# SIMMs for complex data sets

Andrew Parnell
`andrew.parnell@mu.ie`

# Learning outcomes

▶ Understand how covariates might be included in a SIMM
▶ Understand how to compare models using DIC

# Revision: our best SIMM in JAGS

```
modelstring ='
model {
  for (i in 1:N) {
    for (j in 1:J) {
      y[i,j] ~ dnorm(inprod(p[i,]*q[,j], s_mean[,j]+c_mean[,j]) / inprod(p[i,],q[,j]), 1/var_y[i,j])
      var_y[i,j] <- inprod(pow(p[i,]*q[,j],2),s_sd[,j]^2+c_sd[,j]^2)/pow(inprod(p[i,],q[,j]),2)
        + pow(sigma[j],2)
    }
  }
  for(i in 1:N) {
    p[i,1:K] <- expf[i,]/sum(expf[i,])
    for(k in 1:K) {
      expf[i,k] <- exp(f[i,k])
      f[i,k] ~ dnorm(mu_f[k],sigma_f[k]^-2)
    }
  }
  for(k in 1:K) {
    mu_f[k] ~ dnorm(0,1)
    sigma_f[k] ~ dgamma(2,1)
  }
  for(j in 1:J) { sigma[j] ~ dunif(0,10) }
}
'
```

# Key features

- ▶ The key features of our current SIMM are:
    1. We are accounting for uncertainty in sources and TEFs
    2. Individual dietary proportions are provided for each consumer, arising from an overall mean
    3. The dietary proportions are linked to normal distributions via the *centralised log-ratio* (CLR) transform
- ▶ Some remaining restrictions
    1. We still haven't seen code that incorporates covariates
    2. We don't know how to compare between different model structures
    3. The sources, TEFs and consumers, are all assumed to be independent across isotope (i.e. 'circular' on an isospace plot). Could there be covariance between them?

# Adding covariates

- ▶ Let's now create a model with covariates on the dietary proportions
- ▶ Recall that we can do this using the CLR transform on the $p$:

$$[p_{i1}, \ldots, p_{iK}] = \left[ \frac{\exp(f_{i1})}{\sum_j \exp(f_{ij})}, \ldots, \frac{\exp(f_{iK})}{\sum_j \exp(f_{ij})} \right]$$

- ▶ The prior goes on $f$, e.g.

$$f_{ik} = \alpha_k + \beta_k x_i$$

where $x_i$ is the covariate for observation $i$

- ▶ Much of the detailed maths for this work is in our 2013 Environmetrics paper
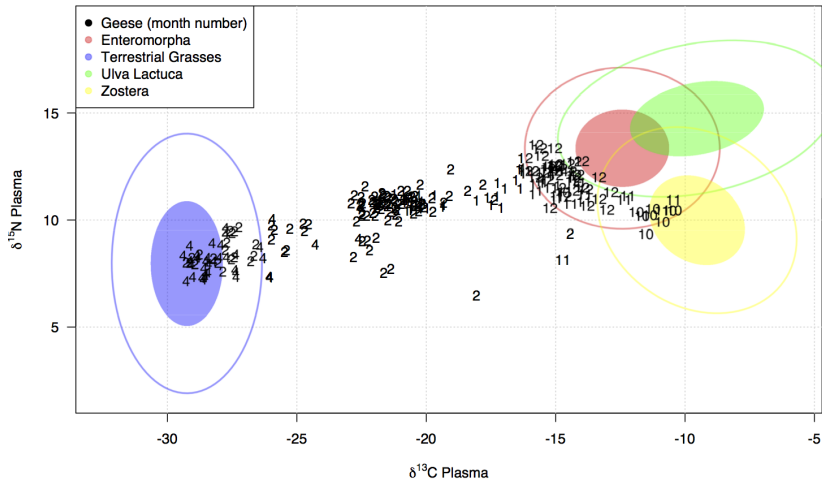- ▶ The CLR allows for much more complex relationships between the dietary proportions

# Why the CLR?

- There is quite a bit of research to show that the Dirichlet distribution is not a good distribution to use for proportions because it suffers from a very rigid correlation structure
- The CLR doesn't suffer from this, but does have an extra restriction that all the $f$s must sum to zero. You can get round this by setting an informative prior
- There are others used too, including the additive log ratio (ALR) and the isometric log ratio (ILR). We recommend the CLR with an informative prior (different to MixSIAR)
- Lots of other distributions are widely used but often inappropriate: e.g. Poisson, $\chi^2$, normal(!)

