

Space is the Place: Effects of Continuous Spatial Structure on Analysis of Population Genetic Data

C.J. Battey^{*,†}, Peter L. Ralph^{*,†} and Andrew D. Kern^{*,†}

*University of Oregon Dept. Biology, Institute for Ecology Evolution

ABSTRACT Real geography is continuous, but standard models in population genetics are based on discrete, well-mixed populations. As a result many methods of analyzing genetic data assume that samples are a random draw from a well-mixed population, but are applied to clustered samples from populations that are structured clinally over space. Here we use simulations of populations living in continuous geography to study the impacts of dispersal and sampling strategy on population genetic summary statistics, demographic inference, and genome-wide association studies. We find that most common summary statistics have distributions that differ substantially from that seen in well-mixed populations, especially when Wright's neighborhood size is less than 100 and sampling is spatially clustered. The combination of low dispersal and clustered sampling causes demographic inference from the site frequency spectrum to infer more turbulent demographic histories, but averaged results across multiple simulations were surprisingly robust to isolation by distance. We also show that the combination of spatially autocorrelated environments and limited dispersal causes genome-wide association studies to identify spurious signals of genetic association with purely environmentally determined phenotypes, and that this bias is only partially corrected by regressing out principal components of ancestry. Last, we discuss the relevance of our simulation results for inference from genetic variation in real organisms.

KEYWORDS Space; Population Structure; Demography; Haplotype block sharing; GWAS

25

26 **Introduction**

28 The inescapable reality that biological organisms live, move, and reproduce in continuous
29 geography is usually omitted from population genetic models. However, mates tend to live
30 near to one another and to their offspring, leading to a positive correlation between genetic

Manuscript compiled: Wednesday 11th September, 2019

¹301 Pacific Hall, University of Oregon Dept. Biology, Institute for Ecology and Evolution. cbattey2@uoregon.edu.

[†]these authors co-supervised this project

31 differentiation and geographic distance. This pattern of “isolation by distance” (?) is one of
32 the most widely replicated empirical findings in population genetics (???). Despite a long
33 history of analytical work describing the genetics of populations distributed across continuous
34 geography (e.g., ??????), much modern work still describes geographic structure as a set of
35 discrete populations connected by migration (e.g., ?????). For this reason, most population
36 genetics statistics are interpreted with reference to discrete, well-mixed populations, and most
37 empirical papers analyze variation within clusters of genetic variation inferred by programs
38 like *STRUCTURE* (?) with methods that assume these are randomly mating units.

39 The assumption that populations are “well-mixed” has important implications for down-
40 stream inference of selection and demography. Methods based on the coalescent (??) assume
41 that the sampled individuals are a random draw from a well-mixed population that is much
42 larger than the sample (?). The key assumption is that the individuals of each generation
43 are *exchangeable*, so that there is no correlation between the fate or fecundity of a parent and
44 that of their offspring (?). If dispersal or mate selection is limited by geographic proximity,
45 this assumption can be violated in many ways. For instance, if mean viability or fecundity
46 is spatially autocorrelated, then limited geographic dispersal will lead to parent–offspring
47 correlations. Furthermore, nearby individuals will be more closely related than an average
48 random pair, so drawing multiple samples from the same area of the landscape will represent
49 a biased sample of the genetic variation present in the whole population (?).

50 Two areas in which spatial structure may be particularly important are demographic infer-
51 ence and genome-wide association studies (GWAS). Previous work has found that discrete
52 population structure can create false signatures of population bottlenecks when attempting
53 to infer demographic histories from microsatellite variation (?), statistics summarizing the
54 site frequency spectrum (SFS) (???), or runs of homozygosity in a single individual (?). The
55 increasing availability of whole-genome data has led to the development of many methods
56 that attempt to infer detailed trajectories of population sizes through time based on a variety
57 of summaries of genetic data (????). Because all of these methods assume that the populations
58 being modeled are approximately randomly mating, they are likely affected by spatial biases
59 in the genealogy of sampled individuals (?), which may lead to incorrect inference of popula-
60 tion changes over time (?). However, previous investigations of these effects have focused on
61 discrete rather than continuous space models, and the level of isolation by distance at which

62 inference of population size trajectories become biased by structure is not well known. Here
63 we test how two methods suitable for use with large samples of individuals – stairwayplot (?)
64 and SMC++ (?) – perform when applied to populations evolving in continuous space with
65 varying sampling strategies and levels of dispersal.

66 Spatial structure is also a major challenge for interpreting the results of genome-wide asso-
67 ciation studies (GWAS). This is because many phenotypes of interest have strong geographic
68 differences due to the (nongenetic) influence of environmental or socioeconomic factors,
69 which can therefore show spurious correlations with spatially patterned allele frequencies (??).
70 Indeed, two recent studies found that previous evidence of polygenic selection on human
71 height in Europe was confounded by subtle population structure (??), suggesting that existing
72 methods to correct for population structure in GWAS are insufficient. However we have little
73 quantitative idea of the population and environmental parameters that can be expected to
74 lead to biases in GWAS.

75 Last, some of the most basic tools of population genetics are summary statistics like F_{IS} and
76 Tajima's D , which are often interpreted as reflecting the influence of selection or demography
77 on sampled populations (?). Statistics like Tajima's D are essentially summaries of the site
78 frequency spectrum, which itself reflects variation in branch lengths and tree structure of the
79 underlying genealogies of sampled individuals. Geographically limited mate choice distorts
80 the distribution of these genealogies (??), which can affect the value of Tajima's D (?). Similarly,
81 the distribution of tract lengths of identity by state among individuals contains information
82 about not only historical demography (??) and selection (?), but also dispersal and mate choice
83 (??). We are particularly keen to examine how such summaries will be affected by models that
84 incorporate continuous space, both to evaluate the assumptions underlying existing methods
85 and to identify where the most promising signals of geography lie.

86 To study this, we have implemented an individual-based model in continuous geography
87 that incorporates overlapping generations, local dispersal of offspring, and density-dependent
88 survival. We simulate chromosome-scale genomic data in tens of thousands of individuals
89 from parameter regimes relevant to common subjects of population genetic investigation such
90 as humans and *Drosophila*, and output the full genealogy and recombination history of all
91 final-generation individuals. We use these simulations to test how sampling strategy interacts
92 with geographic population structure to cause systematic variation in population genetic

93 summary statistics typically analyzed assuming discrete population models. We then examine
94 how the fine-scale spatial structures occurring under limited dispersal impact demographic
95 inference from the site frequency spectrum. Last, we examine the impacts of continuous
96 geography on genome-wide association studies (GWAS) and identify regions of parameter
97 space under which the results from GWAS may be misleading.

98 Materials and Methods

99 *Modeling Evolution in Continuous Space*

100 The degree to which genetic relationships are geographically correlated depends on the
101 chance that two geographically nearby individuals are close relatives – in modern terms, by
102 the tension between migration (the chance that one is descended from a distant location)
103 and coalescence (the chance that they share a parent). A key early observation by Wright
104 (?) is that this balance is often nicely summarized by the “neighborhood size”, defined
105 to be $N_W = 4\pi\rho\sigma^2$, where σ is the mean parent–offspring distance and ρ is population
106 density. This can be thought of as proportional to the average number of potential mates for
107 an individual (those within distance 2σ), or the number of potential parents of a randomly
108 chosen individual. Empirical estimates of neighborhood size vary hugely across species – even
109 in human populations, estimates range from 40 to over 5,000 depending on the population
110 and method of estimation (Table ??).

111 The first approach to modeling continuously distributed populations was to endow indi-
112 viduals in a Wright-Fisher model with locations in continuous space. However, since the
113 total size of the population is constrained, this introduces interactions between arbitrarily
114 distant individuals, which (aside from being implausible) was shown by ? to eventually
115 lead to unrealistic population clumping if the range is sufficiently large. Another method
116 for modeling spatial populations is to assume the existence of a grid of discrete randomly
117 mating populations connected by migration, thus enforcing regular population density by
118 edict. Among many other important results drawn from this class of “lattice” or “stepping
119 stone” models, ? showed that the slope of the linear regression of genetic differentiation (F_{ST})
120 against the logarithm of spatial distance is an estimate of neighborhood size. Although these
121 grid models may be good approximations of continuous geography in many situations, they
122 do not model demographic fluctuations, and limit investigation of spatial structure below

123 the level of the deme, assumptions whose impacts are unknown. An alternative method for
124 dealing with continuous geography is a new class of coalescent models, the Spatial Lambda
125 Fleming-Viot models (??).

126 To avoid questionable assumptions, we here used forward-time, individual-based simu-
127 lations. By scaling the probability of survival in each timestep to local population density
128 we shift reproductive output towards low-density regions, which prevents populations from
129 clustering. Such models have been used extensively in ecological modeling but rarely in
130 population genetics, where to our knowledge previous implementations of continuous space
131 models have focused on a small number of genetic loci, which limits the ability to investi-
132 gate the impacts of continuous space on genome-wide genetic variation as is now routinely
133 sampled from real organisms. By simulating chromosome-scale sequence alignments and
134 complete population histories we are able to treat our simulations as real populations and
135 replicate the sampling designs and analyses commonly conducted on real genomic data.

136 **A Forward-Time Model of Evolution in Continuous Space**

137 We simulated populations using the non-Wright-Fisher module in the program SLiM v3.1
138 (?). Each time step consists of three stages: reproduction, dispersal, and mortality. To reduce
139 the parameter space we use the same parameter, denoted σ , to modulate the spatial scale
140 of interactions at all three stages by adjusting the standard deviation of the corresponding
141 Gaussian functions. As in previous work (??), σ is equal to the mean parent-offspring distance.

142 At the beginning of the simulation individuals are distributed uniformly at random on
143 a continuous, square landscape. Individuals are hermaphroditic, and each time step, each
144 produces a Poisson number of offspring with mean $1/L$ where L is the expected lifespan.
145 Offspring disperse a Gaussian-distributed distance away from the parent with mean zero and
146 standard deviation σ in both the x and y coordinates. Each offspring is produced with a mate
147 selected randomly from those within distance 3σ , with probability of choosing a neighbor at
148 distance x proportional to $\exp(-x^2/2\sigma^2)$.

149 To maintain a stable population, mortality increases with local population density. To do
150 this we say that individuals at distance d have a competitive interaction with strength $g(d)$,
151 where g is the Gaussian density with mean zero and standard deviation σ . Then, the sum
152 of all competitive interactions with individual i is $n_i = \sum_j g(d_{ij})$, where d_{ij} is the distance

153 between individuals i and j and the sum is over all neighbors within distance 3σ . Since g is a
154 probability density, n_i is an estimate of the number of nearby individuals per unit area. Then,
155 given a per-unit carrying capacity K , the probability of survival until the next time step for
156 individual i is

$$p_i = \min \left(0.95, \frac{1}{1 + n_i / (K(1 + L))} \right). \quad (1)$$

157 We chose this functional form so that the equilibrium population density per unit area is close
158 to K , and the mean lifetime is around L .

159 An important step in creating any spatial model is dealing with range edges. Because local
160 population density is used to model competition, edge or corner populations can be assigned
161 artificially high fitness values because they lack neighbors within their interaction radius but
162 outside the bounds of the simulation. We approximate a decline in habitat suitability near
163 edges by decreasing the probability of survival proportional to the square root of distance to
164 edges in units of σ . The final probability of survival for individual i is then

$$s_i = p_i \min(1, \sqrt{x_i/\sigma}) \min(1, \sqrt{y_i/\sigma}) \min(1, \sqrt{(W - x_i)/\sigma}) \min(1, \sqrt{(W - y_i)/\sigma}) \quad (2)$$

165 where x_i and y_i are the spatial coordinates of individual i , and W is the width (and height) of
166 the square habitat. This buffer roughly counteracts the increase in fitness individuals close to
167 the edge would otherwise have.

168 To isolate spatial effects from other components of the model such as overlapping gener-
169 ations, increased variance in reproductive success, and density-dependent fitness, we also
170 implemented simulations identical to those above except that mates are selected uniformly
171 at random from the population, and offspring disperse to a uniform random location on the
172 landscape. We refer to this model as the “random mating” model, in contrast to the first,
173 “spatial” model.

174 We stored the full genealogy and recombination history of final-generation individuals as
175 tree sequences (?), as implemented in SLiM (?). Scripts for figures and analyses are available
176 at <https://github.com/petrelharp/spaceness>.

177 We ran 400 simulations for the spatial and random-mating models on a square landscape
178 of width $W = 50$ with per-unit carrying capacity $K = 5$ (census $N \approx 10,000$), average lifetime
179 $L = 4$, genome size = 10^8 , recombination rate = 10^{-9} , and drawing σ values from a uniform
180 distribution between 0.2 and 4. To speed up the simulations and limit memory overhead

181 we used a mutation rate of 0 in SLiM and later applied mutations to the tree sequence with
182 msprime’s `mutate` function (?). Because msprime applies mutations proportionally to elapsed
183 time, we divided the mutation rate of 10^{-8} mutations per site per generation by the average
184 generation time estimated for each value of σ (see ‘Demographic Parameters’ below) to
185 convert the rate to units of mutations per site per unit time. (We verified that this procedure
186 produced the correct number of mutations by comparing to a subset of simulations with
187 SLiM-generated mutations, which are applied only at meiosis.) Simulations were run for 1.6
188 million timesteps (approximately $30N$ generations).

189 ***Demographic Parameters***

190 Our demographic model includes parameters that control population density (K), mean life
191 span (L), and dispersal distance (σ). However, nonlinearity of local demographic stochasticity
192 causes actual realized averages of these demographic quantitites to deviate from the specified
193 values in a way that depends on the neighborhood size. Therefore, to properly compare to
194 theoretical expectations, we empirically calculated these demographic quantities in simula-
195 tions. We recorded the census population size in all simulations. To estimate generation times,
196 we stored ages of the parents of every new individual born across 200 timesteps, after a 100
197 generation burn-in, and took the mean. To estimate variance in offspring number, we tracked
198 the number of offspring for all individuals for 100 timesteps following a 100-timestep burn-in
199 period, subset the resulting table to include only the last timestep recorded for each individual,
200 and calculated the variance in number of offspring across all individuals in timesteps 50-100.
201 All calculations were performed with information recorded in the tree sequence, using pyslim
202 (<https://github.com/tskit-dev/pyslim>).

