

Demographic inference in a spatially-explicit ecological model from genomic data: a proof of concept for the Mojave Desert Tortoise

Jaime Ashander¹, Peter Ralph², Evan McCartney-Melstad^{1,3}, H. Bradley Shaffer^{1,3}

2018-06-22

¹: Department of Ecology and Evolutionary Biology
610 Charles E. Young Drive East
University of California, Los Angeles
Los Angeles, CA 90095
USA

²: Department of Mathematics and Institute of Ecology and Evolution
University of Oregon
Eugene, OR 97403
USA

³: La Kretz Center for California Conservation Science
Institute of the Environment and Sustainability
University of California, Los Angeles
Los Angeles, CA 90095
USA

Abstract: In this paper, we study the general problem of extracting information from spatially explicit genomic data to inform inference of ecologically and geographically realistic population models. We describe methods and apply them to simulations motivated by the demography of the Mojave desert tortoise (*Gopherus agassizii*). The tortoise is an example of a long-lived, threatened species for which we have an excellent understanding of range, habitat preference, and certain aspects of demography, but inadequate information on other life history components that are important for conservation management. We use an individual-based model on a discretized geographic landscape with overlapping generations and age and sex-specific dispersal, fecundity, and mortality to develop and test a method that uses genomic data to infer demographic parameters. We do this by seeking parameters that best match a set of spatial statistics of genomes, which we introduce and discuss. We find that for inferring only overall population density and mean migration distance, a simple statistical learning method performs well using simulated training data, inferring parameters to within 10% accuracy. In the process, we introduce spatial analogues of common population genetics statistics, and discuss how and why they are expected to contain signal about the geography of population dynamics that are key for ecological modeling generally and conservation of endangered taxa.

Keywords: population genomics; inference; landscape genomics; forward-time simulations; individual based model; demography

1 Introduction

Mechanistic population models are a key tool used by basic and applied ecologists to understand the history and dynamics of natural populations. Population models inform fisheries management (Quinn and Deriso 1999), conservation of endangered species (Caswell 2001), and understanding of emerging infectious diseases (Diekmann and Heesterbeek 2000). Population models are well suited to address both fundamental questions (e.g., how population regulation occurs in spatially extensive, age-structured populations) and applied concerns (e.g., movement of disease vectors across political geographies).

43 Often, a population model is used to project future abundance under a variety of scenarios that affect one
44 or more parameters (e.g., Prates et al. 2016; Benson et al. 2016). For example, a model that includes the
45 effect of temperature on egg hatching rate could be used to project the impact of climate change on future
46 population survival. To provide meaningful information for management, model parameters must be known
47 with enough certainty that one can realistically distinguish the future effects of management or landscape
48 modification scenarios.

49 Demographic quantities such as survival, growth, and fecundity can often be estimated by direct field
50 observation (generally via mark-recapture studies) particularly for short-lived species where data can be
51 collected over several generations. However, the degree to which current demographic parameters accurately
52 reflect long term values is often unknown, particularly when those parameters may fluctuate substantially
53 across time scales or geography. These demographic quantities determine abundance fluctuations and gene
54 flow across the landscape, two processes with conceptually well-understood effects on patterns of genetic
55 relatedness. Therefore, genomic data provide a promising source of additional information to bridge this gap.
56 However, there has thus far been relatively little use of genomic data in fitting mechanistic ecological models,
57 even though there must be a direct relationship between population dynamics and the geographical patterns
58 of standing genetic variation observed in nature.

59 A major barrier to integrating genomic data in ecological models is a lack of analytical results that describe
60 genetic patterns expected under geographically explicit population models. Genomic data are often used for
61 descriptive models – most commonly, either clustering-based methods that seek to identify substructure in a
62 population (e.g., Pritchard, Stephens, and Donnelly 2000; Bradburd, Coop, and Ralph 2017), “resistance”
63 methods that depict genetic similarity using a landscape descriptor of gene flow (e.g., McRae 2006; Petkova,
64 Novembre, and Stephens 2016; Shaffer et al. 2017), or least-cost path analysis to find most likely routes
65 of gene flow (e.g., Wang, Savage, and Shaffer 2009). Although these approaches can provide information
66 about migration rates among a set of discrete populations (Greenwald 2010), none of these methods provide
67 estimates in units that are directly interpretable as describing population dynamics in a generative model of
68 continuous space, such as mean distance traveled by dispersing individuals per year, or number of adults
69 per square kilometer. It is extremely well-understood how demography determines genetic patterns in
70 large, randomly mating populations (e.g., the Wright-Fisher model), but when realistic geography and its
71 idiosyncratic effects are introduced, few analytical predictions are available (but see Ringbauer, Coop, and
72 Barton 2017).

73 Simulations have proven useful in bridging this gap between ecological models and genomic data. A variety
74 of simulation approaches can shed light on evolutionary and some ecological processes (Hoban, Bertorelle,
75 and Gaggiotti 2012). Such studies generally use likelihood free (e.g., Approximate Bayesian Computation)
76 approaches that require choosing summary statistics that describe the high-dimensional outputs (genomic
77 or genome-like data). These simulation tools allow for inference of migration among complex, but discrete,
78 spatially structured populations. For example Vallée, Luciani, and Cox (2016) used a general-purpose
79 individual-based modelling (IBM) library to model the dynamics of 37 recombining markers across all human
80 chromosomes during the Neolithic expansion among Southeast Asian islands. Alves et al. (2016) inferred
81 short- and long-distance dispersal in Eurasian Neolithic expansions of humans using SPLATCHE2 (which
82 combines forward simulations of population sizes with coalescent simulations of genetic data, Ray et al. 2010).
83 Similarly, Prates et al. (2016) inferred past demography of neotropical forest lizards, and S. E. Harris et
84 al. (2016) inferred recent population structure (and correlating it with urbanization) in white-footed mouse
85 (*Peromyscus leucopus*) in the Northeastern USA, both using *fastsimcoal2* (Excoffier and Foll 2011).

86 Here, we use an ecologically realistic, individual-based model to simulate whole genomes of a closed population
87 across a heterogeneous landscape. By simulating across a range of parameters and comparing results from
88 our model to genomes obtained from real populations, we can make inferences about the parameter values
89 corresponding to these real populations. Actual population-scale landscape simulations of individuals with
90 gigabase-sized genomes pushes the limits of current computational feasibility. However, our method is made
91 computationally feasible with large population sizes by recent advances in simulation methods (Kelleher et
92 al. 2018), that allow more rapid simulation and simultaneously record the genealogical relationships across
93 generations for all individuals across the simulated landscape.

94 This study is motivated by the Mojave desert tortoise, *Gopherus agassizii* and the need to create realistic spatial
95 models to guide its conservation. The species lives across much of the Mojave desert in the Southwestern USA,
96 and is threatened due to a combination of habitat destruction, mortality due to human-subsidized predators
97 (Kristan and Boarman 2003; Esque et al. 2010), disease (M. B. Brown et al. 1994), vehicle-associated
98 mortality (W. Boarman and Sasaki 2006), and other factors (Berry 1986; USFWS 2011). A substantial body
99 of ecological fieldwork now characterizes desert tortoise habitat suitability (K. E. Nussear et al. 2009; USFWS
100 2011), and several sizeable demographic studies have estimated sex- and age-specific mortality and fecundity
101 (e.g., Doak, Kareiva, and Klepetka 1994; Karl 1998; Reed, Fefferman, and Averill-Murray 2009). However,
102 certain aspects of tortoise life history – in particular, the effects of juvenile and long-distance dispersal –
103 remain relatively unknown. Since dispersal-mediated gene flow should leave strong signals across the genome,
104 we can reasonably hope that genomic data could inform a mechanistic understanding of tortoise movement
105 across the landscape.

106 In this paper, we (1) Develop a landscape-scale individual-based model (IBM) simulation that maps ecological
107 parameters to a population pedigree; (2) Introduce a general class of spatial population genetic statistics
108 and motivate their use for inference problems such as those modeled here; (3) Develop a statistical method
109 to estimate dispersal and population density by comparing patterns of relatedness on the landscape to
110 simulations of expected relatedness; and (4) Use simulated data to show that our method can simultaneously
111 estimate dispersal and density to within 10% percent of their true values, as long as the dispersal scale is not
112 too large.

113 We then infer two parameters from data produced under the model used for inference, presenting a limited
114 test of the method. This reflects a common scenario where a great deal is known about certain ecological and
115 demographic processes and the goal is to add information from genomic data. We establish a geographically
116 and ecologically realistic population model, develop methods to produce consistent discretizations of the
117 model, explore a class of spatial genetic statistics, determine procedures to match these spatial statistics
118 between datasets, and assess our statistical power across a range of model parameters.

119 The statistical problem we face here is an *inverse problem*, conceptually similar to estimating the migration
120 rate between two randomly mating populations of size N by inverting the analytic relationship $F_{ST} = \frac{1}{4Nm+1}$
121 (Wright 1951), where F_{ST} can be computed from genetic data, and N and m are the effective population
122 size and migration fraction, respectively. However, the functional relationship between genetic statistics
123 and parameters in models of continuous, heterogeneous geography is generally unknown (but see N. H.
124 Barton, Depaulis, and Etheridge 2002; Ringbauer, Coop, and Barton 2017). Much like Approximate Bayesian
125 Computation or sequential Monte Carlo (Marjoram 2013), we use simulation to bypass this issue. We simulate
126 under a range of values of parameters for dispersal and density, calculate a large number of spatial population
127 genetics statistics of the resulting data from each, and use general-purpose statistical learning to approximate
128 the inverse map from statistics back to the (originally unknown) two parameters. In so doing, we demonstrate
129 how spatially-explicit landscape genomic data can be used parameterize population biology models that
130 should be useful across a wide range of applications and taxa.

131 2 Materials and methods

132 Although there are a great number of possible aspects of an ecological model that could be inferred, we focus
133 here on a simple case. Suppose we have a set of georeferenced whole genomes sampled from individuals
134 across a population range, and that features of this spatial range have been recorded and interpreted to
135 yield a measure of habitat suitability for the population of interest. However, both the rate of movement of
136 individuals across geography (dispersal) and the population density remain unknown. We assume (i) that a
137 local measure of carrying capacity is equal to ρ individuals per hectare, multiplied by local habitat suitability,
138 and (ii) that the yearly movement of individuals is Gaussian with a standard deviation of σ meters. This
139 leaves only two scalar parameters to be estimated: ρ and σ .

140 We model a landscape consisting of two large areas of high-quality habitat, which taper to low-quality
141 habitat at their edges, connected by a narrow isthmus of low quality habitat (Figure 1A). Locations of

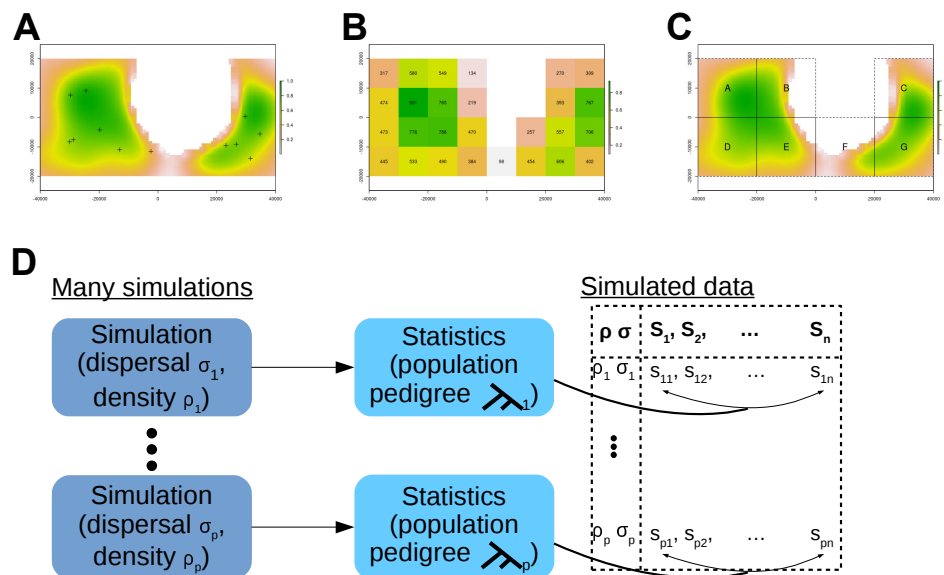


Figure 1: A) Spatial setting and example data, with axes labelled in meters and colors indicating habitat quality (1.0, or green, corresponds to highest-quality habitat, 0.0, or white, corresponds to impassable terrain) across the continuous geography. Individual samples are marked with '+'. For clarity only 12 samples are shown. B) The landscape discretized into patches used in simulations, defined by aggregating the fine scale map (panel A), and labelled with their carrying capacity of mature individuals (for $\rho = 0.1$). C) The landscape map partitioned into regions, labeled with letters; population genetic statistics are computed on groups of samples with group membership determined by the region in which samples occur. D) Simulating data for inferring dispersal σ and density ρ . We calculate n statistics for each of p independent simulations.

142 individuals whose genomes have been sampled are marked by '+' (in a real dataset, these would be fixed by
 143 the sampling location). From these individual samples, we seek to compute statistics that are informative of
 144 the population's demographic parameters, ρ and σ .

145 2.1 Geographic genetic statistics

146 To generate a wide class of potentially informative statistics, we use several population genetic statistics as
 147 *spatial statistics*, including Patterson's F -statistics (Reich et al. 2009; Peter 2016). The F -statistics were
 148 originally used to compare variation among discrete, randomly-mating populations. In that context, the
 149 statistics convey information about admixture and shared branch lengths in the "population phylogeny"
 150 (Moorjani et al. 2013; Reich et al. 2009; Peter 2016) that describes how the populations are related to each
 151 other. Since we use these statistics in a nonstandard way, we now define them and motivate their use as
 152 informative spatial statistics.

