |

Warning message:

"Relative effective sample sizes ('r_eff' argument) not specified.

For models fit with MCMC, the reported PSIS effective sample sizes and MCSE estimates will be over-optimistic."

Computed from 3000 by 1269 log-likelihood matrix

| | Estimate | SE |
|----------|----------|-------|
| elpd_loo | -4442.1 | 54.2 |
| p_loo | 6.2 | 0.8 |
| looic | 8884.2 | 108.4 |

Monte Carlo SE of elpd_loo is 0.0.

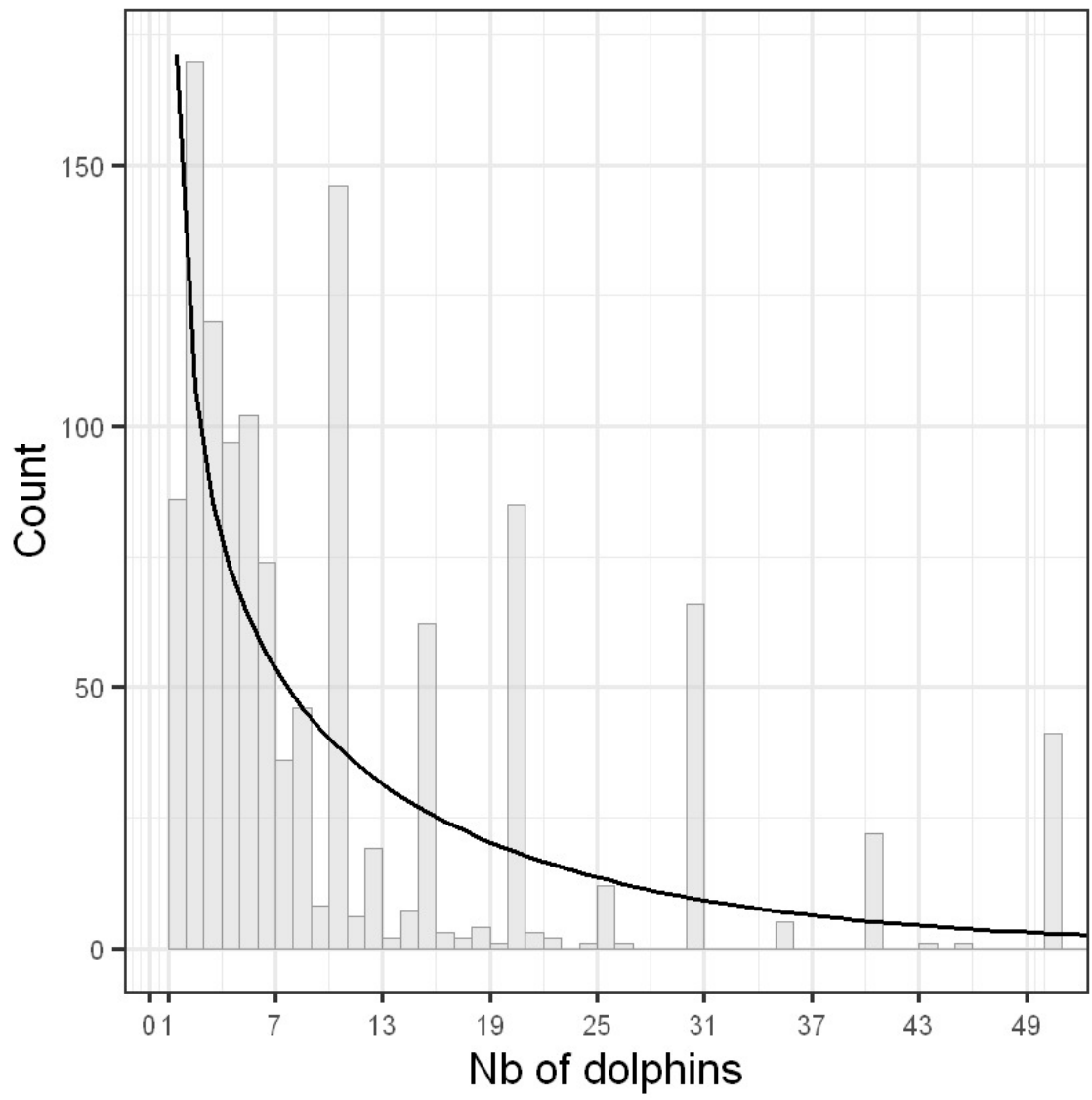
All Pareto k estimates are good (k < 0.5).

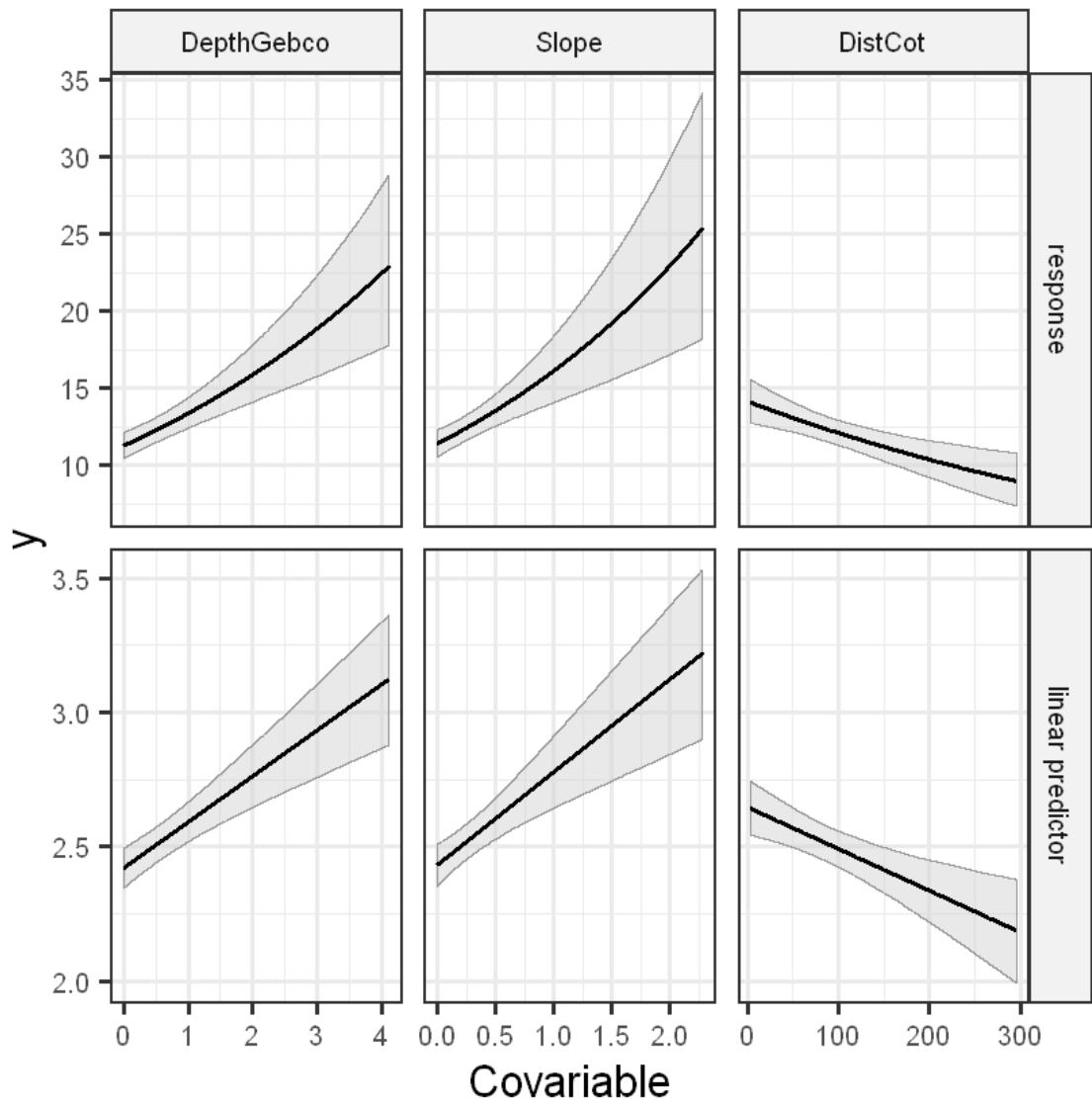
See help('pareto-k-diagnostic') for details.

Que dire de l'ajustement global du modèle aux données?

```
In [24]: rootogram_model_3 <- rootogram(countdata = obs$nombre, y_rep = y_rep_3)
rootogram_model_3 +
  coord_cartesian(xlim = c(1, 50))

plot_relationships(data_df = obs, jagsfit = out3, cov_name = c("DepthGebco", "Sl
ope", "DistCot"))
```



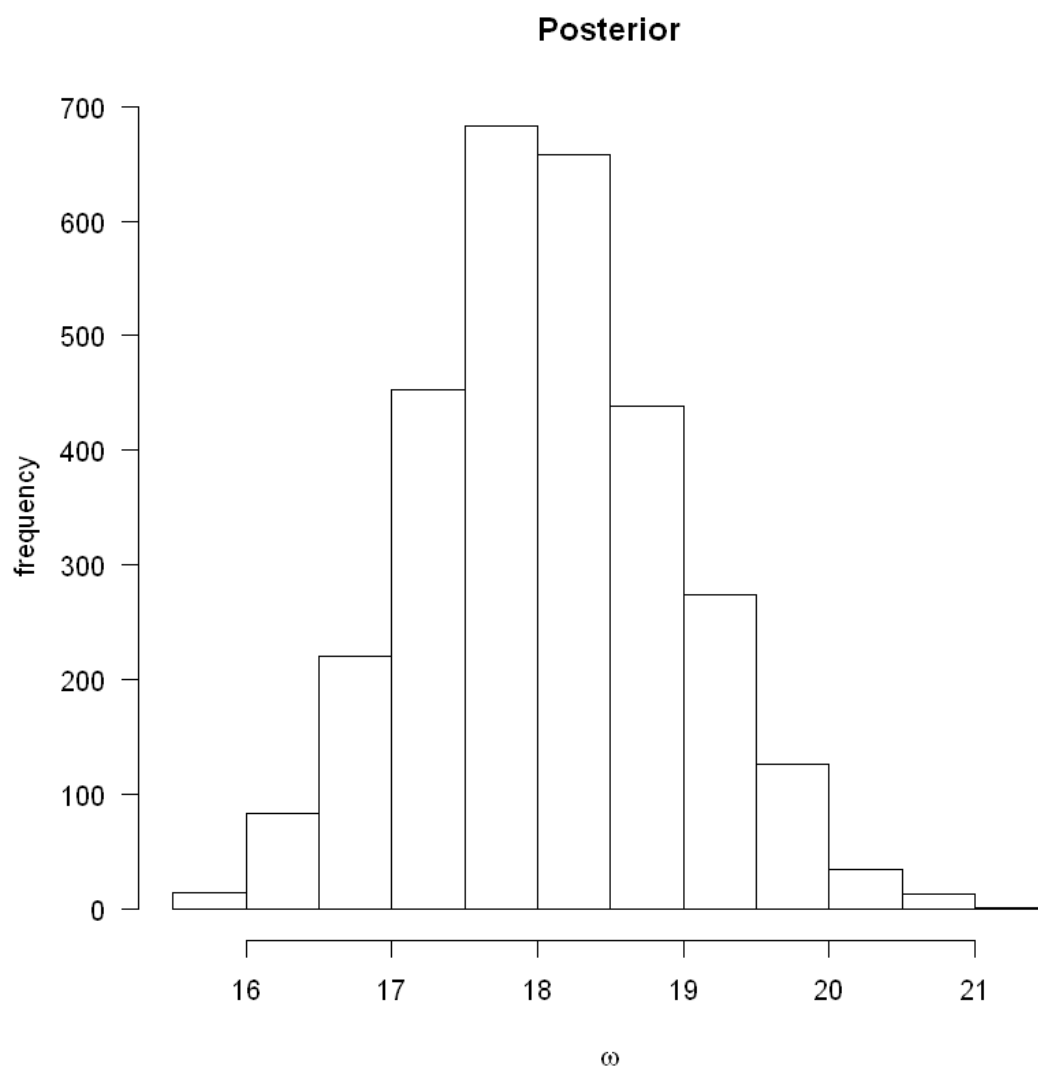
Une première satisfaction

Comme nous le disait le WAIC, l'ajustement est bien meilleur et les zones de crédibilité nous semblent plus réalistes.

- exercice: Est-ce que cela se reflète dans le DIC?
- exercice: tracer les lois a posteriori des inconnues

Qu'en est il du paramètres de surdispersion?

```
In [25]: hist(out3$sims.list$overdispersion, las = 1, bty = 'n',  
             xlab = quote(omega), ylab = "frequency", main = "Posterior")
```

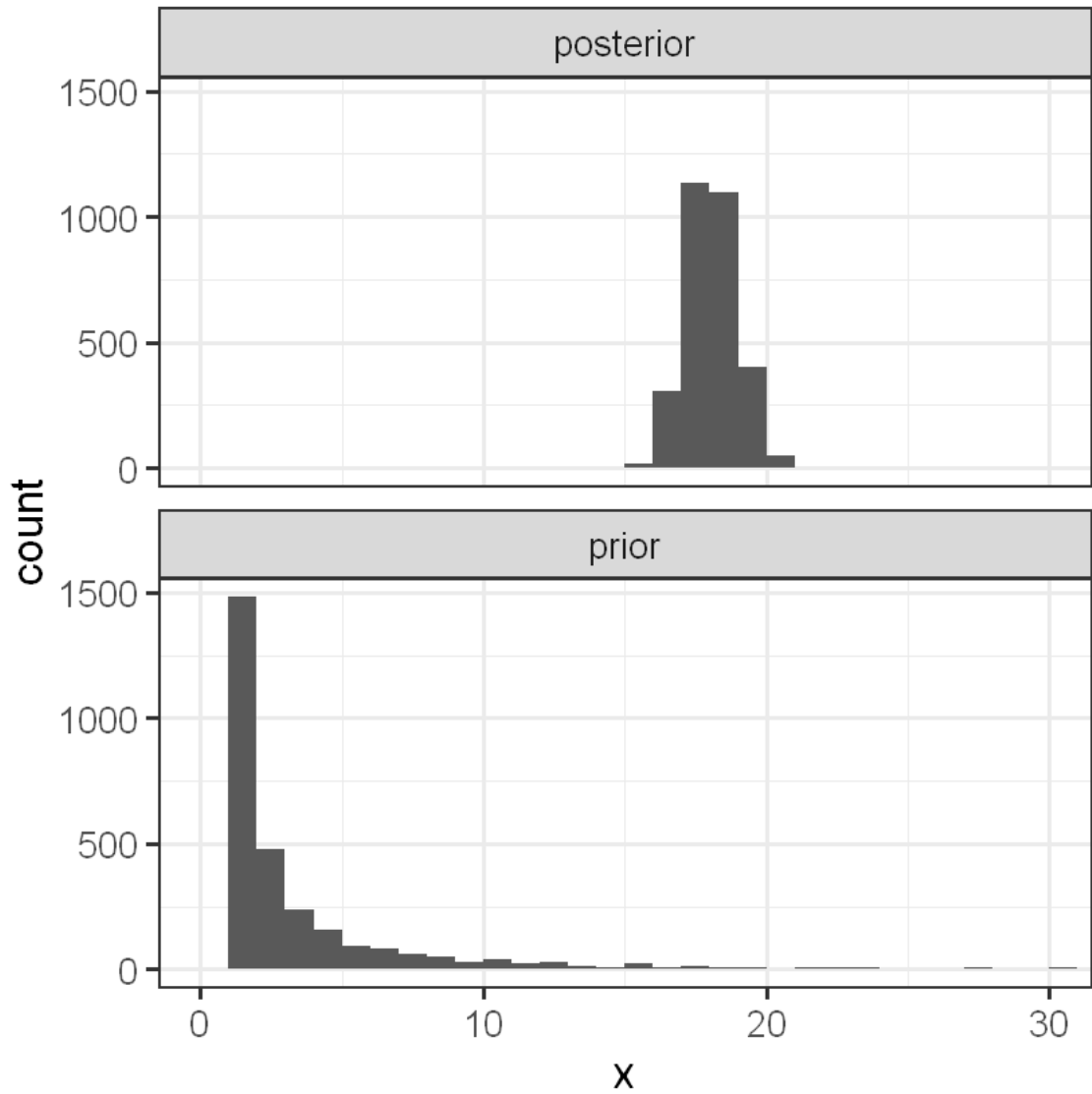


On peut comparer cette distribution a posteriori à la distribution a priori :

```

In [26]: ggplot(data = data.frame(x = c(1 / runif(length(out3$sims.list$overdispersion),
0.0, 1.0),
                                out3$sims.list$overdispersion
                                ),
                distribution = rep(c("prior", "posterior"), each = length(out3$sims.list$overdispersion))
                ),
          aes(x = x, group = distribution)
          ) +
  geom_histogram(breaks = c(0:100)) +
  facet_wrap(~distribution, ncol = 1) +
  coord_cartesian(xlim = c(0, 30))

```



Modéliser la surdispersion de manière plus explicite

Il existe plusieurs processus écologiques susceptibles de générer de la surdispersion dans les données de comptage ici, par exemple des structures d'aggrégation des proies dont dépendent les dauphins et qui ne seraient que partiellement prises en compte par les covariables, *etc.* On a jusqu'à présent ignoré certaines structures dans les données, notamment le fait que celles-ci correspondent à des années différentes. Or, il est possible que certaines années, pour des raisons que l'on comprend mal, on ait vu plus ou moins de dauphins. En d'autres termes, il est possible qu'il y ait une dépendance temporelle dans les données qui n'a pas été prise en compte.