203 ***Sampling***

204 Our model records the genealogy and sequence variation of the complete population, but in
205 real data, genotypes are only observed from a relatively small number of sampled individuals.
206 We modeled three sampling strategies similar to common data collection methods in empirical
207 genetic studies (Figure ??). “Random” sampling selects individuals at random from across
208 the full landscape, “point” sampling selects individuals proportional to their distance from
209 four equally spaced points on the landscape, and “midpoint” sampling selects individuals in
210 proportion to their distance from the middle of the landscape. Downstream analyses were

211 repeated across all sampling strategies.

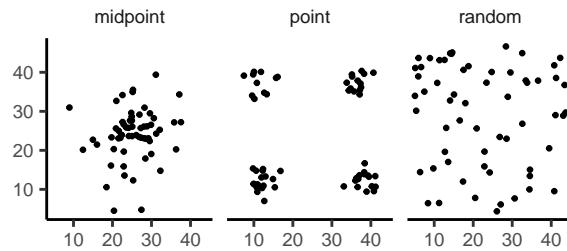


Figure 1 Example sampling maps for 60 individuals on a 50×50 landscape for midpoint, point, and random sampling strategies, respectively.

212 **Summary Statistics**

213 We calculated the site frequency spectrum and a set of 18 summary statistics (Table ??) from
214 60 diploid individuals sampled from the final generation of each simulation using the python
215 package scikit-allel (?). Statistics included common single-population summaries including
216 mean pairwise divergence (π), inbreeding coefficient (F_{IS}), and Tajima's D , as well as an
217 isolation-by-distance regression of genetic distance (D_{xy}) against the logarithm of geographic
218 distance analogous to ?'s approach, which we summarized as the correlation coefficient
219 between the logarithm of the spatial distance and the proportion of identical base pairs across
220 pairs of individuals.

221 Following recent studies that showed strong signals for dispersal and demography in the
222 distribution of shared haplotype block lengths (??), we also calculated various summaries of
223 the distribution of pairwise identical-by-state (IBS) block lengths among sampled chromo-
224 somes. The full distribution of lengths of IBS tracts for each pair of chromosomes was first
225 calculated with a custom python function. We then calculated the first three moments of this
226 distribution (mean, variance, and skew) and the number of blocks over 10^6 base pairs both for
227 each pair of individuals and for the full distribution across all pairwise comparisons.

228 We then estimated correlation coefficients between spatial distance and each moment of
229 the pairwise IBS tract distribution. Because more closely related individuals on average share
230 longer haplotype blocks we expect that spatial distance will be negatively correlated with
231 mean haplotype block length, and that this correlation will be strongest (i.e., most negative)
232 when dispersal is low. The variance, skew, and count of long haplotype block statistics are

meant to reflect the relative length of the right (upper) tail of the distribution, which represents the frequency of long haplotype blocks, and so should reflect recent demographic events (?). For a subset of simulations, we also calculated cumulative distributions for IBS tract lengths across pairs of distant (> 48 map units) and nearby (< 2 map units) individuals. Last, we examined the relationship between allele frequency and the spatial dispersion of an allele by calculating the average distance among individuals carrying each derived allele in a set of simulations representing a range of neighborhood sizes.

The effects of sampling on summary statistic estimates were summarized by testing for differences in mean (ANOVA, (?)) and variance (Levene's test, (?)) across sampling strategies for each summary statistic.

Demographic Modeling

To assess the impacts of continuous spatial structure on demographic inference we inferred population size histories for all simulations using two approaches: stairwayplot (?) and SMC++ (?). Stairwayplot fits its model to a genome-wide estimate of the SFS, while SMC++ also incorporates linkage information. For both methods we sampled 20 individuals from all spatial simulations using random, midpoint, and point sampling strategies.

As recommended by its documentation, we used stairwayplot to fit models with multiple bootstrap replicates drawn from empirical genomic data, and took the median inferred N_e per unit time as the best estimate. We calculated site frequency spectra with scikit-allel (?), generated 100 bootstrap replicates per simulation by resampling over sites, and fit models for all bootstrap samples using default settings.

For SMC++, we first output genotypes as VCF with msprime and then used SMC++'s standard pipeline for preparing input files assuming no polarization error in the SFS. We used the first individual in the VCF as the "designated individual" when fitting models, and allowed the program to estimate the recombination rate during optimization. We fit models using the 'estimate' command rather than the now recommended cross-validation approach because our simulations had only a single contig.

To evaluate the performance of these methods we binned simulations by neighborhood size, took a rolling median of inferred N_e trajectories across all model fits in a bin for each method and sampling strategy. We also examined how varying levels of isolation by distance

263 impacted the variance of N_e estimates by calculating the standard deviation of N_e from each
264 best-fit model and plotting these against neighborhood size.

265 **Association Studies**

266 To assess the degree to which spatial structure confounds GWAS we simulated four types of
267 nongenetic phenotype variation for 1000 randomly sampled individuals in each spatial SLiM
268 simulation and conducted a linear regression GWAS with principal components as covariates
269 in PLINK (?). SNPs with a minor allele frequency less than 0.5% were excluded from this
270 analysis. Phenotype values were set to vary by two standard deviations across the landscape
271 in a rough approximation of the variation seen in height across Europe (???). Conceptually
272 our approach is similar to that taken by ?, though here we model fully continuous spatial
273 variation and compare GWAS output across a range of dispersal distances.

274 In all simulations, the phenotype of each individual is determined by adding independent
275 Gaussian noise with mean zero and standard deviation 10 to a mean that may depend on
276 spatial position. We adjust the geographic pattern of mean phenotype to create spatially
277 autocorrelated environmental influences on phenotype. In the first simulation of *nonspatial*
278 environments, the mean did not change, so that all individuals' phenotypes were drawn
279 independently from a Gaussian distribution with mean 110 and standard deviation 10. Next,
280 to simulate *clinal* environmental influences on phenotype, we increased the mean phenotype
281 from 100 on the left edge of the range to 120 on the right edge (two phenotypic standard
282 deviations). Concretely, an individual at position (x, y) in a 50×50 landscape has mean
283 phenotype $100 + 2x/5$. Third, we simulated a more concentrated "*corner*" environmental
284 effect by setting the mean phenotype for individuals with both x and y coordinates below 20
285 to 120 (two standard deviations above the rest of the map). Finally, in "*patchy*" simulations we
286 selected 10 random points on the map and set the mean phenotype of all individuals within
287 three map units of each of these points to 120.

288 We performed principal components analysis (PCA) using scikit-allel (?) on the matrix of
289 derived allele counts by individual for each simulation. SNPs were first filtered to remove
290 strongly linked sites by calculating LD between all pairs of SNPs in a 200-SNP moving window
291 and dropping one of each pair of sites with an R^2 over 0.1. The LD-pruned allele count matrix
292 was then centered and all sites scaled to unit variance when conducting the PCA, following

293 recommendations in ?.

294 We ran linear-model GWAS both with and without the first 10 principal components as
295 covariates in PLINK and summarized results across simulations by counting the number of
296 SNPs with p -value below 0.05 after adjusting for an expected false positive rate of less than 5%
297 (?). We also examined p values for systematic inflation by estimating the expected values from
298 a uniform distribution (because no SNPs were used when generating phenotypes), plotting
299 observed against expected values for all simulations, and summarizing across simulations by
300 finding the mean σ value in each region of quantile-quantile space. Results from all analyses
301 were summarized and plotted with the “ggplot2” (?) and “cowplot” (?) packages in R (?).

302 **Results**

303 **Demographic Parameters**

304 Adjusting the spatial dispersal and interaction distance, σ , has a surprisingly large effect on
305 demographic quantities that are usually fixed in Wright-Fisher models – the generation time,
306 census population size, and variance in offspring number. These are shown in Figure ???. This
307 occurs because, even though the parameters K and L that control population density and
308 mean lifetime respectively were the same in all simulations, the strength of stochastic effects
309 depends strongly on σ . For instance, the population density near to individual i (denoted n_i
310 above) is computed by averaging over roughly $N_W = 4\pi K\sigma^2$ individuals, and so has standard
311 deviation proportional to $1/\sqrt{N_W}$ – it is more variable at lower densities. (Recall that N_W
312 is Wright’s neighborhood size.) Since the probability of survival is a nonlinear function of
313 n_i , actual equilibrium densities and lifetimes differ from K and L . This is the reason that we
314 included *random mating* simulations – where mate choice and offspring dispersal are both
315 nonspatial – since this should preserve the random fluctuations in local population density
316 while destroying any spatial genetic structure. We verified that random mating models
317 retained no geographic signal by showing that summary statistics did not differ significantly
318 between sampling regimes (Table ??), unlike in spatial models (discussed below).

319 There are a few additional things to note about Figure ???. First, all three quantities are
320 non-monotone with neighborhood size. Census size largely declines as neighborhood size
321 increases for both the spatial and random mating models. However, for spatial models this
322 decline only begins for neighborhood size ≥ 10 . By a neighborhood sizes larger than 100, the

323 spatial and random mating models are indistinguishable from one another, a sign that our
324 simulations are performing as expected. Census sizes range from $\approx 14,000$ at low σ in the
325 random mating model to $\approx 10,000$ for both models when neighborhood sizes approach 1,000.

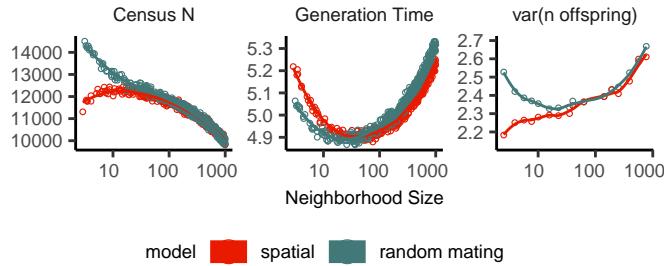


Figure 2 Genealogical parameters from spatial and random mating SLiM simulations, by neighborhood size.

326 Generation time similarly shows complex behavior with respect to neighborhood sizes,
327 and varies between 5.2 and 4.9 timesteps per generation across the parameter range explored.
328 Under both the spatial and random mating models, generation time reaches a minimum at a
329 neighborhood size of around 50. Interestingly, under the range of neighborhood sizes that we
330 examined, generation times between the random mating and spatial models are never quite
331 equivalent – presumably this would cease to be the case at neighborhood sizes higher than
332 we simulated here.

333 Last, we looked at the variance in number of offspring – a key parameter determining the
334 effective population size. Surprisingly, the spatial and random mating models behave quite
335 differently: while the variance in offspring number increases nearly monotonically under the
336 spatial model, the random mating model actually shows a decline in the variance in offspring
337 number until a neighborhood size ≈ 10 before it increases and eventually equals what we
338 observe in the spatial case.

339 **Impacts of Continuous Space on Population Genetic Summary Statistics**

340 Even though certain aspects of population demography depend on the scale of spatial inter-
341 actions, it still could be that population genetic variation is well-described by a well-mixed
342 population model. Indeed, mathematical results suggest that genetic variation in some spatial
343 models should be well-approximated by a Wright-Fisher population if neighborhood size
344 is large and all samples are geographically widely separated (??). However, the behavior of

345 most common population genetic summary statistics other than Tajima's D (?) has not yet
346 been described in realistic geographic models. Moreover, as we will show, spatial sampling
347 strategies can affect summaries of variation at least as strongly as the underlying population
348 dynamics.

349 **Site Frequency Spectra and Summaries of Diversity** Figure ?? shows the effect of varying
350 neighborhood size and sampling strategy on the site frequency spectrum (Figure ??A) and
351 several standard population genetic summary statistics (Figure ??B). Consistent with findings
352 in island and stepping stone simulations (?), the SFS shows a significant enrichment of
353 intermediate frequency variants in comparison to the nonspatial expectation. This bias is
354 most pronounced below neighborhood sizes ≤ 100 and is exacerbated by midpoint and point
355 sampling of individuals (depicted in Figure ??). Reflecting this, Tajima's D is quite positive in
356 the same situations (Figure ??B). Notably, the point at which Tajima's D approaches 0 differs
357 strongly across sampling strategies – varying from a neighborhood size of roughly 50 for
358 random sampling to at least 1000 for midpoint sampling.

359 One of the most commonly used summaries of variation is Tajima's summary of nucleotide
360 divergence, θ_π , calculated as the mean density of nucleotide differences averaged across pairs
361 of samples. As can be seen in Figure ??B, θ_π in the spatial model is inflated by up to three-fold
362 relative to the random mating model. This pattern is opposite the expectation from census
363 population size (Figure ??), because the spatial model has *lower* census size than the random
364 mating model at neighborhood sizes less than 100. Differences between these models likely
365 occur because θ_π is a measure of mean time to most recent common ancestor between two
366 samples, and at small values of σ , the time for dispersal to mix ancestry across the range
367 exceeds the mean coalescent time under random mating. (For instance, at the smallest value
368 of $\sigma = 0.2$, the range is 250 dispersal distances wide, and since the location of a diffusively
369 moving lineage after k generations has variance $k\sigma^2$, it takes around $250^2 = 62500$ generations
370 to mix across the range, which is roughly ten times larger than the random mating effective
371 population size). θ_π using each sampling strategy approaches the random mating expectation
372 at its own rate, but by a neighborhood size of around 100 all models are roughly equivalent.
373 Interestingly, the effect of sampling strategy is reversed relative to that observed in Tajima's
374 D – midpoint sampling reaches random mating expectations around neighborhood size 50,
375 while random sampling is inflated until around neighborhood size 100.

