

## Background

Traditionally cone crop prediction modeling has focused on calculating correlations between various predictors and the annual cone count observed for a stand of trees. Usually, predictor variables are functions of the average temperature in the years preceding the observed cone crop during the growing season, since that is when the various biological processes involved in reproduction occur.

Kelly et al explored several such predictor variables and compared their correlations across a broad dataset involving multiple plant families and over a long period of time. The predictor variables considered were:

- T1 model: the mean summer temperature in the previous year ( $T_{n-1}$ ).
- T2 model: the mean summer temperature 2 years previously ( $T_{n-2}$ ).
- $\Delta T$  model: the change in mean summer temperature over the two preceding years ( $T_{n-1} - T_{n-2}$ ).
- 2T model: both mean summer temperature in the previous year ( $T_{n-1}$ ) and mean summer temperature 2 years previously ( $T_{n-2}$ ).

Correlations were calculated using linear fits between the predictor variable and log of the annual seedfall  $c$ . Kelly et al compared the p-values calculated with the null hypothesis being that no correlation exists, and found that for most plant families the  $\Delta T$  model produced the lowest p-values, and concluded that  $\Delta T$  is an ideal cue for seed crop prediction.

## Motivation

### Disadvantages of correlations

There are a number of reasons to suspect both alternative models and alternative mathematical approaches could be better motivated than the proposed approach of Kelly et al.

1. **p-values:** Although the p-values calculated are small, any analysis in which many p-values are calculated is bound to have statistically significant p-values *somewhere*. No effort to control for the look-elsewhere effect was made. Furthermore the null hypothesis - that the models analyzed by Kelly have 0 correlation with seed crop - doesn't capture what is already known about plant reproduction; for example temperatures *must* have an effect on reproductive processes because plants can't reproduce in temperatures inhospitable to life. It therefore isn't surprising that there's a correlation between some predictor involving temperature and the observed seed crop.
2. **Linearity:** There's no reason *a priori* to think that the relationship between any of the predictors put forth by Kelly et al should be linear with  $\log(c)$ , except that any continuous and differential function can be approximated to first order as a line (Taylor series). No discussion of this choice is made, even though it may be valid, although the implications of this assumption mean that the calculations become easier.
3. **Homoskedasticity and normality:** In the supplemental material, Kelly et al argue that empirically the measured seedfall appears to be homoskedastic, i.e. that the variance doesn't change with the number of seeds observed. This is part of a larger implied argument that the observed log-seedfall  $\log(c)$  is a normally distributed random variable with a fixed variance  $\sigma^2$  and a mean  $\Delta T$ . In short, homoskedasticity is used as an argument that  $\log(c) \sim \mathcal{N}(\Delta T, \sigma)$ .

But is seed production really a homoskedastic process? Naively plants that have more resources available for seed production should have a greater *variation* in seed production. Furthermore the assumption of normality is not justified by the fact that the normal distribution has a nonzero probability on the support  $(-\infty, \infty)$ , while  $\log(c)$  is only defined (for real values) on the interval  $(0, \infty)$ , meaning that *any* normally-distributed predictor can't possibly explain the range of seed

crops observed in nature. Most importantly, log-transforming the seedfall is certainly unjustified for the simple fact that this approach cannot be used to correlate with years in which 0 seeds are observed ( $\log(0) = -\infty$ ). These data points are arguably the *most* important observations in species that exhibit masting because in non-mast years, few to no seeds are often observed.

4. **Data Preprocessing:** As a result of the log-transformation that is carried out on the observed seedfall, years in which no seeds are observed cannot be used in the correlation without utterly dominating the linear fit used to produce the correlation ( $\log(0) = -\infty$ ). Kelly et al directly modify these observations before the log transformation, replacing them with values that are half the smallest nonzero value. Similar procedures are carried out elsewhere in the literature, but in all cases this approach is mathematically unjustified; any valid model of cone production must be able to predict 0-seedfall years.
5.  **$\Delta T$  as a model:** Intuition tells us that plants that experience freezing conditions for multiple years will not reproduce as much as the same plants that experience ideal growing conditions for multiple years because they simply do not have the same resources available for reproduction. However, two subsequent years of ideal reproductive temperatures may have the same small value of  $\Delta T$  as two subsequent years of terrible reproductive temperatures. Kelly et al discuss this as an interesting consequence of this model - that as a general rule plant reproduction will be insensitive to changes in global temperatures because  $\Delta T$  is unaffected by average changes in temperature. Furthermore the number of cones produced is not taken into account in the  $\Delta T$  model, although it is known that multiple sequential masting events are only rarely observed.

