

1002

Supplementary Materials for

1003 Gene-environment interactions govern early regeneration in fir and 1004 beech: evidence from participatory provenance trials across Europe

1005 Katalin Csilléry^{1,*}, Justine Charlet de Sauvage¹, Madleina Caduff^{2,3}, Johannes Alt¹, Marjorie
1006 Bison¹, Mert Celik¹, Nicole Ponta^{1,5}, Daniel Wegmann^{2,3}

1007 ¹Swiss Federal Research Institute WSL, Switzerland

1008 ²University of Fribourg, Switzerland

1009 ³Swiss Institute of Bioinformatics, Switzerland

1010 ⁵Ufficio della Caccia e della Pesca, Switzerland

1011 *Corresponding author: katalin.csillery@wsl.ch

1012 Description of a hidden Markov model of seed germination, 1013 phenological development, and mortality, including the 1014 inference and importance sampling scheme

1015 Let $z_i(t)$ denote the life stage of seed i on day t . We consider three life stages: i) a viable seed
1016 that is dormant and has not yet germinated ($z_i(t) = \mathcal{L}_S$), ii) a seed that germinated into a
1017 currently growing seedling ($z_i(t) = \mathcal{L}_G$) and iii) non-viable seed or a previously germinated
1018 seedling that has died ($z_i(t) = \mathcal{L}_D$). On each day, four events may happen in this order: 1) a
1019 seedling will grow, 2) a seedling may die, 3) a seed may germinate and 4) a seed or seedling
1020 may be observed. We denote by $t_i^{(g)}$ and $t_i^{(d)}$ the days on which seedling i germinated and died,
1021 respectively.

1022 We assume all seedlings can be partitioned into K classes (e.g. combination of provenance
1023 and growing location) and let $c_i \in \mathcal{C} = \{\mathcal{C}_1, \dots, \mathcal{C}_K\}$ denote the class of seedling i .

1024 Germination

1025 Let g_c be the probability that a seed of class $c \in \mathcal{C}$ is viable and can germinate. Thus, $z_i(0) = \mathcal{L}_S$
1026 with probability g_{c_i} and $z_i(0) = \mathcal{L}_D$ with probability $1 - g_{c_i}$.

1027 We assume the probability that a seed germinates, i.e. that it transitions from $z_i(t-1) = \mathcal{L}_S$
1028 to $z_i(t) = \mathcal{L}_G$, to depend on the growing degree days $h(t)$ up to day t . We model this dependence
1029 using the logistic function

$$1027 \quad \mathbb{P}(z_i(t) = \mathcal{L}_G | z_i(t-1) = \mathcal{L}_S, \mathbf{G}_{c_i}) = \frac{1}{1 + \exp[-\gamma_{c_i}(h(t) - \bar{h}_{c_i})]},$$

1030 where $\mathbf{G}_{c_i} = (g_{c_i}, \gamma_{c_i}, \bar{h}_{c_i})$ denotes the vector of parameters relevant for germination, of
 1031 which the logistic growth rate γ_{c_i} determines the variation in germination among seeds of a
 1032 class and \bar{h}_{c_i} the growing degree days necessary to reach a germination probability of $\frac{1}{2}$.

1033 The probability that seed i germinated on day $t_i^{(g)}$ is given by

$$\begin{aligned} \mathbb{P}(t_i^{(g)} | \mathbf{G}_{c_i}) &= g_{c_i} \mathbb{P}(z_i(t_i^{(g)}) \\ &= \mathcal{L}_{\mathcal{G}} | z_i(t_i^{(g)} - 1) = \mathcal{L}_{\mathcal{S}}, \mathbf{G}_{c_i}) \prod_{t=1}^{t_i^{(g)}-1} \mathbb{P}(z_i(t) = \mathcal{L}_{\mathcal{S}} | z_i(t-1) = \mathcal{L}_{\mathcal{S}}, \mathbf{G}_{c_i}), \end{aligned} \quad (1)$$