153 For a set of genomes, denoted A , we write the *genetic diversity* of A , i.e., the mean density of nucleotide
 154 differences between a randomly chosen pair of genomes from A , as $\pi(A) = \mathbb{E}[|a_1 - a_2|]$, where a_1 and a_2 are
 155 alleles at a random site in the genome, coded as 0 or 1, from genomes randomly chosen without replacement
 156 from A . We also denote the *genetic divergence* between two groups, A and B , as the mean density of
 157 nucleotide differences between randomly chosen individuals from the two groups, which can be written as
 158 $\pi(A, B) = \mathbb{E}[|a - b|]$, where a and b now come from randomly chosen genomes in A and B respectively (if
 159 any individual is in both groups, sample without replacement so $\pi(A) = \pi(A, A)$). Patterson's F statistics
 160 can be written in the same way, but using *four* genomes: $F_4(A, B; C, D) = \mathbb{E}[(a - b)(c - d)]$, where a , b ,
 161 c , and d are now alleles of randomly chosen genomes from the four groups A , B , C , and D respectively.
 162 Similarly, $F_3(A; B, C) = \mathbb{E}[(a_1 - b)(a_2 - c)]$ and $F_2(A, B) = \mathbb{E}[(a_1 - b_1)(a_2 - b_2)]$, where now a_1 and a_2
 163 are randomly chosen without replacement from A (and likewise for b_1 and b_2). Note that we can write

164 $F_3(A; B, C) = F_4(A, B; A, C)$ if we interpret the latter as sampling without replacement, as we did for
165 $\pi(A, B)$; for this reason, subsequently we write all F statistics using this format (and dropping the subscript
166 ‘4’).

167 We also introduce analogous statistics that depend on choices of *three* genomes. We will write these in
168 terms of $y(A; B, C) = \mathbb{E}[a(1-b)(1-c) + (1-a)bc]$, the probability that a sample from A differs from
169 two other samples, one from B and one from C . The three-point statistics we use derive from y , but are
170 modified to be zero in a randomly mating population: $Y(A; B, C) = y(a; b, c) - (1/2)(y(b; a, c) + y(c; a, b))$,
171 and $Y_2(A; B) = Y(A; B, B) = y(a; b_1, b_2) - (1/2)(y(b_1; a, b_2) + y(b_2; a, b_1))$.

172 An alternative way to think of these statistics is as estimates of weighted averages of branch length across
173 all genealogical trees relating individuals selected from two or more groups, and scaled by mutation rate
174 (Ralph 2015; Peter 2016). Averaging over marginal gene-trees under an infinite sites model of mutation, and
175 omitting a scaling factor of the mutation rate, the corresponding “branch length” quantities (denoted with a
176 superscript (b)) are:

- 177 • $\pi^{(b)}(A, B)$: the average sum of the lengths of the two branches going from a and b back to their most
178 recent common ancestor (MRCA), averaged over trees and choices of a and b .
- 179 • $\pi^{(b)}(A)$: the same as $\pi^{(b)}(A, B)$ but with both genomes chosen from A .
- 180 • $Y^{(b)}(A; B, C)$: the difference between (the average length of any branches that separate a from b and c)
181 and (one-half of the sum of the lengths of any branches that would separate either b or c from the other
182 two), averaged over trees and choices of a, b , and c .
- 183 • $F^{(b)}(A, B; C, D)$: the difference between (the average length of any branches that separate a and c from
184 b and d) and (the average length of any branches that separate a and d from b and c), averaged over
185 choices of a, b, c, d .

186 To avoid scaling factors and to make this correspondence exact, we measure branch lengths in *expected*
187 *number of mutations*, i.e., scaling branches by the mutation rate per unit time. For instance, since $\pi(A, B)$ is
188 the average number of mutations per site that have occurred between a or b and their MRCA, this makes
189 the expected value of $\pi(A, B)$ equal to $\pi^{(b)}(A, B)$ under an infinite-sites model of neutral mutations. These
190 relationships between statistics computed using genotypes and the summaries of branch lengths are shown in
191 Figure 2.

192 The F and Y statistics are defined so that they have expected values of zero if samples are exchangeable (e.g.,
193 if they all come from a single randomly mating population), because in this case each topology is equally
194 frequent and has the same distribution of branch lengths, so the contributions of each topology cancel. Figure
195 2 also shows formulas for the statistics in terms of divergence.

196 To help develop an intuition for how the statistics work in continuous geography, consider the situation where
197 population density is constant, and movement in any direction is equally likely. Define groups of individuals
198 by whether they fall in different geographic regions. Then, since regions of equal area have equal population
199 size, distance between regions determines how closely connected they are by dispersal – i.e., how likely is a
200 recent MRCA – analogous to the migration rate for discrete populations.

201 This makes it possible to consider the relative frequencies with which various marginal genealogies occur, and
202 hence in which settings Y or F are expected to have positive, zero, or negative values. The frequency with
203 which the first coalescence occurs provides a particularly good heuristic. For example, if group A ’s region lies
204 physically between the regions of groups B and C , then $Y(A; B, C)$ will be negative due to a deficit of trees
205 with topology $(A, (B, C))$, which has overall positive weight, as shown in Figure 2B. Conversely, if B or C lie
206 between the other two, then $Y(A; B, C)$ will be positive. Figure 3 shows several geographic configurations of
207 regions and corresponding effects on F and Y . Across sites where none of the samples are closely related –
208 i.e., there is no *recent* coalescence – all three rooted topologies in Figure 2B are roughly equally likely and
209 have equal branch lengths on average (Wilkins 2004), and will therefore cancel out.

210 Population size also affects the probability of first coalescence. For example, $F_2(1, 2) = F(1, 2; 1, 2)$, which
211 corresponds to Figure 2(C) with $A = C = 1$ and $B = D = 2$. If within-group coalescence is more likely than
212 between-group (as is often the case), then unrooted topology $(AC)(BD) = (11)(22)$ is the most common,
213 implying that $F_2(1, 2) > 0$. However, an increased population density in one region reduces the magnitude of

214 F_2 , since it makes it more likely that two lineages in that region trace back to ancestors outside the region
215 before they coalesce, decreasing this bias. Thus, increased population density is expected to decrease the
216 value of the statistic. Figure 4 shows how varying region size (i.e., group population size) affects F and Y .

217 **How many statistics?** If we have divided our samples into k groups we can compute a large number of
218 statistics. The statistics we consider here are averages over choices of two, three, or four genomes chosen
219 from varying numbers of groups. We call a statistic that depends on k groups a “ k -point statistic” since we
220 imagine each group standing in for a spatial location. “Groups” may be arbitrary: single genomes, single
221 diploids, or larger collections. The total number of statistics of each type is:

- 222 • diversity, $\pi(A) = \pi(A, A)$, a one-point statistic, k .
- 223 • divergence, $\pi(A, B)$, a 2-point symmetric statistic, $k(k - 1)/2$.
- 224 • $F_2(A, B)$, a 2-point, symmetric statistic, $k(k - 1)/2$.
- 225 • $Y_2(A; B)$, a 2-point statistic, $k(k - 1)$.
- 226 • $F_3(A; B, C)$ and $Y_3(A; B, C)$, both 3-point statistics symmetric in B and C , $k(k - 1)(k - 2)/2$.
- 227 • $F_4(A, B; C, D)$, a 4-point statistic with one symmetry and two anti-symmetries, $k(k - 1)(k - 2)(k - 3)/8$.

228 As the number of groups grows, it becomes impossible to compute all possible statistics in reasonable time.
229 For example, with groups of size 20 each statistic takes 15 seconds to compute for a 10Mb region of a human
230 genome (using Python tools from Kelleher, Etheridge, and McVean 2016), so for $k = 30$ it would take roughly
231 451 hours to compute all 108,345 statistics. For a fixed amount of computing power we must either keep
232 the number of groups reasonably small and compute all statistics, or choose sets of statistics to compute,
233 motivated by the biology or geography of the landscape.

234 **Branch lengths or sequence?** The theory outlined above equates each statistic to a corresponding
235 summary of branch lengths in marginal genealogies. Our simulations actually record all marginal genealogies
236 that relate sampled individuals to one another at every point on the genome. As described in Kelleher et al.
237 (2018), this is done for speed, but it has the benefit that we have access to the underlying marginal gene-trees.
238 This means that we can directly compute expected values of the statistics on branch lengths. The alternative
239 is to compute them using genome sequence, generated for simulations by placing mutations on the marginal
240 genealogies. (Since neutral mutations do not, by definition, affect genealogies, placing mutations on the trees
241 *post hoc* is equivalent to generating them as the simulation progresses.) For efficiency we take the former
242 path, working directly with statistics calculated using the branch lengths of the underlying genealogies from
243 simulations (e.g., using $F^{(b)}$ rather than F). We expect performance with sequence-based statistics to be
244 identical, because the deviation of the sequence-based statistic from the underlying tree-based statistic has
245 mean zero with standard deviation inversely proportional to the square root of the sequence length (Ralph
246 2015) – in practice, they will be quite close for large data sets. A set of simple simulations confirms this
247 predicted close agreement between the two methods (See Supplement B and Figure B.1). With empirical
248 data, the statistics *must* be computed using sequence differences, but these are directly comparable to the
249 branch length statistics after a rescaling.

250 2.2 A statistical method to infer population density and dispersal

251 Our overall goal is to estimate dispersal (σ) and population density (ρ) based on the statistics computed
252 from a given genomic dataset with known geographical coordinates. Since the focus of this work is on the
253 establishment of an ecologically realistic model and computation of informative spatial statistics, we do this
254 in a relatively simple way. First, we simulate from an individual-based model across a grid of parameter
255 values. Second, we compute statistics on each output. After this procedure (shown in Figure 1D) we obtain
256 a table containing inputs (values of ρ and σ) and corresponding outputs, i.e., the statistics we calculate,
257 denoted generically here as $(s_1 \dots s_n)$. We seek to infer the relationship between inputs and outputs and do
258 so using inverse interpolation.

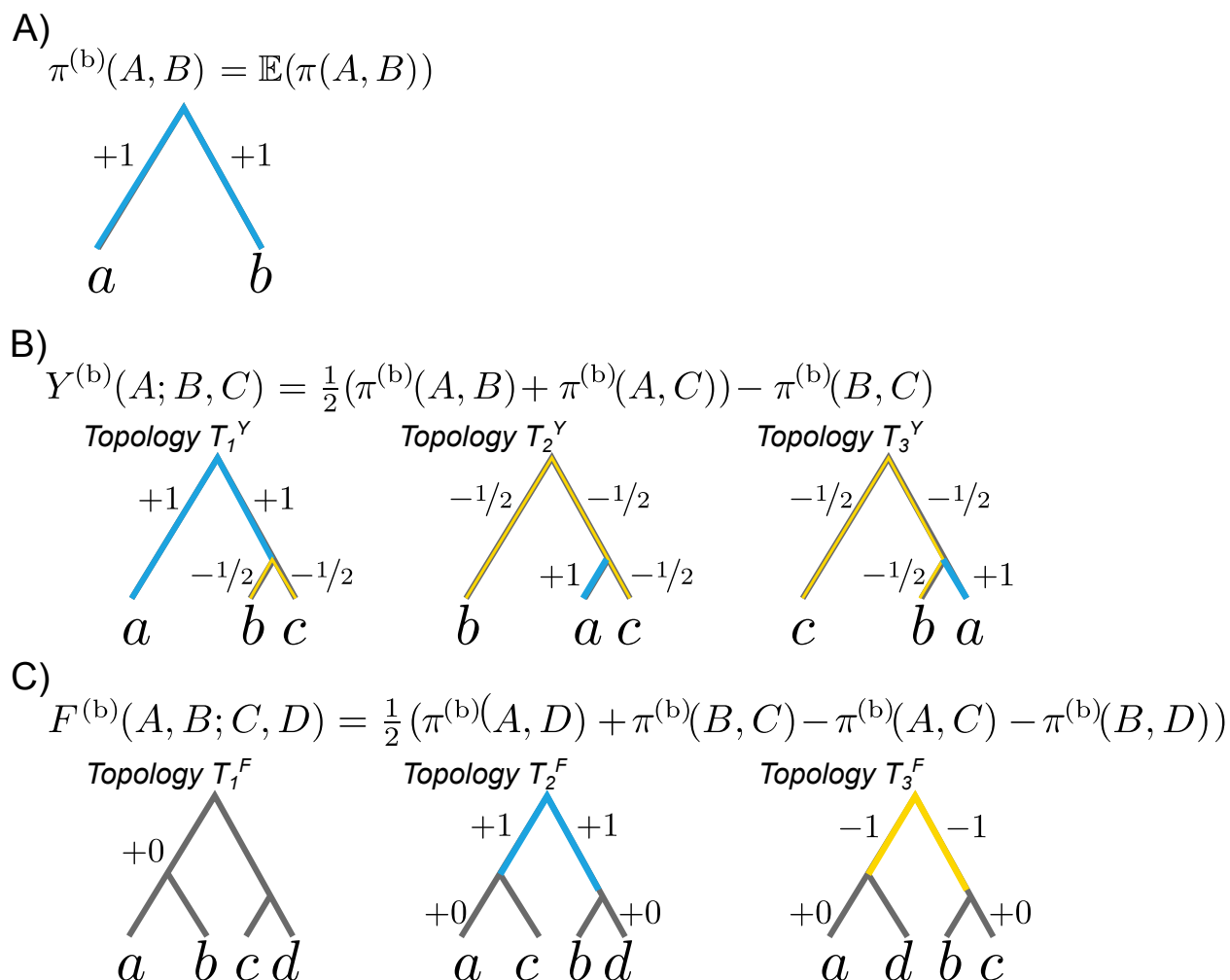


Figure 2: Genomic statistics between groups can be computed as a sum of weighted branch lengths across all gene trees; shown are the weights for the three types of statistic: (A) genetic diversity π , (B) the three-point statistic Y , the four-point statistic F . Branches are measured in number of mutations (for the usual statistics) or in units of expected mutations (for the 'branch length' versions). In the formulas, $\pi(A, B) = \pi^{(b)}(A, B)$ denotes the mean tree distance measured in terms of expected mutations from a random sample in A to a random sample in B , averaged over trees and choices of samples. Weights may depend on the tree topology and are marked on each branch; positive contributions are shown in blue and negative contributions are shown in yellow; gray is zero contribution. The equations in (B) and (C) show how the branch weights for $Y^{(b)}$ and $F^{(b)}$ in depend on those for $\pi^{(b)}$ in (A). Equating $\pi^{(b)}(A, B)$ with the path between a sample from A and a sample from B yields the weights. For instance, in T_1^Y the weight of $-1/2$ on the branch above b (ii) is obtained by adding $+1/2$ (because it is on the path from a to b) to -1 (since it is on the path from b to c).

