

Estimating the Lambda measure in multiple-merger coalescents

Verónica Miró Pina^{1,2,3}, Émilien Joly⁴, and Arno Siri-Jégousse^{3*}

¹Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and
Technology, Barcelona, Spain.

²Universitat Pompeu Fabra (UPF), Barcelona, Spain

³Instituto de Investigaciones en Matemáticas Aplicadas y Sistemas, Universidad
Nacional Autónoma de México, CDMX, México.

³Centro de Investigación en Matemáticas, AC (CIMAT), Guanajuato, México

*e-mail: arno@sigma.iimas.unam.mx

Abstract

Multiple-merger coalescents, also known as Λ -coalescents, have been used to describe the genealogy of populations that have a skewed offspring distribution or that undergo strong selection. Inferring the characteristic measure Λ , which describes the rates of the multiple-merger events, is key to understand these processes. So far, most inference methods only work for some particular families of Λ -coalescents that are described by only one parameter, but not for more general models. This article is devoted to the construction of a non-parametric estimator of the density of Λ that is based on the observation at a single time of the so-called Site Frequency Spectrum (SFS), which describes the allelic frequencies in a present population sample. First, we produce estimates of the multiple-merger rates by solving a linear system, whose coefficients are obtained by appropriately subsampling the SFS. Then, we use a technique that aggregates the information extracted from the previous step through a kernel type of re-construction to give a non-parametric estimation of the measure Λ . We give a consistency result of this estimator under mild conditions on the behavior of Λ around 0. We also show some numerical examples of how our method performs.

1 Introduction

Inferring the evolutionary history of a sample of a present-time population is a very challenging task for both biological and mathematical communities. If no information is available on the past of the population, it is natural to consider Markov processes to model its genealogical tree. The family of Λ -coalescents, also called multiple-merger coalescents [Donnelly and Kurtz, 1999, Pitman, 1999, Sagitov, 1999], provides a rich class of processes to this purpose. They appear naturally as limit genealogies in population models, with possibly highly variable reproductive success [Sagitov, 1999]. Their flexibility (and richness) comes from the fact that their dynamics are fully described by a measure Λ on the interval $[0, 1]$, in the sense that Λ determines all the rates of multiple-merger coalescences. The intrinsic exchangeability of the branches at a given time in Λ -coalescents makes them particularly suited for

35 modeling genealogies in a very general framework. The family of Λ -coalescents contains the Kingman
36 coalescent [Kingman, 1982b,a], which is considered as the neutral model for random genealogies. We
37 refer to Berestycki [2009], Gnedin et al. [2014] for a mathematical survey that contains most of the
38 justifications of the classical formulas that we use in this paper.

39 Some subclasses of Λ -coalescent, such as Beta-coalescents [Schweinsberg, 2003] are intensively used
40 in practice since they are quite easily interpretable and are characterized by only one or two parameters.
41 The inference task is then reduced to the estimation of this parameter, that represents the level of repro-
42 ductive skewness in the population. This family is now well admitted to fit to populations with highly
43 skewed offspring distribution, e.g. marine species [Sigurgíslason and Árnason, 2003, Niwa et al., 2016],
44 tuberculosis cells [Menardo et al., 2020] or cancer cells [Kato et al., 2017]. Additionally, the class of
45 Beta-coalescents contains the Bolthausen-Sznitman coalescent that provides a genealogical process for
46 modeling populations under strong selection effects [Neher and Hallatschek, 2013, Desai et al., 2013,
47 Schweinsberg, 2017, Cortines and Mallein, 2017, Schertzer and Wences, 2023]. Another class of coa-
48 lescents that has been used for modelling populations with skewed offspring distribution are the Dirac-
49 coalescents [Eldon and Wakeley, 2006], which are characterized by two parameters.

50 Most of the previous inference work on Λ -coalescents has focused on estimating the characterizing
51 parameters of the one parameter families cited above [Birkner et al., 2013, Eldon et al., 2015, Koskela,
52 2018, Hobolth et al., 2019]. But the inference of more general Λ measures is still a challenging issue. To
53 our knowledge, it was only treated in Koskela et al. [2018], with Bayesian non-parametric methods. The
54 authors consider a forward in time model of type frequencies, more adapted to data sampled at multiple
55 time points, that are difficult to obtain from natural populations. In this paper we present an estimator that
56 uses only contemporaneous observations.

57 Note that the family of Λ -coalescents does not contain more general genealogical models for popula-
58 tions with varying size [Freund, 2020], under bottlenecks [González Casanova et al., 2021], with diploid
59 reproduction mechanism [Birkner et al., 2018], or recombination, which are part of a larger class of coa-
60 lescents, offering more possibilities for inference and model selection (see for example [Eldon et al.,
61 2015, Blath et al., 2016, Koskela, 2018, Sainudiin and Véber, 2018, Koskela and Wilke Berenguer, 2019,
62 Freund and Siri-Jégousse, 2021]).

63 In this work, we use approximations via Bernstein polynomials to establish a non-parametric estima-
64 tor of the density of Λ , which is built from estimates of the multiple merger rates of the Λ -coalescent
65 process. These estimates are the solutions of a linear system where the coefficients are the expected
66 branch lengths (or SFS) estimated from different resamplings of the observed genomic data of the sam-
67 ple. This technique relies on generalizations of classical recursions for functionals of the coalescent (see
68 Proposition 2.1). By reconstructing the density function of the measure Λ , and by visualizing it, our
69 technique should help to detect some unexpected behavior of the population and to reject a whole class
70 of coalescent models, e.g. Beta coalescents.

71 The type of observations that we use are SNP¹ matrices that summarize mutation data: at every
72 locus (column) where a mutation is observed, the individuals (lines) are assigned a 1 if they carry the
73 mutation and a 0 otherwise (see Table 2 as an example). This means that the method can only work for
74 polarized data i.e. data where derived and ancestral states can be tell apart. We assume that mutations
75 at each locus occur independently. Notice that this does not mean that the columns of the SNP matrix
76 are independent of each other, since the observed mutations all depend on the evolutionary history of the

¹Single Nucleotide Polymorphisms (SNP) are mutations that affect one single base position in the DNA. For the purpose of this paper, SNPs can be substitutions or indels.

77 sampled population. From the SNP matrix we can compute the SFS, which is the vector such that its
78 i -th component contains the number of observed mutations shared by exactly i individuals in the sample.
79 This is a statistic widely used in population genetics, both for model selection [Eldon et al., 2015, Freund
80 and Siri-Jégousse, 2021] and for parameter inference [Fu, 1995, Kersting et al., 2021]. In this paper, we
81 use an extension of the SFS, the so-called *weighted SFS*, which can also be read from the SNP matrix
82 by subsampling the individuals (lines), as explained in section 2.3. In the case of diploid (or polyploid)
83 species, our method requires phased data so that each line of the SNP matrix would represent a haplotype.

84 In addition, we assume given a collection of SNP matrices S_1, \dots, S_n governed by the same subjacent
85 Λ -coalescent. We also assume that those matrices are independent (in the probabilistic sense) so that no
86 coalescence in one genealogical branch of one matrix S_i have an influence on the other matrices. For
87 example, one can think of SNP matrices obtained from different replicates in an evolutionary experiment
88 or from multiple isolated populations, that have evolved independently, under similar conditions (see for
89 example Lenski et al. [1991] for *E. coli*, Johnson et al. [2021] for yeast or Teotónio et al. [2017] for *C.*
90 *elegans*). For applications to real populations, it is not possible to obtain independent SNP matrices from
91 the same population. However, ignoring dependencies between genomic regions that are at sufficient
92 recombination distance is a natural approximation and it has been used for example in McVean and
93 Cardin [2005] to define the sequentially Markov coalescent. In the case of multiple merger coalescent
94 models however, this assumption is less acceptable since the multiple merger events can in general affect
95 long-range dependence within chromosomes (see Koskela [2018], Korfmann et al. [2023]). By computing
96 the weighted SFS from this collection of SNP matrices, we can finely estimate rates of multiple merger
97 events in the Λ -coalescent.s

98 **Summary of the method**

99 In order to obtain a non-parametric estimator of the characteristic measure Λ , one has to follow these
100 steps, that will be explained in detail in the next section.

- 101 1. Using a subsampling technique, that is explained in section 2.3, compute the weighted SFS for
102 each of the n SNP matrices.
- 103 2. By averaging over the weighted SFS obtained from the different SNP matrices, obtain an estimator
104 of the branch lengths of the tree, as given in (10).
- 105 3. Solving the linear system in Proposition 2.1, obtain estimates of the rates of multiple merger events
106 in the Λ -coalescent.
- 107 4. Use a kernel based reconstruction (see (7)) to estimate the density of the characterizing measure of
108 the coalescent. Consistency of this estimator is proven in Theorem 2.1.

109 **2 Results**

110 Table 1 summarizes the notations introduced in this section and used through the paper.

111 **2.1 Definition of the estimator**

112 The Λ -coalescents are exchangeable continuous-time Markov chains taking their values in the partitions
113 of the integers. The chain jumps from one state to another by coagulating some blocks of the partition to

<u>Integers/ indices:</u>	
m	number of individuals (or haplotypes) in the sample
n	number of (almost) independent observations
<u>Multiple merger coalescents:</u>	
Λ, ν	characteristic measure ($\Lambda(dx) = x^2 \nu(dx)$)
$\lambda_{k,j}$	rates of multiple merger events (or coalescent rates)
$r_{k,j}$	combinatorial coalescent rates
d_k	sum of the combinatorial coalescent rates
<u>Density reconstruction:</u>	
$(B_j^m)_{0 \leq j \leq m}$	Bernstein polynomials of degree m
h	window parameter
ϕ	regularizing function
K_h	kernel of window h
K_h^x	kernel of window h around x
$P_{m,h}^x$	sequence of Bernstein polynomials used to approximate K_h^x
$\nu_{m,h}$	estimator of the density ν
<u>Estimator of the coalescent rates:</u>	
$\hat{r}_{k,j}$	estimator of the combinatorial coalescent rates
\bar{c}	any partition of the integer m
\bar{c}_0	partition made of singletons (denotes the usual initial configuration of the coalescent process)
$L_j(\bar{c}_0)$	partial length of order j
X_j	genotype of individual j
$X_{j,k}$	allele carried by individual j at locus k (0: wild-type, 1: mutated)
$(M_1(\bar{c}_0), \dots, M_{m-1}(\bar{c}_0))$	site frequency spectrum of a sample of size m
$(M_1(\bar{c}), \dots, M_{m-1}(\bar{c}))$	weighted site frequency spectrum with weights \bar{c}