1034 where

$$\mathbb{P}(z_i(t) = \mathcal{L}_{\mathcal{S}} | z_i(t-1) = \mathcal{L}_{\mathcal{S}}, \mathbf{G}_{c_i}) = 1 - \mathbb{P}(z_i(t) = \mathcal{L}_{\mathcal{G}} | z_i(t-1) = \mathcal{L}_{\mathcal{S}}, \mathbf{G}_{c_i}).$$

1035 Development

Let $x_i(t)$ quantify the size of a currently growing ($z_i(t) = \mathcal{L}_{\mathcal{G}}$) seedling i on day t . We assume that newly germinated seedlings have a size of zero and grow linearly after germination with class-specific growth rate δ_c such that for a seed i with class c_i

$$x_i(t) = \begin{cases} \delta_{c_i}(t - t_i^{(g)}) & \text{if } z_i(t) = \mathcal{L}_{\mathcal{G}} \text{ and } t_i^{(g)} < t \leq t_i^{(d)} \\ 0 & \text{otherwise} \end{cases}. \quad (2)$$

1036 Note that the variable $x_i(t)$, while referred to as “size”, does not mean an actual measurable
 1037 size (e.g. tree height), but is used here just as a measure of how a seedling is moving along
 1038 development stages. Note also that while modeled linearly, we allow for a non-linear translation
 1039 to life stages (see below).

1040 Death

A currently growing seedling i dies with a probability that is dependent on its current life stage. We assume the highest death rate occurs at germination and then declines exponentially with size. For a seedling i of class c_i we thus have

$$\mathbb{P}(z_i(t) = \mathcal{L}_{\mathcal{D}} | z_i(t-1) = \mathcal{L}_{\mathcal{G}}, \mathbf{D}_{c_i}) = \alpha_{c_i} \exp[-\beta_{c_i} x_i(t)],$$

where $\mathbf{D}_{c_i} = (\alpha_{c_i}, \beta_{c_i})$ denotes the vector of parameters relevant for death and

$$\mathbb{P}(z_i(t) = \mathcal{L}_{\mathcal{G}} | z_i(t-1) = \mathcal{L}_{\mathcal{G}}, \mathbf{D}_{c_i}) = 1 - \mathbb{P}(z_i(t) = \mathcal{L}_{\mathcal{D}} | z_i(t-1) = \mathcal{L}_{\mathcal{G}}, \mathbf{D}_{c_i}).$$

1041 The probability that seed i died on day $t_i^{(d)}$ is given by

$$\mathbb{P}(t_i^{(d)}|t_i^{(g)}, \mathbf{D}_{c_i}) = \prod_{t=t_i^{(g)}+1}^{t_i^{(d)}-1} \left[\mathbb{P}(z_i(t) = \mathcal{L}_G | z_i(t-1) = \mathcal{L}_G, \mathbf{D}_{c_i}) \right] \mathbb{P}(z_i(t_i^{(d)}) = \mathcal{L}_D | z_i(t_i^{(d)}-1) = \mathcal{L}_G, \mathbf{D}_{c_i}). \quad (3)$$

1042 Note that since $x_i(t) = x_i(t-1) + \delta_{c_i}$, these probabilities can be calculated iteratively for
1043 efficiency:

$$\begin{aligned} \exp[-\beta_{c_i} x_i(t)] &= \exp[-\beta_{c_i}((x_i(t-1) + \delta_{c_i}))] \\ &= \exp[-\beta_{c_i} x_i(t-1)] \exp[-\beta_{c_i} \delta_{c_i}]. \end{aligned}$$

1044 This iteration only involves the constant term $\exp[-\beta_{c_i} \delta_{c_i}]$, and no exponentials need to be
1045 evaluated.