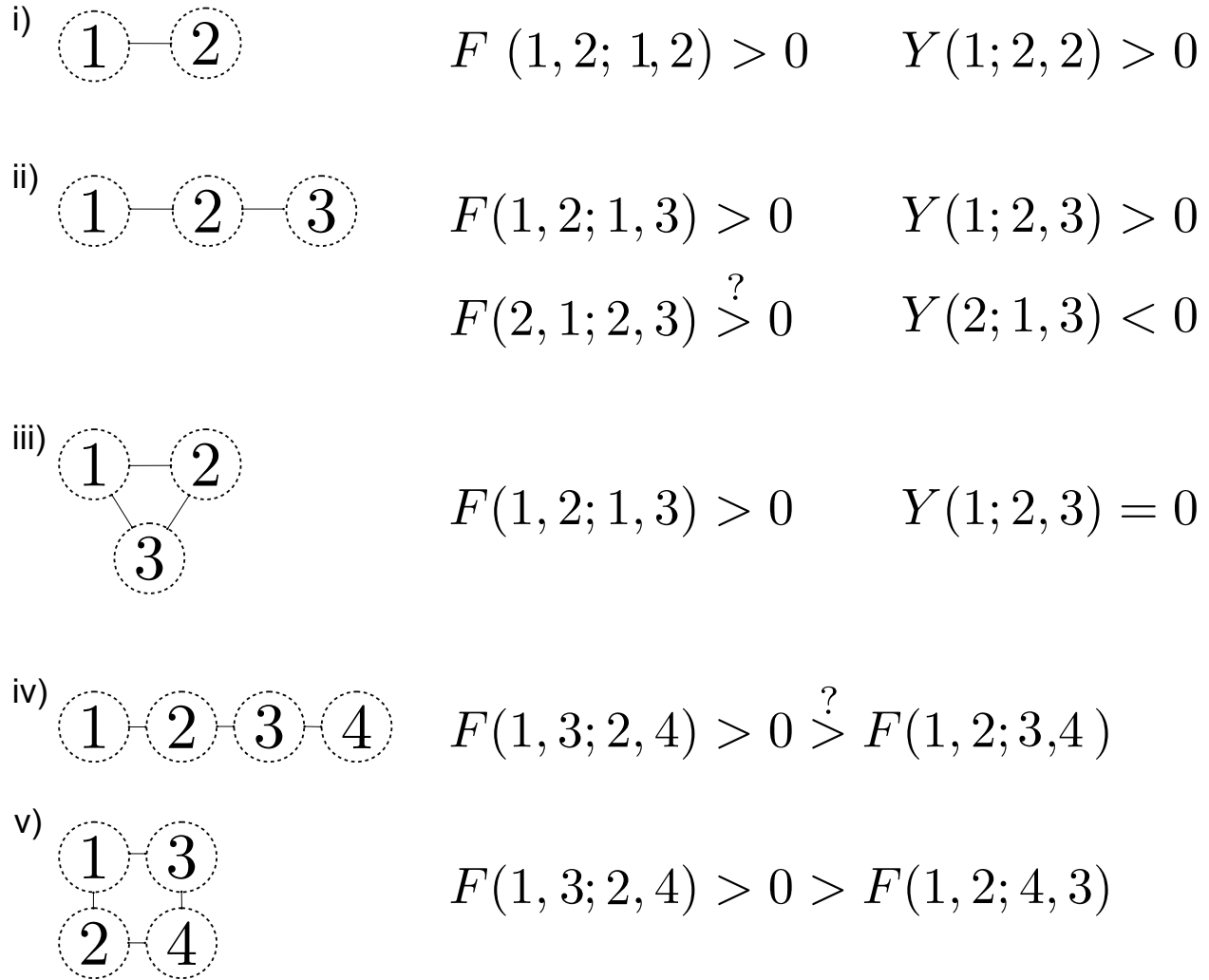


Figure 3: Effects of differing spatial configurations of regions in isotropic space on the signs of Y and F statistics. In all cases the sign is derived by reasoning about the probability of the first coalescence involving two regions and then connecting this to the probability that topologies occur in Figure 2. All population sizes (i.e., areas) are equal. Larger distances between regions in isotropic space result in decreased probability that the first coalescence involves these regions. Equivalent distances between equal-sized regions results in equal probability that the first coalescence involves these regions. If the actual sign is in doubt, the comparator is marked by ‘?’. Recall that $F_3(1; 2, 3) = F(1, 2; 1, 3)$ and that $F_2(1, 2) = F(1, 2; 1, 2)$. See Supplement A for justification.

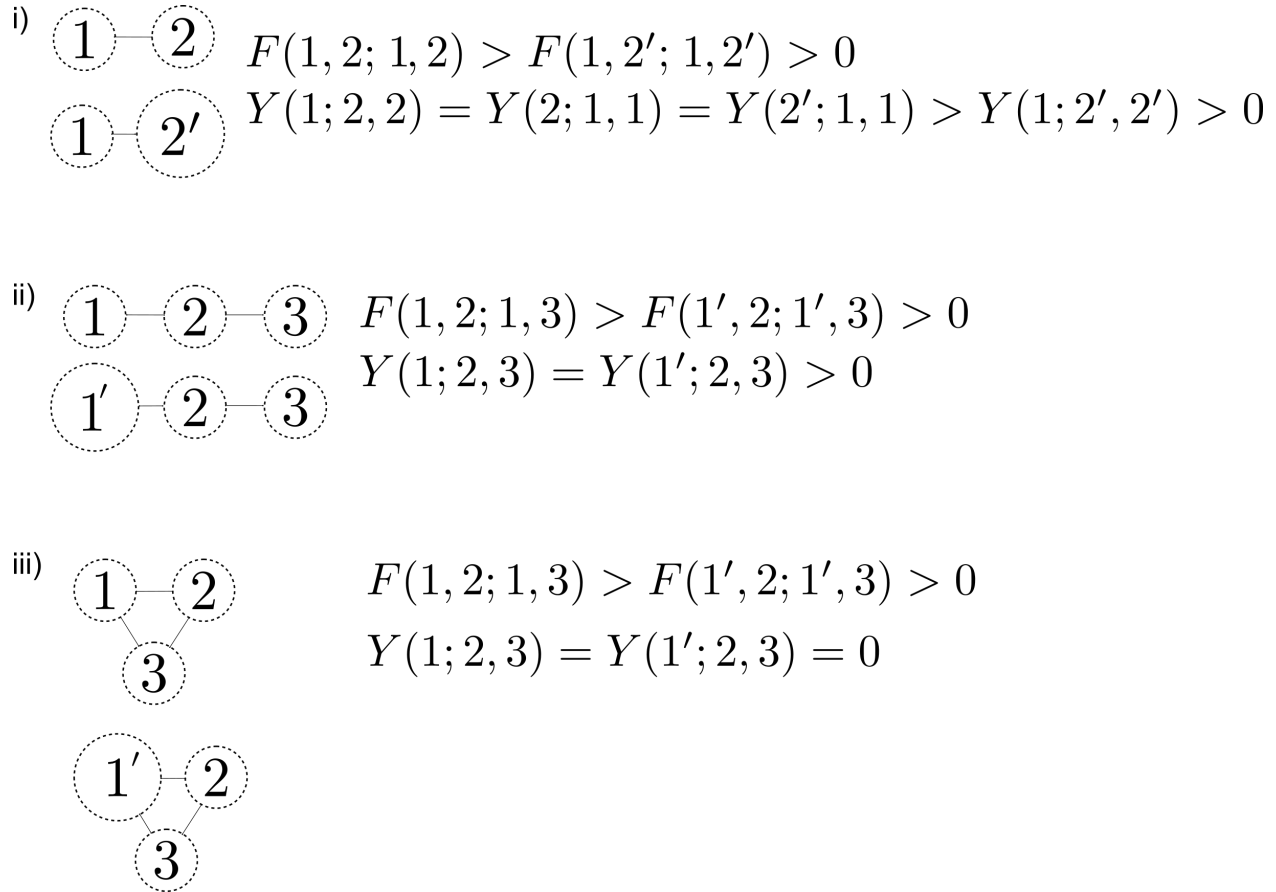


Figure 4: Effects of differing sizes and configurations of regions in isotropic space on Y and F statistics. In all cases the sign is derived by reasoning about the probability of the first coalescence involving two regions and then connecting this to the probability that topologies occur in Figure 2. All distances between centroids of adjacent populations are assumed equal. Regions denoted by larger circles have twice the population size (e.g., twice the population density) of those denoted by smaller circles. The probability of the first coalescence involving two members from a larger population is less than that of the first coalescence involving two members of a smaller population. Recall that $F_3(1; 2, 3) = F(1, 2; 1, 3)$ and that $F_2(1, 2) = F(1, 2; 1, 2)$. See Supplement A for justification.

259 2.2.1 Inference via inverse interpolation

260 Consider using real genomic data to infer migration. In a classic Wright-Fisher model there is a clean
261 parametric dependence of a genetic statistic, F_{ST} , on migration rate m and population size N . For our more
262 complex model there is an unknown analytic relationship between σ , ρ , and the statistics we introduce above.
263 More generally, simulations give us noisy observations of an unknown function $f(\theta)$ that maps parameters
264 (here, $\theta = (\sigma, \rho)$) to the n statistics: for each simulation, run with parameters θ_i , we can think of the resulting
265 set of statistics as

$$S_i = (s_{i1}, \dots, s_{in}) = f(\theta_i) + \epsilon_i \\ \epsilon_i \sim N(0, \Sigma),$$

266 where Σ is an unknown covariance matrix defining how the noise ϵ_i is correlated across observations i . Given
267 a new set of statistics \tilde{S} , we then seek to estimate the corresponding parameters, $\tilde{\theta}$.

268 For our purposes, it suffices to take an average over the known parameter values θ_i , each weighted according
269 to the proximity of its associated statistics S_i to the observed \tilde{S} :

$$\tilde{\theta} = \sum_i \theta_i \frac{\exp(-\|\tilde{S} - S_i\|^2/2\omega^2)}{\sum_j \exp(-\|\tilde{S} - S_j\|^2/2\omega^2)}$$

270 We choose the bandwidth, ω , by k -fold crossvalidation. To do this, we randomly divide the data S and θ into
271 k blocks; denote by $S^{(i)}$ the i -th block and S^{-i} the rest of the data. Then, for each bandwidth ω and each
272 $1 \leq i \leq k$, use the other data S^{-i} and θ^{-i} to predict parameters for every entry in the i th block $\tilde{\theta}^{(i)}$, and
273 then compute the mean relative error in the i th as:

$$RE_i = 1/2 \sum_j |\tilde{\theta}_j^{(i)} - \theta_j^{(i)}|/|\theta_j^{(i)}|.$$

274 We choose the bandwidth ω that minimizes the median relative error across all k validation blocks.

275 2.2.2 Individual-based simulations of a desert tortoise population

276 Our focal species, the Mojave desert tortoise (*Gopherus agassizii*) is widespread and exists as a collection
277 of semi-discrete populations spread across the landscape. For species like this, one of the most important
278 parameters for any spatial population model—dispersal—is also one of the most challenging to estimate
279 empirically. In *G. agassizii* adults are generally dependent on burrows and therefore relatively stationary.
280 But tortoises are also long-lived, so rare adult dispersal may provide significant, but rarely observed, genetic
281 connectivity. At the same time, juveniles are secretive and have low survival, and movement is presumably
282 much more common but also much more difficult to observe. Thus the long life, sparse distribution, and
283 low juvenile survival rates all make it very challenging to accurately estimate the lifetime-averaged dispersal
284 rate directly. Data from direct observations can provide some insights, particularly on short-term movement
285 (e.g., radio telemetry, Nafus et al. 2017), but analyses of genetic samples from across the population’s range
286 can potentially tell us much more about the recent history of movement in the population. The genomes
287 of individuals encode information about their relatedness both across space and back through time. With
288 genomic data from many individuals it is possible to accurately estimate statistics that describe these patterns
289 of genetic relatedness and thus reflect movement and post-migration breeding of these elusive animals.

290 We developed a landscape-scale individual-based model (IBM) simulation of a closed population of inter-
291 breeding individuals, including relatedness across the genome. The model includes several demographic
292 complexities, including sex (individuals are diploid, with sex fixed at birth), age structure (i.e., varying
293 survival and fecundity by age class), density-dependence, movement in space, and maturity (which increases
294 survival). Our IBM simulates a closed population, which could be an entire species.

295 We implement the model by discretizing continuous space into a grid of discrete patches, where the patches
296 are contiguous and each represents a subpopulation (Figure 1B). In the simulation below, each of these

297 includes approximately 50-500 individuals. Within each patch, we assume that individuals mate randomly.
298 We also implement population regulation within the patches.

