



Matthias Birkner

Weierstrass Institute for Applied Analysis and Stochastics

Inference for coalescents with multiple collisions  
Ongoing project with J. Blath and M. Steinrücken, TU Berlin

WIAS Institutskolloquium, 17 December 2007

## Outline

---

1. Some 'classical' mathematical population genetics
2. Coalescents with multiple collisions
3. Combinatorics of the infinitely-many sites mutation model
4. A Monte Carlo method for likelihood estimation
5. Illustration

Genetic variability at a 250bp piece of the mitochondrial cytochrome *b*-gene in a sample of 117 atlantic cod (a random subsample from the dataset described in E. Árnason, *Genetics* 2004)

	468	481	487	488	490	496	508	523	562	601	631	643	649	685	691
66	t	a	a	c	a	a	t	g	a	t	g	a	c	c	g
17	-	-	-	-	-	-	c	-	-	-	-	-	-	-	-
14	-	-	-	-	-	-	-	a	-	-	-	-	-	t	-
8	-	-	-	-	-	-	-	-	-	-	-	-	-	-	t
1	-	-	-	-	-	-	-	-	-	-	-	-	-	t	-
2	-	-	-	t	-	-	-	-	-	-	-	-	-	-	-
1	-	-	-	-	-	-	-	a	-	-	-	g	-	t	-
1	-	-	-	-	-	-	-	-	-	-	-	-	t	-	-
1	-	-	-	-	g	-	c	-	-	-	-	-	-	-	-
1	-	-	-	-	-	g	-	-	-	-	-	-	-	-	-
1	-	-	g	-	-	-	-	-	-	-	a	-	-	-	t
1	-	-	-	-	-	-	c	-	g	-	-	-	-	-	-
1	g	-	-	-	-	-	-	-	-	-	-	-	-	-	-
1	-	-	-	-	-	-	-	-	-	c	-	-	-	-	-
1	-	c	-	-	-	-	c	-	-	-	-	-	-	-	-

John Gillespie's 'Great obsession' of population genetics:

“What evolutionary forces could have lead to such divergence between individuals in the same species?”

John Gillespie's 'Great obsession' of population genetics:

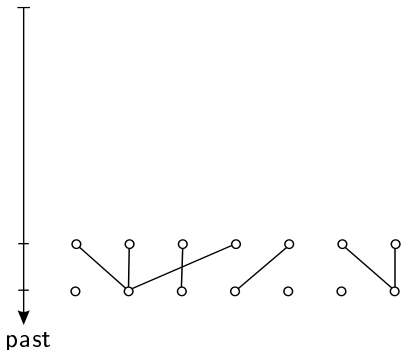
“What evolutionary forces could have lead to such divergence between individuals in the same species?”

A more humble obsession:

How can stochastic models help to understand genetic variability inside populations?

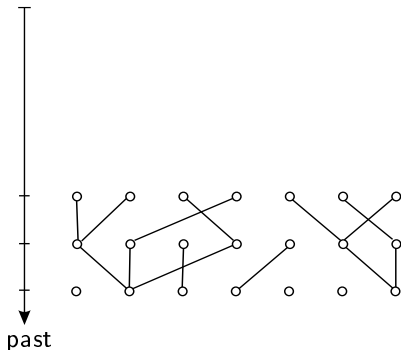
## Wright-Fisher model: A fundamental model for 'genetic drift'

- ▷ A (haploid) population of  $N$  individuals per generation,
- ▷ each individual in the present generation picks a 'parent' at random from the previous generation,
- ▷ genetic types are inherited (possibly with a small probability of mutation).



