

BAYESIAN INFERENCE IN GENETICS

Jaromir Sant, Paul Jenkins, Jere Koskela, Dario Spanò

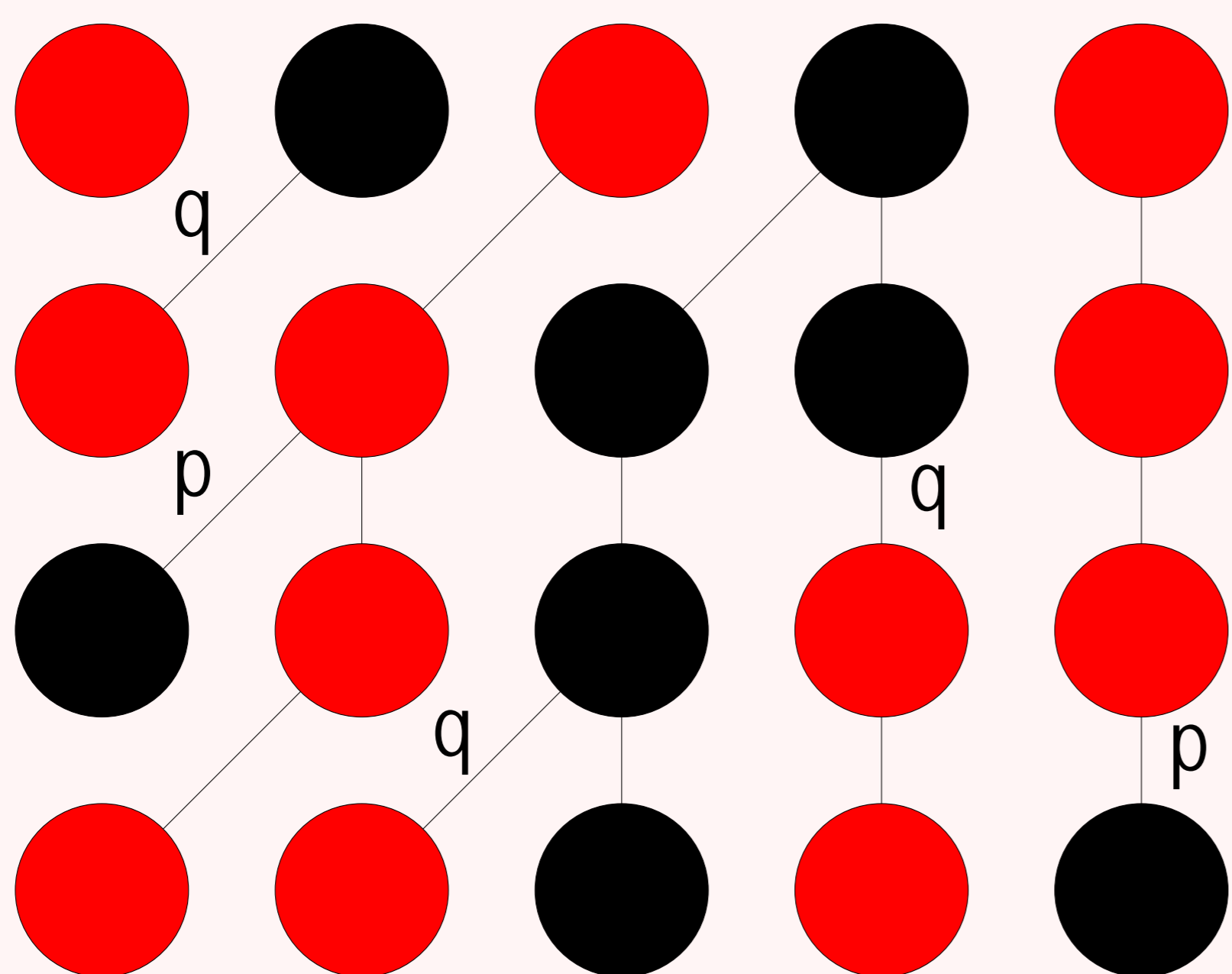
The Wright-Fisher Model

Consider a haploid population of fixed size N and consider only two alleles a and A . In the neutral Wright-Fisher model without mutation, parents are chosen uniformly at random and offspring inherit their type. To add mutation and selection:

Mutation - after choosing parent, flip a coin. If heads then the offspring mutates, otherwise it retains type of their parent