1046 Emission probabilities: the observed life stages

1047 Let τ_1, \dots, τ_M denote M times at which seeds were measured, and let $\tau_0 < \tau_1$ denote the
1048 beginning of the experiment. Let $\mathbf{s}_i = (s_i(\tau_1), \dots, s_i(\tau_M))$ denote the life stage identified for
1049 seedling i identified at the observation times with $s_i(\tau_m) \in \mathcal{S} = \{\mathcal{S}_0, \mathcal{S}_1, \dots, \mathcal{S}_S\}$. Of these, \mathcal{S}_0
1050 denotes the unobserved state, i.e. the seedling has not yet germinated or has already died. For
1051 all other stages $\mathcal{S}_s, s > 0$, we model emission probabilities as normalized gamma densities. Let
1052 $w_s(x) = \text{Gamma}(x; m_s, \sigma_s^2)$ be the weight of stage s for seedling size x , where $\mathbf{m} = (m_1, \dots, m_S)$
1053 and $\boldsymbol{\sigma}^2 = (\sigma_1^2, \dots, \sigma_S^2)$ are the modes and variances of the gamma distribution relevant for
1054 stages $s = 1, \dots, S$. Denoting by $\mathbf{E} = (\mathbf{m}, \boldsymbol{\sigma}^2)$ the vector of emission parameters, the emission
1055 probabilities are then given by

$$\mathbb{P}(s_i(\tau_m) = \mathcal{S}_s | z_i(\tau_m), x_i(\tau_m), \mathbf{E}) = \begin{cases} 1 & \text{if } z_i(\tau_m) = \mathcal{L}_S \text{ and } s = 0 \\ 1 & \text{if } z_i(\tau_m) = \mathcal{L}_D \text{ and } s = 0 \\ \frac{w_s(x_i(\tau_m))}{\sum_l w_l(x_i(\tau_m))} & \text{if } z_i(\tau_m) = \mathcal{L}_G \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

1056 We assume that the stage $s_D \in \mathcal{S}$ with the highest mortality is known and set $m_{\mathcal{S}_S} = x_D$.
1057 Note that the model presented is not identifiable, as multiple combinations of modes m_s and
1058 growth rates δ_c can lead to the same emissions. To avoid these non-identifiability issues, we
1059 will arbitrarily assume $m_{\mathcal{S}_S} = x_D = 100$.

The emission probability of the full vector of observations is given by

$$\mathbb{P}(\mathbf{s}_i|t_i^{(g)}, \mathbf{z}_i, \delta_{c_i}, \mathbf{E}) = \prod_{m=1}^M \mathbb{P}(s_i(\tau_m)|z_i(\tau_m), \delta_{c_i}(\tau_m - t_i^{(g)}), \mathbf{E})$$

1060 Probability of full realizations

1061 Given the above model, the seeds follow a Markov model with class-specific parameters
 1062 $\boldsymbol{\theta}_c = (\mathbf{G}_c, \delta_c, \mathbf{D}_c)$ as well as the emission parameters \mathbf{E} . For seed i , we distinguish three
 1063 different types of realizations $\mathbf{z}_i = (z_i(0), \dots, z_i(\tau_M))$ until time τ_M :

- 1064 1. The seed may be viable, germinate between time points $m - 1$ and m and be observed
 1065 for the remainder of the experiment at time points τ_m, \dots, τ_M . We obtain the probability
 1066 of the observations by integrating out the specific day of germination $\tau_{m-1} \leq t_i^{(g)} < \tau_m$:

$$\begin{aligned} \mathbb{P}(\mathbf{s}_i, z_i(\tau_0, \dots, \tau_{m-1}) = \mathcal{L}_S, z_i(\tau_m, \dots, \tau_M) = \mathcal{L}_G | \boldsymbol{\theta}_{c_i}, \mathbf{E}) = \\ \sum_{t^{(g)}=\tau_{m-1}}^{\tau_m} \mathbb{P}(t^{(g)} | \mathbf{G}_{c_i}) \prod_{t=t^{(g)}+1}^{\tau_M} [\mathbb{P}(z_i(t) = \mathcal{L}_G | z_i(t-1) = \mathcal{L}_G, \mathbf{D}_{c_i})] \mathbb{P}(\mathbf{s}_i | t^{(g)}, \mathbf{z}_i, \mathbf{E}, \delta_{c_i}), \end{aligned}$$

1067 where we used the notation $z_i(t_1, t_2, \dots) = \mathcal{L}$ to indicate that the life stage was \mathcal{L} at all
 1068 times t_1, t_2, \dots specified.

