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

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 $T \to \infty$. We define

$$ilde{s}_T = rg\min_{ar{s}_T} \int_{\mathcal{S}} \mathbb{E}_{\mathbb{P}^{(T)}_s} \left[\ell \left(\sqrt{T} \left(ar{s}_T - s
ight)
ight)
ight] p(s) ds$$

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**



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

$$\mathbb{P}\left[Y_{k+1}^N=jert Y_k^N=i
ight]=inom{N}{j}\psi_i^j(1-\psi_i)^{N-j}$$

with

$$_{N'} = rac{i(1+rac{s}{N})(1-rac{1}{N} heta_{A
ightarrow a}) + (N-i)rac{1}{N} heta_{a
ightarrow A}}{N}$$

$$Z_{T,s}(u):=rac{d\mathbb{P}^{(T)}_{s+rac{u}{\sqrt{T}}}}{d\mathbb{P}^{(T)}_s}(X^T)$$

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 $\sup_{s\in K} \mathbb{E}_{\mathbb{P}^{(T)}_s} \left| \left| Z_{T,s}(u_2)^{rac{1}{2}} - Z_{T,s}(u_1)^{rac{1}{2}}
ight|^2
ight|^2
ight|$ $\leq B(1+R^a)|u_2-u_1|^q$ **C2**: $\forall K \subset \Theta$ compact, $\exists g_T(\cdot)$ a suitable monotonically

increasing continuous function s.t.



$i(1 + \frac{s}{N}) + N - i$

The Wright-Fisher Diffusion

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

$$rac{1}{N}Y^N_{\lfloor tN
floor}
ightarrow X_t$$

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

$$egin{aligned} dX_t = &rac{1}{2} \left(sX_t(1-X_t) - heta_{a o A}X_t + heta_{A o a}(1-X_t)
ight) dt \ &+ \sqrt{X_t(1-X_t)} dW_t \end{aligned}$$

 $s \in \mathcal{S}$ is the selection parameter we wish to infer and $\theta_{a \to A}, \theta_{A \to a} > 0$ are the corresponding mutation parameters which we assume to be known.

Sample paths from a Wright-Fisher Diffusion with s=1, $heta_{a ightarrow A}=0.5, heta_{A ightarrow a}=0.8$

$$egin{aligned} orall u \in \mathbb{U}_{T,s} &:= \{u: s + rac{u}{\sqrt{T}} \in \Theta\} \ &\sup_{s \in K} \mathbb{E}_{\mathbb{P}^{(T)}_s} \left[Z_{T,s}(u)^rac{1}{2}
ight] \leq e^{-g_T(|u|)} \end{aligned}$$

- **C3**: The random functions $Z_{T,s}(u)$ have marginal distributions which converge uniformly in $s \in K$ as $T \to \infty$ to those of the random function $Z_s(u)$
- **C4** : The random function

$$\psi(v) = \int_{\mathbb{R}} \ell(v-u) rac{Z_s(u)}{\int_{\mathbb{R}} Z_s(y) dy} du$$

attains its minimum value at the unique point $ilde{u}(s) = ilde{u}$ 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.







Pioneering research and skills