

# glm to MCMCglmm

*Kevin Healy*

*24 May 2017*

This is a short example of linear modeling starting with a glm and working up towards a MCMCglmm. However, we will mainly focus on the basics needed to run and check MCMCglmm models with some extra more complex models including adding random effects.

## Installation

First we need to install some packages including the lme4 to run mixed models using a frequentist approach and MCMCglmm to run a Bayesian approach to mixed models.

```
if(!require(MCMCglmm)) install.packages("MCMCglmm")
if(!require(lme4)) install.packages("lme4")
if(!require(MCMCpack)) install.packages("MCMCpack")
```

We will also install from GitHub the MulTree package which is still under development (so watch out for BUGS) but contains some handy data and also will allow us to use MCMCglmm and later to include the error associated with building phylogenies. To do so we need to get them from GitHub and so we need to run.

```
if(!require(devtools)) install.packages("devtools")
library(devtools)
install_github("TGuillerme/mulTree", ref = "master")
```

If you have not heard of GitHub it is a great way to share code, build packages and use version control to back up your work. For more check out here

Next we load up the packages we just installed from the library and we are good to go.

```
library(MCMCglmm)
library(lme4)
library(mulTree)
library(MCMCpack)
```

## Data

For the duration of the tutorials we will use some data that is part of a MulTree package. To get it just run

```
data(lifespan)
```

This data file contains a subset of the data used in an analysis on the role of flying (volant) in the evolution of maximum lifespan in birds and mammals Link to paper. Note that these data have been log transformed, mean centered and expressed in units of standard deviation. The original lifespan data were taken from the Anage database. We will come back to why it is often useful to transform data into z-scores later but for now we will simply assume our data is well behaved. Lets have a look at it now.

```
#data have been log transformed, mean centered and
#expressed in units of standard deviation.
head(lifespan_volant)
```

```
##           species    class longevity      mass    volant
## 1 Dolichotis_patagonum Mammalia -0.1490041  1.0875446 nonvolant
## 2      Eidolon_helvum Mammalia  0.4686111 -0.2748337    volant
## 3      Elephas_maximus Mammalia  2.1071286  3.1220340 nonvolant
## 4      Equus_asinus Mammalia  1.6128024  2.0352764 nonvolant
## 5      Equus_burchellii Mammalia  1.2962194  2.2295299 nonvolant
## 6      Equus_caballus Mammalia  1.9001076  2.2548716 nonvolant
```

## Lets run some models

Let's first start off running a simple glm for a subset of data for mammals

```
#subset for mammals
lifespan_mammals <- lifespan_volant[lifespan_volant$class == "Mammalia",]

#### and run a simple glm
glm_mod <- glm(formula = longevity ~ mass + volant, family = "gaussian", data = lifespan_mammals)
summary(glm_mod)
```

```
##
## Call:
## glm(formula = longevity ~ mass + volant, family = "gaussian",
##      data = lifespan_mammals)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.02569  -0.38641  -0.07613   0.41818   1.46075
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.02144    0.09241   0.232 0.816885
## mass          0.45655    0.05844   7.812 1.6e-12 ***
## volantvolant  0.95325    0.25568   3.728 0.000286 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 0.423457)
##
##      Null deviance: 81.317  on 133  degrees of freedom
## Residual deviance: 55.473  on 131  degrees of freedom
## AIC: 270.09
##
## Number of Fisher Scoring iterations: 2
```

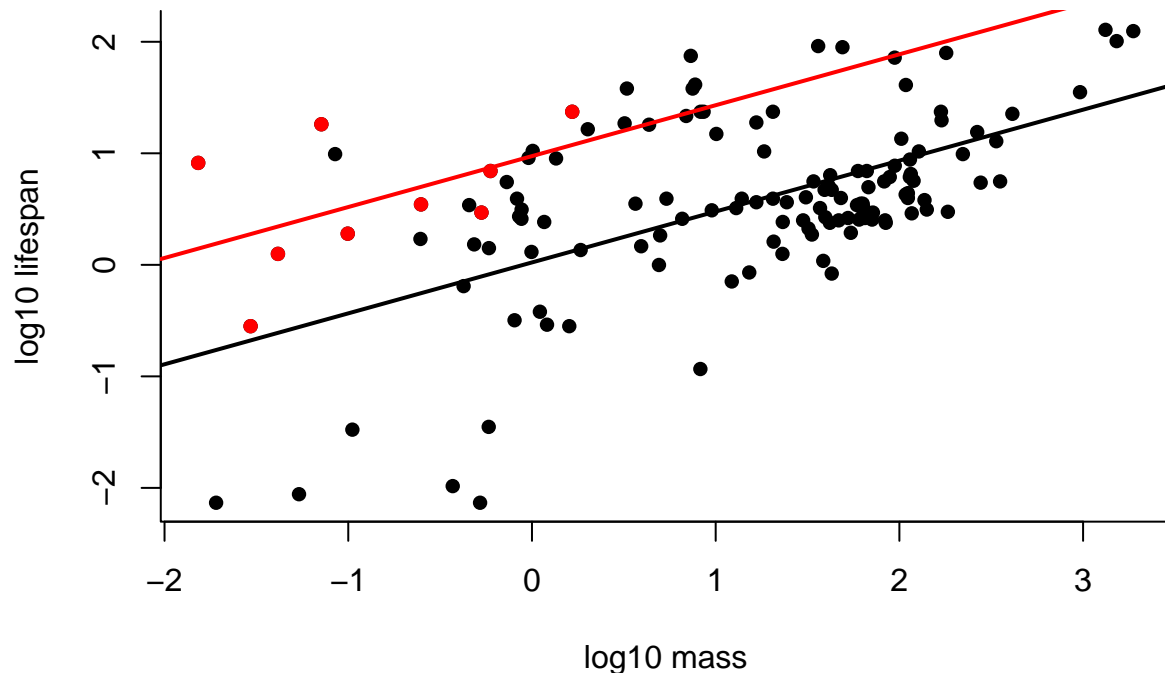
We can plot the results

```
#simple plots
plot(longevity ~ mass, data = lifespan_mammals, pch = 16, bty = "l",
      xlab = "log10 mass", ylab = "log10 lifespan")

#add in the volant species as red
points(lifespan_mammals[lifespan_mammals$volant == "volant", "longevity"] ~
       lifespan_mammals[lifespan_mammals$volant == "volant", "mass"],
       col = "red", pch = 16)

#add in the nonvolant regression line
```

```
abline(glm_mod, lwd = 2)
#add in the volant regression line
abline(glm_mod$coefficients[1] + glm_mod$coefficients[3],
       glm_mod$coefficients[2], lwd = 2, col = "red")
```



Most people will be familiar with `lm` and `glm` models so we won't spend any more time here. One thing that we might do is account for some of the structure in the error term. To do so we would need to include a random term.

## glmm

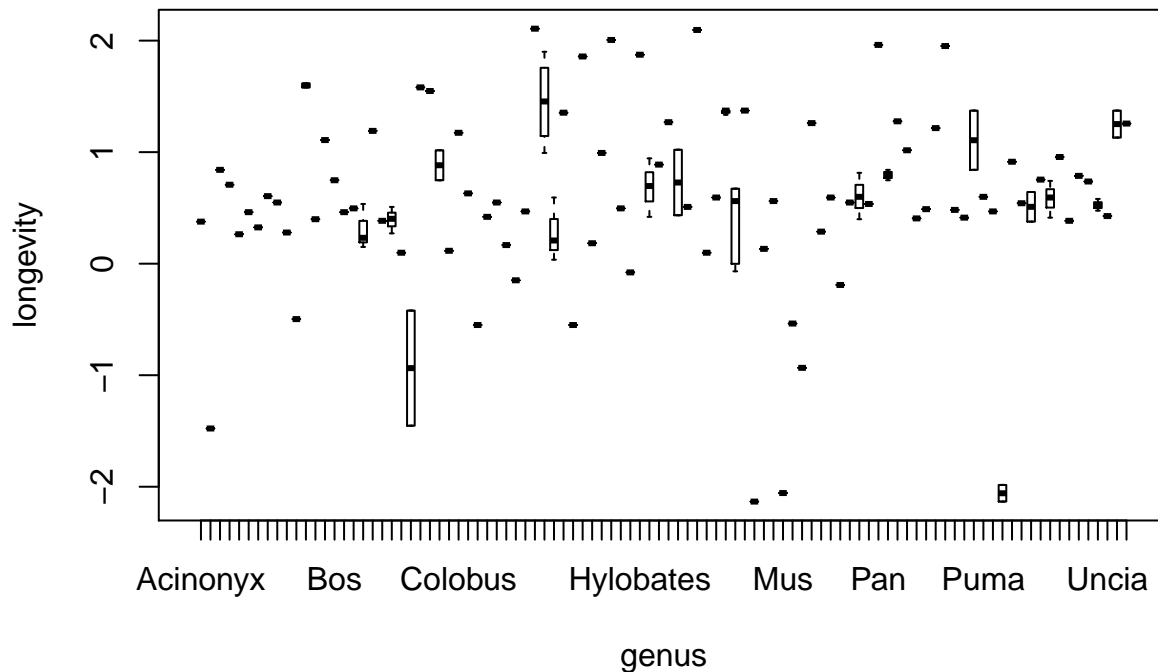
While our linear model looks good we might also want to try to control for the structure of our data. For example, maybe all the data points are not fully independent with some data points more likely to have values closer to other ones. In this case we might want to add a random term to control for this. We can imagine such a case in our data above where two species from the same genus might show more similar values of lifespan in comparison to other species.