On va donc tenir compte de cette effet année, et le modéliser en incluant dans le modèle l'année comme effet aléatoire, c'est-à-dire on va supposer que les années sont échangeables : il n'est pas possible *a priori* de classer les années en fonction de l'abondance d'animaux. On va supposer une distribution commune aux effets années, et par convenance on va supposer une distribution normale :

$$\forall t, \text{year}_t \sim \mathcal{N}(0, \sigma_{\text{year}})$$

Le paramètre $\sigma_{\text{year}} (> 0)$ quantifie la variabilité inter-annuelle. On suppose ici que l'effet moyen d'une année prise au hasard est centrée sur 0, mais peut s'écarter de 0 en fonction de ce qui est attendu sous une loi normale d'écart-type σ_{year} . En d'autres termes, en plus d'estimer les paramètres year_t , nous allons estimer l'hyperparamètre σ_{year} qui va gouverner la distribution normale. Il faut donc spécifier des priors pour cet hyperparamètre.

```

In [27]: model.jags <- '
model{
  # Prior
  intercept ~ dnorm(0.0, hyperparam_inter)
  slope[1] ~ dnorm(0.0, hyperparam_slope[1])
  slope[2] ~ dnorm(0.0, hyperparam_slope[2])
  slope[3] ~ dnorm(0.0, hyperparam_slope[3])
  inv_overdispersion ~ dunif(0, 1)
  overdispersion <- 1/inv_overdispersion
  sigma_year ~ dnorm(0.0, 1.0)T(0.0,)
  tau_alpha <- pow(sigma_year, -2)
  for(j in 1:n_year) {
    alpha[j] ~ dnorm(intercept, tau_alpha)
  }

  # Likelihood
  for (i in 1:n_obs) {
    lambda[i] <- exp(alpha[YEAR[i]] + slope[1] * X1[i] + slope[2] * X2[i] + slope[3] * X3[i])
    r[i] <- lambda[i] / (overdispersion - 1)
    Y[i] ~ dnegbin(inv_overdispersion, r[i])
  }
}
'

data.jags <- list(n_obs = nrow(obs),
                 n_year = length(unique(obs$an)),
                 Y = obs$nombreShift,
                 X1 = as.vector(scale(obs$DepthGebco)),
                 X2 = as.vector(scale(obs$Slope)),
                 X3 = as.vector(scale(obs$DistCot)),
                 YEAR = obs$an - min(obs$an) + 1,
                 hyperparam_inter = 1.0,
                 hyperparam_slope = 1/rep(log(2)/2, 3)^2
                )

# Inits function
inits.jags <- function(idchain, glm) {
  list(intercept = rnorm(1, glm$coefficients[1], sqrt(10 * diag(vcov(glm))[1])),
        slope = rnorm(3, glm$coefficients[2:4], sqrt(10 * diag(vcov(glm))[2:4])),
        inv_overdispersion = runif(1, 0, 1),
        sigma_year = abs(rnorm(1)),
        alpha = rnorm(length(unique(obs$an)))
        )
}

params.jags <- c("intercept", "slope", "overdispersion", "alpha", "sigma_year")

# Start Gibbs sampler
temp <- jags(data = data.jags,
             inits = lapply(1:nc, inits.jags, glm = dauphin2.glm),
             param = params.jags,
             model.file = textConnection(model.jags),
             n.chains = nc,
             n.burnin = nb,
             n.iter = ni,
             n.thin = nt
            )

out4 = temp$BUGSoutput
# Inferences
# After Burn in period let's work inference assuming ergodic regime is reached
print(out4$summary, dig = 2)

```

Compiling model graph
 Resolving undeclared variables
 Allocating nodes
 Graph information:
 Observed stochastic nodes: 1269
 Unobserved stochastic nodes: 19
 Total graph size: 11788

Initializing model

| | mean | sd | 2.5% | 25% | 50% | 75% | 97.5% | Rha |
|----------------|----------|-------|----------|----------|----------|----------|----------|-----|
| t | | | | | | | | |
| alpha[1] | 2.482 | 0.224 | 1.995 | 2.348 | 2.493 | 2.634 | 2.893 | |
| 1 | | | | | | | | |
| alpha[2] | 2.998 | 0.189 | 2.615 | 2.874 | 3.004 | 3.127 | 3.356 | |
| 1 | | | | | | | | |
| alpha[3] | 2.748 | 0.188 | 2.365 | 2.625 | 2.753 | 2.877 | 3.100 | |
| 1 | | | | | | | | |
| alpha[4] | 2.578 | 0.248 | 2.058 | 2.422 | 2.587 | 2.738 | 3.062 | |
| 1 | | | | | | | | |
| alpha[5] | 2.711 | 0.099 | 2.506 | 2.645 | 2.713 | 2.781 | 2.900 | |
| 1 | | | | | | | | |
| alpha[6] | 2.748 | 0.092 | 2.564 | 2.687 | 2.749 | 2.809 | 2.919 | |
| 1 | | | | | | | | |
| alpha[7] | 2.747 | 0.066 | 2.619 | 2.703 | 2.749 | 2.794 | 2.873 | |
| 1 | | | | | | | | |
| alpha[8] | 2.602 | 0.084 | 2.439 | 2.544 | 2.604 | 2.661 | 2.761 | |
| 1 | | | | | | | | |
| alpha[9] | 2.626 | 0.068 | 2.489 | 2.581 | 2.628 | 2.673 | 2.755 | |
| 1 | | | | | | | | |
| alpha[10] | 2.225 | 0.089 | 2.052 | 2.165 | 2.225 | 2.286 | 2.396 | |
| 1 | | | | | | | | |
| alpha[11] | 2.381 | 0.091 | 2.191 | 2.319 | 2.383 | 2.444 | 2.553 | |
| 1 | | | | | | | | |
| alpha[12] | 2.294 | 0.111 | 2.070 | 2.223 | 2.297 | 2.368 | 2.500 | |
| 1 | | | | | | | | |
| alpha[13] | 2.133 | 0.065 | 2.007 | 2.090 | 2.132 | 2.177 | 2.262 | |
| 1 | | | | | | | | |
| deviance | 8781.529 | 5.870 | 8772.066 | 8777.285 | 8780.823 | 8785.192 | 8794.405 | |
| 1 | | | | | | | | |
| intercept | 2.539 | 0.098 | 2.338 | 2.479 | 2.541 | 2.603 | 2.726 | |
| 1 | | | | | | | | |
| overdispersion | 16.715 | 0.821 | 15.208 | 16.140 | 16.685 | 17.267 | 18.400 | |
| 1 | | | | | | | | |
| sigma_year | 0.308 | 0.084 | 0.181 | 0.247 | 0.295 | 0.355 | 0.501 | |
| 1 | | | | | | | | |
| slope[1] | 0.104 | 0.030 | 0.045 | 0.084 | 0.104 | 0.125 | 0.160 | |
| 1 | | | | | | | | |
| slope[2] | 0.097 | 0.030 | 0.034 | 0.077 | 0.098 | 0.118 | 0.153 | |
| 1 | | | | | | | | |
| slope[3] | -0.067 | 0.030 | -0.127 | -0.087 | -0.067 | -0.047 | -0.011 | |
| 1 | | | | | | | | |
| | n.eff | | | | | | | |
| alpha[1] | 1200 | | | | | | | |
| alpha[2] | 2400 | | | | | | | |
| alpha[3] | 3000 | | | | | | | |
| alpha[4] | 3000 | | | | | | | |
| alpha[5] | 540 | | | | | | | |
| alpha[6] | 920 | | | | | | | |
| alpha[7] | 2700 | | | | | | | |
| alpha[8] | 3000 | | | | | | | |
| alpha[9] | 3000 | | | | | | | |
| alpha[10] | 2100 | | | | | | | |
| alpha[11] | 3000 | | | | | | | |
| alpha[12] | 550 | | | | | | | |
| alpha[13] | 1300 | | | | | | | |
| deviance | 240 | | | | | | | |
| intercept | 1900 | | | | | | | |
| overdispersion | 800 | | | | | | | |

On peut également vérifier la convergence des différents paramètres ...

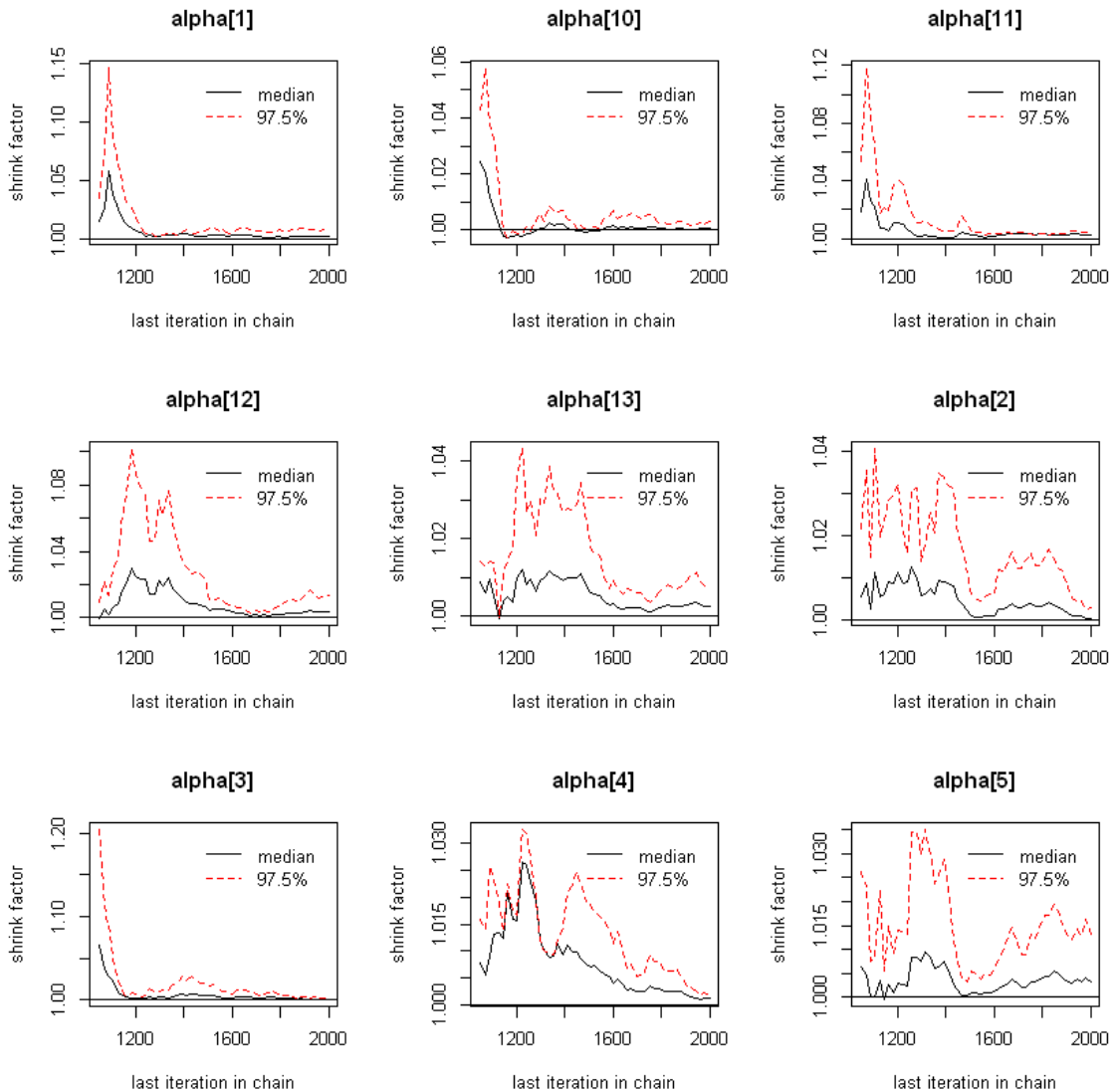

```
In [28]: library(coda)
gelman.diag(x=out4, confidence = 0.95, transform=FALSE, autoburnin=TRUE)
gelman.plot(x=out4)
```

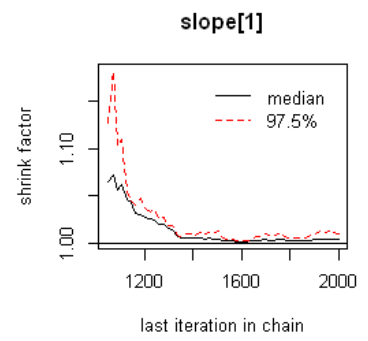
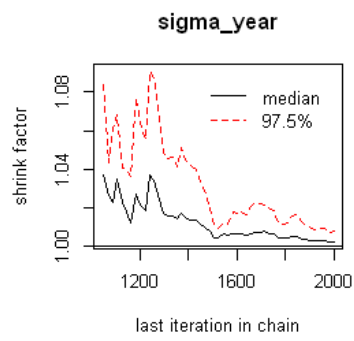
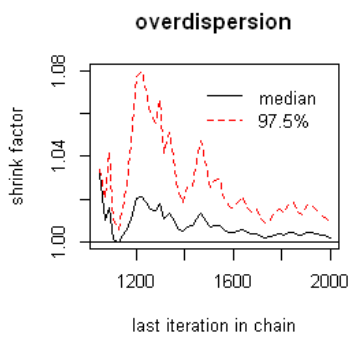
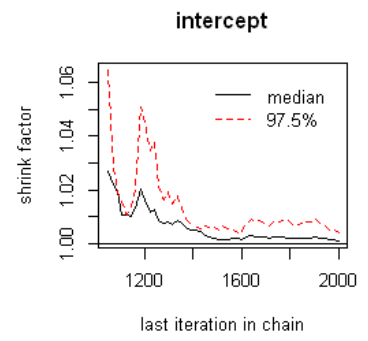
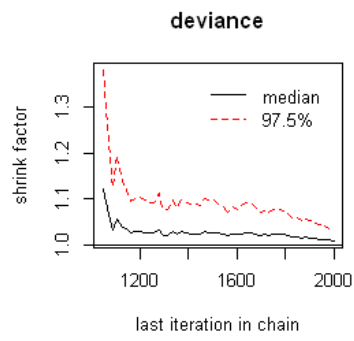
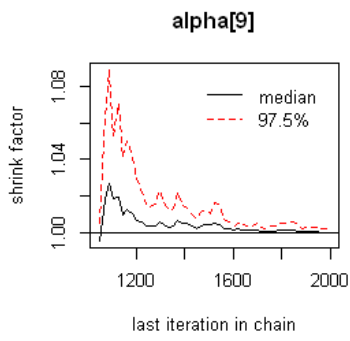
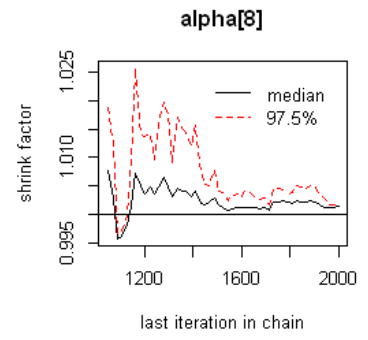
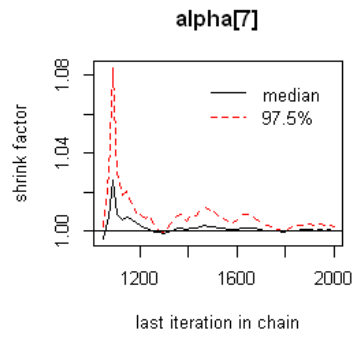
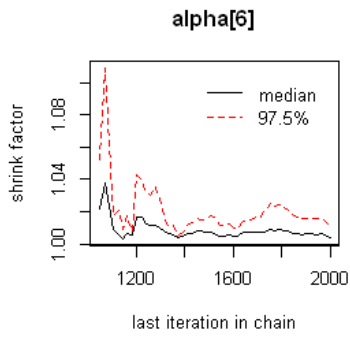
Potential scale reduction factors:

