

# 1 Hierarchical Generalized Additive Models: an 2 introduction with mgcv

## 3 ABSTRACT

4 This is just placeholder text until we write a proper abstract

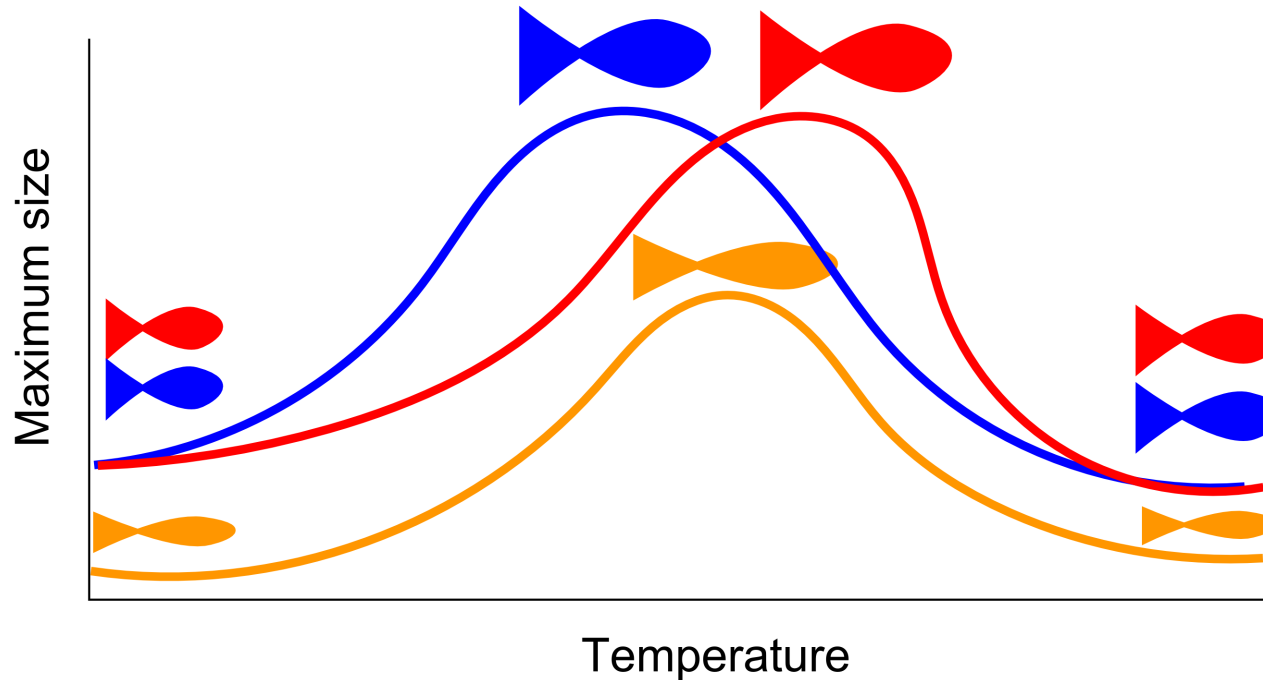
## 5 I: INTRODUCTION

6 As ecology has progressed as a quantitative discipline and the questions ecologists ask have become  
7 more complicated, the statistical techniques ecologists use have increased in their flexibility to  
8 model complex relationships. Two of the more popular and powerful techniques now in use  
9 are generalized additive models (GAMs; Simon N Wood 2006) for modelling flexible regression  
10 functions, and generalized linear mixed models (“hierarchical generalized linear models” HGLMs  
11 or simply “hierarchical models”; Bolker et al. 2009, A. Gelman et al. (2013)) for modelling  
12 between-group variability in regression relationships.

13 At first glance, GAMs and HGLMs are very different tools. GAMs are used to estimate  
14 smooth functional relationships between predictor variables and the response, assuming that  
15 the phenomena under investigation is not linear (in the GLM sense) but is a smooth function  
16 of the predictor variables. Example of such relationships would be the vertical distribution of  
17 abundance of a population as a function of depth (Stanley, Pedersen, and Snelgrove 2016) or  
18 the swimming speed of snakes as function of temperature (Vickers, Aubret, and Coulon 2017).  
19 HGLMs, on the other hand, are used to estimate linear relationships between predictor variables  
20 and response, but impose a structure where predictors are organized into groups (often referred  
21 to as “blocks”) and the relationships between predictor and response may differ between those  
22 groups. Either or both slope and intercept may be subject to grouping. A typical example of  
23 HGLM use might be to include site-specific effects in a model of counts, or to model individual  
24 level heterogeneity in a study with repeated observations of multiple individuals.

25 Both GAMs and HGLMs can be used to fit potentially highly variable models by “pooling”  
26 parameter estimates towards one another. The connection between the two methods is quite  
27 deep and a GAMs may be interpreted (and fitted) as HGLMs and vice-versa. Given this  
28 connection, the obvious extension to the standard GAM framework is to allow the smooth  
29 functional relationship between predictor and response to vary between different grouping levels,  
30 but in such a way that the different functions are in some sense pooled toward each other. We  
31 often want to know both how the functional relationship between varies between groups, and if  
32 there is a strong relationship on average across groups. We will refer to this type of model as a  
33 *hierarchical GAM*, or HGAM.

34 There are many potential uses for HGAMs. For example, to estimate how the maximum size  
35 different fish species reach varies along a common temperature gradient (figure 1). Each species  
36 will typically have its own response function, but since the species overlap in range, they should  
37 have similar responses over at least some of the temperature gradient; figure 1 shows all three  
38 species reach their largest maximum sizes in the center of the temperature gradient. Estimating  
39 a separate function for each species throws away a lot of shared information and could result in  
40 highly noisy function estimates if there were only a few data points for each species. Estimating  
41 a single average relationship could result in an average function that did not predict any specific  
42 group well. In our example, using a single global temperature-size relationship would miss the  
43 three species distinct temperature optima, and that the orange species is significantly smaller at



**Figure 1.** Hypothetical example of functional variability between different group levels. Each line indicates how the maximum possible body size for different species of fish in a community might vary as a function of average water temperature. While the orange species shows lower maximum size at all temperatures, and the red and blue species differ in what temperature they can achieve the maximum possible size, all three curves are similarly smooth, and peak close to one another, relative to the entire range of tested temperatures.

all temperatures than the other two (figure 1). We prefer a hierarchical model that fit a single global temperature-size curve plus species-specific curves that were penalized to be close to the mean function.

The capability to fit HGAMs already exists in the popular *mgcv* package for the R statistical programming language. There are many different options available representing different model assumptions with corresponding trade-offs. This paper will cover the different approaches to group-level smoothing, the options for each one and why a user might choose it, and demonstrate the different approaches across a range of case studies.

This paper is divided into six sections. Part II is a brief (and friendly) review of how generalized additive models work and their relation to hierarchical models. In part III, we discuss different ways of modelling hierarchical generalized additive models, what assumptions each model makes about how information is shared between groups, and different ways of specifying these models in *mgcv*. In part IV, we discuss some of the tools available for plotting model output and assessing model goodness of fit. In part V, we discuss some of the computational and statistical issues involved in fitting HGAMs in *mgcv*. Finally, in part VI, we work through a few examples of analyses using this approach, to demonstrate the modelling process and how hierarchical GAMs can be incorporated into the quantitative ecologist's toolbox.

## II: AND INTRODUCTION TO GENERALIZED ADDITIVE MODELS

One of the most common model formulations in statistics is the generalized linear model (McCullagh and Nelder 1989) — models that relate some response ( $y$ ) to linear combinations of

explanatory variables. We may allow the response to be distributed according to some exponential family distribution (e.g., letting the response be a trial, a count or a strictly positive number – binomial, Poisson or Gamma distributions, respectively). The generalized additive modelling (GAM) framework (Hastie and Tibshirani 1990; Ruppert, Wand, and Carroll 2003; Simon N Wood 2006) allows the relationships between the explanatory variables (henceforth covariates) and the response to be described by smooth terms (usually *splines* (Boor 1978), but potentially other structures). In general we are then talking about models of the form:

$$\mathbb{E}(y) = g^{-1} \left( \beta_0 + \sum_{j=1}^J f_j(x_j) \right),$$

where  $y$  is the response (with an appropriate distribution and link function  $g$ ),  $f_j$  is a smooth function of the covariate  $x_j$ ,  $\beta_0$  is an intercept term and  $g^{-1}$  is the inverse link function. Here there are  $J$  smooths and each is a function of only one covariate, though it is possible to construct smooths of multiple variables.

Each smooth  $f_j$ s is represented by a sum of simpler *basis functions* ( $b_k$ ) multiplied by corresponding coefficients ( $\beta_k$ ), which need to be estimated to be estimated:

$$f_j(x_j) = \sum_{k=1}^K \beta_k b_k(x_j),$$

The size of  $K$  of each smooth will determine the flexibility of the resulting term (referred to as “basis size”, “basis complexity” or “basis richness”). Though it seems like the basis can be overly complex (“how big should I make  $K$ ?”) and lead to overfitting, we need not worry about this as we use a penalty to ensure that the functions complexity is appropriate; hence the basis only need to be “large enough” and we let the penalty deal with excess wigglyness.

The penalty for a term is usually based on derivatives of that term – as the derivatives give the wigglyness of the function and hence its flexibility. We trade-off the fit of the model against the wigglyness penalty to obtain a model that both fits the data well but does not overfit. To control this trade-off we estimate a *smoothing parameter*. Figure 2 shows optimal smoothing (where the smoothing parameter is estimated to give a parsimonious model) in the first plot; the second plot shows what happens when the smoothing parameter is set to zero, so the penalty has no effect (interpolation); the right plot shows when the smoothing parameter is set to a very large value, giving a straight line. Smooths of this kind are often referred to as a *basis-penalty smoothers*.

**say something about knots!**

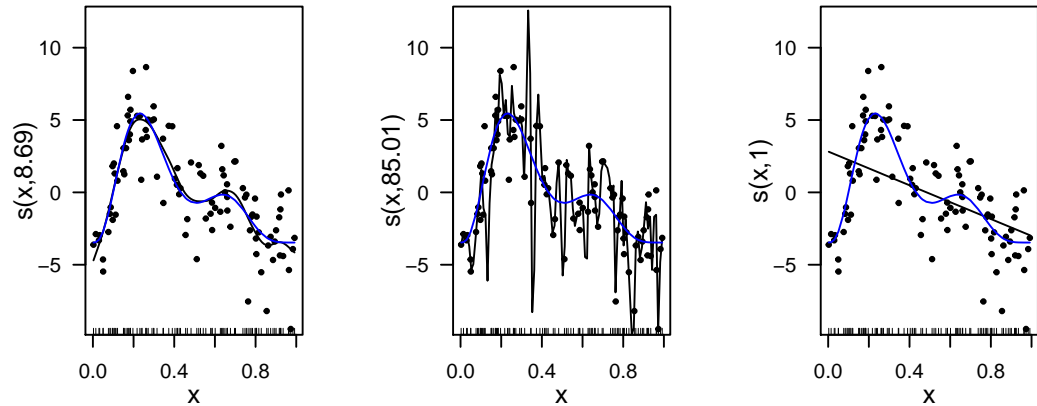
**say something about penalty matrices and also smoothing parameters**

The number of basis functions,  $K$ , limits the maximum basis complexity for a given smooth term. To measure the wigglyness of a given term, we use the *effective degrees of freedom* (EDF) which, at a maximum is the number of coefficients to be estimated in the model, minus any constraints. The EDF can take non-integer values and a larger value indicates a more wiggly term. See Simon N Wood (2006) Section 4.4 for further details.

There are many possible basis functions that can be used to model the  $b_k$ s. Here we’ll use thin plate regression splines, which have the appealing property that knot choice is somewhat automatic (the best approximation to including knots at each data point is used; Simon N. Wood (2003)).

**DLM:: put the cubic splines back in here**

TPRS are also defined for any number of predictors, so multivariate smoothers can be constructed easily. The basis is *isotropic* so smoothing is treated the same in all directions. So if one had, a bivariate smooth of temperature and time, a one degree change in temperature would equate to

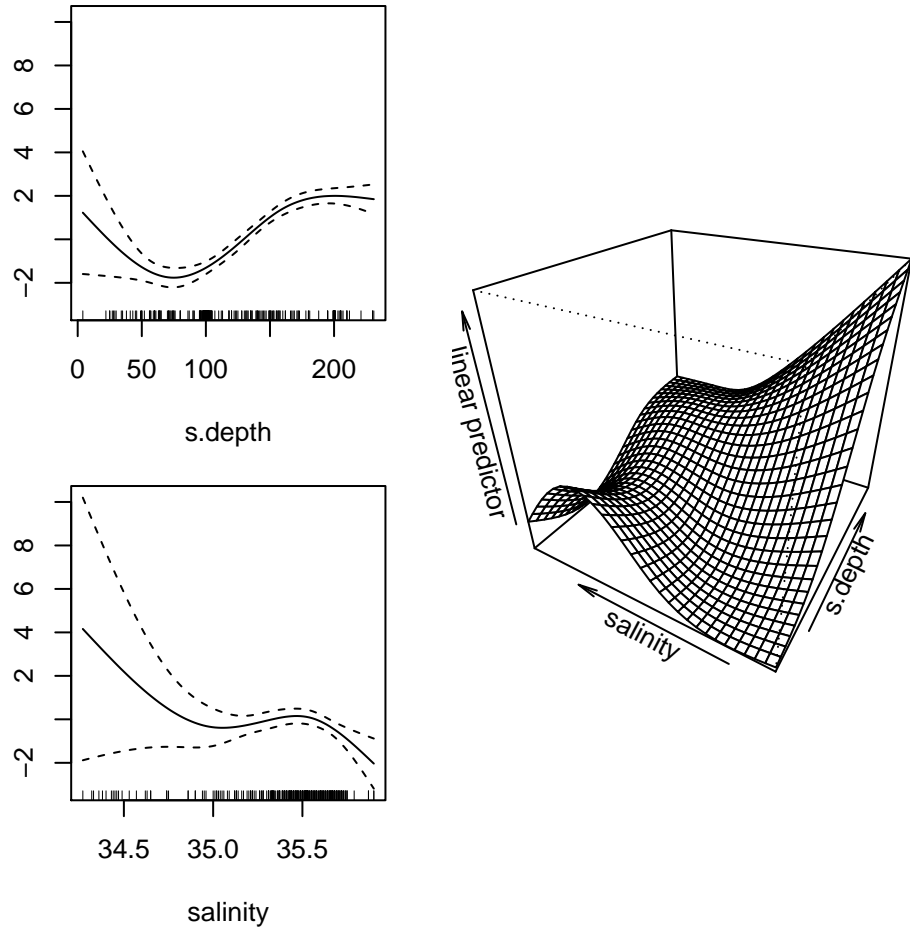


**Figure 2.** Examples of how different choices of the smoothing parameter effect the resulting function. Data (points) were generated from the blue function and noise added to them. In the left plot the smoothing parameter was estimated to give a good fit to the data, in the middle plot the smoothing parameter was set to zero, so the penalty has no effect and the function interpolates the data, the right plot shows when the smoothing parameter is set to a very large value, so the penalty removes all terms that have any wigglyness, giving a straight line. Numbers in the  $y$  axis labels show the estimated degrees of freedom for the term.