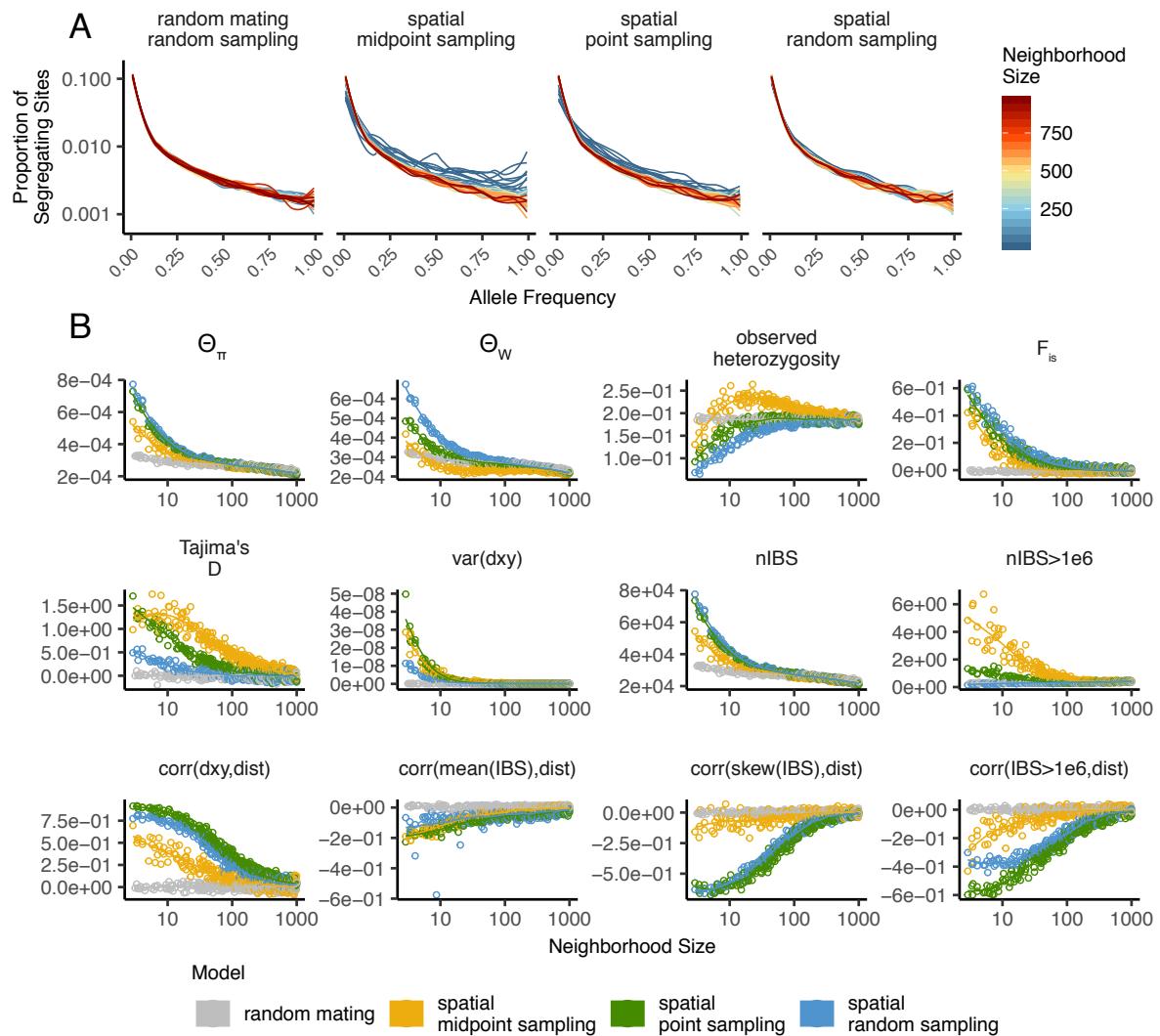


Figure 3 Site frequency spectrum (A) and summary statistic distributions (B) by sampling strategy and neighborhood size.

376 Values of observed heterozygosity and its derivative F_{IS} also depend heavily on neighbor-
377 hood size under spatial models as well as the sampling scheme. F_{IS} is inflated above the
378 expectation across most of the parameter space examined and across all sampling strategies.
379 This effect is caused by a deficit of heterozygous individuals in low-dispersal simulations –
380 a continuous-space version of the Wahlund effect (?). Indeed, for random sampling under
381 the spatial model, F_{IS} does not approach the random mating equivalent until neighborhood
382 sizes of nearly 1000. On the other hand, the dependency of raw observed heterozygosity
383 on neighborhood size is not monotone. Under midpoint sampling observed heterozygosity
384 is inflated even over the random mating expectation, as a result of the a higher proportion
385 of heterozygotes occurring in the middle of the landscape (Figure ??). This echoes a report
386 from ? who observed a similar excess of heterozygosity in the middle of the landscape when
387 simulating under a lattice model.

388 **IBS tracts and correlations with geographic distance** We next turn our attention to the effect
389 of geographic distance on haplotype block length sharing, summarized for sets of nearby and
390 distant individuals in Figure ???. There are two main patterns to note. First, nearby individuals
391 share more long IBS tracts than distant individuals (as expected because they are on average
392 more closely related). Second, the difference in the number of long IBS tracts between nearby
393 and distant individuals decreases as neighborhood size increases. This reflects the faster
394 spatial mixing of populations with higher dispersal, which breaks down the correlation
395 between the IBS tract length distribution and geographic distance. This can also be seen in the
396 bottom row of Figure ??B, where the correlation coefficients between the summaries of the
397 IBS tract length distribution (the mean, skew, and count of tracts over 10^6 bp) and geographic
398 distance approaches 0 as neighborhood size increases.

399 The patterns observed for correlations of IBS tract lengths with geographic distance are
400 similar to those observed in the more familiar regression of allele frequency measures such
401 as D_{xy} (i.e., “genetic distance”) or F_{ST} against geographic distance (?). D_{xy} is positively
402 correlated with the geographic distance between the individuals, and the strength of this
403 correlation declines as dispersal increases (Figure ??B), as expected (??). This relationship
404 is very similar across random and point sampling strategies, but is weaker for midpoint
405 sampling, perhaps due to a dearth of long-distance comparisons. In much of empirical
406 population genetics a regression of genetic differentiation against spatial distance is a de-facto

407 metric of the significance of isolation by distance. The similar behavior of moments of the
408 pairwise distribution of IBS tract lengths shows why haplotype block sharing has recently
409 emerged as a promising source of information on spatial demography through methods
410 described in ? and ?.

411 **Spatial distribution of allele copies** Mutations occur in individuals and spread geographically
412 over time. Because low frequency alleles generally represent recent mutations (??), the
413 geographic dispersion of an allele may covary along with its frequency in the population. To
414 visualize this relationship we calculated the average distance among individuals carrying a
415 focal derived allele across simulations with varying neighborhood sizes, shown in Figure ??.
416 On average we find that low frequency alleles are the most geographically restricted, and that
417 the extent to which geography and allele frequency are related depends on the amount of
418 dispersal in the population. For populations with large neighborhood sizes we found that
419 even very low frequency alleles can be found across the full landscape, whereas in populations
420 with low neighborhood sizes the relationship between distance among allele copies and their
421 frequency is quite strong. This is the basic process underlying ?'s (?) method for estimating
422 dispersal distances based on the distribution of low frequency alleles, and also generates the
423 greater degree of bias in GWAS effect sizes for low frequency alleles identified in ?.

424 **Effects of Space on Demographic Inference**

425 One of the most important uses for population genetic data is inferring demographic history
426 of populations. As demonstrated above, the site frequency spectrum and the distribution
427 of IBS tracts varies across neighborhood sizes and sampling strategies. Does this variation
428 lead to different inferences of past population sizes? To ask this we inferred population
429 size histories from samples drawn from our simulated populations with two approaches:
430 stairwayplot (?), which uses a genome-wide estimate of the SFS, and SMC++ (?), which
431 incorporates information on both the SFS and linkage disequilibrium across the genome.

432 Figure ??A shows the median inferred population size histories from each method across all
433 simulations, grouped by neighborhood size and sampling strategy. In general these methods
434 tend to slightly overestimate ancient population sizes and infer recent population declines
435 when neighborhood sizes are below 20 and sampling is spatially clustered (Figure ??A, Figure
436 ??). The overestimation of ancient population sizes however is relatively minor, averaging

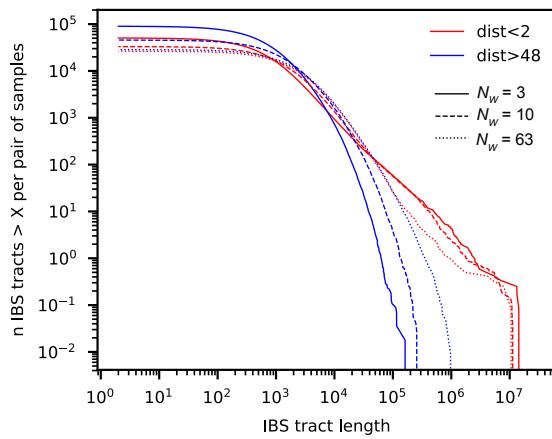


Figure 4 Cumulative distributions for IBS tract lengths per pair of individuals at different geographic distances, across three neighborhood sizes (N_W).

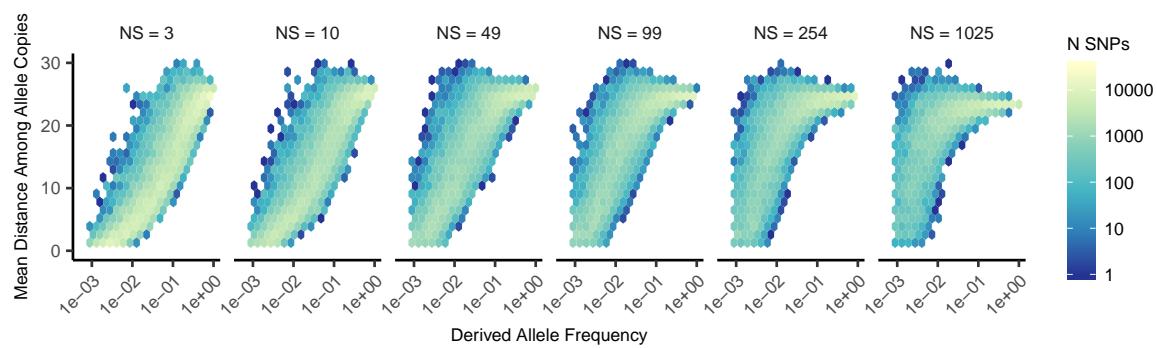


Figure 5 Trends in the distance among allele copies at varying derived allele frequencies and neighborhood sizes.

437 around a two-fold inflation at 10,000 generations before present in the worst-affected bins.
438 For stairwayplot we found that many runs infer dramatic population bottlenecks in the last
439 1,000 generations when sampling is spatially concentrated, resulting in ten-fold or greater
440 underestimates of recent population sizes. However SMC++ appeared more robust to this
441 error, with runs on point- and midpoint-sampled simulations at the lowest neighborhood
442 sizes underestimating recent population sizes by roughly half and those on randomly sampled
443 simulations showing little error. Above neighborhood sizes of around 100, both methods
444 performed relatively well when averaging across results from multiple simulations.

445 However, individual model fits from both methods frequently reflected turbulent demo-
446 graphic histories (Figure ??), with the standard deviation of inferred N_e across time points
447 often exceeding the expected N_e for both methods (Figure ??B). That is, despite the constant
448 population sizes in our simulations, both methods tended to infer large fluctuations in popu-
449 lation size over time, which could potentially result in incorrect biological interpretations. On
450 average the variance of inferred population sizes was elevated at the lowest neighborhood
451 sizes and declines as dispersal increases, with the strongest effects seen in stairwayplot model
452 fits with for clustered sampling and neighborhood sizes less than 20 (Figure ??B).

453 **GWAS**

454 To ask what confounding effects spatial genetic variation might have on genome-wide associa-
455 tion studies we performed GWAS on our simulations using phenotypes that were determined
456 solely by the environment – so, any SNP showing statistically significant correlation with
457 phenotype is a false positive. As expected, spatial autocorrelation in the environment causes
458 spurious associations across much of the genome if no correction for genetic relatedness
459 among samples is performed (Figures ?? and ??). This effect is particularly strong for clinal
460 and corner environments, for which the lowest dispersal levels cause over 60% of SNPs in
461 the sample to return significant associations. Patchy environmental distributions, which
462 are less strongly spatially correlated (Figure ??A), cause fewer false positives overall but
463 still produce spurious associations at roughly 10% of sites at the lowest neighborhood sizes.
464 Interestingly we also observed a small number of false positives in roughly 3% of analyses on
465 simulations with nonspatial environments, both with and without PC covariates included in
466 the regression.