## A different approach

Instead of carrying out linear correlations, let's focus on a different approach. Besides the reasons I've previously discussed for trying something different than Kelly et al, there are also other reasons why one might want to take a Bayesian approach. [Here's a good discussion about this if you want to learn more.](#)

Here's the general idea:

1. We start by writing down a probability distribution which describes the probability of observing an individual data point  $c_i$  given some model. I propose that the probability is given by a Poisson distribution:

$$P(c_i | c_{\mu,i}) = \frac{c_{\mu,i}^{c_i} e^{-c_{\mu,i}}}{c_i!} \quad (1)$$

where  $c_i$  is the number of cones observed for the stand on day  $i$ , and  $c_{\mu,i}$  is the Poisson rate parameter - the expected number of cones on day  $i$ .

The Poisson distribution is a *counting distribution*, a distribution that gives the probability of a number of events happening in a given amount of time. It has some nice properties: it's discrete (i.e. it is only defined for integer cone counts) and is defined on the support  $[0, \infty)$ , as is required by our data.

2. From the probability of observing an individual day's cone crop  $c_i$ , we can write down the probability of observing the cone crop for all  $N$  days (the entire dataset) by taking the product of all the individual probabilities:

$$P(\{c_i\} | \{c_{\mu,i}\}) = \prod_{i=0}^N \frac{c_{\mu,i}^{c_i} e^{-c_{\mu,i}}}{c_i!} \quad (2)$$

This is the *likelihood*; in the literature it is usually written as

$$\mathcal{L}(D | \theta) \quad (3)$$

where for us the data  $D = \{c_i\}$  and the model  $\theta = \{c_{\mu,i}\}$ .

3. From here, we can write down an expression for the *posterior* probability distribution  $\mathcal{P}$ , which is a distribution over our model parameters conditioned on our data. We make use of Bayes's theorem, and up to a normalization constant

$$\mathcal{P}(\theta | D) \propto \mathcal{L}(D | \theta) P(\theta) \quad (4)$$

where  $P(\theta)$  is a prior distribution; this distribution characterizes the epistemic uncertainty in our model *prior* to observing any data.

In the Bayesian approach, getting the posterior probability distribution is the entire goal. We could spend a long time talking about how this differs from frequentist approaches which usually focus on metrics related to the likelihood (for example, least squares fitting is equivalent to maximum likelihood estimation with the implicit and extremely restrictive assumption that there is no uncertainty in the independent variable, and normally distributed uncertainty in the observed dependent variable), but for now I'm going to just point to [Jake VanderPlas's blog](#) and [this talk by Chris Fonnesbeck](#).

As you'll see with models discussed below, usually the posterior can only be written down in a nice closed form for certain special cases where the likelihood and priors are *conjugate*. In that case you

can actually calculate the posterior by hand, but in most cases you really need to turn to computers to sample from the posterior distribution numerically, and that's what we'll do here.

## Modeling

This part focuses on coming up with reasonable expressions of the rate parameter  $c_{\mu,i}$ . Fundamentally I think that the expected number of cones should be determined entirely by what resources are available to produce them. For the discussion below, I only attempted to fit data for site 1 (10ABAM\_OR).

### Three years preceding model

This model assumes that the energy available to produce cones is mostly gathered from average temperatures in specific windows in the three years preceding the cone crop, with no contribution from years before that. As a proxy for Photosynthetically Active Radiation (PAR) we instead use temperature (since that's the data we have), and make the assumption that the two are proportional. Under that assumption we can write down the expected number of cones:

$$c_{\mu,i} = c_0 + \underbrace{\alpha_0 \langle T \rangle_{i-l_0, w_0} + \alpha_1 \langle T \rangle_{i-l_1, w_1} + \alpha_2 \langle T \rangle_{i-l_2, w_2}}_{\text{PAR}} - c_{i-l_3} \quad (5)$$

where  $\langle T \rangle_{i-l_k, w_k}$  denotes the moving average of the temperature  $T$  over a window of size  $2w_k + 1$  days surrounding the day  $i - l_k$ . Here, the  $\alpha_k$  are fit parameters which determine the relative importance of each year's sunlight contribution to the stand's energy reserves.  $c_0$  is the initial energy reserves of the stand at the beginning of our observations,  $c_{i-l_3}$  is the lagged cone count from  $l_3$  days in the past.