Table 1: Main notations

114 form one new block. The exchangeability implies that the jump rates only depend on the number of blocks
 115 of the partitions. Thus, the dynamics of this chain are totally described by the array $(\lambda_{k,j})_{k \geq j \geq 2}$ where
 116 $\lambda_{k,j}$ stands for the rate at which a given j -tuple, among k present blocks, coalesces into one new block.
 117 Moreover, the class of exchangeable coalescents is consistent. This means that a coalescent starting from
 118 a partition π made of n blocks, and restricted to the set of partitions of $[m]$, with $m \leq n$, has the same law
 119 as a coalescent starting from the partition π restricted to $[m]$. This consistency property implies that, for
 120 $k \geq j \geq 2$,

$$\lambda_{k,j} = \lambda_{k+1,j} + \lambda_{k+1,j+1} \quad (1)$$

121 which is equivalent to

$$\lambda_{k,j} = \int_0^1 x^{j-2} (1-x)^{k-j} \Lambda(dx) = \int_0^1 x^j (1-x)^{k-j} \nu(dx) \quad (2)$$

122 where ν is a measure on $[0, 1]$ satisfying $\int_0^1 x^2 \nu(dx) < \infty$, and $\Lambda(dx) = x^2 \nu(dx)$ (see Pitman [1999]).
 123 The measure Λ (or ν) fully describes the coalescent, in the sense that all the multiple merger rates can
 124 be computed by integrating against these measures, as shown by equation (2). One has to note that
 125 both measures Λ and ν characterize the coalescent, nevertheless the measure ν is easier to estimate than
 126 the measure Λ . Besides this fact, the estimation of ν can only insure a good estimation outside of an
 127 open interval containing 0 due to the presence of x^{-2} in its definition. In this work, we suppose that the
 128 measure ν is absolutely continuous with respect to the Lebesgue measure over the interval $(0, 1]$, that is
 129 $\nu(dx) = \nu(x)dx$. This assumption includes Beta coalescents, but excludes Dirac coalescents as well as
 130 the Kingman coalescent.

131 The rate at which the number of blocks of the coalescent jumps from k to j is

$$r_{k,j} = \binom{k}{j} \lambda_{k,j}. \quad (3)$$

132 These new rates are then referred to as the combinatorial coalescent rates and the total coalescent rate
 133 from a state with k blocks is given by

$$d_k = \sum_{j=2}^k r_{k,j}.$$

134 This section is focused on the practical reconstruction of the measure ν when estimations $(\hat{r}_{m,j})_{2 \leq j \leq m}$
 135 of the $(r_{m,j})_{2 \leq j \leq m}$ are given. Note that the number m is fixed so that the only information needed is a
 136 good estimation of the coalescent rates of a group of m individuals. The reconstruction of the density
 137 of the measure ν depends on a hyperparameter h that can be understood as a parameter controlling the
 138 smoothing. Mathematically, this reconstruction can be seen as a map $(\hat{r}_{m,j})_j \mapsto \hat{\nu}_{m,h}$ taking values in
 139 the continuous functions on $(0, 1]$, where h holds for a window hyperparameter in a kernel estimation
 140 technique that we explain in details below.

141 The particular form of the coefficient $\lambda_{m,j}$ obtained from (2) allows us to use Bernstein polynomials
 142 as a technical tool for the approximation step. The Bernstein polynomial basis is defined as a double
 143 index sequence given by

$$B_j^m(x) = \binom{m}{j} x^j (1-x)^{m-j}$$

144 for all $0 \leq j \leq m$. In the following, we call Bernstein polynomial of order m a polynomial that belongs
 145 to the vector space spanned by the $(B_j^m)_{0 \leq j \leq m}$ where m is fixed. The ideas behind the construction of the
 146 estimator $\hat{v}_{m,h}$ of v are the following.

147 1. Choose a non-negative even function ϕ that is continuous, has a compact support containing 0
 148 and such that $\int_{\mathbb{R}} \phi(x) dx = 1$. Such a function is called a regularizing function. Without loss of
 149 generality, we assume that the support of ϕ is included in a segment of the form $[-\tau, \tau]$. For $h > 0$,
 150 define $K_h(x) = h^{-1} \phi(h^{-1}x)$. The function K_h is called a kernel of window h . Classical functional
 151 theorems show that a continuous function f is well approximated by its smoothed version given by
 152 the convolution $f \star K_h$ defined, for any fixed x , by

$$f \star K_h(x) = \int_{\mathbb{R}} K_h(x-y) f(y) dy \xrightarrow{h \rightarrow 0} f(x).$$

153 The convolution will apply on the function v that has support included in $(0, 1]$ so the integration
 154 of the function $K_h^x : y \rightarrow K_h(x-y)$ is only on the segment $(0, 1]$. The parameter h is a smoothing
 155 parameter and it has to be chosen carefully to avoid underfitting (h too small) or overfitting (h too
 156 large). Then, the kernel K_h^x is a function that selects - in a smooth way - data points around x . The
 157 kernel itself is used to give weight to those points depending on their distance to the center point
 158 x . In general, a gaussian kernel (see section 3) is usually a good choice for practical reasons but
 159 many more kernels can be used in this step. In the formal theorem given below (see Theorem 2.1),
 160 the bounded support condition is important for technical reasons but the gaussian kernel mimics
 161 this assumption and then works in practice. One can see [Györfi et al., chapter 5] for a detailed
 162 discussion on the consistency of kernel estimators.

163 2. Since the function K_h^x is smooth on its compact support, it can be approximated by a sequence of
 164 Bernstein polynomials $P_{m,h}^x$ so that

$$\|K_h^x - P_{m,h}^x\|_{\infty} \xrightarrow{m \rightarrow \infty} 0.$$

165 Those polynomials have the explicit definition

$$P_{m,h}^x(y) = \sum_{j=0}^m \binom{m}{j} K_h^x\left(\frac{j}{m}\right) y^j (1-y)^{m-j}.$$

166 3. According to points 1. and 2., the quantity

$$\int_0^1 P_{m,h}^x(y) v(y) dy \tag{4}$$

167 is a good candidate for the estimator of the density of $v(x)$. From the definition in (4), it is not clear
 168 to see that it is completely defined by the values of the $r_{m,j}$. But it can be rewritten in the following
 169 way

$$\int_0^1 P_{m,h}^x(y) v(y) dy = \sum_{j=0}^m \binom{m}{j} K_h^x\left(\frac{j}{m}\right) \int_0^1 y^j (1-y)^{m-j} v(y) dy. \tag{5}$$

170 In the sum of the right hand side, the terms for $j = 0$ or $j = 1$ won't lead to quantities that only
 171 depend on $r_m = (r_{m,j})_{2 \leq j \leq m}$. This is not a problem if one takes x to be outside of the support of K_h
 172 for the terms $j = 0, 1$. This leads to the restriction $x - m^{-1} \geq h\tau$. Then, we define

$$v_{m,h}(x) = \sum_{j=2}^m K_h^x \left(\frac{j}{m} \right) r_{m,j}, \quad (6)$$

173 if $x \in [m^{-1} + \tau h, 1]$ and 0 elsewhere.

174 4. One has to note that $v_{m,h}(x)$ is still not an estimator of $v(x)$ since it is not computable from the data.
 175 In the next section, we explain how to compute estimators of the coalescent rates $\hat{r}_m := (\hat{r}_{m,j})_{2 \leq j \leq m}$.
 176 From (5) and (4), for any $x \in [m^{-1} + \tau h, 1]$, we define the final estimator $\hat{v}_{m,h}$ by a plug in of the
 177 estimated values $(\hat{r}_{m,j})_{2 \leq j \leq m}$,

$$\hat{v}_{m,h}(x) = \sum_{j=2}^m K_h^x \left(\frac{j}{m} \right) \hat{r}_{m,j} \quad (7)$$

178 and we define $\hat{v}_{m,h}(x) = 0$ for $x < m^{-1} + \tau h$. In particular, we see that the interval $[m^{-1} + \tau h, 1]$
 179 tends to $(0, 1]$ when we let $m \rightarrow \infty$ and $h \rightarrow 0$.

180 We establish in Theorem 2.1 that this estimator is consistent and we provide some control on the error.
 181 For that purpose we need to define the modulus of continuity. For a function $f : [0, 1] \rightarrow \mathbb{R}$ we define, for
 182 any $\delta > 0$,

$$\omega_f(x, \delta) = \sup_{y: |x-y| \leq \delta} |f(x) - f(y)|.$$

183 **Theorem 2.1.** *Let $\tau > 0$, $h > 0$ and $m \geq 2$. Let ϕ be a non-negative α -Hölder function of support $[-\tau, \tau]$
 184 and such that $\phi(x) \leq \mathbb{1}_{[-\tau, \tau]}(x)$. Define $K_h(x) = h^{-1} \phi(h^{-1}x)$ and for $x \in [m^{-1} + \tau h, 1]$, let $\hat{v}_{m,h}(x)$ as in
 185 (7). Then,*

$$\begin{aligned} \|\hat{v}_{m,h} - v\|_{1, [m^{-1} + h\tau, 1]} &\leq \sup_{x \in [m^{-1} + \tau h, 1]} \omega_v(x, \tau h) + C\Lambda[0, 1/m] \\ &+ 1.5c \frac{v[1/m, 1]}{h^{1+\alpha} m^{\alpha/2}} + \|\hat{r}_m - r_m\|_1 \end{aligned}$$

186 where $\|\cdot\|_{1, [m^{-1} + h\tau, 1]}$ holds for the L_1 norm on the interval $[m^{-1} + h\tau, 1]$, and C, c are some positive
 187 constants.

188 The proof of this Theorem can be found in Appendix A.

189 In words, there are four different sources of error in the reconstruction of v . The first one comes from
 190 the modulus of continuity of v . Since the measure v is continuous, this quantifies how much it varies
 191 in a small interval of size τh . This is the most general way of stating the theorem. In most of practical
 192 cases, the function v is sufficiently regular to weaken this assumption. In particular, if v is Lipschitz
 193 of constant k , we can replace the term $\omega_v(x, \delta)$ by $k\delta$ when x and δ are such that $x - \delta \geq m^{-1} + \tau h$ to
 194 insure to avoid any problem of discontinuity. Taking $\delta = \tau h$, the window parameter h can be seen as a
 195 smoothing effect parameter. As every smoothing parameter it will have to be balanced in such a way to
 196 avoid over-estimation or under-estimation. This first term in the bound drives us in taking h the smallest
 197 possible whereas the third term in the bound is governing h the other way around. The second source of

198 error comes from the mass of Λ close to 0. The third one to the mass of ν in $[1/m, 1]$. Observe that, since
 199 $\int_0^1 x^2 \nu(dx) < \infty$, we have that $\nu[1/m, 1] = o(m^2)$. This implies that we can always find a Hölder exponent
 200 α such that the term $\nu[1/m, 1]/(h^{1+\alpha} m^{\alpha/2})$ converges to 0. This also reflects the trade-off between the
 201 regularity of ϕ and the behavior of ν near 0. The fourth source of error is error in the estimation of the
 202 $r_{m,j}$'s. The next section is devoted to providing a method to estimate these coalescent rates.