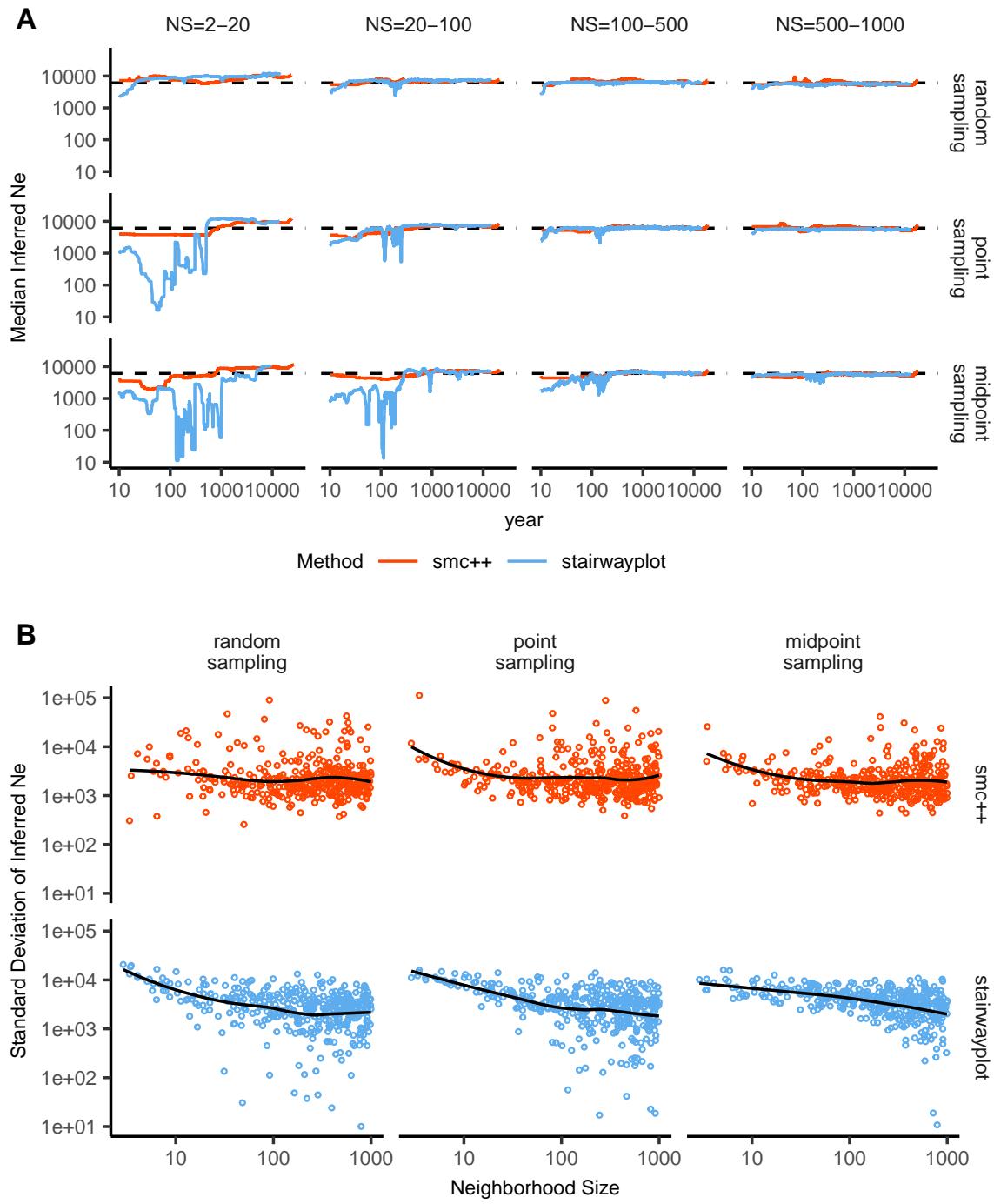


Figure 6 A: Rolling median inferred N_e trajectories for `stairwayplot` and `smc++` across sampling strategies and neighborhood size bins. The dotted line shows the mean N_e of random-mating simulations. B: Standard deviation of individual inferred N_e trajectories, by neighborhood size and sampling strategy. Black lines are loess curves. Plots including individual model fits are shown in Figure ??.

467 The confounding effects of geographic structure are well known, and it is common practice
468 to control for this by including principal components (PCs) as covariates to control for these
469 effects. This mostly works in our simulations – after incorporating the first ten PC axes as
470 covariates, the vast majority of SNPs no longer surpass a significance threshold chosen to
471 have a 5% false discovery rate (FDR). However, a substantial number of SNPs – up to 1.5% at
472 the lowest dispersal distances – still surpass this threshold (and thus would be false positives
473 in a GWAS), especially under “corner” and “patchy” environmental distributions (Figure ??C).
474 At neighborhood sizes larger than 500, up to 0.31% of SNPs were significant for corner and
475 clinal environments. Given an average of 132,000 SNPs across simulations after MAF filtering,
476 this translates to up to 382 false-positive associations; for human-sized genomes, this number
477 would be much larger. In most cases the p values for these associations were significant after
478 FDR correction but would not pass the threshold for significance under the more conservative
479 Bonferroni correction (see example Manhattan plots in figure ??).

480 Clinal environments cause an interesting pattern in false positives after PC correction:
481 at low neighborhood sizes the correction removes nearly all significant associations, but at
482 neighborhood sizes above roughly 250 the proportion of significant SNPs increases to up to
483 0.4% (Figure ??). This may be due to a loss of descriptive power of the PCs – as neighborhood
484 size increases, the total proportion of variance explained by the first 10 PC axes declines from
485 roughly 10% to 4% (Figure ??B). Essentially, PCA seems unable to effectively summarize the
486 weak population structure present in large-neighborhood simulations, but these populations
487 continue to have enough spatial structure to create significant correlations between genotypes
488 and the environment. A similar process can also be seen in the corner phenotype distribution,
489 in which the count of significant SNPs initially declines as neighborhood size increases and
490 then increases at approximately the point at which the proportion of variance explained by
491 PCA approaches its minimum.

492 Figure ??D shows quantile-quantile plots that show the degree of genome-wide inflation of
493 test statistics in PC-corrected GWAS across all simulations and environmental distributions.
494 For clinal environments, $-\log_{10}(p)$ values are most inflated when neighborhood sizes are
495 large, consistent with the pattern observed in the count of significant associations after
496 PC regression. In contrast corner and patchy environments cause the greatest inflation in
497 $-\log_{10}(p)$ at neighborhood sizes less than 100, which likely reflects the inability of PCA to

498 account for fine-scale structure caused by very limited dispersal. Finally, we observed that PC
499 regression appears to overfit to some degree for all phenotype distributions, visible in Figure
500 ??D as points falling below the 1:1 line.

501 **Discussion**

502 In this study, we have used efficient forward time population genetic simulations to describe
503 the myriad influence of continuous geography on genetic variation. In particular, we examine
504 how three main types of downstream empirical inference are affected by unmodeled spatial
505 population structure – 1) population genetic summary statistics, 2) inference of population
506 size history, and 3) genome-wide association studies (GWAS). As discussed above, space often
507 matters (and sometimes dramatically), both because of how samples are arranged in space,
508 and because of the inherent patterns of relatedness established by geography.

509 **Effects of Dispersal**

510 Limited dispersal inflates effective population size, creates correlations between genetic and
511 spatial distances, and introduces strong distortions in the site frequency spectrum that are
512 reflected in a positive Tajima's D (Figure ??). At the lowest dispersal distances, this can
513 increase genetic diversity threefold relative to random-mating expectations. These effects
514 are strongest when neighborhood sizes are below 100, but in combination with the effects
515 of nonrandom sampling they can persist up to neighborhood sizes of at least 1000 (e.g.,
516 inflation in Tajima's D and observed heterozygosity under midpoint sampling). If samples
517 are chosen uniformly from across space, the general pattern is similar to expectations of the
518 original analytic model of ?, which predicts that populations with neighborhood sizes under
519 100 will differ substantially from random mating, while those above 10,000 will be nearly
520 indistinguishable from panmixia.

521 The patterns observed in sequence data reflect the effects of space on the underlying
522 genealogy. Nearby individuals coalesce rapidly under limited dispersal and so are connected
523 by short branch lengths, while distant individuals take much longer to coalesce than they
524 would under random mating. Mutation and recombination events in our simulation both
525 occur at a constant rate along branches of the genealogy, so the genetic distance and number
526 of recombination events separating sampled individuals simply gives a noisy picture of the
527 genealogies connecting them. Tip branches (i.e., branches subtending only one individual)

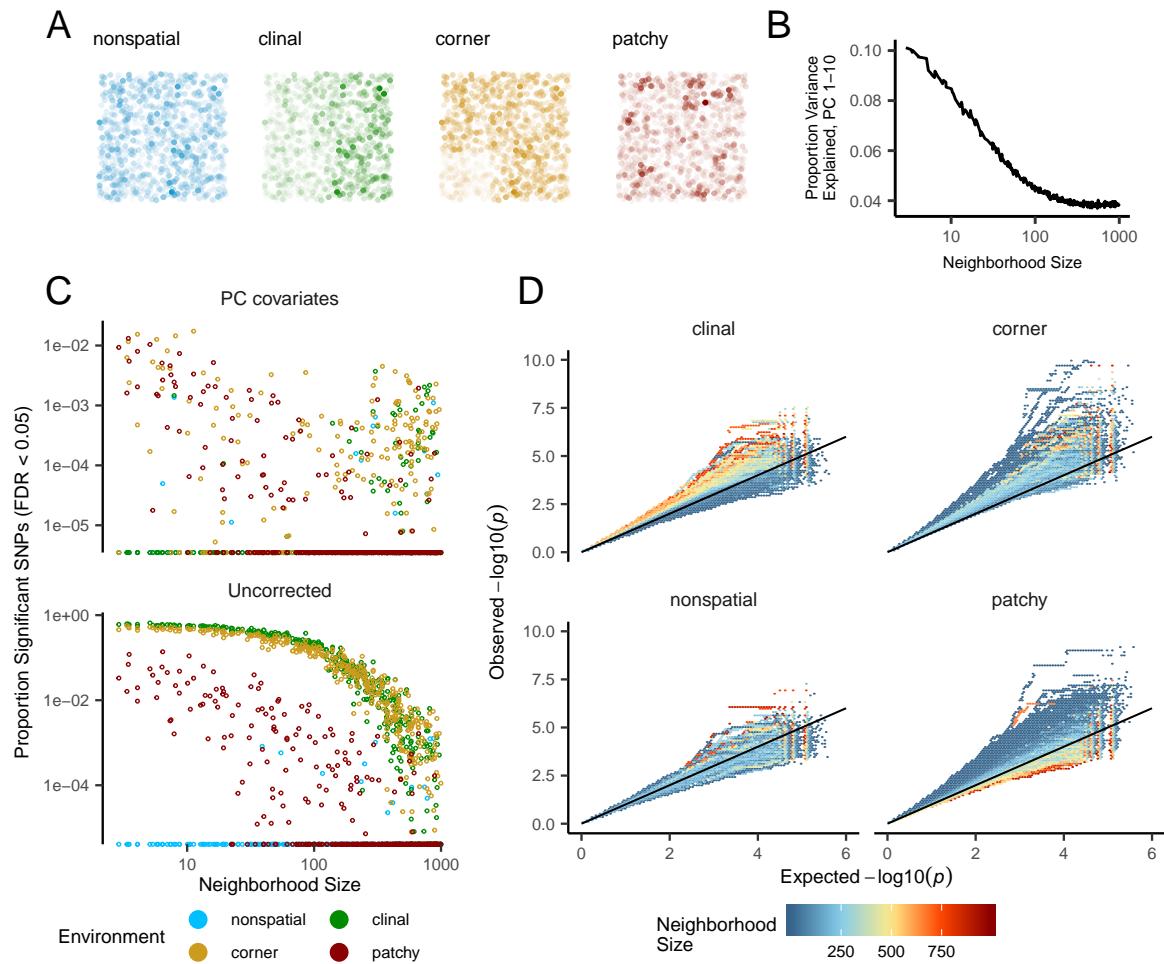


Figure 7 Impacts of spatially varying environments and isolation by distance on linear regression GWAS. Simulated quantitative phenotypes are determined only by an individual's location and the spatial distribution of environmental factors. In **A** we show the phenotypes and locations of sampled individuals under four environmental distributions, with transparency scaled to phenotype. As neighborhood size increases a PCA explains less of the total variation in the data (**B**). Spatially correlated environmental factors cause false positives at a large proportion of SNPs, which is partially but not entirely corrected by adding the first 10 PC coordinates as covariates (**C**). Quantile-quantile plots in **D** show inflation of $-\log_{10}(p)$ after PC correction across all simulations and environmental distributions, with colors scaled by the median neighborhood size in each region of q-q space.

528 are then relatively short, and branches in the middle of the genealogy connecting local groups
529 of individuals relatively long, leading to the biases in the site frequency spectrum shown in
530 Figure ??.

531 The genealogical patterns introduced by limited dispersal are particularly apparent in
532 the distribution of haplotype block lengths (Figure ??). This is because identical-by-state
533 tract lengths reflect the impacts of two processes acting along the branches of the underlying
534 genealogy – both mutation and recombination – rather than just mutation as is the case
535 when looking at the site frequency spectrum or related summaries. This means that the
536 pairwise distribution of haplotype block lengths carries with it important information about
537 genealogical variation in the population, and correlation coefficients between moments of the
538 this distribution and geographic location contain signal similar to the correlations between
539 F_{ST} or D_{xy} and geographic distance (?). Indeed this basic logic underlies two recent studies
540 explicitly estimating dispersal from the distribution of shared haplotype block lengths (??).
541 Conversely, because haplotype-based measures of demography are particularly sensitive to
542 variation in the underlying genealogy, inference approaches that assume random mating
543 when analyzing the distribution of shared haplotype block lengths are likely to be strongly
544 affected by spatial processes.

545 **Effects of Sampling**

546 One of the most important differences between random mating and spatial models is the
547 effect of sampling: in a randomly mating population the spatial distribution of sampling
548 effort has no effect on estimates of genetic variation (Table ??), but when dispersal is limited
549 sampling strategy can compound spatial patterns in the underlying genealogy and create
550 pervasive impacts on all downstream genetic analyses (see also ?). In most species, the
551 difficulty of traveling through all parts of a species range and the inefficiency of collecting
552 single individuals at each sampling site means that most studies follow something closest
553 to the “point” sampling strategy we simulated, in which multiple individuals are sampled
554 from nearby points on the landscape. For example, in ornithology a sample of 10 individuals
555 per species per locality is a common target when collecting for natural history museums. In
556 classical studies of *Drosophila* variation the situation is considerably worse, in which a single
557 orchard might be extensively sampled.

558 When sampling is clustered at points on a landscape and dispersal is limited, the sampled
559 individuals will be more closely related than a random set of individuals. Average coalescence
560 times of individuals collected at a locality will then be more recent and branch lengths shorter
561 than expected by analyses assuming random mating. This leads to fewer mutations and
562 recombination events occurring since their last common ancestor, causing a random set of
563 individuals to share longer average IBS tracts and have fewer nucleotide differences. For some
564 data summaries, such as Tajima’s D , Watterson’s Θ , or the correlation coefficient between
565 spatial distance and the count of long haplotype blocks, this can result in large differences in
566 estimates between random and point sampling (Figure ??). Inferring underlying demographic
567 parameters from these summary statistics – unless the nature of the sampling is somehow
568 taken into account – will be subject to bias if sampling is not random across the landscape.