## Priors

I chose some prior probability distributions based on what I know about cone production. These characterize the epistemic uncertainty about our model parameters prior to observing any data:

Parameter	Prior	Unit of measure	Comment
$c_0$	HalfNorm(100)	# of cones	Initial energy reserves at start of dataset
$\alpha_0$	HalfNorm(20)	cones/ $^{\circ}$ K	Weakly informative choice of half-normal distribution. This is probably a small number
$\alpha_1$	HalfNorm(20)	cones/ $^{\circ}$ K	Weakly informative choice of half-normal distribution. This is probably a small number
$\alpha_2$	HalfNorm(20)	cones/ $^{\circ}$ K	Weakly informative choice of half-normal distribution. This is probably a small number
$w_0$	Uniform(1, 100)	days	Window size used to calculate the average temperature in the first year. Probably in the range of 1-100 days long
$w_1$	Uniform(1, 100)	days	Window size used to calculate the average temperature in the second year. Probably in the range of 1-100 days long
$w_2$	Uniform(1, 100)	days	Window size used to calculate the average temperature in the second year. Probably in the range of 1-100 days long
$l_0$	Uniform(185, 545)	days	Lag time of the moving average of the temperature in the first year; constrained to be 0.5 to 1.5 years before the measured crop
$l_1$	Uniform(545, 910)	days	Lag time of the moving average of the temperature in the second year; constrained to be 1.5 to 2.5 years before the measured crop
$l_2$	Uniform(910, 1275)	days	Lag time of the moving average of the temperature in the second year; constrained to be 2.5 to 3.5 years before the measured crop
$l_3$	Uniform(910, 1275)	days	Lag time used to get the last cone crop, constrained to be 2.5 to 3.5 years before the measured crop

## Posterior

I used Markov Chain Monte Carlo to sample from the posterior distribution using a software package called [emcee](#). You can think of this like a black box which proposes values of the model parameters  $\theta = (c_0, \alpha_0, \alpha_1, \alpha_2, w_0, w_1, w_2, l_0, l_1, l_2, l_3)$ ; I wrote code that takes those proposed values and uses them to compute the prior and likelihood, which are then multiplied to get a posterior probability, which I then pass back to the black box. The [emcee](#) sampler then uses this posterior probability to generate new proposed values of  $\theta$ , and the process repeats many times. After an initial burn-in period, if the probabilities are defined correctly and the sampler is able to efficiently explore the parameter space, the sampler will eventually converge to certain values of  $\theta$ , and [due to some clever underlying math](#) we can guarantee that the samples being generated will [match the posterior distribution](#).

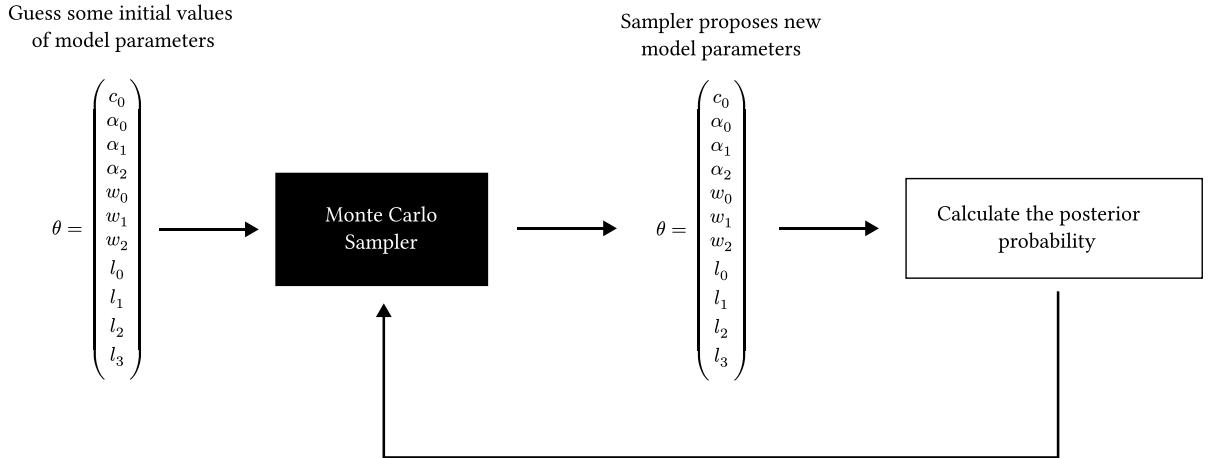


Figure 1: The MCMC sampler takes in an initial  $\theta$ , and proposes new values of  $\theta$ . For each proposal, we calculate the posterior probability of the proposed value, and pass it back to the sampler; it uses this value to generate new proposed values of  $\theta$ . After some point, the new proposed values will converge to posterior probability distribution.

Looking at the values of  $\theta$  generated by the sampler, we can see they never really settle to one specific value. There's probably an issue with the way I've parameterized the problem.