# A SIMM with covariates

```
modelstring ='
model {
  ...
  for(i in 1:N) {
    p[i,1:K] <- expf[i,]/sum(expf[i,])
    for(k in 1:K) {
      expf[i,k] <- exp(f[i,k])
      f[i,k] ~ dnorm(mu_f[i,k],sigma_f[k]^-2)
    }
  }
  for(k in 1:K) {
    for(i in 1:N) { mu_f[i,k] <- alpha[k] + beta[k]*x[i] }
    sigma_f[k] ~ dgamma(2,1)
    alpha[k] ~ dnorm(0,1)
    beta[k] ~ dnorm(0,1)
  }
  ...
```

# The Geese data

# A Fourier basis

- ▶ For the Geese data we don't want to use a linear covariate
- ▶ Instead we want to use a *Fourier* covariate which measures how the dietary proportions change periodically
- ▶ We're going to structure our mean as:

$$f_{ik} = \alpha_k + \beta_k sin\left(\frac{2\pi x_i}{365}\right) + \gamma_k cos\left(\frac{2\pi x_i}{365}\right)$$
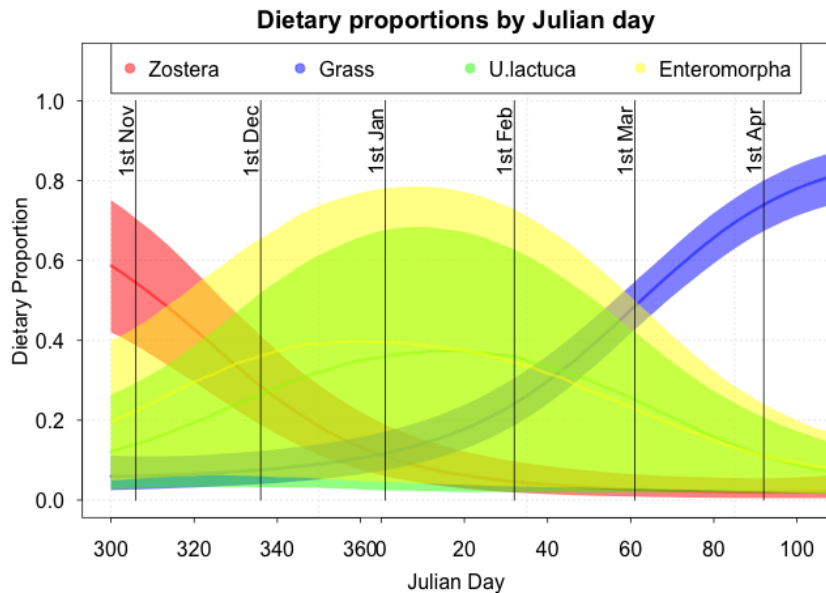
where $x_i$ is Julian day

- ▶ This will allow for periodic behaviour over 365 days. The sign and magnitude of the parameters $\alpha$ and $\beta$ will determine the shape of the periodic behaviour

# Model setup

```
modelstring ='
model {
  ...
  for(i in 1:N) {
    p[i,1:K] <- expf[i,]/sum(expf[i,])
    for(k in 1:K) { expf[i,k] <- exp(f[i,k]) }
  }
  for(k in 1:K) {
    f[1:N,k] <- X[1:N,1:L]%*%beta[1:L,k]
  }
  for(l in 1:L) {
    for(k in 1:K) { beta[l,k] ~ dnorm(0,1) }
  }
  ...
}
'
```