569 However, we observed the largest sampling effects using “midpoint” sampling. This
570 model is meant to reflect a bias in sampling effort towards the middle of a species’ range.
571 In empirical studies this sampling strategy could arise if, for example, researchers choose to
572 sample the center of the range and avoid range edges to maximize probability of locating
573 individuals during a short field season. Because midpoint sampling provides limited spatial
574 resolution it dramatically reduces the magnitude of observed correlations between spatial
575 and genetic distances. More surprisingly, midpoint sampling also leads to strongly positive
576 Tajima’s D and an inflation in the proportion of heterozygous individuals in the sample –
577 similar to the effect of sampling a single deme in an island model as reported in (?). This
578 increase in observed heterozygosity appears to reflect the effects of range edges, which are
579 a fundamental facet of spatial genetic variation. If individuals move randomly in a finite
580 two-dimensional landscape then regions in the middle of the landscape receive migrants from
581 all directions while those on the edge receive no migrants from at least one direction. The
582 average number of new mutations moving into the middle of the landscape is then higher
583 than the number moving into regions near the range edge, leading to higher heterozygosity
584 and lower inbreeding coefficients (F_{IS}) away from range edges. Though here we used only a
585 single parameterization of fitness decline at range edges we believe this is a general property
586 of non-infinite landscapes as it has also been observed in previous studies simulating under
587 lattice models (??).

588 In summary, we recommend that empirical researchers collect individuals from across as

589 much of the species' range as practical, choosing samples separated by a range of spatial scales.
590 Many summary statistics are designed for well-mixed populations, and so provide different
591 insights into genetic variation when applied to different subsets of the population. Applied
592 to a cluster of samples, summary statistics based on segregating sites (e.g., Watterson's Θ
593 and Tajima's D), heterozygosity, or the distribution of long haplotype blocks, can be expected
594 to depart significantly from what would be obtained from a wider distribution of samples.
595 Comparing the results of analyses conducted on all individuals versus those limited to single
596 individuals per locality can provide an informative contrast. Finally we wish to point out
597 that the bias towards intermediate allele frequencies that we observe may mean that the
598 importance of linked selection, at least as is gleaned from the site frequency spectrum, may be
599 systematically underestimated currently.

600 **Demography**

601 Previous studies have found that population structure and nonrandom sampling can create
602 spurious signals of population bottlenecks when attempting to infer demographic history
603 with microsatellite variation, summary statistics, or runs of homozygosity (????). Here we
604 found that methods that infer detailed population trajectories through time based on the SFS
605 and patterns of LD across the genome are also subject to this bias, with some combinations
606 of dispersal and sampling strategy systematically inferring deep recent population bottle-
607 necks and overestimating ancient N_e by a around a factor of 2. We were surprised to see
608 that both stairwayplot and SMC++ can tolerate relatively strong isolation by distance – i.e.,
609 neighborhood sizes of 20 – and still perform well when averaging results across multiple
610 simulations. Inference in populations with neighborhood sizes over 20 was relatively un-
611 biased unless samples were concentrated in the middle of the range (Figure ??). Although
612 median demography estimates across many independent simulations were fairly accurate,
613 empirical work has only a single estimate to work with, and individual model fits (Figure
614 ??) suggest that spuriously inferred population size changes and bottlenecks are common,
615 especially at small neighborhood sizes. As we will discuss below, most empirical estimates
616 of neighborhood size, including all estimates for human populations, are large enough that
617 population size trajectories inferred by these approaches should not be strongly affected by
618 spatial biases created by dispersal in continuous landscapes. In contrast, ? found that varying

619 migration rates through time could create strong biases in inferred population trajectories
620 from an n -island model with parameters relevant for human history, suggesting that changes
621 in migration rates through time are more likely to drive variation in inferred N_e than isolation
622 by distance.

623 We found that SMC++ was more robust to the effects of space than stairwayplot, under-
624 estimating recent populations by roughly half in the worst time periods rather than nearly
625 10-fold as with stairwayplot. Though this degree of variation in population size is certainly
626 meaningful in an ecological context, it is relatively minor in population genetic terms. A
627 more worrying pattern was the high level of variance in inferred N_e trajectories for individual
628 model fits using these methods, which was highest in simulations with the smallest neighbor-
629 hood size (Figure ??, Figure ??). This suggests that, at a minimum, researchers working with
630 empirical data should replicate analyses multiple times and take a rolling average if model
631 fits are inconsistent across runs. Splitting samples and running replicates on separate subsets –
632 the closest an empirical study can come to our design of averaging the results from multiple
633 simulations – may also alleviate this issue.

634 Our analysis suggests that many empirical analyses of population size history using meth-
635 ods like SMC++ are robust to error caused by spatial structure within continuous landscapes.
636 Inferences drawn from static SFS-based methods like stairwayplot should be treated with
637 caution when there are signs of isolation by distance in the underlying data (for example, if
638 a regression of F_{ST} against the logarithm of geographic distance has a significantly positive
639 slope), and in particular an inference of population bottlenecks in the last 1000 years should
640 be discounted if sampling is clustered, but estimates of deeper time patterns are likely to be
641 fairly accurate. The biases in the SFS and haplotype structure identified above (see also ???)
642 are apparently small enough that they fall within the range of variability regularly inferred by
643 these approaches, at least on datasets of the size we simulated.

644 **GWAS**

645 Spatial structure is particularly challenging for genome-wide association studies, because the
646 effects of dispersal on genetic variation are compounded by spatial variation in the environ-
647 ment (?). Spatially restricted mate choice and dispersal causes variation in allele frequencies
648 across the range of a species. If environmental factors affecting the phenotype of interest also

649 vary over space, then groups of individuals in different regions will allele frequencies and
650 environmental exposures will covary over space. In this scenario an uncorrected GWAS will
651 infer genetic associations with a purely environmental phenotype at any site in the genome
652 that is differentiated over space, and the relative degree of bias will be a function of the degree
653 of covariation in allele frequencies and the environment (i.e., Figure ??C, bottom panel). This
654 pattern has been demonstrated in a variety of simulation and empirical contexts (??????????).

655 Incorporating PC positions as covariates in a linear-regression GWAS (?) is designed to
656 address this challenge by regressing out a baseline level of “average” differentiation. In
657 essence, a PC-corrected GWAS asks “what regions of the genome are more associated with
658 this phenotype than the average genome-wide association observed across populations?” In
659 our simulations, we observed that this procedure can fail under a variety of circumstances. If
660 dispersal is limited and environmental variation is clustered in space (i.e., corner or patchy
661 distributions in our simulations), PCA positions fail to capture the fine-scale spatial structure
662 required to remove all signals of association. Conversely, as dispersal increases, PCA loses
663 power to describe population structure before spatial mixing breaks down the relationship
664 between genotype and the environment. These effects were observed with all spatially
665 correlated environmental patterns, but were particularly pronounced if environmental effects
666 are concentrated in one region, as was also found by ?. Though increasing the number of PC
667 axes used in the analysis may reduce the false-positive rate, this may also decrease the power
668 of the test to detect truly causal alleles (?).

669 In this work we simulated a single chromosome with size roughly comparable to one
670 human chromosome. If we scale the number of false-positive associations identified in our
671 analyses to a GWAS conducted on whole-genome data from humans, we would expect to see
672 several thousand weak false-positive associations after PC corrections in a population with
673 neighborhood sizes up to at least 1000 (which should include values appropriate for many
674 human populations). Notably, very few of the spurious associations we identified would be
675 significant at a conservative Bonferroni-adjusted p -value cutoff (see Figure ??). This suggests
676 that GWAS focused on finding strongly associated alleles for traits controlled by a limited
677 number of variants in the genome are likely robust to the impacts of continuous spatial struc-
678 ture. However, methods that analyze the combined effects of thousands or millions of weakly
679 associated variants such as polygenic risk scores (?) are likely to be affected by subtle popula-

tion structure. Indeed as recently identified in studies of genotype associations for human height in Europe (??), PC regression GWAS in modern human populations do include residual signal of population structure in large-scale analyses of polygenic traits. When attempting to make predictions across populations with different environmental exposures, polygenic risk scores affected by population structure can be expected to offer low predictive power, as was shown in a recent study finding lower performance outside European populations (?).

In summary, spatial covariation in population structure and the environment confounds the interpretation of GWAS *p*-values, and correction using principal components is insufficient to fully separate these signals for polygenic traits under a variety of environmental and population parameter regimes. Other GWAS methods may be less sensitive to this confounding, but there is no obvious reason that this should be so. One approach to estimating the degree of bias in GWAS caused by population structure is LD score regression (?). Though this approach appears to work well in practice, its interpretation is not always straightforward and it is likely biased by the presence of linked selection (?). In addition, we observed that in many cases the false-positive SNPs we identified appeared to be concentrated in LD peaks similar to those expected from truly causal sites (Figure ??), which may confound LD score regression.

We suggest a straightforward alternative for species in which the primary axes of population differentiation is space (note this is likely not the case for some modern human populations): run a GWAS with spatial coordinates as phenotypes and check for *p*-value inflation or significant associations. If significant associations with sample locality are observed after correcting for population structure, the method is sensitive to false positives induced by spatial structure. This is essentially the approach taken in our “clinal” model (though we add normally distributed noise to our phenotypes). Of course, it is possible that genotypes indirectly affect individual locations by adjusting organismal fitness and thus habitat selection across spatially varying environments, but we believe that this hypothesis should be tested against a null of stratification bias inflation rather than accepted as true based on GWAS results.

707 **Where are natural populations on this spectrum?**

708 For how much of the tree of life do spatial patterns circumscribe genomic variation? In
709 Table ?? we gathered estimates of neighborhood size from a range of organisms to get an

Table 1 Neighborhood size estimates from empirical studies.

Species	Description	Neighborhood Size	Method	Citation
<i>Ipomopsis aggregata</i>	flowering plant	12.60 - 37.80	Genetic	(?)
<i>Borreria frutescens</i>	salt marsh plant	20 - 30	Genetic+Survey	(?)
<i>Oreamnos americanus</i>	mountain goat	36 - 100	Genetic	(?)
<i>Homo sapiens</i>	Gainj- and Kalam-speaking people, Papua New Guinea	40 - 213	Genetic	(?)
<i>Formica sp.</i>	colonial ants	50 - 100	Genetic	(?)
<i>Astrocaryum mexicanum</i>	palm tree	102 - 895	Genetic+survey	(?)
<i>Spermophilus mollis</i>	ground squirrel	204 - 480	Genetic+Survey	(?)
<i>Sceloporus olivaceus</i>	lizard	225 - 270	Survey	(?)
<i>Dieffenbachia longispatha</i>	beetle-pollinated colonial herb	227 - 611	Survey	(?)
<i>Aedes aegypti</i>	Yellow-fever mosquito	268	Genetic	(?)
<i>Homo sapiens</i>	Gainj- and Kalam-speaking people, Papua New Guinea	410	Survey	(?)
<i>Quercus laevis</i>	Oak tree	> 440	Genetic	(?)
<i>Drosophila pseudoobscura</i>	fruit fly	500 - 1,000	Survey+Crosses	(?)
<i>Homo sapiens</i>	POPRES data NE Europe	1,342 - 5,425	Genetic	(?)
<i>Bebicium vittatum</i>	intertidal snail	240,000	Survey	(?)
<i>Bebicium vittatum</i>	intertidal snail	360,000	Genetic	(?)

idea of how likely dispersal is to play an important role in patterns of variation. Though this sample is almost certainly biased towards small-neighborhood species (because few studies have quantified neighborhood size in species with very high dispersal or population density), we find that neighborhood sizes in the range we simulated are fairly common across a range of taxa. At the extreme low end of empirical neighborhood size estimates we see some flowering plants, large mammals, and colonial insects like ants. Species such as this have neighborhood size estimates small enough that spatial processes are likely to strongly influence inference. These include some human populations such as the Gainj- and Kalam-speaking people of Papua New Guinea, in which the estimated neighborhood sizes in (?) range from 40 to 410 depending on the method of estimation. Many more species occur in a middle range of neighborhood sizes between 100 and 1000 – a range in which spatial processes play a minor role in our analyses under random spatial sampling but are important when sampling of individuals in space is clustered. Surprisingly, even some flying insects with huge census population sizes fall in this group, including fruit flies (*D. melanogaster*) and mosquitoes (*A. aegypti*). Last, many species likely have neighborhood sizes much larger than we simulated, including modern humans in northeastern Europe (?). For these species demographic inference and summary statistics are likely to reflect minimal bias from spatial effects as long as dispersal is truly continuous across the landscape. While that is so we caution that association studies in which the effects of population structure are confounded with spatial variation in the environment are still sensitive to dispersal even at these large neighborhood sizes.

731 ***Future Directions and Limitations***

As we have shown, a large number of population genetic summary statistics contain information about spatial population processes. We imagine that combinations of such summaries might be sufficient for the construction of supervised machine learning regressors (e.g., ?) for the accurate estimation of dispersal from genetic data. Indeed, ? found that inverse interpolation on a vector of summary statistics provided a powerful method of estimating dispersal distances. Expanding this approach to include the haplotype-based summary statistics studied here and applying machine learning regressors built for general inference of nonlinear relationships from high-dimensional data may allow precise estimation of spatial parameters

740 under a range of complex models.

741 One facet of spatial variation that we did not address in this study is the confounding of
742 dispersal and population density implicit in the definition of Wright's neighborhood size. Our
743 simulations were run under constant densities, but ?'s approach to demographic inference
744 in space suggests that density and dispersal can in some cases be estimated separately from
745 genetic data. Much additional work remains to be done to better understand how these
746 parameters interact to shape genetic variation in continuous space, which we leave to future
747 studies.

748 Though our simulation allows incorporation of realistic demographic and spatial processes,
749 it is inevitably limited by the computational burden of tracking tens or hundreds of thousands
750 of individuals in every generation. In particular, computations required for mate selection
751 and spatial competition scale approximately with the product of the total census size and
752 the neighborhood size and so increase rapidly for large populations and dispersal distances.
753 The reverse-time model of continuous space evolution described by ? and implemented by ?
754 allows exploration of parameter regimes with population and landscape sizes more directly
755 comparable to empirical cases like humans. Alternatively, implementation of parallelized
756 calculations may allow progress with forward-time simulations.

757 Finally, we believe that the difficulties in correcting for population structure in continuous
758 populations using principal components analysis or similar decompositions is a difficult
759 issue, well worth considering on its own. How can we best avoid spurious correlations while
760 correlating genetic and phenotypic variation without underpowering the methods? Perhaps
761 optimistically, we posit that process-driven descriptions of ancestry and/or more generalized
762 unsupervised methods may be able to better account for carry out this task.