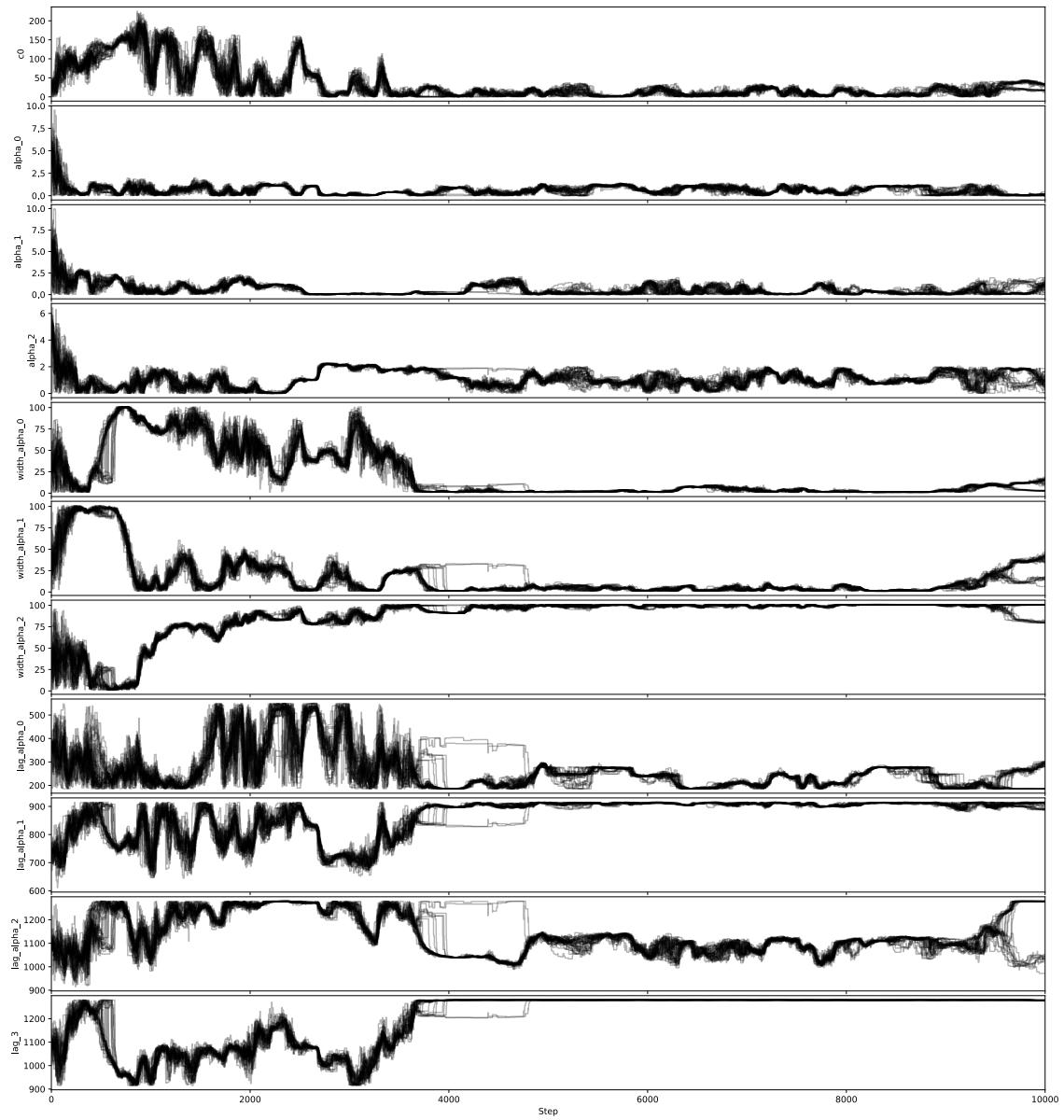


Figure 2: Markov chains for each fit parameter generated by emcee. Initially the chains vary as the MCMC sampler searches the parameter space of the problem; eventually they converge but fall into a region of instability, indicating that the model probably needs to be reparameterized.