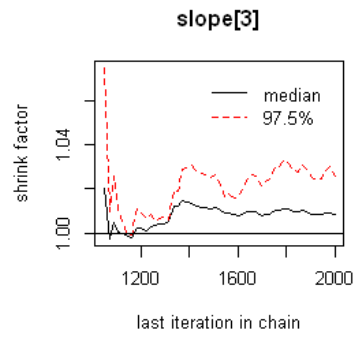
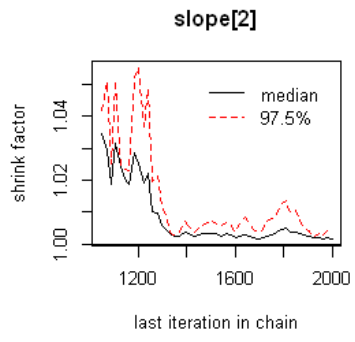
| | Point est. | Upper C.I. |
|----------------|------------|------------|
| alpha[1] | 1.00 | 1.01 |
| alpha[10] | 1.00 | 1.00 |
| alpha[11] | 1.00 | 1.00 |
| alpha[12] | 1.00 | 1.01 |
| alpha[13] | 1.00 | 1.01 |
| alpha[2] | 1.00 | 1.00 |
| alpha[3] | 1.00 | 1.00 |
| alpha[4] | 1.00 | 1.00 |
| alpha[5] | 1.00 | 1.01 |
| alpha[6] | 1.00 | 1.01 |
| alpha[7] | 1.00 | 1.00 |
| alpha[8] | 1.00 | 1.00 |
| alpha[9] | 1.00 | 1.00 |
| deviance | 1.01 | 1.03 |
| intercept | 1.00 | 1.00 |
| overdispersion | 1.00 | 1.01 |
| sigma_year | 1.00 | 1.01 |
| slope[1] | 1.00 | 1.01 |
| slope[2] | 1.00 | 1.00 |
| slope[3] | 1.01 | 1.03 |

Multivariate psrf

1.02



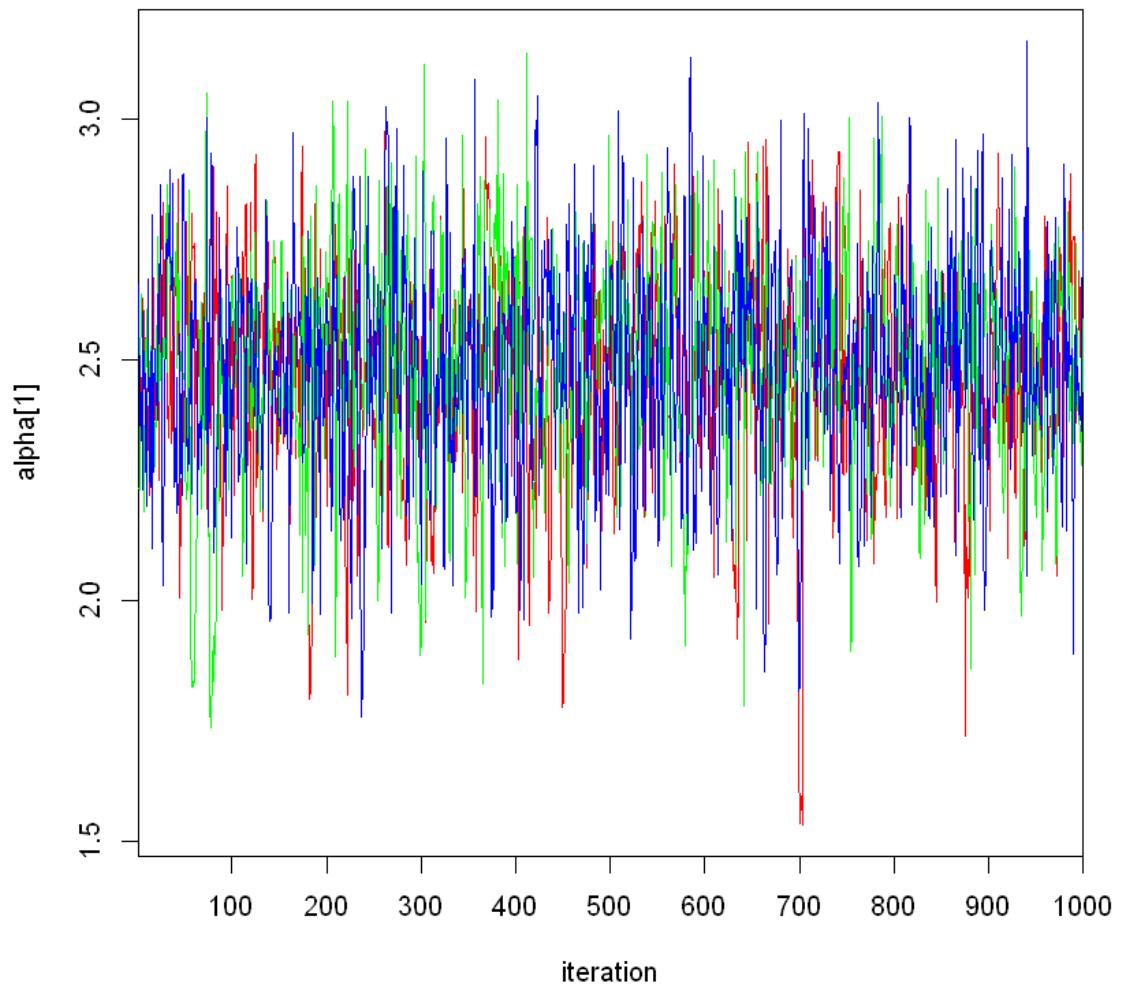




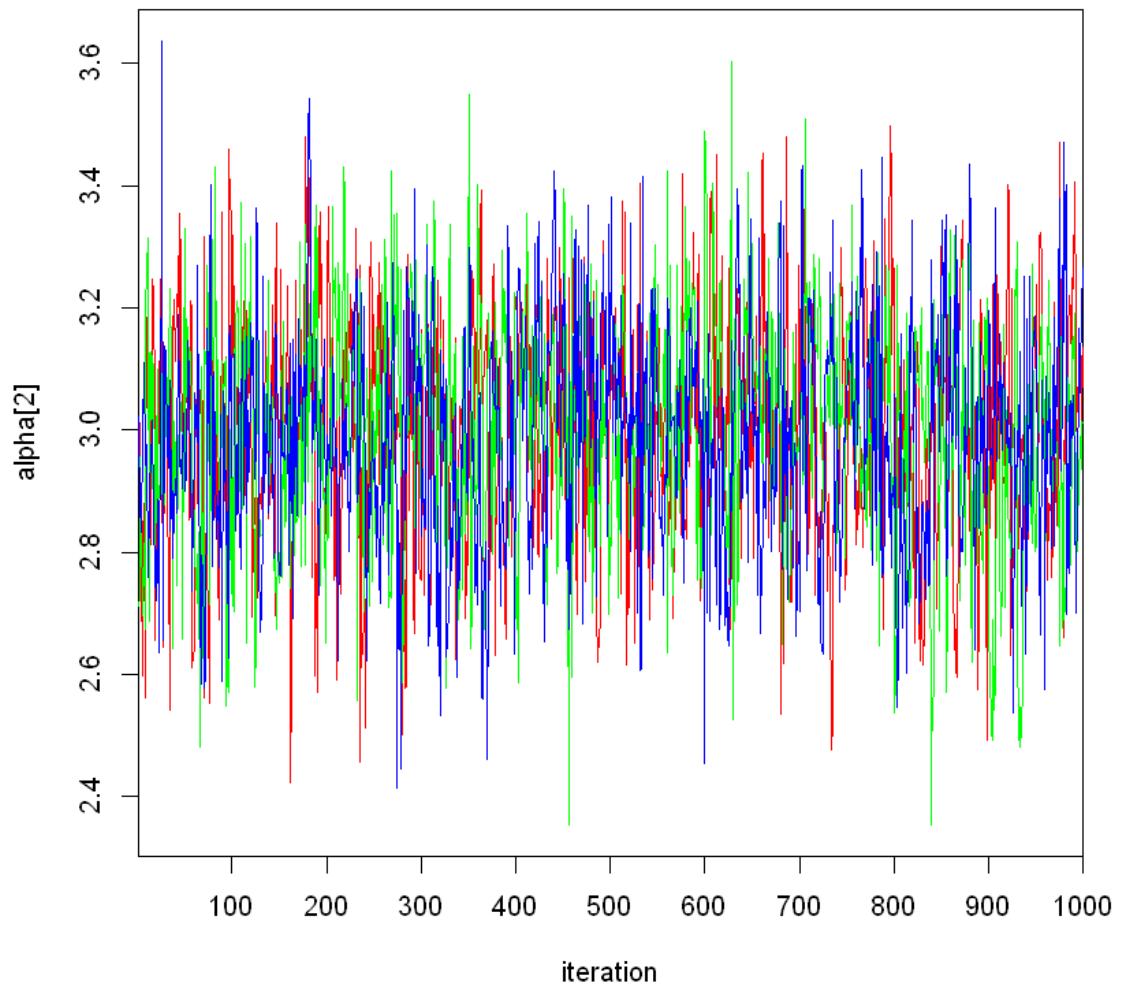
et leur trace

```
In [29]: traceplot(out4)
```

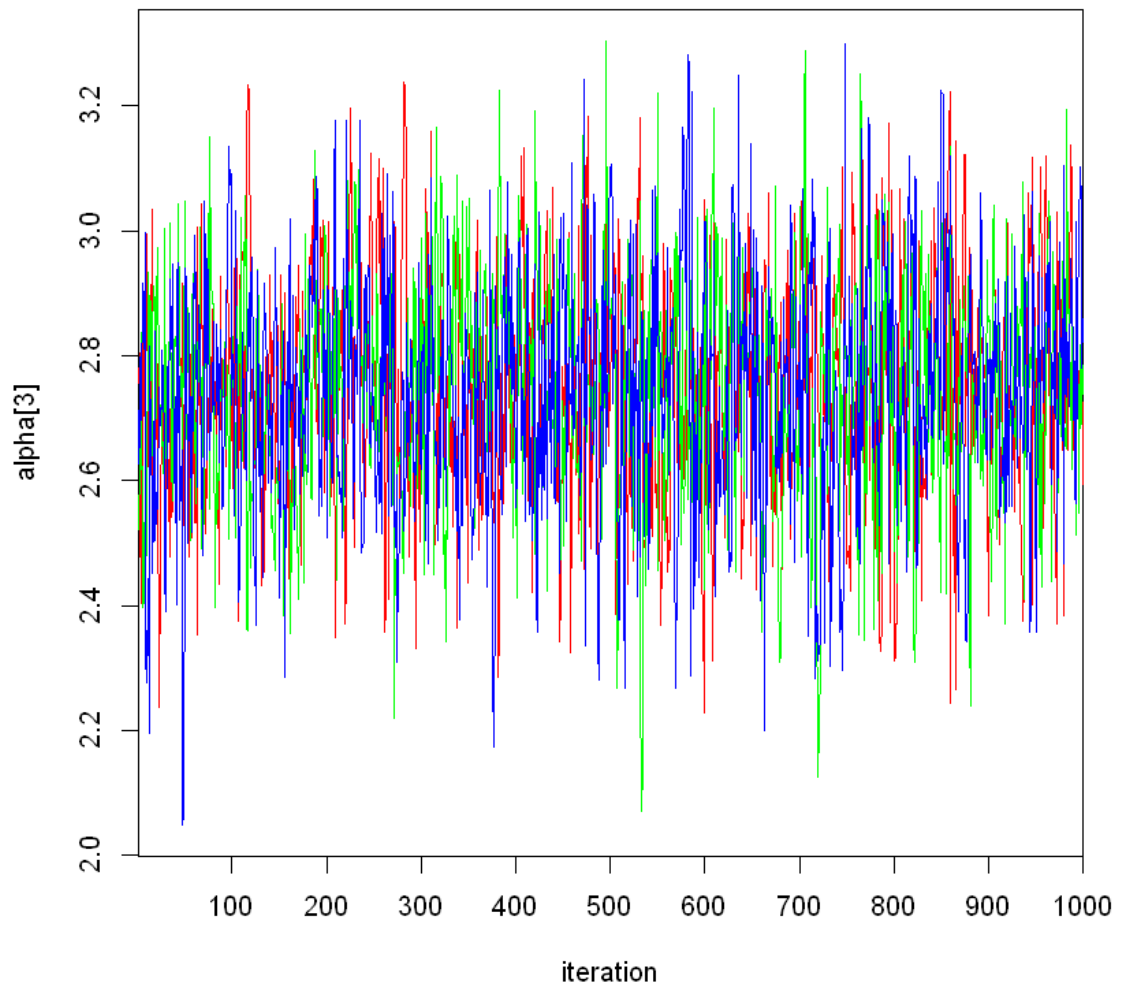
alpha[1]



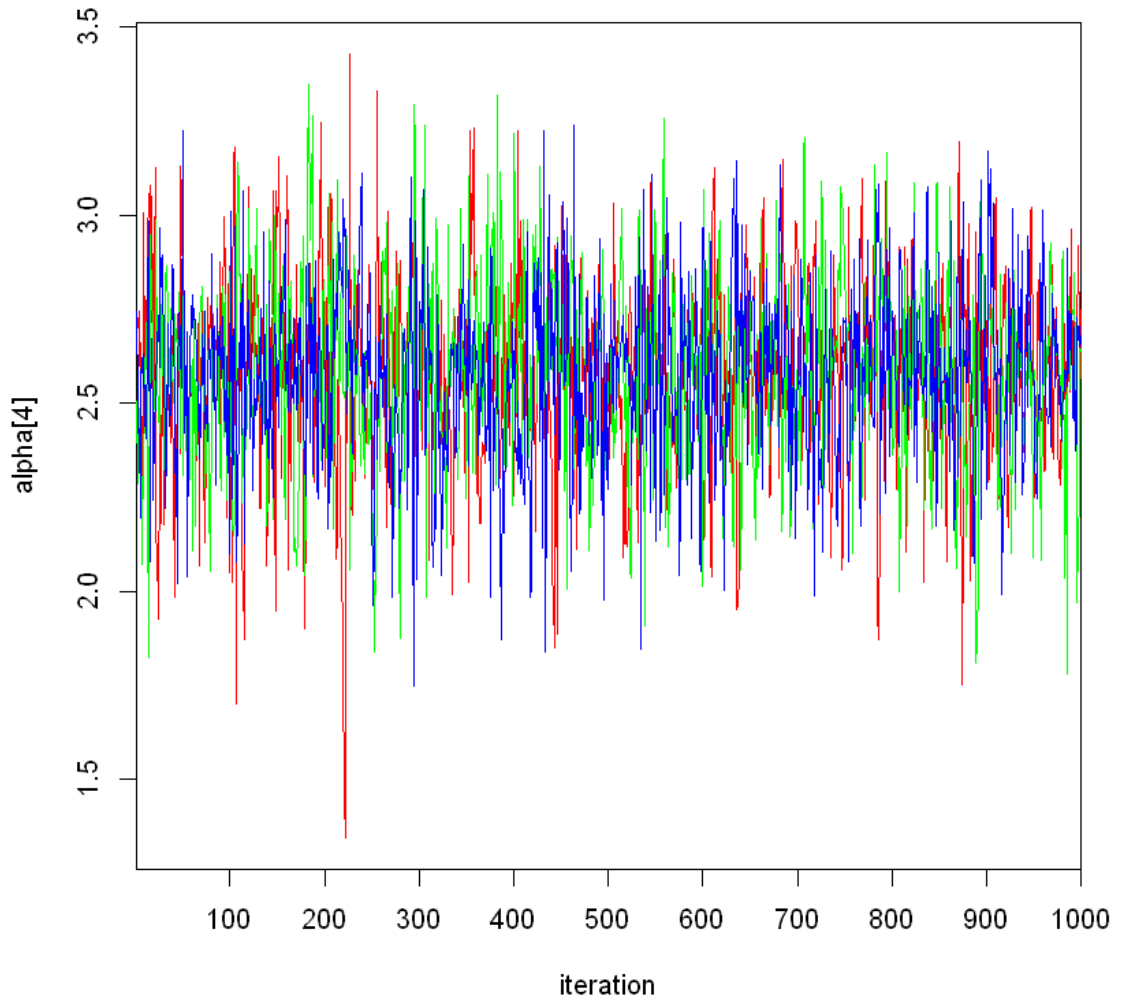
alpha[2]