Full script in `run_geese_harmonic.R` file

**Dietary proportions by Julian day**

# Comparing models

- How do we know that this model fitted the data better than the model without covariates?
- How do we choose between models generally?
- These are very hard and open questions in the field of statistics

# A rough primer on model comparison

- $p$-values. The traditional way. These tell you very little about whether a parameter is important in the model
- Likelihood ratio tests (with $p$-values). A bit better. These compare how 'likely' the data is under one hypothesis vs the other.
- Information criteria. Idea here is to penalise the likelihood by some measure of 'model complexity', so as to choose models which fit the data well and are not too complex. We will use the *Deviance Information Criterion* (DIC), which is already part of JAGS
- Bayes Factors. These are theoretically the gold standard in Bayesian hypothesis testing. However, they can be very sensitive to the choice of prior distribution
- Cross-validation. Obtained by removing portions of the data, fitting to the remainder, and then predicting values for the missing portion. Very useful for larger data sets

# The Deviance Information Criterion

- The DIC is defined as:

$$DIC = -2 \log L + 2p_D$$

  where $L$ is the likelihood and $p_D$ is the *effective number of parameters*

- A smaller DIC indicates a 'better' model. DIC doesn't give any estimate of uncertainty so there is no way to discern exactly how small a jump is required to choose a model

- $p_D$ is approximately calculated as the difference between how well the model fits the data at the mean value of the parameters, and how well the model fits the data at the mean of the likelihood

- From JAGS, we can get DIC by running the extra command `dic.samples`. Note that the DIC sometimes takes much longer to converge than the parameters

# DIC example

- ► Compare the Geese model with the time series structure to a model without including the sine or cosine terms

- ► Procedure
    1. Set up each model in JAGS as normal
    2. Additionally run the dic.samples function
    3. Extract the DIC and $p_D$ value
    4. Choose the model that has the smallest value

- ► Can also do this amongst multiple different models, providing they all use the same data

- ► Need to make sure that each model has converged

# DIC example - code

```
X = cbind(1,sin(2*pi*con$julianday/365),cos(2*pi*con$julianday/365))
data=list(y=con[,2:3],s_mean=sources[,c(1,3)],s_sd=sources[,c(2,4)],
          c_mean=tefs[,c(1,3)],c_sd=tefs[,c(2,4)],
          q=cd,N=nrow(con),K=nrow(sources),
          J=ncol(con[,2:3]),X=X,L=ncol(X))
model=jags.model(textConnection(modelstring), data=data, n.chains=3)
dic.samples(model,n.iter=2000)
# Mean deviance:  1793
# penalty 9.512
# Penalized deviance: 1802
X2 = cbind(1,con$julianday/365)
data2=list(y=con[,2:3],s_mean=sources[,c(1,3)],s_sd=sources[,c(2,4)],
          c_mean=tefs[,c(1,3)],c_sd=tefs[,c(2,4)],
          q=cd,N=nrow(con),K=nrow(sources),
          J=ncol(con[,2:3]),X=X2,L=ncol(X2))
model2=jags.model(textConnection(modelstring), data=data2, n.chains=3)
dic.samples(model2,n.iter=2000)
# Mean deviance:  1879
# penalty 8.922
# Penalized deviance: 1888
```

# Some final notes on the DIC

▶ DIC can be quite sensitive to model focus. If you move parameters around (like in the centered random effects model) you can get different DIC values

▶ The value of $p_D$ is useful. For simple models, it should be roughly the true number of parameters (e.g. 5 and 4 in the previous models). However, in hierarchical models it can be non-integer representing the fact that the parameters are shared between groups. In some cases $p_D$ can be negative!

▶ JAGS contains the option to create a 'superior' version of $p_D$ called popt which you specify via the type argument. This penalises extra parameters more harshly but isn't quite as interpretable. However, it is often more stable in more complicated models.

# Summary

- We can add in rich covariate behaviour through the CLR, though need to be careful with priors
- DIC can help us choose between models. More complex models always fit the data better, but can often over-fit yielding poor predictive and explanatory performance