203 2.2 Techniques for the estimation of the coalescent rates

204 Let us turn to the problem of estimating the $m - 1$ combinatorial coalescent rates $(r_{m,j})_{2 \leq j \leq m}$. The
 205 estimator $\hat{r}_m = (\hat{r}_{m,j})_{2 \leq j \leq m}$ is defined as the solution of a linear system of $m - 1$ equations. Set $e_j =$
 206 $(0, \dots, 0, 1, 0, \dots, 0)$ where the value 1 is at j -th coordinate.

207 To do so, we consider the set of partitions of the integer m , $\mathcal{C}_m = \{\bar{c} = (c_1, \dots, c_m), \sum_j j c_j = m\}$.
 208 We denote the number of blocks of $\bar{c} \in \mathcal{C}_m$ by $|\bar{c}| = \sum_j c_j$. Then, we define the multi-dimensional block
 209 counting process of the coalescent with initial configuration \bar{c} as $(\Pi(t, \bar{c}))_{t \geq 0} = (N_1(t, \bar{c}), \dots, N_m(t, \bar{c}))_{t \geq 0}$,
 210 where $N_j(t, \bar{c})$ holds for the number of blocks of size j in the coalescent at time t . The total number of
 211 blocks in the coalescent process at time t is then given by $|\Pi(t, \bar{c})| = \sum_j N_j(t, \bar{c})$.

212 The partial length of order j ($1 \leq j \leq m - 1$) of the coalescent starting from \bar{c} is defined by

$$L_j(\bar{c}) = \int_0^{D(\bar{c})} N_j(t, \bar{c}) dt$$

213 where $D(\bar{c})$ is the depth of the tree, also called time to the most recent common ancestor (TMRCA),
 214 $D(\bar{c}) = \inf\{t \geq 0, \Pi(t, \bar{c}) = e_m\} = \inf\{t \geq 0, |\Pi(t, \bar{c})| = 1\}$.

215 **Proposition 2.1.** Let $\bar{c}_0 = (c_1, \dots, c_m) = m.e_1 = (m, 0, 0, \dots)$ and $\bar{b}_j = (m - j).e_1 + e_j, j \in \{2, \dots, m\}$.
 216 Then we have, for $i \in \{2, \dots, m\}$,

$$d_m \mathbb{E}[L_i(\bar{c}_0)] = c_i + \sum_{j=2}^{m-1} r_{m,j} \mathbb{E}[L_i(\bar{b}_j)], \quad i \in \{1, \dots, m - 1\}. \quad (8)$$

217 Equations (8) provide a linear system of $m - 1$ equations for $m - 1$ unknown values $(r_{m,j})_{2 \leq j \leq m}$ which
 218 can be rewritten as

$$AR = B,$$

219 where A is a $(m - 1) \times (m - 1)$ matrix such that the last column contains the expected branch lengths,
 220 $(\mathbb{E}[L_i(\bar{c}_0)])_{1 \leq i \leq m-1}^T$ and for $j \in \{1, m - 2\}$, the j -th column contains the difference between the expected
 221 branch lengths and partial lengths starting from b_j , i.e. $(\mathbb{E}[L_i(\bar{c}_0)] - \mathbb{E}[L_i(\bar{b}_j)])_{1 \leq i \leq m-1}^T$; B is the vector
 222 $(m, 0, \dots, 0)$ and $R = (r_{m,j})_{2 \leq j \leq m}$.

223 The precision in the estimation of the vector r_m will thus depend on the precision of the estimation of
 224 the expectations in (8). The proof of Proposition 2.1 is a consequence of the Markov property applied at
 225 the first jump time (see Appendix B).

226 2.3 Obtaining the partial lengths from the SNP matrix

227 In this section we provide a method to estimate the expected values in (8). We assume that mutations
 228 happen at rate θ , as a Poisson process along the branches of the genealogical tree resulting from the

229 coalescent process considered in the previous section.² This mutation rate can also be seen as a scaling
 230 parameter since the law of the mutations is not changed if one multiplies this parameter θ and the multiple
 231 merger rates are all multiplied by a common constant. Since the data is only considering the mutations
 232 affecting individuals in the sample, this scaling effect is untraceable, resulting in the model being not
 233 identifiable. To avoid dealing with those issues, we fix θ to be 1 in the sequel, which corresponds to
 234 measuring time in units of mutation rate. In the infinite sites model, we suppose that each mutation affects
 235 a distinct site so that they can all be observed in present-time data. This is a reasonable assumption since
 236 in most species the mutation rate is low, i.e. mutations are in nature rare events. In addition, multiple
 237 mergers generally reduce the time to the most recent common ancestor, which reduces the number of
 238 generations at which mutation events can occur.

239 We assume provided a sample of m individuals, whose genotypes are encoded in vectors X_1, \dots, X_m .
 240 Formally, the vector X_j describing the genotype of individual j is a vector $X_j = (X_{j,k})_k \in \{0, 1\}^p$, $p \in \mathbb{N}$
 241 of binary values where $X_{j,k} = 1$ if and only if the individual j carries the mutation at site k (also called
 242 the SNP k). All this information is stored in a SNP matrix, where the vectors X_j are written as rows.
 243 See Table 2 for an example. Observe that the size of the vectors is not relevant here as the matrix can be
 244 adjusted by withdrawing columns of 0's, which are the sites where no genetic variation is detected.

245 We also define the site frequency spectrum $(M_1(\bar{c}_0), M_2(\bar{c}_0), \dots, M_{m-1}(\bar{c}_0))$ as the vector such that
 246 $M_i(\bar{c}_0)$ counts the number of mutations carried by exactly i individuals, that is $M_i(\bar{c}_0) = \text{Card}\{k, \sum_{j=1}^m X_{j,k} =$
 247 $i\}$. Also, the total number of mutations observed in the sample is $M(\bar{c}_0) = \text{Card}\{k, \sum_{j=1}^m X_{j,k} \geq 1\}$.
 248 Notice that these quantities depend on the initial configuration (\bar{c}_0) . In fact, considering a sample
 249 of m different individuals corresponds to starting the coalescent process from m singletons, i.e. from
 250 $\bar{c}_0 = m \cdot e_1 = (m, 0, \dots, 0)$. Our assumption on $\theta = 1$ leads to

$$\mathbb{E}[M_i(\bar{c}_0)] = \mathbb{E}[L_i(\bar{c}_0)].$$

251 As a result, any estimator of the expected value $\mathbb{E}[M_i(\bar{c}_0)]$ is also an estimator of the expected value of
 252 the partial length of order i .

Ind.	SNP 1	SNP 2	SNP 3	SNP 4	SNP 5	SNP 6	SNP 7
X_1	1	0	0	0	0	0	1
X_2	1	0	0	0	0	1	1
X_3	1	0	0	0	0	1	1
X_4	0	0	1	1	0	0	0
X_5	0	0	1	1	1	0	0
X_6	0	1	0	0	0	0	0

Table 2: Example of a SNP matrix, where 0 denotes the ancestral type and 1 denotes the derived type. Here the sample size is $m = 6$ and the total number of mutations is $M(\bar{c}_0) = 7$.

253 In order to make use of the linear system (8), one has to be able to estimate the partial length $\mathbb{E}[L_i(\bar{b}_j)]$
 254 of a coalescent tree started at any partition \bar{b}_j of the integer m of the form $\bar{b}_j = (m - j)e_1 + e_j$, $j \in$
 255 $\{2, \dots, m\}$. This can be done by a trick that consists in observing the site frequency spectrum differently

² θ is the population size rescaled mutation rate, i.e. θ is proportional to $N_e \times \mu$, where N_e is the effective population size and μ is the mutation rate of the sequence, per generation. In humans, the mutation rate per nucleotide per generation is of the order of 10^{-8} and N_e is of the order of 10^4 , which would correspond to $\theta = 1$ if the sequence length is of order 10^4 .

256 and that is based on the consistency and the Markov properties of the coalescent. Indeed, the law of
 257 the coalescent tree starting from a partition \bar{c} is the same as that of the coalescent tree starting from $|\bar{c}|$
 258 individuals where individuals are associated to weights given by the partition \bar{c} . We will therefore define
 259 a *weighted SFS*, by subsampling m' individuals with replacement. The weight of each individual is given
 260 by the number of times it appears in the subsample. A mutation carried by an individual with weight j
 261 will count as a mutation carried by j individuals. The weighted SFS corresponding to a subsample with
 262 weights \bar{c} is denoted by $(M_1(\bar{c}), \dots, M_{m'-1}(\bar{c}))$.

263 Let us give an example using the SNP matrix of Table 2. In this matrix $m = 6$, $M(\bar{c}_0) = 7$ and the
 264 observed SFS is $(2, 3, 2, 0, 0)$. Now, think of $\bar{c} = \bar{b}_3 = 3e_1 + e_3$, i.e. a sample where three individuals have
 265 weight 1 and one individual has weight 3. We chose at random one individual to which we associate the
 266 block of size 3, e.g. X_1 and three individuals to which we associate weight 1 (X_2, X_3 and X_4). Then the
 267 weighted version of X_1 is $(3, 0, 0, 0, 0, 3)$. The corresponding weighted SFS is $(2, 1, 0, 0, 2)$. Recall that,
 268 due to the exchangeability of the coalescent process, the law of the weighted SFS does not depend on the
 269 labels of the individuals which are assigned to each weight.

270 Formally, for a general partition \bar{c} of any integer $m' \geq m$ such that $|\bar{c}| \leq m$, let (h_1, \dots, h_m) be any
 271 partition vector associated to \bar{c} , such that $\sum h_j = m'$, $\text{Card}\{j, h_j = i\} = c_i$ and $\text{Card}\{j, h_j = 0\} = m - |\bar{c}|$.
 272 Then,

$$M_i(\bar{c}) = \text{Card}\left\{k, \sum_{j=1}^m h_j X_{j,k} = i\right\}. \quad (9)$$

273 Observe that $m' = m$ in the example above.

274 Now, suppose that we can observe n i.i.d. SNP matrices of m individuals and that from each of these
 275 n matrices we obtain a realization of $(M_1^k(\bar{c}), \dots, M_{m'-1}^k(\bar{c}))$ for $k \in \{1, \dots, n\}$ using the subsampling
 276 procedure described above. The estimator of $\mathbb{E}[L_i(\bar{c})]$ is given by

$$\hat{L}_i(\bar{c}) = \frac{1}{n} \sum_{k=1}^n M_i^k(\bar{c}). \quad (10)$$

277 3 A numerical example

278 We apply our method to simulated data, where the coagulation measure Λ is known, so we can compare
 279 the estimated measure $\hat{v}_{m,h}$ to the true measure v . The goal of this section is not to quantify the error in
 280 the reconstruction, which is already bounded in Theorem 2.1, but rather to illustrate our method.