- 1069 2. The seed may be viable, germinate between time points $m - 1$ and m , be observed at
 1070 time points τ_m, \dots, τ_n and die prior to time point $\tau_{n+1} \leq \tau_M$. We obtain the probability
 1071 of the observations by integrating out the specific days of germination $\tau_{m-1} \leq t_i^{(g)} < \tau_m$
 1072 and death $\tau_n < t_i^{(d)} \leq \tau_{n+1}$:

$$\begin{aligned} \mathbb{P}(\mathbf{s}_i, z_i(\tau_0, \dots, \tau_{m-1}) = \mathcal{L}_S, z_i(\tau_m, \dots, \tau_n) = \mathcal{L}_G, z_i(\tau_{n+1}, \dots, \tau_M) = \mathcal{L}_D | \boldsymbol{\theta}_{c_i}, \mathbf{E}) = \\ \sum_{t^{(g)}=\tau_{m-1}}^{\tau_m} \mathbb{P}(t^{(g)} | \mathbf{G}_{c_i}) \mathbb{P}(\mathbf{s}_i | t^{(g)}, \mathbf{z}_i, \mathbf{E}, \delta_{c_i}) \sum_{t^{(d)}=\tau_n+1}^{\tau_{n+1}} \mathbb{P}(t^{(d)} | t^{(g)}, \mathbf{D}_{c_i}) \end{aligned}$$

- 1073 3. The seed was never observed at any of the time points τ_1, \dots, τ_M . This may occur due to
 1074 three distinct reasons:

- (a) The seed may be non-viable and stay in state $z_i(t) = \mathcal{L}_D$ for all $t = 0, \dots, \tau_M$. This
 occurs with probability

$$\mathbb{P}(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0, z_i(\tau_1, \dots, \tau_M) = \mathcal{L}_D | \boldsymbol{\theta}_{c_i}, \mathbf{E}) = 1 - g_{c_i}$$

(b) The seed may be viable and never germinate. This occurs with probability

$$\mathbb{P}\left(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0, z_i(\tau_1, \dots, \tau_M) = \mathcal{L}_S | \boldsymbol{\theta}_{c_i}, \mathbf{E}\right) = g_c \prod_{t=1}^{\tau_M} \mathbb{P}(z_i(t) = \mathcal{L}_S | z_i(t-1) = \mathcal{L}_S, \mathbf{G}_{c_i})$$

(c) The seed may be viable, germinate between time points $m - 1$ and m and die before timepoint m . We obtain the probability of the observations by integrating out the specific days of germination $\tau_{m-1} \leq t_i^{(g)} < \tau_m - 1$ and death $t_i^{(g)} < t_i^{(d)} \leq \tau_m$:

$$\begin{aligned} & \mathbb{P}\left(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0, z_i(\tau_0, \dots, \tau_{m-1}) = \mathcal{L}_S, z_i(\tau_m, \dots, \tau_M) = \mathcal{L}_D | \boldsymbol{\theta}_{c_i}, \mathbf{E}\right) \\ &= \sum_{t^{(g)}=\tau_{m-1}}^{\tau_m-1} \mathbb{P}(t^{(g)} | \mathbf{G}_{c_i}) \sum_{t^{(d)}=t^{(g)}+1}^{\tau_m} \mathbb{P}(t^{(d)} | t^{(g)}, \mathbf{D}_{c_i}). \end{aligned}$$