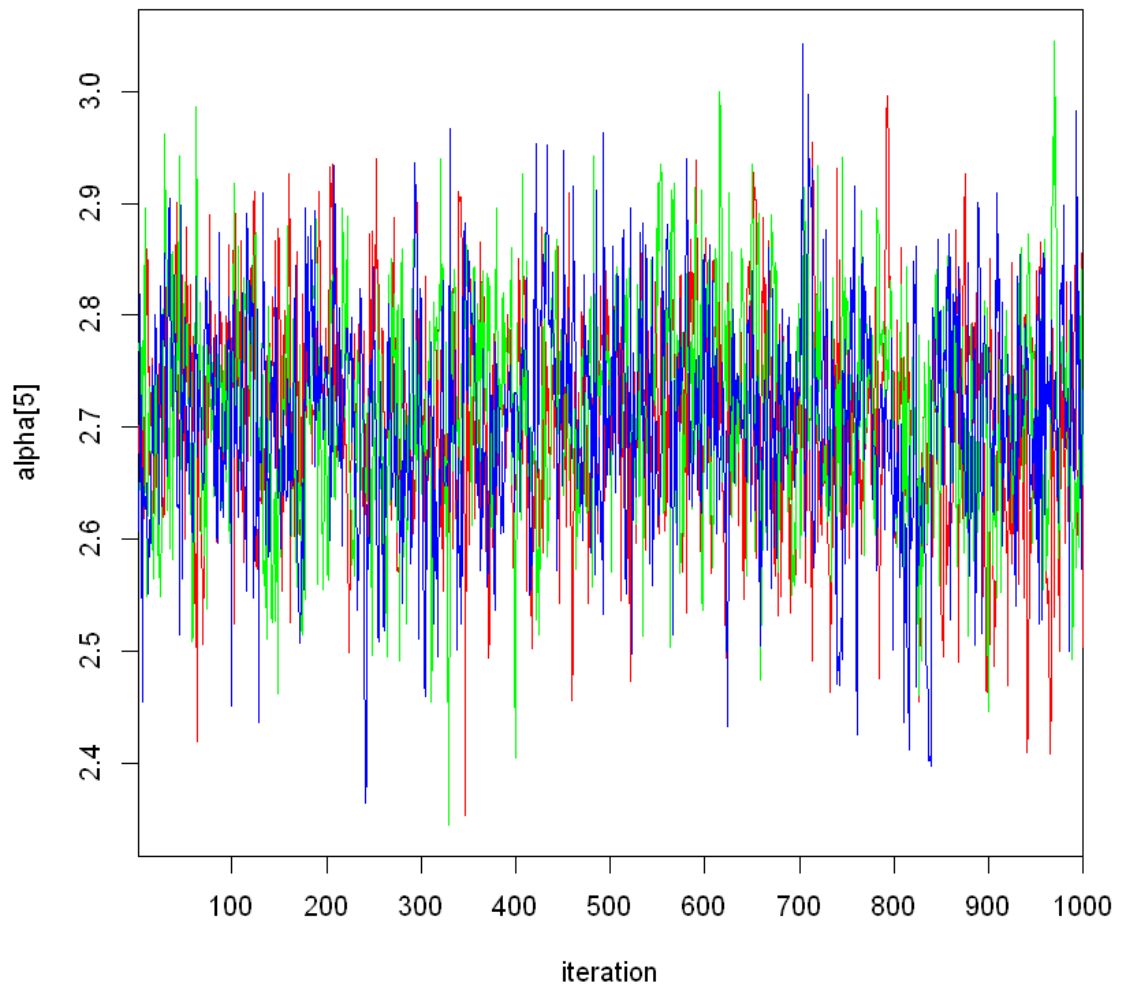
alpha[3]



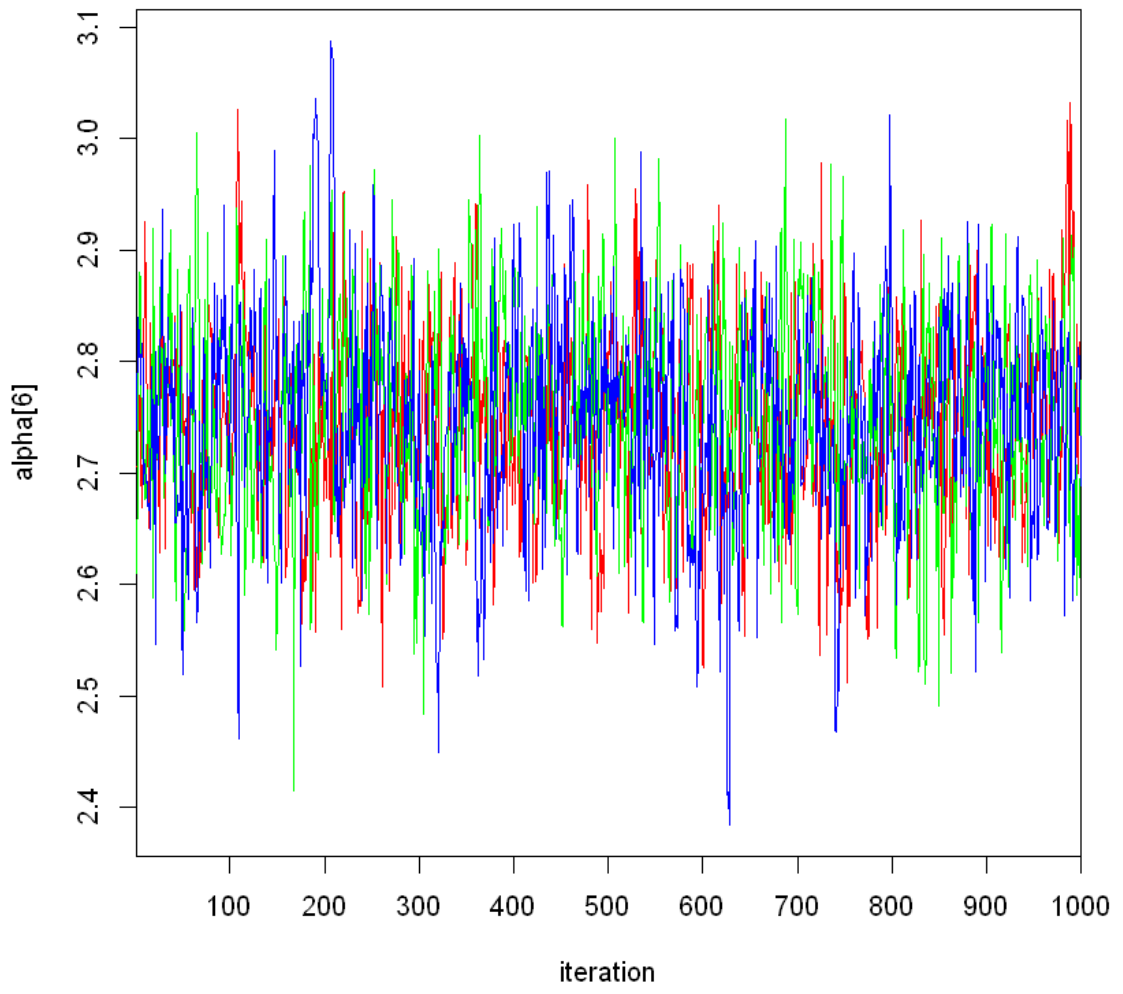
alpha[4]



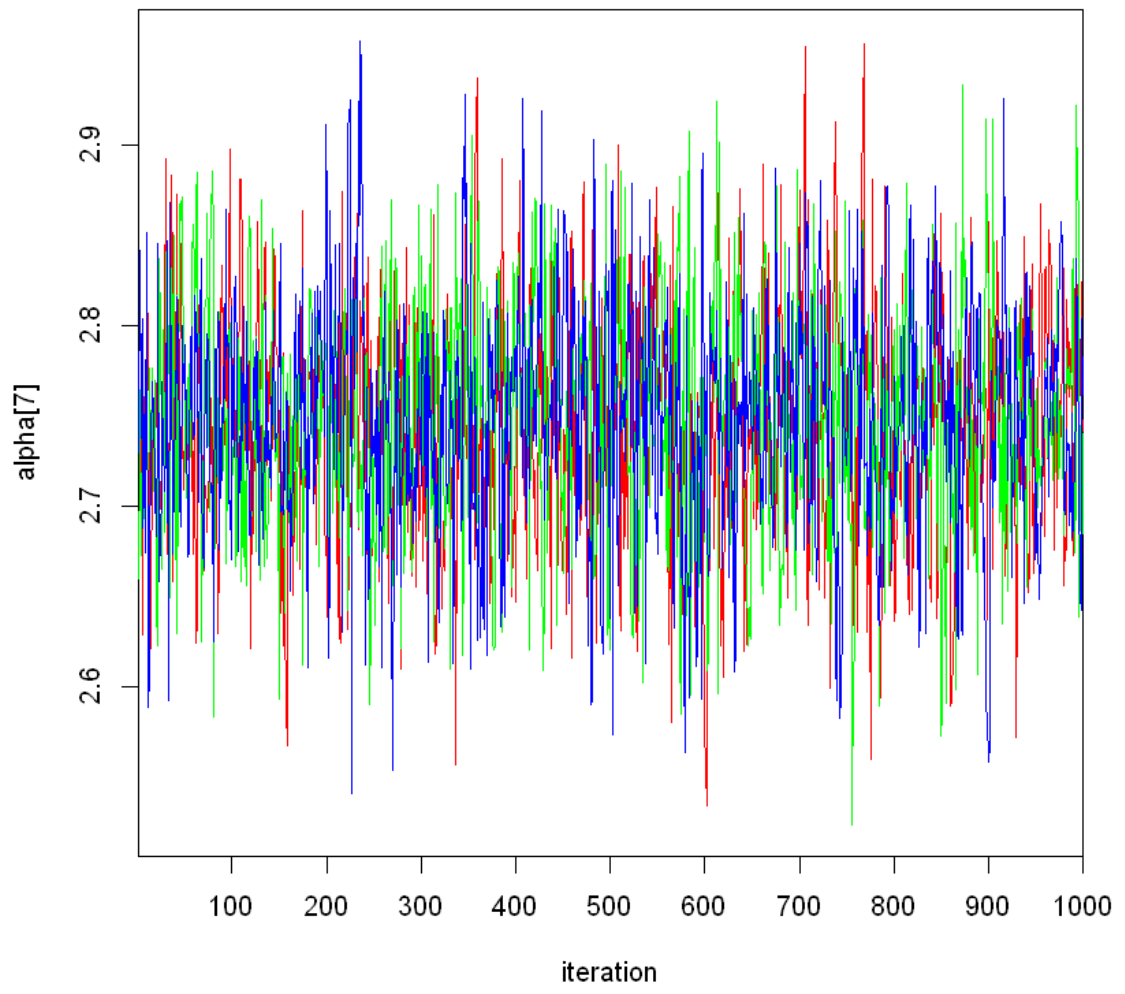
alpha[5]

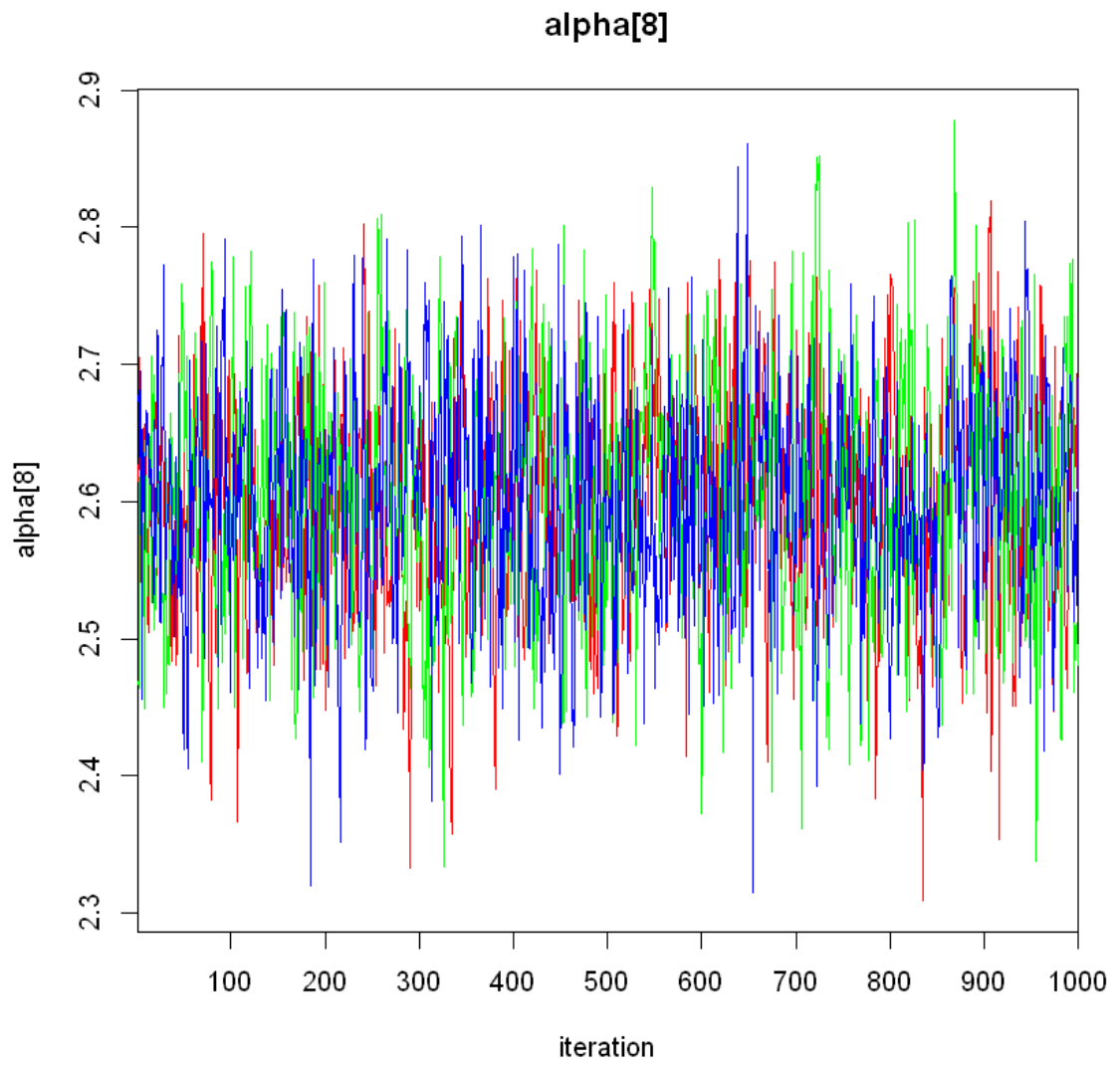


alpha[6]