281 We used the population genetics simulator `msprime` [Baumdicker et al., 2022, Kelleher et al., 2016].
 282 There are two Λ -coalescent models available in `msprime`: the Beta-coalescent [Schweinsberg, 2003]
 283 and the Dirac-coalescents [Eldon and Wakeley, 2006]. Since we have assumed that the Λ measure has
 284 a (Lebesgue) density, we chose the Beta-coalescent model. It is a one parameter family in which the
 285 characterizing parameter α denotes the degree of skewness of the offspring distribution of the population.
 286 Smaller values of α correspond to a more skewed offspring distribution, while $\alpha \rightarrow 2$ corresponds to the
 287 Kingman coalescent. For this example we fixed $\alpha = 1.2$, but we obtain qualitatively similar results for
 288 other values such that $1 < \alpha < 2$, which is the range of values allowed by `msprime` (Supplementary
 289 Figures S2 and S3). Observe that in this model the characteristic measure $v(x)$ behaves like $x^{-1-\alpha}$
 290 near 0, which implies it has no first moment, providing a "worst-case scenario" for the non-parametric
 291 estimation.

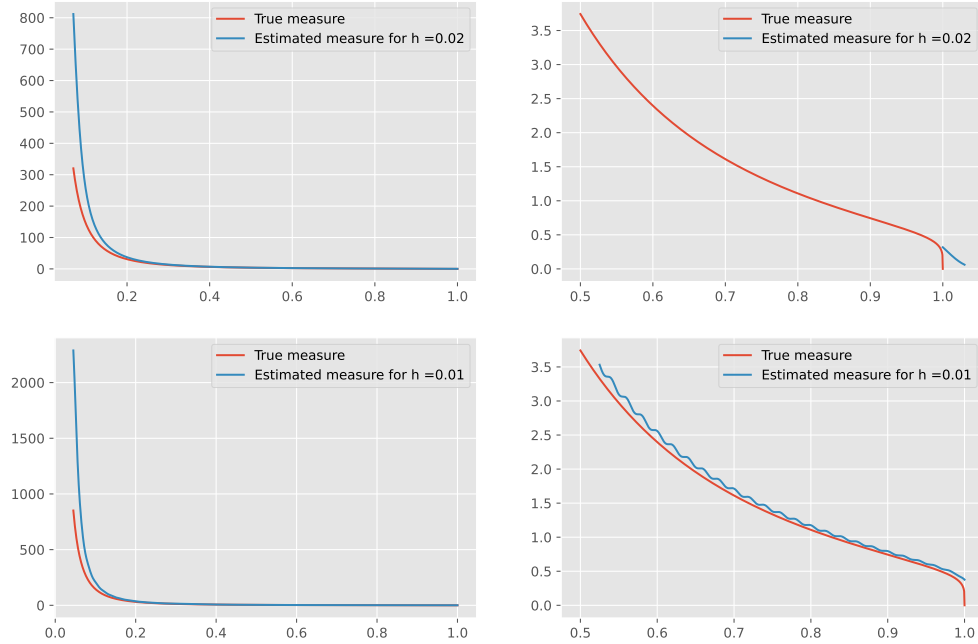


Figure 1: **Density of ν** . Here the true value of the coalescent rates ($r_{m,j}$'s) is used to reconstruct the density of ν using equation (6) (blue curves). We compare against the true density of the measure ν (red curves). We used a Beta-coalescent model with $\alpha = 1.2$ and $m = 50$ and we show two different values of h . Panels on the right are zoomed versions of the panels on the left (notice the difference in the x-axis).

292 We simulated the SNP matrix of a sample of size m . We estimated the partial lengths by computing
 293 the SFS and the weighted SFSs, as explained in the previous section. We simulated n independent SNP
 294 matrices, that we used for inference. We explored different settings for m and n . Then we solved the
 295 linear system (8) in Proposition 2.1. Since some of the coalescent rates (the $r_{m,j}$'s) are small, sometimes
 296 solving the linear system $AR = B$ with the estimated values of the partial lengths $\hat{L}_i(\bar{c})$'s yielded negative
 297 values of the $r_{m,j}$'s, so, instead we used `scipy.optimize`, to minimize $AR - B$ with the constraint that
 298 the $r_{m,j}$'s must be positive.

299 To reconstruct the measure ν , we followed our method using a Gaussian kernel for K_h . Figure 1
 300 compares the true density of ν and the density of ν reconstructed using the true values of the $r_{m,j}$'s (using
 301 equation (6)). We can observe that there is a rather large error close to 0. This is because the family of Λ
 302 measures that we were able to simulate using `msprime` have a lot of mass close to 0 (which corresponds
 303 to the second term of the error computed in Theorem 2.1). We show results for two different values of
 304 h . We observe that the error in the reconstruction depends on h , and that for $m = 50$, $h = 0.02$ is a good
 305 parameter choice.

306 In a second step, we applied the whole method to our simulated data, i.e. we first estimated the $r_{m,j}$'s

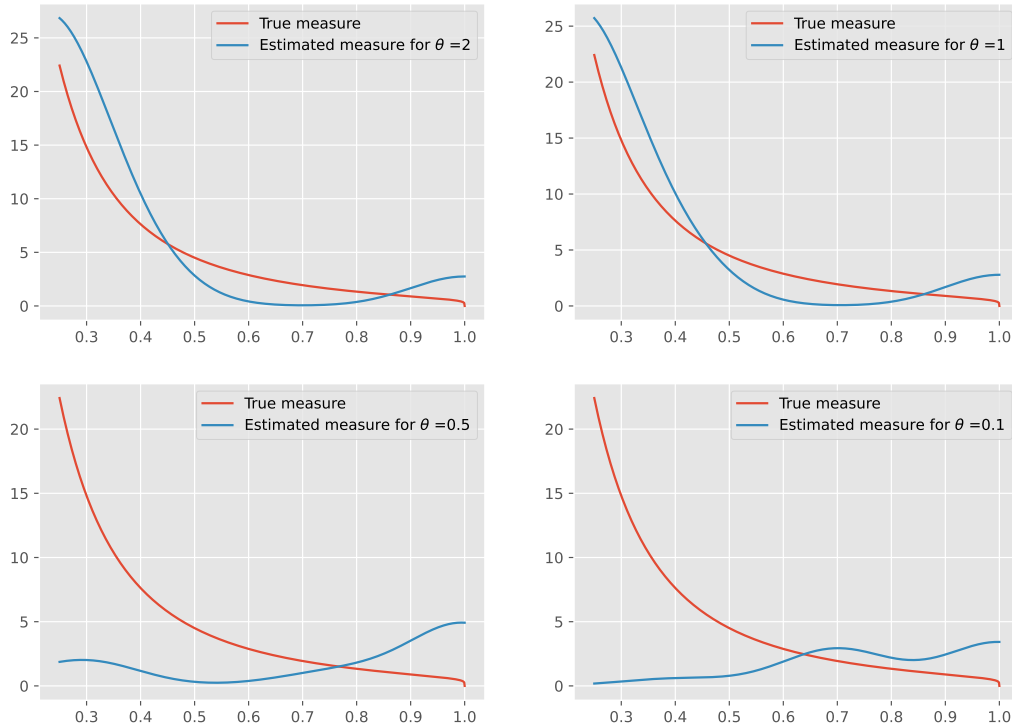


Figure 2: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 10$, $h = 0.1$, $\tau = 2$, $\alpha = 1.2$.

307 from our simulated SNP matrices and then we reconstructed the measure v .

308 We simulated different values of the sample size m and (population-size rescaled) mutation rate θ .
 309 Our method does not allow to estimate the mutation rate θ , but as discussed in the previous section,
 310 changing θ corresponds to changing the time-scale, i.e. multiplying the $r_{m,j}$'s by some constant. Since
 311 the aim of this project is not to get this constant, in order to visualize the results, we normalized the
 312 measures by the integral of Λ in $[m^{-1} + h\tau, 1]$, which is an approximation for the total coalescence rate,
 313 i.e. the constant we need.

314 We fixed the number of independent replicates to $n = 20$. This would correspond, for a species with
 315 20 chromosomes, to consider that each chromosome is an almost independent replicate. Figure 2 shows
 316 the obtained results, which strongly depend on the mutation rate. For θ larger than 1 (which corresponds
 317 more or less to 500 mutations per sample in our model), the estimation of v is quite good, in the sense
 318 that it has a similar behavior to the true v . However, for lower mutation rates, the estimated measure does
 319 not match with the real measure.

320 If we increase the number of independent samples, the quality of our estimates improves (Figure 3).

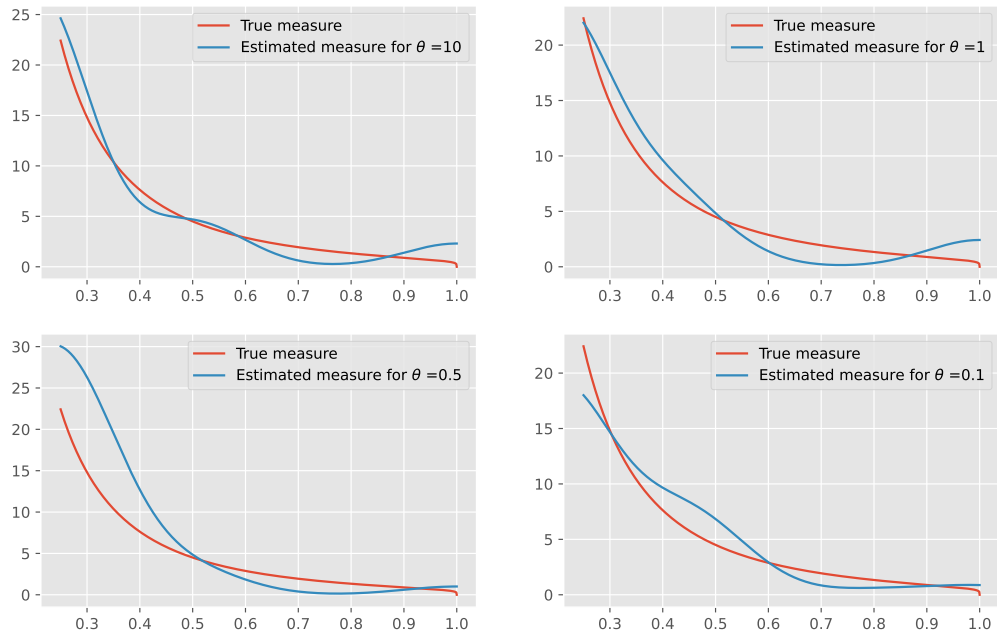


Figure 3: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 100$, $m = 10$, $h = 0.1$, $\tau = 2$, $\alpha = 1.2$.