105 a one second change in time, which is an odd assumption to make. In the more general case  
 106 where units are not alike, we can use *tensor products* to combine two or more univariate smooths  
 107 into a more complex basis. Each component can be made up from a different basis, playing to  
 108 their particular strengths.

109 In the linear modelling literature we can specify a single interaction between terms (in R, `a:b`)  
 110 or a “full interaction”, which includes the marginal terms (`a*b` in R, which is equivalent to `a +`  
 111 `b + a:b`). There are parallels for smooths too, allowing us to separate-out the main effect terms  
 112 from the interactions (in R `te(a, b)` specifies the tensor product which is equivalent to `ti(a) +`  
 113 `ti(b) + ti(a, b)`). The ability to separate out the interactions and main effects will become  
 114 very useful in the next section, once we start looking at group-level smooths. For an example  
 115 of a tensor product interaction, see figure 3, which illustrates marginal smooths and a tensor  
 116 product interaction of depth and salinity predicting mackerel egg counts from a marine survey.

117 We represent the terms in our model as basis functions, which end up as additional columns  
 118 in our design matrix and parameter vector, and penalties, which penalize the likelihood and  
 119 stop our model from being too wiggly. Taking a pragmatic Bayesian approach to the problem,  
 120 the penalty is really a prior on how we think the model should act. In which case the penalty  
 121 matrix itself is a prior precision matrix (inverse variance) for the term. With that in mind,  
 122 we can think about random effects as “smooths” in our model, albeit ones with ridge penalties  
 123 (Kimeldorf and Wahba 1970; Simon N. Wood 2017). For instance, to include a random effect  
 124 modelling between group variation in intercepts there will be one basis function for each level  
 125 of the grouping variable, that takes a value of 1 for any observation in that group and 0 for  
 126 any observation not in the group. The penalty matrix for these terms is a  $n_g$  by  $n_g$  identity  
 127 matrix, where  $n_g$  is the number of groups. This means that each group-level coefficient will be  
 128 penalized in proportion to its squared deviation from zero. This is equivalent to how random  
 129 effects are estimated in standard mixed effect models. The penalty term here is proportionate to  
 130 the inverse of the variance of the fixed effect estimated by standard hierarchical model solvers  
 131 (Verbyla et al. 1999 does this contradict what’s above??).



**Figure 3.** Tensor product of depth and salinity with data taken from a 1992 survey of mackerel eggs. The two left plots show the marginal smooths of each term in the model ( $ti(s.depth)$  above and  $ti(salinity)$  below), the right plot shows the interaction effect ( $ti(s.depth, salinity)$ ). Data are from Simon N Wood (2006).

132 This connection between random effects and basis function smooths extends beyond the varying-  
133 intercept case. Any basis-function representation of a smooth function can be transformed so  
134 that it can be represented as a combination of a random effect with an associated variance, and  
135 possibly one or more fixed effects, corresponding to functions in the null space of the original  
136 basis-function (see below). While this is beyond the scope of this paper, see Verbyla et al. (1999)  
137 or Simon N. Wood, Scheipl, and Faraway (2013) for a more detailed discussion on the connections  
138 between these approaches.

#### 139 **Smoothing penalties vs. shrinkage penalties**

140 **does this go above??**

141 **EJP: I think this makes sense here**

142 Penalties can have two effects on how well a model fits: they can penalize how wiggly a given  
143 term is (smoothing) and they can penalize the absolute size of the function (shrinkage). The  
144 penalty can only effect the components of the smooth that have derivatives (the *range space*),  
145 not the other parts (the *nullspace*). For 1-dimensional thin plate regression splines, this means  
146 that there is a linear term left in the model, even when the penalty is in full force (as  $\lambda \rightarrow \infty$ ),  
147 as shown in figure 4 **BLAH**. It is often useful to be able to remove nullspace functions as well,  
148 to be able to shrink them to zero if they do not contribute significantly to a given model fit.  
149 This can be done either by tweaking the penalty matrix so that it both smooths and shrinks  
150 as the single penalty term increases, or by adding a new penalty term that just penalizes the  
151 null space for the model. Figure 4 shows an example of what the basis functions (Fig. 4A),  
152 and smoothing penalties and shrinkage penalties (Fig. 4B) look like for a 6-basis function cubic  
153 spline and for a 6-basis function thin-plate spline. The random effects smoother we discussed  
154 earlier is an example of a pure shrinkage penalty; it penalizes all deviations away from zero, no  
155 matter the pattern of those deviations. This will come into play in section III, when we use  
156 random effects smoothers as one of the components in fitting HGAMs.

157 **DLM: say something about cs basis here**

158 **EDIT this figure**

159 **Should this figure include the intercept basis function as well for each? Right now**  
160 **I've excluded as mgcv drops it automatically, but it is part of the basis, and comes**  
161 **into play with tensor products...**

162 **Do we actually need the cubic spline here? It would simplify the presentation to**  
163 **just show the thin plate spline, and it gets the idea across.**

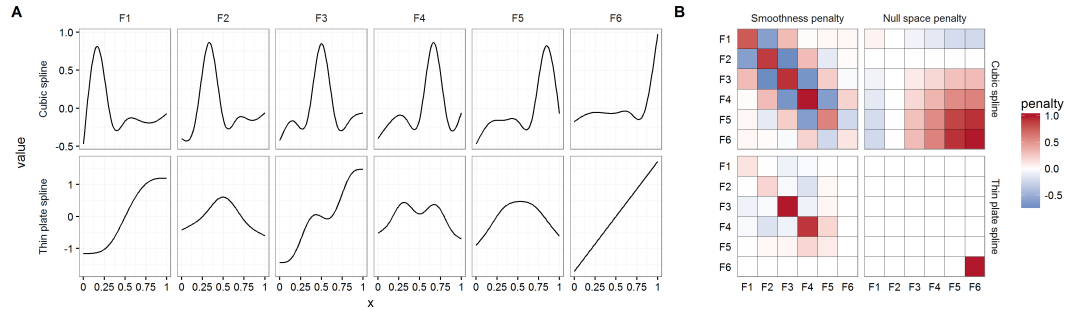
#### 164 **Comparison to hierarchical linear models**

165 Generalized linear mixed effect models (GLMMs; also referred to as hierarchical generalized linear  
166 models, multilevel models etc; e.g., Bolker et al. 2009; Andrew Gelman 2006) are an extension  
167 of regression modelling that allow the modeller to include structure in the data – the structure is  
168 usually of the form of a nesting of the observations. For example individuals are nested within  
169 sample sites, sites are nested within forests and forests within states. The depth of the nesting is  
170 limited by the fitting procedure and number of parameters to estimate.

171 HGLMs are a highly flexible way to think about groupings in the data, the groupings used in  
172 the models often refer to the spatial or temporal scale of the data (McMahon and Diez 2007)  
173 though can be based on any useful grouping.

174 We would like to be able to think about the groupings in our data in a simple way, even when  
175 the covariates in our model are related to the response in a non-linear way. The next section  
176 investigates the extension of the smoothers we showed above to the case where each observation  
177 is in a group, with a group-level smooth.

178 **\*\* this last section should say something like: these are both latent gaussian thingos, the model**  
179 **structure is the difference – we want to be able to heirarchically structure our smoothers\*\***



**Figure 4.** a) Examples of the basis functions associated with a six basis function thin plate spline (top) and a cubic spline (bottom) calculated for  $x$  data spread evenly between  $x = 0$  and  $x = 1$ . Each line represents a single basis function. To generate a given smooth function, each basis function would be multiplied by its own coefficient. b) The smoothing and shrinkage penalty matrices for the thin plate and cubic smoothers shown on left. Red entries indicate positive values and blue indicate negative values. For example, for the thin plate spline, functions f3 and f4 would have the greatest proportionate effect on the total penalty (as they have the largest values on the diagonal), whereas function f6 would not contribute to the penalty at all (all the values in the 6th row and column of the penalty matrix are zero). This means function f6 is in the null space of this basis, and would be treated as completely smooth.

### III: WHAT ARE HIERARCHICAL GAMS?

#### What do we mean by hierarchical smooths?

The smoothers in section II allowed us to model flexible relationships between our response and predictor variables. In this section, we will describe how to model model inter-group variability using smooth curves and how to fit these models in `mgcv`. Model structure is key in this framework, so we start with three choices:

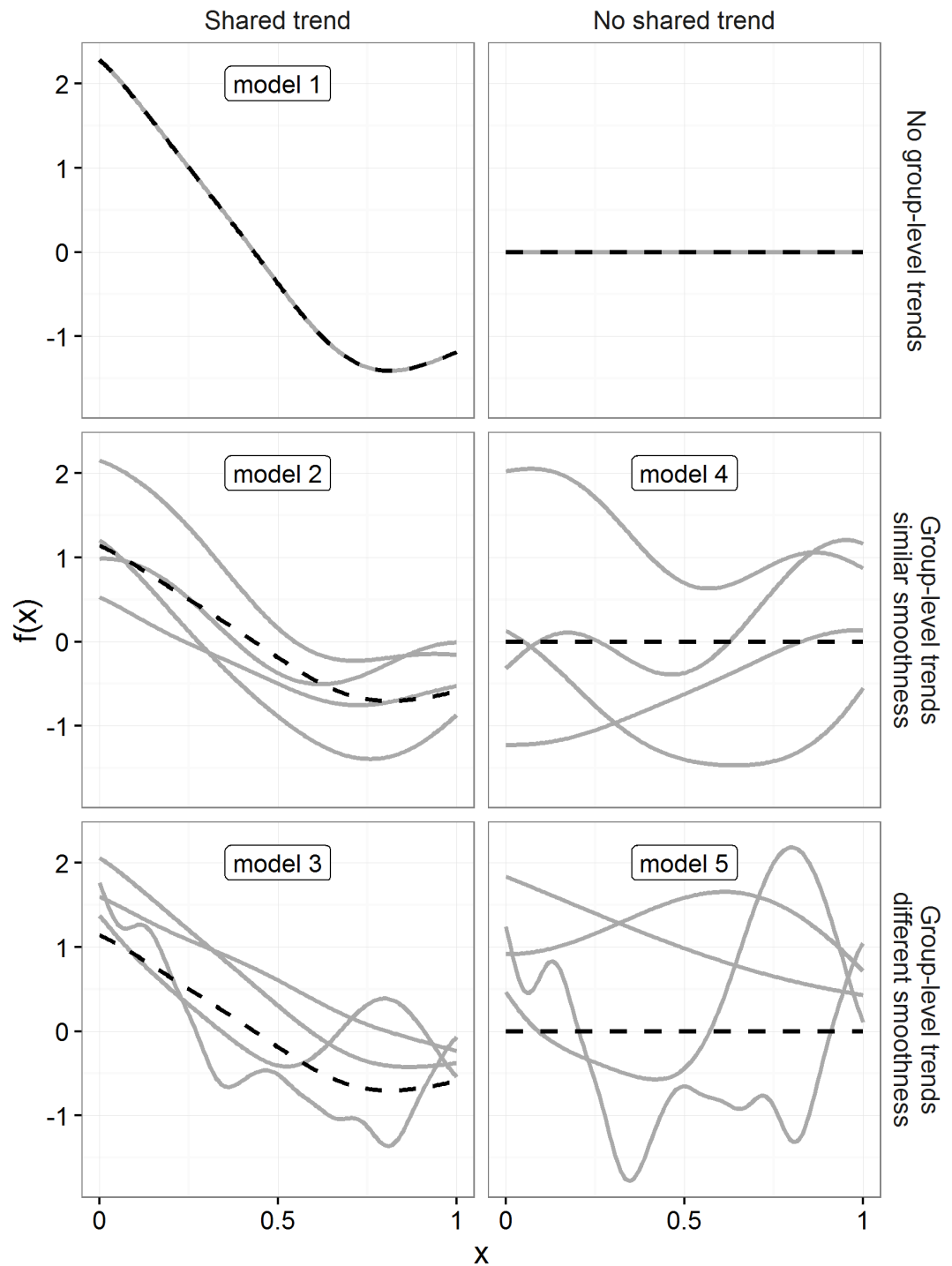
1. Should each group have its own smooth, or will a global smooth term suffice?
2. Do all of the group-specific curves have the same smoothness, or should each group have its own smoothing parameter?
3. Will the smooths for each group have a similar shape to one another – a shared average curve<sup>1</sup>?

These three choices result in five possible models (figure 5), beyond the null model of “no relation between response and predictor(s)”:

1. A single common smooth for all observations.
2. A single common smooth plus group-level smooths that have the same wigglyness.
3. A single common smooth plus group-level smooths with differing wigglyness.
4. Group-specific smooths without an average trend, but with all smooths having the same wigglyness.
5. Group-specific smooths with different wigglyness.

It is important to note that “similar wigglyness” and “similar shape” are distinct ideas; two functions can have very similar wigglyness but very different shapes. Wigglyness simply measures

<sup>1</sup>For this paper, we consider two functions to have similar shape if the average squared distance between the functions is small (assuming the functions have been scaled to have a mean value of zero across their ranges). This definition is somewhat restricted; for instance, a cyclic function would not be considered to have the same shape as a phase-shifted version of that function, nor would two normal distributions with the same mean but different standard deviations. The benefit of this definition of shape, however, is that it is straightforward to translate into quadratic penalties as we have been using.