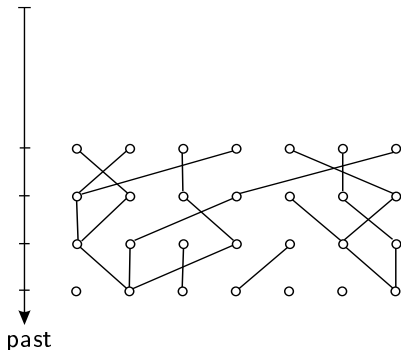
## Wright-Fisher model: A fundamental model for 'genetic drift'

- ▷ A (haploid) population of  $N$  individuals per generation,
- ▷ each individual in the present generation picks a 'parent' at random from the previous generation,
- ▷ genetic types are inherited (possibly with a small probability of mutation).



## Wright-Fisher model: A fundamental model for 'genetic drift'

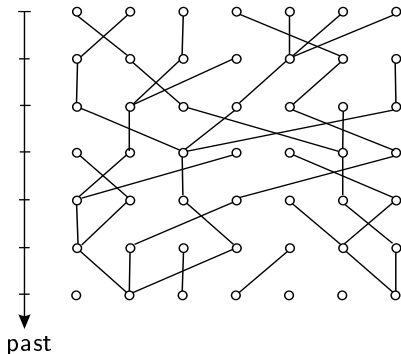
- ▷ A (haploid) population of  $N$  individuals per generation,
- ▷ each individual in the present generation picks a 'parent' at random from the previous generation,
- ▷ genetic types are inherited (possibly with a small probability of mutation).





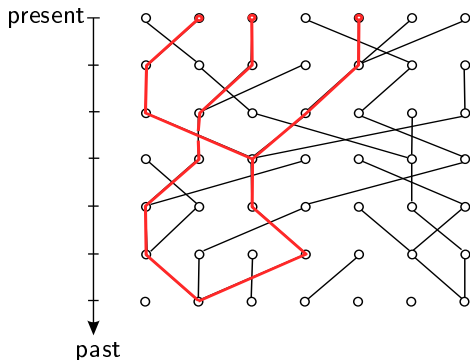
## Wright-Fisher model: A fundamental model for 'genetic drift'

- ▷ A (haploid) population of  $N$  individuals per generation,
- ▷ each individual in the present generation picks a 'parent' at random from the previous generation,
- ▷ genetic types are inherited (possibly with a small probability of mutation).



## Genealogical point of view

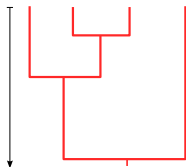
Sample  $n (\ll N)$  individuals from the 'present generation'



## Kingman's coalescent

### Theorem (Kingman, 1982)

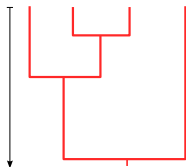
In the limit  $N \rightarrow \infty$ , the genealogy of an  $n$ -sample, measured in units of  $N$  generations, is described by a continuous-time Markov chain where each pair of lineages merges at rate 1.



## Kingman's coalescent

### Theorem (Kingman, 1982)

In the limit  $N \rightarrow \infty$ , the genealogy of an  $n$ -sample, measured in units of  $N$  generations, is described by a continuous-time Markov chain where each pair of lineages merges at rate 1.



The same limit appears for any *exchangeable* offspring vectors

$$(\nu_1, \dots, \nu_N), \quad (\text{independent over generations}),$$

if time is measured in  $\frac{N}{\sigma^2}$  generations, where  $\sigma^2 = \lim_{N \rightarrow \infty} \text{Var}(\nu_1)$ .

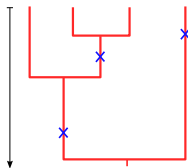
## Kingman's coalescent: superimposing neutral types

Assume that the considered genetic types do not affect their bearer's reproductive success.

If as population size  $N \rightarrow \infty$ ,

$\frac{N}{\sigma^2} \times$  mutation prob. per ind. per generation  $\rightarrow r$ ,

the type configuration in the sample can be described by putting mutations with rate  $r$  along the genealogy.