321 However, obtaining such a large number of almost independent samples, with enough mutation on each
 322 requires either a high-mutation high-recombination regime that is difficult to observe from real data. This
 323 parameter regime seems more likely to be achieved if the independent samples are replicates from an
 324 evolutionary experiment.

325 Supplementary Figure S1 shows the real values of the multiple merger rates $r_{m,j}$'s and their estimated
 326 value $\hat{r}_{m,j}$'s. We can observe that, although the error in the estimated multiple merger rates is large
 327 (especially when r is close to m , the density of the reconstructed measure still shows the same trend as
 328 the true density.

329 Finally, increasing m does not necessarily provide better estimates (Supplementary Figures S6 and
 330 S7). This is because the effect of increasing the sample size is counteracted by an increase in the size of
 331 the linear system given by (8) and the errors in the estimation of the partial lengths propagate through the
 332 system. Decreasing m (Supplementary Figures S4 and S5) can also increase the reconstruction error (see
 333 Theorem 2.1).

334 4 Discussion

335 We have developed the first non-parametric estimator of the density function of the measure Λ of Λ -
336 coalescents from observations taken at a single time point. Our method can infer any measure Λ , in a
337 way that is agnostic to biology, since we are not making any assumption on how the population evolves.
338 It could be used to hint or exclude already known coalescent classes such as Beta coalescents or the
339 Bolthausen-Sznitman coalescent which are linked to biological phenomena such as skewed offspring
340 distribution or certain modes of selection. It provides a well-fitting model for observed genetic diversity,
341 which can be used to find species or genomic regions where other processes additionally, on top of
342 those modelled by classical Λ -coalescents (Beta, Dirac or or the Bolthausen-Sznitman) influence genetic
343 diversity. We are also able to quantify the consistency of this estimator in our main theorem. Note that
344 our method fits better models with a measure Λ having a non-explosive behavior close to 0, providing a
345 convincing rejection tool of the family of Beta coalescents.

346 The method works in three steps. First, we estimate the branch lengths of the coalescent from SNP
347 matrices thanks to a new observable that we called the weighted SFS. Second, we obtain the coalescent
348 rates from the branch lengths by inverting the linear system (8). Third, we estimate the coalescent mea-
349 sure from the coalescent rates thanks to a non-parametric method based on Bernstein polynomials, as
350 explained in the first paragraph of the Results section. Our (new) mathematical results focus only on
351 step 3. Notice that steps 1 and 2 may be useful for other statistical techniques, in particular if SFS based
352 statistics are needed, or if coalescent rates are needed.

353 The error generated at step 1 does only come from the central limit theorem so it is minimized if we
354 can observe data from a large number of independent replicates or if the mutation rate is high. It can
355 increase a lot if those assumptions do not hold. We tried the weighted SFS method by by simulating
356 genome-wide data from a single sample and treating distant genomic regions as independent (pseudo-
357)samples. To do so, we first simulated the ancestry of a very large population using `msprime` and then
358 subsampled individuals with or without replacement. This yielded large errors in the estimation of the
359 expected branch lengths. This is because of the shared evolutionary history between individuals of a
360 same population, that makes the ancestry trees of the different samples highly correlated (especially for
361 the deeper nodes).

362 The error produced at step 2 is due to the propagation of the error in step 1 at the moment of inverting
363 a linear system. This part might be improved by using some more efficient numerical tools, or by adding
364 constraints or equations in the system, e.g. recursive equations derived from (1). A more promising track
365 to follow is that we may need not all the rates $r_{m,j}$ since they get too small when j goes close to m and do
366 not bring substantial information for the estimation of Λ .

367 Finally, we managed to bound the error generated at step 3 in Theorem 2.1 thanks to the theory of
368 non-parametric statistics. The error naturally explodes close to 0 when the measure Λ accumulates too
369 much mass at 0. This is typically the case in Beta-coalescents with parameter α between 1 and 2. Thus,
370 obtaining a non-explosive measure close to 0 could also mean that the evolution of the population does
371 not fall into the domain of attraction of Beta coalescents.

372 In this work we did not focus on the estimation of the mutation rate θ . Recall that it is not possible
373 to discriminate between the couple (θ, Λ) and $(c\theta, c\Lambda)$ for some constant c from one date set taken at a
374 single time. We refer to Eldon et al. [2015], Freund and Siri-Jégousse [2021] for some further techniques
375 involving estimation of the mutation rate and model selection.

376 5 Code availability

377 Code for the numerical example is available at : https://github.com/vmiropina/bootstrap_coalescent/.

378 6 Funding

379 Most of this work was conducted when VMP was a postdoc at UNAM, funded by the DGAPA-UNAM
380 postdoctoral program. VMP also acknowledges support of the Spanish Ministry of Science and Inno-
381 vation to the EMBL partnership, the Centro de Excelencia Severo Ochoa and the CERCA Programme /
382 Generalitat de Catalunya. The research of EJ was partially funded by CONACYT, Ciencia de Frontera
383 project 1043167. The research of ASJ was partially funded by UNAM-DGAPA PAPIIT grant IN104722.

384 References

- 385 F. Baumdicker, G. Bisschop, D. Goldstein, G. Gower, A. P. Ragsdale, G. Tsambos, S. Zhu, B. Eldon,
386 E. C. Ellerman, J. G. Galloway, A. L. Gladstein, G. Gorjanc, B. Guo, B. Jeffery, W. W. Kretzschmar,
387 K. Lohse, M. Matschiner, D. Nelson, N. S. Pope, C. D. Quinto-Cortés, M. F. Rodrigues, K. Saunack,
388 T. Sellinger, K. Thornton, H. van Kemenade, A. W. Wohns, Y. Wong, S. Gravel, A. D. Kern, J. Koskela,
389 P. L. Ralph, and J. Kelleher. Efficient ancestry and mutation simulation with msprime 1.0. *Genetics*,
390 220(3):iyab229, 2022.
- 391 N. Berestycki. Recent progress in coalescent theory. *Ensaïos Matematicos*, 16(1):1–193, 2009.
- 392 M. Birkner, J. Blath, and B. Eldon. Statistical properties of the site-frequency spectrum associated with
393 lambda-coalescents. *Genetics*, 195(3):1037–1053, 2013. ISSN 0016-6731. doi: 10.1534/genetics.113.
394 156612. URL <https://www.genetics.org/content/195/3/1037>.
- 395 M. Birkner, H. Liu, and A. Sturm. Coalescent results for diploid exchangeable population models. *Elec-*
396 *tron. J. Probab.*, 23(49):44, 2018.
- 397 J. Blath, M. C. Cronjäger, B. Eldon, and M. Hammer. The site-frequency spectrum associated with
398 ξ -coalescents. *Theor. Pop. Biol.*, 110:36–50, 2016.
- 399 A. Cortines and B. Mallein. A N -branching random walk with random selection. *ALEA, Lat. Am. J.*
400 *Probab. Math. Stat.*, 14:17–137, 2017.
- 401 M. M. Desai, A. M. Walczak, and D. S. Fisher. Genetic diversity and the structure of genealogies in
402 rapidly adapting populations. *Genetics*, 193(2):565–585, 2013.
- 403 P. Donnelly and T. G. Kurtz. Particle representations for measure-valued population models. *Ann.*
404 *Probab.*, 27(1):166–205, 1999.
- 405 B. Eldon and J. Wakeley. Coalescent processes when the distribution of offspring number among indi-
406 viduals is highly skewed. *Genetics*, 172(4):2621–2633, 2006. doi: 10.1534/genetics.105.052175.
- 407 B. Eldon, M. Birkner, J. Blath, and F. Freund. Can the site-frequency spectrum distinguish exponential
408 population growth from multiple-merger coalescents? *Genetics*, 199(3):841–856, 2015.

- 409 F. Freund. Cannings models, population size changes and multiple-merger coalescents. *J. Math. Biol.*,
410 80(5):1497–1521, 2020.
- 411 F. Freund and A. Siri-Jégousse. The impact of genetic diversity statistics on model selection between
412 coalescents. *Comput. Stat. Data Anal.*, 156:107055, 2021.
- 413 Y.-X. Fu. Statistical properties of segregating sites. *Theor. Pop. Biol.*, 48(2):172–197, 1995.
- 414 A. Gnedin, A. Iksanov, and A. Marynych. Λ -coalescents: a survey. *J. Appl. Probab.*, 51A(Celebrating 50
415 Years of The Applied Probability Trust):23–40, 2014. ISSN 0021-9002. doi: 10.1239/jap/1417528464.
416 URL <http://dx.doi.org/10.1239/jap/1417528464>.
- 417 A. González Casanova, V. Miró Pina, and A. Siri-Jégousse. The symmetric coalescent and Wright-Fisher
418 models with bottlenecks. *Ann. Appl. Probab.*, 32(2):235–268, 2021.
- 419 L. Györfi, A. Krzyżak, and H. Walk. *A distribution-free theory of nonparametric regression*, volume 1.
420 Springer.
- 421 A. Hobolth, A. Siri-Jégousse, and M. Bladt. Phase-type distributions in population genetics. *Theor. Pop.*
422 *Biol.*, 127:16–32, 2019.
- 423 M. S. Johnson, S. Gopalakrishnan, J. Goyal, M. E. Dillingham, C. W. Bakerlee, P. T. Humphrey,
424 T. Jagdish, E. R. Jerison, K. Kosheleva, K. R. Lawrence, J. Min, A. Moulana, A. M. Phillips,
425 J. C. Piper, R. Purkanti, A. Rego-Costa, M. J. McDonald, A. N. Nguyen Ba, and M. M. De-
426 sai. Phenotypic and molecular evolution across 10,000 generations in laboratory budding yeast
427 populations. *eLife*, 10:e63910, jan 2021. ISSN 2050-084X. doi: 10.7554/eLife.63910. URL
428 <https://doi.org/10.7554/eLife.63910>.
- 429 M. Kato, D. A. Vasco, R. Sugino, D. Narushima, and A. Krasnitz. Sweepstake evolution revealed by
430 population-genetic analysis of copy-number alterations in single genomes of breast cancer. *Royal Soc.*
431 *Open Sci.*, 4(9), 2017. doi: 10.1098/rsos.171060. URL [http://rsos.royalsocietypublishing.](http://rsos.royalsocietypublishing.org/content/4/9/171060)
432 [org/content/4/9/171060](http://rsos.royalsocietypublishing.org/content/4/9/171060).
- 433 J. Kelleher, A. M. Etheridge, and G. McVean. Efficient coalescent simulation and genealogical analysis
434 for large sample sizes. *PLoS Comput. Biol.*, 12(5):e1004842, 2016.
- 435 G. Kersting, A. Siri-Jégousse, and A. H. Wences. Site frequency spectrum of the bolthausen-sznitman
436 coalescent. *ALEA. Lat. Am. J. Prob. Math. Stat.*, 18:1483–1505, 2021.
- 437 J. F. Kingman. On the genealogy of large populations. *J. Appl. Probab.*, pages 27–43, 1982a.
- 438 J. F. Kingman. The coalescent. *Stochastic Process. Appl.*, 13(3):235–248, 1982b.
- 439 K. Korfmann, T. Sellinger, F. Freund, M. Fumagalli, and A. Tellier. Simultaneous inference of past
440 demography and selection from the Ancestral Recombination Graph under the Beta Coalescent. *bioRxiv*
441 *preprint doi.org/10.1101/2022.09.28.508873*, 2023.
- 442 J. Koskela. Multi-locus data distinguishes between population growth and multiple merger coalescents.
443 *Stat. Appl. Genet. Mol. Biol.*, 17(3), 2018.