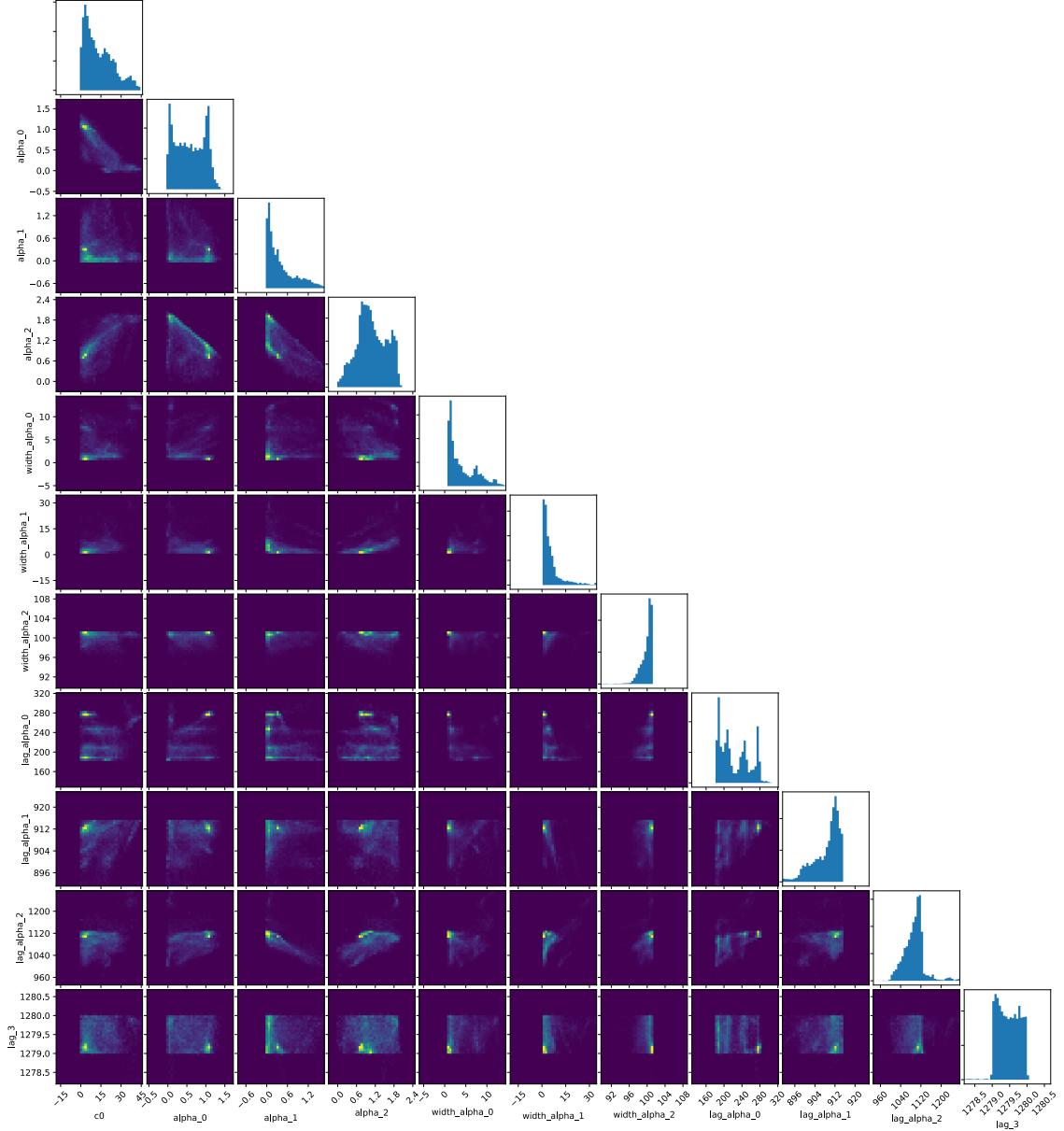


Figure 3: Usually we discard the initial values produced by the sampler as it takes time for the Markov chains to converge to the posterior distribution. In this case the sampler never converged; even so, for illustrative purposes we can throw out the first 6000 samples from each Markov chain, and then plot the histograms of the values of the model parameters. You can think of this plot as a series of 2D projections of the posterior probability which help us see how the model parameters depend on one another.

I tried a number of different transformations on the data but haven't found a stable parameterization of this model yet. Part of the difficulty is due to the fact that some of the variables here are meant to be discrete variables: the width of the moving average windows and the lag period for each window must all be integers, whereas  $\alpha$  can be floating point. This introduces computational complexities that are not well handled by some of the major Bayesian statistics packages. So to even get this model to run I've had to compromise:

- I've needed to use `emcee`, a package with greater flexibility but quite slow sampling rate

- The sampler chooses floating point samples for *all* parameters, but my model casts these to integers before computing the expected number of cones  $c_{\mu,i}$ . This discretization means that numerically speaking the sampler isn't as effective at exploring the parameter space, as multiple floating point values of e.g.  $w_0$  will be truncated to the same integer value:

$$\begin{aligned}\langle T \rangle_{i-20.2,25.4} &\rightarrow \langle T \rangle_{i-20,25} \\ \langle T \rangle_{i-20.5,25.7} &\rightarrow \langle T \rangle_{i-20,25}\end{aligned}\tag{6}$$