Let's plot it out and have a look

```
#use the gsub function to create a vector of the genus for each species
genus <- (gsub("_.*", "", lifespan_mammals$species))

#bind it back into the database
lifespan_mammals <- data.frame(lifespan_mammals, genus = genus)

#plot out lifespan ~ genus to get an idea of whether genus is structured
#randomly
plot(longevity ~ genus, data = lifespan_mammals)
```



Genus is something I would like to control for, however I am not really interested in which groups are different and for practical reasons I don't want to fit every single group as a fixed factor. Hence I could include it as a random term using a lmer model which is used to fit mixed models using a maximum likelihood approach.

```
# Lets fit our model. For lmer the random term is fitted using (1|genus)
# to indicate that you want to fit separate intercepts
# to each of group in your random term.
lmer_mod <- lmer(longevity ~ mass + volant + (1|genus), data = lifespan_mammals)
summary(lmer_mod)
```

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: longevity ~ mass + volant + (1 | genus)
## Data: lifespan_mammals
##
## REML criterion at convergence: 220.9
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -2.0997 -0.2667 -0.0240  0.3312  1.8262
##
## Random effects:
## Groups   Name            Variance Std.Dev.
## genus    (Intercept)  0.37772   0.6146
## Residual                0.05953   0.2440
## Number of obs: 134, groups: genus, 98
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept)  -0.01585    0.10366  -0.153
## mass          0.48637    0.06365   7.641
## volantvolant  1.00386    0.27950   3.592
##
## Correlation of Fixed Effects:
```

```
##          (Intr) mass
## mass      -0.746
## volantvolnt -0.534  0.496
```

Like before we see estimates for each of the fixed effects which look similar to our glm model. The next section that we are interested in is the random effects where we see the terms genus and Residual. Our residual term is acting as in the glm model telling how much of the variance is unexplained while the genus term tells us how much of the variance is due to variance associated with the genus groupings. We can see here for example that genus accounts for more variation than our random term after accounting for our fixed effects.

I will leave lmer models here, however if you are interested in more on mixed effects models using this approach check out [this tutorial](#). For now let's move on to the main event with MCMCglmm.

## MCMCglmm

So far we have fitted a very simple glm and a glmm model, now we will fit a linear model and a mixed model using the MCMCglmm package. We will start first with a model similar to our glm.

First things first, since we are using a Bayesian approach we will need to set up the priors. In most cases we want to use a non-informative prior that doesn't influence the estimated posterior distribution. We are basically saying that we don't know anything about the expected values for our parameters. That is we have no prior information of what the intercepts and slopes will be.

To give priors for MCMCglmm we need to make an object that is in a list format that includes terms of B (fixed effect), R (residual terms) and G (random terms which we will come to later).

For now let's build a prior with just a fixed term and a residual term.

```
prior <- list(B = list(mu= diag(3)*0, V=diag(3)*1e+10),
              R = list(nu=0.002, V=1))
```

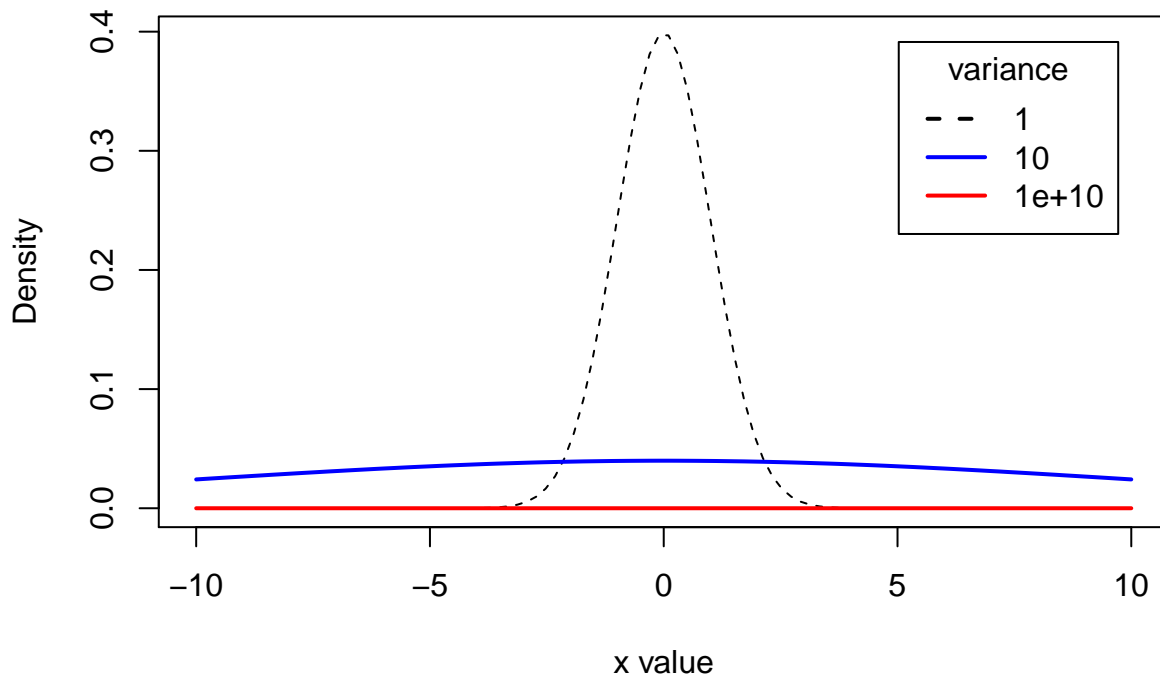
For fixed effects (B) the terms V and mu give the variance and mean of a normal distribution. Here we set mu as 0 and the variance as a large number to make these priors uninformative. Since we have three fixed terms (two intercepts and one slope) we can use the diag function to create a matrix to store a prior for each. Normally we don't need to set this as MCMCglmm will set non-informative priors automatically for fixed terms.

If we plot it out, we can see that higher values for V give less informative priors.

```
#create some normal distributions over some range x
x <- seq(-10, 10, length=100)
#V = 1
hx_1 <- dnorm(x, 0,1 )
#V = 10
hx_10 <- dnorm(x, 0,10 )
#V = 1e+10
hx_1e10 <- dnorm(x, 0,1e+10 )
plot(x, hx_1, type="l", lty=2, xlab="x value",
     ylab="Density", main="fixed effects prior")

lines(x, hx_10, lwd=2, col="blue")
lines(x, hx_1e10, lwd=2, col="red")
labels <- c("1","10","1e+10")
legend("topright", inset=.05, title="variance",labels,
      lwd=2, lty=c(2, 1, 1), col=c("black", "blue", "red"))
```

## fixed effects prior

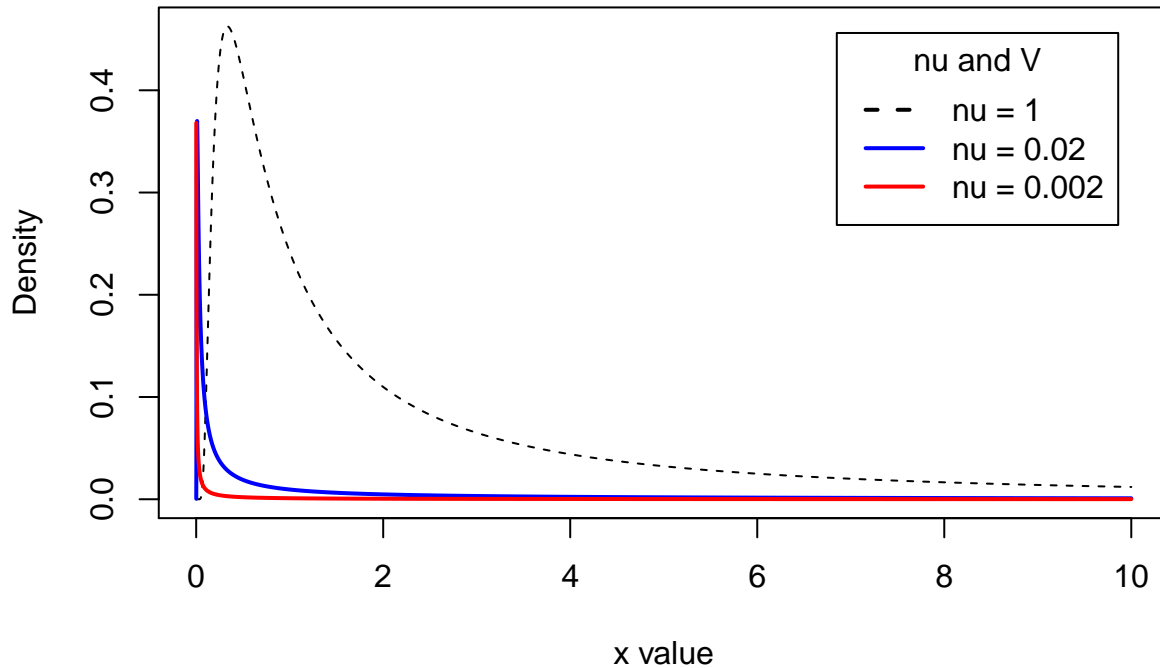


For the variance terms we need to make sure that the distribution is bounded at zero as the variance term needs to be positive. To do this MCMCglmm uses an inverse-Gamma distribution. In MCMCglmm this is described by two parameters  $\nu$  and  $V$ . These terms are related to the shape (alpha) and scale (beta) parameters on an inverse-Gamma with  $\alpha = \nu/2$ , and  $\beta = (\nu \cdot V)/2$ . As we don't want our estimates to be heavily influenced by our prior we will use weakly informative prior values such as described as  $V = 1$  and  $\nu = 0.002$ . (For more on priors the see course notes)

```
#create some normal distributions over some range x
x <- seq(0, 10, length=10000)
#nu = 0.002 V = 1
hx_0.002 <- dinvgamma(x, 0.002/2, (0.002*1)/2)
#nu = 0.02 V = 1
hx_0.02 <- dinvgamma(x, 0.02/2, (0.02*1)/2)
#nu = 0.002 V = 2
hx_0.002_2 <- dinvgamma(x, 1/2, (1*1)/2)
plot(x, hx_0.002_2, type="l", lty=2, xlab="x value",
     ylab="Density", main="fixed effects prior")

lines(x, hx_0.02, lwd=2, col="blue")
lines(x, hx_0.002, lwd=2, col="red")
labels <- c("nu = 1", "nu = 0.02", "nu = 0.002")
legend("topright", inset=.05, title="nu and V", labels,
     lwd=2, lty=c(2, 1, 1), col=c("black", "blue", "red"))
```

## fixed effects prior



Next we need to decide on the parameters relating to running the mcmc chain in the model. We need to include how many iterations we want to run the MCMC chain for (`nitt`), the burnin we want to discard at the start of the chain (`burnin`) and also how often we want to sample and store from the chain (`thin`). We discard a burnin as we don't want the starting point of the chain to over-influence our final estimates. For now let's just use a burnin of 1/6 of the `nitt`, just to be safe. The thinning is used to help reduce autocorrelation in our sample, how much you use often depends on how much autocorrelation you find.

To save time we will only run this model over 12000 iterations (However, much larger `nitt` is often required).