763 **Data Availability**

764 Scripts used for all analyses and figures are available at <https://github.com/petrelharp/spaceness>.

765 **Acknowledgements**

766 We thank Brandon Cooper, Matt Hahn, Doc Edge, and others for reading and thinking about
767 this manuscript. CJB and ADK were supported by NIH award R01GM117241.

768 **Appendix 1**

769 ***Comparisons with Alternate Spatial Models***

770 Stepping-stone models have been extremely useful in population genetics and, when sim-
771 ulated in reverse time, remain much faster than our implementation of continuous-space
772 evolution in forward time. Similarly, the Wright-Malécot models have lead to analytic re-
773 sults that have proved reasonable approximations for estimating neighborhood sizes from
774 empirical data (?). When are these models "good enough"?

775 First we address a two-dimensional Wright-Malécot model. SLiM includes an implemen-
776 tation of a similar model in recipe 14.1. This model simulates a Wright-Fisher population
777 with discrete generations and fixed total population size. Individuals are initially distributed
778 randomly across the landscape, and each generation

779 **Supplementary Figures and Tables**

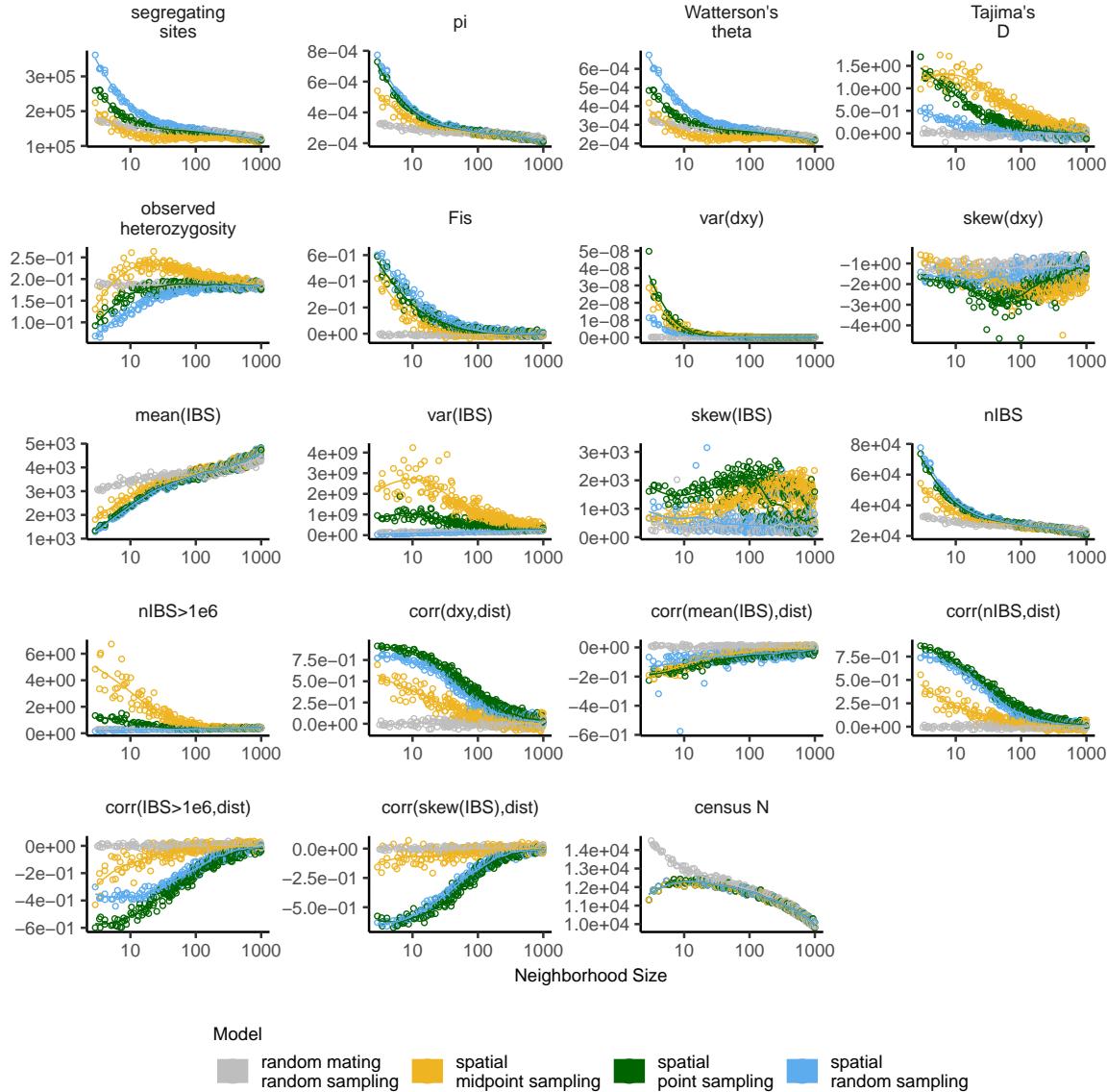


Figure S1 Change in summary statistics by neighborhood size and sampling scheme calculated from simulated sequence data of 60 individuals.

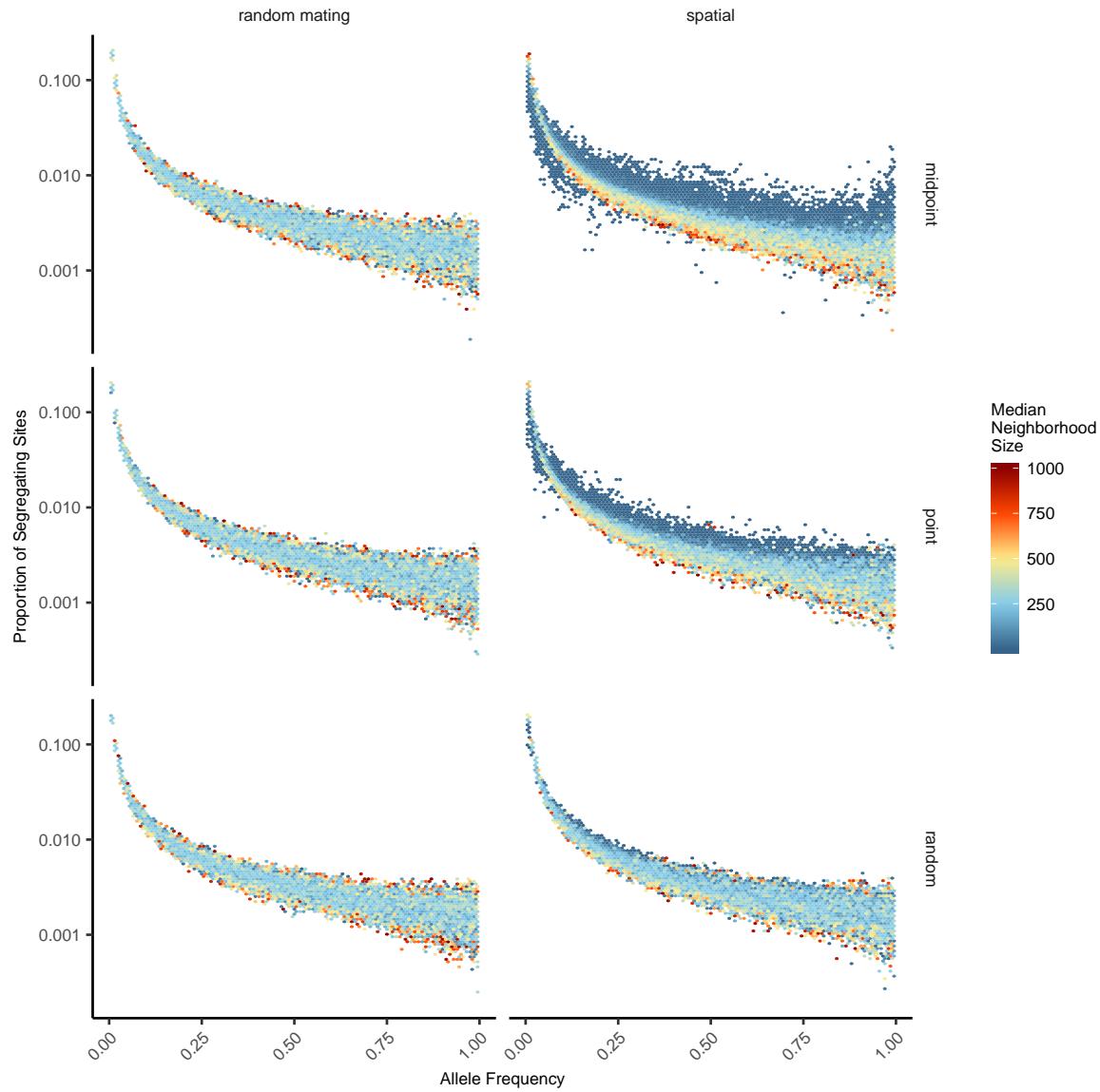


Figure S2 Site frequency spectra for random mating and spatial SLiM models under all sampling schemes.

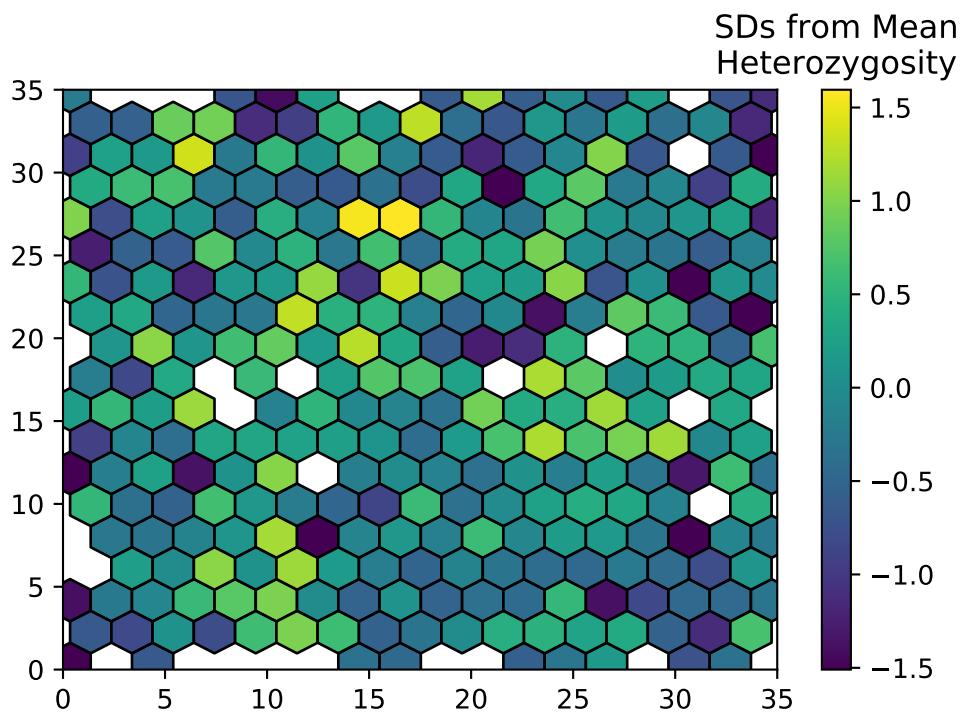


Figure S3 Normalized mean observed heterozygosity by location across 200 randomly sampled individuals

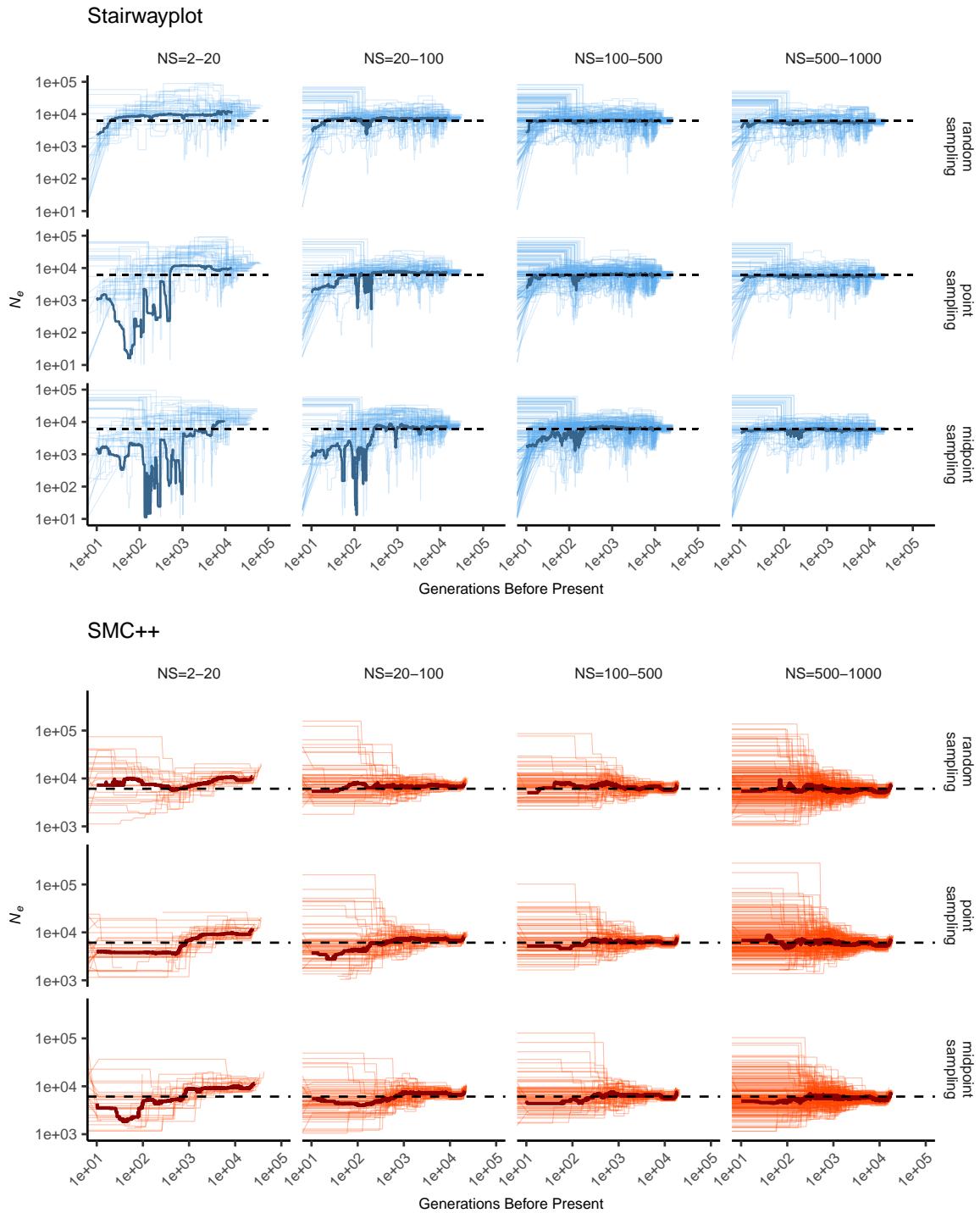


Figure S4 Inferred demographic histories for spatial SLiM simulations, by sampling scheme and neighborhood size (NS) range. Thick lines are rolling medians across all simulations in a bin and thin lines are best fit models for each simulation. Dashed horizontal lines are the average N_e across random-mating SLiM models estimated from θ_π .