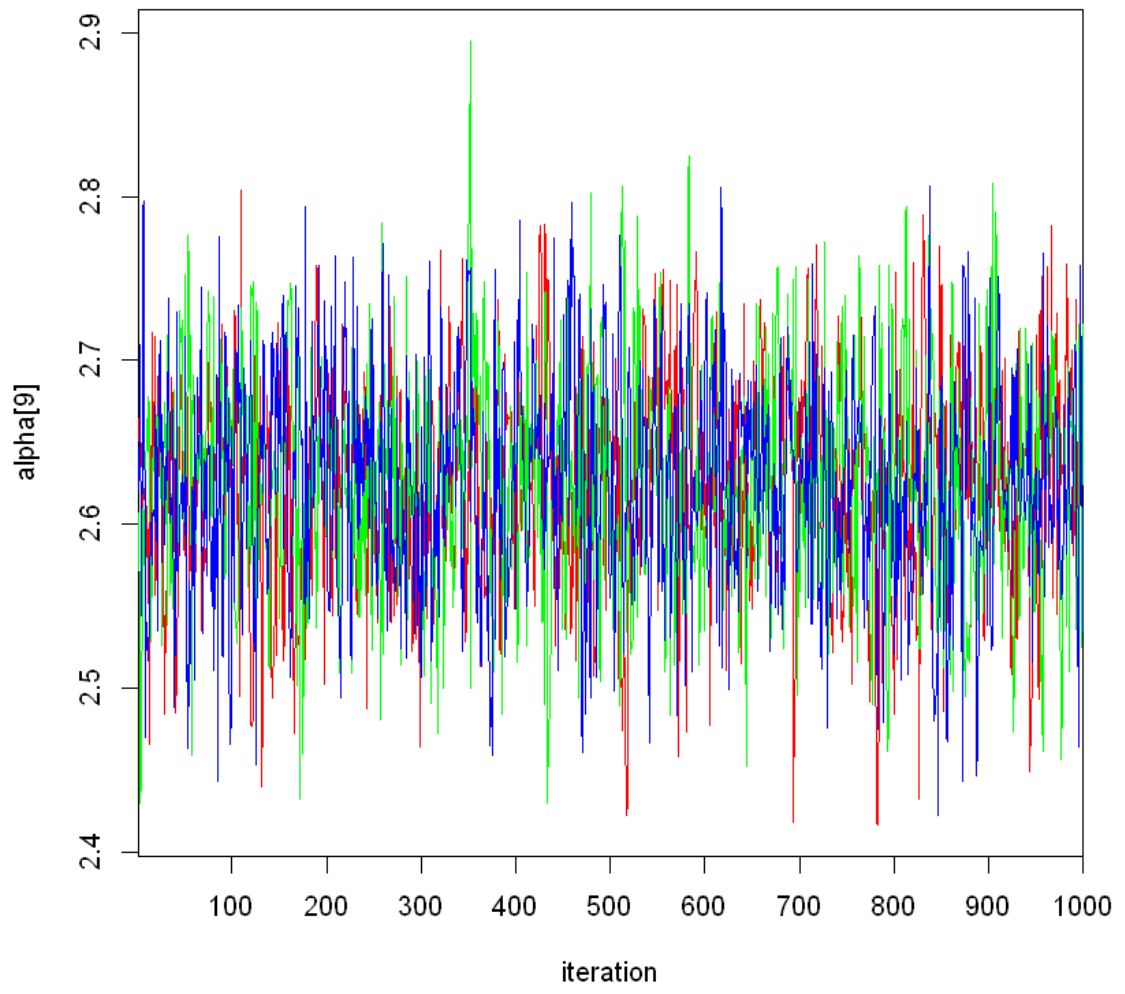


alpha[7]

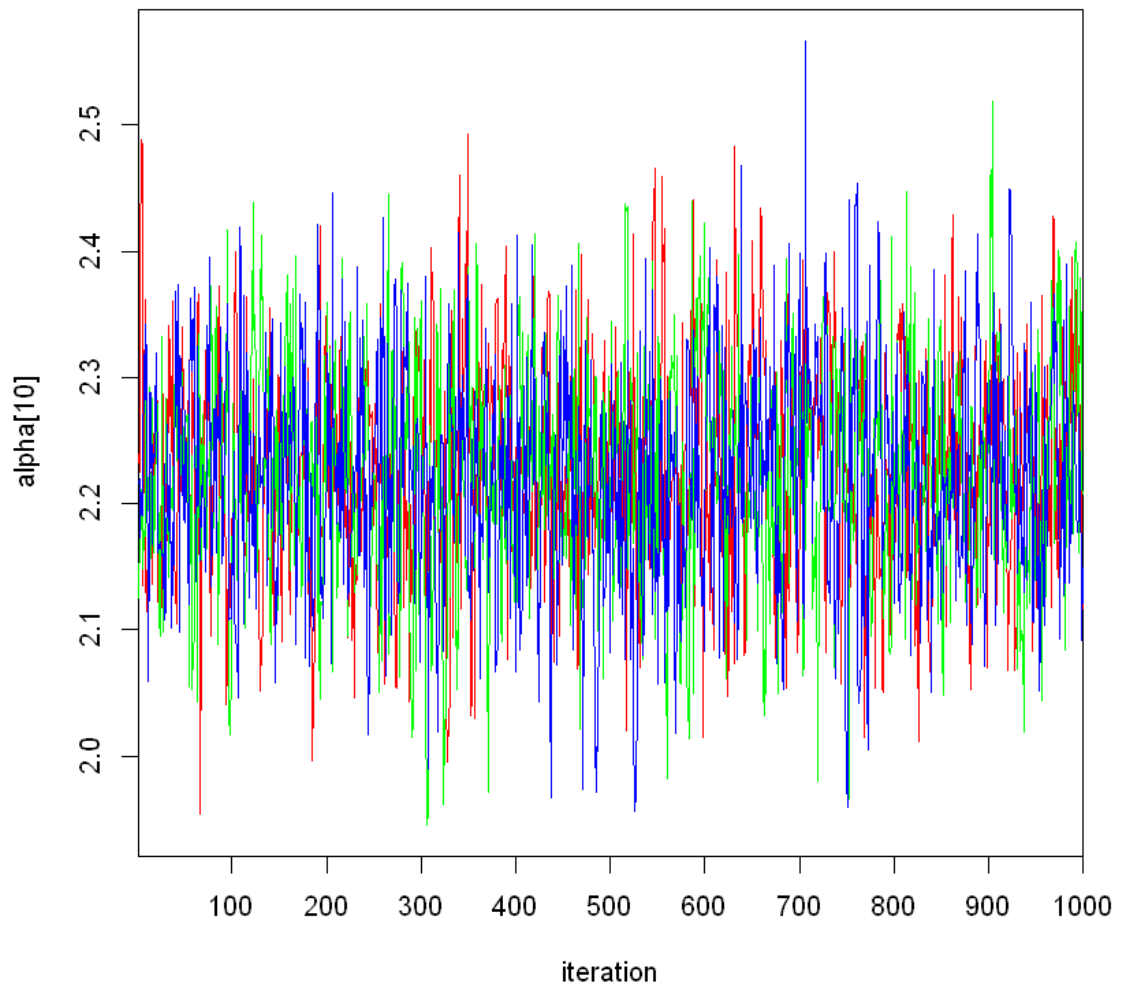




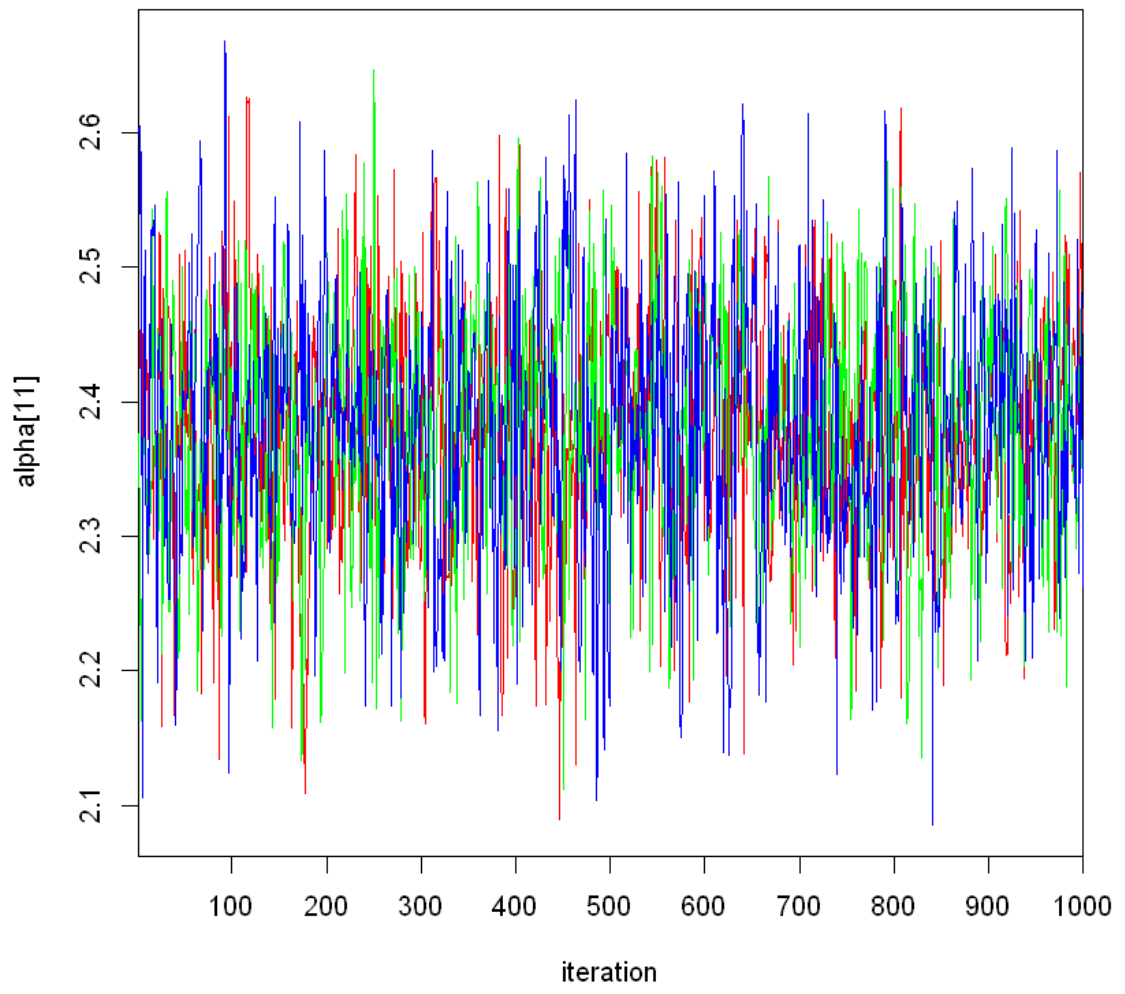
alpha[9]



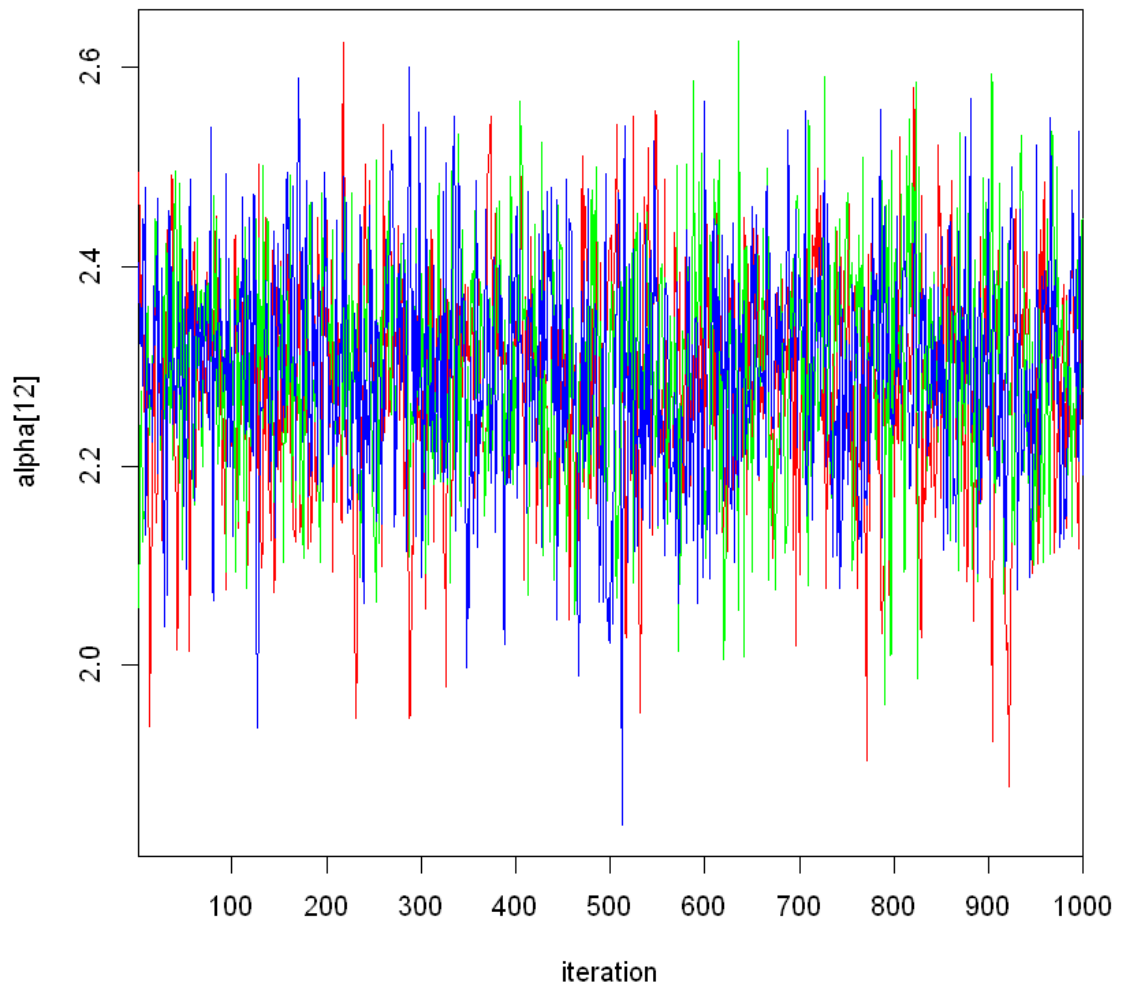
alpha[10]



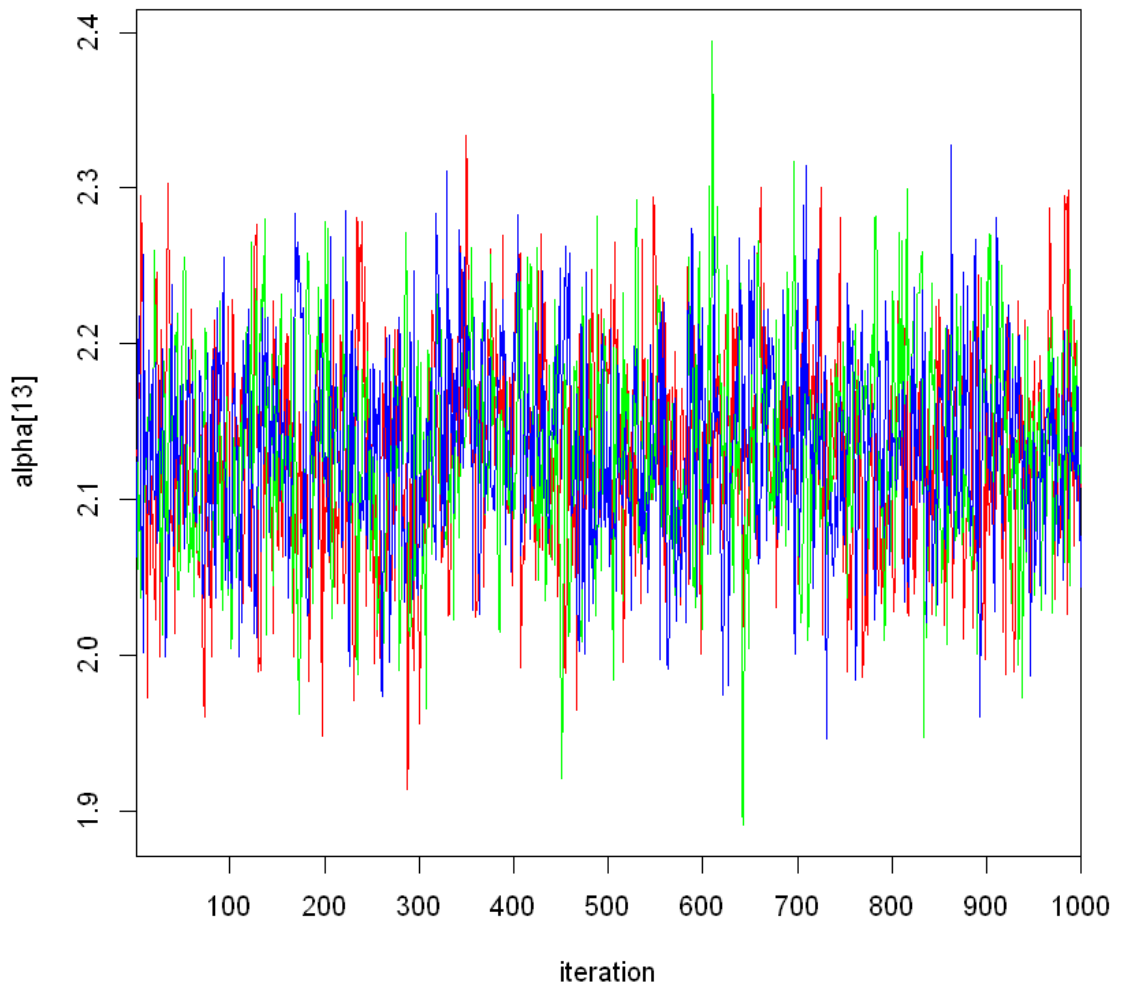
alpha[11]



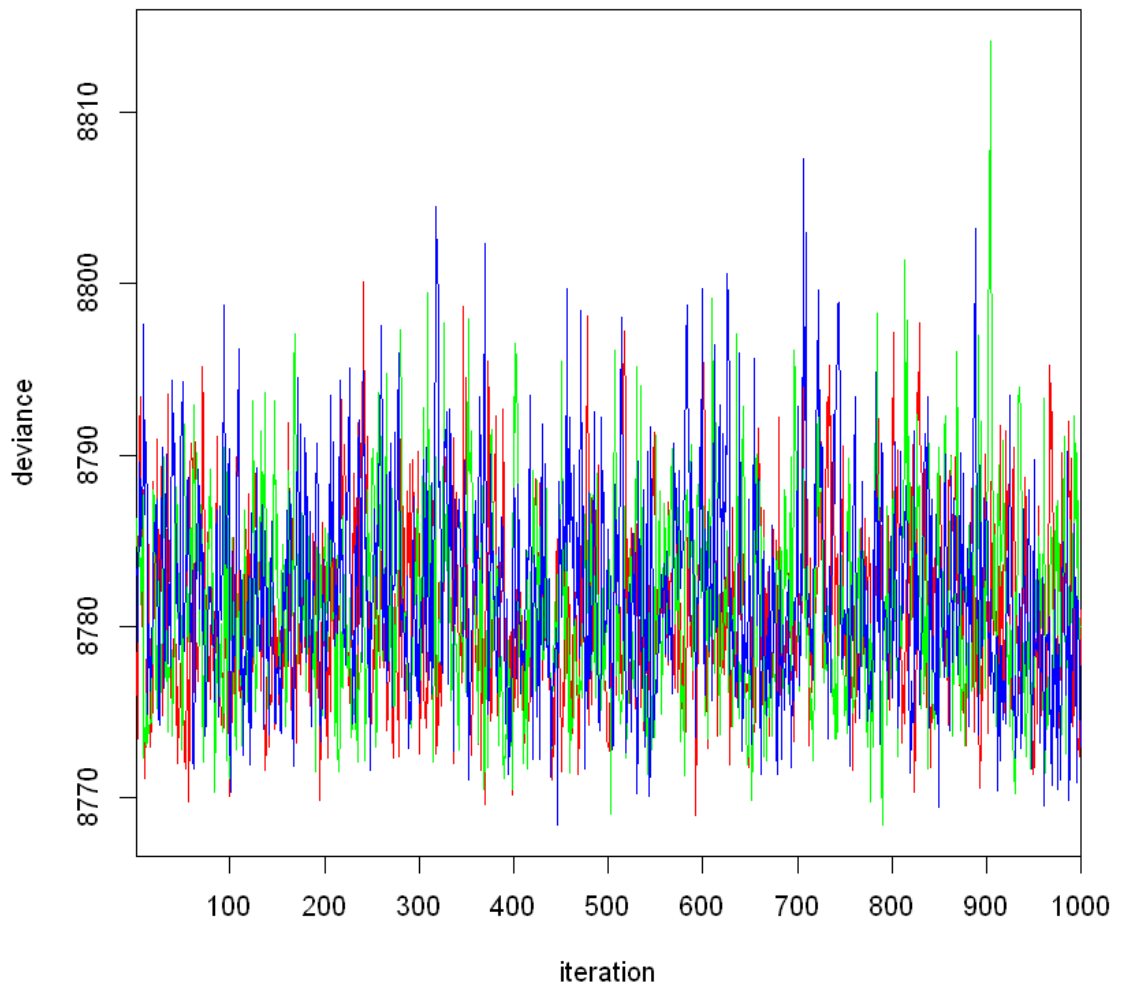
alpha[12]



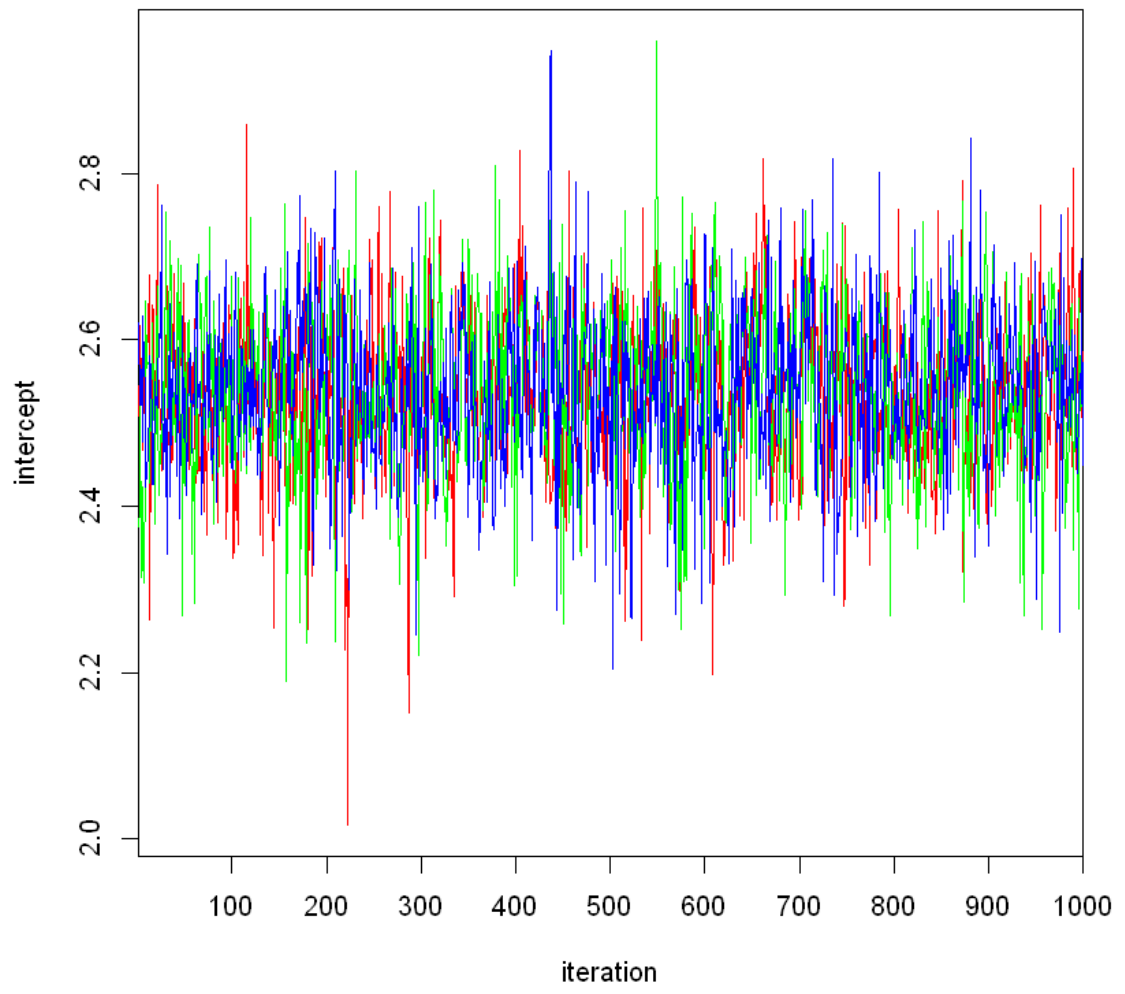
alpha[13]