```
#no. of iterations
nitt <- c(12000)
#length of burnin
burnin <- c(2000)
#amount of thinning
thin <- c(5)
```

Now we can run the model using our data. Let's run a model similar to our first glm

```
mod_mcmc_fix <- MCMCglmm(fixed = longevity ~ mass + volant,
  family="gaussian",
  data = lifespan_mammals,
  nitt = nitt,
  burnin = burnin,
  thin = thin,
  prior = prior)
```

```
##
##                               MCMC iteration = 0
##
##                               MCMC iteration = 1000
##
##                               MCMC iteration = 2000
```

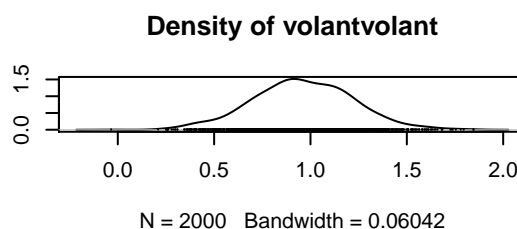
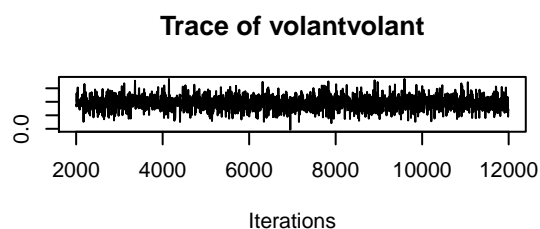
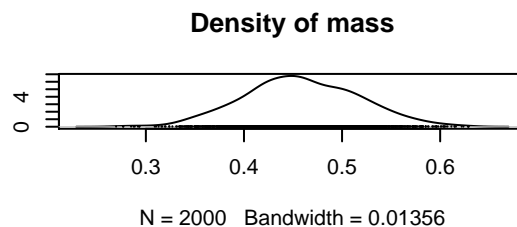
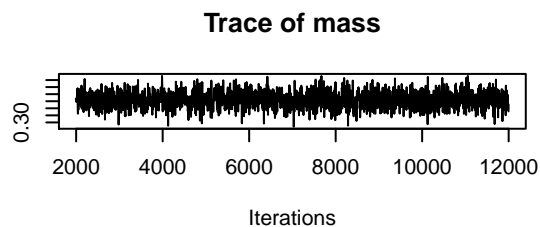
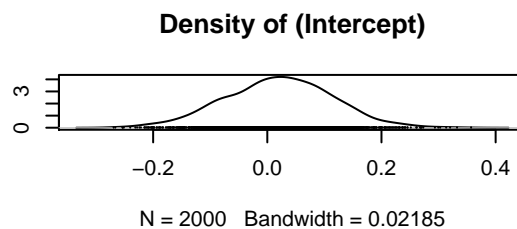
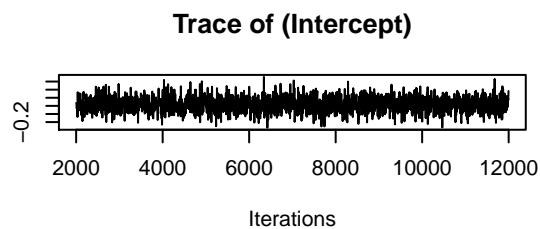
```
##
##           MCMC iteration = 3000
##
##           MCMC iteration = 4000
##
##           MCMC iteration = 5000
##
##           MCMC iteration = 6000
##
##           MCMC iteration = 7000
##
##           MCMC iteration = 8000
##
##           MCMC iteration = 9000
##
##           MCMC iteration = 10000
##
##           MCMC iteration = 11000
##
##           MCMC iteration = 12000
```

As the model runs we see the iterations print out. These chains can take some time to run, depending on the model, however, since we only ran our chains for 12000 iterations it doesn't take long here.

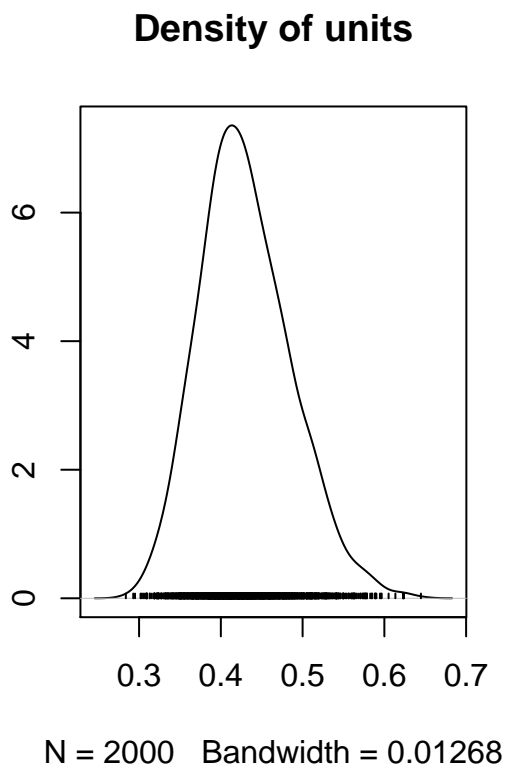
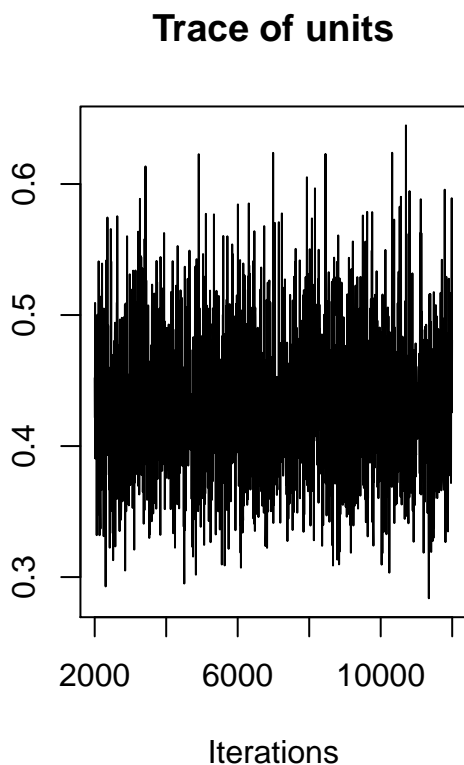
Before we even look at our model we need to check if the model ran appropriately. We can do this by visually inspecting the chains to make sure there has been no unruly behavior! We can extract the full chains using `model$Sol` for the fixed effects and `model$VCV` for the variance terms. So `Sol[,1]` will give you the first fixed term, in this case the intercept, and `VCV[,1]` will give you the first random term, which is just the residual term here. As our model is an mcmc object when we use the plot function we get a trace plot.

```
#plot the first fixed term, the intercept.
plot(mod_mcmc_fix$Sol)
```





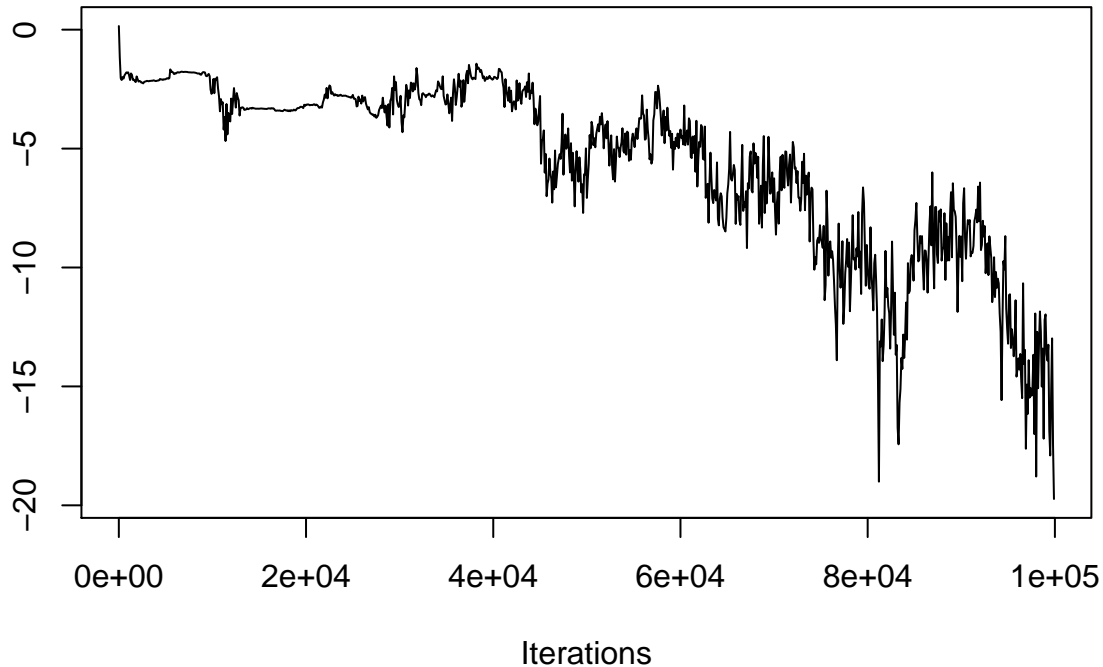
*#plot the first variance term, the residual error term.*  
`plot(mod_mcmc_fix$VCV)`



On the right hand side of the plots is the posterior distributions for each of the terms. On the left side of

these plots are the traces of the mcmc chain for each estimate. What we want to see in these trace plots is “hairy caterpillars” (not my phrase!). That is a trace with no obvious trend that is bouncing around some stable point.

What we don’t want to see in the trace plots can be demonstrated if we only run a model over a very short chain (itt == 10000) or more difficult model fit (we will see these more difficult models later). Notice in the example below that without a burnin the start of trace is well outside the area that the chain is converging towards.



So far in our simple mod\_mcmc\_fix model everything looks good visually, however we also want to check the level of auto-correlation in these traces. We can do this using autocorr.diag() which gives the level of correlation along the chain between some lag sizes.