### Posterior Predictive Distribution

Although the model is not well behaved, for illustrative purposes we can also pretend it is, and discuss how we can use Bayesian models to make predictions. The way that this is done in practice is again using a rigorous statistical treatment - we calculate the **posterior predictive distribution**. This is the probability of observing some predicted dataset given the data we've already observed:

$$P(\tilde{D} | D) = \int_{\theta} P(\tilde{D} | \theta)P(\theta | D) d\theta\tag{7}$$

In practice we do this by asking the computer to generate Poisson-distributed random numbers (Equation 1) using samples of the values of the model parameters  $\theta$ . For each sample of the model parameters, we generate a new dataset; here, I've generated 10000 such datasets:

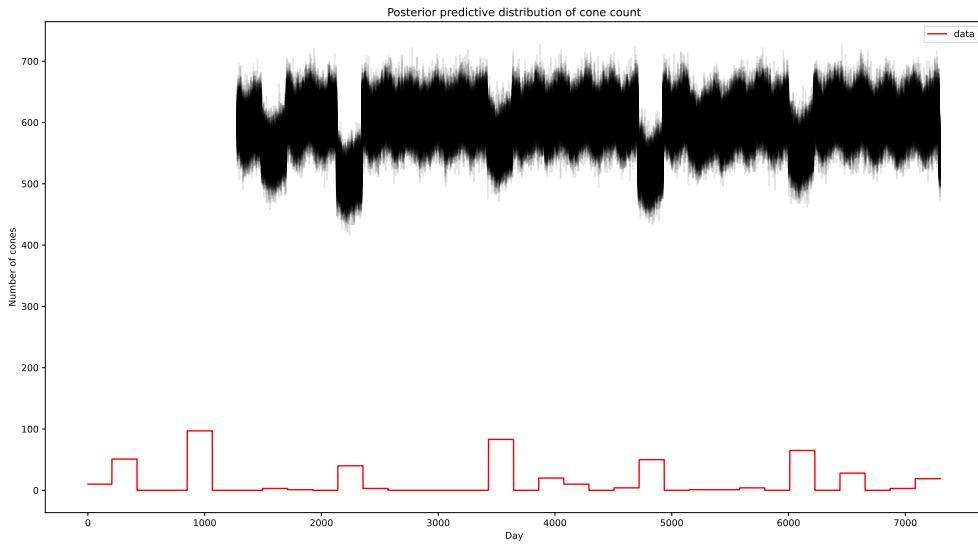


Figure 4: 10000 samples from the posterior predictive distribution, all plotted on the same graph. Each of these can be thought of as a prediction. The actual data is not close here because the sampler never converged.

Of course, these are just samples from the posterior predictive distribution; if we want to get an idea of what the distribution itself looks like, we can histogram the cone counts on each day across all predictions:

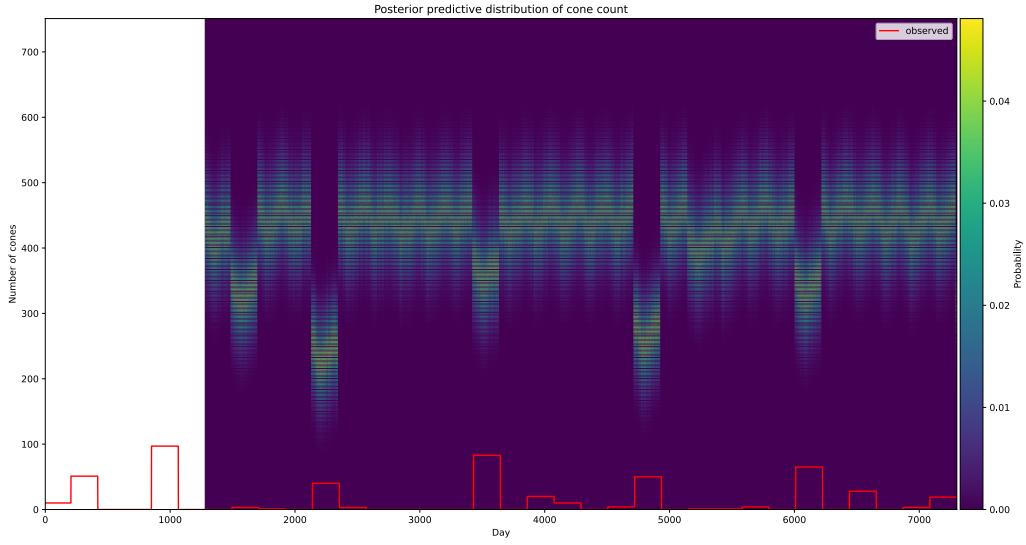


Figure 5: 10000 samples from the posterior predictive distribution histogrammed together. The actual data is not close to the distribution here because the sampler never converged.