299 Movement between patches is parameterized by the standard deviation of the dispersal kernel, σ . We compute
300 the *migration matrix* M whose (i, j) th entry is the probability that an individual in patch i moves to patch j
301 in a given year. This is computed as the probability that an individual uniformly located in patch i moves a
302 random, Gaussian-distributed distance with mean zero and standard deviation σ and ends up in patch j ; if
303 the corresponding geographic regions are denoted A_i and A_j then this can be computed as

$$M_{ij} = \int_{A_i} \int_{A_j} \frac{1}{2\pi\omega^2} e^{-\frac{|x-y|^2}{2\omega^2}} dx dy.$$

304 This computation is done by numerical integration using the R package `landsim` (Ralph 2017).

305 Our simulation performs one step in the **life cycle** for each *year*, as follows:

- 306 1. Dispersal: Females each choose to disperse with probability 1/2; males are more vagile and disperse
307 every year. Each dispersing individual in patch i independently chooses a new location, moving to j
308 with probability M_{ij} . Since $M_{ii} > 0$, this may result in no movement.
- 309 2. Maturation: Newly born individuals are *immature*, and to mature they need to find available resources.
310 The probability per immature individual of maturing is $K/(S + K)$, where S is the local number of
311 already mature individuals, and K is the local “carrying capacity”. Each subpopulation’s carrying
312 capacity is equal to the product of ρ and the integral of habitat quality over the corresponding geographic
313 patch.
- 314 3. Birth: If there are available mates, every mature female produces offspring in a *single clutch*, mating
315 with a randomly chosen male of reproductive age (at least 15 years old) from the same population as
316 the mother, if any males are available. (If none are available, she produces no offspring.) The new
317 offspring have age 0, and the number of these produced per clutch is Poisson with a mean that depends
318 on age (see Supplement D), derived from (Reed, Fefferman, and Averill-Murray 2009).
- 319 4. Growth: increment all ages by one year.
- 320 5. Survival: kill individuals (including new ones) with probability depending on their age determined as in
321 (Reed, Fefferman, and Averill-Murray 2009) and listed in Supplement D, Table D.2.

322 As our model structure is inspired by *G. agassizii*, parameters for survival and fecundity are both drawn
323 directly from literature on tortoises. Potentially critical aspects of tortoise biology, however, are missing from
324 the model including variation in dispersal rates by age.

325 Further details and parameters of the model are given in Supplement D.

326 **Implementation with pedigree recording.** We implemented our IBM in `simuPOP` (Peng and Kimmel
327 2005), a flexible individual-based simulation library with a `python` interface. Patches were implemented
328 as subpopulations, with sex structure to allow different dispersal probabilities for males and females. Age-
329 dependence of survival and fecundity were both implemented via `python` functions passed to `simuPOP`.

330 To record the population pedigree (actually the embellished pedigree recording all relationships between all
331 genomic segments in the entire population, or ‘nedegree’) from our forward simulations into a ‘tree sequence’
332 data structure we used the efficient pedigree recording method described in Kelleher et al. (2018) and
333 implemented in `python` for `simuPOP` in `ftprime` 0.0.6rc. In this method, haploid genomes correspond to
334 nodes in the tree. Using data from every recombination, we record the tree structure on every genomic
335 interval. Periodically during a forward simulation, this data structure is simplified, and its size reduced in
336 memory, to the genomes of living individuals and their ancestors. At the end of the simulation, this tree
337 sequence can be queried to describe underlying genealogies or genome sequence, and enables extremely fast
338 computation of statistics.

339 2.3 Evaluating the method

340 We simulated single chromosomes of length 10^8 base pairs, with a recombination rate of 10^{-8} per base pair
341 per generation on the landscape of Figure 1A. We ran one simulation for 15,000 years at each of the 225
342 parameter combinations from 15 values of ρ evenly spaced in the range 0.05 to 0.2 individuals per hectare and
343 15 values of σ logarithmically spaced in the range 10 to 1000m per year. The range of density ρ corresponds
344 to 500-2000 individuals in a patch with optimal habitat (value 1.0 on the landscape of Figure 1A). Although
345 historic estimates of tortoise density are elusive (USFWS 2011), desert tortoises have been reported at
346 densities on the order of 10 per km^2 (0.1 per hectare) (Allison and McCoy 2014), and yearly home range
347 sizes span 0.01-88.6 hectares (since 1980 in CA/NV from Table 11.1 of Berish and Medica 2014). These
348 ranges are consistent with our ranges of σ and ρ . Simulations were initialized with the results of a coalescent
349 simulation in a very small population, as for instance might result from a rapid expansion. This means that
350 the populations are *not* at demographic equilibrium (at least for larger population sizes), which makes the
351 inference problem both more difficult and more realistic.

352 **Geographic regions for statistical computation.** We merged adjacent patches into seven regions based
353 on geography (Figure 1C). For each simulation, we sampled 500 individuals (see next paragraph for how these
354 individuals are chosen) and grouped them according to the seven regions, so that individuals whose location
355 fell in a region were all in the same group. We then computed all statistics between these seven groups of
356 individuals. By using only seven groups, we can compute all 406 possible statistics among the groups in a
357 reasonable amount of time.

358 **Choice of sampled individuals.** Our next step was to compare statistics computed from each simulation
359 to those obtained from data. However, these patterns of statistics can be sensitive to the geographic positions
360 of the sampled individuals, even within the geographic regions defining each group. To remove this source of
361 noise, we developed a scheme to choose, in each simulation, a set of individuals closely matching the spatial
362 positions in our dataset. To achieve this, we first defined a reference simulation (with values $\rho \approx 0.104$ and
363 $\sigma \approx 100$), and chose 500 individuals from the reference simulation. (For inference from non-simulated data
364 the reference would be the empirical samples.) Our goal was then to choose individuals in other simulations
365 that are geographically close to these. Suppose, however that we chose 50 reference individuals from a given
366 patch; we are not guaranteed in a different simulation to have 50 individuals (total) in that same patch,
367 a problem that becomes worse as the spatial discretization becomes finer. Therefore, for each reference
368 individual we assigned weights to each patch corresponding to the distribution of that individual's location
369 after 800 migration steps (by applying the migration matrix 800 times to a vector that indicates the initial
370 location). This yields 500 weightings of the patches to sample from, each corresponding to a reference sample.
371 To choose individuals from other simulations, then, we choose a patch based on the weighting and then
372 sample an individual uniformly from it. (If it is empty or all individuals have been sampled, we choose from
373 the next patch.) This process is illustrated in Figure 5.

374 **Evaluation by crossvalidation.** We evaluated the model using k -fold crossvalidation. This is a similar
375 procedure to how we chose $\hat{\omega}$ but now we use it to evaluate the model's performance in terms of relative
376 error. We randomly divided the simulated data into validation blocks. Then for every block we used the
377 simulations not in that block to infer parameters for each simulation in that block. For each validation block
378 we computed mean relative error. We did this using all possible statistics ($n = 406$): each of the six types
379 computed across all seven groups.

380 In many other supervised learning contexts, having a large number of predictors (here, the statistics) relative
381 to observations (here, the number of simulations) can cause overfitting and biased predictions (Hastie,
382 Tibshirani, and Friedman 2009). Thus, for comparison, we did the same thing with a much smaller set
383 ($n = 7$) of transformations computed on individual or groups of the statistics. These were designed based on
384 the geographic and biological setting to give us more nearly independent information. We refer to these as
385 "custom" or "biologically-motivated" statistics (see Supplement C for definitions of these).

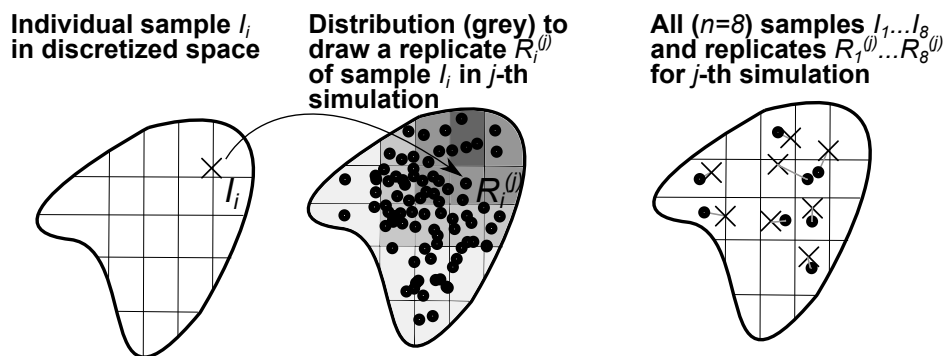


Figure 5: A method to choose simulated samples that match n reference samples: 1. identify patch locations of samples (I_i) on the discretized landscape, 2. for each sample I_i construct a distribution to draw a replicate in the j -th simulation ($R_i^{(j)}$) based on the location of I_i (solid arrow), 3. for each of j simulations draw the n replicates (connected by gray lines). Samples I_i could be empirical data or individuals in a reference simulation, in which case their location is simply their patch.

3 Results

3.1 Model behavior

The IBM simulations produced a population with age-structure and dynamics similar to that seen in the models of tortoise demography (e.g. Reed, Fefferman, and Averill-Murray 2009; Doak, Kareiva, and Klepetka 1994) from which the age-structure in the model was largely parameterized. For example, the shape of mean realized lifetime fitness (number of offspring) versus age for females (see Supplement D, Figure D.4) agrees qualitatively with that of reproductive value seen in other models.

3.2 The simulated data

Plots of all 406 statistics are difficult to interpret (see Supplement E). However, it is clear that many of the statistics carry significant signal regarding both density and dispersal. For example, separating relative divergence comparisons between regions into three categories — within regions, between neighbors, and between non-neighbors (Figure 6) — results in interpretable patterns. In self-comparisons, divergence dips well below the mean as dispersal declines, especially for less-dense populations. Meanwhile comparisons between neighbors display varying patterns that relate to geography. For example, the E-F divergence, which involves neighbors on the low-quality habitat isthmus (Figure 1A), behaved differently than the other statistics. Further intuition can be gained from examining other biologically-motivated combinations of statistics (see Appendix C, Figure C.2).

3.3 Performance

Using all possible statistics, the median relative error across all 5 cross-validation blocks was 0.104 (10.4%) with standard deviation of 0.016. In contrast, using only the few biologically motivated “custom” statistics the median relative error was 0.169 with standard deviation of 0.028. These results indicate the work needed to develop and specify statistics for a given geography was counterproductive: using only “custom” statistics had almost twice the prediction error as the method using all statistics. Figure 7 shows predicted versus actual values for the inferred parameters, density and dispersal scale.

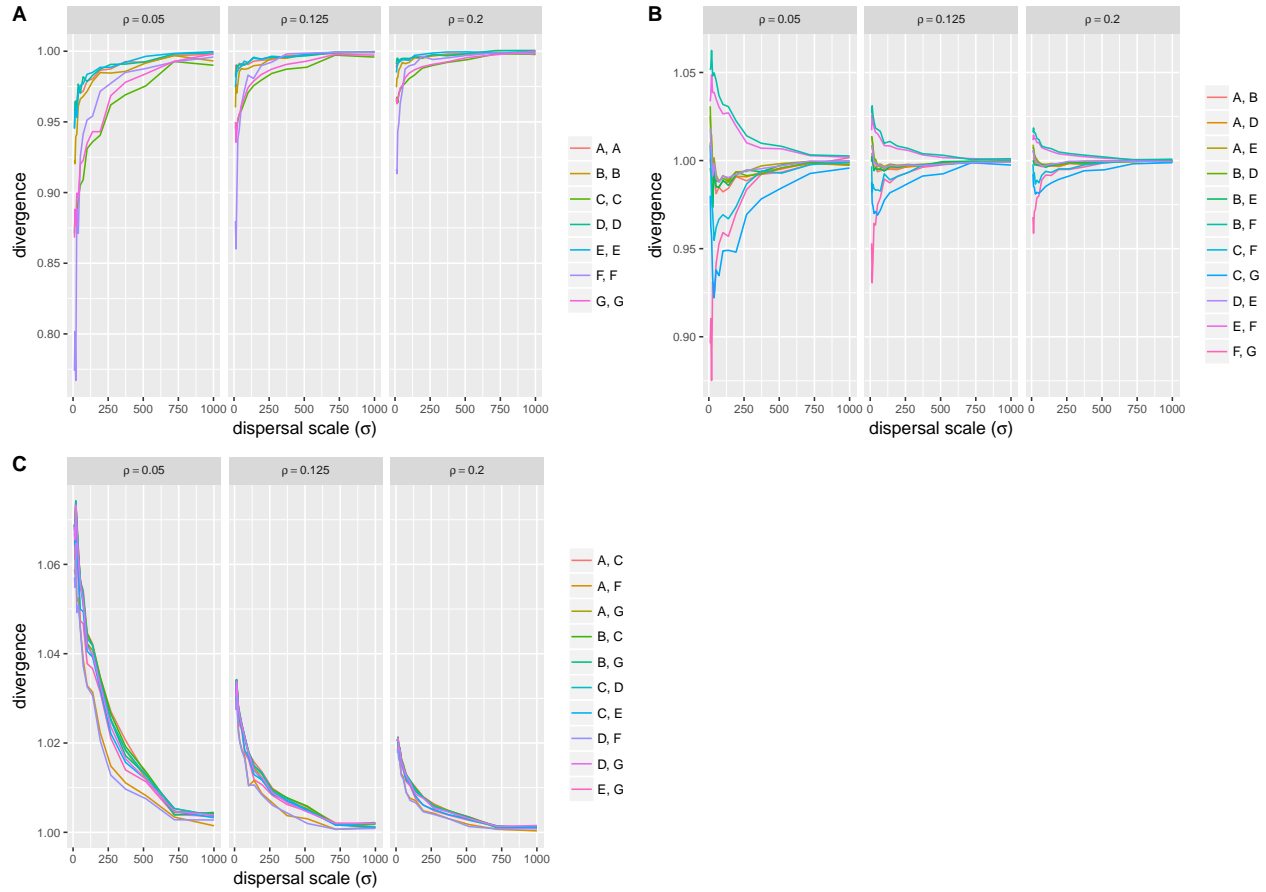


Figure 6: All pairwise divergences (scaled to mean divergence) for varying dispersal scale (σ) for comparisons (A) within regions, (B) between neighbors, and (C) between non-neighbors. Neighbors are defined based on King's neighborhood; see Supplement C. Within each panel there are three subpanels labelled with the population density ρ (in individuals per hectare).