```
autocorr.diag(mod_mcmc_fix$Sol)
```

```
##           (Intercept)           mass volantvolant
## Lag 0      1.000000000    1.000000000    1.000000000
## Lag 5     -0.004186761    0.0004268486   0.008269119
## Lag 25    -0.011355300    0.0173513524  -0.013247957
## Lag 50    -0.013273705   -0.0231914532  -0.013357976
## Lag 250   -0.043002922    0.0036305563  -0.034266359
```

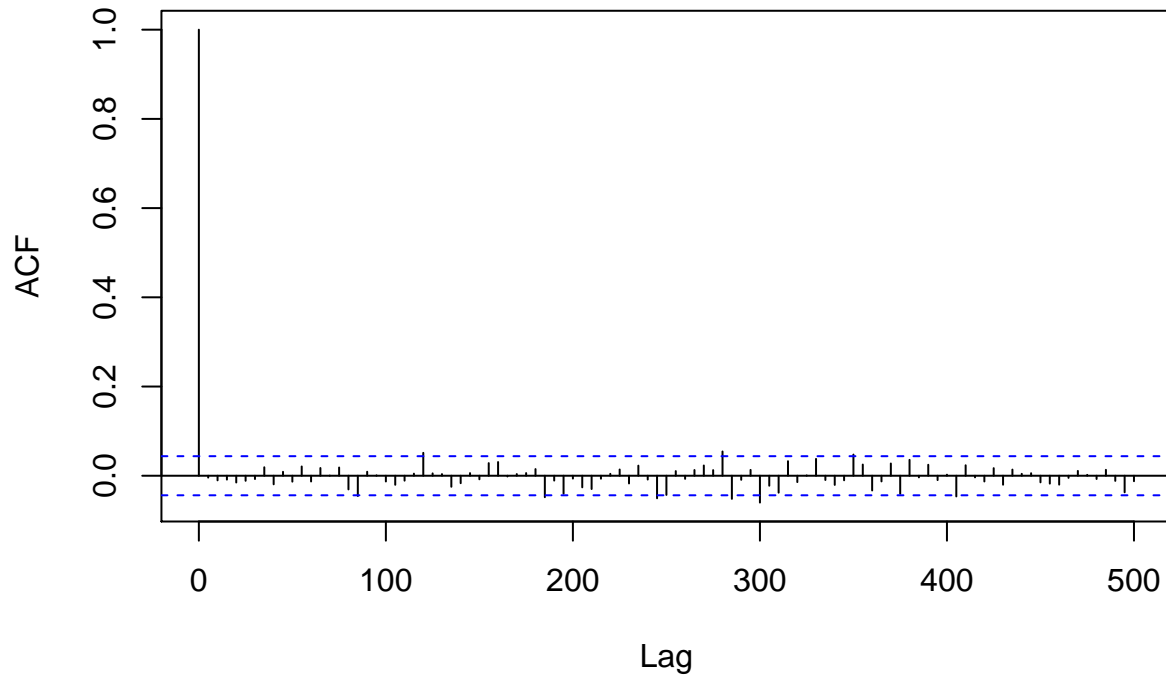
```
autocorr.diag(mod_mcmc_fix$VCV)
```

```
##           units
## Lag 0      1.00000000
## Lag 5     -0.02387519
## Lag 25     0.04774305
## Lag 50    -0.03675542
## Lag 250    0.06677416
```

or we can look at autocorrelation plots for each of the traces. For example, let’s check the auto-correlation in the intercept chain using the acf function

```
#acf plot for the first fixed estimate in our model (the intercept)
acf(mod_mcmc_fix$Sol[,1], lag.max =100)
```

### Series mod\_mcmc\_fix\$Sol[, 1]



For our intercept the auto-correlation plot looks good. However if there is bad auto-correlation one quick way to deal with this is to simply increase the thinning. While we don't need to in our case an example would be running something like

```
nitt2 <- 240000
burnin2 = 40000
thin2 = 100
mod_mcmc_long <- MCMCglmm(fixed = longevity ~ mass + volant,
                          family="gaussian",
                          data = lifespan_mammals,
                          nitt = nitt2,
                          burnin = burnin2,
                          thin = thin2,
                          prior = prior,
                          verbose=FALSE)
```

Noticed I also increased the number of iterations. One rough and ready rule that I like to use is to aim for an effective sample size of somewhere between 1000-2000 for all my estimates. The effective sample size is the number of samples in the posterior after the burnin, thinning and autocorrelation are accounted for.

```
#acf plot for the first fixed estimate in our model (the intercept)
effectiveSize(mod_mcmc_long$Sol)
```

```
## (Intercept)      mass volantvolant
##      2000      2000      2000
```

```
effectiveSize(mod_mcmc_long$VCV)
```

```
## units
## 1678.764
```

\*One thing to note is that while thinning might help autocorrelation it won't solve it and you might have

to use parameter expanded priors. These are priors that help weight the chain away from zero, a common problem when variance is low or with certain phylogenetic structures. They work by splitting the prior into 2 components with one component weighing the chain away from zero. We will come back to such a prior shortly.

One last thing to check is that our MCMC chain has properly converged and that our estimate is not the result of some type of transitional behaviour. That is have our chains “found” the optimum or do we need to let them run longer before they settle around some estimate. To check this we will run a second model and see if it converges on the same estimates as our first model.

```
mod_mcmc_2 <- MCMCglmm(fixed = longevity ~ mass + volant,
                      family="gaussian",
                      data = lifespan_mammals,
                      nitt = nitt2,
                      burnin = burnin2,
                      thin = thin2,
                      prior = prior,
                      verbose=FALSE)
```

We can now check the convergence of the two chains using the Gelman and Rubin Multiple Sequence Diagnostic. This calculates the within-chain and between-chain variance of the chains and then gives a scale reduced factor, (for more see [here](#). When this number is close to one (something below 1.1 is usually good) the chains are indistinguishable and hence can be considered to be converged.

```
#checking convergence for our fixed factors
gelman.diag(mcmc.list(mod_mcmc_long$Sol, mod_mcmc_2$Sol))
```

```
## Potential scale reduction factors:
##
##               Point est. Upper C.I.
## (Intercept)           1           1
## mass                  1           1
## volantvolant          1           1
##
## Multivariate psrf
##
## 1
```

```
#checking convergence for our random terms
gelman.diag(mcmc.list(mod_mcmc_long$VCV, mod_mcmc_2$VCV))
```

```
## Potential scale reduction factors:
##
##               Point est. Upper C.I.
## units              1           1.02
```

Since everything looks good, we will finally look at the results of our model.

```
summary(mod_mcmc_long)
```

```
##
## Iterations = 40001:239901
## Thinning interval = 100
## Sample size = 2000
##
## DIC: 270.1205
##
## R-structure: ~units
```

```
##
##      post.mean l-95% CI u-95% CI eff.samp
## units    0.4287  0.3305    0.537    1679
##
## Location effects: longevity ~ mass + volant
##
##      post.mean l-95% CI u-95% CI eff.samp pMCMC
## (Intercept)   0.02416 -0.15679  0.20066    2000  0.808
## mass          0.45561  0.34378  0.57506    2000 <5e-04 ***
## volantvolant  0.95038  0.44186  1.45539    2000 <5e-04 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

First off we can find the estimates for the fixed factors are under the Location effects section (Notice the similarity to our glm and nlme model). Each parameter has a measure of the effect size under post.mean and a lower and higher 95% credible interval (CI). These are simply calculated from the posterior distributions we looked at in the above plots, so if you would rather calculated the median instead of using the mean we can simple use

```
median(mod_mcmc_long$Sol[,1])
```

```
## [1] 0.02297672
```

We also have the effective sample size (`eff.samp`) and the `pMCMC` which calculated as two times the probability that the estimate is either  $>$  or  $< 0$ , using which ever one is smaller. However, since our data has been mean centred and expressed in units of standard deviation we can simply look at what proportion of our posterior is on either side of zero. This mean centering and expressing our data in units of standard deviation hence allows us to use a cut off point like a p-value but without boiling down the whole distribution to one value.

We also have the DIC which is a Bayesian version of AIC. Like AIC it is a measure of the trade-off between the “fit” of the model and the number of parameters, with a lower number better.

## Putting the m in MCMCglmm

Since we can run a simple linear MCMCglmm lets add a random term of genus in the example above. Like before we need to set up the prior however we will let the model estimate the fixed effects this time. To add a random term we now add a `G structure` that acts just like the other random variance terms and is defined using `nu` and `V`.

```
prior <- list(G = list(G1 = list(nu=0.002, V=1)),
             R = list(nu=0.002, V=1))
```