For both figures above, the predictions start 3 years into the dataset because that's how much data we need to calculate a posterior probability; the model needs 3 years of temperature data before it can start making predictions.

This is the kind of graph I'd eventually like to make once the bugs in the model/code get sorted out. For now I've just included it to give you an idea of how I want to use these tools, to give you time to think about this approach which is very different to what has historically been done, and to hear your feedback.

### Resource-Accumulation Model

Instead of looking at windows of time in some number of years preceding a cone crop, we could instead take into account all energy expenditure of the stand starting from the beginning of the data.

$$c_{\mu,i} = c_0 + \underbrace{\alpha \int_0^{t_i} T(t) dt}_{\text{PAR}} - \int_0^{t_i} c(t) dt \quad (8)$$

Here the resources accumulated by the tree over time are considered: the PAR received each day is approximated as being proportional to the temperature that day; a potentially dubious approximation. The available resources of the stand include all the PAR absorbed since the beginning of the dataset less any spent on cone production.

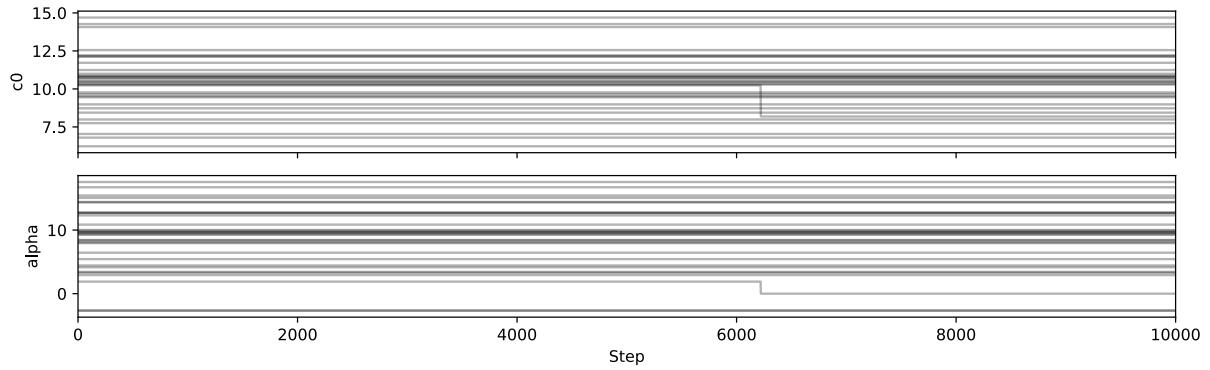


Figure 6: Markov chains for each fit parameter generated by emcee for the RAModel. The walkers never seem to explore the parameter space.

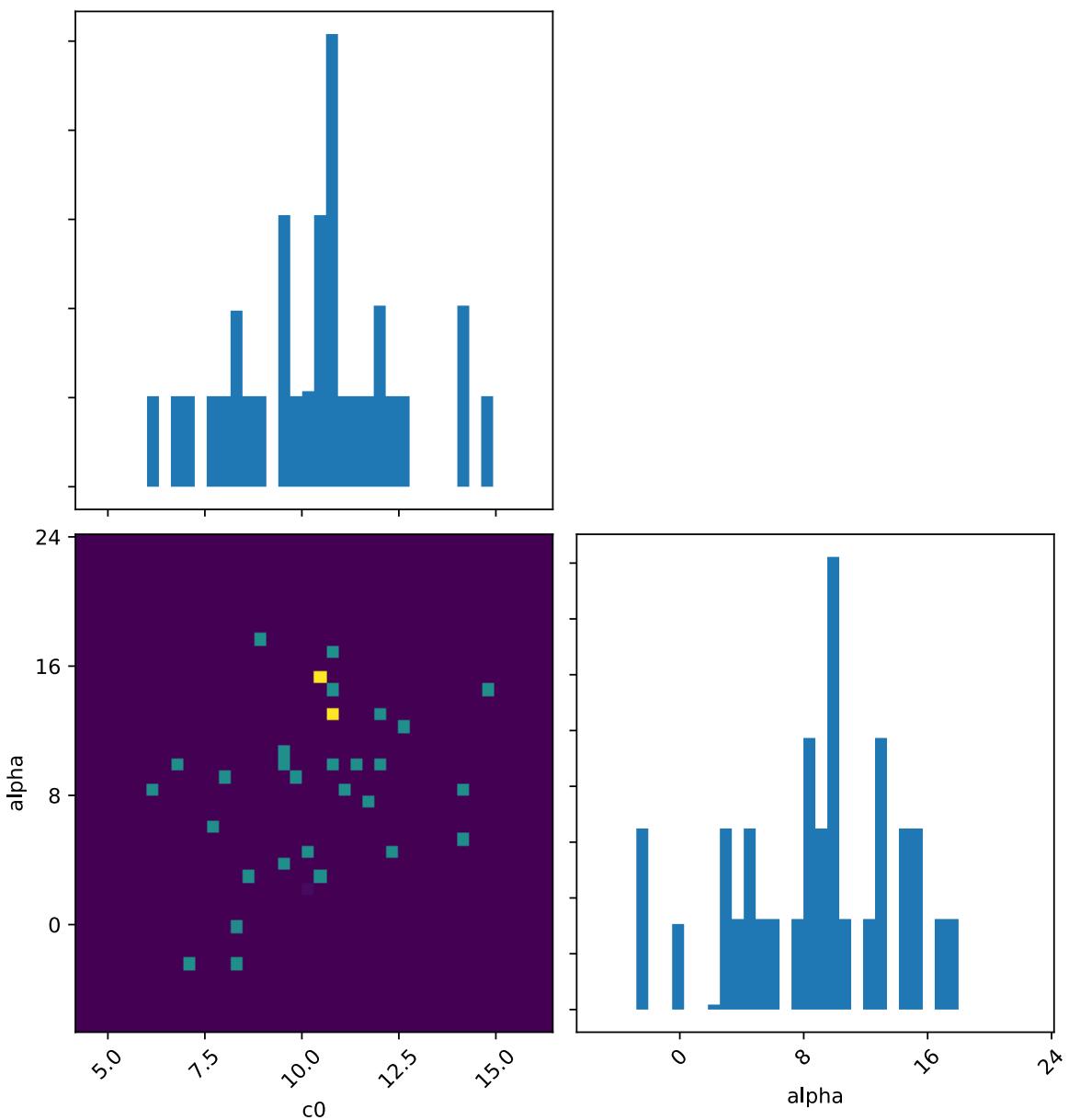


Figure 7: The posterior probability distribution sampled by emcee for the RAModel. The sampler never explores the parameter space; additional debugging is required.

I haven't bothered computing the posterior predictive distribution because I clearly can't sample effectively with the way I've written this model. So I've got to debug this before moving on to posterior predictive checks.

## Discussion

It took a number of months to realize that we don't need to entertain models that are similar to  $\Delta T$  if we don't want to, although that's what we had initially. We have the freedom to explore whatever model makes sense to us. Right now something like the resource-accumulation model makes the most sense to me, in part because it actually fully accounts for the energy flowing into the tree as PAR and flowing out of the tree from cone production, and partly because it is simpler than the n-years preceding model. There are fewer free parameters and computationally it's easier to model. This means I don't need to stick to emcee, I can switch to the much faster pymc, which may allow more models and parametrizations to be explored.

### A clue about where to go next

Some samples produced by the sampler result in a PAR term which is smaller than the cone contribution term for certain days; this is particularly true in the immediate aftermath of masting events, in which the cone-contribution term overwhelms the PAR term. This of course results in a negative expected cone crop rate parameter, which doesn't make physical sense. Moreover this results in nan-valued probabilities after masting events, which makes posterior predictive checks impossible in these regions. Up to this point I've been dealing with this by essentially omitting their contribution to the posterior probability.

One alternative is to set the probability for these data points to 0, which should push the sampler to higher values of  $\alpha$  to compensate. So far I've had absolutely no luck with this - it just seems like the Markov chains don't vary when I do this. Sorry I don't have more interesting news yet, but I'm still debugging, reparameterizing, and thinking about this.

Finally, the observed model instabilities could be due to

- Poor parameterization. The sampler most effectively explores parameter space when the steps it needs to take are roughly equal for each model parameter.
- A bug
- Poor initial guesses for the values of the parameters
- Some other problem with the model. I've tried to strike a balance so that the models I was running were not overly complex, but it's possible I need to include a term for the energy spent by the stand on things that aren't reproducing, e.g. growing wood, respirating, etc. I don't see this as an intrinsic reason why this approach won't work, I just haven't yet explored it.
- Unjustified priors. One of the things I like about the Bayesian approach is that all the assumptions are explicit. There are almost certainly better choices to be made for the priors here, so it's possible that could really affect the sampling.