Statistical Inferences in Population Genetics

Rasmus Nielsen 6/7/2022

Inferences of

- Demography (migration rates, divergence times, changes in population size, admixture proportions)
- Selection (identifying loci under positive or negative selection, estimating selection coefficients, time of selection)

Genetic Drift

1

Change of allele frequencies through time.



Coalescence theory



Past

The *n*-coalescent



The *n*-coalescent

For a sample of size n, the probability that all gene copies have distinct parental gene copies in the previous generation is

$$\left(1-\frac{1}{2N}\right)\left(1-\frac{2}{2N}\right)\times...\times\left(1-\frac{n-1}{2N}\right)=1-\frac{\binom{n}{2}}{2N}+O(N^{-2})$$

P(n gene copies do not find a common ancestor in rgenerations) is then approximately

$$\left[1 - \binom{n}{2} (2N)^{-1} + O(N^{-2})\right]^{r}$$

Set r = 2Nt and let $N \rightarrow \infty$

$$\left[1-\binom{n}{2}(2N)^{-1}+O(N^{-2})\right]^r\approx e^{-\binom{n}{2}t}$$

In general, the time in which there a *j* ancestors in the sample is exponentially distributed with mean

 $\binom{j}{2}^{-1}$



- (1) Time is scaled in terms of 2N. The smaller the population the faster the coalescence process.
- (2) The coalescence process is faster when there are many lineages. The further in the past you get in the tree – the slower the coalescence process. Branches deep in the tree will be longer



Sequence comparisons



Segregating sites with Single Nucleotide Polymorphisms (SNPs)

Infinite Sites Model



Assume there are 2N individuals (we model a haploid population with N individuals as a diploid population with 2N individuals).

Probability of having same parent in one generation = 1/2N

Mean number of generations until they have the same parent (time to the most recent common ancestor, tMRCA) = 2N

Mean number of mutations separating two individuals = $2N \times 2 \times \mu = 4N\mu = \theta$

Tajima's estimator of $4Nu = \theta$



$$\pi = (1 + 1 + 2)/3 = 4/3$$

 π (= $\hat{\theta}_T$) is an *estimator* of $4Nu = \theta$ (In this case the total value for the entire region)

Coalescence Trees



Felsenstein's Equation

$$p(X \mid \Theta) = \int_{G \in \psi} p(X \mid G) p(G \mid \Theta) dG$$

•Can only be evaluated directly in very simple cases.

•Simulation based approaches are typically used for real data.

MCMC

Set up a Markov chain on state space on all supported values of Θ and *G* and with stationary distribution $p(\Theta, G | X)$. Now since

$p(\Theta, G | X) \propto \Pr(X | G) p(G | \Theta) p(\Theta)$

this can easily be done using Metropolis-Hastings sampling, i.e. updates to Θ and *G* are proposed from a proposal distribution $q(\Theta, G \to \Theta', G')$ and accepted with probability

 $\frac{P(X \mid G')P(G' \mid \Theta')q(\Theta', G' \rightarrow \Theta, G)}{P(X \mid G)P(G \mid \Theta)q(\Theta, G \rightarrow \Theta', G')}$

MCMC algorithms





Demographic history of eastern pacific stickleback vs. western pacific stickleback.

Data: DNA sequences

Eastern	1	CGAAAAGTATCCATCTCGCAGTGCTGAGCTAGACA
Eastern	2	CGAAAAGTATCCATCTCGCAGTGCTGAGCTAGACA
Eastern	3	CGAAAAGTATCCATCTCGCAGTGCTGAGCTAGACA
Eastern	4	CGAAAGGTATCCATCTCGCAGTGCTGAGTTAGACA
Eastern	5	CAAAAAGTATCCATCTCGCAGTGCTGAGTTAGACA
Eastern	6	CAAAAAGTATCCATCTCGCAGTGCTGAGTTAGACA
Eastern	7	CGAAAAGTATCCATCTCGCAGTGCTGAACTAGACA
Eastern	8	CGAAAAGTATCCATCTCGCAGTGCTAAGCTAGACA
Eastern	9	TGAAAAGTATCCATCTCGCAGTGCTAAGCTAGACA
Western	1	CGAAAAGTATCCATCTCGCAGTGCTGAGCTAGACA
Western	2	CGCGAAGCACTTGCCCCATAGCGCTAAGCCGCGTT
Western	3	CGCGTAGCACTTGCCCCATAGCGCTAAGCCGCGTT
Western	4	CGCAAAGCGCTTGCCCCATAACGCTAAGCCGCGTT
Western	5	CGCAAAGCGCTTGCCCCATAACGCTAAGCCGCGTT
Western	6	CGCAAAGCACTTGCCCCATAACGCTAAGCCGCGTT
Western	7	CGCAAAGCACTTGCCCCATAACGCTAAGCCGCGTT





FIGURE 8.—The integrated likelihood surfaces for M_1 (dots) and M_2 (solid lines) estimated from the data by ORTI *et al.* 1994.







Some MCMC approaches

- BEAST (Rambaut)
- Various version of IM (Hey)
- MIGRATE (Beerli)

Approximate Bayesian Computation

- 1. Draw $\theta_i \sim \pi(\theta)$.
- 2. Simulate $x_i \sim p(x | \theta_i)$.
- 3. Reject θ_i if $\rho(S(x_i), S(y)) > \epsilon$.



Figure 3 | Best-supported demographic models inferred by approximate Bayesian computation model-selection. a, Eurasia; b, North America. Grey

Then this happened...



Problem 1: Each position in the genome might have a unique tree – there are millions of trees!







Challenge

- Full likelihood inference considered intractable in samples of n>2 (but see PSMC).
- Instead, rely on summary statistics/various features extracted from the data.
- Examples: π , *S*, *F*_{*ST*}, Patterson's *D*
- More complex examples: The Site-Frequency-Spectrum, principle components, features extracted by a convolutional neural network.

Principle Component Analysis (PCA)



Structure analysis of 1056 individuals from 52 populations for 377 microsatelite loci. Rosenberg et al. (2002)



Frequency spectrum



Africans Site Frequency Spectrum (SFS)



Composite likelihood function

$$L(\Theta) \equiv \prod_{j \in \Phi} \left(p_j(\Theta) \right)^{n_j}$$

Sampling probability of a SNP with site pattern *j* (a binary vector)

Number of SNPs with pattern *j* in the data

SNPs within a gene are correlated. But estimator is consistent. The estimate has the same properties as a real likelihood estimator except that it converges slightly slower because of the correlation (Wiuf 2006).

Marginal likelihood for a single site Nielsen (2000)

$$\begin{split} L(\Omega|X_i) &= \lim_{\theta \to 0} \Pr\left(X_i \mid \Omega, \theta, S_i > 0\right) \\ &= \lim_{\theta \to 0} \frac{\int (\theta/2)^{-1} \sum_{j:b_{ij} \in B_i} (1 - e^{-\theta T_{ij}/2}) e^{-\theta(T_i - T_{ij})/2} dF(G|\Omega)}{\int (\theta/2)^{-1} (1 - e^{-\theta T_i/2}) dF(G|\Omega)} \\ &= \frac{\int t_i dF(G|\Omega)}{\int T_i dF(G|\Omega)} = \frac{\mathrm{E}(t_i \mid \Omega)}{\mathrm{E}(T_i \mid \Omega)}. \end{split}$$
Estimation

(Nielsen 2000)

$$p_j(\Theta) = \frac{E[t]}{E[T]}$$



$$x = \{3, 2\}$$

$$T = 5\tau_5 + 4\tau_4 + 3\tau_3 + 2\tau_2.$$

$$t = \tau_4 + \tau_3 + 2\tau_2.$$

Sampling distribution can be calculated analytically or using simple simulation schemes.

Question

Consider the polarized SFS. Under the standard (Kingman's coalescent) with *n*=3, what is the probability of the pattern {1, 2} meaning one ancestral and two derived alleles?

Hint: the expected time in the tree with *i* lineages is 2/(i(i - 1))

Question

Consider the polarized SFS. Under the standard (Kingman's coalescent) with n=3, what is the probability of the pattern $\{1, 2\}$ meaning one ancestral and two derived alleles?



Expected total tree length:

 $3 \times \left(\frac{1}{3}\right) + 2 \times 1 = 3$ Expected length of blue edge: 1 P(x = {1,2}) = 1/3

Site Frequency Spectrum



Data

Directly sequenced polymorphism data from 20 European-Americans, 19 African-Americans and one chimpanzee from 9,316 protein coding genes (Bustamante et al. 2005).

Objectives

To detect natural selection in individual genes using the frequency spectrum.

To account for demography by estimating parameters of a demographic model.

Demographic model





African-Americans



European-Americans





Alternative...

(e.g. Williamson et al. 2005)

Use
$$\Pr(X | \Theta) = \int \binom{n}{x} p^x (1-p)^{n-x} f(p | \Theta) dp$$

The density $f(p \mid \Theta)$ can be obtained by numerical solution of diffusion equations (e.g. Crank–Nicolson approximation).

SFS based infrences

- Simulation (e.g., Fastsimcoal2)
- Numerical solution of diffusion equation (e.g., dadi).







Genome-Wide Inference of Ancestral Recombination Graphs



Matthew D. Rasmussen^{1#}*, Melissa J. Hubisz¹, Ilan Gronau¹, Adam Siepel^{1,2}*

1 Department of Biological Statistics and Computational Biology, Comell University, Ithaca, New York, United States of America, 2 European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambs, United Kingdom

Abstract

The complex correlation structure of a collection of orthologous DNA sequences is uniquely captured by the "ancestral recombination graph" (ARG), a complete record of coalescence and recombination events in the history of the sample. However, existing methods for ARG inference are computationally intensive, highly approximate, or limited to small

Selective Sweeps





Selective Sweeps





Coalescence Tree



Coalescence Models



Importance Sampling

$$\int_{\Psi} \Pr(X \mid G) p(G \mid \Theta) dG = \int_{\Psi} \Pr(X \mid G) p(G \mid \Theta) \frac{h(G)}{h(G)} dG$$
$$= E \left[\frac{\Pr(X \mid G) p(G \mid \Theta)}{h(G)} \right]$$

So

$$\Pr(X \mid \Theta) \approx \frac{1}{k} \sum_{i=1}^{k} \frac{\Pr(X \mid G_i) p(G_i \mid \Theta)}{h(G_i)}$$

where $G_{i,i} = 1, 2, ..., k$, has been simulated from h(G).



Stern et al. 2019. bioRxiv http://dx.doi.org/10.1101/592675.

Aaron Stern







Stern et al. 2019. bioRxiv http://dx.doi.org/10.1101/592675.



Aaron Stern



Lactase rs4988235



Stern et al. 2019. bioRxiv http://dx.doi.org/10.1101/592675.



Stern et al. 2019. bioRxiv http://dx.doi.org/10.1101/592675.

Allele frequency







Problem 2: DNA is chopped into short 'reads' with high error rates

Ļ	Ļ	Ļ	Ļ

Reverse problem

1	1	1	1



Frequency spectrum



10X Next Gen. Sequencing data



Using a fixed cut-off for SNP calling can never produce unbiased allele frequency estimates for all values of *p*.

Maximum Likelihood Estimation

- Assume genotype likelihoods, $p(X_i^{(v)} | G_i^{(v)})$, can be calculated in site *v*, individual *d*, for all *v* and *d*.
- Parameterize the frequency spectrum as $\mathbf{P} = (p_0, p_1, ..., p_{2k})$.
- Likelihood function:

$$L(\mathbf{P}) = \prod_{v} \left(\sum_{j=0}^{2k} p_{j} \left[\sum_{G_{1}^{(v)}} \dots \sum_{G_{k}^{(v)}} c(j, G^{(v)}) \prod_{d=0}^{k} p(X_{d}^{(v)} \mid G_{d}^{(v)}) \right] \right)$$
$$c(j, G^{(v)}) = \begin{cases} \binom{2k}{j}^{-1} 2^{\sum_{d=1}^{k} I\left(G_{d}^{(v)} = 1\right)} & \text{if } \sum_{d=1}^{k} G_{d}^{(v)} \neq j \\ 0 & \text{else} \end{cases}$$

$$G^{(v)} = (G_1^{(v)}, G_2^{(v)}, ..., G_k^{(v)}), \quad G_d^{(v)} \in \{0, 1, 2\}$$

Nielsen et al. 2012. PloS One 10.1371

Dynamic Programming Algorithm

Initialization:

Set $h_0 = p(X_d | G_d = 0)$, $h_1 = 2p(X_d | G_d = 1)$, $h_2 = p(X_d | G_d = 2)$, and $h_j = 0$ for j = 3, 4, ..., 2k.

Recursion

For *d* = 2, 3,..., *k*:

For j = 2d, 2d-1,...,2: Set $h_j = p(X_d | G_d = 2)h_{j-2} + 2p(X_d | G_d = 1)h_{j-1} + p(X_d | G_d = 0)h_j$ Set $h_{j-1} = p(X_d | G_d = 0)h_{j-1} + 2p(X_d | G_d = 1)h_{j-2}$ Set $h_{j-2} = p(X_d | G_d = 0)h_{j-2}$

Termination

Set
$$h_{j,=} h_j \binom{2k}{j}^{-1}$$
 for $j=0,1,2,...,2k$.

Nielsen et al. 2012. PloS One 10.1371
Maximum Likelihood Estimation

• The likelihood function can then be expressed as

$$L(\mathbf{P}) = \prod_{v} \left(\sum_{j=0}^{2k} p_j h_j^{(v)} \right)$$

- Derivatives can be calculated analytically.
- Algorithm which is quadratic in the number of individuals and linear in the number of sites.
- Optimization to eight decimals precision takes < 1 minute for 1 GB and 60 individuals on one desktop computer.



The distribution of true and estimated folded SFS in a sample from 50 MB 10 diploid individuals, where 1% of all SNPs are variable in the population and follow a distribution of allele frequencies, p, proportional to 1/p. An error rate of 0.5% is assumed. The mean sequencing depths is 5X.

Nielsen et al. 2012. PloS One 10.1371



The distribution of true and estimated folded SFS in a sample from 50 MB 10 diploid individuals, where 1% of all SNPs are variable in the population and follow a distribution of allele frequencies, p, proportional to 1/p. An error rate of 0.5% is assumed. The mean sequencing depths is 1X

Nielsen et al. 2012. PloS One 10.1371

50 Tibetan Exomes



Yi, X. et al. 2010. Science 329: 75-78.

Tibetans have approx. constant hemoglobin concentrations as a function of increasing altitude up to 4000 meters (Beall *et al.* 2006).

Photo by (

02007

50 Tibetan Exomes



Yi, X. et al. 2010. Science 329: 75-78.

		1.16		
		2	EPAS1 EPAS1 Variant frequency Variant frequer Tibet Han Chinese	ncy
1	A A A A A A A A A A A A A A A A A A A		87% 8%	
X	EPAS1	hemoglobin	THE CONTRACT	
	Genotype	concentration		
	CC	178		
	CG	178.9		
	GG	167.5		5

Yi, X. et al. 2010. Science 329: 75-7



Huerta Sanchez et al. 2014. Nature 512: 194–197

Denisovans





Denisovan!



Huerta Sanchez et al. 2014. Nature 512: 194–197

Introgression

Species 1

Species 2



Admixture

Population 1



Population 2



Assignment and Admixture

- Multi loci data?
- -To which population does an individual belong?
- How many groupings should we choose?

Structure analysis of 1056 individuals from 52 populations for 377 microsatelite loci. Rosenberg et al. (2002)



Transition to the Neolithic



Figure 1. Isochrone map. The spread of the Neolithic transition, obtained by interpolating the dates in calibrated years before present (BP) of 918 Early Neolithic sites (circles) in Europe and the Near East (see the electronic supplementary material for details on the dataset and interpolation). Map created with ArcGIS 10.

Transition to the Neolithic











Contemporary Eurasians





STEPPE IN TIME

An ancient-DNA study links the Corded Ware culture of northern Europe with the Yamnaya culture of the Eurasian steppe. It points to a mass migration northwest that would support the Steppe hypothesis, one of two theories that compete to explain the origins of the Indo-European family of languages.



From Haak et al. (2015)



Individual samples ...



Map data © 2016 Google, INEGI · Phylogenetic trees: http://www.jade-cheng.com/graphs/





Admixture tracts



Recombination Clock!

Ancestry painting (23andme.com)

the second se

African American Man		
	Europe 63%	
	Africa 35%	
	Asia 2%	
	Not Genotyped	

Predictions assuming independence





Kostenki



Andaine Seguin-Orlando et al. 2014. Science: 346: 1113-1118

Kostenki



Age ~54,000 years

Andaine Seguin-Orlando *et al.* 2014. Science: 346: 1113-1118

Rapanui

Polynesian Expansion



Rapanui



Rapanui







Rapanui