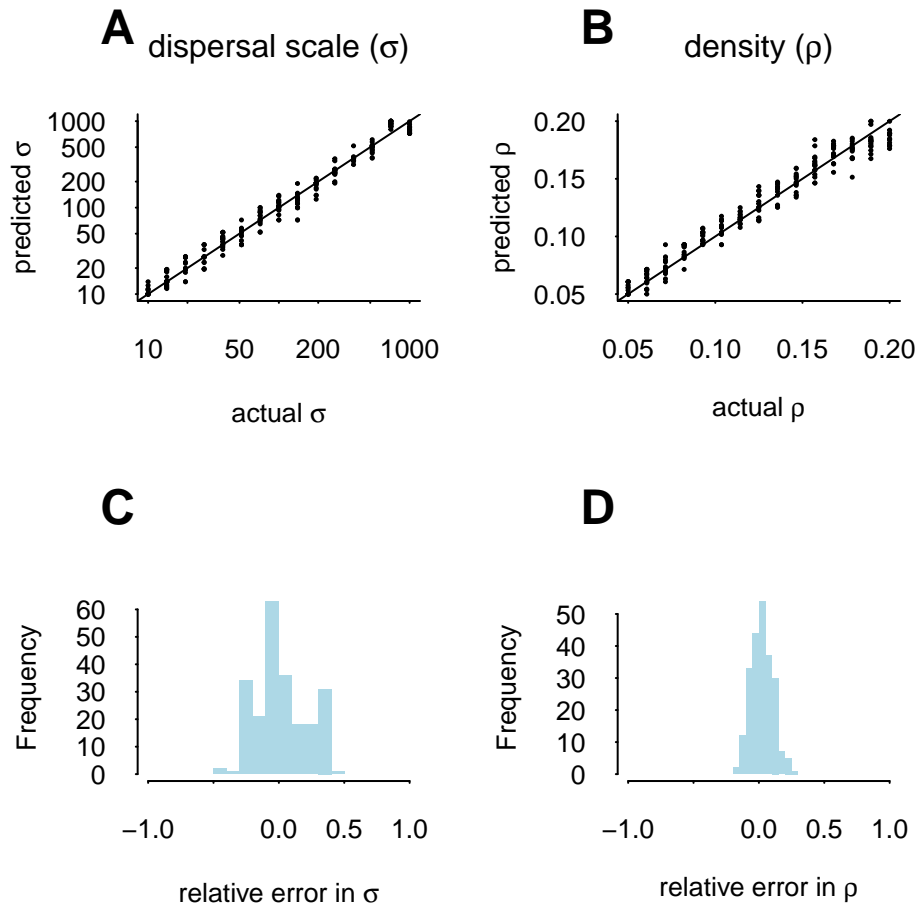


Figure 7: A, B) Inferred versus actual simulated parameters from 5-fold crossvalidation for dispersal (σ , A) and density (ρ , B). In each panel the solid line is 1:1. C, D) Relative error, i.e., the difference between the predicted value and the actual value divided by the actual value, for both dispersal (σ , C) and density (ρ , D).

410 4 Discussion

411 The methods described here enable simultaneous estimation of dispersal distance (σ) and carrying capacity
412 (ρ) in a spatially explicit model of an age-structured population. Using data simulated from a model for
413 which other parameters (e.g., age-specific survival) are fixed based on empirical values for *G. agassizii*, we
414 show that our methods can estimate dispersal and carrying capacity parameters to within 10% of their true
415 values. Thus, these methods provide a way to estimate parameters for which it is difficult to obtain data
416 within the context of relatively well-known parameters.

417 In addition, we have introduced spatial analogues of common population genetic statistics and showed how
418 and why they contain signal about geographic dynamics. Over the past decade, the four-point statistic
419 has been widely applied in population genetics (Reich et al. 2009; Peter 2016) but its utility in continuous
420 geography has not been appreciated. Here we show that the information derived from calculating all possible
421 statistics for a partitioning of continuous space provides sufficient information to recover spatial population
422 dynamics. However, as the equations in Figure 2 show, the three-point Y and four-point F are actually linear
423 combinations of the two-point statistic π . Thus, to learn about geographic dynamics it would be sufficient
424 to use only the pairwise, two-point statistic π in a statistical learning method that discovers correlations
425 between outcome variables and arbitrary linear combinations of inputs. General purpose machine learning
426 methods that are capable of such discovery include Random Forests (Breiman 2001) and techniques from
427 deep learning. However, this assumes perfect data. Because the four-point statistic F_2 is most robust to
428 sequencing error as it is not affected by singleton mutations (and the other F -statistics can be written as
429 linear combination of F_2 ; Peter 2016), in practice it may be wise to also use all pairwise F_2 s as input data for
430 inference.

431 We use inverse interpolation to estimate parameters from our simulations, and cross-validation to quantify the
432 uncertainty in parameters inferred from simulated data. For empirical data, where the parameters are truly
433 unknown, cross-validation cannot be used to quantify the parameter uncertainty. Therefore the bootstrap
434 or jackknife, which like cross-validation are based on sampling, will be necessary to estimate uncertainty in
435 estimates from empirical data.

436 These results show that integrating genomic data into structured ecological models is a feasible way to
437 estimate parameters— in our case, carrying capacity and dispersal scale— for which it is difficult to gather
438 sufficient data to estimate directly. There are additional caveats to keep in mind. First, although our method
439 incorporates prior knowledge by fixing parameters based on literature values, it is not formally Bayesian.
440 Because of this, the results here do not account for uncertainty in these estimates from the literature. Second,
441 the model we used does not account for temporal or spatial variation in survival rates or fecundity.

442 Using our method requires spatially-resolved samples of individual genomes. In practice these locations only
443 need be resolved to the scale of discrimination of the model. However, the spatial scale of discretization affects
444 inference and additional work is needed to understand these effects. Further, the distribution of these samples
445 on the landscape likely affects inference. Exactly how remains unclear, but because pairwise comparisons
446 contain information about both density and space we suggest matching the distribution of samples to the
447 population density.

448 As noted above, a lack of analytical results is a barrier to performing demographic inference with geographically
449 explicit populations. We overcame this barrier with individual-based forward simulations. Another issue,
450 however, is the possible influence of deep time: genetic variation can be influenced by both modern geography
451 and existing variation present long before the present landscape became recognizable. To address this issue
452 fully is beyond the goals of the current paper. However, as discussed in Kelleher et al. (2018) our simulation
453 framework permits combining coalescent simulations in deep time with our forwards simulations. This
454 approach (which Liu, Athanasiadis, and Weale 2008 called “sideways”) would permit analyzing multiple
455 scenarios for deep time dynamics in tandem with a set of forward simulations as we have used here.

456 **5 Acknowledgements**

457 We gratefully acknowledge helpful discussions, assistance with software, or feedback from Roy Averill-Murray,
458 Linda Allison, Jerome Kelleher, Bo Peng, Dick Tracy, and Chava Weitzman. This project was funded in part
459 by a grant from the USFWS and by the NSF (to HBS). This work benefited from access to the University of
460 Oregon high performance computer, Talapas.

Supplemental Information for

“Demographic inference in a spatially-explicit ecological model from genomic data: a proof of concept for the Mojave Desert Tortoise”

Jaime Ashander, Peter Ralph, Evan McCartney-Melstad, H. Bradley Shaffer

A Supplementary explanations of spatial statistics

We lay out the reasoning behind the statements on Figures 3 and 4 below with reference to the diagrams in Figure 2 depicting Y -statistics (with tips A, B, C) and F -statistics (with tips A, B, C, D). In each case we describe the labelling of the tips with a dictionary-like notation: {key : value}. We also denote the probability that the first coalescence involves individuals from groups A and B , for example, as \bar{AB} .

For Y statistics our heuristic is $\bar{BC} - \frac{1}{2}(\bar{AC} + \bar{AB})$. For F -statistics our heuristic is $\bar{AC} + \bar{BD} - (\bar{AD} + \bar{BC})$. We conjecture that for both of these heuristics, if the value is positive, so is the statistic. For example, if $\bar{AC} + \bar{BD} - (\bar{AD} + \bar{BC}) > 0$ then $F > 0$.

A.1 Figure 3: Spatial configuration and statistics

A.1.1 Y -statistics:

- i): Y has tips $\{A: 1, B: 2, C: 2\}$ and because $B = C$ then $\bar{BC} > \bar{AC} = \bar{AB}$ and the statistic is positive;
- ii): Y has tips $\{A: 1, B: 2, C: 3\}$ and the geometry means $\bar{BC} = \bar{AB} > \bar{AC}$ and the statistic is positive;
- iii): Y has tips $\{A: 2, B: 1, C: 3\}$ and the geometry means $\bar{BC} < \bar{AB} = \bar{AC}$ and the statistic is negative;
- iv): Y has tips $\{A: 1, B: 2, C: 3\}$ and the geometry means $\bar{BC} = \bar{AB} = \bar{AC}$ and the statistic is zero.

A.1.2 F -statistics:

- i): F has tips $\{A: 1, B: 2, C: 1, D: 2\}$ then because $A = D$ and $B = C$, $\bar{BC} = \bar{AD} > \bar{AC} = \bar{BD}$ and the statistic is positive;
- ii): F has tips $\{A: 1, B: 2, C: 1, D: 3\}$ then because $A = C$ and the geometry, $\bar{AC} > \bar{AD} = \bar{BD} > \bar{AD}$ and this statistic is positive;
- iii): If F has tips $\{A: 1, B: 2, C: 1, D: 3\}$ then because of the geometry, $\bar{AC} > \bar{AD} = \bar{BC} = \bar{BD}$ and the statistic is positive; If F has tips $\{A: 2, B: 1, C: 2, D: 3\}$ then because $A = C$ and the geometry, $\bar{AC} > \bar{AD} = \bar{BC} > \bar{BD}$ and this statistic has ambiguous sign (but if coalescences within the same group are greater than between groups, it's positive);
- iv): if F has tips $\{A: 1, B: 3, C: 2, D: 4\}$ then because of the geometry, $\bar{AC} = \bar{BD} = \bar{BC} > \bar{AD}$, the statistic is positive; if F has tips $\{A: 1, B: 2, C: 3, D: 4\}$ then because of the geometry, $\bar{BC} > \bar{AC} = \bar{BD} > \bar{AD}$ and the statistic has ambiguous sign that depends on how the chance of first coalescence involving individuals from any two groups decays with the distance between the groups;
- v): if F has tips $\{A: 1, B: 2, C: 4, D: 3\}$ then because of the geometry, $\bar{AD} = \bar{BC} > \bar{AC} = \bar{BD}$, the statistic is negative; if F has tips $\{A: 1, B: 3, C: 2, D: 4\}$ then because of the geometry, $\bar{BC} = \bar{AD} < \bar{AC} = \bar{BD}$ and this statistic is positive.

A.2 Figure 4: population size and statistics

A.2.1 Y -statistics:

- i): If Y has tips $\{A: 1, B: 2, C: 2\}$ and because $B = C$ then $\bar{BC} > \bar{AC} = \bar{AB}$ and the statistic is positive; the tips $\{A: 2, B: 1, C: 1\}$ are equal by symmetry. When 2 is replaced with $2'$, then $\bar{BC} > \bar{BC}'$

500 and so without other changes under the first labelling statistic must be *less* positive. However under
501 the second labelling, coalescence within group $2'$ has no effect and so it should be equal to the statistic
502 on equal-sized groups.

- 503 • **ii)**: For Y , coalescence within the group $1'$ has no effect and so it should be equal to the statistic on
504 equal-sized groups.
- 505 • **iii)**: For Y , coalescence within the group $1'$ has no effect and so it should be equal to the statistic on
506 equal-sized groups.

507 A.2.2 F -statistics:

- 508 • **i)**: As in Figure 2, **line 1** but when 2 is replaced with $2'$ and F has tips $\{A: 1, B: 2', C: 1, D: 2'\}$ then
509 \overline{BC} declines and the statistic is *less* positive.
- 510 • **ii)**: As in Figure 2, **line 2** but when 1 is replaced with $1'$ and F has tips $\{A: 1', B: 2, C: 1', D: 3\}$ then
511 \overline{AC} declines and the statistic is *less* positive.
- 512 • **iii)**: As in Figure 2, **line 4** but when 1 is replaced with $1'$ and F has tips $\{A: 1', B: 2, C: 1', D: 3\}$
513 then \overline{AC} declines and the statistic is *less* positive.