The probability of $\mathbb{P}(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0) | \boldsymbol{\theta}_{c_i}, \mathbf{E}$ is thus given by

$$\begin{aligned} & \mathbb{P}\left(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0 | \boldsymbol{\theta}_{c_i}, \mathbf{E}\right) = \\ & \mathbb{P}\left(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0, z_i(\tau_1, \dots, \tau_M) = \mathcal{L}_D | \boldsymbol{\theta}_{c_i}, \mathbf{E}\right) \\ &+ \mathbb{P}\left(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0, z_i(\tau_1, \dots, \tau_M) = \mathcal{L}_S | \boldsymbol{\theta}_{c_i}, \mathbf{E}\right) \\ &+ \sum_{m=1}^{M-1} \mathbb{P}\left(s_i(\tau_1, \dots, \tau_M) = \mathcal{S}_0, z_i(\tau_0, \dots, \tau_{m-1}) = \mathcal{L}_S, z_i(\tau_m, \dots, \tau_M) = \mathcal{L}_D | \boldsymbol{\theta}_{c_i}, \mathbf{E}\right) \end{aligned}$$

1079 Micro-gardens in which seeds were not individually tracked

1080 In our micro-garden experiment, seeds were not individually tracked and hence only the counts
 1081 of life-stages across a set of seedlings of a specific class grown together as one set are available.
 1082 Let $\mathbf{n}_{j\tau_m} = (n_{j\tau_m 1}, \dots, n_{j\tau_m S})$ denote the vector of stage counts for set j of size $|\mathbf{n}_j|$ at time
 1083 τ_m and let $\mathbf{n}_j = (\mathbf{n}_{j\tau_1}, \dots, \mathbf{n}_{j\tau_M})$. Let us further denote by \mathbf{S} the matrix of stage assignments
 1084 to seeds in set j with elements $S_{[im]}$ denotes the stage of seed i at time τ_m . Rows \mathbf{s}_i of \mathbf{S}
 1085 thus correspond to the observations of seed i . To calculate the probability $\mathbb{P}(\mathbf{n}_j | \boldsymbol{\theta}_{c_i}, \mathbf{E})$ of all
 1086 observations of \mathbf{n}_j of that set across all time points, we need to integrate over all assignments
 1087 of stages to seeds \mathbf{S} :

$$\mathbb{P}(\mathbf{n}_j | \boldsymbol{\theta}_{c_i}, \mathbf{E}) = \sum_{\mathbf{S}} \mathbb{P}(\mathbf{n}_j | \mathbf{S}) \mathbb{P}(\mathbf{S} | \boldsymbol{\theta}_{c_i}, \mathbf{E}),$$

1088 where

$$\mathbb{P}(\mathbf{S}|\boldsymbol{\theta}_{c_i}, \mathbf{E}) = \prod_i \mathbb{P}(\mathbf{s}_i|\boldsymbol{\theta}_{c_i}, \mathbf{E})$$

1089 and

$$\mathbb{P}(\mathbf{n}_j|\mathbf{S}) = \begin{cases} 1 & \text{if } n(\mathbf{S}) = \mathbf{n}_j \\ 0 & \text{otherwise} \end{cases},$$

1090 that is, $\mathbb{P}(\mathbf{n}_j|\mathbf{S}) = 1$ if the stage counts per time point $n(\mathbf{S})$ match the observed counts \mathbf{n}_j and
1091 zero otherwise.

1092 The space of potential stage assignments \mathbf{S} is very large even for a small number of seeds
1093 and cannot be explored in full. We therefore use a Markov approximation to integrate over
1094 stage assignments. Since a naive sampling of state assignments $\mathbf{S} \sim \mathbb{P}(\mathbf{S}|\boldsymbol{\theta}_{c_i}, \mathbf{E})$ results in a
1095 larger number of samples that are incompatible with the data, we generate samples of \mathbf{S} using
1096 importance sampling:

$$\begin{aligned} \mathbb{P}(\mathbf{n}_j|\boldsymbol{\theta}_{c_i}, \mathbf{E}) &\approx \sum_{k=1}^K w_k \mathbb{P}(\mathbf{n}_j|\mathbf{S}_k) \\ \mathbf{S}_k &\sim \mathbb{P}_I(\mathbf{S}|\boldsymbol{\theta}_{c_i}, \mathbf{E}) \\ w_k &= \frac{\mathbb{P}(\mathbf{S}|\boldsymbol{\theta}_{c_i}, \mathbf{E})}{\mathbb{P}_I(\mathbf{S}|\boldsymbol{\theta}_{c_i}, \mathbf{E})}. \end{aligned}$$