- 444 J. Koskela and M. Wilke Berenguer. Robust model selection between population growth and multiple
445 merger coalescents. *Math. Biosci.*, 311:1–12, 2019.
- 446 J. Koskela, P. A. Jenkins, and D. Spanò. Bayesian non-parametric inference for Λ -coalescents:
447 Posterior consistency and a parametric method. *Bernoulli*, 24(3):2122–2153, 2018.
- 448 R. Lenski, M. Rose, S. Simpson, and S. Tadler. Long-term experimental evolution in escherichia coli. i.
449 adaptation and divergence during 2,000 generations. *Am. Nat.*, 138(6):1315–41, 1991.
- 450 G. McVean and N. Cardin. Approximating the coalescent with recombination. *Philos. Trans. R. Soc.*
451 *Lond. B Biol. Sci.*, 1459(360):1387–93, 2005.
- 452 F. Menardo, S. Gagneux, and F. Freund. Multiple merger genealogies in outbreaks of mycobacterium
453 tuberculosis. *Mol. Biol. Evol.*, 38(1):290–306, 2020.
- 454 R. A. Neher and O. Hallatschek. Genealogies of rapidly adapting populations. *Proc. Natl. Acad. Sci.*
455 *U.S.A.*, 110(2):437–442, 2013.
- 456 H. S. Niwa, K. Nashida, and T. Yanagimoto. Reproductive skew in japanese sardine inferred from dna
457 sequences. *ICES J. Mar. Sci.*, 73(9):2181–2189, 2016.
- 458 J. Pitman. Coalescents with multiple collisions. *Ann. Probab.*, 27(4):1870–1902, 1999. ISSN 0091-1798.
459 doi: 10.1214/aop/1022677552. URL <http://dx.doi.org/10.1214/aop/1022677552>.
- 460 H. Queffelec and C. Zuily. *Analyse pour l'agrégation-Agrégation/Master Mathématiques*. Dunod, 2020.
- 461 S. Sagitov. The general coalescent with asynchronous mergers of ancestral lines. *J. Appl. Probab.*, 36(4):
462 1116–1125, 1999.
- 463 R. Sainudiin and A. Véber. Full likelihood inference from the site frequency spectrum based on the
464 optimal tree resolution. *Theor. Pop. Biol.*, 124:1–15, 2018.
- 465 E. Schertzer and A. H. Wences. Relative vs absolute fitness in a population genetics model. how stronger
466 selection may promote genetic diversity. *arXiv preprint arXiv:2301.07762*, 2023.
- 467 J. Schweinsberg. Coalescent processes obtained from supercritical galton–watson processes. *Stochastic*
468 *Process. Appl.*, 106(1):107–139, 2003.
- 469 J. Schweinsberg. Rigorous results for a population model with selection ii: genealogy of the population.
470 *Electron. J. Probab.*, 22, 2017.
- 471 H. Sigurgíslason and E. Árnason. Extent of mitochondrial dna sequence variation in atlantic cod from
472 the faroe islands: a resolution of gene genealogy. *Heredity*, 91(6):557–564, 2003. doi: 10.1038/sj.hdy.
473 6800361. URL <https://doi.org/10.1038/sj.hdy.6800361>.
- 474 H. Teotónio, S. Estes, P. Phillips, and C. Baer. Experimental evolution with caenorhabditis nematodes.
475 *Genetics*, 206(2):697–716, 2017.

476 7 Appendix A: Proof of Theorem 2.1

477 In this section we prove the consistency of the proposed estimator. One source of error comes from the
 478 uncertainty on the terms $\hat{f}_{m,j}$. There are also two other sources of approximation error given by the fact
 479 that we approximate v by $v \star K_h$ and that we approximate K_h^x by $P_{m,h}^x$. To prove the theorem, we need to
 480 tackle these three sources of error. The first results that we present deals with the smooth approximation
 481 of v by $v \star K_h$.

482 **Proposition 7.1.** *Let ϕ be such that for all $x \in \mathbb{R}$, $\phi(x) \leq \mathbb{1}_{\{[-\tau, \tau]\}}(x)$, for some $\tau > 0$. Then, for any*
 483 *$\delta > 0$ and any $x \in (0, 1]$ such that $x - \delta \geq m^{-1} + \tau h$,*

$$|v(x) - v \star K_h(x)| \leq \omega_v(x, \delta) + 4 \sup_{x \in [m^{-1} + \tau h, 1]} |v(x)| \cdot \left(\tau - \frac{\delta}{h} \vee 0 \right).$$

484 *Proof.* Using the fact that K_h is of integral 1, we see that for any $\delta > 0$ and any x ,

$$\begin{aligned} |v(x) - v \star K_h(x)| &\leq \int_{\mathbb{R}} |v(x) - v(x+y)| K_h(y) dy \\ &= \int_{-\delta}^{\delta} |v(x) - v(x+y)| K_h(y) dy + \int_{\mathbb{R} \setminus [-\delta, \delta]} |v(x) - v(x+y)| K_h(y) dy. \end{aligned}$$

485 Observe that the last term is well defined since the support of K_h is included in $[-\tau h, \tau h]$. Then,

$$\begin{aligned} |v(x) - v \star K_h(x)| &\leq \omega_v(x, \delta) + 2 \sup_{x \in [m^{-1} + \tau h, 1]} |v(x)| \int_{\mathbb{R} \setminus [-\delta, \delta]} K_h(y) dy \\ &= \omega_v(x, \delta) + 4 \sup_{x \in [m^{-1} + \tau h, 1]} |v(x)| \int_{\delta}^{+\infty} K_h(y) dy. \end{aligned}$$

486 This last integral is obviously 0 if the support of K_h and $[\delta, +\infty)$ are disjoint i.e., if $\tau h < \delta$. Otherwise,

$$\int_{\delta}^{+\infty} K_h(y) dy \leq \int_{\delta}^{h\tau} \frac{1}{h} \phi\left(\frac{y}{h}\right) dy \leq \int_{\delta}^{h\tau} \frac{1}{h} dy = \tau - \frac{\delta}{h}.$$

487

□

488 Next, we deal with the approximation of K_h^x by the Bernstein polynomials $P_{m,h}^x$.

489 **Proposition 7.2.** *Let f be a continuous function on $[0, 1]$. Then*

$$\|f - P_m(f)\|_{\infty} \leq 1.5 \|\omega_f(\cdot, m^{-1/2})\|_{\infty}$$

490 where $P_m(f) = \sum_{j=0}^m \binom{m}{j} f\left(\frac{j}{m}\right) y^j (1-y)^{m-j}$.

491 The proof can be found in Theorem II, p. 518 of Queffélec and Zuily [2020] where the constant can
 492 be taken equal to 1.5 from the proof page 519. Proposition 7.2 will be applied on the kernel function K_h .
 493 It is obvious to see that for any $\delta > 0$,

$$\omega_{K_h^x}(y, \delta) = \omega_{K_h}(y, \delta) = h^{-1} \omega_{\phi}(y, h^{-1} \delta).$$

494 Thus the regularity of K_h^x is directly inherited from the regularity of ϕ . For example, if ϕ is α -Hölder,

$$\|\omega_{K_h^x}(\cdot, m^{-1/2})\|_\infty \leq \frac{c}{h^{1+\alpha} m^{\alpha/2}} \quad (11)$$

495 for some constant $c > 0$.

496 The following simple result is the verification that the plug in strategy of (7) is valid when one controls
497 the L_1 error of estimation on the vector $r_m = (r_{m,j})_{2 \leq j \leq m}$.