**Figure 5.** Alternate types of functional variation  $f(x)$  that can be fitted with HGAMs. The dashed line indicates the average function value for all groups, and each solid line indicates the functional value at a given predictor value for an individual group level.



how quickly a function changes across its range, and it is easy to construct two functions that differ in shape but have the same wiggleness. For example, a logistic curve might have the same squared total second derivative between -1 and +1 as a sine curve, but they have very different shapes. Figure 5, model 4 illustrates this case. Similarly, two curves could have very similar overall shape, but differ in their wiggleness. For instance, if one function was equal to the second function plus a high-frequency oscillation. Figure 5 model 3 illustrates this.

We will discuss the trade-offs between different models and guidelines about when each of these models is appropriate in section IV. The remainder of this section will focus on how to specify each of these five models using `mgcv`.

## Coding hierarchical GAMs in R

### **EJP: Going with canned and simulated data for the examples rather than real as it's a bit less messy**

Each of these models can be coded straightforwardly in `mgcv`. To help illustrate this throughout the section when describing how to set these models up, we will refer to the response variable as  $y$ , continuous predictor variables as  $x$  (or  $x_1$  and  $x_2$ , in the case multiple predictors), and `fac` to designate the discrete grouping factor whose variation we are interested in understanding.

We will also use two example datasets to demonstrate how to code these models (see the appendix for code to generate these examples):

A. The `C02` dataset, available in R in the `datasets` package. This data is from an experimental study by Potvin, Lechowicz, and Tardif (1990) of  $\text{CO}_2$  uptake in grasses under varying concentrations of  $\text{CO}_2$ , measuring how concentration-uptake functions varied between plants from two locations (Mississippi and Quebec) and two temperature treatments (chilled and warm). A total of 12 plants were measured, and uptake measured at 7 concentration levels for each plant (figure 6a). Here we will focus on how to use these techniques to estimate inter-plant variation in functional responses.

B. A hypothetical study of what bird movement might look like along a migration corridor, sampled throughout the year. We have simulated this data for this paper (see supplemental code XX). This dataset consists of records of numbers of observed locations of 100 tagged individuals each from six species of bird, at ten locations along a latitudinal gradient, with one observation taken every four weeks. Not every bird was observed at each time point, so counts vary randomly between location and week. The data set (`bird_move`) consists of the variables `count`, `latitude`, `week` and `species` (figure 6b). This example will allow us to demonstrate how to fit these models with interactions and with non-normal (count) data.

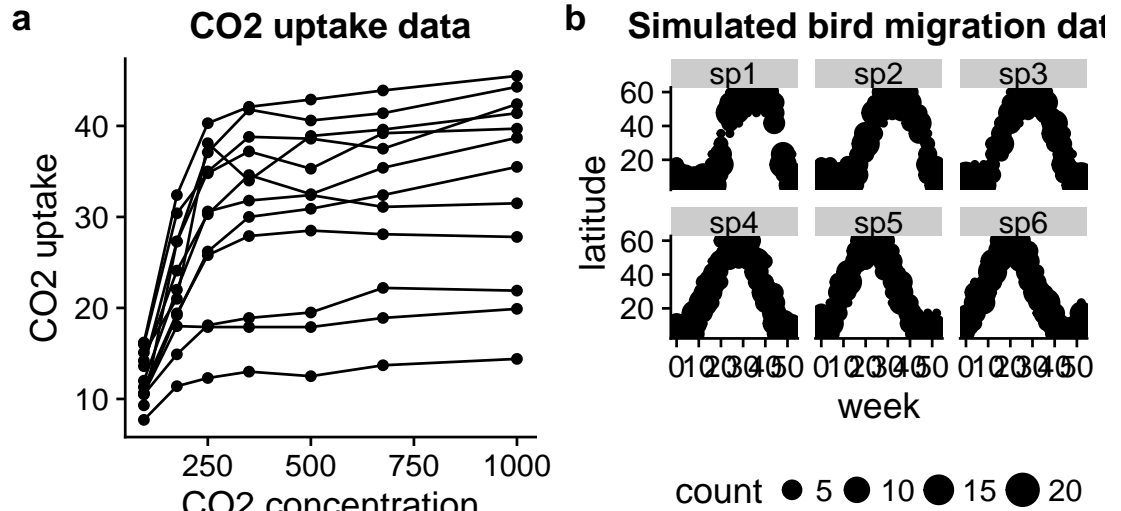
It is important to note that the grouping variable should be coded in R as an unordered factor – a character will raise an error and numeric will lead to a completely different model specification. Whether the factor is ordered or not will not matter for most of the smoothers we use here. However, for models 3&5 order will matter (see below for further details).

Throughout the examples we use Restricted Maximum Likelihood (REML) to estimate model coefficients and smoothing parameters. We strongly recommend using either REML or marginal likelihood (ML) when fitting GAMs for the reasons outlined in (Simon N. Wood 2011).

In each case some data processing and manipulation has been done to obtain the graphics and results below. We recommend readers take a look at the source RMarkdown [CITE] document for this paper to get the full code.

### **A single common smooth for all observations (Model 1)**

We start with the simplest model we can in our framework and include many details here to ensure that readers are comfortable with the terminology and R functions we are going to use later.



**Figure 6.** Example data sets used throughout section III. a) Grass CO<sub>2</sub> uptake versus CO<sub>2</sub> concentration for 12 individual plants (black lines). b) Simulated data set of bird migration, with point size corresponding to weekly counts of 6 species along a latitudinal gradient (zeros excluded for clarity).

For our CO<sub>2</sub> data set, we will model  $\log_e(\text{uptake})$  as a function of two smooths: a thin plate regression spline of  $\log$  concentration, and a random effect for species to model species-specific intercepts.<sup>2</sup> Mathematically:

$$\log_e(\text{uptake}_i) = f(\log_e(\text{conc}_i)) + \zeta_{\text{Plant\_uo}} + \epsilon_i$$

where  $\zeta_{\text{Plant\_uo}}$  is the random effect for plant and  $\epsilon_i$  is a Gaussian error term. We assume that  $\log_e(\text{uptake}_i)$  is normally distributed.

**DLM:** not sure if this note is necessary...

**DLM:** need to justify why we use  $\log$  concentration not just concentration? (could just cite?)

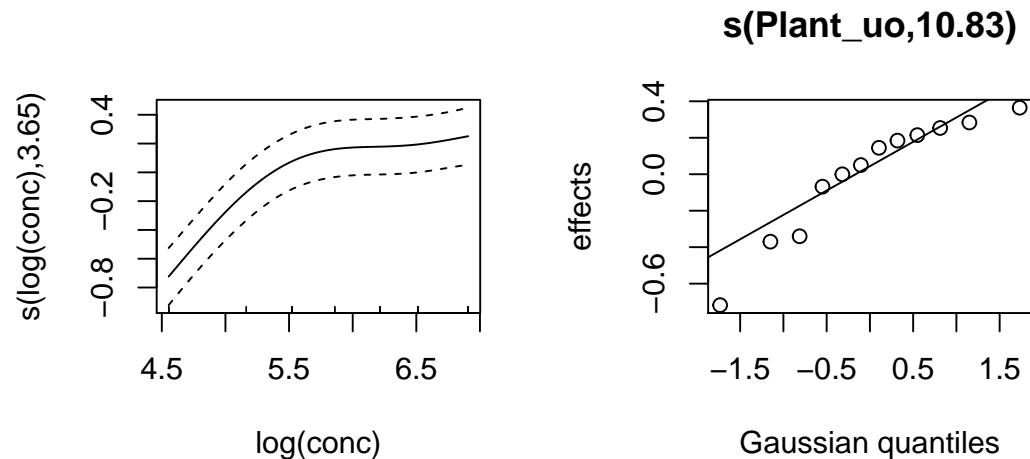
I think it's useful here as it shows one approach to deal with multiplicative functional variation. I've simplified it to remove the note about  $\log$ -transforming concentration though. The reasons for that should be obvious

In R we can write our model as:

```
C02_mod1 <- gam(log(uptake) ~ s(log(conc), k = 5, bs = "tp") +
                    s(Plant_uo, k = 12, bs = "re"),
                    data=C02, method="REML")
```

This is the typical GAM setup, with a single smooth term for each variable. Specifying the model is similar to specifying a `glm` in R, with the addition of `s()` terms to include one-dimensional or isotropic multidimensional smooths. The first argument to `s()` is the terms to be smoothed, the type of smooth to be used for the term is specified by the `bs=...` argument, and the number of basis functions is specified by `k=...`

<sup>2</sup>Note that we're actually modelling  $\ln(\text{uptake})$ ; this can be a useful approach when dealing with estimating multiple functional relationships as it means that functions that differ from each other by a multiplicative constant (so  $f_1(x) = \alpha \cdot f_2(x)$ ) will differ by an additive constant when  $\log$ -transformed (which can be estimated by simple random effects):  $\ln(f_1(x)) = \ln(\alpha) + \ln(f_2(x))$ .



**Figure 7.** mgcv plotting output for model 1 applied to the CO2 dataset.

Figure 7 illustrates `mgcv`'s default plotting out for `CO2_mod1`: the left panel shows the estimated global functional relationship, and the right shows a quantile-quantile plot of the estimates effects vs Gaussian quantiles, which can be used to check our model.

#### DLM: add more to plot description!

Looking at the effects by term is useful but we are often interested in fitted values or predictions our models. This can be useful to construct plots (like those in Figure ??). The next block of code shows how you could plot this to illustrate inter-plant variation in the functional response, plotting untransformed uptake and concentration to make the figure easier to interpret.

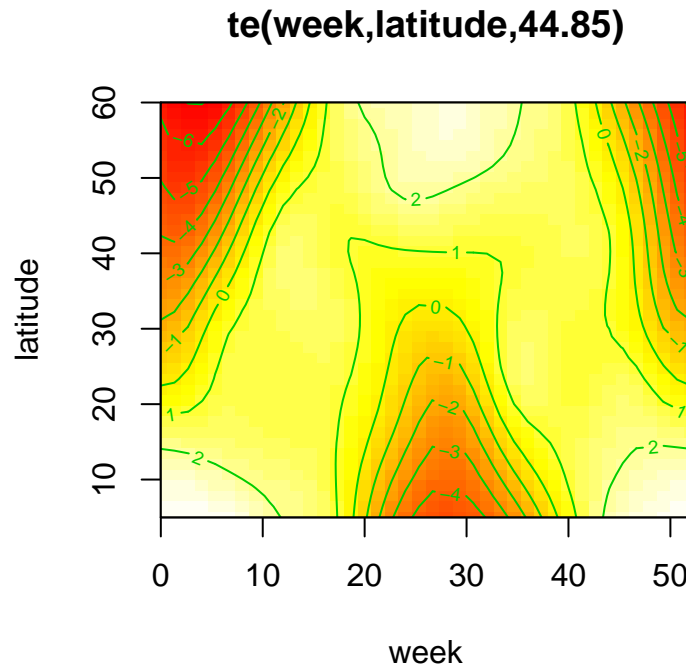
```
“{r co2_mod1_ggplot, fig.width=6, fig.height=3, fig.cap=“\label{fig:co2_mod1_predict}”
# setup prediction data CO2_mod1_pred <- with(CO2, expand.grid(conc=seq(min(conc),
max(conc), length=100), Plant_uo=levels(Plant_uo))) # make the prediction, add this
and a column of standard errors to the prediction # data.frame. Predictions are on the log
scale. CO2_mod1_pred <- cbind(CO2_mod1_pred, predict(CO2_mod1, CO2_mod1_pred,
se.fit=TRUE))
```

#### MAKE THE PLOT

```
ggplot(data=CO2, aes(x=conc, y=uptake, group=Plant_uo)) + facet_wrap(~Plant_uo) +
geom_point() + geom_line(aes(y=exp(fit)), data=CO2_mod1_pred) + geom_ribbon(aes(ymin=exp(fit
- 2se.fit), ymax=exp(fit + 2se.fit), x=conc), data=CO2_mod1_pred, alpha=0.3, in-
herit.aes=FALSE) “
```

We can include interactions in an `s()` term via isotropic smooths such as thin plate regression splines or we can use the tensor product (`te()`) function, if we don't believe the composite terms are isotropic. In this case `bs` and `k` can be specified as a single value (in which case each marginal smooth has the same basis or complexity) or as a vector of basis types or complexities. For example, `y~te(x1,x2, k=c(10,5), bs=c("tp", "cs"))`, would specify a non-isotropic smooth of `x1` and `x2`, with the marginal basis for `x1` being a thin plate regression spline with 10 basis functions, and the smooth of `x2` being a cubic regression spline with a penalty on the null space.

For our bird example, we want to look at the interaction between location and time, so for this we setup the model as:



**Figure 8.** The default plot for this GAM illustrates the average log-abundance of all bird species at each latitude for each week, with yellow colours indicating more individuals and red colours fewer.

$$\text{count}_i = \exp(f(\text{week}_i, \text{latitude}_i))$$

where we assume that  $\text{count}_i \sim \text{Poisson}$ . For the smooth term,  $f$ , we employ a tensor product of latitude and week, using a thin plate regression spline for the marginal latitude effects, and a cyclic cubic spline for the marginal week effect to account for the cyclic nature of weekly effects (we expect week 1 and week 52 to have very similar values), both splines had basis complexity ( $k$ ) of 10. We will also assume the counts of individuals at each location in each week follow a Poisson distribution, and we will ignore species-specific variability.

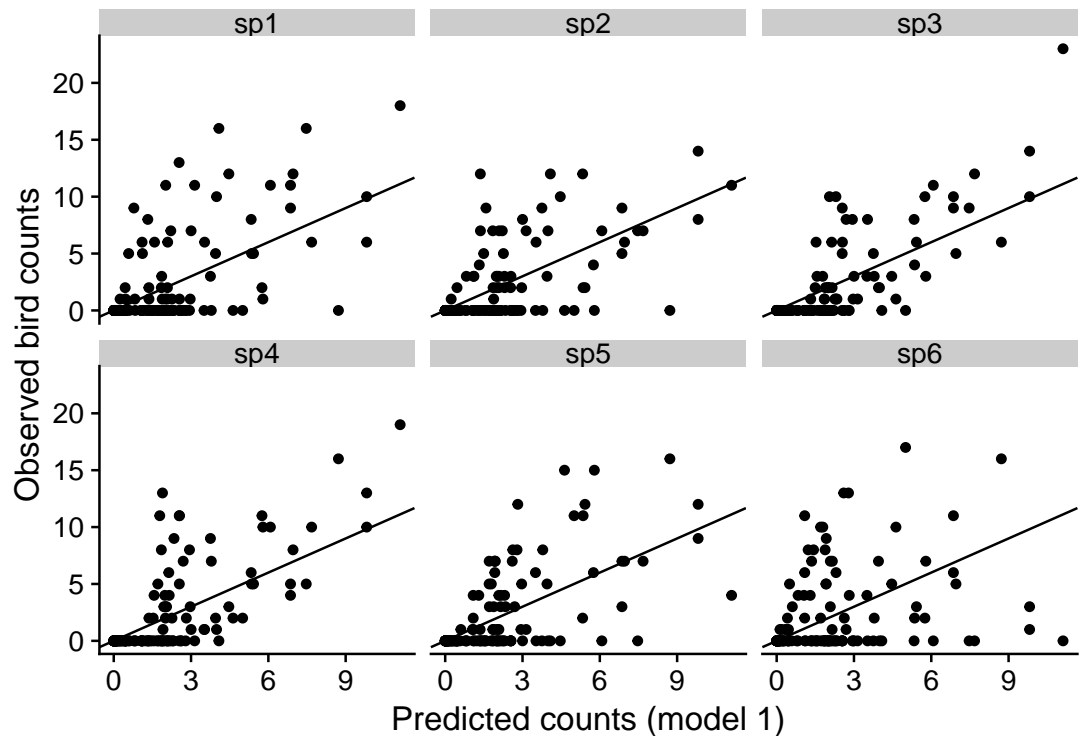
```
library(tidyr)
library(viridis) # for color plotting

bird_move <- read.csv("../data/bird_move.csv") # load data

bird_mod1 <- gam(count ~ te(week, latitude, bs=c("cc", "tp"), k=c(10, 10)),
  data=bird_move, method="REML", family=poisson)

plot(bird_mod1, pages=1, scheme=2, rug=FALSE)
box()
```

Figure 8 shows birds starting at low latitudes in the winter then migrating to high latitudes from the 10th to 20th week, staying there for 15-20 weeks, then migrating back. However, the plot also indicates a large amount of variability in the timing of migration. The source of this variability is apparent when looking at the specifics of migration timing of each species (figure 6b).



**Figure 9.** Observed counts by species versus predicted counts from `bird_mod1` (1-1 line added as reference). If our model fitted well we would expect that all species should show similar dispersions around the 1-1 line. Instead we see that variance around the predicted is much higher for species 1 and 6.

304 All six species in figure 6b) show relatively precise migration patterns, but they differ in the  
 305 timing of when they leave their winter grounds and the time they spend at their summer grounds.  
 306 Averaging over all of this variation results in a relatively imprecise (diffuse) average estimate of  
 307 migration timing (figure 8, 9), and viewing species-specific plots of observed versus predicted  
 308 values (figure 9), it is apparent that the model fits some of the species better than others. This  
 309 model could potentially be improved by adding inter-group variation in migration timing. The  
 310 rest of this section will focus on how to model this type of variation.

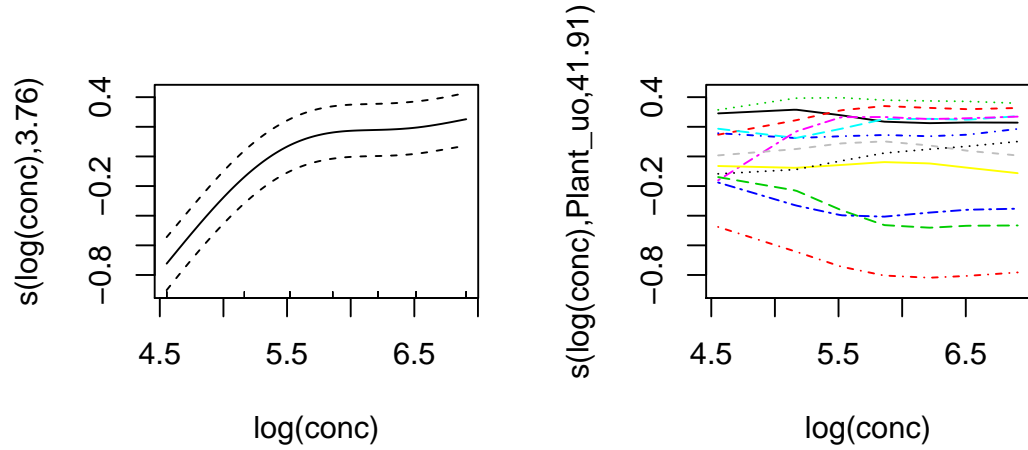
```
bird_move$mod1 = predict(bird_mod1, type="response")

ggplot(bird_move, aes(x=mod1, y= count))+
  facet_wrap(~species)+
  geom_point()+
  geom_abline()+
  labs(x="Predicted counts (model 1)", y= "Observed bird counts")+
  cowplot::theme_cowplot()
```

### 311 **A single common smooth plus group-level smooths that have the same wigglyness (Model 2)**

312 Model 2 is a close analogue to a GLMM with varying slopes: all groups have similar functional  
 313 responses, but allows for inter-group variation in responses. This approach works by allowing  
 314 each grouping level to have its own functional response, but penalizing functions that are too far  
 315 from the average.

316 This can be coded in `mgcv` by explicitly specifying one term for the global smooth (as in model  
 317 1 above) then adding a second smooth term specifying the group level smooth terms, using



**Figure 10.** Global function (left) and group-specific deviations from the global function (right) for CO2\_mod2

a penalty term that tends to draw these group-level smooths to zero. For one-dimensional smooths, `mgcv` provides an explicit basis type to do this, the factor smooth or “fs” basis (see `?smooth.construct.fs.smooth.spec` for detailed notes). This smoother creates a copy of each set of basis functions for each level of the grouping variable, but only estimates one set of smoothing parameters for all groups. The penalty is also set up so each component of its null space is given its own penalty (so that all components of the smooth are penalized towards zero)<sup>3</sup>. As there can be issues of co-linearity between the global smooth term and the group-specific terms (see section V for more details), it is generally necessary to use a smoother with a more restricted null space than the global smooth; for thin plate splines this can be done by setting `m=2` for the global smooth and `m=1` for the group smooth (Baayen et al. 2016) [also cite Wieling paper here]. e.g.: `y~s(x,bs="tp",m=2)+s(x,fac,bs="fs",m=1,xt=list(bs="tp"))`.

We modify our previous CO<sub>2</sub> model as follows:

$$\log_e(\text{uptake}_i) = f(\log_e(\text{conc}_i)) + f_{\text{Plant\_uo}_i}(\log_e(\text{conc}_i)) + \epsilon_i$$

where  $f_{\text{Plant\_uo}_i}(\log_e(\text{conc}_i))$  is the smooth of concentration for the given plant. In R we then have:

```
C02_mod2 <- gam(log(uptake) ~ s(log(conc), k=5, m=2, bs="tp") +
                    s(log(conc), Plant_uo, k=5, bs="fs", m=1),
                    data=C02, method="REML")

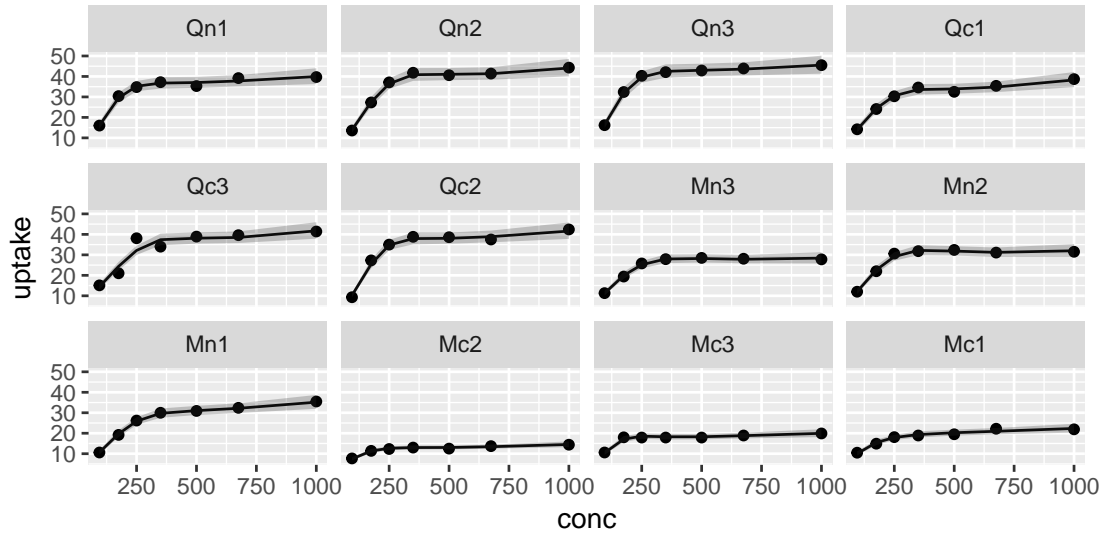
source("../code/functions.R")

C02_mod2 <- gam(log(uptake) ~ s(log(conc), k=5, m=2, bs="tp") +
                    s(log(conc), Plant_uo, k=5, bs="fs", m=1),
                    data=C02, method="REML")

plot(C02_mod2, page=1, seWithMean=TRUE)
```

Figure 10 shows the fitted smoothers for C02\_mod2. The plots of group-specific smooths indicate that plants differ not only in average log-uptake (which would correspond to each plant having a

<sup>3</sup>As part of the penalty construction, each group will also have its own intercept (part of the penalized null space), so there is no need to add a separate term for group specific intercepts as we did in model 1.



**Figure 11.** Predicted uptake values (lines) versus observed uptake for each plant, based on CO2 model 2.

straight line at different levels for the group-level smooth), but differ slightly in the shape of their functional responses. Figure 11 shows how the global and group-specific smooths combine to predict uptake rates for individual plants:

The “fs”-based approach mentioned above does not work for higher-dimensional tensor product smooths (if one is willing to use thin plate regression splines for the multivariate smooth then one can use “fs”). Instead, the group-specific term can be specified with a tensor product of the continuous smooths and a random effect for the grouping parameter. This term will again have a separate set of basis functions for each group, one penalty for the smooth term, and a second penalty drawing all basis functions toward zero<sup>4</sup>. e.g.:  $y \sim \text{te}(x_1, x_2, \text{bs} = \text{"tp"}, m = 2) + \text{te}(x_1, x_2, \text{fac}, \text{bs} = \text{c("tp", "tp", "re")}, m = 1)$ . We illustrate this approach below on the bird migration data.

```
bird_mod2 <- gam(count ~ te(week, latitude, bs=c("cc", "tp"),
                        k=c(10, 10), m=c(2, 2)) +
                te(week, latitude, species, bs=c("cc", "tp", "re"),
                  k=c(10, 10, 6), m=c(1, 1, 1)),
                data=bird_move, method="REML", family=poisson)
```

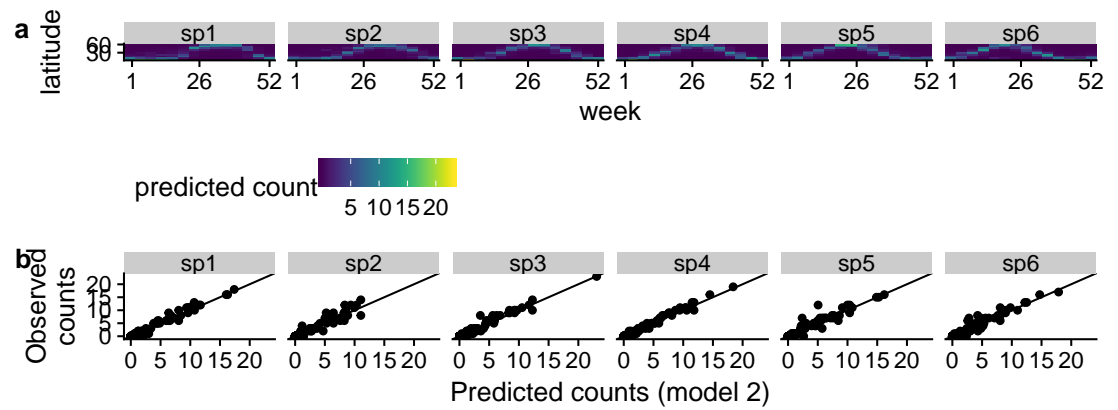
Model 2 is able to effectively capture the observed patterns of interspecific variation in migration behaviour (figure 12a), shows a much tighter fit between observed and predicted values, as well as less evidence of over-dispersion in some species compared to model 1 (figure 12b).

### A single common smooth plus group-level smooths with differing wigglyness (Model 3)

This model class is very similar to model 2, but we now allow each group-specific smooth to have its own smoothing parameter and hence it’s own level of wigglyness. This increases the computational cost of the model, and means that the only information shared between groups is through the global smoothing term. This is useful if different groups differ substantially in how variable they are.

Fitting a separate smooth term (with its own penalties) can be done in mgcv by using the `by=fac` argument in the `s()` function. Therefore, we can code this model as:  $y \sim s(x, \text{bs} = \text{"tp"}) + s(x, \text{by} = \text{fac}, m = 1, \text{bs} = \text{"ts"}) + s(\text{fac}, \text{bs} = \text{"re"})$ . Note two major differences from how model 2

<sup>4</sup>Note that this differs from the “fs” penalty, which assigned one penalty per null space term.



**Figure 12.** a) Predicted migration paths for each species based on `bird_mod2`, with lighter colors corresponding to higher predicted counts. b) Observed counts versus predictions from `bird_mod2`.

was specified: 1., we explicitly include a random effect for the intercept (the `bs="re"` term), as group-specific intercepts are not incorporated into these smooth terms automatically (as would be the case with `bs="fs"` or a tensor product random effect); 2., we explicitly use a basis with a fully penalized null space for the group-level smooth (`bs="ts"`, for “thin plate with shrinkage”), as this method does not automatically penalize the nullspace, so there is potential for co-linearity issues between unpenalized components of the global and group-level smoothers.

Our C02 model is then modified as follows:

```
C02_mod3 = gam(log(uptake) ~ s(log(conc), k=5, m=2, bs="tp") +
                  s(log(conc), by= Plant_uo, k=5, bs="ts", m=1) +
                  s(Plant_uo, bs="re", k=12),
                  data= C02, method="REML")
```

Figure 13 shows a subsample of the group-specific smooths from this model, to prevent crowding. It is apparent from this that some groups (e.g. Qc1) have very similar shapes to the global smooth (differing only in intercept), others do differ from the global trend, with higher uptake at low concentrations and lower uptake at higher concentrations (e.g. Mc1, Qn1), or the reverse pattern (e.g. Mn1).

Using model 3 with higher-dimensional data is also straightforward; `by=fac` terms work as well in tensor-product smooths as they do with isotrophic smooths. We can see this with our bird model:

```
bird_mod3 <- gam(count ~ te(week, latitude, bs=c("cc", "tp"),
                           k=c(10, 10), m=c(2, 2)) +
                  te(week, latitude, bs= c("cc", "tp"),
                           k=c(10, 10), m=c(1, 1), by=species),
                  data=bird_move, method="REML", family=poisson)
```

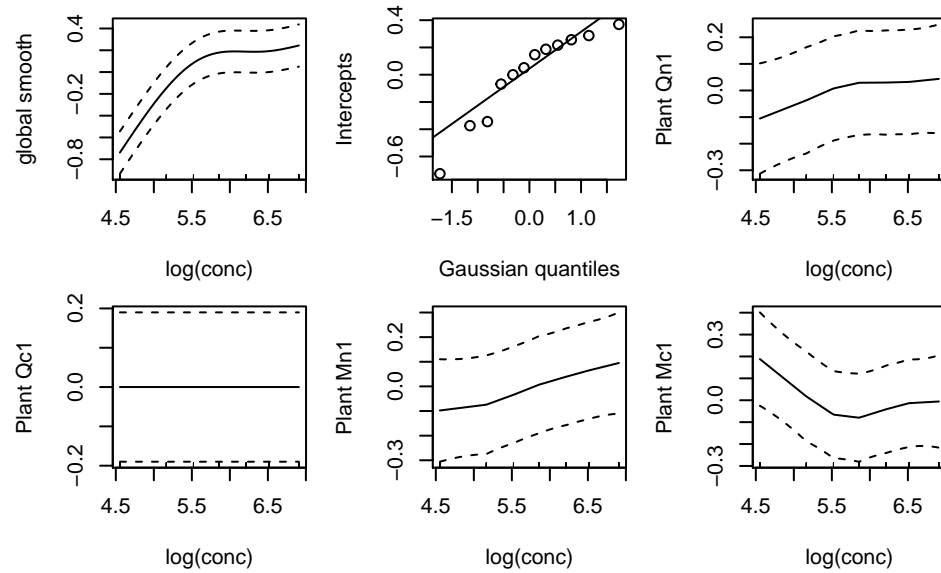
The fitted model for `bird_mod3` is visually indistinguishable from `bird_mod2` (figure 12) so we do not illustrate it here.

#### **Models without global smooth terms (models 4 and 5)**

We can modify the above models to exclude the global term (which is generally faster; see section V). When we don't model the global term, we are allowing each factor to be different, though there may be some similarities in the shape of the functions.

#### **Model 4:**





**Figure 13.** Functional relationships for the CO<sub>2</sub> data estimated for model 3. Top left: the global smooth; Top middle: species-specific random effect intercepts. The remaining plots are a selected subset of the plant-specific smoothers, indicating how the functional response of that plant differs from the global smooth.

Model 4 (shared smooths) is simply model 2 without the global smooth term: `y~s(x,fac,bs="fs")` or `y~te(x1,x2,fac,bs=c("tp","tp","re"))`. This model assumes all groups have the same smoothness, but that the individual shapes of the smooth terms are not related. (Plots are very similar to model 2.)

```
C02_mod4 <- gam(log(uptake) ~ s(log(conc), Plant_uo, k=5, bs="fs", m=2),
  data=C02, method="REML")

bird_mod4 <- gam(count ~ te(week, latitude, species, bs=c("cc", "tp", "re"),
  k=c(10, 10, 6), m=2),
  data=bird_move, method="REML", family=poisson)
```

### Model 5:

Model 5 is simply model 3 without the first term: `y~s(x,by=fac)` or `y~te(x1,x2, by=fac)`. (Plots are very similar to model 3.)

```
C02_mod5 <- gam(log(uptake) ~ s(log(conc), by=Plant_uo, k=5, bs="tp", m=2) +
  s(Plant_uo, bs="re", k=12), data= C02, method="REML")

bird_mod5 <- gam(count ~ te(week,latitude, by=species, bs= c("cc", "tp"),
  k=c(10, 10), m = 2),
  data=bird_move, method="REML", family=poisson)
```

Where group-level smooths are coded using the `by=fac` argument in the `s()` function, ; if the factor is unordered, `mgcv` will set up a model with one smooth for each grouping level. If the factor is ordered, `mgcv` will not set the basis functions for the first grouping level to zero. In model 3 (with an ungrouped smooth included) the ungrouped smooth will then correspond to the first grouping level, rather than the average functional response, and the group-specific smooths will correspond to deviations from the first group. In model 5, using an ordered factor will result in the first group not having a smooth term associated with it at all.

## V: MODELLING ISSUES

Which of the five models should you choose for a given data set? There are two major trade-offs to take into account. The first is the bias-variance trade-off: more complex models can account for more fluctuations in the data, but also tend to give more variable predictions, and can overfit. The second tradeoff is model complexity versus computer time: more complex models can include more potential sources of variation and give more information about a given data set, but will generally take more time and computational resources to fit and debug. We will discuss both of these trade-offs in this section.

### Bias-variance tradeoffs

The bias-variance tradeoff is a fundamental concept in classical statistical analysis. When trying to estimate any value (in the cases we are focusing on, a smooth functional relationship between predictors and data), bias measures how on average an estimate is from the true value of the thing we are trying to estimate, and the variance of an estimator corresponds to how much that estimator would fluctuate if applied to multiple different samples taken from the same population. These two properties tend to be traded off when fitting models; for instance, rather than estimating a population mean from data, we could simply use a fixed value regardless of the observed data. This estimate would have no variance (as it is always the same) but would have high bias unless the true population mean happened to equal zero.<sup>5</sup> The core insight into why penalization is useful is that the penalty term slightly increases the bias but can substantially decrease the variance of an estimator, relative to its unpenalized version (Efron and Morris 1977).

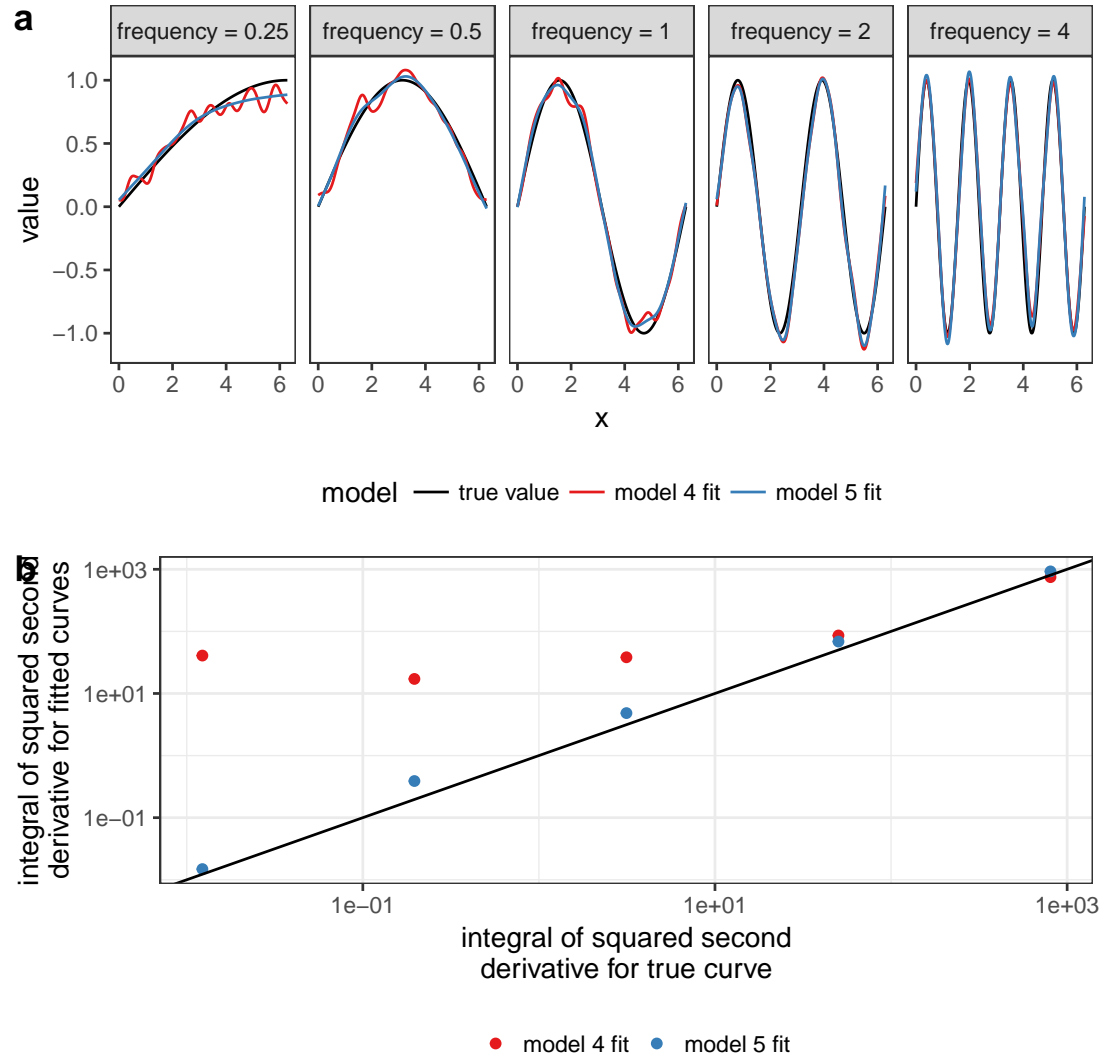
In GAMs and HGLMs, the bias-variance tradeoff is managed by the penalty terms (random effect variances in HGLM terminology). Larger penalties correspond to lower variance, as the estimated function is unable to wiggle a great deal, but also correspond to higher bias unless the true function is close to the null space for a given smoother (e.g. a straight line for thin-plate splines with 2nd derivative penalties, or zero for a standard random effect). The computational machinery used by mgcv to fit smooth terms is designed to find penalty terms that best trade off bias for variance to find a smooth that can effectively to predict new data.

The bias-variance tradeoff comes into play with HGAMs when choosing whether to fit separate penalties for each group level or assign a common penalty for all group levels (i.e. deciding between models 2 & 3 or models 4 & 5). If the functional relationships we are trying to estimate for different group levels actually vary in how wiggly they are, setting the penalty for all group-level smooths equal (models 2&4) will either lead to overly variable estimates for the least variable group levels, overly smoothed (biased) estimates for the most wiggly terms, or a mixture of these two, depending on how the fitting criteria used (ML, REML, or GCV) determines where the optimal smoothing parameter should be set.

We developed a simple numerical experiment to determine whether mgcv fitting criteria tend to set estimated smoothness penalties high or low in the presence of among-group variability in smoothness. We simulated data from five different groups, with all groups having the same levels of the covariate  $x$ , ranging from 0 to  $2\pi$ . For each group, the true function relating  $x$  to  $y$  was a sine wave, but the frequency varied from 0.25 (equal to half a cycle across the range of  $x$ ) to 4 (corresponding to 4 full cycles across the range). We added normally distributed error to all  $y$ -values, with a standard deviation of 0.2. We then fit both model 4 (where all curves were assumed to be equally smooth) and model 5 (with varying smoothness) to the entire data set, using REML criteria to estimate penalties. For this example (Fig. 14a), requiring equal smoothness for all group levels resulted mgcv underestimating the penalty for the lowest frequency (most smooth) terms, but accurately estimating the true smoothness of the highest frequency terms as measured by the squared second derivative of the smooth fit versus that of the true function (Fig. 14b).

---

<sup>5</sup>While this example may seem contrived, this is exactly what happens when we assume a given fixed effect is equal to zero (and thus exclude it from a model).



**Figure 14.** plotting example

441 This implies that assuming equal smoothness will tend to lead to underestimating the true  
 442 smoothness of low-variability terms, and thus leading to more variable estimates of these terms.  
 443 If this is a potential issue, we recommend fitting both models and using the model evaluation  
 444 techniques discussed in Section IV to determine if there is evidence for among-group variability  
 445 in smoothness. For instance, the AIC for model 4 fit to this data is -180, whereas it is -211,  
 446 implying a substantial improvement in fit by allowing smoothness to vary. However, it may  
 447 be the case that there is too few data points per group to estimate separate smoothness levels,  
 448 in which case model 2 or model 4 may still be the better option even in the face of varying  
 449 smoothness.

450 The ideal case would be to assume that among group penalties follow their own distribution  
 451 (estimated from the data), to allow variation in smoothness while still getting the benefit of  
 452 pooling information on smoothness between groups. However, this is currently not implemented  
 453 in mgcv (and would be difficult to set up via mgcv's method of structuring penalties). It is  
 454 possible to set up this type of varying penalty model in flexible Bayesian modelling software such  
 455 as Stan or INLA (see below for a discussion of these tools), but how to set up this type of model  
 456 has not been well studied, and is beyond the scope of this paper.

It may seem like there is also a bias–variance tradeoff between choosing to use a single global smoother (model 1) or a global smoother plus group-level terms (model 2-3), as in model 1, all the data is used to estimate a single smooth term, and thus should have lower variance than models 2-3, but higher bias for any given group in the presence of inter-group functional variability. However, in practice, this tradeoff will already be handled by `mgcv` via estimating penalties; if there are no average differences between functional responses, `mgcv` will penalize the group-specific functions toward zero, and thus toward the global model. The choice between using model 1 versus models 2-3 should generally be driven by computational costs; model 1 is typically much faster to fit than models 2-3, even in the absence of among-group differences, so if there is no need to estimate inter-group variability, model 1 will typically be more efficient.

A similar issue exists when choosing between models 2/3 and 4/5; if all group levels have very different functional shapes, the global term will get penalized toward zero in models 2/3, so they will reduce to models 4/5. Again, the choice to include a global term or not should be made based on scientific considerations (is the global term of interest to estimate) and computational considerations (which we will discuss next).

## Complexity – computation tradeoffs

GAMs and GLMMs have substantially increased the range of flexible models available to the average researcher, and the HGAM models we discussed in section III extend on this broad base. However, the more flexible a model is, the larger an effective parameter space any fitting software has to search to find parameters that can predict the observed data. While numerical algorithms for solving complex models are always improving, it can still be surprisingly easy to use up massive computational resources trying to fit a model to even relatively small datasets. While we typically want to choose a model based on model fit (see above and section IV) and our goals for what the model will be used for, computing resources can often act as an effective upper limit on possible model complexity. Fitting an HGAM means adding extra computational complexity on top of either a GAM model with only global terms or a GLMM without smooth terms. For a given data set (with a fixed number `n` data points) and assuming a fixed family and link function, the time it takes to compute a given HGAM will depend, roughly, on four factors: the number of basis functions to be estimated, the number of smooth penalties to be estimated, whether the model needs to estimate both a global smooth and groupwise smooths, and the algorithm used to estimate parameters and fitting criteria used.

The most straightforward factor that will affect the amount of computational resources is the number of parameters in the model. Adding group-level smooths (moving from model 1 to 2-5) means that there will be more regression parameters to estimate, since each grouping level needs a separate coefficient for each basis function in the smooth. For a dataset with `g` different groups and `n` data points, fitting a model will just a global smooth,  $y \sim s(x, k=k)$  will require only `k` coefficients, and takes  $\mathcal{O}(nk^2)$  operations<sup>6</sup> to evaluate, but fitting the same data using a group-level smooth (model 4,  $y \sim s(x, fac, bs="fs", k=k)$ ) will require  $\mathcal{O}(nk^2g^2)$  operations to evaluate; in effect, adding a group-level smooth will increase computational time by an order of the number of groups squared<sup>7</sup>. The effect of this is visible in the examples we fit in section III when comparing the number of coefficients and relative time it takes to compute model 1 versus the other models (Table 1). One way to deal with this issue would be to reduce the number of basis functions (`k`) used when fitting group-level smooths when the number of groups is large; in effect, this would increase the flexibility of the model to accommodate inter-group differences, while reducing its ability to model variance within any given group. It can also make sense to use more computationally efficient basis functions when fitting large data sets, such as p-splines or

<sup>6</sup>To understand the effects of these terms, we will use “big-O” notation; when we say a given computation is of order  $\mathcal{O}(n \log n)$ , it means that, for that computation, as  $n$  gets large, the amount of time the computation will take will grow proportionally to  $n \log n$ , so more quickly than linearly with  $n$ , but not as fast as  $n$  squared.

<sup>7</sup>Including a global smooth (models 2-3) or not (models 4-5) will not generally substantially affect the number of coefficients needed to estimate (compare the number of coefficients in Table 1, model 2 vs. model 4, or model 3 versus model 5). Adding a global term will only add at most `k` extra terms, and it actually ends up being less than that, as `mgcv` drops basis functions from co-linear smooths to ensure that the model matrix is full rank.

**Table 1.** Relative computational time and model complexity for different HGAM formulations of the two example data sets from section III. All times are scaled relative to the length of time model 1 takes to fit to that data set. The # of coefficients measures the total number of model parameters (including intercepts). The # of smooths is the total number of unique penalty values estimated by the model.

model	relative time	# of terms	
		coefficients	penalties
A. CO2 data			
1	1	17	2
2	5	65	3
3	14	65	14
4	4	61	3
5	9	61	13
B. bird movement data			
1	1	90	2
2	110	540	5
3	140	624	14
4	100	541	3
5	66	535	12

cubic splines, rather than thin-plate splines, as thin-plate splines can take a substantial amount of overhead to compute the actual basis functions to use [CITE].

Adding additional smoothing parameters (moving from model 2 to model 3, or moving from model 4 to 5) is even more costly than increasing the number of coefficients to estimate, as estimating smoothing parameters is computationally intensive (Simon N. Wood 2011). This means that models 2 and 4 will generally be substantially faster than 3 and 5 when the number of groups is reasonably large, as models 3 and 5 fit a separate set of penalties for each group level. The effect of this is visible in comparing the time it takes to fit model 2 to model 3 (which has a smooth for each group) or models 4 and 5 for the example data (Table 1). Note that this will not hold for every model, though; for instance, model 5 takes less time to fit the bird movement data than model 4 does (Table 1B).

#### Alternative formulations: bam, gamm, and gamm4 (with a brief foray into Bayes)

When fitting models with large numbers of different group levels, it is often possible to speed up computation substantially by using one of the alternative fitting algorithms available through mgcv.

The first tool available, that requires the least changes to your code compared to the base `gam` function, is the `bam` function. This function is designed to improve performance when fitting data sets with large amounts of data. It uses two tools to do this. First, it saves on the amount of memory needed to compute a given model by using a random subset of the data to calculate the basis functions for the smoothers, then breaking the data up into blocks and updating model fit within each block (Simon N. Wood, Goude, and Shaw 2015). While this is primarily designed to reduce the amount of memory needed to fit these models, it can also substantially reduce computation time. Second, the `bam` function, when fitting using its default “fREML” (for “Fast REML”) method, you can use the `discrete` when fitting the model. This option causes `bam` to simplify each covariate to a set of discrete levels (instead of a continuous range), substantially reducing the amount of computation needed. Setting “discrete = TRUE” lets `bam` estimate the number of bins to use for each covariate. It is also possible to manually specify the number of bins by passing `discrete` a vector of values. See `?mgcv::bam` for more details.

It also takes more computational overhead compared to `gam` to set a `bam` model up, so for small numbers of groups, it can actually be slower than `gam` (Figure 15), however, as the number

of groups increases, computational time for `bam` increases more slowly than for `gam`; in our simulation tests, when the number of groups is greater than 16, `bam` can be upward of an order of magnitude faster (Figure 15). Note that `bam` can be somewhat less computationally stable when estimating these models (i.e. less likely to converge) so it does typically make sense to still use `gam` for smaller data sets.

The second option is to fit these models using one of two dedicated mixed effect model estimation packages, `nlme` and `lme4`. The `mgcv` package includes the function `gamm` that allows you to call `nlme` to solve a given GAM, automatically handling the transformation of smooth terms into random effects (and back into basis-function representations for plotting and other statistical analyses). To use `lme4`, you will have to install the `gamm4` package, and use the `gamm4` function from this package. Using `gamm` or `gamm4` to fit models rather than `gam` can substantially speed up computation when the number of groups is large, as both `nlme` and `lme4` take advantage of the sparse structure of the random effects, where most basis functions will be zero for most groups (i.e. any group-specific basis function will only take a non-zero value for observations in that group level). As with `bam`, `gamm` and `gamm4` are generally slower than `gam` for fitting HGAMs when the number of group levels is small (in our simulations, <8 group levels), however they do show substantial speed improvements even with a moderate number of groups, and were as fast as or faster to calculate than `bam` for all numbers of grouping levels we tested (Figure 15)<sup>8</sup>.

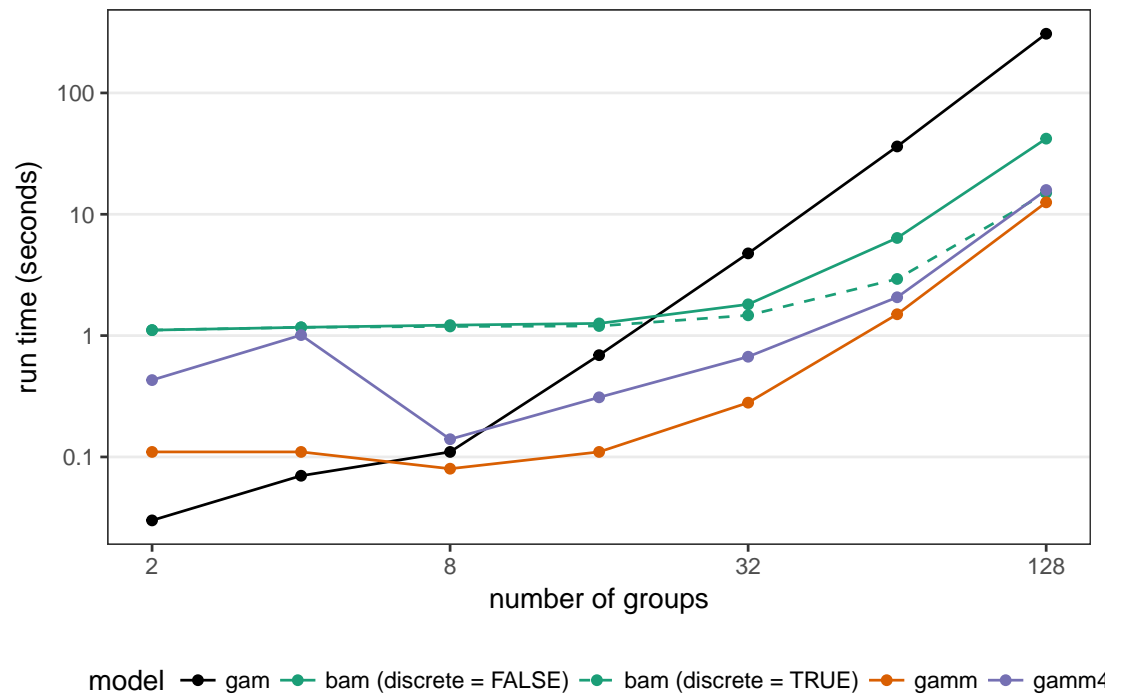
Setting up models 1-5 in `bam` uses the same code as we have previously covered; the only difference is that you use the `bam` instead of `gam` function, and have the additional option of discretizing your covariates. The advantage of this approach is that `bam` allows you to use almost all of the same families available to the `gam` function, and `bam` model output can be evaluated using the same functions (e.g. `summary`, `AIC`, `plot`, etc.) so it is simple to substitute for `gam` if you need to speed a model up.

Both `gamm` and `gamm4` require at least a few changes to how you code models. First, there are a few limitations on how you are able to specify models 1-5 in both frameworks. Factor smooth (`bs="fs"`) basis setups work in both `gamm` and `gamm4`. However, as the `nlme` package does not support crossed random effects, it is not possible to have two “fs” terms for the same grouping variable in `gamm` models (e.g. `y~s(x1, grp,bs="fs")+s(x2, grp, bs="fs")`). These type of crossed random effects are allowed in `gamm4`. The use of `te` and `ti` terms are not possible in `gamm4` however, due to issues with how random effects are specified in the `lme4` package, making it impossible to code models where multiple penalties apply to a single basis function. Instead, for multidimensional group-level smooths, the alternate function `t2` needs to be used to generate these terms, as it creates tensor products with only a single penalty for each basis function (see `?mgcv::t2` for details on these smoothers, and Simon N. Wood, Scheipl, and Faraway (2013) for the theoretical basis behind this type of tensor product). So for instance, model 2 for the bird movement data we discussed in section III would need to be coded as:

```
bird_mod4_gamm4 = gamm4(count ~ t2(week,latitude,species, bs= c("cc", "tp","re"),
                                k=c(10,10,6),m = 2),
                        data= bird_move, family= poisson)
```

These packages also do not support the same range of families for the dependent variable; `gamm` only supports non-Gaussian families by using a fitting method called penalized quasi-likelihood (PQL) that is slower and not as numerically stable as the methods used in `gam`, `bam`, and `gamm4`. Non-Gaussian families are well supported by `lme4` (and thus `gamm4`), but can only fit them using marginal likelihood (ML) rather than REML, so may tend to over-smooth relative to `gam` using REML estimation. Further, neither `gamm` nor `gamm4` supports several of the extended families

<sup>8</sup>It is also possible to speed up both `gam` and `bam` by using multiple processors in parallel, whereas this is not currently possible for `gamm` and `gamm4`. For large numbers of grouping levels, this should speed up computation as well, at the cost of using more memory. However, computation time will likely not decline linearly with the number of cores used, since not all model fitting sets are parallizable, and performance of the cores can vary. As parallel processing can be complicated and dependent on the type of computer you are using to configure properly, we do not go into how to use these methods here. The help file `?mgcv::mgcv.parallel` explains how to use parallel computations for `gam` and `bam` in detail.



**Figure 15.** Elapsed time to estimate the same model using each of the four approaches. Each data set was generated with 20 observations per group using a unimodal global function and random group-specific functions consisting of an intercept, a quadratic term, and logistic trend for each group. Observation error was normally distributed. Models were fit using model 2:  $y \sim s(x, k=10, bs='cp') + s(x, fac, k=10, bs='fs', xt=list(bs='cp'), m=1)$ . All models were run on a single core.

579 available through `mgcv`, such as zero-inflated, negative binomial, or ordered categorical and  
580 multinomial distributions.

## 581 Estimation issues when fitting both global and groupwise smooths

582 When fitting models with separate global and groupwise smooths (models 2 and 3), one issue to  
583 be aware of is concurvity between the global smooth and groupwise terms. Concurvity measures  
584 how well one smooth term can be approximated by some combination of the other smooth terms  
585 in the model (see `?mgcv::concurvity` for details). For models 2 and 3, the global term is entirely  
586 concure with the groupwise smooths. This is because, in the absence of the global smooth term,  
587 it would be possible to recreate that average effect by shifting all the groupwise smooths so they  
588 were centered around the global mean. In practical terms, this has the consequence of increasing  
589 uncertainty around the global mean relative to a model with only a global smooth. In some  
590 cases, it can result in the estimated global smooth being close to flat, even in simulated examples  
591 with a known strong global effect. This concurvity issue may also increase the time it takes  
592 to fit these models (for example, compare the time it takes to fit models 3 and 5 in Table 1).  
593 That these models can still be estimated is because of the penalty terms; all of the methods we  
594 have discussed for fitting model 2 (“fs” terms or random effect tensor products) automatically  
595 create a penalty for the null space of the group-level terms, so that only the global term has  
596 its own unpenalized null space, and both the REML and ML criteria work to balance penalties  
597 between nested smooth terms (this is why nested random effects can be fitted). We have noted,  
598 however, that `mgcv` still occasionally finds degenerate solutions with simulated data where the  
599 fitted global term ends up over-smoothed.

600 What we recommend to avoid this issue is to use a combination of smoother choice and setting  
601 model degrees of freedom so that the groupwise terms are either slightly less flexible or have  
602 a smaller null space. For instance, in the examples in section III, we used smoothers with an  
603 unpenalized null space (standard thin-plate splines) for the global smooth and ones with no null  
604 space for the groupwise terms<sup>9</sup>. When using thin-plate splines, it may also help to use splines  
605 with a lower order of derivative penalized in the groupwise smooths than the global smooths,  
606 as lower-order “tp” splines have fewer basis functions in the null space. For example, we used  
607 `m=2` (penalizing squared second derivatives) for the global smooth, and `m=1` (penalizing squared  
608 first derivatives) for groupwise smooths in models 2 and 3. Another option would be to use a  
609 lower number of basis functions (`k`) for groupwise relative to global terms, as this will reduce  
610 the maximum flexibility possible in the groupwise terms. We do caution that these are just  
611 rules of thumb. As of this writing, there is no published work looking what the effect of adding  
612 groupwise smooths has on the statistical properties of estimating a global smooth. In cases  
613 where an accurately estimated global smooth is essential, we recommend either fitting model 1,  
614 or using Markov Random Fields (Appendix A) and calculate the global smooth by averaging  
615 across grouping levels.

## 616 A brief foray into the land of Bayes

## 617 VI: EXAMPLES

618 *EJP: I think we should aim for 3-4 good examples here, highlighting different as-*  
619 *pects of the model fitting problem. The example I’ll be showing is a set of zooplank-*  
620 *ton community time series data, where multiple species were tracked throughout the*  
621 *year for a period of roughly 20 years. This example can highlight both testing mod-*  
622 *els 4/5 (for comparing different species’ season cycles) and models 2/3 (for testing*  
623 *for differences between years for a single species). I think we need at least one*  
624 *example showing how to use these methods for multivariate regression (e.g. spa-*  
625 *tial analysis), and potentially an example showing how to these models work for*

<sup>9</sup>For model 2, the “fs” smoother, and tensor products of random effect (“re”) and other smooth terms do not have a penalized null space by construction (they are full rank), as noted above. For model 3 groupwise terms, we used basis types that had a penalty added to the null space: `bs=“tp”`, `“cs”`, or `“ps”` have this property.



626 *non-normal data, and for including other covariates. In all examples, I think we*  
627 *should focus on how to fit each data set, visualize models, and compare different*  
628 *model fits.*

629 In this final section, we will go through a few example analyses, to highlight how to use these  
630 models in practice, and to illustrate how to fit, test, and visualize each model.

### 631 **Example 1: Inter- and intra-specific variation in zooplankton seasonal cycles over time**

632 This first example will demonstrate how to use these models to fit community data, to show  
633 when using a global trend may or may not be justified, and to illustrate how to use these  
634 models to fit seasonal time series. Here, we are using data from the Wisconsin Department  
635 of Natural Resources collected by Richard Lathrop from a chain of lakes (Mendota, Menona,  
636 Kegonsa, and Waubesa) in Wisconsin, to study long-term patterns in the seasonal dynamics of  
637 zooplankton. This data consists of roughly bi-weekly samples (during open-water conditions) of  
638 the zooplankton communities, taken from the deepest point of each lake via vertical tow collected  
639 every year from 1976 to 1994 (the collection and processing of this data is fully described in  
640 Lathrop (2000)). We will use this data estimate variability in seasonality among species in the  
641 community, and between lakes for the most abundant taxon in the sample (*Daphnia mendotae*).  
642 As we are focusing on seasonal cycles rather than average or maximum abundances, we have scaled  
643 all densities by log-transforming them then scaling by the within year species- and lake-specific  
644 mean and standard deviation (so all species in all lake-years will have a mean scaled density of  
645 zero and standard deviation of one).

646 This is what the data looks like:

```
647 zooplankton = read.csv("../data/zooplankton_example.csv")
648
649 str(zooplankton)
650
651 ## 'data.frame':    5848 obs. of  6 variables:
652 ## $ day           : int  10 10 10 10 10 10 10 10 12 12 ...
653 ## $ year          : int  1980 1980 1980 1980 1980 1980 1980 1980 1984 1984 ...
654 ## $ lake          : Factor w/ 4 levels "Kegonsa","Mendota",...: 2 2 2 2 2 2 2 2 1 1 ...
655 ## $ taxon         : Factor w/ 8 levels "Calanoida copepodites",...: 1 2 3 4 5 6 7 8 1 2 .
656 ## $ density       : num  5000 28000 3000 6000 1867000 ...
657 ## $ density_scaled: num  -1.4173 -0.0714 -2.3837 0.1945 1.2613 ...
658
659 levels(zooplankton$taxon)
660
661 ## [1] "Calanoida copepodites" "Chydorus sphaericus"
662 ## [3] "Cyclopoida copepodites" "Daphnia mendotae"
663 ## [5] "Diacyclops thomasi"    "Keratella cochlearis"
664 ## [7] "Leptodiaptomus siciloides" "Mesocyclops edax"
665
666 levels(zooplankton$lake)
667
668 ## [1] "Kegonsa" "Mendota" "Menona" "Waubesa"
```

659 We will split the data into testing and training sets, so we can evaluate how well our models fit  
660 out of sample. As there are multiple years of data here, we will use data from the even years to  
661 fit (train) models, and that from the odd years to test the fit:

```
662 library(mgcv)
663 zoo_train = subset(zooplankton, year%%2==0)
664 #the modulus (%%) finds the remainder after division by the right.
665 #here we use it to find even numbers
666
667 zoo_test = subset(zooplankton, year%%2==1)
```

Our first exercise here will be to demonstrate how to model community-level variability in seasonality, by regressing scaled density on day of year, with species-specific curves. As we are not interested here in average seasonal dynamics, we will focus on models 4&5<sup>10</sup>. As this is seasonal data, we will use cyclic smoothers as the basis for seasonal dynamics.

```

zoo_comm_mod4 = gam(density_scaled~s(day, taxon, bs="fs",k=10,xt=list(bs="cc")),
  data=zoo_train,
  #we need to specify the start and end knots for day
  knots = list(day =c(1,365)),
  #We'll use ML as we are comparing models that differ in fixed effects
  method = "ML"
)

summary(zoo_comm_mod4)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## density_scaled ~ s(day, taxon, bs = "fs", k = 10, xt = list(bs = "cc"))
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.02597    0.02621  -0.991   0.322
##
## Approximate significance of smooth terms:
##             edf Ref.df    F p-value
## s(day,taxon) 54.52     71 18.44 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.308   Deviance explained = 32.1%
## -ML = 3615.4   Scale est. = 0.64168    n = 2947

# as all of the model features except the formula are the same in model 4&5,
# we just use the update function to refit the model with the new formula
zoo_comm_mod5 = update(zoo_comm_mod4,
  formula = density_scaled~s(day, by=taxon, k=10,bs="cc"))

summary(zoo_comm_mod5)

##
## Family: gaussian
## Link function: identity
##
## Formula:
## density_scaled ~ s(day, by = taxon, k = 10, bs = "cc")
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.001563   0.014736  -0.106   0.916
##
## Approximate significance of smooth terms:

```

<sup>10</sup>Here we are focusing on only the most common species in the data set. If we wanted to estimate the seasonal dynamics for rarer species, adding a global smooth term might be useful, so we could borrow information from the more common species.

```

697 ##
698 ## s(day):taxonCalanoida copepodites      6.968      8 44.248 < 2e-16 ***
699 ## s(day):taxonChydorus sphaericus      5.595      8  8.933 2.47e-15 ***
700 ## s(day):taxonCyclopoida copepodites    5.806      8 21.013 < 2e-16 ***
701 ## s(day):taxonDaphnia mendotae         6.941      8 15.977 < 2e-16 ***
702 ## s(day):taxonDiacyclops thomasi        6.663      8 38.303 < 2e-16 ***
703 ## s(day):taxonKeratella cochlearis     3.904      8  3.527 1.35e-06 ***
704 ## s(day):taxonLeptodiaptomus siciloides 6.223      8  5.891 3.36e-09 ***
705 ## s(day):taxonMesocyclops edax          5.137      8 27.439 < 2e-16 ***
706 ## ---
707 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
708 ##
709 ## R-sq.(adj) =  0.31   Deviance explained = 32.1%
710 ## -ML = 3601.5   Scale est. = 0.63978   n = 2947

```

We can see that both models have very similar fits, with an adjusted  $R^2$  of 0.308 for model 4 and 0.31 for model 5. Model 5 has a somewhat lower AIC (AIC(zoo\_comm\_mod4) = 7115, AIC(zoo\_comm\_mod5) = 7104), implying a better overall fit. However, the two models show very similar fit to the data:

```

library(ggplot2)
library(dplyr)

#Create synthetic data to use to compare predictions
zoo_plot_data = expand.grid(day = 1:365, taxon = factor(levels(zoo_train$taxon)))

zoo_mod4_fit = predict(zoo_comm_mod4, zoo_plot_data, se.fit = T)
zoo_mod5_fit = predict(zoo_comm_mod5, zoo_plot_data, se.fit = T)

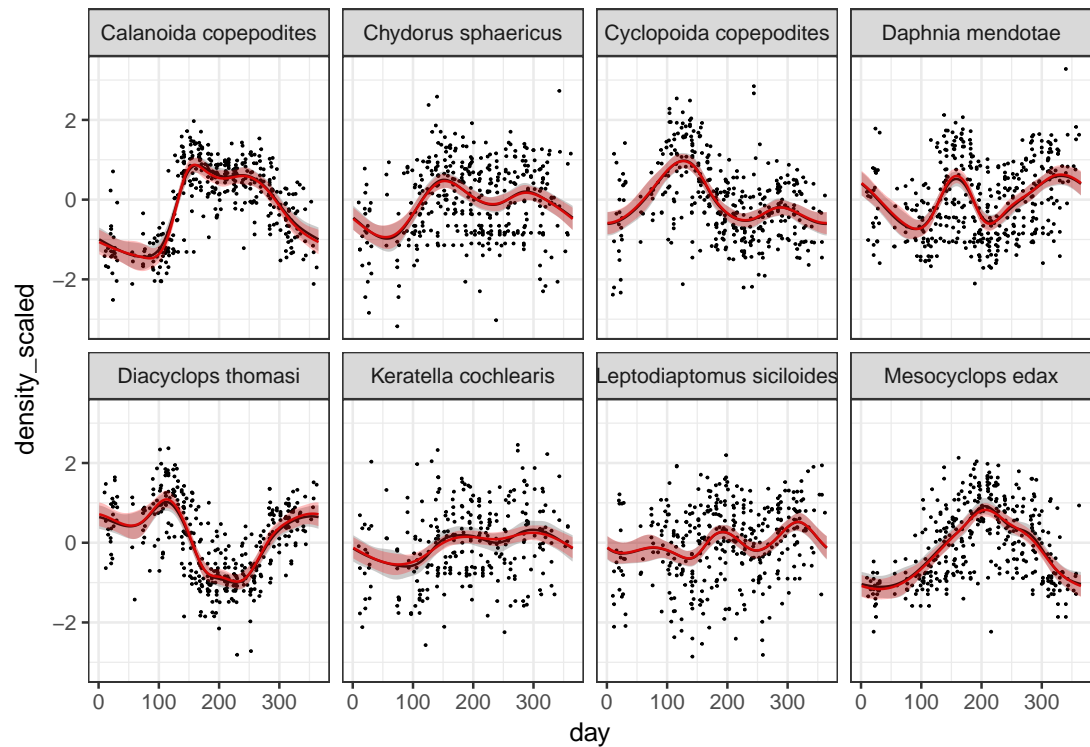
zoo_plot_data$mod4_fit = as.numeric(zoo_mod4_fit$fit)
zoo_plot_data$mod5_fit = as.numeric(zoo_mod5_fit$fit)

zoo_plot_data$mod4_se = as.numeric(zoo_mod4_fit$se.fit)
zoo_plot_data$mod5_se = as.numeric(zoo_mod5_fit$se.fit)

zoo_plot = ggplot(zoo_plot_data, aes(x=day))+
  facet_wrap(~taxon, nrow = 2)+
  geom_point(data= zoo_train, aes(y=density_scaled),size=0.1)+
  geom_line(aes(y=mod4_fit))+
  geom_line(aes(y=mod5_fit),color="red")+
  geom_ribbon(aes(ymin = mod4_fit - 2*mod4_se,
                 ymax = mod4_fit + 2*mod4_se),
             alpha=0.25)+
  geom_ribbon(aes(ymin = mod5_fit - 2*mod5_se,
                 ymax = mod5_fit + 2*mod5_se),
             alpha=0.25, fill="red")+
  theme_bw()

print(zoo_plot)

```



715

716 The two curves are very close for all species, but the differences in smoothness that resulted in  
 717 model 5 having an higher AIC than model 4 seem to be driven by the low seasonality of *Keratella*  
 718 *cochlearis* and *Leptodiaptomus siciloides* relative to the other species. Still, both models show  
 719 very similar fits to the training data, model 5 is only slightly better at predicting out of sample  
 720 fits for *K. cochlearis*, and not at all better for *L. siciloides*:

```
#Getting the out of sample predictions for both models:
zoo_test$mod4 = as.numeric(predict(zoo_comm_mod4,zoo_test))
zoo_test$mod5 = as.numeric(predict(zoo_comm_mod5,zoo_test))

#Correlations between fitted and observed values for all species:
zoo_test_summary = zoo_test %>%
  group_by(taxon)%>%
  summarise(mod4_cor = round(cor(density_scaled, mod4),2),
            mod5_cor = round(cor(density_scaled, mod5),2))

print(zoo_test_summary)
```

```
721 ## # A tibble: 8 x 3
722 ##   taxon                mod4_cor mod5_cor
723 ##   <fct>                <dbl>    <dbl>
724 ## 1 Calanoida copepodites 0.690    0.700
725 ## 2 Chydorus sphaericus  0.360    0.350
726 ## 3 Cyclopoida copepodites 0.540    0.530
727 ## 4 Daphnia mendotae     0.340    0.350
728 ## 5 Diacyclops thomasi    0.810    0.810
729 ## 6 Keratella cochlearis 0.230    0.200
730 ## 7 Leptodiaptomus siciloides 0.330    0.330
731 ## 8 Mesocyclops edax     0.570    0.570
```

732 Now let's look at how to fit inter-lake variability in dynamics for just *Daphnia mendotae*. Here,  
 733 we will compare models 1,2, and 3, to determine if a single global function is appropriate for

734 all four lakes, or if we can effectively model variation between lakes with a shared smooth or  
735 lake-specific smooths.

```
daphnia_train = subset(zoo_train,taxon=="Daphnia mendotae")
daphnia_test = subset(zoo_test,taxon=="Daphnia mendotae")

zoo_daph_mod1 = gam(density_scaled~s(day, bs="cc",k=10),
                    data=daphnia_train,
                    knots = list(day =c(1,365)),
                    method = "ML"
                    )

summary(zoo_daph_mod1)
```

```
736 ##
737 ## Family: gaussian
738 ## Link function: identity
739 ##
740 ## Formula:
741 ## density_scaled ~ s(day, bs = "cc", k = 10)
742 ##
743 ## Parametric coefficients:
744 ##             Estimate Std. Error t value Pr(>|t|)
745 ## (Intercept) 2.309e-16  4.245e-02      0      1
746 ##
747 ## Approximate significance of smooth terms:
748 ##             edf Ref.df      F p-value
749 ## s(day) 6.824      8 13.96 <2e-16 ***
750 ## ---
751 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
752 ##
753 ## R-sq.(adj) =  0.217   Deviance explained = 23.1%
754 ## -ML = 519.41   Scale est. = 0.7262      n = 403
```

```
zoo_daph_mod2 = update(zoo_daph_mod1,
                      formula = density_scaled~s(day, bs="cc",k=10) +
                                s(day,lake, k=10, bs="fs",
                                xt=list(bs="cc")))

summary(zoo_daph_mod2)
```

```
755 ##
756 ## Family: gaussian
757 ## Link function: identity
758 ##
759 ## Formula:
760 ## density_scaled ~ s(day, bs = "cc", k = 10) + s(day, lake, k = 10,
761 ##      bs = "fs", xt = list(bs = "cc"))
762 ##
763 ## Parametric coefficients:
764 ##             Estimate Std. Error t value Pr(>|t|)
765 ## (Intercept) 0.007772  0.043919  0.177    0.86
766 ##
767 ## Approximate significance of smooth terms:
768 ##             edf Ref.df      F p-value
769 ## s(day)      6.797      8 10.881 <2e-16 ***
770 ## s(day,lake) 5.348     35  0.432  0.0039 **
```

```

771 ## ---
772 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
773 ##
774 ## R-sq.(adj) =  0.246   Deviance explained = 26.9%
775 ## -ML = 516.48   Scale est. = 0.69947    n = 403

zoo_daph_mod3 = update(zoo_daph_mod1,
                      formula = density_scaled~s(day, bs="cc",k=10) +
                                s(day,by=lake, k=10, bs="cc"))

summary(zoo_daph_mod3)

776 ##
777 ## Family: gaussian
778 ## Link function: identity
779 ##
780 ## Formula:
781 ## density_scaled ~ s(day, bs = "cc", k = 10) + s(day, by = lake,
782 ##      k = 10, bs = "cc")
783 ##
784 ## Parametric coefficients:
785 ##              Estimate Std. Error t value Pr(>|t|)
786 ## (Intercept) -0.0007017  0.0416980  -0.017   0.987
787 ##
788 ## Approximate significance of smooth terms:
789 ##              edf Ref.df      F p-value
790 ## s(day)          6.8552361      8 13.985 < 2e-16 ***
791 ## s(day):lakeKegonsa 0.0488015      8  0.006 0.34599
792 ## s(day):lakeMendota 0.0005693      8  0.000 0.73813
793 ## s(day):lakeMenona  0.9270931      8  0.198 0.16129
794 ## s(day):lakeWaubesa 2.2353408      8  1.454 0.00116 **
795 ## ---
796 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
797 ##
798 ## R-sq.(adj) =  0.246   Deviance explained = 26.5%
799 ## -ML = 514.72   Scale est. = 0.69938    n = 403

800 The AIC values indicate that both model 2 and 3 are better fits than model 1, but models 2
801 and 3 have similar fits to one another. There does not seem to be a large amount of inter-lake
802 variability (all three models have similar adjusted  $R^2$ ), and model 3 indicates that only Lake
803 Waubesa deviates substantially from the overall dynamics. The plots for all three models (model
804 1 as dashed line, model 2 in black and model 3 in red) show that Medota and Menona lakes are
805 very close to the average and to one another for both models (which is unsurprising, as they are
806 very closely connected by a short river) but both Kegons and Waubesa show evidence of a more
807 pronounced spring bloom and lower winter abundances. While this is stronger in Lake Waubesa,
808 model 2 (in black) shows that it is still detectable in Lake Kegonsa if we do not need to fit a
809 separate penalty for each lake.

library(ggplot2)
library(dplyr)

#Create synthetic data to use to compare predictions
daph_plot_data = expand.grid(day = 1:365, lake = factor(levels(zoo_train$lake)))

daph_mod1_fit = predict(zoo_daph_mod1, daph_plot_data, se.fit = T)
daph_mod2_fit = predict(zoo_daph_mod2, daph_plot_data, se.fit = T)

```

```

daph_mod3_fit = predict(zoo_daph_mod3, daph_plot_data, se.fit = T)

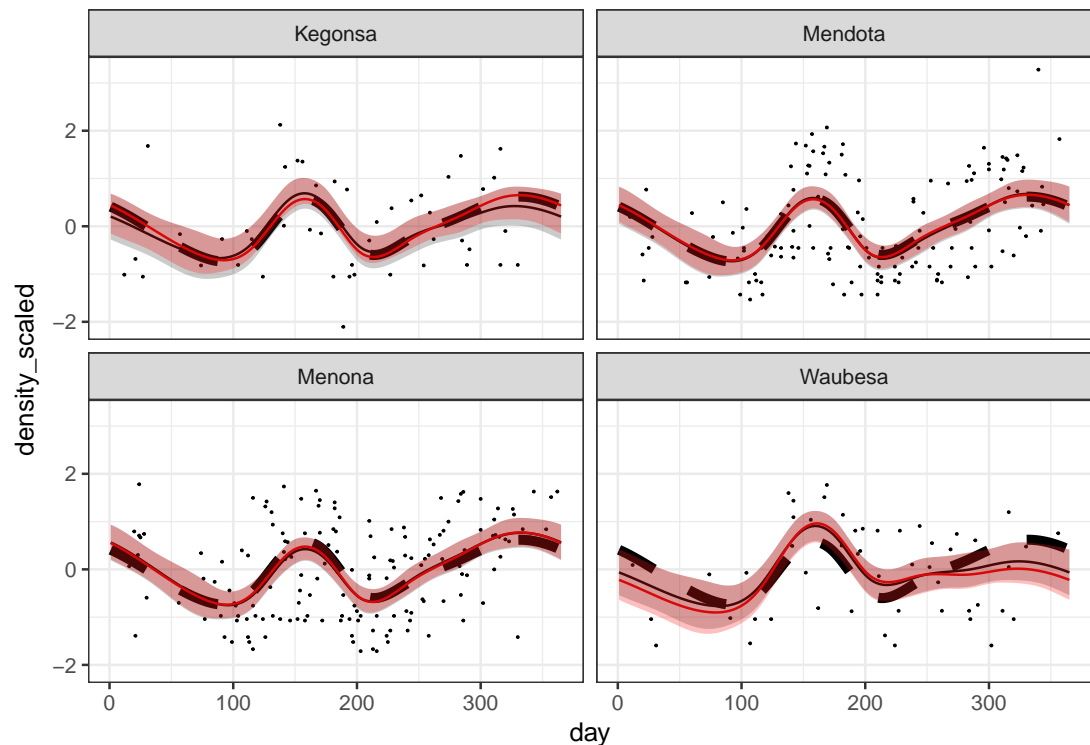
daph_plot_data$mod1_fit = as.numeric(daph_mod1_fit$fit)
daph_plot_data$mod2_fit = as.numeric(daph_mod2_fit$fit)
daph_plot_data$mod3_fit = as.numeric(daph_mod3_fit$fit)

daph_plot_data$mod1_se = as.numeric(daph_mod1_fit$se.fit)
daph_plot_data$mod2_se = as.numeric(daph_mod2_fit$se.fit)
daph_plot_data$mod3_se = as.numeric(daph_mod3_fit$se.fit)

daph_plot = ggplot(daph_plot_data, aes(x=day))+
  facet_wrap(~lake, nrow = 2)+
  geom_point(data= daphnia_train, aes(y=density_scaled),size=0.1)+
  geom_line(aes(y=mod1_fit),linetype=2, size=2)+
  geom_line(aes(y=mod2_fit),color="black")+
  geom_line(aes(y=mod3_fit),color="red")+
  geom_ribbon(aes(ymin = mod2_fit - 2*mod2_se,
                ymax = mod2_fit + 2*mod2_se),
            alpha=0.25)+
  geom_ribbon(aes(ymin = mod2_fit - 2*mod3_se,
                ymax = mod2_fit + 2*mod2_se),
            alpha=0.25, fill="red")+
  theme_bw()

print(daph_plot)

```



810

811 In this case, model 2 is able to predict as good or better out of sample as model 1 or 3, indicating  
 812 that jointly smoothing the lake together improved model prediction. However, None of the  
 813 models did well in terms of predicting Lake Kegonsa dynamics out of sample (with a correlation

of only 0.11 between predicted and observed densities), indicating that this model may be missing substantial year-to-year variability in *D. mendotae* dynamics:

```
#Getting the out of sample predictions for both models:
daphnia_test$mod1 = as.numeric(predict(zoo_daph_mod1,daphnia_test))
daphnia_test$mod2 = as.numeric(predict(zoo_daph_mod2,daphnia_test))
daphnia_test$mod3 = as.numeric(predict(zoo_daph_mod3,daphnia_test))

# We'll look at the correlation between fitted and observed values for all species:
daph_test_summary = daphnia_test %>%
  group_by(lake)%>%
  summarise(mod1_cor = round(cor(density_scaled, mod1),2),
            mod2_cor = round(cor(density_scaled, mod2),2),
            mod3_cor = round(cor(density_scaled, mod3),2))

print(daph_test_summary)
```

```
## # A tibble: 4 x 4
##   lake      mod1_cor mod2_cor mod3_cor
##   <fct>      <dbl>    <dbl>    <dbl>
## 1 Kegonsa    0.0800    0.110    0.0900
## 2 Mendota    0.420     0.430    0.430
## 3 Menona     0.350     0.400    0.390
## 4 Waubesa    0.290     0.290    0.270
```

## BIBLIOGRAPHY

- Baayen, R. Harald, Jacolien van Rij, Cecile de Cat, and Simon N. Wood. 2016. “Autocorrelated Errors in Experimental Data in the Language Sciences: Some Solutions Offered by Generalized Additive Mixed Models.” *arXiv:1601.02043 [Stat]*, January.
- Bolker, Benjamin M, Mollie E Brooks, Connie J Clark, Shane W Geange, John R Poulsen, M Henry H Stevens, and Jada-Simone S White. 2009. “Generalized linear mixed models: a practical guide for ecology and evolution.” *Trends in Ecology & Evolution* 24 (3): 127–35.
- Boor, Carl de. 1978. *A Practical Guide to Splines*. Springer.
- Efron, Bradley, and Carl Morris. 1977. “Stein’s Paradox in Statistics.” *Scientific American* 236 (5): 119–27.
- Gelman, A., J.B. Carlin, H.S. Stern, D.B. Dunson, A. Vehtari, and D.B. Rubin. 2013. *Bayesian Data Analysis, Third Edition*. Chapman & Hall/Crc Texts in Statistical Science. Taylor & Francis.
- Gelman, Andrew. 2006. “Multilevel (Hierarchical) Modeling: What It Can and Cannot Do.” *Technometrics* 48 (3): 432–35.
- Hastie, T J, and Robert J Tibshirani. 1990. *Generalized Additive Models*. Monographs on Statistics and Applied Probability. Taylor & Francis.
- Kimeldorf, George S., and Grace Wahba. 1970. “A Correspondence Between Bayesian Estimation on Stochastic Processes and Smoothing by Splines.” *The Annals of Mathematical Statistics* 41 (2): 495–502. doi:10.1214/aoms/1177697089.
- Lathrop, R. C. 2000. “Madison Wisconsin Lakes Zooplankton 1976 - 1994.” *Environmental Data Initiative*. <http://dx.doi.org/10.6073/pasta/ec3d0186753985147d4f283252388e05>.
- McCullagh, P, and John A Nelder. 1989. *Generalized Linear Models, Second Edition*. CRC Press.
- McMahon, Sean M, and Jeffrey M Diez. 2007. “Scales of association: hierarchical linear models



- and the measurement of ecological systems.” *Ecology Letters* 10 (6): 437–52.
- Potvin, C, M.J. Lechowicz, and S. Tardif. 1990. “The Statistical Analysis of Ecophysiological Response Curves Obtained from Experiments Involving Repeated Measures.” *Ecology*, 1389–1400.
- Ruppert, David, M P Wand, and R J Carroll. 2003. *Semiparametric Regression*. Cambridge University Press.
- Stanley, R.R.E., E. J. Pedersen, and P.V.R. Snelgrove. 2016. “Biogeographic, Ontogenetic, and Environmental Variability in Larval Behaviour of American Lobster (*Homarus Americanus*).” *Marine Ecology Progress Series* 553: 125–46.
- Verbyla, Arūnas P., Brian R. Cullis, Michael G. Kenward, and Sue J. Welham. 1999. “The Analysis of Designed Experiments and Longitudinal Data by Using Smoothing Splines.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 48 (3): 269–311. doi:10.1111/1467-9876.00154.
- Vickers, Mathew J., Fabien Aubret, and Aurélie Coulon. 2017. “Using GAMM to Examine Inter-Individual Heterogeneity in Thermal Performance Curves for *Natrix Natrix* Indicates Bet Hedging Strategy by Mothers.” *Journal of Thermal Biology* 63 (January): 16–23. doi:10.1016/j.jtherbio.2016.11.003.
- Wood, Simon N. 2006. *Generalized Additive Models*. An Introduction with R. CRC Press.
- Wood, Simon N. 2003. “Thin Plate Regression Splines.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 65 (1): 95–114.
- . 2011. “Fast Stable Restricted Maximum Likelihood and Marginal Likelihood Estimation of Semiparametric Generalized Linear Models.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 73 (1): 3–36. doi:10.1111/j.1467-9868.2010.00749.x.
- . 2017. *Generalized Additive Models: An Introduction with R, 2nd Edition*. 2nd ed. Boca Raton, FL: CRC Press.
- Wood, Simon N., Yannig Goude, and Simon Shaw. 2015. “Generalized Additive Models for Large Data Sets.” *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 64 (1): 139–55. doi:10.1111/rssc.12068.
- Wood, Simon N., Fabian Scheipl, and Julian J. Faraway. 2013. “Straightforward Intermediate Rank Tensor Product Smoothing in Mixed Models.” *Statistics and Computing* 23 (3): 341–60. doi:10.1007/s11222-012-9314-z.