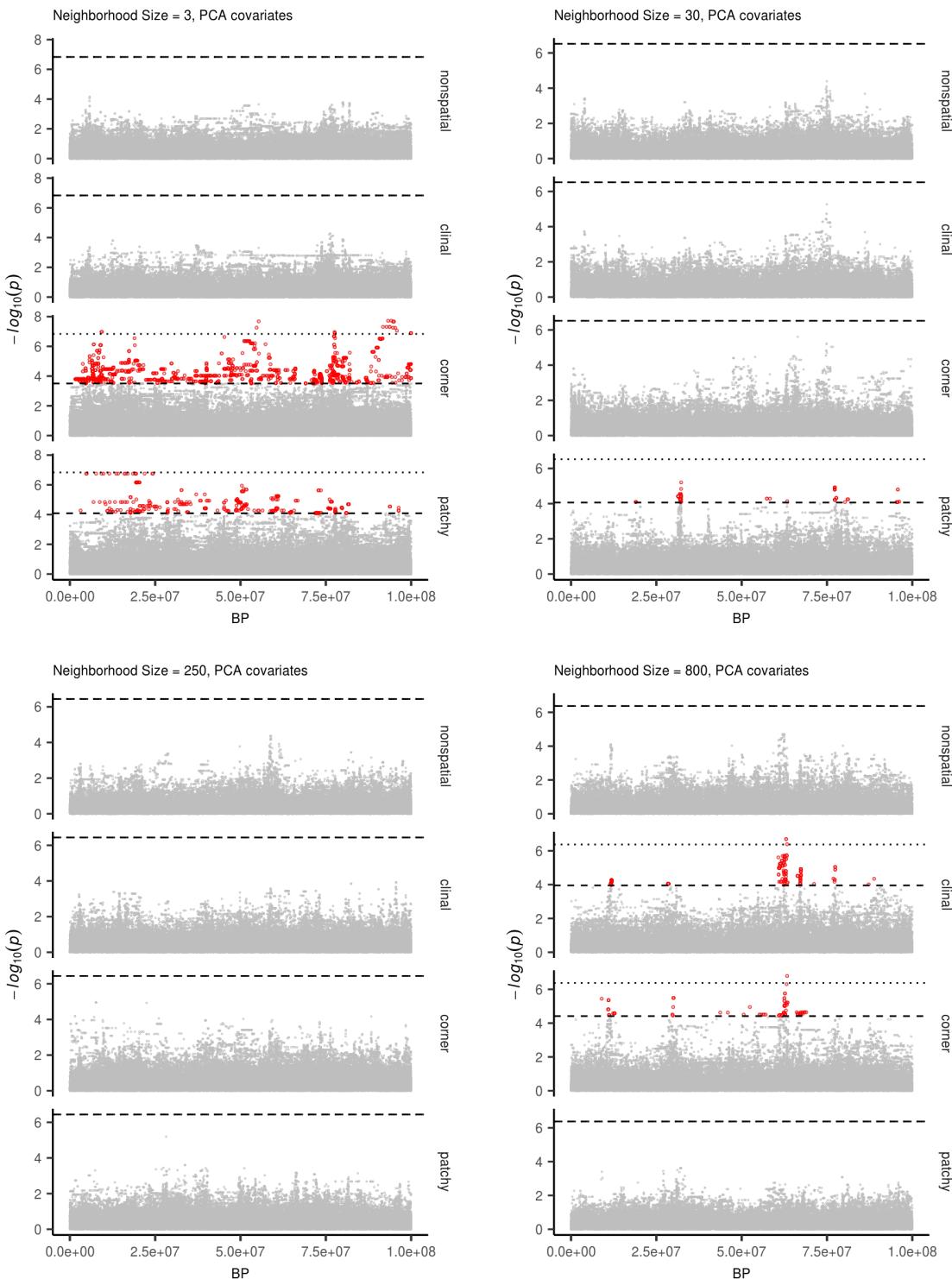


Figure S5 Manhattan plots for a sample of simulations at varying neighborhood sizes. Labels on the right of each plot describes the spatial distribution of environmental factors (described in the methods section of the main text). Points in red are significantly associated with a nongenetic phenotype using a 5% FDR threshold (dashed line). For runs with significant associations the dotted line is a Bonferroni-adjusted cutoff for $p = 0.05$.

Table S1 Summary statistics calculated on simulated genotypes.

Statistic	Description
Θ_{pi}	Mean of the distribution of pairwise genetic differences
Θ_W	Effective population size based on segregating sites
Segregating Sites	Total number of segregating sites in the sample
Tajima's D	Difference in Θ_{pi} and Θ_W over its standard deviation
Observed Heterozygosity	Proportion of heterozygous individuals in the sample
F_{IS}	Wright's inbreeding coefficient $1 - H_e / H_o$
$var(D_{xy})$	Variance in the distribution of pairwise genetic distances
$skew(D_{xy})$	Skew of the distribution of pairwise genetic distances
$mean(IBS)$	Mean of the distribution of pairwise identical-by-state (IBS) tract lengths taken over all pairs.
$var(IBS)$	Variance of the distribution of pairwise identical-by-state (IBS) tract lengths taken over all pairs.
$skew(IBS)$	Skew of the distribution of pairwise identical-by-state (IBS) tract lengths taken over all pairs.
$nIBS$	Mean number of IBS tracts with length > 2bp across all pairs in the sample.
$nIBS > 1e6$	Mean number of IBS tracts over 1×10^6 bp per pair across all pairs in the sample.
$corr(D_{xy}, dist)$	Pearson correlation between genetic distance and $\log_{10}(spatial\ distance)$
$corr(mean(IBS), dist)$	Pearson correlation between the mean of the IBS tract distribution for each pair of samples and $\log_{10}(spatial\ distance)$
$corr(nIBS, dist)$	Pearson correlation between the number of IBS tracts for each pair of samples and $\log_{10}(spatial\ distance)$
$corr(IBS > 1e6, dist)$	Pearson correlation between the number of IBS tracts > 1×10^6 bp for each pair of samples and $\log_{10}(spatial\ distance)$
$corr(skew(IBS), dist)$	Pearson correlation between the skew of the distribution of pairwise haplotype block lengths for each pair of samples and $\log_{10}(spatial\ distance)$

Table S2 Anova and Levene's test p values for differences by sampling strategy. Bolded values are rejected at $\alpha = 0.05$

variable	model	p(equal means)	p(equal variance)
segsites	random mating	0.998190	0.980730
Θ_π	random mating	0.997750	0.996450
Θ_W	random mating	0.998190	0.980730
Tajima's D	random mating	0.879690	0.188770
observed heterozygosity	random mating	0.531540	0.433230
F_{IS}	random mating	0.474790	0.785730
$mean(D_{xy})$	random mating	0.997770	0.996510
$var(D_{xy})$	random mating	0.283630	0.647240
$skew(D_{xy})$	random mating	0.958320	0.260750
$corr(D_{xy}, dist)$	random mating	0.601980	0.000000
$mean(IBS)$	random mating	0.997960	0.997730
$var(IBS)$	random mating	0.486450	0.399490
$skew(IBS)$	random mating	0.117980	0.069770
$nIBS$	random mating	0.997680	0.996570
$nIBS > 1e6$	random mating	0.834870	0.888730
$corr(mean(IBS), dist)$	random mating	0.073270	0.308420
$corr(IBS > 1e6, dist)$	random mating	0.268440	0.002100
$corr(skew(IBS), dist)$	random mating	0.396920	0.000620
$corr(nIBS, dist)$	random mating	0.581090	0.000000
segsites	spatial	0.000000	0.000000
Θ_π	spatial	0.026510	0.013440
Θ_W	spatial	0.000000	0.000000
Tajima's D	spatial	0.000000	0.000000
observed heterozygosity	spatial	0.000000	0.000000
F_{IS}	spatial	0.000000	0.000120
$mean(D_{xy})$	spatial	0.025390	0.012910
$var(D_{xy})$	spatial	0.004970	0.006230
$skew(D_{xy})$	spatial	0.000000	0.000000
$corr(D_{xy}, dist)$	spatial	0.000000	0.000000
$mean(IBS)$	spatial	0.272400	0.114250
$var(IBS)$	spatial	0.000000	0.000000
$skew(IBS)$	spatial	0.000000	0.000000
$nIBS$	spatial	0.033920	0.016640
$nIBS > 1e6$	spatial	0.000000	0.000000
$corr(mean(IBS), dist)$	spatial	0.000000	0.590540
$corr(IBS > 1e6, dist)$	spatial	0.000000	0.000000
$corr(skew(IBS), dist)$	spatial	0.000000	0.000000
$corr(nIBS, dist)$	spatial	0.000000	0.000000

Resubmission Cover Letter
Genetics

C. J. Battey,
Peter Ralph,
and Andrew Kern
Wednesday 11th September, 2019

To the Editor(s) –

We are writing to submit a revised version of our manuscript, “Space is the Place: Effects of Continuous Spatial Structure on Analysis of Population Genetic Data”.

Sincerely,

C. J. Battey, Peter Ralph, and Andrew Kern

Reviewer AE:

The manuscript admirably explores a lot of consequences of isolation-by-distance in the context of a novel model that is easily amenable to forward simulation; however, given that this model may be used in a lot of future studies based on the precedent set here, there is some concern about the model and its support. Reviewers 2 and 3 highlight this in particular (it underlies the main 2 points of reviewer 2's review, and the core of Reviewer 3's comment), and I agree. Whatever can be done to strengthen the standing of this model, and/or connect it to more thoroughly studied models, will be helpful for the manuscript. The concern would be that there are peculiarities of this model that do not generalize well. A new supplemental section or opener to the results section establishing the model more thoroughly would make the strongest response.

(I would generally cut down the quoted bits like the above to only what's essential, but haven't done that yet.)

(IMPORTANT: don't reorder or delete "points" below - it messes up the automatic numbering!)

(AE.1) Line 35: Also cite Wilkins and Wakeley, Genetics 2002; Wilkins 2004

Reply:

(AE.2) (p. ??, l. ??) "Such models have been used extensively in ecological modeling but rarely in population genetics " Detailing these previous uses via citations and elaboration may help alleviate the major concern about the provenance of this model and its unique behaviors (see general comments above and R2 and R3 comments).

Reply:

(AE.3) (p. ??, l. ??) Please describe computation time needed per replicate

Reply:

(AE.4) (p. ??, l. ??) I read the acknowledgement to the Hearth and Creative Sky Brewing with a sense of familiarity in feeling of gratitude to my own favorite cafes and breweries, but I it's not a great precedent for Acknowledgements to be filled this way. Please cut.

Reply: Good point; we have done this.

(AE.5) Figure 4: Show random-mating expectation

Reply:

(AE.6) Figure 3A, S2: Perhaps more revealing to show on log-log scale?

Reply:

(AE.7) Figure S3: Caption seems to be missing detail

Reply:

Reviewer 1:

This study explores biases arising in population-based inference when 1) real population samples are coming from spatial habitat with various degree of structuring while inference is made assuming random mating population; 2) imperfect sampling in practice that fails to represent full diversity across entire population habitat; 3) phenotypes that vary across geography and create spurious associations with genotypes. While earlier studies explored the effect of strong structure on population genetic inference and GWAS, this work focuses on less extreme scenarios of structuring that arises in populations evolving in continuous habitat. By using non-Wright-Fisher model, authors simulated chromosome-scale samples from populations that evolved in continuous space, and that can model environmental factors to create phenotypes varying over space. As a result, this study identified spatial structuring scenarios (small neighbourhood size 10-100) that coupled with imperfect sampling strategies lead to a biased inference of widely used population genetic statistics (altogether 18 statistics) such as pi (average pairwise sequence differences), heterozygosity (and inbreeding coefficient), and IBS tract sharing. Accordingly, inference of the effective population size history was also strongly affected under these parameter ranges. Finally, the authors use their spatial modelling to demonstrate that typical GWAS with PC-based correction cannot entirely remove spurious signals of genotype-phenotype associations arising from purely environmental factors. Overall, the authors explore an important but

often neglected source of bias that can affect inference in many population-based studies (in medical genetics, evolutionary biology and ecology). This study can be of interest to a broader audience of readership, and I have only minor comments to improve clarity and increase accessibility for readers:

(1.1) When neighbourhood size is small (10-100), the mean number of IBS tracts $> 2\text{bp}$ (n_{IBS} as in Table S1) is elevated similar to Wright's inbreeding coefficient, but mean of the distribution of pairwise IBS (mean(IBM)) is decreased. What could be the source of this discrepancy? How exactly mean(IBM) was calculated?

Reply:

(1.2) The authors use K to denote both carrying capacity (p. ??, l. ??) and population density (p. ??, l. ??). It might be better to use a different notation for these quantities since carrying capacity is fixed while density is an emergent quantity in the non-Wright-Fisher model. Use of K to denote carrying capacity and density is a bit confusing. For example, on (p. ??, l. ??) it is said that 'the "population density" (K) and "mean lifetime" (L) parameters were the same in all simulations'. Here K seems to indicate carrying capacity rather than density? The latter is an emergent quantity and varies across simulation runs?