498 **Proposition 7.3.** *For any $h > 0$ and any $m \geq 2$ one has that*

$$\|\hat{v}_{m,h} - v_{m,h}\|_1 \leq \|\hat{r}_m - r_m\|_1.$$

499 *Proof.* The proof of the proposition is obtained by noticing that the terms $K_h(x - j/m)$ have integral equal
500 to 1. More formally, notice that the support of $\hat{v}_{m,h}$ or $v_{m,h}$ is included in $[m^{-1} + \tau h, 1]$ and so

$$\begin{aligned} \|\hat{v}_{m,h} - v_{m,h}\|_1 &\leq \int_{m^{-1} + \tau h}^1 \left| \sum_{j=2}^m K_h^x \left(\frac{j}{m} \right) r_{m,j} - \sum_{j=2}^m K_h^x \left(\frac{j}{m} \right) \hat{r}_{m,j} \right| dx \\ &\leq \int_{m^{-1} + \tau h}^1 \sum_{j=2}^m K_h^x \left(\frac{j}{m} \right) |r_{m,j} - \hat{r}_{m,j}| dx \\ &\leq \sum_{j=2}^m |r_{m,j} - \hat{r}_{m,j}| \int_{\mathbb{R}} K_h^x \left(\frac{j}{m} \right) dx = \|\hat{r}_m - r_m\|_1, \end{aligned}$$

501 which concludes the proof. □

502 When one sums up the previous results we end with the following result of consistency.

503 *Proof of Theorem 2.1.* Noticing that for any x ,

$$|\hat{v}_{m,h}(x) - v(x)| \leq |\hat{v}_{m,h}(x) - v_{m,h}(x)| + |v_{m,h}(x) - v \star K_h(x)| + |v \star K_h(x) - v(x)|$$

504 and using the previous results, we see that it is enough to show that

$$|v_{m,h}(x) - v \star K_h(x)| \leq C\Lambda[0, 1/m] + 1.5c \frac{v[1/m, 1]}{h^{1+\alpha} m^{\alpha/2}}.$$

505 But since x belongs to $[m^{-1} + h\tau, 1]$, we have that $v_{m,h}(x) = \int_0^1 P_{m,h}^x(y) v(y) dy$ and so

$$\begin{aligned} |v_{m,h}(x) - v \star K_h(x)| &\leq \int_0^1 |P_{m,h}^x(y) - K_h^x(y)| v(y) dy \\ &= \int_0^{\frac{1}{m}} \frac{|P_{m,h}^x(y)|}{y^2} y^2 v(y) dy + \int_{\frac{1}{m}}^1 |P_{m,h}^x(y) - K_h^x(y)| v(y) dy. \end{aligned}$$

506 On the one hand, the polynomial $|P_{m,h}^x(y)|/y^2$ is uniformly bounded so there exists a constant C such that

$$\int_0^{\frac{1}{m}} \frac{|P_{m,h}^x(y)|}{y^2} y^2 v(y) dy \leq C\Lambda[0, 1/m].$$

507 On the other hand, by Proposition 7.2 and (11), we have that

$$\int_{\frac{1}{m}}^1 |P_{m,h}^x(y) - K_h^x(y)| v(y) dy \leq \|K_h^x - P_{m,h}^x\|_\infty \int_{\frac{1}{m}}^1 v(y) dy \leq 1.5c \frac{v[1/m, 1]}{h^{1+\alpha} m^{\alpha/2}}$$

508

□

509 8 Appendix B: Proof of Proposition 2.1

510 In this section we prove Proposition 2.1 by establishing a more general recursion result. Consider any
511 initial configuration $\bar{c} \in \mathcal{C}_m$. By the strong Markov property (applied at the time of the first jump of the
512 process), we obtain that, for $i \in \{1, \dots, m-1\}$,

$$d_{|\bar{c}|} \mathbb{E}[L_i(\bar{c})] = c_i + \sum_{\bar{b}: \bar{b} \preceq \bar{c}} S_{\bar{c}, \bar{b}} \mathbb{E}[L_i(\bar{b})] \quad (12)$$

where $S_{\bar{c}, \bar{b}}$ is the rate at which the multi-dimensional block counting process jumps from \bar{c} to \bar{b} , i.e.

$$S_{\bar{c}, \bar{b}} = \lambda_{|\bar{c}|, \Sigma(c_j - b_j) \mathbf{1}_{b_j < c_j}} \prod_{j, b_j < c_j} \binom{c_j}{c_j - b_j},$$

513 already defined in Eq. (20) of Hobolth et al. [2019]. In the summation in (12), the notation $\bar{b} \preceq \bar{c}$ holds for
514 the vectors \bar{b} obtained by a single step merging of a multiple number of lineage into one. This imposes,
515 in particular, that there exists a unique index $i_0 \leq m$ such that $b_{i_0} > c_{i_0}$ and in this case, $b_{i_0} = c_{i_0} + 1$. This
516 last fact is actually a characterization of the partial order \preceq . An additional fact is that

$$d_{|\bar{c}|} = \sum_{\bar{b}: \bar{b} \preceq \bar{c}} S_{\bar{c}, \bar{b}}.$$

517 Proposition 2.1 is obtained in the particular example where the initial configuration only contains single-
518 tons, that is $\bar{c}_0 = m.e_1 = (m, 0, 0, \dots)$. In this case we recover the classical functional leading to the site
519 frequency spectrum. After the first merge, the $m-2$ attainable partitions in (12) are of the form

$$\bar{b}_j = (m-j)e_1 + e_j, \quad j \in \{2, \dots, m\}.$$

520 An important observation is that, in this case, $S_{\bar{c}_0, \bar{b}_j} = r_{m,j}$.

521 **Remark 8.1.** *There is a classical relation between the rates $(r_{m,j})_{2 \leq j \leq m}$ and the expectation of the total
522 branch lengths of the associated coalescents starting from $j \in \{2, \dots, m\}$ individuals, denoted by $B(j)$.
523 The total branch length is the sum of all the lengths of the coalescent tree from the leaves until the root.
524 It is also obtained thanks to the Markov property applied at the first jump time:*

$$d_m \mathbb{E}[B(m)] = m + \sum_{j=2}^{m-1} r_{m, m-j+1} \mathbb{E}[B(j)]. \quad (13)$$

525 *A first approach to estimate the coalescence rates would consist in considering this equation and add a
526 number of recursive equations (1) to obtain a linear system of $m-1$ equations. However this very rigid
527 strategy, based only on the total number of blocks of the coalescent, leads to uncontrolled error as the
528 error resulting from the estimate of $r_{m,2}$ will spread additively in the rest of the estimates. The strategy
529 based on recursive equations for the partial lengths of the coalescent, as proposed in Proposition 2.1,
530 leads to a more stable system.*

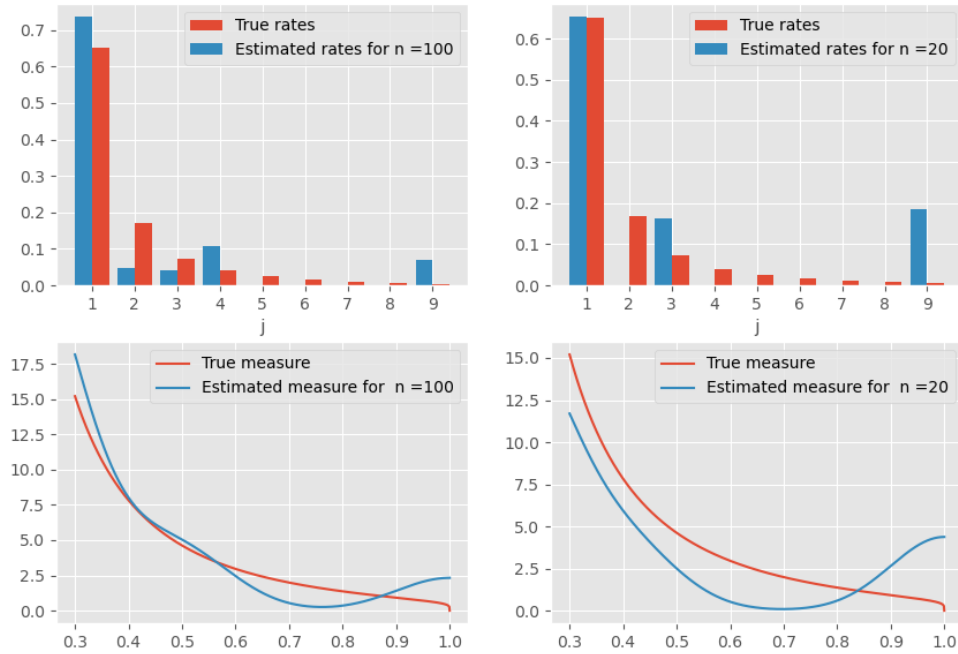


Figure S1: **Multiple merger rates and density of v .** In the upper panels, we compare the multiple merger rates $r_{m,j}$'s to the estimated values $\hat{r}_{m,j}$'s. In the bottom panels, we compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 50$, $h = 0.01$, $\tau = 2$, $\alpha = 1.2$.

531 9 Appendix C: Supplementary Figures

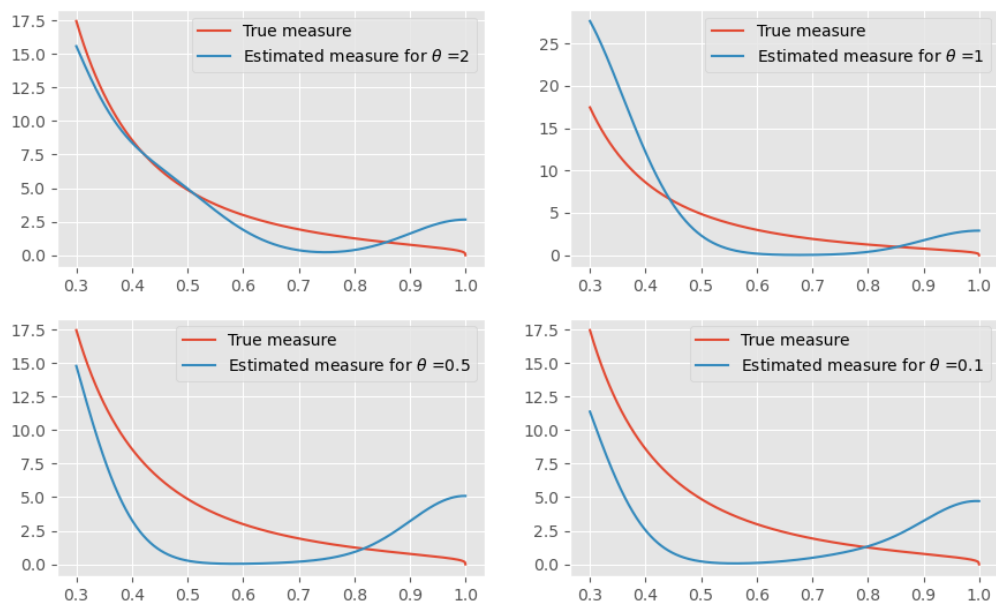


Figure S2: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 50$, $h = 0.01$, $\tau = 2$, $\alpha = 1.3$.