Selection - weight individuals by the relative fitness when choosing parents

Example of WF Model with Mutation and Selection



If we let Y_k^N be the number of individuals having allele A , and we assume that a and A have relative fitness $1 : 1 + \frac{s}{N}$, and mutation probabilities $\frac{1}{N} a \rightarrow A; \frac{1}{N} A \rightarrow a$, then

$$P[Y_{k+1}^N = j | Y_k^N = i] = \binom{N}{j} \frac{j_i (1 - i)^{N-j}}{i(1 + \frac{s}{N})(1 - \frac{1}{N} A \rightarrow a) + (N - i)\frac{1}{N} a \rightarrow A}$$

with

$$i := \frac{i(1 + \frac{s}{N})(1 - \frac{1}{N} A \rightarrow a) + (N - i)\frac{1}{N} a \rightarrow A}{i(1 + \frac{s}{N}) + N - i}$$

The Wright-Fisher Diffusion

Rescaling the Wright-Fisher model leads to the Wright-Fisher diffusion:

$$\frac{1}{N} Y_{[tN]}^N \rightarrow X_t$$

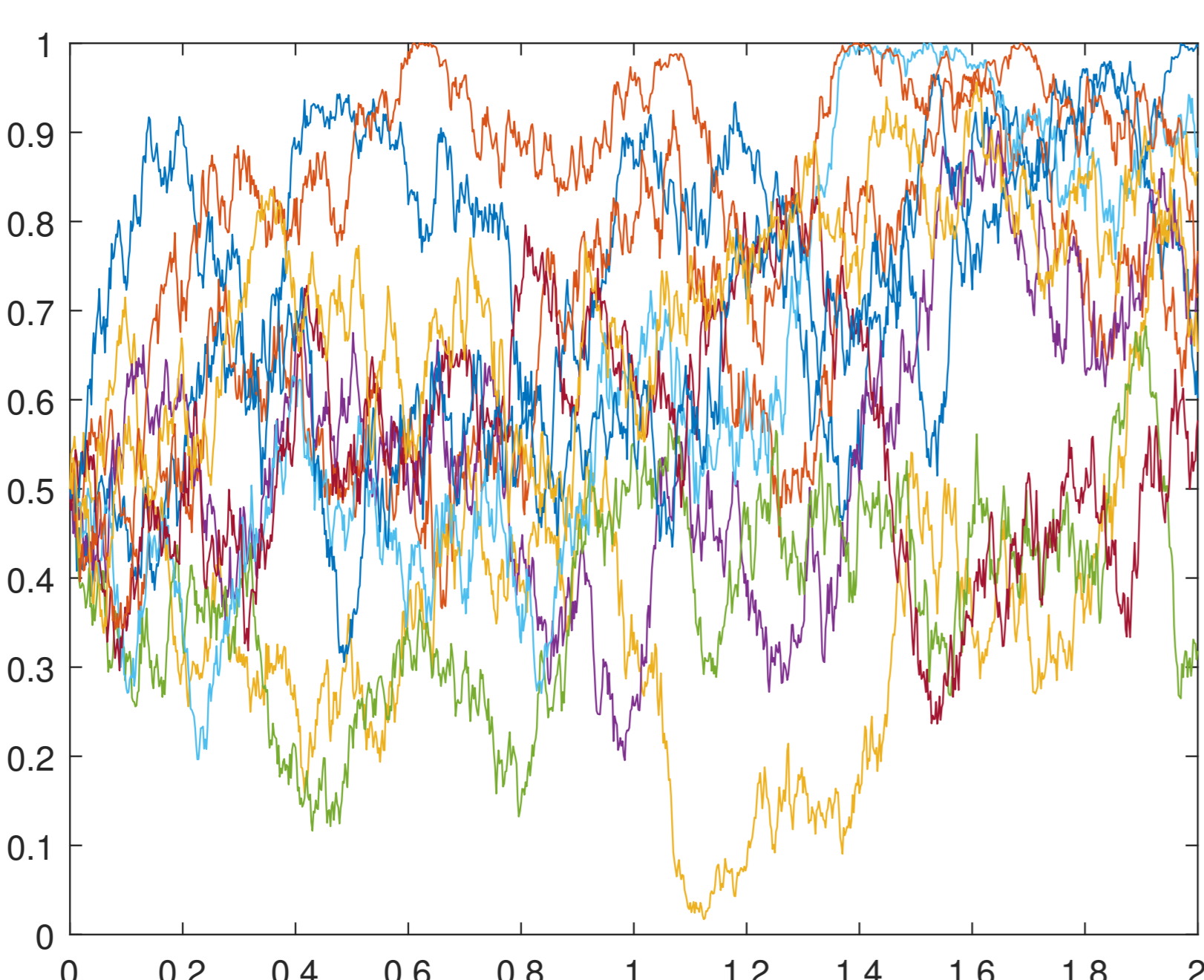
where the convergence is pathwise in $D_{[0,1]}([0; \infty))$, and X satisfies the SDE

$$dX_t = \frac{1}{2} (sX_t(1 - X_t) - a \rightarrow A X_t + A \rightarrow a (1 - X_t)) dt + \sqrt{X_t(1 - X_t)} dW_t$$

$s \in \mathcal{S}$ is the selection parameter we wish to infer and

$a \rightarrow A; A \rightarrow a > 0$ are the corresponding mutation parameters which we assume to be known.

Sample paths from a Wright-Fisher Diffusion with $s = 1, a \rightarrow A = 0.5; A \rightarrow a = 0.8$



The Inferential Setting

We observe one whole sample path $(X_t)_{t \in [0; T]}$ continuously through time, and consider the properties of the Bayesian estimator \tilde{s}_T for s in the asymptotic limit $T \rightarrow \infty$. We define

$$\tilde{s}_T = \arg \min_{s_T} \int_{\mathcal{S}} E_{P_s^{(T)}} \left[\ell \left(\sqrt{T} (\tilde{s}_T - s) \right) \right] p(s) ds$$

where p and ℓ are a suitably chosen prior and loss function respectively. The idea is to obtain bounds on what we can learn from the data in the absence of observational error, as done in [1] & [2]. The object of interest is the **likelihood ratio function**

$$Z_{T;s}(u) := \frac{dP_{s+\frac{u}{\sqrt{T}}}^{(T)}(X^T)}{dP_s^{(T)}(X^T)}$$

where we look at an order $\frac{1}{\sqrt{T}}$ perturbation around a fixed s .

Properties of the Bayesian Estimator

The Ibragimov-Has'minskii Conditions & Theorem

C1 : $\forall K \subset \Theta$ compact, $\exists a; B \in \mathbb{R}$ s.t. $\forall R > 0, \forall u_1; u_2$ s.t. $|u_1| < R, |u_2| < R$, and $q > 0$

$$\sup_{s \in K} E_{P_s^{(T)}} \left[\left| Z_{T;s}(u_2)^{\frac{1}{2}} - Z_{T;s}(u_1)^{\frac{1}{2}} \right|^2 \right] \leq B(1 + R^a) |u_2 - u_1|^q$$

C2 : $\forall K \subset \Theta$ compact, $\exists g_T(\cdot)$ a suitable monotonically increasing continuous function s.t.

$$\forall u \in U_{T;s} := \left\{ u : s + \frac{u}{\sqrt{T}} \in \Theta \right\}$$

$$\sup_{s \in K} E_{P_s^{(T)}} \left[Z_{T;s}(u)^{\frac{1}{2}} \right] \leq e^{-g_T(|u|)}$$

C3 : The random functions $Z_{T;s}(u)$ have marginal distributions which converge uniformly in $s \in K$ as $T \rightarrow \infty$ to those of the random function $Z_s(u)$

C4 : The random function

$$(\nu) = \int_{\mathbb{R}} \ell(\nu - u) \frac{Z_s(u)}{\int_{\mathbb{R}} Z_s(y) dy} du$$

attains its minimum value at the unique point $\tilde{u}(s) = \tilde{u}$ with probability 1

Theorem : If \tilde{s}_T is the Bayesian estimator and C1-C4 hold, then we have that \tilde{s}_T :

is uniformly consistent in $s \in K$

is uniformly asymptotically normal

displays moment convergence for any $p > 0$ uniformly on compacts $K \subset \Theta$

is asymptotically efficient for a suitable choice of loss function

References

- [1] Y. A. Kutoyants, *Statistical inference for ergodic diffusion processes*. Springer Series in Statistics. London, 2004
- [2] G. A. Watterson. Estimating and testing selection: the two-alleles, genetic selection diffusion model. *Adv. in Appl. Probab.*, 11(1):14-30, 1979.