1097 We simulate a fraction π of samples under the full model, i.e. by simulating for each seed
1098 whether it germinated or not with probability g_c , and if it did germinate, the day of germination
1099 according to (1) and the day of death according to (3). These samples have weight $w_k = \frac{1}{\pi}$.

1100 We generate the remaining fraction $1 - \pi$ of samples as follows:

- Let $\tilde{\mathbf{g}}_k$ denote a sorted vector of sampled days at which a seeds germinated. The same day may occur multiple times and if all seeds of a set germinated, its length matches that of the number if seeds of a specific class. Let $\tilde{\mathbf{d}}_k$ analogously be a sorted vector of the intervals (prior to the first, after the last or between consecutive observation times) during which death occurred. We fill these containers as follows: We first identify the minimal number of germination events per interval and pick a random germination day within that interval. If three seeds were observed at time observation time m but four seeds at time $m + 1$, for instance, at least one germination event must have occurred between these observation times and we pick a random day according to (1), truncated to that interval:

$$\mathbb{P}(t_i^{(g)}|\mathbf{G}_{c_i}, \tau_m \leq g < \tau_{m+1}) = \frac{\mathbb{P}(t_i^{(g)}|\mathbf{G}_{c_i})}{\sum_{g'=\tau_m}^{\tau_{m+1}-1} \mathbb{P}(t_i^{(g')}|\mathbf{G}_{c_i})}.$$

1101 We analogously identify the minimal number of death events that occurred in each interval
1102 and add those intervals to $\tilde{\mathbf{d}}_k$.

1103 While it is unknown if and when the remaining seeds germinated, the number of additional
1104 germination and death events must be equal in each interval. For each of the unaccounted
1105 seeds we therefore i) first simulate if the seed is germinating or not with probability g_c ,
1106 ii) then, in case the seed germinated, simulate a random germination day (without any
1107 interval restriction) according to (1) and add a death event to the interval in which the
1108 germination was simulated.

- 1109 • We next simulate sample pairs of germination dates and death intervals given the constraint
1110 that germination must predate death. We start with the earliest death interval \tilde{d}_{k1} of
1111 $\tilde{\mathbf{d}}_k$ (the first since the vector is sorted) and randomly pick a germination date from $\tilde{\mathbf{g}}_k$
1112 that is compatible with a death in that interval, i.e. from $\tilde{\mathbf{g}}'_k = (\tilde{g}_{ki} | \tilde{g}_{ki} < \tau_{\tilde{d}_{k1}})$. Once
1113 the germination data was simulated, we simulate a valid death date in the chosen death
1114 interval that postdates the simulated germination date. We repeat that process for each
1115 death interval in sequence.
- 1116 • We calculate the importance weights w_k as $\frac{1}{1-\pi}$ times the ration between the probability of
1117 the simulated \mathbf{S}_k under the full model and its probability under the importance sampling
1118 scheme outlined above.

In each iteration, we initially generate 20,000 samples and determine the effective sample size ESS of the samples generated as

$$ESS = \frac{(\sum_k w_k)^2}{\sum_k w_k^2}.$$

1119 We then generate an additional batches of 20,000 samples until the ESS of all samples combined
1120 exceeds 2,000.

1121 Implementation

1122 We implemented the above model and inference scheme as a command-line C++ program
1123 `tree_growth` using the library `statools`. The implementation, along with a brief user manual,
1124 is available at https://bitbucket.org/wegmannlab/tree_growth/. All estimations in this
1125 paper were done with commit `a9ae485`.