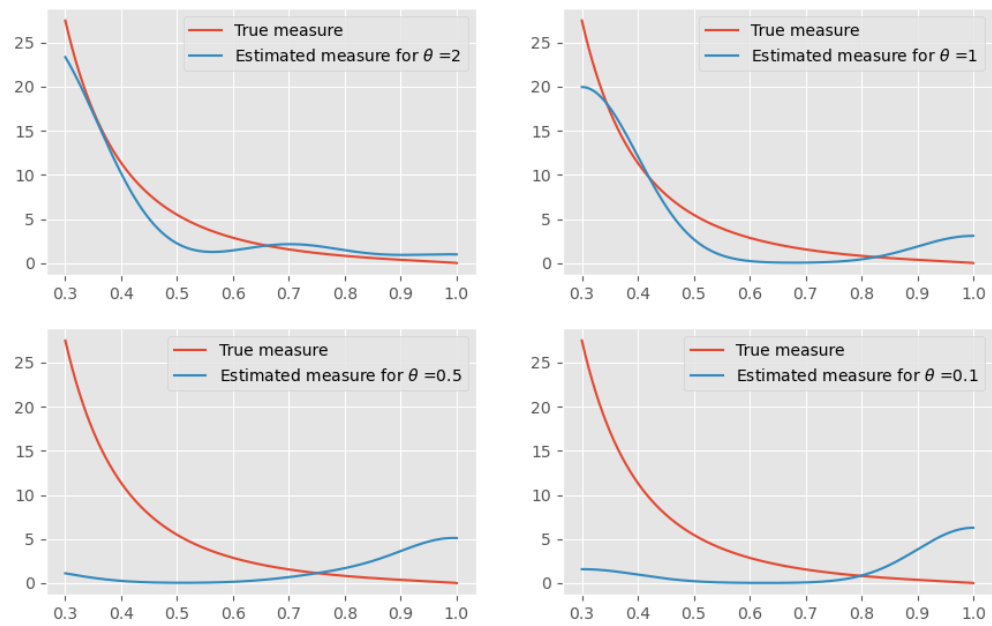


Figure S3: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 50$, $h = 0.01$, $\tau = 2$, $\alpha = 1.7$.

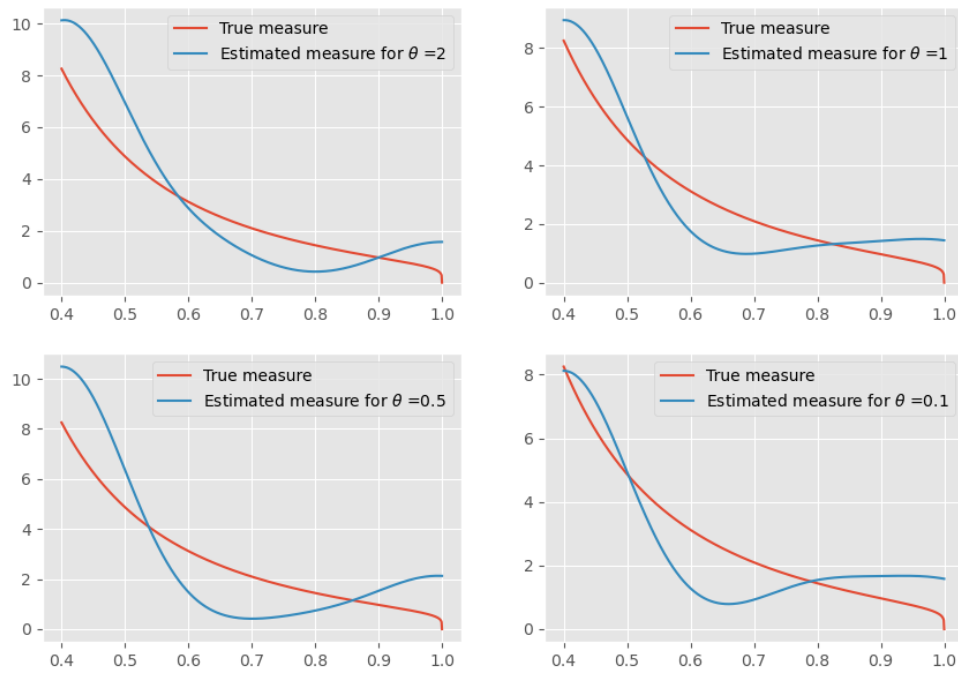


Figure S4: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 5$, $h = 0.1$, $\tau = 2$, $\alpha = 1.2$.

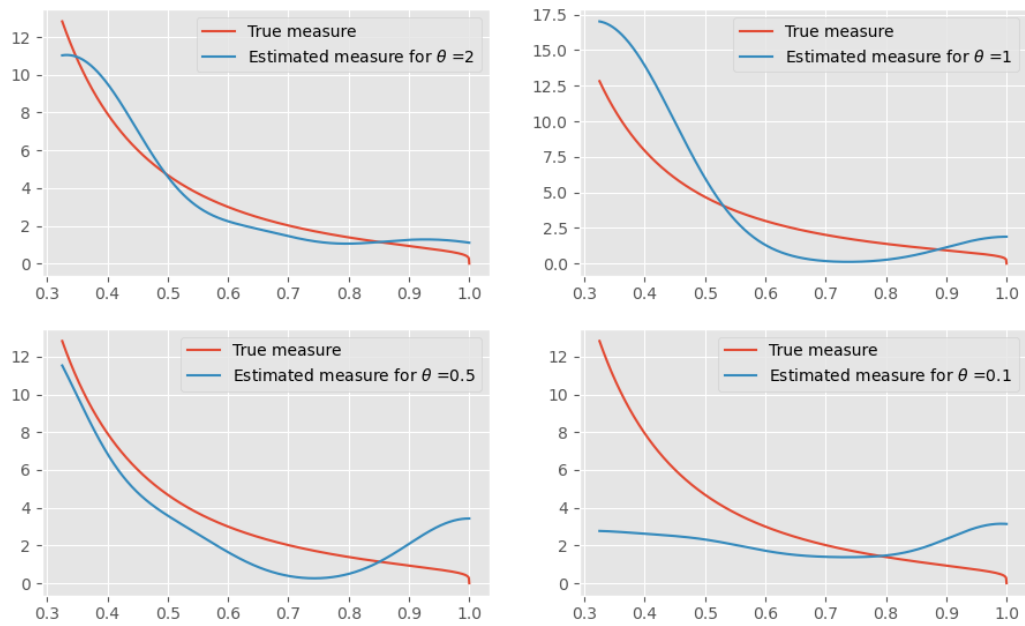


Figure S5: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 8$, $h = 0.1$, $\tau = 2$, $\alpha = 1.2$.

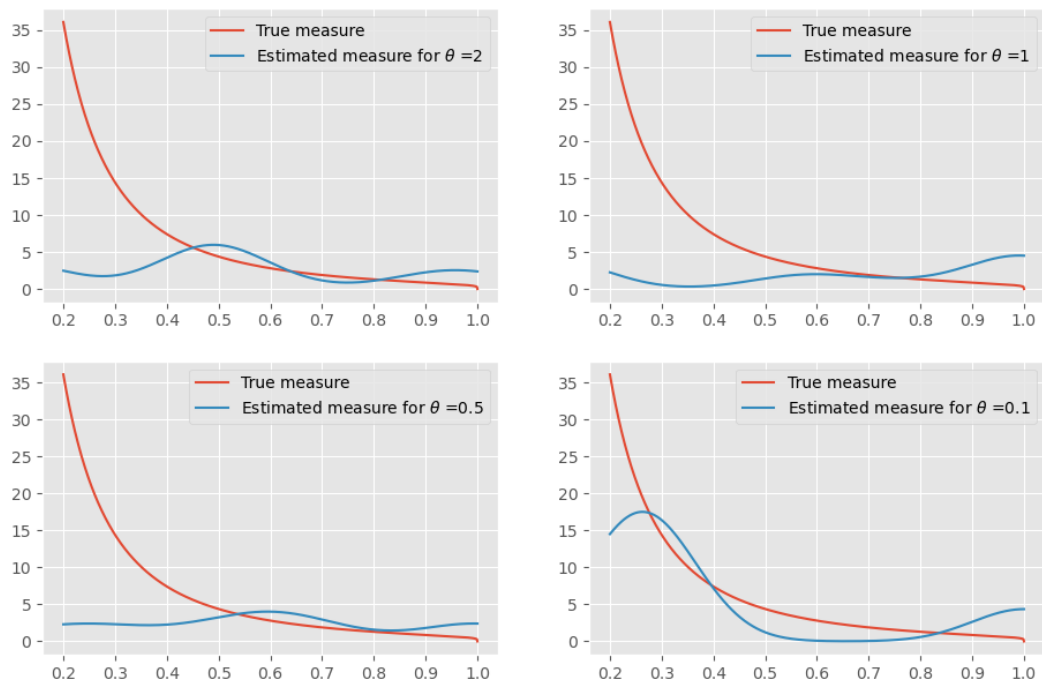


Figure S6: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 20$, $h = 0.1$, $\tau = 2$, $\alpha = 1.2$.

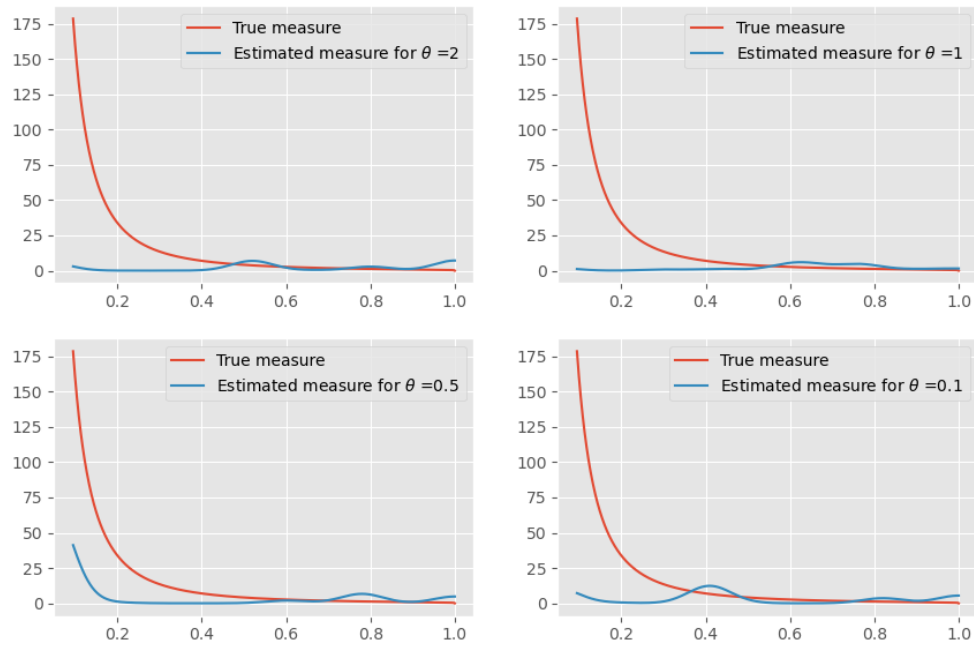


Figure S7: **Density of v .** We compare the true density (red curves) to the density estimated from simulated SNP matrices using equation (7) (blue curves). Here $n = 20$, $m = 50$, $h = 0.02$, $\tau = 2$, $\alpha = 1.2$.