514 B Branch lengths or sequence details

515 Using msprime 0.5.0 we simulated populations on a 10-subpopulations stepping stone with the total
516 migration proportion out of each subpopulation set to 10^{-4} per generation. For each of 5 replicate populations
517 we performed coalescent simulations of 500 samples for a genome length of 10^6 base pairs with recombination
518 rate and mutation rate both set to 10^{-8} per base pair. Code for the simulations is shown below. For both
519 F_4 and Y_3 we chose 22 random statistics (among all possible using the 10 subpopulations as groups) and
520 computed both using site-based and branch length-based methods. The results are shown in Figure B.1.

```
import msprime
import numpy as np
import time

def run_sim(M, nsamples, **kwargs):
    assert(M.shape[0] == M.shape[1])
    n = M.shape[0]
    print("Simulating on a landscape of {} patches.".format(n))
    # sample a total of nsamples, uniformly spread
    sample_sizes = [int(np.ceil(nsamples/n)) for _ in range(n)]
    while sum(sample_sizes) < nsamples:
        sample_sizes[np.random.choice(range(len(sample_sizes)))] += 1
    population_configurations = [
        msprime.PopulationConfiguration(sample_size=k)
        for k in sample_sizes]
    # run the simulation
    begin_time = time.time()
    ts = msprime.simulate(
        population_configurations=population_configurations,
        migration_matrix=M,
        **kwargs)
    end_time = time.time()
    print("Simulation took {} seconds.".format(end_time - begin_time))
    return ts
```

```
np.random.seed(111)
n = 10
# stepping stone migration
# outmigration rate
outmig = 0.0001
M = np.zeros(shape=(n, n))
for ii in range(n):
    M[ii, ii] = 0
    if ii > 0:
        M[ii, ii - 1] = 1 * outmig / 2
    if ii < n - 1:
        M[ii, ii + 1] = 1 * outmig / 2

ts = tuple(run_sim(M,
                  nsamples=500,
                  length=1e6, Ne=1e4,
                  recombination_rate=1e-8,
                  mutation_rate=1e-8,
                  num_replicates=5))
tsl = tuple(run_sim(M,
                  nsamples=500,
                  length=1e6, Ne=1e4,
                  recombination_rate=1e-8,
                  mutation_rate=1e-10,
                  num_replicates=5))
nstats = int(n * (n-1) / 2 / 2)
```

521 C Biologically-motivated custom statistics

522 The first custom statistic aims to capture the overall timescale over which sampled alleles find a common
523 ancestor; for this we use divergence averaged across all comparisons. The next three custom statistics aim to
524 quantify the timescale over which individual alleles sampled from opposite sides of the population range find
525 a common ancestor; for this we apply the three-point statistics to groups spanning the range, with the focal
526 group in one corner and the two comparison groups on the opposite end of the landscape. For example, the Y
527 statistic with a focal group in the northwest corner and the comparison groups on the east side is $Y_3(A; C, G)$.
528 Specifically, we computed Y_3 with focal groups in the west, $Y_3^W = \frac{1}{2}(Y_3(A; C, G) + Y_3(D; C, G))$; and east
529 $Y_3^E = \frac{1}{2}(Y_3(C; A, D) + Y_3(G; A, D))$. We also computed $\bar{F}_3^{\text{corners}}$, the average of the four F_3 statistics having
530 a corner population as focal and the two most distant corner populations as the other two arguments.

531 The final three custom statistics aim to quantify differences in timescales over which alleles find common
532 ancestors depending on whether they are sampled from i) within the same region, ii) neighboring regions, or
533 iii) non-neighboring regions. To do this, we categorized every divergence statistic $\pi(A, B)$ according to these
534 three categories, i.e., whether i) $A = B$, ii) A and B were neighbors, or iii) otherwise, and averaged mean
535 divergences within each of these three categories, denoting these quantities π_w , π_n , and π_{nn} , respectively.
536 Then, we took differences of these average divergences to create three statistics: between neighbors and
537 within-region, $\pi_n - \pi_w$; between non-neighbors and within-region, $\pi_{nn} - \pi_w$; and between neighbors and
538 non-neighbors, $\pi_n - \pi_{nn}$. The neighbors and non-neighbors for each group in Figure 1C (using a King's
539 neighborhood) are listed in Table C.1.

540 We compare performance of interpolating parameter values from two different sets of statistics. First, **all** the
541 statistics meaning every one of the six types of statistics computed across all 7 groups ($n = 406$). Second,
542 **custom** statistics only, meaning the biologically motivated predictors ($n = 7$): mean divergence $\bar{\pi}$, neighbor

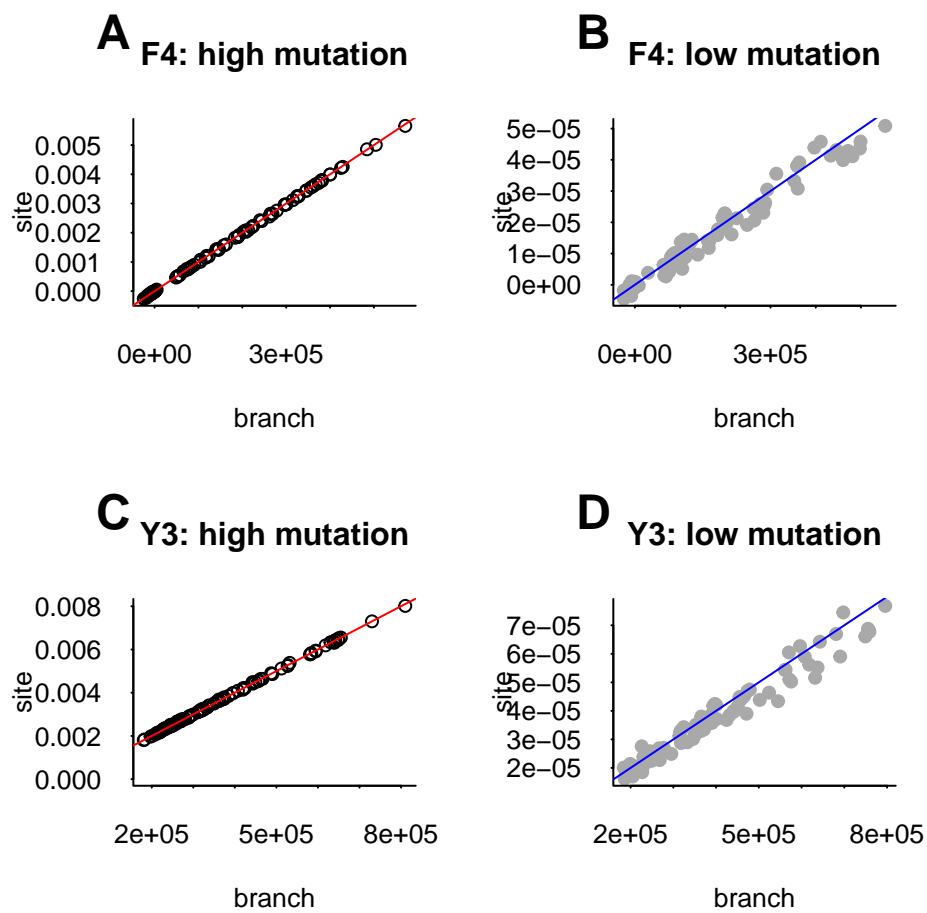


Figure B.1: Comparing F_4 and Y_3 calculated using sequence data (y-axes) and using branch lengths of marginal genealogies (x-axes) for high (circles and red lines) and low (gray dots and blue lines) mutation rates. The slopes of the lines in each plot are the mutation rates.

Table C.1: Regions, their neighbors, and their non-neighbors.

region	neighboring regions	non-neighboring regions
A	B, D, E	C, F, G
B	A, D, E, F	C, G
C	G	A, B, D, E, F
D	A, B, E	C, F, G
E	A, B, D, F	C, G
F	B, E, F, G	A, C, D
G	C, F	A, B, D, E

vs within-region divergence $\pi_n - \pi_w$, non-neighbor vs within-region divergence $\pi_{nn} - \pi_s$, non-neighbor vs neighbor divergence $\pi_{nn} - \pi_n$, $\bar{F}_3^{\text{corners}}$, Y_3^W , and Y_3^E .

Intuition from custom statistics. For example, the difference between Y_3^W and Y_3^E likely relates to the relative strength of migration and population sizes on the east and west parts of the landscape. Between-subpopulation migration is higher and population sizes are lower on the eastern part of the landscape. This means that when the focal group is in the west, the other two groups are likely to have coalesced earlier and thus Y_3 has a larger value. This is similar to Figure 4i where the three-point statistic decreases in magnitude when the non-focal population is larger: $Y(1; 2, 2) > Y(1; 2', 2')$ because group $2'$ has a larger population size than 2.

D Simulation parameters

Simulations were run for 15000 years, starting from the age distribution shown in Figure D.3, which is roughly at equilibrium. Table D.2 gives age-specific survival and fecundity. Genetics were specified as a single chromosome of length of 10^8 base pairs with 10 recombining loci and a recombination rate of 10^{-6} per generation.

Table D.2: Life table of age-based female fecundity and survival after Reed *et al.* (2009) but with fecundity (r_0) doubled (life tables model only female offspring). Males are non-reproductive below age 15 and of equal fitness above this cutoff. Immature individuals survive at 0.5 the rate of immature individuals.

min age	max age	female fecundity	survival	immature survival
0	0	0.000	1.0000	0.50000
1	1	0.000	0.7645	0.38225
2	2	0.000	0.7711	0.38555
3	3	0.000	0.7793	0.38965
4	4	0.000	0.7878	0.39390
5	5	0.000	0.7960	0.39800
6	6	0.000	0.8044	0.40220
7	7	0.000	0.8131	0.40655
8	8	0.000	0.8214	0.41070
9	9	0.000	0.8300	0.41500
10	10	0.000	0.8383	0.41915
11	11	0.000	0.8478	0.42390
12	12	1.608	0.8618	0.43090

13	13	1.864	0.8745	0.43725
14	14	3.176	0.8874	0.44370
15	15	4.642	0.9000	0.45000
16	16	5.178	0.9086	0.45430
17	17	6.152	0.9173	0.45865
18	18	6.568	0.9238	0.46190
19	19	6.750	0.9276	0.46380
20	20	7.356	0.9317	0.46585
21	21	7.478	0.9348	0.46740
22	22	7.680	0.9365	0.46825
23	23	7.754	0.9382	0.46910
24	24	7.820	0.9391	0.46955
25	25	7.886	0.9404	0.47020
26	26	7.946	0.9414	0.47070
27	27	8.006	0.9420	0.47100
28	28	8.064	0.9439	0.47195
29	29	8.128	0.9444	0.47220
30	30	8.184	0.9452	0.47260
31	31	8.238	0.9463	0.47315
32	32	8.292	0.9467	0.47335
33	33	8.346	0.9484	0.47420
34	34	8.400	0.9482	0.47410
35	35	8.450	0.9493	0.47465
36	36	8.500	0.9508	0.47540
37	37	8.550	0.9498	0.47490
38	38	8.600	0.9518	0.47590
39	39	8.648	0.9526	0.47630
40	40	8.696	0.9537	0.47685
41	41	8.742	0.9532	0.47660
42	42	8.790	0.9547	0.47735
43	43	8.836	0.9545	0.47725
44	44	8.880	0.9565	0.47825
45	45	8.924	0.9545	0.47725
46	46	8.966	0.9569	0.47845
47	47	9.010	0.9573	0.47865
48	48	9.052	0.9579	0.47895
49	49	9.092	0.9587	0.47935
50	50	9.132	0.9596	0.47980
51	51	9.172	0.9579	0.47895
52	52	9.212	0.9619	0.48095
53	53	9.248	0.9604	0.48020
54	54	9.286	0.9587	0.47935
55	55	9.322	0.9636	0.48180
56	56	9.360	0.9588	0.47940
57	57	9.398	0.9642	0.48210
58	58	9.432	0.9628	0.48140
59	59	9.468	0.9614	0.48070
60	60	9.502	0.9639	0.48195
61	61	9.536	0.9625	0.48125
62	62	9.568	0.9654	0.48270

63	63	9.598	0.9641	0.48205
64	64	9.630	0.9674	0.48370
65	65	9.660	0.9615	0.48075
66	66	9.690	0.9650	0.48250
67	67	9.718	0.9689	0.48445
68	68	9.748	0.9626	0.48130
69	69	9.780	0.9667	0.48335
70	70	9.806	0.9713	0.48565
71	71	9.834	0.9645	0.48225
72	72	9.860	0.9693	0.48465
73	73	9.888	0.9684	0.48420
74	74	9.910	0.9673	0.48365
75	Inf	9.934	0.0500	0.02500

557 The mean realized lifetime fitness (number of offspring) versus age for both males and females is shown in
558 Figure D.4. This quantity is computed from simulated data by sweeping forward in time and recording the
559 number of offspring produced after an individual reaches a given age (recorded on the x-axis); the y-axis
560 shows the mean of this quantity versus age.

561 E Extended Results

562 E.1 All statistics

563 With all possible statistics shown, the biological meaning in the patterns is difficult to discern (Figure E.5).

564 E.2 Crossvalidation results

565 Figure E.6 shows the median relative error (black dots) across all k folds of cross-validation for our method,
566 either with “custom” statistics or all possible statistics. In all cases, the method inferred parameters within
567 a few percent of their true (simulated) values. The number of folds in cross-validation did not much affect
568 performance.

569 References

- 570 Allison, Linda J, and Earl D McCoy. 2014. “Abundance of North American Tortoises: Chapter 14.” In
571 *Biology and Conservation of North American Tortoises*, edited by H.R. Mushinsky, E. D. McCoy, and D.C.
572 Rostal, 96–101. USA: Johns Hopkins University Press.
- 573 Alves, Isabel, Miguel Arenas, Mathias Currat, Anna Sramkova Hanulova, Vitor C. Sousa, Nicolas Ray, and
574 Laurent Excoffier. 2016. “Long-Distance Dispersal Shaped Patterns of Human Genetic Diversity in Eurasia.”
575 *Molecular Biology and Evolution* 33 (4): 946–58. doi:10.1093/molbev/msv332.
- 576 Barton, Nick H., Frantz Depaulis, and Alison M. Etheridge. 2002. “Neutral Evolution in Spatially Continuous
577 Populations.” *Theoretical Population Biology* 61 (1): 31–48. <http://dx.doi.org/10.1006/tpbi.2001.1557>.
- 578 Benson, John F, Peter J Mahoney, Jeff A Sikich, Laurel EK Serieys, John P Pollinger, Holly B Ernest, and
579 Seth PD Riley. 2016. “Interactions Between Demography, Genetics, and Landscape Connectivity Increase

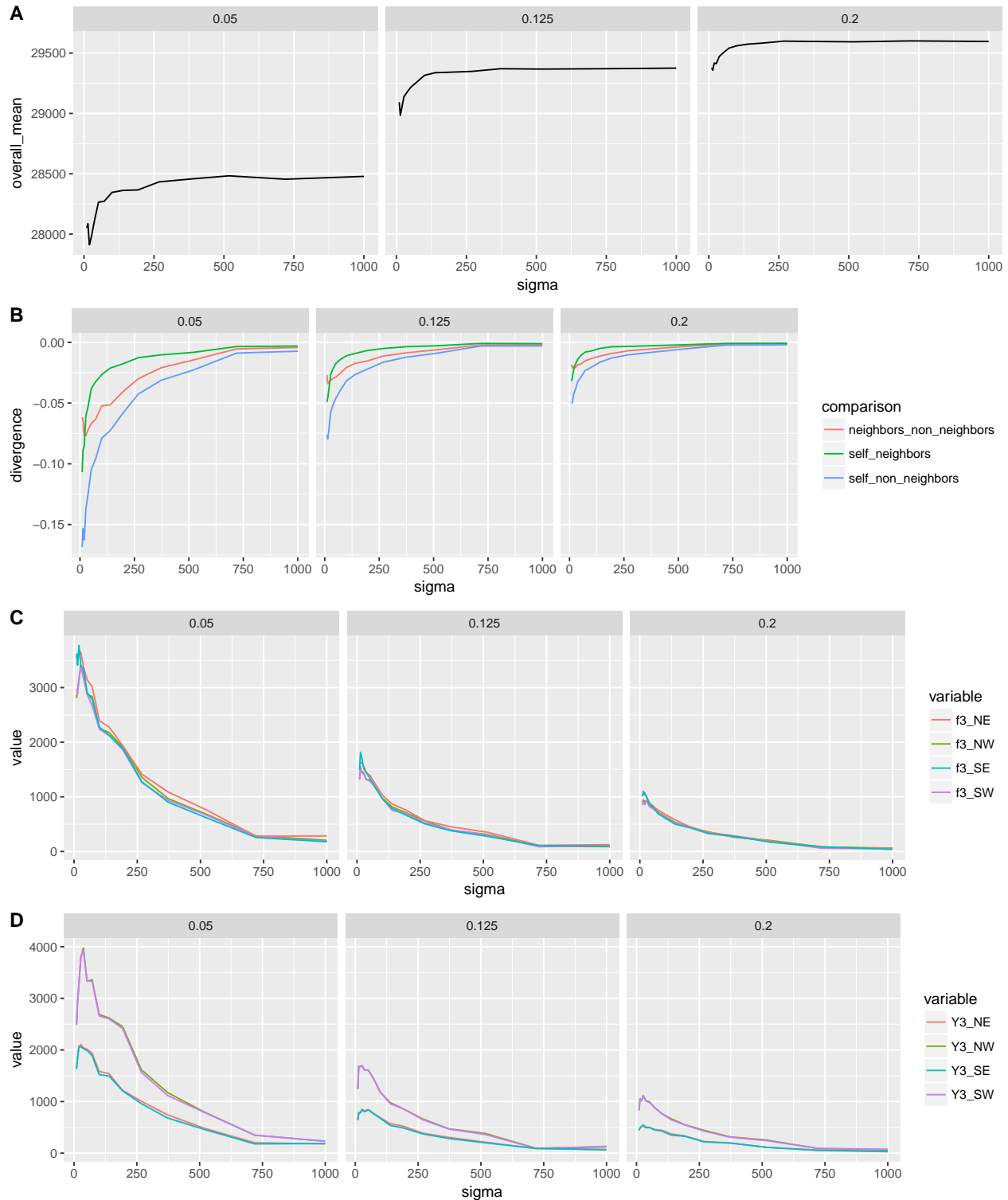


Figure C.2: Custom statistics for varying dispersal scale (σ). Within each panel there are three subpanels labelled with the population density ρ in individuals per hectare (0.05, 0.125, 0.2). A) Divergence averaged across all populations. B) divergence_{n-*nn*} scaled by mean divergence for neighbors defined in Table C.1. Three-point statistics with focal population in NW, SW, SE or NE corner (see legend): C) F_3 , in the main text F_3^{corners} is the average of the lines shown, and D) Y_3 , in the main text Y_3^W is an average of the values for W focal populations and Y_3^E is an average of the values for E focal populations.

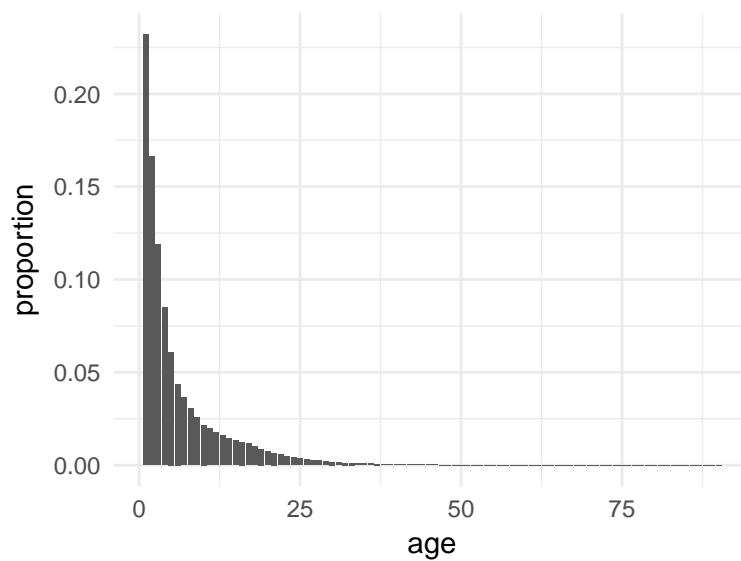


Figure D.3: Initial age distribution

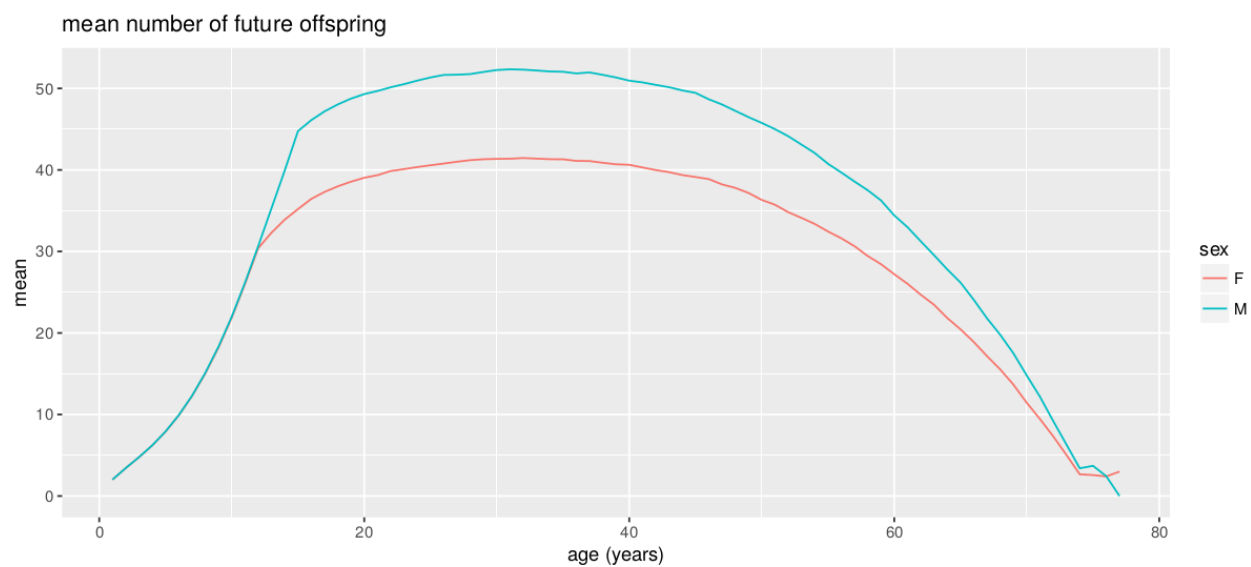


Figure D.4: Mean number of future offspring versus age (analog of reproductive value) observed in the model.

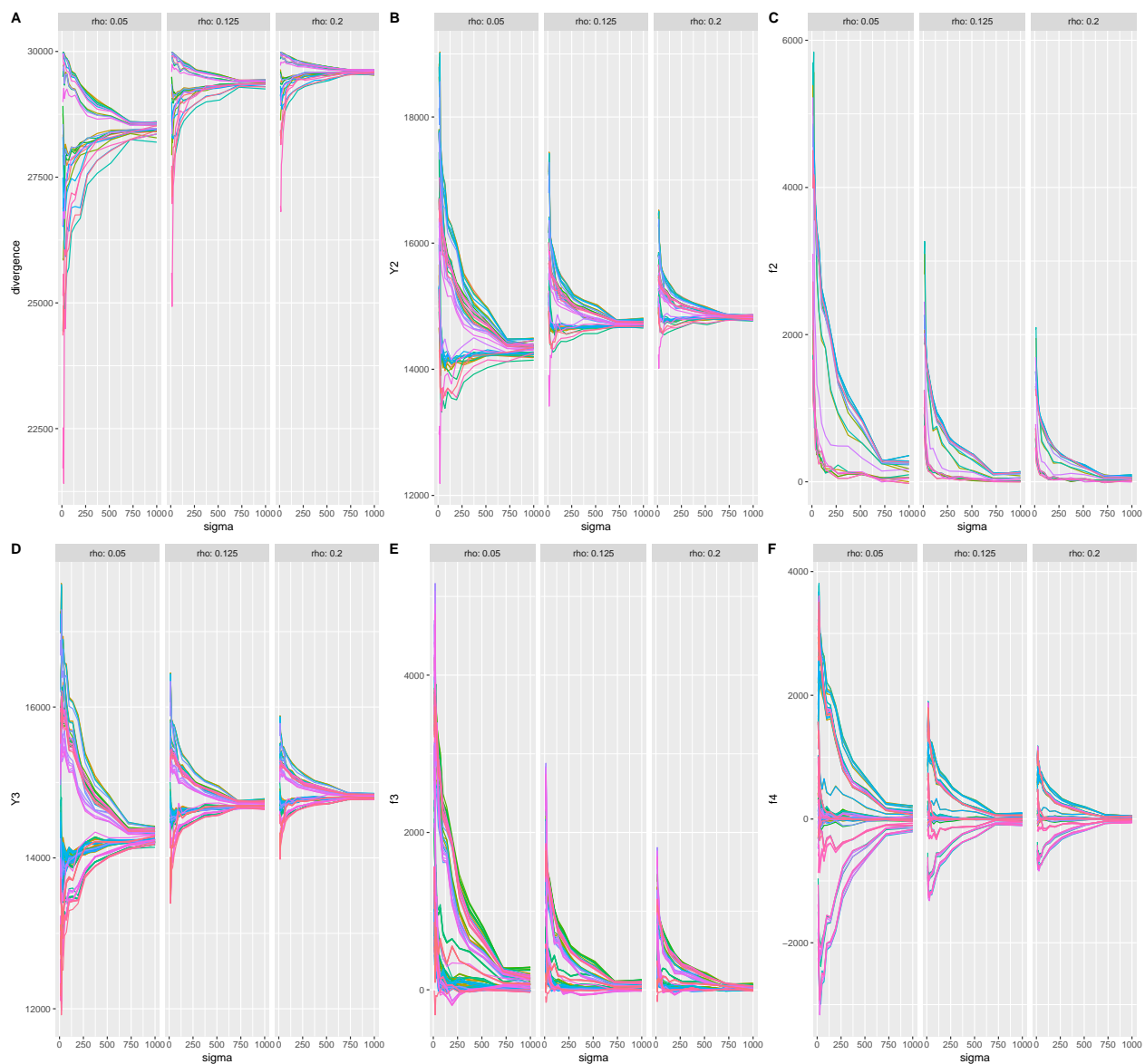


Figure E.5: All values of (A) divergence, (B) Y_2 , (C) F_2 , (D) y_3 (recall $Y_3(A; B, C) = y_3(a; b, c) - 1/2(y_3(b; a, c) + y_3(c; a, b))$), (E) F_3 , and (F) F_4 across all combinations of the 7 groups for varying dispersal scale σ . Within each panel there are three subpanels labelled with the population density ρ in individuals per hectare (0.05, 0.125, 0.2). Note differing y -axis scales.