We can now include the random term in the model in the section `random= ~.`

```
nitt_m <- 1000
burnin_m <- 1
thin_m <- 1

mod_mcmc_mixed <- MCMCglmm(fixed = longevity ~ mass + volant,
                          rcov=~ units,
                          random= ~ genus,
                          family="gaussian",
                          data = lifespan_mammals,
                          nitt = nitt_m,
                          burnin = burnin_m,
                          thin = thin_m,
                          prior = prior,
```

```

                                verbose=FALSE)
summary(mod_mcmc_mixed)

##
## Iterations = 2:1000
## Thinning interval = 1
## Sample size = 999
##
## DIC: 95.83104
##
## G-structure: ~genus
##
##      post.mean l-95% CI u-95% CI eff.samp
## genus      0.374    0.259    0.509    308.6
##
## R-structure: ~units
##
##      post.mean l-95% CI u-95% CI eff.samp
## units      0.0667  0.03654   0.1025    131.8
##
## Location effects: longevity ~ mass + volant
##
##      post.mean l-95% CI u-95% CI eff.samp pMCMC
## (Intercept)  -0.01674 -0.21382  0.17345    999.0  0.847
## mass          0.48669  0.35980  0.59756    999.0 <0.001 **
## volantvolant  1.00390  0.41189  1.52475    803.6 <0.001 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

As an exercise check the output of this model.

Once you are happy with the above model you have essentially got the basics of running a Bayesian mixed effects model. Building on the complexity involves adding more terms which are covered in good detail in the course notes and vignette.

## Phylogenetic models

We have already ran a model with genus as a random term in order to attempt to deal with pseudoreplication. However, life is not ordered in catagories but as a continuous branching tree meaning we are potentially throwing away a lot of useful information. To include the full phylogeny we will use both a Phylogenetic generalized linear model and a MCMCglmm animal model to run the Bayesian version of the phylogenetic comparative analysis.

First we will need ape and caper, two packages that are very useful for phylogenetic analysis.

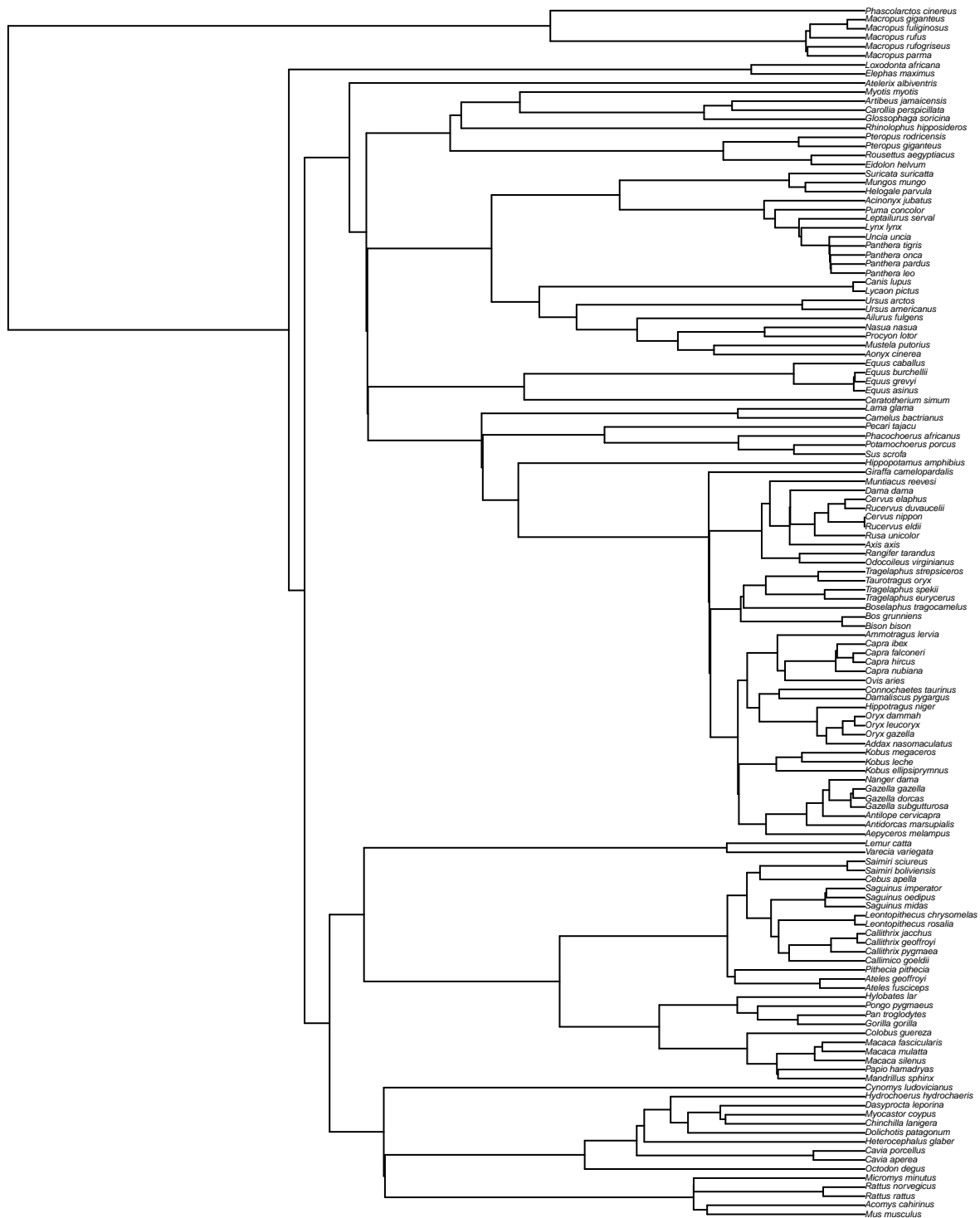
```

if(!require(ape)) install.packages("ape")
if(!require(caper)) install.packages("caper")
library(ape)
library(caper)

```

We will use a phylogeny of mammals constructed in Kuhn et al 2011, were they produce 10,000 trees with each individual tree comprising one resolution of the polytomies of a previously published supertree. For now we will just use the first tree in the list.

```
# The 10Ktrees from Khun et al (2011) gives a set of trees to represent different polytomies.
# For now let just take one.
mammal_tree <- trees_mammalia[[1]]
plot(mammal_tree, cex = 0.3)
```



```
#the number of species
Ntip(mammal_tree)
```

```
## [1] 134
#we can also check that its ultrametric
is.ultrametric(mammal_tree)
```

```
## [1] TRUE
```

## PGLS

In phylogenetic comparative methods we try to deal with this non-independence by using the structure of a phylogeny to weight the error term so that our model fit is no longer blind to the non-independence of our data.

The first step is to make sure that our data and phylogeny match up correctly. We will use the `comparative.data` function from the `caper` package to make sure that the species in the dataset are matched up to the species in the phylogeny. To do so we need to give the phylogeny (`phy`); the dataset and the name of the column in the dataset with the species names (`names.col`). As we want to calculate a variance-covariance matrix of the phylogeny we set `vcv = true`.

```
comp_data <- comparative.data(phy = mammal_tree,
                             data = lifespan_volant,
                             names.col = species,
                             vcv=TRUE)

head(comp_data$data)

##           class longevity      mass    volant
## Mus_musculus   Mammalia -2.056791 -1.2680281 nonvolant
## Acomys_cahirinus Mammalia -1.477934 -0.9776073 nonvolant
## Rattus_rattus   Mammalia -1.984124 -0.4313382 nonvolant
## Rattus_norvegicus Mammalia -2.133185 -0.2824073 nonvolant
## Micromys_minutus Mammalia -2.133185 -1.7193276 nonvolant
## Octodon_degus   Mammalia -0.190961 -0.3721030 nonvolant

##notice in the comp_data$data that there are now no birds
###these have been dropped
head(comp_data$dropped)

## $tips
## character(0)
##
## $unmatched.rows
## [1] "Buteo_jamaicensis"      "Gyps_fulvus"
## [3] "Haliaeetus_leucocephalus" "Parabuteo_unicinctus"
## [5] "Anas_acuta"             "Anser_anser"
## [7] "Aythya_marila"          "Branta_bernicle"
## [9] "Cygnus_olor"            "Dendrocygna_viduata"
## [11] "Uria_aalge"             "Leptoptilos_crumeniferus"
## [13] "Colius_striatus"        "Caloenas_nicobarica"
## [15] "Columba_livia"          "Dacelo_novaeguineae"
## [17] "Falco_tinnunculus"      "Gallus_gallus"
## [19] "Pavo_cristatus"         "Phasianus_colchicus"
## [21] "Grus_virgo"             "Balearica_regulorum"
## [23] "Cinclus_cinclus"        "Petrochelidon_pyrrhonota"
## [25] "Riparia_riparia"        "Ficedula_hypoleuca"
## [27] "Sericornis_frontalis"    "Parus_major"
```



```
## [29] "Sturnus_vulgaris"      "Zosterops_lateralis"
## [31] "Pelecanus_crispus"     "Pelecanus_onocrotalus"
## [33] "Geronticus_eremita"    "Threskiornis_aethiopicus"
## [35] "Phoenicopterus_roseus" "Phoebastria_immutabilis"
## [37] "Phoebastria_nigripes"  "Amazona_aestiva"
## [39] "Ara_ararauna"          "Ara_macao"
## [41] "Ara_militaris"         "Cacatua_galerita"
## [43] "Cacatua_moluccensis"   "Cyanoliseus_patagonus"
## [45] "Eclectus_roratus"      "Cacatua_roseicapilla"
## [47] "Melopsittacus_undulatus" "Myiopsitta_monachus"
## [49] "Nymphicus_hollandicus" "Poicephalus_senegalus"
## [51] "Psittacula_krameri"    "Psittacus_erithacus"
## [53] "Aptenodytes_patagonicus" "Spheniscus_demersus"
## [55] "Bubo_bubo"             "Dromaius_novaehollandiae"
## [57] "Struthio_camelus"      "Sula_dactylatra"
```

We can now run some models, first lets run two models with lambda fixed to 1 and something close to 0.

```
####lets define our fixed factors
formula_a <- longevity ~ mass + volant

#object comp_data which contains but the phylogeny and the data.
#Lets set the lambda in this case to 1.

pgls_l1 <- pgls(formula = formula_a, data = comp_data, lambda = c(1))
pgls_l0 <- pgls(formula = formula_a, data = comp_data, lambda = c(0.01))
summary(pgls_l1)
```

```
##
## Call:
## pgls(formula = formula_a, data = comp_data, lambda = c(1))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.29647 -0.08937 -0.02065  0.06000  0.67568
##
## Branch length transformations:
##
## kappa  [Fix]   : 1.000
## lambda [Fix]   : 1.000
## delta  [Fix]   : 1.000
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.053484   0.792826  1.3288   0.1862
## mass        -0.433933   0.078131 -5.5539 1.495e-07 ***
## volantvolant -1.059141   0.812158 -1.3041   0.1945
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1252 on 131 degrees of freedom
## Multiple R-squared:  0.1906, Adjusted R-squared:  0.1782
## F-statistic: 15.42 on 2 and 131 DF, p-value: 9.663e-07
```

```
summary(pglis_10)
```

```
##
## Call:
## pglis(formula = formula_a, data = comp_data, lambda = c(0.01))
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.226378 -0.021377  0.002155  0.017977  0.290707
##
## Branch length transformations:
##
## kappa  [Fix]  : 1.000
## lambda [Fix]  : 0.010
## delta  [Fix]  : 1.000
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.0011739  0.0991748   0.0118 0.9905740
## mass         0.4760918  0.0586018   8.1242 2.931e-13 ***
## volantvolant 0.9950322  0.2533605   3.9273 0.0001382 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.05256 on 131 degrees of freedom
## Multiple R-squared:  0.335,    Adjusted R-squared:  0.3249
## F-statistic:    33 on 2 and 131 DF,  p-value: 2.474e-12
```

Next lets run a model were lambda is no longer fixed. We can do this by specifying lambda = “ML” which tells the model to estimate it using maximum likelihood.

```
#Finally we also need to set the lambda in this case to ML.
#This means the we will using Maximum Likelihood
#to calculate the lambda.
pglis_mod <- pglis(formula = formula_a, data = comp_data, lambda = "ML")
summary(pglis_mod)
```

```
##
## Call:
## pglis(formula = formula_a, data = comp_data, lambda = "ML")
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.194491 -0.038027  0.004578  0.031555  0.151222
##
## Branch length transformations:
##
## kappa  [Fix]  : 1.000
## lambda [ ML]  : 0.933
##   lower bound : 0.000, p = < 2.22e-16
##   upper bound : 1.000, p = < 2.22e-16
##   95.0% CI    : (0.866, 0.968)
## delta  [Fix]  : 1.000
##
## Coefficients:
```

```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept) 0.076286  0.394227  0.1935  0.84686
## mass        0.416752  0.070962  5.8729 3.334e-08 ***
## volantvolant 0.988252  0.427932  2.3094  0.02249 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06309 on 131 degrees of freedom
## Multiple R-squared:  0.2084, Adjusted R-squared:  0.1963
## F-statistic: 17.25 on 2 and 131 DF, p-value: 2.246e-07
```

Now under Branch length transformations we also now get the estimated branch transformation under maximum likelihood. As we are only interested in fitting only lambda for now the other types of transformations, (kappa and delta), are held fixed.

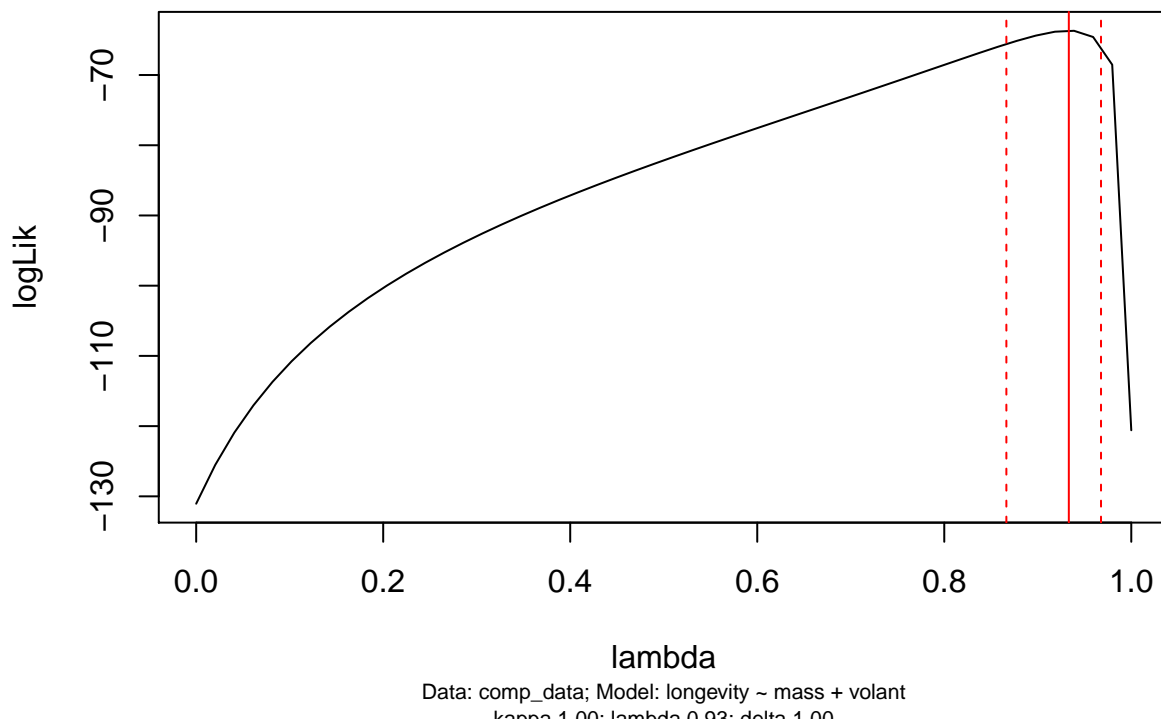
Lambda here estimated as

```
pgls_mod$param["lambda"]
```

```
##      lambda
## 0.9331958
```

As it is close to 1 the traits in this model are correlated under Brownian motion. If our value was 0 it would indicate that our data points are essentially independent. We can then go and check various elements of the model such as the likelihood profile for lambda.

```
mod_profile <- pgls.profile(pgl_mod)
plot(mod_profile)
```



which looks good as the profile shows a nice clear peak.

We would also then go ahead and check our residuals etc but for now we will assume everything is good and move onto running a similar model using MCMCglmm

## MCMCglmm

So far we have fitted a pgls model that included phylogeny to account for non-independence. Now we will use a Bayesian approach where we include phylogeny as a random term using the animal model in the MCMCglmm package.

As we are using a Bayesian approach we will first set up the priors. In most cases we want to use a non-informative prior that doesn't influence the estimated posterior distribution. For the random effect prior we will use an inverse-Gamma distribution. In MCMCglmm this is described by two parameters nu and V. These terms are related to the shape (alpha) and scale (beta) parameters on an inverse-Gamma with  $\alpha = \text{nu}/2$ , and  $\text{Beta} = (\text{nu} * V)/2$ . As we don't want our estimates to be heavily influenced by our prior we will use weakly informative prior values such as described as  $V = 1$  and  $\text{nu} = 0.002$ . (For more on priors for the animal model see course notes)

```
prior <- list(R = list(V=1, nu=0.002),
             G = list(G1 = list(V=1, nu=0.002)))
```

We describe our prior as above for the random (G) and residual variances (R) each of them as a list, which we will in turn put within a list. If we wanted to include more random terms we would include a G2, G3 etc for each additional random term within the G list. We could also specify priors for the fixed terms using B, however MCMCglmm will automatically do that for us and as it usually does a good job at it we will ignore it here.

Next we need to decide on the parameters relating to running the mcmc chain in the model. We need to include how many iterations we want to run the chain for (nitt), the burnin we want to discard at the start of the chain (burnin) and also how often we want to sample and store from the chain (thin). We discard a burnin as we don't want the starting point of the chain to over-influence our final estimates. For now let's just use a burnin of 1/6 of the nitt, just to be safe. The thinning is used to help reduce autocorrelation in our sample, how much you use often depends on how much autocorrelation you find.

To save time we will only run this model over 12000 iterations (However, much larger nitt is often required).

```
#no. of iterations
nitt <- c(12000)
#length of burnin
burnin <- c(2000)
#amount of thinning
thin <- c(5)
```

Now we need to set up the data. We have already cleaned and matched up our data earlier using the comparative.data function but we need to now add an extra column into our dataset called "animal" which contains the species matched between the tree and the data.

```
#Matched data
mcmc_data <- comp_data$data
#As MCMCglmm requires a column named animal for it to identify it
#as a phylo model we include an extra column with the species names in it.
mcmc_data <- cbind(animal = rownames(mcmc_data), mcmc_data)
mcmc_tree <- comp_data$phy
```

MCMCglmm reserves the random variable "animal" to call a model that includes the phylogeny as an additive genetic effect. If we name it something else, like say "species", MCMCglmm will either throw an error looking for "animal", or if we do not provide a phylogeny under pedigree it will run "species" like a standard random term. Now we can run the model.

```
mod_mcmc <- MCMCglmm(fixed = formula_a,
                    random= ~ animal,
                    family="gaussian",
```

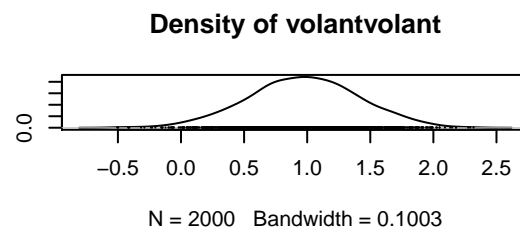
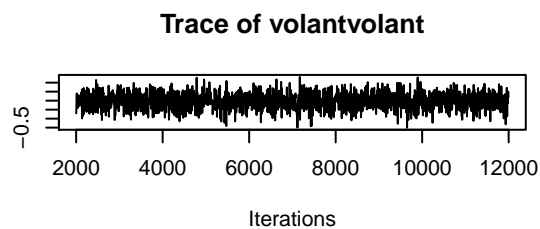
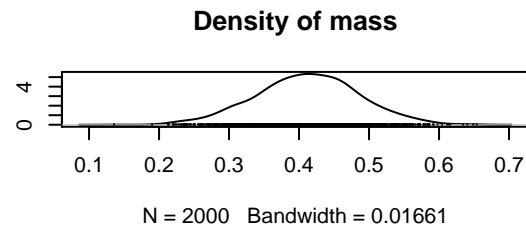
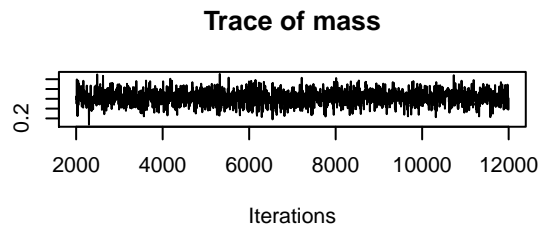
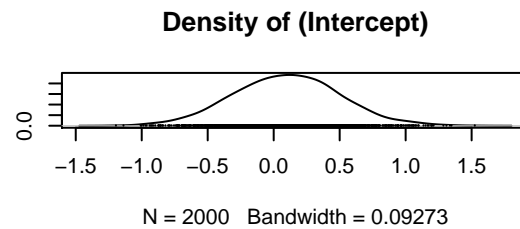
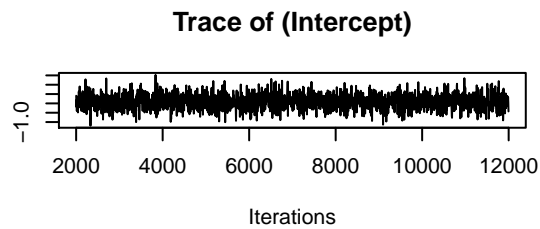
```
pedigree = mcmc_tree,  
data = mcmc_data,  
nitt = nitt,  
burnin = burnin,  
thin = thin,  
prior = prior)
```