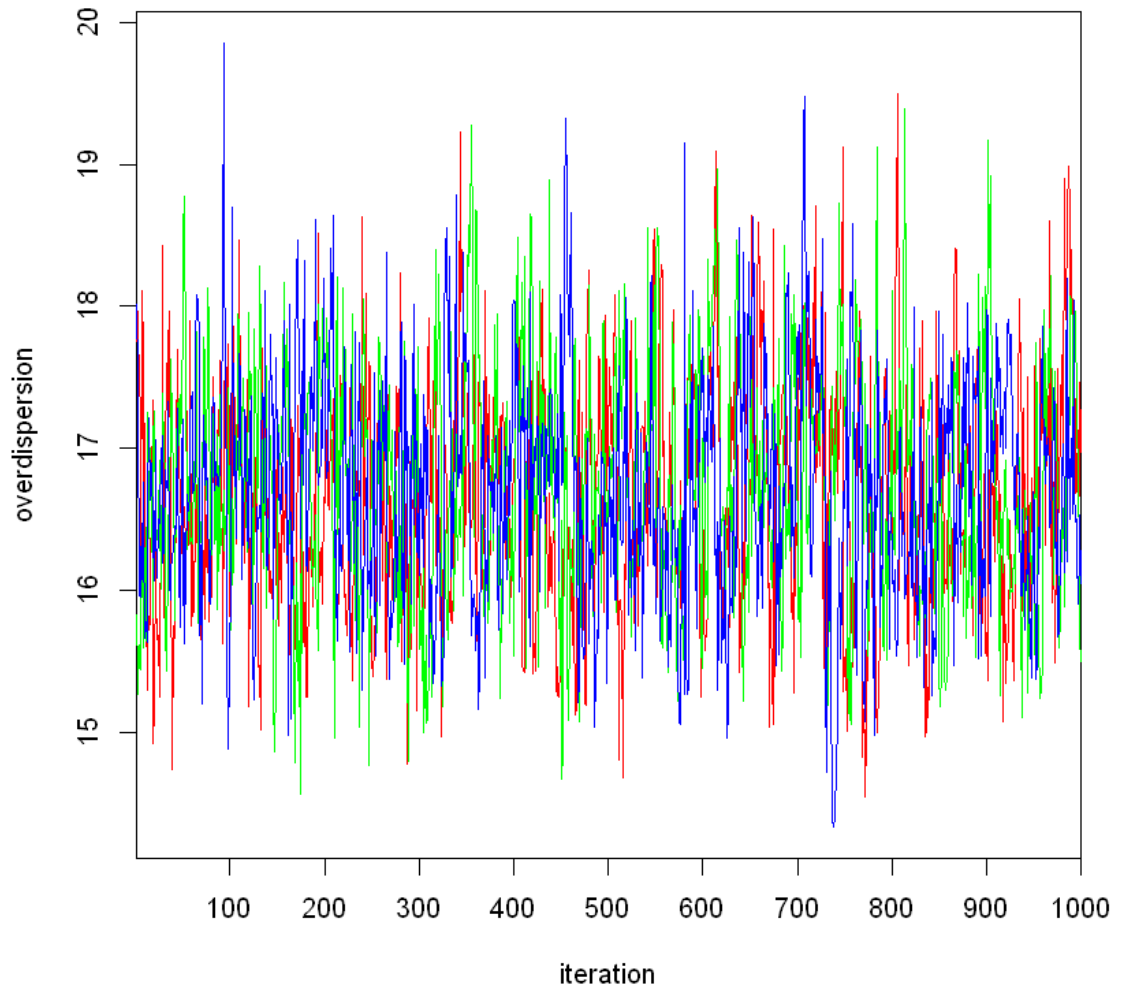
deviance



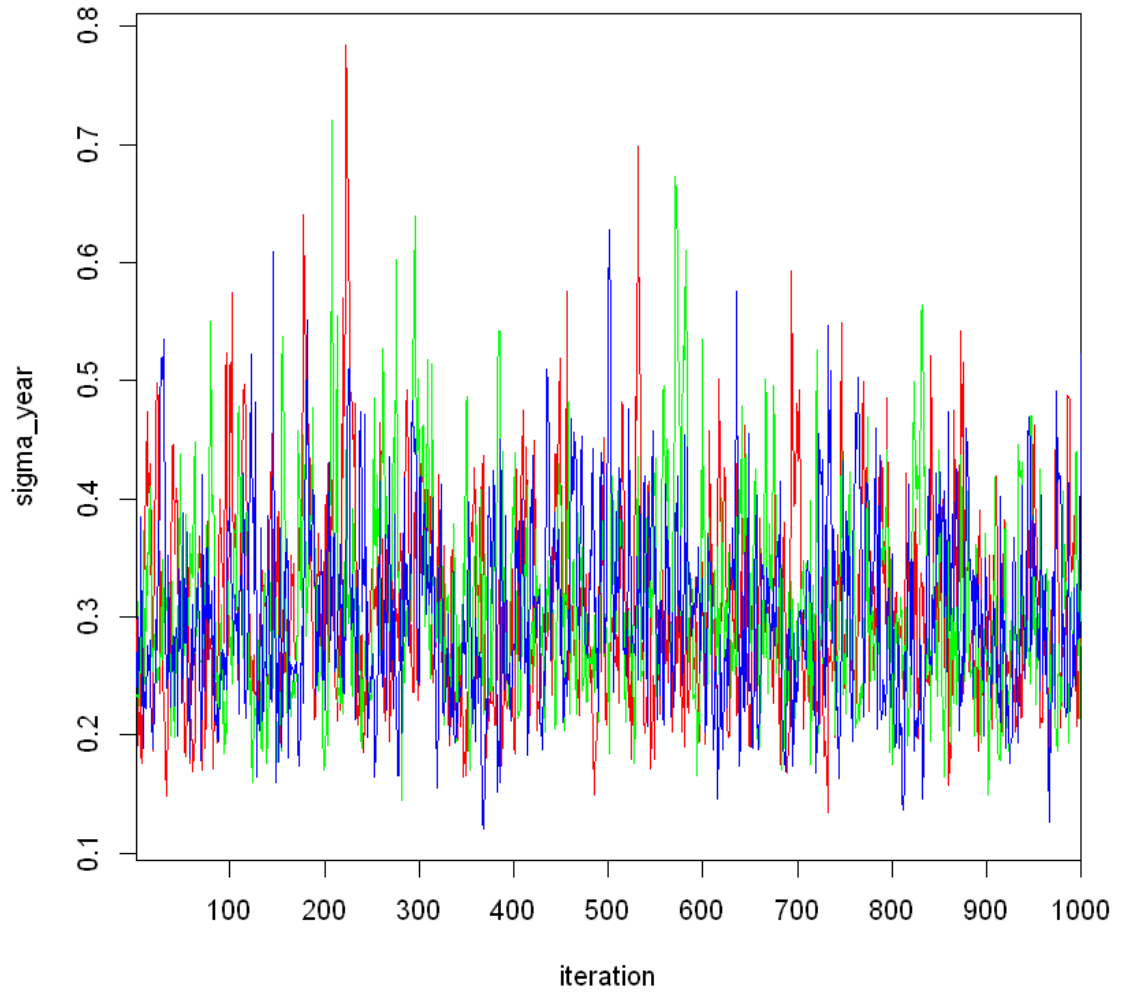
intercept

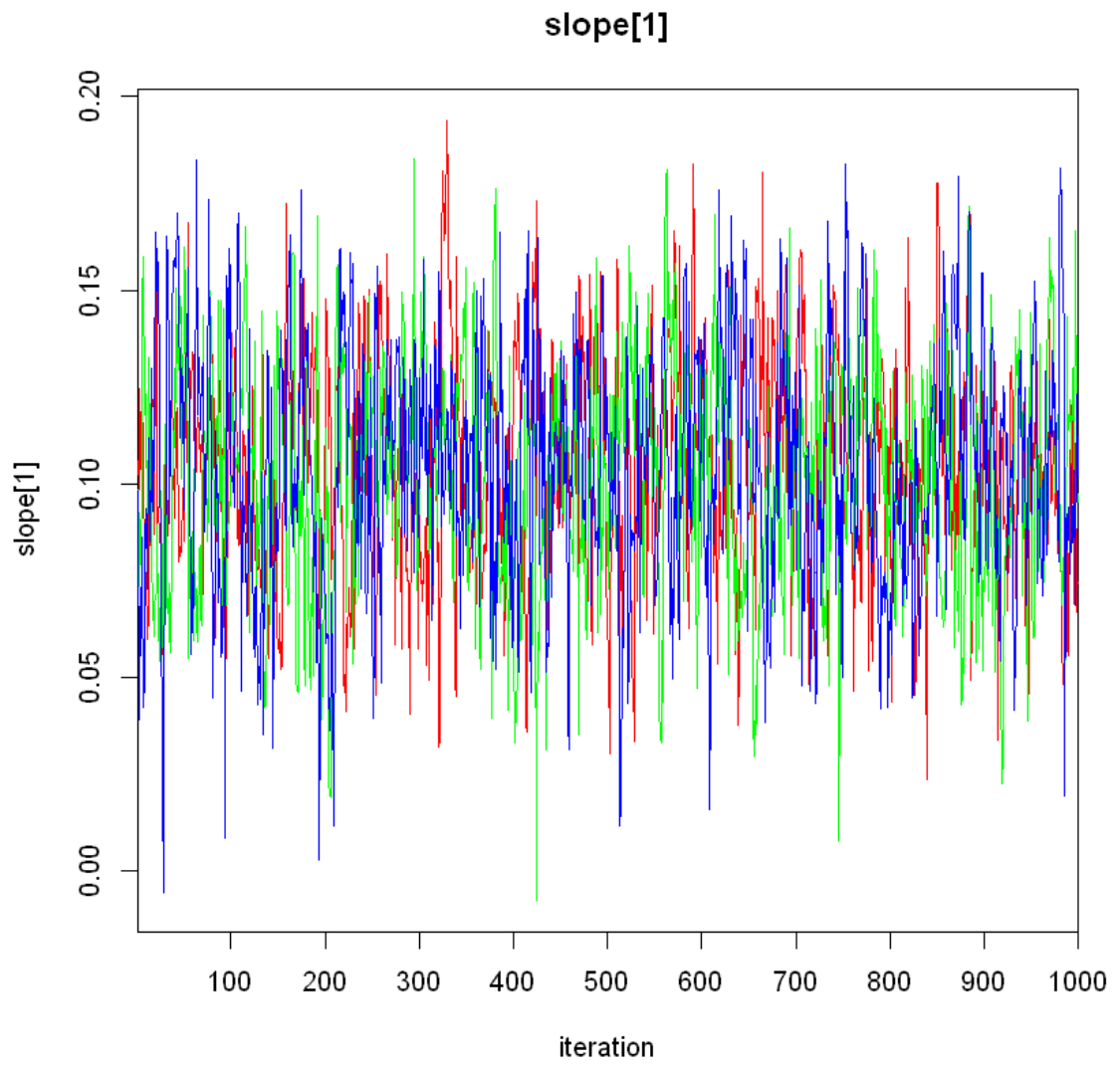


overdispersion

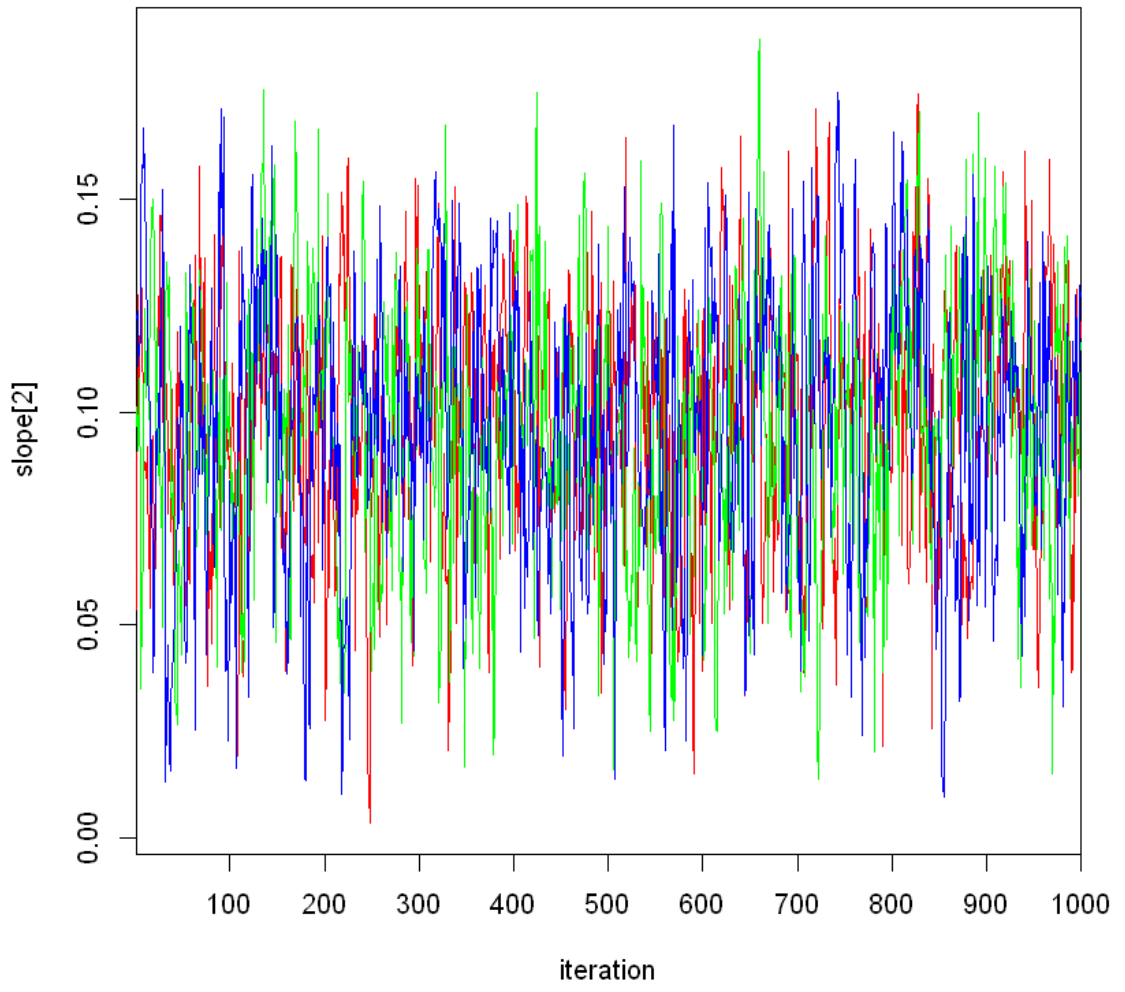


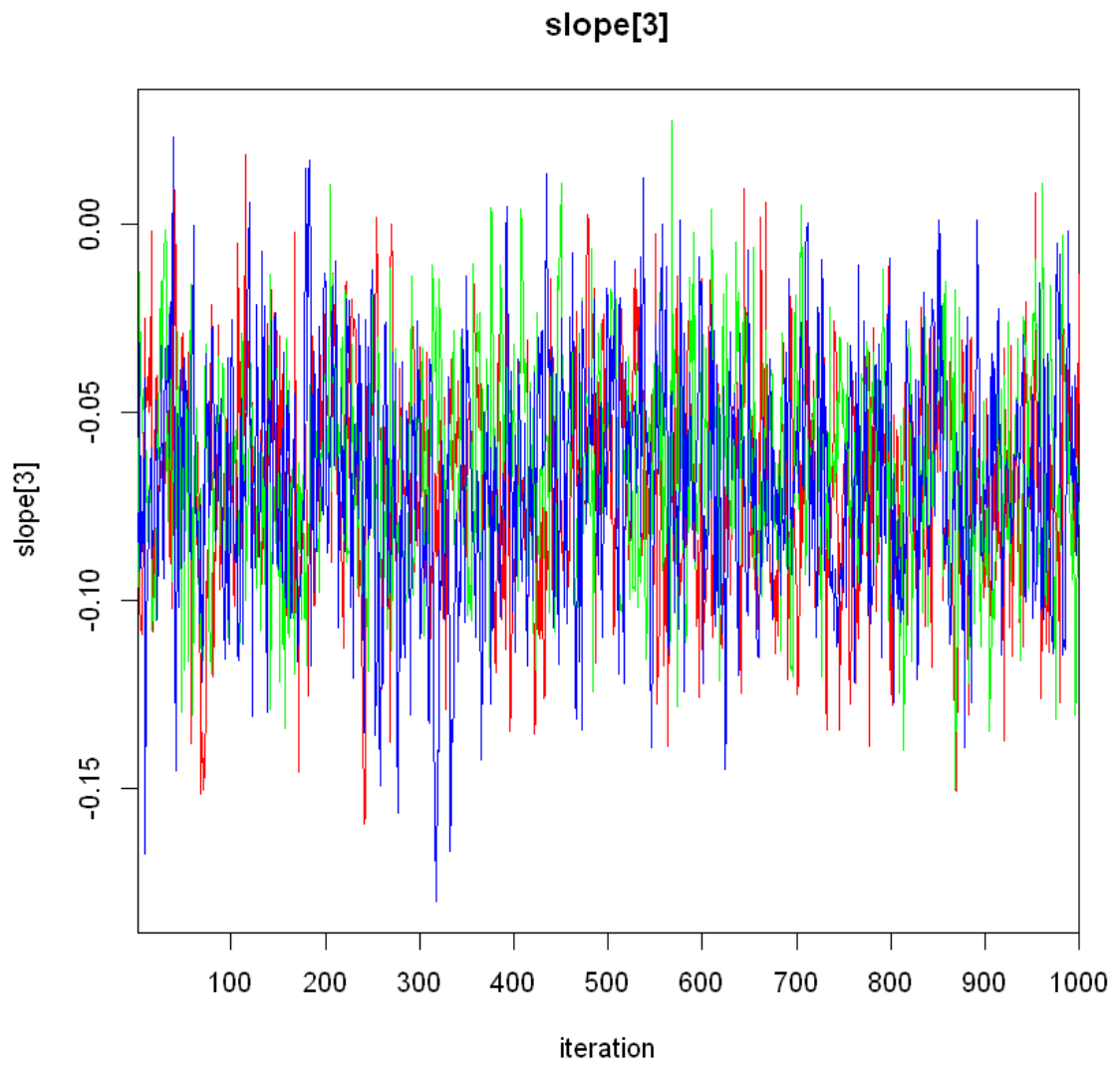
sigma_year





slope[2]





WAIC

On peut calculer le WAIC, du point de vue de ce critère le model est il meilleur que les précédents?

```
In [30]: ### calcul du WAIC
## besoin de calculer la log vraisemblance sous le modele
y_rep_4 <- log_lik4 <- array(NA, dim = c(nc * (ni - nb)/nt, nrow(obs)))

for(i in 1:nrow(obs)) {
  overdispersion <- drop(out4$sims.list$overdispersion)
  linpred <- out4$sims.list$slope[, 1] * data.jags$X1[i] +
    out4$sims.list$slope[, 2] * data.jags$X2[i] +
    out4$sims.list$slope[, 3] * data.jags$X3[i] +
    out4$sims.list$alpha[, data.jags$YEAR[i]]
  log_lik4[, i] <- dnbinom(data.jags$Y[i],
    prob = 1 / overdispersion,
    size = exp(linpred) / (overdispersion - 1),
    log = TRUE
  )
  y_rep_4[, i] <- rnbinom(nrow(y_rep_4),
    prob = 1 / overdispersion,
    size = exp(linpred) / (overdispersion - 1)
  ) + 1
}; rm(i, linpred, overdispersion)

loo::waic(log_lik4)
loo::loo(log_lik4)
```

Warning message:

"1 (0.1%) p_waic estimates greater than 0.4. We recommend trying loo instead."

Warning message:

"1 (0.1%) p_waic estimates greater than 0.4. We recommend trying loo instead."

Computed from 3000 by 1269 log-likelihood matrix

| | Estimate | SE |
|-----------|----------|-------|
| elpd_waic | -4398.3 | 56.6 |
| p_waic | 14.8 | 1.2 |
| waic | 8796.6 | 113.2 |

Warning message:

"Relative effective sample sizes ('r_eff' argument) not specified.
For models fit with MCMC, the reported PSIS effective sample sizes and
MCSE estimates will be over-optimistic."

Computed from 3000 by 1269 log-likelihood matrix

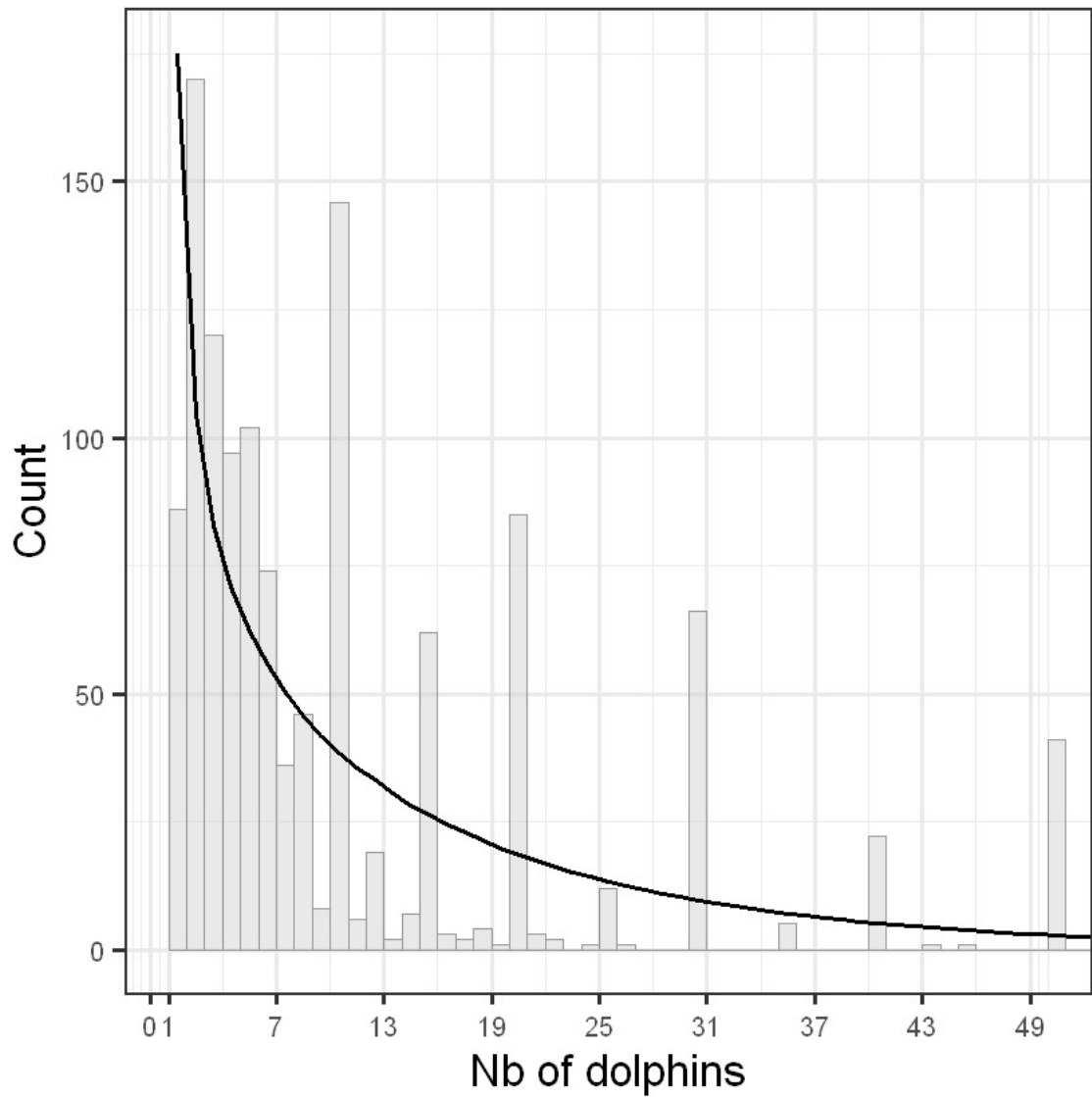
| | Estimate | SE |
|----------|----------|-------|
| elpd_loo | -4398.4 | 56.6 |
| p_loo | 14.9 | 1.2 |
| looic | 8796.7 | 113.2 |

Monte Carlo SE of elpd_loo is 0.1.