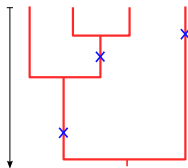
## Kingman's coalescent: superimposing neutral types

Assume that the considered genetic types do not affect their bearer's reproductive success.

If as population size  $N \rightarrow \infty$ ,

$\frac{N}{\sigma^2} \times$  mutation prob. per ind. per generation  $\rightarrow r$ ,

the type configuration in the sample can be described by putting mutations with rate  $r$  along the genealogy.



Kingman's coalescent is the *standard model* of mathematical population genetics.

## Coalescents with multiple collisions, aka ' $\Lambda$ -coalescents'

While  $n$  lineages, any  $k$  coalesce at rate

$$\lambda_{n,k} = \int_{[0,1]} x^{k-2} (1-x)^{n-k} \Lambda(dx),$$

where  $\Lambda$  is a finite measure on  $[0, 1]$ .

(Sagitov, 1999; Pitman, 1999).



## Coalescents with multiple collisions, aka ' $\Lambda$ -coalescents'

While  $n$  lineages, any  $k$  coalesce at rate

$$\lambda_{n,k} = \int_{[0,1]} x^{k-2} (1-x)^{n-k} \Lambda(dx), \text{ where } \Lambda \text{ is a finite measure on } [0,1].$$

(Sagitov, 1999; Pitman, 1999).

Interpretation:

re-write  $\lambda_{n,k} = \int_{[0,1]} x^k (1-x)^{n-k} \frac{1}{x^2} \Lambda(dx)$  to see:

at rate  $\frac{1}{x^2} \Lambda([x, x+dx])$ , an ' $x$ -resampling event' occurs.

Thinking forwards in time, this corresponds to an event in which the fraction  $x$  of the total population is replaced by the offspring of a single individual.





## Coalescents with multiple collisions, aka ' $\Lambda$ -coalescents'

While  $n$  lineages, any  $k$  coalesce at rate

$$\lambda_{n,k} = \int_{[0,1]} x^{k-2} (1-x)^{n-k} \Lambda(dx), \text{ where } \Lambda \text{ is a finite measure on } [0, 1].$$

(Sagitov, 1999; Pitman, 1999).

Interpretation:

re-write  $\lambda_{n,k} = \int_{[0,1]} x^k (1-x)^{n-k} \frac{1}{x^2} \Lambda(dx)$  to see:

at rate  $\frac{1}{x^2} \Lambda([x, x+dx])$ , an ' $x$ -resampling event' occurs.

Thinking forwards in time, this corresponds to an event in which the fraction  $x$  of the total population is replaced by the offspring of a single individual.

Note:  $\Lambda = \delta_0$  corresponds to Kingman's coalescent.



## Cannings' models in the 'domain of attraction of a $\Lambda$ -coalescent'

Fixed population size  $N$ , *exchangeable* offspring numbers in one generation

$$(\nu_1, \nu_2, \dots, \nu_N).$$

Sagitov (1999), Möhle & Sagitov (2001) clarify under which conditions the genealogies of a sequence of exchangeable finite population models are described by a  $\Lambda$ -coalescent:

- ▷  $c_N :=$  pair coalescence probability over one generation  $\rightarrow 0$   
(  $c_N = \frac{1}{N-1} \mathbb{E}[\nu_1(\nu_1 - 1)]$  )
- ▷ two double mergers asymptotically negligible compared to one triple merger
- ▷  $Nc_N \mathbb{P}(\text{a given family has size} \geq Nx) \sim \int_x^1 y^{-2} \Lambda(dy)$

Time is measured in  $1/c_N$  generations (in general  $\neq 1/\text{pop. size}$ )

## A 'heavy-tailed' Cannings model

---

Haploid population of size  $N$ . Individual  $i$  has  $X_i$  *potential offspring*,  
 $X_1, X_2, \dots, X_N$  are i.i.d. with mean  $m := \mathbb{E}[X_1] > 1$ ,  
 $\mathbb{P}(X_1 \geq k) \sim \text{Const.} \times k^{-\alpha}$  with  $\alpha \in (1, 2)$ .

Note: infinite variance.

Sample  $N$  without replacement from all potential offspring to form the next generation.

## A 'heavy-tailed' Cannings model

Haploid population of size  $N$ . Individual  $i$  has  $X_i$  potential offspring,  $X_1, X_2, \dots, X_N$  are i.i.d. with mean  $m := \mathbb{E}[X_1] > 1$ ,  $\mathbb{P}(X_1 \geq k) \sim \text{Const.} \times k^{-\alpha}$  with  $\alpha \in (1, 2)$ .

Note: infinite variance.

Sample  $N$  without replacement from all potential offspring to form the next generation.

### Theorem (Schweinsberg, 2003)

Let  $c_N = \text{prob. of pair coalescence one generation back in } N\text{-th model}$ .  $c_N \sim \text{const. } N^{1-\alpha}$ , measured in units of  $1/c_N$  generations, the genealogy of a sample from the  $N$ -th model is approximately described by a  $\Lambda$ -coalescent with  $\Lambda = \text{Beta}(2 - \alpha, \alpha)$ .

$$\left( \text{Beta}(2 - \alpha, \alpha)(dx) = \mathbf{1}_{[0,1]}(x) \frac{1}{\Gamma(2-\alpha)\Gamma(\alpha)} x^{1-\alpha} (1-x)^{\alpha-1} dx \right)$$

## Why $\Lambda = \text{Beta}(2 - \alpha, \alpha)$ ?

---

Heuristic argument:

Probability that first individual's offspring provides more than fraction  $y$  of the next generation, given that the family is substantial (i.e. given  $X_1 \geq \varepsilon N$ )

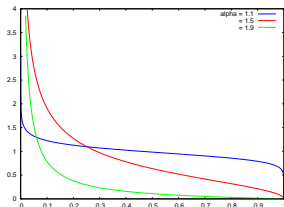
## Why $\Lambda = \text{Beta}(2 - \alpha, \alpha)$ ?

Heuristic argument:

Probability that first individual's offspring provides more than fraction  $y$  of the next generation, given that the family is substantial (i.e. given  $X_1 \geq \varepsilon N$ )

$$\begin{aligned}
 &\approx \mathbb{P}\left(\frac{X_1}{X_1 + (N-1)m} \geq y \mid X_1 \geq \varepsilon N\right) \\
 &= \mathbb{P}\left(X_1 \geq (N-1)m \frac{y}{1-y} \mid X_1 \geq \varepsilon N\right) \\
 &\sim \text{const.} \frac{(1-y)^\alpha}{y^\alpha} = \text{const.}' \text{Beta}(2-\alpha, \alpha)([y, 1]).
 \end{aligned}$$

## The family $\text{Beta}(2 - \alpha, \alpha)$ , $\alpha \in (1, 2]$



- ▷ Kingman's coalescent is included as a boundary case:  
 $\text{Beta}(2 - \alpha, \alpha) \rightarrow \delta_0$  weakly as  $\alpha \rightarrow 2$ .
- ▷ Smaller  $\alpha$  means tendency towards more extreme resampling events.
- ▷ For  $\alpha \leq 1$ , corresponding coalescents *do not* come down from infinity.
- ▷  $\text{Beta}(2 - \alpha, \alpha)$ -coalescents appear as genealogies of  $\alpha$ -stable continuous mass branching process (via a time-change).

## 'Meta-mathematic' associations





## Asymptotics of the frequency spectrum

Consider an  $n$ -Beta( $2 - \alpha, \alpha$ )-coalescent, mutations at rate  $r$  according to the *infinitely-many-sites* model (assuming known ancestral types). Let

$M(n) :=$  #total number of mutations in the sample,

$M_k(n) :=$  #number of mutations affecting exactly  $k$  samples,

$k = 1, 2, \dots, n - 1$ .

**Theorem** (Berestycki, Berestycki & Schweinsberg, 2005–)

$$\frac{M(n)}{n^{2-\alpha}} \rightarrow r \frac{\alpha(\alpha-1)\Gamma(\alpha)}{2-\alpha}, \quad \frac{M_k(n)}{n^{2-\alpha}} \rightarrow r\alpha(\alpha-1)^2 \frac{\Gamma(k+\alpha-2)}{k!}$$

in probability as  $n \rightarrow \infty$ .

## Asymptotics of the frequency spectrum

Consider an  $n$ -Beta( $2 - \alpha, \alpha$ )-coalescent, mutations at rate  $r$  according to the *infinitely-many-sites* model (assuming known ancestral types). Let

$$M(n) := \# \text{total number of mutations in the sample,}$$

$$M_k(n) := \# \text{number of mutations affecting exactly } k \text{ samples,}$$

$$k = 1, 2, \dots, n - 1.$$

**Theorem** (Berestycki, Berestycki & Schweinsberg, 2005–)

$$\frac{M(n)}{n^{2-\alpha}} \rightarrow r \frac{\alpha(\alpha-1)\Gamma(\alpha)}{2-\alpha}, \quad \frac{M_k(n)}{n^{2-\alpha}} \rightarrow r\alpha(\alpha-1)^2 \frac{\Gamma(k+\alpha-2)}{k!}$$

in probability as  $n \rightarrow \infty$ .

Thus  $M_1(n)/M(n) \approx 2 - \alpha$  for  $n$  large, which suggests

$$\hat{\alpha}_{\text{BBS}} := 2 - \frac{M_1(n)}{M(n)}$$

as an estimator for  $\alpha$ .

## A likelihood approach

---

If the observations had been generated by putting mutations at rate  $r > 0$  on a realisation of a certain  $\Lambda$ -coalescent (from some class, e.g.,  $\text{Beta}(2 - \alpha, \alpha)$ ), for which  $(\hat{\Lambda}, \hat{r})$  is

$\mathbb{P}_{\Lambda, r}(\text{observations})$  maximal?

## Infinitely-many-sites model

An infinite sequence of completely linked sites, mutations always hit a new site

Example:

Seq.	segr. site			
	1	2	3	4
1	1	0	0	0
2	1	1	0	0
3	0	0	1	1
4	0	0	1	1
5	0	0	1	0

(0=wild type, 1=mutant

assume known ancestral types)

Obs. fit IMS  $\iff$  no sub-matrix  $\begin{matrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{matrix}$  (nor row permutation).

## Infinitely-many-sites model, II

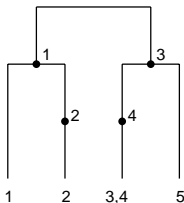
If the infinitely-many-sites model applies, the observations correspond to a unique rooted perfect phylogeny (or 'genetree').

Sequences,

Genetree,

obs. types

Seq.	segr. site			
	1	2	3	4
1	1	0	0	0
2	1	1	0	0
3	0	1	1	
4	0	0	1	1
5	0	0	1	0



type	multiplicity
(1, 0)	1
(2, 1, 0)	1
(4, 3, 0)	2
(3, 0)	1

Construct e.g. using Gusfield's (1991) algorithm.

Note: purely combinatorial, does not depend on a probabilistic model for the observations.

## “Naive approach”

---

We have

$$p_{\Lambda,r}(T, \mathbf{n}) = \sum_{T \in \mathcal{C}_{T,\mathbf{n}}} \mathbb{P}_{\Lambda,r}(\text{marked geneal. tree of } n\text{-sample} = T),$$

where  $\mathcal{C}_{T,\mathbf{n}}$  are all marked coalescent trees compatible with the observations.

Problem: Too many trees!

## Recursions for tree probabilities (B. & Blath, 2007)

$T$  a tree of (ordered) types, type multiplicity vector  $\mathbf{n}$ .

$$\begin{aligned}
 p_{\Lambda,r}(T, \mathbf{n}) &= \frac{1}{r_n} \sum_{i: n_i \geq 2} \sum_{k=2}^{n_i} \binom{n_i}{k} \lambda_{n,k} \frac{n_i - k + 1}{n - k + 1} p_{\Lambda,r}(T, \mathbf{n} - (k-1)\mathbf{e}_i) \\
 &+ \frac{r}{r_n} \sum_{\substack{k: n_k=1, x_{k0} \text{ distinct} \\ s(x_k) \neq x_j \forall j}} p_{\Lambda,r}(s_k(T), \mathbf{n}) \\
 &+ \frac{r}{r_n} \sum_{\substack{k: n_k=1, x_{k0} \\ \text{distinct}}} \sum_{j: s(x_k)=x_j} (n_j + 1) p_{\Lambda,r}(r_k(T), r_k(\mathbf{n} + \mathbf{e}_j)).
 \end{aligned}$$

where  $\mathbf{e}_j$ :  $j$ -th unit vector,  $s_k(T)$ : removes first coordinate of  $k$ -th sequence in  $\mathbf{n}$ ,  $r_k(T)$ : removes  $k$ -th sequence from  $T$ ,  $x_{k0}$  'distinct':  $\iff$

$x_{k0} \neq x_{ij}, \forall (i, j) \neq (k, 0)$ ,  $r_n = rn + \sum_{k=2}^n \binom{n}{k} \lambda_{n,k}$ .

Extends Ethier & Griffiths (1987) to  $\Lambda$ -case and Möhle (2005) to IMS.

Note: true recursion in *complexity* of  $(T, \mathbf{n})$ .

## Markov chains and linear equations

---

$|S| < \infty$ ,  $(q_{xy})$  transition kernel  $S$ ,  $f : S \rightarrow \mathbb{R}$ .

$$u(x) = f(x) \sum_{y \in S} q_{xy} u(y), \quad x \in S' \subset S$$

with given boundary values on  $S \setminus S'$ .

$X$  a  $q$ -Markov chain,  $\tau := \min\{k : X_k \notin S'\}$ .

If  $\tau \leq K$  for a fixed  $K < \infty$ ,

$$u(x) = \mathbb{E}_x \left[ \prod_{i=0}^{\tau} f(X_i) \right].$$



## A Monte-Carlo method

Using this and the recursion for  $p_{\Lambda,r}$ :

$$p_{\Lambda,r}(T, \mathbf{n}) = \mathbb{E}_{(T, \mathbf{n})} \left[ \prod_{i=0}^{\tau} f(X_i) \right]$$

for a suitable Markov chain on type trees with multiplicities (analogous to Griffiths & Tavaré, 1994).

- ▶ Unbiased estimate  $\hat{p}_{\Lambda,r}(T, \mathbf{n})$  via independent runs
- ▶ Finite runtime: complexity of  $(T, \mathbf{n})$  ( $:= \# \text{mutations} + \text{sample size}$ ) decreases in each step
- ▶ Can view chain as an integral on “ $(\Lambda)$ -coalescent histories”

## The Monte-Carlo method (details to be glossed over)

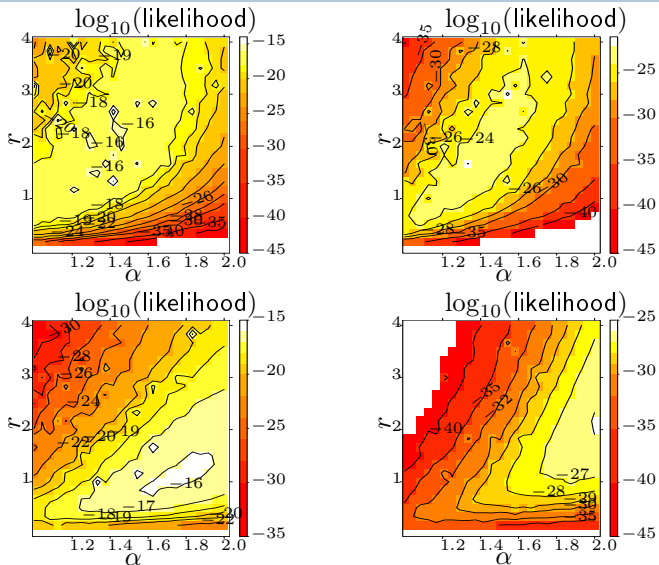
Transition mechanism

$$(T, \mathbf{n}) \rightarrow \begin{cases} (s_k(T), \mathbf{n}) & \text{w. p. } \frac{1}{r_n f(T, \mathbf{n})} r \text{ if} \\ & n_k = 1, x_{k0} \text{ distinct, } s(x_k) \neq x_j \forall j, \\ (r_k(T), r_k(\mathbf{n} + \mathbf{e}_j)) & \text{w. p. } \frac{1}{r_n f(T, \mathbf{n})} r(n_j + 1) \text{ if} \\ & n_k = 1, x_{k0} \text{ distinct, } s(x_k) = x_j, \\ (T, \mathbf{n} - (k-1)\mathbf{e}_i) & \text{w. p. } \frac{1}{r_n f(T, \mathbf{n})} \binom{n}{k} \lambda_{n,k} \frac{n_i - k + 1}{n - k + 1} \text{ if } 2 \leq k \leq n_i, \end{cases}$$

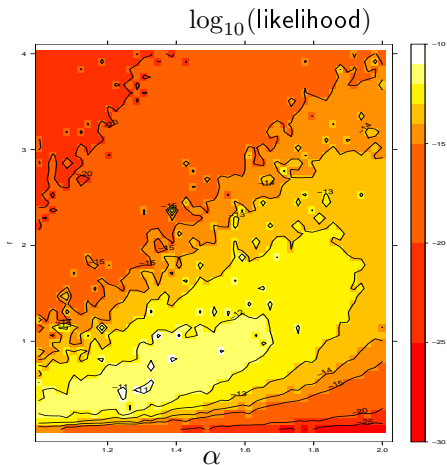
where

$$\begin{aligned} r_n f(T, \mathbf{n}) = & \sum_{\substack{k: n_k=1, x_{k0} \text{ distinct} \\ s_k(x_k) \neq x_j \forall j}} r + \sum_{\substack{k: n_k=1, x_{k0} \\ \text{distinct}}} \sum_{j: s_k(x_k) = x_j} r(n_j + 1) \\ & + \sum_{1 \leq i \leq d: n_i \geq 2} \sum_{k=2}^{n_i} \binom{n}{k} \lambda_{n,k} \frac{n_i - k + 1}{n - k + 1}. \end{aligned}$$

## Simulated datasets: $\alpha = 1.25, 1.5, 1.75, 2, r = 2.0$



# $\hat{P}_{\text{Beta}(2-\alpha, \alpha), r}(\text{data})$ for the cod sample



Maximum at  $\hat{\alpha} = 1.3$ ,  $\hat{r} = 0.7$ .  $\hat{\alpha}_{\text{BBS}} = 2 - 9/14 \approx 1.36$ .

## Further issues

---

- ▷ Reduce variance of estimator via importance sampling?
- ▷ Properties of estimators?
- ▷ Interplay of demographic stochasticity and recombination, “ $\Lambda$ -ancestral recombination graph”?
- ▷ More general mutation models, unknown ancestral types
- ▷ Selection, population substructure
- ▷ ...

---

beta genetree is available (under GNU General public licence) from

<http://www.wias-berlin.de/people/birkner/bgt/>