```
##  
##           MCMC iteration = 0  
##  
##           MCMC iteration = 1000  
##  
##           MCMC iteration = 2000  
##  
##           MCMC iteration = 3000  
##  
##           MCMC iteration = 4000  
##  
##           MCMC iteration = 5000  
##  
##           MCMC iteration = 6000  
##  
##           MCMC iteration = 7000  
##  
##           MCMC iteration = 8000  
##  
##           MCMC iteration = 9000  
##  
##           MCMC iteration = 10000  
##  
##           MCMC iteration = 11000  
##  
##           MCMC iteration = 12000
```

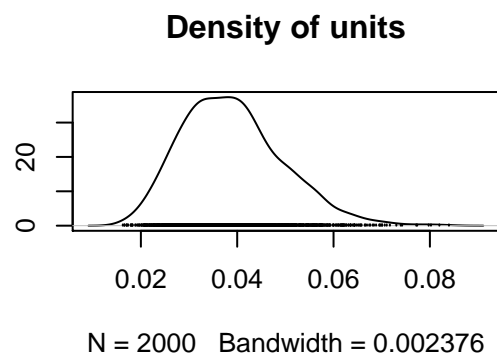
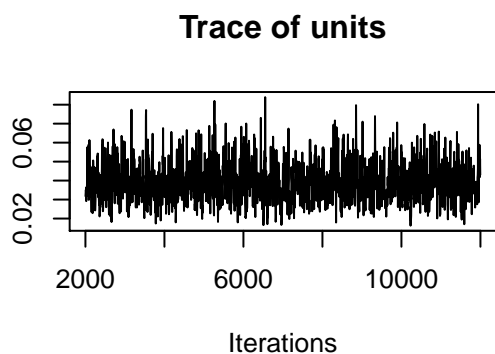
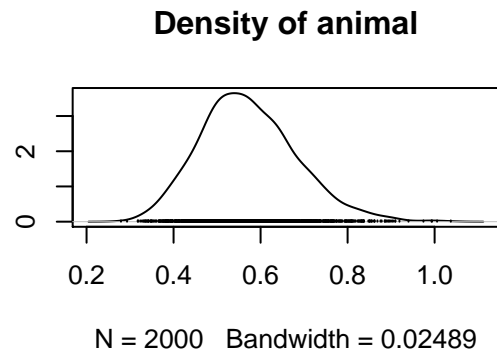
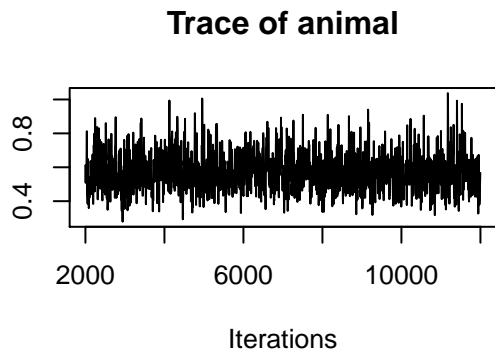
As the model runs we see the iterations print out. These chains can take some time to run, depending on the model, however, since we only ran our chains for 12000 iterations it doesn't take long here.

Before we even look at our model we need to check if the model ran appropriately. We can do this by visually inspecting the chains to make sure there has been no unruly behaviour! We can extract the full chains using `modelSol` for the fixed effects and `modelVCV` for the random effect variances. So `Sol[,1]` will give you the first fixed term, in this case the intercept, and `VCV[,1]` will give you the first random term, which is “animal” and so on. As our model is an `mcmc` object when we use the `plot` function we get a trace plot.

```
plot(mod_mcmc$Sol)
```



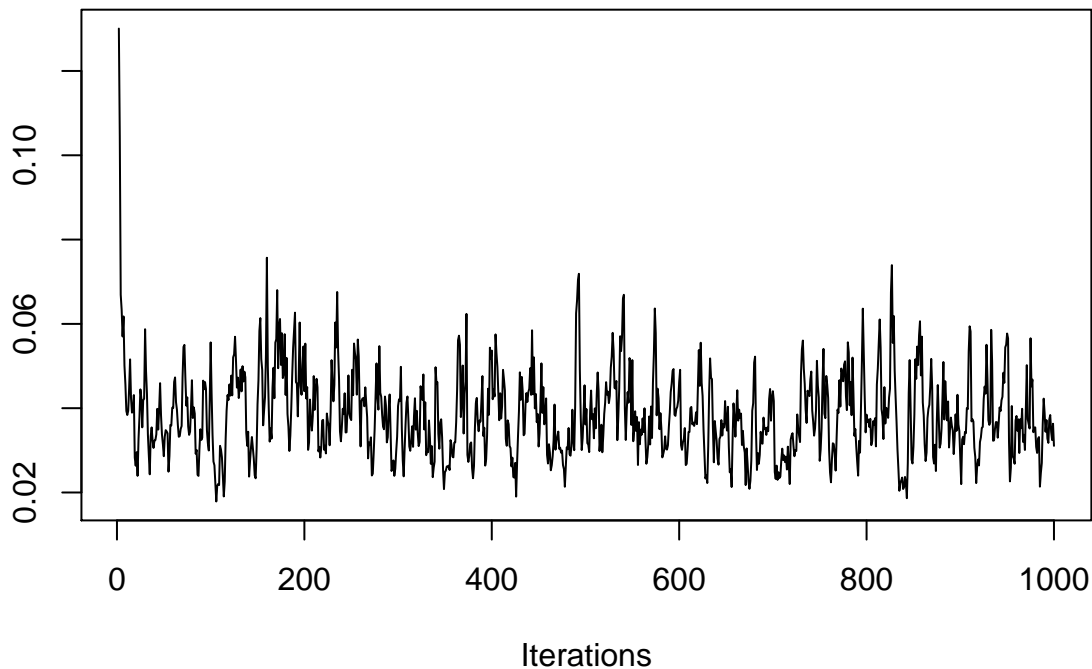
```
plot(mod_mcmc$VCV)
```



What we want to see “hairy caterpillars” and a trace with no obvious trend that is bouncing around some

stable point.

What we don't want to see in the trace plots can be demonstrated if we only run our model over a very short chain (itt == 1000). Notice that without a burnin the start of trace is well outside the area that the chain will converge towards.



So in our longer run model everything looks good visually, however we also want to check the level of autocorrelation in these traces. We can do this using `autocorr.diag()` which gives the level of correlation along the chain between some lag sizes.

```
autocorr.diag(mod_mcmc$Sol)
```

```
##           (Intercept)           mass volantvolant
## Lag 0      1.000000000  1.00000000  1.00000000
## Lag 5     -0.012043947  0.03392189 -0.01423925
## Lag 25    -0.001472627  0.01622992  0.02448349
## Lag 50    -0.021021924 -0.02729474 -0.01965705
## Lag 250   -0.002009103  0.01196044 -0.01705418
```

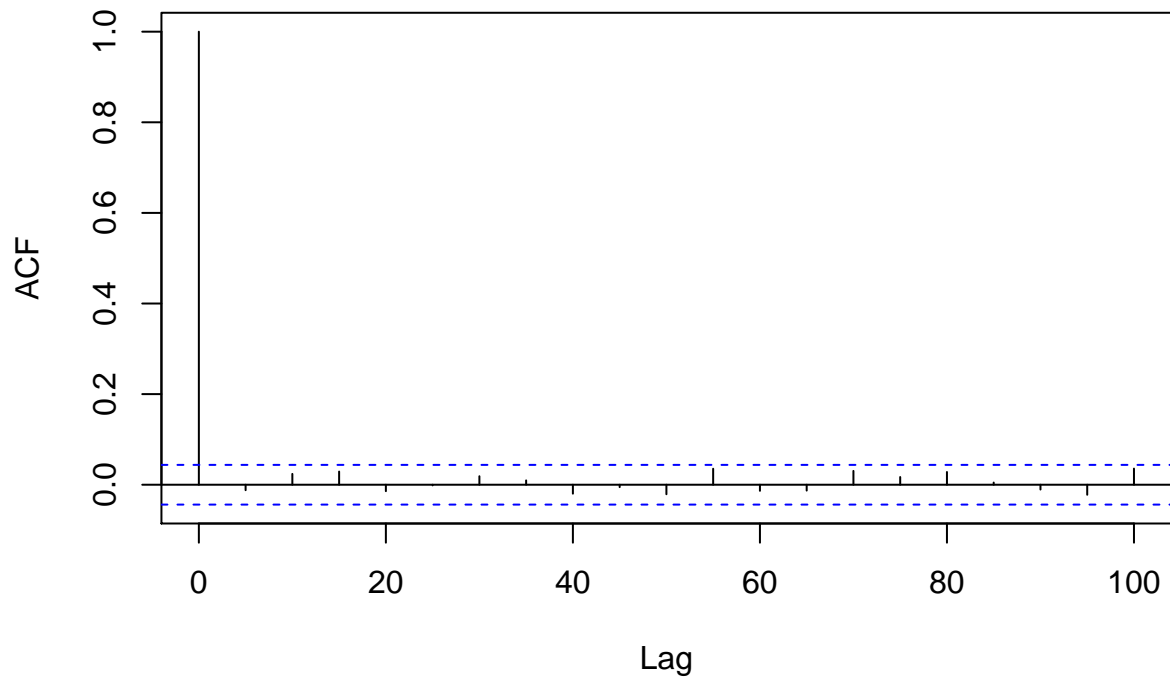
```
autocorr.diag(mod_mcmc$VCV)
```

```
##           animal           units
## Lag 0      1.00000000  1.00000000
## Lag 5      0.33261821  0.35118680
## Lag 25     0.01376557  0.05241862
## Lag 50    -0.01566639 -0.01424567
## Lag 250    0.02959519 -0.03046731
```

or we can look at autocorrelation plots for each of the traces, we'll look at just one using the `acf` function here.

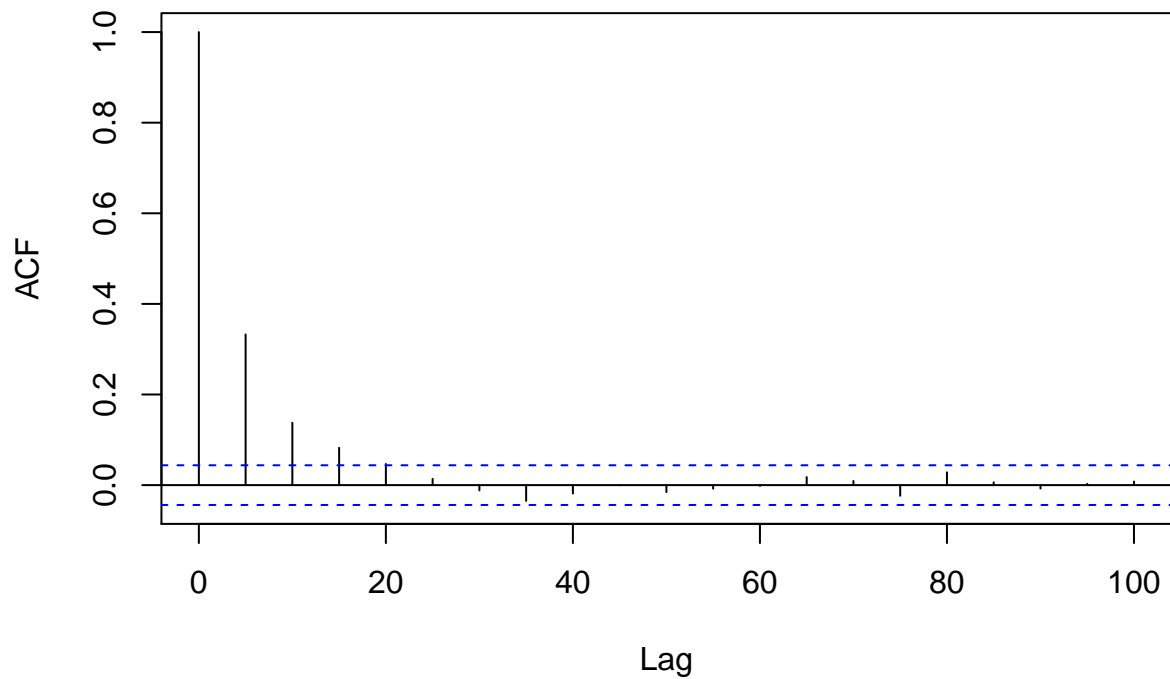
```
#acf plot for the first fixed estimate in our model (the intercept)
acf(mod_mcmc$Sol[,1], lag.max = 20)
```

### Series mod\_mcmc\$Sol[, 1]



```
#acf plot for the first random term in our model (the animal term)  
acf(mod_mcmc$VCV[,1], lag.max = 20)
```

### Series mod\_mcmc\$VCV[, 1]



For our intercept the autocorrelation plot looks good, however the animal term still shows some autocorrelation. One quick way to deal with this is to simply increase the thinning.



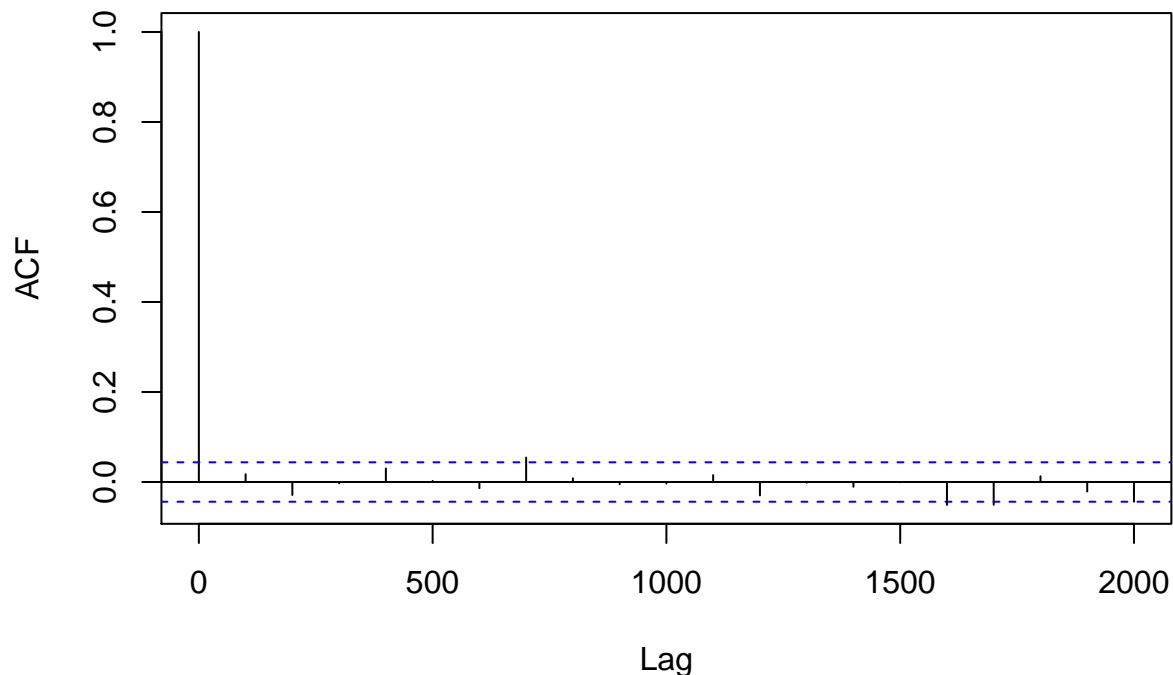
```

nitt2 <- 240000
burnin2 = 40000
thin2 = 100
mod_mcmc_long <- MCMCglmm(fixed = formula_a,
  random= ~ animal,
  family="gaussian",
  pedigree = mcmc_tree,
  data = mcmc_data,
  nitt = nitt2,
  burnin = burnin2,
  thin = thin2,
  prior = prior,
  verbose=FALSE)

acf(mod_mcmc_long$VCV[,1], lag.max =20)

```

**Series mod\_mcmc\_long\$VCV[, 1]**



That looks better now. Noticed I also increased the number of iterations. One rough and ready rule that I like to use is to aim for an effective sample size of my chains, which is the number of iterations used in the posterior after the burnin, thinning and accounting for autocorrelation, somewhere between 1000-2000.

```

#acf plot for the first fixed estimate in our model (the intercept)
effectiveSize(mod_mcmc_long$Sol)

```

```

## (Intercept)      mass volantvolant
##    2000.000    2000.000    1771.895

```

```

effectiveSize(mod_mcmc_long$VCV)

```

```

## animal  units
##    2000    2000

```

One last thing to check is that our MCMC chain has properly converged and that our estimate is not the result of some type of transitional behaviour. That is have our chains “found” the optimum or do we need to let them run longer before they settle around some estimate. To check this we will run a second model and see if it converges on the same estimates as our first model.

```
mod_mcmc_2 <- MCMCglmm(fixed = formula_a,
                      random= ~ animal,
                      family="gaussian",
                      pedigree = mcmc_tree,
                      data = mcmc_data,
                      nitt = nitt2,
                      burnin = burnin2,
                      thin = thin2,
                      prior = prior,
                      verbose=FALSE)
```

We can now check the convergence of the two chains using the Gelman and Rubin Multiple Sequence Diagnostic. This calculates the within-chain and between-chain variance of the chains and then gives a scale reduced factor, (for more see here. When this number is close to one (say below 1.1) the chains are indistinguishable and hence can be considered to be converged.

```
#checking convergence for our fixed factors
gelman.diag(mcmc.list(mod_mcmc_long$Sol, mod_mcmc_2$Sol))
```

```
## Potential scale reduction factors:
##
##           Point est. Upper C.I.
## (Intercept)          1          1.00
## mass                 1          1.00
## volantvolant         1          1.01
##
## Multivariate psrf
##
## 1
```

```
#checking convergence for our random terms
gelman.diag(mcmc.list(mod_mcmc_long$VCV, mod_mcmc_2$VCV))
```

```
## Potential scale reduction factors:
##
##           Point est. Upper C.I.
## animal              1          1
## units                1          1
##
## Multivariate psrf
##
## 1
```

Since everything looks good, we will finally look at the results of our model.

```
summary(mod_mcmc_long)
```

```
##
## Iterations = 40001:239901
## Thinning interval = 100
## Sample size = 2000
##
## DIC: 23.53959
```

```
##
## G-structure: ~animal
##
##      post.mean l-95% CI u-95% CI eff.samp
## animal    0.5765  0.3791  0.7966    2000
##
## R-structure: ~units
##
##      post.mean l-95% CI u-95% CI eff.samp
## units    0.03924  0.01946  0.05986    2000
##
## Location effects: longevity ~ mass + volant
##
##      post.mean l-95% CI u-95% CI eff.samp pMCMC
## (Intercept)  0.07853 -0.64964  0.92337    2000  0.835
## mass         0.41269  0.28189  0.58085    2000 <5e-04 ***
## volantvolant  0.96283  0.12242  1.83950    1772  0.037 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

First off we can find the estimates for the fixed factors are under the Location effects section (Again notice the similarity to our pgl's model). Each parameter has a measure of the effect size using the post.mean and a lower and higher 95% credible interval (CI). These are simply calculated from the posterior distributions we looked at in the above plots, so if you would rather calculated the median instead of using the mean we can simple use

```
median(mod_mcmc_long$Sol[,1])
```

```
## [1] 0.07456271
```

We also have the effective sample size (eff.samp) and the pMCMC which calculated as two times the probability that the estimate is either > or < 0, using which ever one is smaller. However since our data has been mean centred and expressed in units of standard deviation we can look at what proportion of our posterior is on either side of zero.

For the random terms we have the posterior distribution for our G-structure which includes or phylogenetic effect and the R-structure which is our residual variation.

We also have the DIC which is a Bayesian version of AIC. Like AIC it is a measure of the trade-off between the “fit” of the model and the number of parameters, with a lower number better.

Finally, we can also calculate the  $H^2$  which is comparable to pagels lambda as

```
H <- (var(mod_mcmc_long$VCV[, "animal"])) /
      (var(mod_mcmc_long$VCV[, "animal"])
       + var(mod_mcmc_long$VCV[, "units"]))
H
```

```
## [1] 0.9903149
```

Before moving on to the next section try running the above analysis subsetting for birds as opposed to mammals.

```
#aves tree from Jetz et al 2012
aves_tree <- trees_aves[[1]]
```