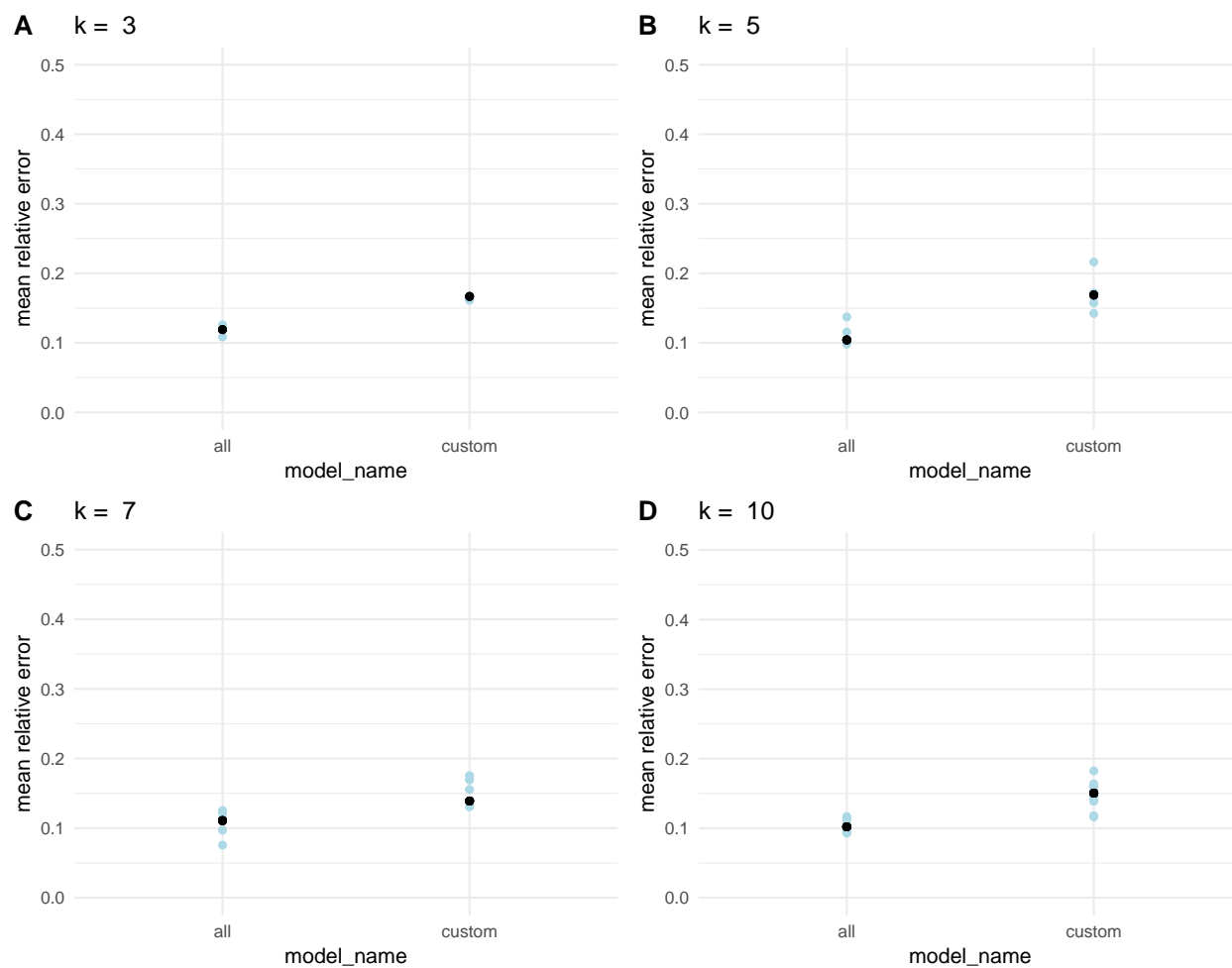


Figure E.6: Model performance under k -fold crossvalidation (for values of k noted in the titles) for both inverse interpolation with all predictors or just a few biologically motivated ones: mean relative error (blue dots) of the k crossvalidation replicates and the median across all replicates (black dot).

- 580 Extinction Probability for a Small Population of Large Carnivores in a Major Metropolitan Area.” *Proc. R.*
581 *Soc. B* 283 (1837). The Royal Society: 20160957.
- 582 Berish, Joan E, and Phil A Medica. 2014. “Home Range and Movements of North American Tortoises:
583 Chapter 11.” In *Biology and Conservation of North American Tortoises*, edited by H.R. Mushinsky, E.D.
584 McCoy, and D.C. Rostal, 96–101. USA: Johns Hopkins University Press.
- 585 Berry, Kristin H. 1986. “Incidence of Gunshot Deaths in Desert Tortoise Populations in California.” *Wildlife*
586 *Society Bulletin (1973-2006)* 14 (2). JSTOR: 127–32.
- 587 Boarman, WI, and M Sazaki. 2006. “A Highway’s Road-Effect Zone for Desert Tortoises (*Gopherus Agassizii*).”
588 *Journal of Arid Environments* 65 (1). Elsevier: 94–101.
- 589 Bradburd, Gideon, Graham Coop, and Peter Ralph. 2017. “Inferring Continuous and Discrete Population
590 Genetic Structure Across Space.” *BioRxiv*. Cold Spring Harbor Laboratory, 189688.
- 591 Breiman, Leo. 2001. “Random Forests.” *Machine Learning* 45 (1). Springer: 5–32.
- 592 Brown, Mary B, Isabella M Schumacher, Paul A Klein, Keith Harris, Terrie Correll, and Elliott R Jacobson.
593 1994. “Mycoplasma Agassizii Causes Upper Respiratory Tract Disease in the Desert Tortoise.” *Infection and*
594 *Immunity* 62 (10). Am Soc Microbiol: 4580–6.
- 595 Caswell, Hal. 2001. *Matrix Population Models: Construction Analysis and Interpretation*. Second. Sunderland,
596 MA, USA: Sinauer Associates.
- 597 Diekmann, Odo, and Johan Andre Peter Heesterbeek. 2000. *Mathematical Epidemiology of Infectious*
598 *Diseases: Model Building, Analysis and Interpretation*. Vol. 5. John Wiley & Sons.
- 599 Doak, Daniel, Peter Kareiva, and Brad Klepetka. 1994. “Modeling Population Viability for the Desert
600 Tortoise in the Western Mojave Desert.” *Ecological Applications* 4 (3). ESA: 446–60.
- 601 Esque, Todd C, Ken E Nussear, K Kristina Drake, Andrew D Walde, Kristin H Berry, Roy C Averill-Murray,
602 A Peter Woodman, et al. 2010. “Effects of Subsidized Predators, Resource Variability, and Human Population
603 Density on Desert Tortoise Populations in the Mojave Desert, Usa.” *Endangered Species Research* 12 (2):
604 167–77.
- 605 Excoffier, Laurent, and Matthieu Foll. 2011. “Fastsimcoal: A Continuous-Time Coalescent Simulator of
606 Genomic Diversity Under Arbitrarily Complex Evolutionary Scenarios.” *Bioinformatics* 27 (9): 1332–4.
607 doi:10.1093/bioinformatics/btr124.
- 608 Greenwald, KR. 2010. “Genetic Data in Population Viability Analysis: Case Studies with Ambystomatid
609 Salamanders.” *Animal Conservation* 13 (2). Wiley Online Library: 115–22.
- 610 Harris, Stephen E., Alexander T. Xue, Diego Alvarado-Serrano, Joel T. Boehm, Tyler Joseph, Michael J.
611 Hickerson, and Jason Munshi-South. 2016. “Urbanization Shapes the Demographic History of a Native
612 Rodent (the White-Footed Mouse, *Peromyscus Leucopus*) in New York City.” *Biology Letters* 12 (4): 20150983.
613 doi:10.1098/rsbl.2015.0983.
- 614 Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. *The Elements of Statistical Learning 2nd*
615 *Edition*. Springer series in statistics New York.
- 616 Hoban, Sean, Giorgio Bertorelle, and Oscar E. Gaggiotti. 2012. “Computer Simulations: Tools for Population
617 and Evolutionary Genetics.” *Nature Reviews Genetics* 13 (2): 110–22. doi:10.1038/nrg3130.
- 618 Karl, Alice E. 1998. “Reproductive Strategies, Growth Patterns, and Survivorship of a Long-Lived Herbivore
619 Inhabiting a Temporally Variable Environment.” PhD thesis, University of California—Davis.
- 620 Kelleher, Jerome, Alison M Etheridge, and Gilean McVean. 2016. “Efficient Coalescent Simulation and
621 Genealogical Analysis for Large Sample Sizes.” *PLoS Computational Biology* 12 (5). Public Library of Science:
622 e1004842.
- 623 Kelleher, Jerome, Kevin Thornton, Jaime Ashander, and Peter Ralph. 2018. “Efficient Pedigree Recording

- 624 for Fast Population Genetics Simulation.” *BioRxiv*. Cold Spring Harbor Laboratory, 248500.
- 625 Kristan, William B, and William I Boarman. 2003. “Spatial Pattern of Risk of Common Raven Predation on
626 Desert Tortoises.” *Ecology* 84 (9). Wiley Online Library: 2432–43.
- 627 Liu, Youfang, Georgios Athanasiadis, and Michael E Weale. 2008. “A Survey of Genetic Simulation Software
628 for Population and Epidemiological Studies.” *Human Genomics* 3 (1). BioMed Central: 79.
- 629 Marjoram, P. 2013. “Approximation Bayesian Computation.” *OA Genet* 1 (3): 853–53. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4297650/>.
- 631 McRae, Brad H. 2006. “Isolation by Resistance.” *Evolution* 60 (8). BioOne: 1551–61.
- 632 Moorjani, Priya, Kumarasamy Thangaraj, Nick Patterson, Mark Lipson, Po-Ru Loh, Periyasamy Govindaraj,
633 Bonnie Berger, David Reich, and Lalji Singh. 2013. “Genetic Evidence for Recent Population Mixture in
634 India.” *The American Journal of Human Genetics* 93 (3). Elsevier: 422–38.
- 635 Nafus, Melia G, Todd C Esque, Roy C Averill-Murray, Kenneth E Nussear, and Ronald R Swaisgood. 2017.
636 “Habitat Drives Dispersal and Survival of Translocated Juvenile Desert Tortoises.” *Journal of Applied Ecology*
637 54 (2). Wiley Online Library: 430–38.
- 638 Nussear, Kenneth E, Todd C Esque, Richard D Inman, Leila Gass, Kathryn A Thomas, Cynthia SA Wallace,
639 Joan B Blainey, David M Miller, and Robert H Webb. 2009. “Modeling Habitat of the Desert Tortoise
640 (*Gopherus agassizii*) in the Mojave and Parts of the Sonoran Deserts of California, Nevada, Utah, and
641 Arizona.” US Geological Survey.
- 642 Peng, Bo, and Marek Kimmel. 2005. “simuPOP: A Forward-Time Population Genetics Simulation Environ-
643 ment.” *Bioinformatics* 21 (18): 3686–7. doi:10.1093/bioinformatics/bti584.
- 644 Peter, Benjamin M. 2016. “Admixture, Population Structure, and F-Statistics.” *Genetics* 202 (4). Genetics
645 Soc America: 1485–1501.
- 646 Petkova, Desislava, John Novembre, and Matthew Stephens. 2016. “Visualizing Spatial Population Structure
647 with Estimated Effective Migration Surfaces.” *Nature Genetics* 48 (1). Nature Publishing Group: 94.
- 648 Prates, Ivan, Alexander T. Xue, Jason L. Brown, Diego F. Alvarado-Serrano, Miguel T. Rodrigues, Michael
649 J. Hickerson, and Ana C. Carnaval. 2016. “Inferring Responses to Climate Dynamics from Historical
650 Demography in Neotropical Forest Lizards.” *Proceedings of the National Academy of Sciences* 113 (29):
651 7978–85. doi:10.1073/pnas.1601063113.
- 652 Pritchard, Jonathan K, Matthew Stephens, and Peter Donnelly. 2000. “Inference of Population Structure
653 Using Multilocus Genotype Data.” *Genetics* 155 (2). Genetics Soc America: 945–59.
- 654 Quinn, Terrance J, and Richard B Deriso. 1999. *Quantitative Fish Dynamics*. Oxford University Press.
- 655 Ralph, Peter. 2017. “Efficient, Forwards-Time Simulation of Populations on Raster-Based Landscapes.”
656 <https://github.com/petrelharp/landsim>.
- 657 Ralph, Peter L. 2015. “An Empirical Approach to Demographic Inference.” *ArXiv Preprint ArXiv:1505.05816*.
- 658 Ray, Nicolas, Mathias Currat, Matthieu Foll, and Laurent Excoffier. 2010. “SPLATCHE2: A Spatially Explicit
659 Simulation Framework for Complex Demography, Genetic Admixture and Recombination.” *Bioinformatics*
660 26 (23): 2993–4. doi:10.1093/bioinformatics/btq579.
- 661 Reed, J Michael, Nina Fefferman, and Roy C Averill-Murray. 2009. “Vital Rate Sensitivity Analysis as a
662 Tool for Assessing Management Actions for the Desert Tortoise.” *Biological Conservation* 142 (11). Elsevier:
663 2710–7.
- 664 Reich, David, Kumarasamy Thangaraj, Nick Patterson, Alkes L Price, and Lalji Singh. 2009. “Reconstructing
665 Indian Population History.” *Nature* 461 (7263). Nature Publishing Group: 489.
- 666 Ringbauer, H, G Coop, and N H Barton. 2017. “Inferring Recent Demography from Isolation by Distance of

- 667 Long Shared Sequence Blocks.” *Genetics* 205 (3): 1335–51. doi:10.1534/genetics.116.196220.
- 668 Shaffer, H Bradley, Evan McCartney-Melstad, Peter L Ralph, Gideon Bradburd, Erik Lundgren, Jannet Vu,
669 Bridgette Hagerty, Fran Sandmeier, Chava Weitzman, and C Richard Tracy. 2017. “Desert Tortoises in the
670 Genomic Age: Population Genetics and the Landscape.” *BioRxiv*. Cold Spring Harbor Laboratory, 195743.
- 671 USFWS. 2011. “Revised Recovery Plan for the Mojave Population of the Desert Tortoise (*Gopherus*
672 *Agassizii*).”
- 673 Vallée, François, Aurélien Luciani, and Murray P. Cox. 2016. “Reconstructing Demography and Social
674 Behavior During the Neolithic Expansion from Genomic Diversity Across Island Southeast Asia.” *Genetics*
675 204 (4): 1495–1506. doi:10.1534/genetics.116.191379.
- 676 Wang, Ian J, Wesley K Savage, and H Bradley Shaffer. 2009. “Landscape Genetics and Least-Cost Path
677 Analysis Reveal Unexpected Dispersal Routes in the California Tiger Salamander (*Ambystoma Californiense*).”
678 *Molecular Ecology* 18 (7). Wiley Online Library: 1365–74.
- 679 Wilkins, Jon F. 2004. “A Separation-of-Timescales Approach to the Coalescent in a Continuous Population.”
680 *Genetics* 168 (4). Genetics Soc America: 2227–44.
- 681 Wright, S. 1951. “The Genetical Structure of Populations.” *Annals of Eugenics* 15 (4): 323.