All Pareto k estimates are good (k < 0.5).
See help('pareto-k-diagnostic') for details.

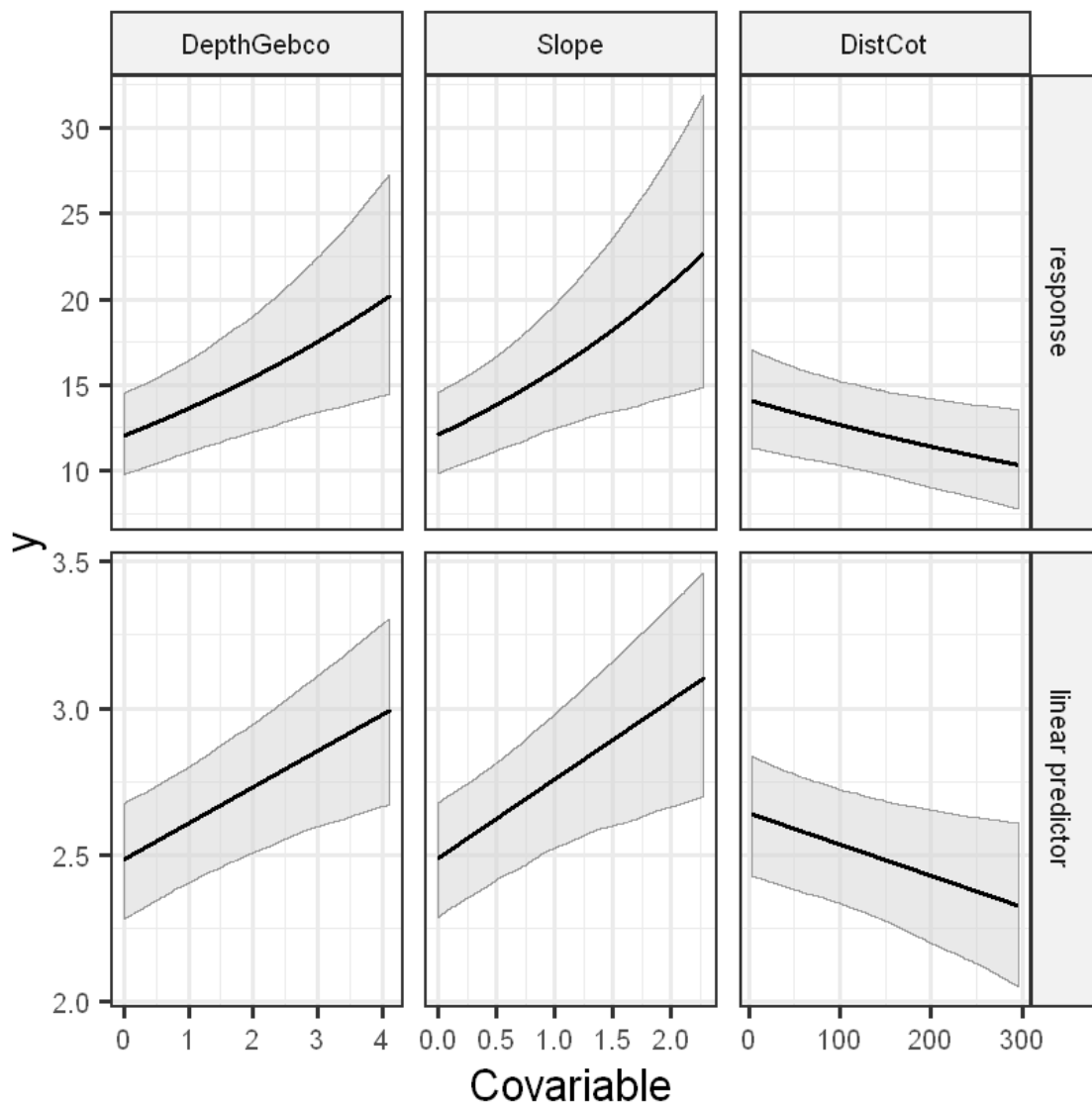
L'histogramme des donnees se compare-t-il à l'histogramme attendu sous le modele?

```
In [31]: ## regarder l'histogramme des donnees et le comparer a l'histogramme attendu sous le modele
rootogram_model_4 <- rootogram(countdata = obs$nombre, y_rep = y_rep_4)
rootogram_model_4 +
  coord_cartesian(xlim = c(1, 50))
```



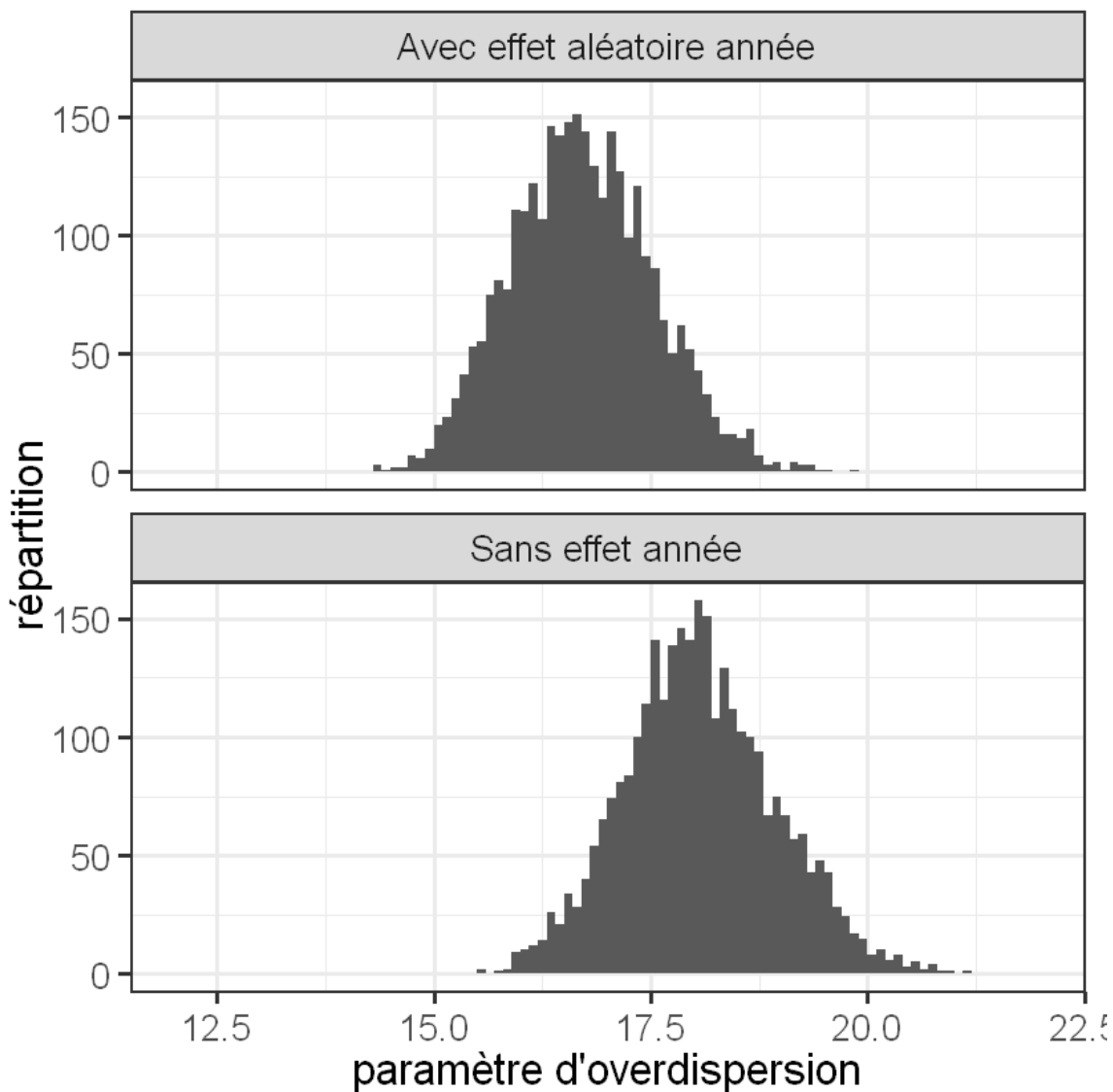
Un ajustement des effets avec des incertitudes plus réalistes

```
In [32]: ### un peu mieux que le modele precedent...  
plot_relationships(data_df = obs, jagsfit = out4, cov_name = c("DepthGebco", "Slope", "DistCot"))
```



On va regarder ce que ce modèle avec effets aléatoires change au niveau de l'estimation du paramètre de surdispersion:

```
In [33]: ggplot(data = data.frame(x = c(out3$sims.list$overdispersion,
                                         out4$sims.list$overdispersion
                                         ),
                                   distribution = rep(c("Sans effet année", "Avec effet aléatoire année"), each=length(out3$sims.list$overdispersion) )
                                   ),
               aes(x = x, group = distribution)
               ) +
  geom_histogram(breaks = seq(10, 25, 0.1)) +
  facet_wrap(~distribution, ncol = 1) +
  coord_cartesian(xlim = c(12, 22))+labs(x="paramètre d'overdispersion", y="répartition")
```



Une partie de la surdispersion est absorbée par l'effet année mais il reste une grande part qui n'est pas expliquée :

```
In [34]: 1 - round(mean(out4$sims.list$overdispersion /
                        out3$sims.list$overdispersion
                        ), 2
                    )
```

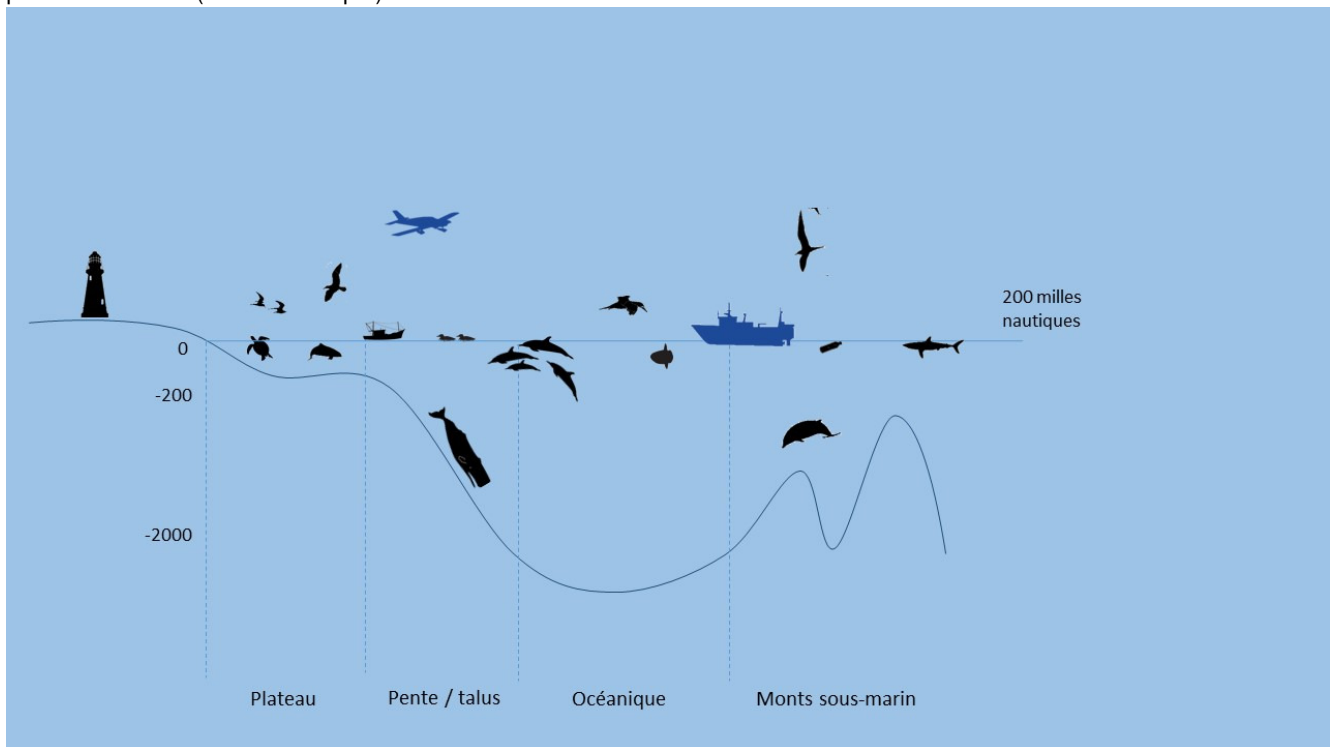
0.07

Une dernière option ici est de tenir compte d'un effet spatial avec un modèle de type Conditional Auto-Regressive (CAR). Ce modèle est facile à ajuster avec **Stan** (https://mc-stan.org/users/documentation/case-studies/icar_stan.html (https://mc-stan.org/users/documentation/case-studies/icar_stan.html)). Il faut cependant utiliser un nouveau bloc pour y définir une nouvelle distribution : la distribution CAR pour l'effet spatial. L'avantage de **Stan** ici est de pouvoir utiliser une paramétrisation qui va éviter des calculs coûteux. Ce modèle est ajustable dans un script supplémentaire mais ne sera pas couvert dans le TD.

On va donc ici considérer que notre meilleur modèle est le modèle surdispersé avec effet aléatoire de l'année . Les relations estimées avec l'environnement nous donnent la crédibilité de ces différentes variables environnementales sur la réponse conditionnellement à cette construction.

Conclusion

Nous avons ajuster une série de modèle à la complexité croissante pour estimer le nombre de dauphins dans le Golfe de Gascogne en fonction de diverses covariables environnementales. On peut voir ici qu'il y a plus de dauphins dans des zones de grandes profondeurs, loin des côtes et aussi des zones où la bathymétrie changent très rapidement sur de petites distances (talus océanique).



Néanmoins, le modèle sélectionné reste mauvais car il ignore un aspect important des données : les observateurs qui font les relevés sur le terrain ont tendance à arrondir les chiffres au delà d'une certaine abondance d'animaux, et cela induit des pics à des multiples de 5 ou 10. Il faudrait donc tenir compte du comportement des observateurs car celui-ci induit une source d'erreur de mesure qui n'est pas anodine, et qui contribue à générer de la surdispersion.

In []: