

1 The main question: Are PMM's appropriate for my data structure?

1.1 What is my desired model and data structure?

I'd like to use a repeat measure PMM as detailed in chapter 11 of Modern Phylogenetic Comparative Methods and the application in Evolutionary Biology (de Villemereuil Nakagawa, 2014).

I am trying to predict the time between flowering and leaf expansion (FLS) for 23 species from Harvard Forest as a function of their pollination syndrome, minimum precipitation across their range, day of flowering with all two way interactions and accounting for phylogeny.

The possible problem (the point of this treatment is to determine whether or not this is a problem) is all of my variable are recorded at different levels of organization.

- FLS- repeat measure across 3-5 individuals per species over 15 year.
- pollination syndrome- one value/species
- minimum precipitation-one value/species
- flowering time: repeat measure across 3-5 individuals per species over 15 year.
- See the head of the data below:

##	species	name	fopn.jd	min_precip	pol	funct.flr
## 1	ACPE Acer_pensylvanicum		136	24	0	9
## 2	ACPE Acer_pensylvanicum		129	24	0	19
## 3	ACPE Acer_pensylvanicum		138	24	0	16
## 4	ACPE Acer_pensylvanicum		130	24	0	12
## 5	ACPE Acer_pensylvanicum		135	24	0	9
## 6	ACPE Acer_pensylvanicum		138	24	0	7

I standardize (z-score) all predictors for comparison.

This model I have been using and would like to continue to do so:

```
FLS = pol+flowing+precip+precip:fl+precip:pol+pol:flowing+(1—name)+(1—tree.id/species),covranef  
= list(name= A))
```

The covariance structure to the random effect "name" accounts for phylogenetic correlations in the inter-specific residuals. The tree.id/species random effect is in this case "an individual/species-specific" term accounting for variability in the estimate not attributed to phylogeny. This all seems well and good until we encounter the following line from de Villemereuil Nakagawa:

" As a comparative biologist, the reader would most likely be interested in the between species slope. If co factor x (ie my predictor), only contains one value per species, then there is no problem... Things are

slightly more complicated using individual measurements in x, but it is still possible to obtain between species and within species slopes using a technique called within group centering (Davis et al 1961, van de Pol and Wright 2009). This separates x into two components, the containing the species level mean, and one containing species level variability.”

Here’s the thing. I am not precisely interested in estimating the within species slope. Rather, I’d like to estimate the between species slope, but have the model “account” for the within species variation but not necessarily estimate it explicitly.

Yet, the quote above would suggest the proper model is:

```
FLS = pol+meanflotime+varflotime+precip+all 2 way interactions +(1—name)+(1—species) cov-  
ranef = list(name= A))
```

When I run this model, it generates huge credible interval for all predictors, and only mean flowering time matters. This is contrary to all other models I’ve ever run on this data set and other hysteryanth data, as well as my biological expectation. (eg I am fairly certain pollination syndrome matters.)

2 So what went wrong?

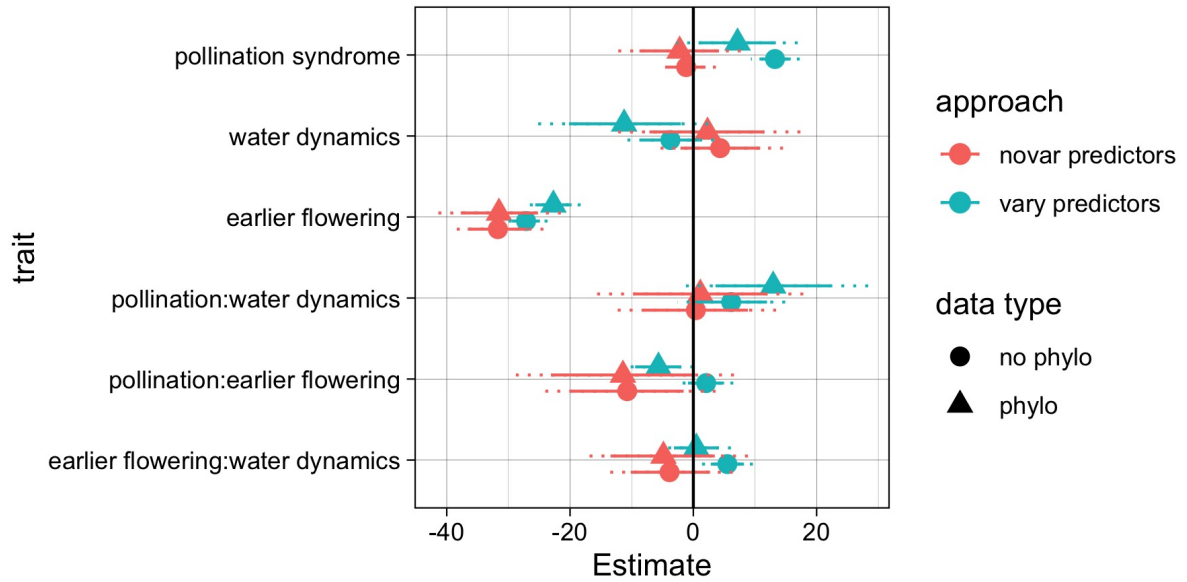
Why are these answers so different? Nacho and I discussed this. Neither of us were sure. Nacho thinks maybe that what seems important is these data are changing to phylogenetic covariance structure in the residuals.

- Suggeston: Because the variables aren’t balanced (ie Acer rubrum has 100 FLS and flowering time estimates and other species have as few as 40), this is biasing the phylogeny through weighting it differently.
- But here’s what I’m thinking: To me this would be a problem with a lack of balance in the response variable more than the predictor because the model is pooling only on intercepts. I don’t see why partitioning within and among species variance in a *predictor*, would help this.

3 A foray into qualitative model comparison

So we don’t really know what is wrong. But if we follow the lead that the phylogenetic inaccuracies are driving the giant differences between the models, it follows accordingly that removing the phylogeny and covariance structures entirely should bring he models closer together.

So I ran 4 models. 1. meanfloweringtime w/ phylo (no vary predictors). 2. meanfloweringtime w/o phylo (no vary predictors). 3. original model w/ phylo (vary predictors). 4. original model w/o phylo (vary predictors. *Note: for simplicitly of comparision I didn’t include the flowering variationterm in this graphic, just mean flowering time per species. When I did run that full model, there was also a much weaker effect of flowering time variation, but the other predictors didn’t change.* Here is a comparative mu plot:



As you can see here, even when ignoring the phylogeny the mean flowering time and original models are very different. Adding the phylogeny changes them a bit, but each one is still more similar to itself with or without the phylogeny .

4 Take-away:

I think it's still important to understand why these two approaches yield such different estimates. However, because my original model with and without the phylogeny still give estimates that are FAR more qualitatively similar to each other than either of the mean flowering time models, I don't feel like my original formulation is wrong per say, just different.

I would love to hear your thoughts on this.