Reply: We agree that this distinction is worth emphasizing! We've adjusted our language to hopefully remind the reader that K is a parameter that controls population density, rather than being equal to it, at (p. ??, l. ??) and (p. ??, l. ??) and (p. ??, l. ??).

(1.3) Concerning the non-Wright-Fisher model used, it would be helpful to emphasize that some of the parameters are emergent in contrast to Wright-Fisher model. For example, on Page 11, lines 306-308, the author's goal was to look at census size variation and variation in other quantities. This would be better understood if to emphasize that these parameters are emergent properties in the non-Wright-Fisher model used.

Reply:

(1.4) Page 9, line 242, Perhaps 'Demographic Inference' might better reflect the content of this section.

Reply:

(1.5) (p. ??, l. ??) This sentence with 'Gaussian noise with mean zero and standard deviation 10' is confusing since it was mentioned earlier that the modelled phenotype must vary as human height across Europe, and human height varies 2 standard deviations. Only after reading the whole paragraph it becomes clear that 'standard deviation 10' here refers to unit of height. Please consider rephrasing this sentence.

Reply:

(1.6) (p. ??, l. ??) In the sentence, 'We also examined p values for systemic inflation' I think the authors meant 'systematic inflation'.

Reply: Whoops; thanks. Fixed.

(1.7) Page 11, Please correct the legend in Figure 2: must be 'spatial model' and 'random mating' model.

Reply:

(1.8) Optional: a dashed line in Figure 2 that shows the total carrying capacity of $50*50*5=12500$ would be helpful.

Reply:

(1.9) Page 13, line 349, The phrase 'affect summaries of variation' is better to replace with 'summaries of genetic variation'.

Reply:

(1.10) Please add or correct references to supplementary figures: For example, Figure S2 was probably meant to accompany Figure 3A, while Figure S1 Figure 3B, but references in the text are absent. In fact, the first reference is made to Figure S3 on page 15.

Reply:

(1.11) There are also several typos and errors in the text. For example, on Page 12, lines 309; Page 27, line 655.

Reply:

Reviewer 2:

Battey et al. use spatially explicit population genetic simulations to analyze the effects of spatial structure on (i) the estimation of key population genetic parameters, in turn used to (ii) make inferences about population history, and on (iii) confounding in genome-wide association studies (GWAS). I Liked the paper a lot. It's interesting, well-written and addresses an important question - the effect of spatial population structure on population genetic statistics and inference-and I enjoyed reading it. The most positive aspects were:

1. It nice to actually see spatially explicit simulations and I'm happy that forward simulation is now fast enough that you can do this sort of thing.
2. The paper is very clear and well-written, easy to understand the motivation and most of the details. That's not always the case for this sort of paper.
3. I felt that the section about the effect on GWAS was the most interesting and novel part of the paper and gave me some intuition that I hadn't had before.

I don't have any major criticisms. There were a few aspects that I thought might warrant some additional discussion, and a few specific questions below. The general questions I had after reading it were:

(2.1) *To what extent are any of the results dependent on the exact method of simulation. There are a number of choices about the exact details of the simulations (e.g. the way the overlapping generations are handled, the edge effects and, particularly, the form of Equation 1 - see below). It's not so much that these are non-standard (since I don't think there is a standard) and they all sort of make sense heuristically, and I was left wondering whether these sorts of choices actually make a difference. Do the authors have some thoughts/intuition/results about that? Given that the results in Fig. 3 seem quite consistent with expectations, I suspect that on some level it doesn't make much difference but then there are intermediate results like Fig. 2 which seem a bit counter-intuitive and I wonder if those aspects depend on the simulation scheme.*

Reply:

(2.2) Related to the first point, to what extent are the results qualitatively different to those that would be obtained in a stepping-stone model? My interpretation is that they are actually very similar, but I didn't see whether that was explicitly discussed. In some sense, it's still easier to do large simulations in a stepping-stone model so it would be nice to be reassured that that's still ok.

Reply:

(2.3) The source of equation (1) is not obvious to me. I sort of see how it makes sense, but a little bit more intuition or a brief derivation or an illuminating either in the main text or the supplement, would be helpful.

Reply:

(2.4) The authors use a scaling factor in equation (2) to counteract the increase in fitness of individuals at the edges. Can they provide a figure showing that this is the case. What does "roughly" mean on line 164. Perhaps a heatmap of the fitness of individuals across the grid with and without the scaling factor?

Reply:

(2.5) It would be helpful provide the figure showing that generating mutations during the forward simulations in SLIM is equivalent to applying mutations using msprime on pre-generated trees (line 185)? It sounds like this procedure would underestimate the variance in the number of mutations, since you remove the effect of random generation time. Is this effect small?

Reply:

(2.6) Can the authors provide a bit more intuition behind the patterns of variation seen in generation time, census population size, and variance in the number of offspring with respect to neighborhood size seen in Figure 2? For example, it is not obvious to me why the census population size, for example, should decline systematically with respect to neighborhood size. Presumably this isn't just due to the local demographic stochasticity. Could the authors briefly interpret the observed patterns or cite appropriate literature?

Reply:

(2.7) Fig. 7D: I am surprised by the extent to which the observed values of $-\log_{10}(p)$ fall below the $y=x$ line. Particularly in the lower right panel for large neighbourhood sizes. I would expect that to be close to panmictic - why are the P-values underdispersed? That seems like a potential bug, or else something weird is going on.

Reply:

(2.8) Lines 706-716, It might be worth citing Haworth et al *Nature Communications* 2019 (<https://doi.org/10.1038/s41467-018-08219-1>) who do the proposed test (GWAS for birth location) in UK Biobank to illustrate the population structure.

Reply:

(2.9) The analysis and discussion around the effect of GWAS is focused on PCA correction. Do mixed models help at all?

Reply:

(2.10) The github link to the code didn't work for me. I assume it will be made public later, but at this point I can't tell whether the code is available/useable.

Reply:

Reviewer 3:

The present study deals with a "hot topic" in spatial population genetics. Most inferential and descriptive methods in statistical spatial population genetic rely on a discrete approximation of space and it is not clear what impact this approximation may have when individuals migrate along a continuum instead. Spatial patterns in sampling is also another major issue which is often simply dismissed, mainly because of the paucity of statistical methods to deal with it. This work touches on these important issues in a timely manner.

Although I was enthusiastic about the topic, I was quite disappointed with the core of the study, i.e., the forward-in-time simulation of populations in continuous space. The field has been struggling with this issue for decades – examples of spectacular failures like the Wright-Malecot model (see Felsenstein's "pain

in the torus' article, 1975) or, more recently, the "mugration" or "discrete trait analysis" model in phylodynamics (see De Maio et al. 2015) have probably mostly harmed our research field – that one cannot make the economy of using a sound probabilistic model for generating geo-referenced genetic data. It does not seem to be the case here unfortunately.

(3.1) *First, the simulation starts with individuals distributed uniformly at random in space. Is there any indication that the three-step algorithm used here maintains this distribution during the course of evolution? If it does not, then is there any stationary regime and how many generations does one need to wait before reaching it? I do appreciate that the competitive interaction term was introduced in order to avoid seeing the "clumping" of individuals that hampers the Wright-Malecot model. Yet, just because there are no such clusters does not mean that the spatial distribution of individuals reaches a stable regime and that the distribution reached, if any, is reasonable from a biological perspective.*

Reply:

Second, the demographic process used here involves birth and death of individuals. Does the population survive asymptotically or, like any birth-death process, eventually dies with probability one? In fact, one needs to know a little about the dynamics of the population size to decide whether the corresponding process is reasonable from a biological standpoint.

(3.2)

Reply:

(3.3) *Third, it is not clear what the relationship between the expected lifespan and the probability of survival is. The expected lifespan, L , is first defined as the inverse of the expected number of offspring produced by a parent. The authors also define the probability of survival of a given individual at a given point in space, p_i . Hence, the expected lifespan at a point in space (and time) is the mean of a geometric distribution with parameter p_i , i.e., $1/p_i$. Now, it is far from being obvious what the relationship between these two approaches for defining the expected lifespan actually is.*

Reply:

(3.4) Also, the web page <https://github.com/petrelharp/spaceness> does not seem to exist so that I was not able to experiment with the forward-in-time generator used here unfortunately.

Reply:

(3.5) All in all, more efforts need to be made here in my opinion to show that the forward-in-time simulations generate sensible outcomes. Sensible in terms of the behavior of the population demography at equilibrium (provided such equilibrium indeed exists) along with that of the spatial distribution of individuals. The authors could provide some guarantee of the good behavior of their model as evidenced from simulations using a broad range of parameter values for generating data. Alternatively, they could elect to use the spatial-Lambda-Fleming-Viot model for their simulations, which, in my opinion would seem the most sensible option given that (1) it is possible to run backward-in-time simulations under this model, thereby saving a lot of computation time and (2) it is a well-studied model with good mathematical and biological properties and (3) it is implemented in a publicly available software program (<https://github.com/jeromekelleher/discsim>)

Reply:

(3.6) Figure 2: I do not understand why the neighborhood size varies to the same extent in the random mating model as it does for the spatial model. For the random mating model, I would have expected the neighborhood size to be equal to the census size since all individuals have the same probability of being a parent of any given offspring. From lines 166->171, it is clear that the spatial model would converge to the random mating model when the mean parent-offspring distance tends to infinity only if we were to ignore the impact of range edges. I am thus wondering whether the variation of neighborhood size one observes in Fig 2 for the random mating model is just a consequence of border effects. If that is the case, then the authors should state it clearly and try to justify it from a biological perspective.

Reply:

(3.7) Line 729-731: "Many more species occur in a middle range of neighborhood sizes between 100 and 1000 - a range in which spatial processes play a minor role in our analyses [...]" Do the authors think that the spatial processes would still play a minor role when neighborhood sizes exceed 100-1000 if the habitat was larger than that taken in the present

simulations? It would also probably be useful to mention that neighborhood sizes given in Table 1 should be compared with extreme caution since the size of the corresponding habitats vary across species. More generally, I suspect that the size of the habitat has a substantial impact on the vast majority of statistics examined in this study. Indeed, the mean parent-offspring distance, which is at the core of the definition of Wright's neighborhood size, is only small or large relative to the size of the habitat.

Reply:

(3.8) Line 753-757: please add a reference to Guindon, Guo and Welch (2016). This study clearly shows that population density and dispersal parameters are identifiable and can indeed be estimated in practice under the spatial Lambda-Fleming-Viot model.

Reply:

Reviewer 4:

The manuscript by Battey et al explores the consequence of a well-known violation to population genetic models: the fact that populations are spatially structured and mate along a geographical cline, rather than randomly. This topic is important, particularly in light of recent working describing how spatially correlated genetic and environmental impacts can confound some population genetic insights, such as positive selection for height in Europe. The analyses and investigations presented here are thorough and sensible, and my comments are primarily intended to broaden accessibility for this interesting topic.

(4.1) Introduction. The discussion is very clear, articulating the three primary goals of the project: the impact of failing to model spatial population structure on 1) population genetic summary statistics, 2) inference on demographic history from population genetic data, and 3) impacts on GWAS summary statistics. I found the discussion a bit easier to follow than the introduction and would suggest streamlining and introducing the topic a bit more. Since the paper follows the flow described in the discussion, it might help orient readers by introducing these topics in the same order.

Reply:

(4.2) I agree that most modern work describes structure as discrete populations connected by migration. However, some methods/studies have explicitly modeled spatial structure, e.g. especially in ecology or using methods like *dadi* (diffusion approximations). Highlighting some examples of previously identified structure not possible to infer without modeling geography would be helpful to contextualize this work.

Reply:

(4.3) There is some reference to spatial models using grids (e.g. Rousset 1997). Some additional discussion contextualizing more recent methods like EEMS that also construct demes and model migration through divergence between neighboring demes would be helpful and interesting.

Reply:

(4.4) Demographic modeling. Both approaches tested, *stairwayplot* and *SMC++*, are most sensitive to older demographic events, and consequently are very noisy and underestimate effect population sizes, especially in smaller neighborhood sizes. Models that consider haplotype structure are much better suited to this time period. It would be helpful to either 1) discuss the varying time sensitivities of different classes of demographic inference methods and how spatial patterns of genetic variation would influence these inferences, or 2) apply a method of this class (many options, e.g. *DoRIS*, *IBDNe*, *Tracts*, *Globetrotter*, etc) and show how it performs.

Reply:

(4.5) GWAS mixed models. To what extent can spatial signals (e.g. corner, patchy) be corrected with mixed models, e.g. with PCs and PC-adjusted GRM as in Conomos et al, 2016 using PC-AiR and PC-Relate? Is patchiness related to dispersal? I'm curious how this relates to the predictive ability of GWAS phenotypes with some spatial association that may or may not be associated with environmental effects.

Reply:

(4.6) Code availability. This github link doesn't work, but is important to be able to evaluate for review: <https://github.com/petrelharp/spaceness>

Reply:

(4.7) Definitions and interpretations. There are quite a large number of metrics discussed in Figure 3B, and it's a lot to take in. It might be helpful to have a table with a reminder of what the metric is, its interpretation, and how it is computed.

Reply:

(4.8) Notation: "Offspring disperse a Gaussian-distributed distance away from the parent with mean zero and standard deviation σ in both the x and y coordinates. Each offspring is produced with a mate selected randomly from those within distance 3σ , with probability of choosing a neighbor at distance x proportional to $\exp(-x^2/2\sigma^2)$." I think x may be overloaded here, or I'm confused. Clarify?

Reply:

(4.9) When introducing the "spatial model" as opposed to this "random model," the more concrete illustration in Figure 1 hasn't yet been referenced, which makes it harder to follow. It would be helpful to introduce this figure with the model. Additionally, when Figure 1 is introduced, the order is from right to left (random, then point, then midpoint). It would be helpful to rearrange the figure to mirror what's in the text.

Reply:

(4.10) Not sure I follow this example: "Concretely, an individual at position (x, y) in a 50×50 landscape has mean phenotype $100 + 2x/5$."

Reply:

(4.11) Minor typo (through vs though): "This occurs because, even through the "population density" (K) and "mean lifetime" (L) parameters..."

Reply:

(4.12) Define NS abbreviation in Figure 5.

Reply: