Likelihood-based demographic inference using the coalescent



Raphael Leblois Centre de Biologie et de Gestion des Populations, CBGP INRA, Montpellier

Master MEME, March 2011

Likelihood-based demographic inference using the coalescent

- 1. Reminder : main coalescence principles
- 2. Simulating coalescent trees and polymorphism data
- 3. Likelihood-based inferences
- 4. Maximum likelihood and Isolation by Distance



In the coalescent theory, we look at the genealogy of a sample of genes going backward in time until the most recent common ancestor (MRCA)



- \rightarrow a new approach in population genetics :
- ✓ Classical approach
 - Population
 - Gene frequencies
 - Forward in time

- ✓ Coalescent approach
 - Sample
 - Gene Genealogies
 - backward in time



Probability of coalescence of 2 genes in one generation = probability that the two genes have a common parental gene

$$P(T_2 = 1) = \frac{1}{N}$$

Coalescence of 2 genes in 2 generations in a haploid population of size N

Past

Present

(Prob. that the 2 gene do not coalesce at *t*=1) *(Prob. that the 2 gene coalesce at t=2)

$$P(T_2 = 2) = \left(1 - \frac{1}{N}\right) \frac{1}{N}$$

Coalescence of two genes in *t* generations in a haploid population of size *N*

Past



Present

(Prob. that the 2 gene do not coalesce in the first *t*-1 generations) *(Prob. that the 2 gene coalesce at *t*)

$$P(T_2 = t) = \left(1 - \frac{1}{N}\right)^{t-1} \frac{1}{N}$$

Coalescence of two genes in *t* generations in a haploid population of size *N*

for x << 1 $(1-x)^t \approx e^{-xt}$

The discret geometric distribution can be approximated by an continuous exponential distribution for large *N*

$$P(T_2 = t) = \left(1 - \frac{1}{N}\right)^{t-1} \frac{1}{N} \approx \frac{1}{N} e^{-Nt}$$

Coalescence times follow an exponential ditribution of rate *N* (it is also its expectation)

T2

Т3

Coalescence of *j* genes in *t* generations in a haploid population of size *N*

Assumption: no multiple coalescence for large N

 $\binom{2}{j} = j^*(j-1)/2$ gene pairs can coalesce with probability 1/N

Pr(two genes among *j* coalesce in one generation) = $\frac{j(j-1)}{2N}$

coalescence times for a sample of *j* genes/lineages follow a geometric distribution with parameter $j^*(j-1)/2N$, and can be approximated by an exponential distribution with expectation $2N/(j^*(j-1))$

$$\Pr(T_j = t) = (1 - \frac{j(j-1)}{2N})^{t-1} (\frac{j(j-1)}{2N}) \approx \frac{j(j-1)}{2N} e^{-\frac{j(j-1)}{2N}t}$$

Coalescence of *j* genes in *t* generations in a haploid population of size *N*





$$E(T_j) = \frac{2N}{j(j-1)}$$

the larger the sample size or lineage number is, the larger the expected coalescence times are

 $\operatorname{var}(T_{j}) = \frac{4N^{2}}{j^{2}(j-1)^{2}}$

Coalescence times have high variance : two independent loci could show very different coalescence times, and thus very different coalescent trees (genealogies)

Coalescence en t générations de j lignées

$$E(T_j) = \frac{2N}{j(j-1)}$$

the larger the sample size or lineage number is, the larger the expected coalescence times are



 $var(T_j) = \frac{4N^2}{j^2(j-1)^2}$ Coalescence times have high two independent loci could show very different coalescence times, and thus very different Coalescence times have high variance : coalescent trees (genealogies)



TMRCA : length of the coalescent trees

TMRCA = Time to the Most Recent Common Ancestor

- = Time of the last node (coalescence) of the tree
- = Tree length



$$E[TMRCA] = \sum_{i=2}^{j} E[T_i] = \sum_{i=2}^{j} \frac{2N}{i(i-1)}$$
$$= 2N \times \sum_{i=2}^{j} (\frac{1}{i-1} - \frac{1}{i})$$
$$= 2N(1 - \frac{1}{j})$$

TMRCA expectation tends to 2N for large samples
 TMRCA of a relatively small sample is close the TMRCA of the whole population

coalescent trees and mutations

Under neutrality assumption, mutations are independent of the genealogy, because genealogical process strictly depends on demographic parameters

First, genealogies are build given the demographic parameters considered (e.g. N),



Then mutation are added a posteriori on each branch of the genealogy, from MRCA to the leaves

We thus obtain polymorphism data under the demographic and mutational model considered

coalescent trees and mutations

The number of mutations on each branch is a function of the mutation rate of the genetic marker (μ) and the branch length (t). μ = mean number of mutation per locus per generation. e.g. 5.10⁻⁴ for microsatellites, 10⁻⁷ per nucleotide for DNA sequences



For a branch of length *t*, the number of mutation thus follows a binomial distribution with parameters (μ ,t). Often approximated by a Poisson distribution with parameter (μ *t).

$$\Pr(k \text{ mut}|t) = \frac{(\mu t)^{k} e^{-\mu t}}{k!}$$

Arbre de coalescence et mutations

Different mutational models for the different genetic markers, e.g. :

✓ For DNA sequences : mutation matrix for different nucleotide transition rates ($Pr[A \rightarrow T]$, $Pr[A \rightarrow C]$, $Pr[T \rightarrow G]$, etc...)

✓ For SNPs : $Pr(Ancestral \rightarrow Derived)=\mu$, $Pr(Anc. \rightarrow Der.)=0$

✓ For microsatellites :stepwise models



SMM (Stepwise mutation model): each mutation add or remove a motif to the parental allele

Main advantages of the coalescent

- The coalescent is a powerful probabilistic model for gene genealogies
- The genealogy of a population genetic sample, and more generally its evolutionary history, is often unknown and cannot be repeated ⇒ the coalescent allows to take this unknown history into account
- The coalescent often simplifies the analyses of stochastic population genetic models and their interpretation
- Genetic data polymorphism largely reflects the underlying genealogy ⇒ the coalescent greatly facilitate the analysis of the observed genetic variability and the understanding of evolutionary processes that shaped the observed genetic polymorphism.

Main advantages of the coalescent

- The coalescent allows extremely efficient simulations of the expected genetic variability under various demo-genetic models (sample vs. entire population)
- The coalescent allows the development of powerful methods for the inference of populational evolutionary parameters (genetic, demographic, reproductive,...), some of those methods uses all the information contained in the genetic data (likelihood-based methods)

Trees and polymorphism data simulation

- reminder :
 - ✓ For neutral markers, the number of offspring is independent of the genetic types of their parents
 - → Demographic processes are thus independent of mutational processes
 - Simulation of polymorphism data can thus be done in two steps :
 - (1) Tree simulation : topology and branch length
 - (2) Addition of mutations on the tree

Coalescent tree simulation

Two main methods :

Coalescent continuous approximations

very fast but approximations only valid for large population sizes, weak mutation and migration rates and "simple" demographic models

Generation by generation

Ok for all demograophic and mutational models but relatively slow

RAPIDITY :

Continuous approximations > Generation by generation

FLEXIBILITY :

Generation by generation > Continuous approximations



- very simple and exact (without any approximations):
 - ✓ Go backward in time generation by generation
 - At each generation, we stochastically draw potential events affecting the genealogy

e.g. coalescence, migration, recombinaison

 Stop at the most recent common ancestor of all sampled genes = MRCA

- Toy example :
 - ✓ 4 gene sample
 - ✓ single neutral locus
 - ✓ panmictic haploid population of size N=10

• Example : 4 genes, neutral, 1 pop N=10

nodes / lineages numbering	1	2	3	4
random number between 1 and <i>N</i> for each lineage				
lineage starting generation	0	0	0	0





• Example : 4 genes, neutral, 1 pop N=10



Prob for a coalescence in *j* lineages in one generation

=j(j-1)/2N

= probability of drawing 2 identical integers in *j* uniform drawings between 1 and *N*

- Example : 4 genes, neutral, 1 pop N=10
 Probability of a coalescence in *j* lineages in one generation
 =j(j-1)/2N
 - = probability of drawing 2 identical integers in *j* uniform drawings between 1 and *N*

in other terms, we randomly and uniformly draw a parent for each gene/lineage among the N potential parents (stable population size)

Genes/lineages sharing the same parent coalesce

• Example : 4 genes, neutral, 1 pop N=10



Coalescence at generation 1 of nodes/lineages 3 and 4





• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	1	3	5
random number between 1 and <i>N</i> for each lineage	2	5	6
lineage starting generation	0	0	1

Coalescence at generation 1 of nodes/lineages 3 and 4

new node 5



• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	1	3	5
random number between 1 and <i>N</i> for each lineage	3	1	7
lineage starting generation	0	0	1

nothing happened at generation 2





• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	1	3	5
random number between 1 and <i>N</i> for each lineage	7	4	8
lineage starting generation	0	0	1

nothing happened at generation 3





• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	1	3	5
random number between 1 and <i>N</i> for each lineage	5	2	5
lineage starting generation	0	0	1

Coalescence at generation 4 of nodes/lineages 1 and 5



• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	3	6
random number between 1 and <i>N</i> for each lineage	2	5
lineage starting generation	0	5

Coalescence at generation 4 of nodes/lineages 1 and 5 new node 6



• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	3	6
random number between 1 and <i>N</i> for each lineage	3	9
lineage starting generation	0	5

nothing at generation 5,6,...



• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	3	6
random number between 1 and <i>N</i> for each lineage	7	7
lineage starting generation	0	5

Coalescence at generation 20 of the two last lineages 3 and 6

Gn=20

• Example : 4 genes, neutral, 1 pop *N*=10

nodes / lineages numbering	3	6
random number between 1 and <i>N</i> for each lineage	7	7
lineage starting generation	0	5

Coalescence at generation 20 of the two last lineages 3 and 6 new node 7 = MRCA of (6) the sample



The coalescence tree (topology and branch lengths) is build.

It is a stochastic process, so if we build many trees, they will all be different but share some common properties.

To get polymorphism data, we need to add mutations on the tree...



Coalescent tree simulation Hudson continuous approximations

• Principle: 2 successive steps

(1) The topology of the tree is build by randomly coalescing lineages

(2) Branch length are simulated using expected coalescence times between two coalescence events
- Example : 4 genes, neutral, 1 pop *N*=10
 - (1) The topology of the tree is build by randomly coalescing lineages
- 1st coalescence = random draw of 2 lineages among the 4 → lineages 2 and 4 coalesce to give lineage 5

• Example : 4 genes, neutral, 1 pop *N*=10

1

- (1) The topology of the tree is build by randomly coalescing lineages
- 2^d coalescence = random draw of 2 lineages among the 3 lineages left → lineages 1 and 5 coalesce to give lineage 6 ③ ③

(5)

(4)

2

3

3

- Example : 4 genes, neutral, 1 pop *N*=10
 - (1) The topology of the tree is build by randomly coalescing lineages
 - 3^d and last coalescence = the last 2 lineages 6 and 3 coalesce to give lineage 7, the MRCA



Example : 4 genes, neutral, 1 pop N=10
 (1) Topology is build



Example : 4 genes, neutral, 1 pop *N*=10
(2) Branch length simulation

there are 3 branch lengths to simulate T_4 , T_3 , T_2



• Example : 4 genes, neutral, 1 pop *N*=10

3 branch lengths to simulate T4, T3, T2

$$\Pr(T_{j} = k) = \frac{j(j-1)}{2N} e^{\frac{-j(j-1)}{2N}k}$$

T₄ drawn from an exponential distribution with parameter (expectation) j(j-1)/2N=4*3/2*10

(algorithms to draw exponential deviates are availables)

• Example : 4 genes, neutral, 1 pop N=10

3 branch lengths to simulate T4, T3, T2

Ex:

T4 drawn from exp. $(j(j-1)/2N=4*3/2*10) \rightarrow 1,2$ T3 drawn from exp. $(3*2/2*10) \rightarrow 2,6$ T2 drawn from exp. $(2*1/2*10) \rightarrow 15,7$

(5)



under the demographic model considered!

General principle (reminder) :

Mutations are distributed on the different branches from the MRAC to the leaves as a function of the mutation rate μ

Each mutation induce a change in the allelic/nucleotidic state of the descending node

This genetic state change is made according to the mutational model considered, which may reflect real mutational processes of some genetic markers



On a branch of length t, the number of mutation follow a binomial with parameters (μ ,t)

Approximated by a Poisson distribution with parameter ($\mu^* t$)

$$\Pr(k \text{ mut}|t) = \frac{(\mu t)^{k} e^{-\mu t}}{k!}$$



Example for microsatellites under a SMM : gain or loss of a motif (repeat) for each mutation

addition of mutation numbers on each branch following the Poisson distribution

$$\Pr(k \text{ mut}|t) = \frac{(\mu t)^{k} e^{-\mu t}}{k!}$$



- Example for microsatellites under a SMM : gain or loss (p=0.5) of a motif (repeat) for each mutation
- Choice of the MRCA type (random): 20
- node 7 to 6 : one time $\pm 1 \rightarrow 21$
- node 6 to 1 : one time $\pm 1 \rightarrow 22$
- node 6 to 5 : 0 time $\pm 1 \rightarrow 21$
- node 5 to 2 : one time $\pm 1 \rightarrow 20$
- node 5 to 4 : 0 time $\pm 1 \rightarrow 21$



Example for microsatellites under a SMM : gain or loss (p=0.5) of a motif (repeat) for each mutation

node 7 to 3 : 3 times $\pm 1 \rightarrow 19$

A polymorphism sample of 4 genes is obtained with allelic states 19, 20, 21, 22



Example on DNA sequence markers (5 bp). Choice of the ancestral sequence (ATTGC) independent mutation on each site

- 7 to 6 : 1 mut on site 1 → TTTGC
- 6 to 1 : 1 mut on site 3 → TTAGC
- 5 to 2 : 1 mut on site 5 → TTTGG
- 7 to 3 : 1 mut on each site 2,3,4 → AAACC



Example on DNA sequence markers (5 bp).

Choice of the ancestral sequence (ATTGC)

independent mutation on each site

The polymorphism sample is then composed of 4 different sequences :

TTAGC,TTTGG,TTTGC,AAACC



• <u>Exploratory approaches</u> : to study the effects of various parameters on the shape of coalescent trees and on the distribution of polymorphism in a sample



- Exploratory approach : demographic effects
 - growing population size (e.g. invasion of a new habitat)
 - There are more ancient coalescences (small *N*) than recent coalescences (large *N*), coalescent trees thus have longer terminal branches



A population size growth induces an excess of low frequency alleles (rare alleles) 53

- Exploratory approach : demographic effects
 - population size contraction (e.g. threatened species)
 - There are more recent coalescences (small *N*) than ancient coalescences (large *N*), coalescent trees thus have shorter terminal branches



A contraction induces a deficit of low frequency alleles

- Exploratory approach : to study the effects of various parameters on the shape of coalescent trees, on the distribution of polymorphism in a sample and on various sumary statistics computed on a genetic sample (e.g. H_e, F_{ST},...)
- Simulation tests : to create simulated data sets to test the precision and robustness of genetic data analysis methods
- Inferential approach : to estimate populational evolutionary parameters (pop sizes, dispersal, demographic history) from polymorphism data

- Inferential approaches are based on the modeling of population genetic processes. Each population genetic model is characterized by a set of demographic and genetic parameters P
- The aim is to infer those parameters from a polymorphism data set (genetic sample)
- The genetic sample is then considered as the realization ("output") of a stochastic process defined by the demo-genetic model

- First, compute or estimate the probability Pr(D | P*) of observing the data D given some parameter values P*, it is the likelihood : L(P* | D)=Pr(D | P*)
- Second, find the set of parameter values that maximize this probability of observing the data (maximum likelihood method)

Maximum likelihood method



!! many parameters \rightarrow large parameter space to explore **!!**

- Problem : Most of the time, the likelihood Pr(D|P) of a genetic sample cannot be computed directly because there is no explicite mathematical expression
- However, the probability Pr(D|P,G_i) of observing the data D given a specific genealogy G_i and the parameter values P can be computed.
- then we take the sum of all genealogy-specific likelihoods on the whole genealogical space, weighted by the probability of the genealogy given the parameters : $L(P|D) = \int_{G} \Pr(D|G;P)\Pr(G|P)dG$

 The likelhood can be written as the sum of Pr(D|P,G_i) over the genealogical space (all possible genealogies) :

$$L(P|D) = \int_{G} \Pr(D|G;P)\Pr(G|P)dG$$

mutational parameters Coalescent theory
demographic parameters

 Genealogies are nuisance parameters (or missing data), they are important for the computation of the likelihood but there is no interest in estimating them

very different from the phylogenetic approaches

$$L(P|D) = \int_{G} \Pr(D|G;P)\Pr(G|P)dG$$

Sum over all possible genealogies
 \Rightarrow usually untractable !!!

Monte Carlo simulations are used : a large number *K* of genealogies are simulated according to Pr(G|P) and the mean over those simulations is taken as the expectation of Pr(D|G;P):

$$L(P|D) = E_{pr(G|P)} \Big[\Pr(D|G;P) \Big] \approx \frac{1}{K} \sum_{k=1}^{K} \Pr(D|G_k;P) \Big]$$

simulation of many genealogies is necessary to get a good estimation of the likelihood

$$L(P|D) = E_{pr(G|P)} \left[\Pr(D|G;P) \right] \approx \frac{1}{K} \sum_{k=1}^{K} \Pr(D|G_k;P)$$

Monte Carlo simulations are often not very efficient because there are too many genealogies giving extremely low probabilities of observing the data, more efficient algorithms are used to explore the genealogical space and focuss on genealogies well supported by the data.

More efficient algorithms :

- IS : Importance Sampling
- MCMC : Monte Carlo Markov chains associated with Metropolis-Hastings algorithm

allows better exploration of the genealogies proportionnaly to their probability of explaining the data P(D|P;G).

Demographic inference under the coalescent the approach of Felsenstein et al. (MCMC)

- Probability of a genealogy given the parameters of the demographic model Pr(G_i|P) can be computed from the continuous time approximations (cf. Hudson approximations to construct coalescent trees)
- then the probability of the data given a genealogy and mutational parameters Pr(D|G_i,P) can be easily computed from the mutation model parameters, the mutation rate and the Poison distribution of mutations.
- From this, an efficient algorithm to explore the genealogical and the parameter spaces should allows the inference of the likelihood over the spaces.

Demographic inference under the coalescent the approach of Felsenstein et al. (MCMC)

 Probability of a genealogy given the parameters of the demographic model (N_i or $\{N_i, m_{ii}\}$ if structured populations)

example for a unique panmictic population

migr

$$Pr(G|P) = \prod_{\tau=1}^{TMRCA} \begin{pmatrix} j_{\tau}(j_{\tau}-1) & j_{\tau}(j_{\tau}-1) \\ 4N & 4N \end{pmatrix}$$
Product over all demographic events (coalescence or migration) affecting the genealogy
lineage number before the event
Time interval between this event and the previous one



Demographic inference under the coalescent the approach of Felsenstein et al. (MCMC)

• Probability of a genealogy given the parameters of the demographic model $\Pr(G|P) = \prod_{\tau=1}^{T} \left(\frac{j_{\tau}(j_{\tau}-1)}{4N} e^{\frac{j_{\tau}(j_{\tau}-1)}{4N}} k_{\tau} \right)$

• Probability of the sample given the genealogy and mutational parameters (mutation rate μ , M_{mut} mutation matrix)

$$Pr(D|G) = \prod_{b=1}^{B} \left((M_{mut})^{i_b} (\mu L_b)^{i_b} e^{\mu L_b} \right)^{i_b}$$

Product over all tree
ranches
Poisson probability of getting i_b
mutations on a time interval L_b
67

Demographic inference under the coalescent the approach of Felsenstein et al. (MCMC)

• Probability of a genealogy given the parameters of the demographic model $\Pr(G|P) = \prod_{\tau=1}^{T} \left(\frac{j_{\tau}(j_{\tau}-1)}{4N} e^{\frac{j_{\tau}(j_{\tau}-1)}{4N}} k_{\tau} \right)$

• Probability of the sample given the genealogy and mutational parameters $B_{B}(\mu = \sum_{i_{b}} \mu_{i_{b}})^{i_{b}} = \mu_{i_{b}}$

$$\Pr(D|G) = \prod_{b=1}^{B} \left(\left(P_{mut} \right)^{i_b} \frac{(\mu L_b)^{i_b}}{i_b!} e^{\mu L_b} \right)$$

• by definition $L(P|D) \approx \frac{1}{K} \sum_{k=1}^{K} \Pr(D|G_k; P) \approx \frac{1}{K} \sum_{k=1}^{K} \Pr(D|G_k) \Pr(G_k|P)$

Demographic inference under the coalescent the approach of Felsenstein et al. (MCMC)

It is a very complexe problem because of the large genealogical and parameter spaces to explore

more parameters more complexe genealogies

Models with more parameters will need more computation times or more efficient algorithms to explore the 2 spaces

→ better to always try to consider simple but robust models

Metropolis-Hastings algorithm for the parameter space

(1) start from a point (vector of parameter values, Θ)

(2) propose a change in the parameter space Θ' from the proposal distribution $q(\Theta \rightarrow \Theta')$

(3) accept the change with probability

$$h = \min\left(1, \frac{L(\Theta'; D)}{L(\Theta; D)} \frac{P(\Theta')}{P(\Theta)} \frac{q(\Theta' \to \Theta)}{q(\Theta \to \Theta')}\right)$$

(4) go back to (1)

This algorithm ensure that the parameter space is explored proportionnaly to the likelihood

Metropolis-Hastings algorithm : an efficient exploration of the space



One example : MsVar (Beaumont 1999)

Demographic model : one population with variable size
 Taille Population contraction or expansion



3+1 parameters N_0 , N_1 et t_g (+ μ) to be estimated using a MCMC Metropolis-Hastings algorithm
- Monte Carlo Markov chains simulation using the Metropolis-Hastings algorithm (MCMC)
 - \checkmark To explore the genealogy space
 - \checkmark and the parameter space

Monte Carlo Markov chains simulation (MCMC)

✓ To explore the genealogies, we then build a new genealogy by a "partial deletion-reconstruction" algorithm from the ourrest one :

current one :



Figure 2 .- Transition from an old genealogy G, to a new genealogy G. Dotted lines are the times of coalescences or migrations. (A) Genealogy G, with migration. The black bar marks a migration from the white population to the black or vice versa; z is the node to be picked. (B) Partial genealogy G, after drawing the coalescent node z at random and dissolving the branch to the next coalescent below, (C) Simulation of the coalescent with migration. One possible outcome with three consecutive steps is shown: (1) using Equation 7 a new time interval is drawn and a migration event from white to black is also drawn (Equation 8) and the lineage is extended down to that event; (2) a new time interval is drawn: it extends too far back, so the lineage advances down to the time j; (3) a new time interval is drawn with a coalescent event at its bottom end. The process stops at time k. (D) the final configuration G.

Monte Carlo Markov chains simulation (MCMC)

✓ To explore the genealogies, we then build a new genealogy by a "partial deletion-reconstruction" algorithm from the current one :



potential problem : Trees are correlated...

Monte Carlo Markov chains simulation (MCMC)

✓ To explore the genealogies, a new genealogy is build by a "partial deletion-reconstruction" algorithm from the current genealogy

✓ in parallel, the parameter space will be explored by modifying parameter values in the MCMC

 \rightarrow at each step of the MCMC:

either the genealogy is modified,

or a parameter value is modified

Results presented by Renaud Vitalis

Les Orangs-Outans et la déforestation



Figure 5. Total estimated orangutan populations since the early Holocene. Data are taken from Rijksen & Meijaard.^48 $\,$



Quelle est la cause de la baisse de taille de population? La génétique peut elle nous aider?

• Les Orangs-Outans et la déforestation : les données



200 individus 14 locus microsatellites





• Les Orangs-Outans et la déforestation :



Fig.1 Population size change

MsVar détecte bien un réduction de taille de population



• Les Orangs-Outans et la déforestation :



MsVar détecte bien un réduction de taille de population



• Les Orangs-Outans et la déforestation :



MsVar détecte bien un réduction de taille de population



et permet d'obtenir une datation

FE : Forest exploitation F: Farmers HG: Hunter-gatherers

• Les Orangs-Outans et la déforestation :



MsVar détecte bien un réduction de taille de population



et permet d'obtenir une datation: l'exploitation de la forêt semble être la cause...

FE : Forest exploitation F: Farmers HG: Hunter-gatherers

• Probability of a sample *D* given mutational and demographic parameters of the model considered can be computed using the probabilities of transition between the different events affecting the genealogy (with mutations), i.e. the different ancestral states H_k .

a genealogy = genealogical history of the sample can be divided into *m* successive events/states H_k (coalescences, mutations, migrations)

 $G_{i} = \{H_{k}; 0 > k > -m\} = \{H_{0}, H_{-1}, \dots, H_{-m+1}, H_{-m}, \}$

 Probability of a sample D given mutational and demographic parameters of the model considered can be computed using the probabilities of transition between the different ancestral states H_k.

$$G_{i} = \{H_{k}; 0 > k > -m\} = \{H_{0}, H_{-1}, \dots, H_{-m+1}, H_{-m}, \}$$

the probability of a given state H_k can be expressed as the probability of all possible ancestral states H_{k-1} multiplied by their associated transition probability possible $Pr(H_k|H_{k-1})$

$$p(H_k) = \sum_{\{H_{k-1}\}} p(H_k \mid H_{k-1}) p(H_{k-1}).$$

- the principle of Griffiths et al. importance sampling approach :
 - = the recurrence between ancestral samples

$$p(H_k) = \sum_{\{H_{k-1}\}} p(H_k \mid H_{k-1}) p(H_{k-1}).$$

exploring all possible ancestral sample configurations is usually impossible,

Monte Carlo simulations are used to explore a given number *K* of possible genealogies by building genealogies backward in time from the initial sample configuration H_0 to the MRCA (Absorbing Markov chains with absorbing state being the MRCA)

the principle of Griffiths et al. importance sampling approach :
 = the recurrence between ancestral samples p(H_k) = ∑ p(H_k | H_{k-1})p(H_{k-1}).
 ➡ Monte Carlo simulations on Z possible genealogies build backward in time from D=H₀ to the MRCA

$$p(D = H_0) \approx E_Z[p(H_{z,0} | H_{z,-1}) p(H_{z,-1} | H_{z,-2}) \dots p(H_{z,-m+1} | H_{z,-m}) p(H_{z,-m})]$$

the probability of the data for a given genealogy *z* is the product of all transition probabilities between ancestral states from H_0 to H_{-m} = the MRCA

- the principle of Griffiths et al. importance sampling approach : = the requirement between encostrol complete $r(H_1) = \sum_{i=1}^{n} r(H_1 + H_{i-1})r(H_{i-1})$
 - = the recurrence between ancestral samples $p(H_k) = \sum_{\{H_{k-1}\}} p(H_k | H_{k-1})p(H_{k-1})$. Monte Carlo simulations

$$p(D = H_0) \approx E_Z[p(H_{z,0} | H_{z,-1}) p(H_{z,-1} | H_{z,-2}) \dots p(H_{z,-m+1} | H_{z,-m}) p(H_{z,-m})]$$

This is the approach of Griffiths & Tavaré 1984, implemented in GeneTree for DNA sequence

- Genealogies / coalescent trees are explored according to the forward transition probabilities $Pr(H_k|H_{k-1})$
- The approach is working but relatively inefficiently because it uses forward transition probabilities to build genealogies backward
 - too many tree simulation needed

• the new importance sampling of Delorio & Griffiths 2004 :

= the recurrence between ancestral samples $p(H_k) = \sum_{\{H_{k-1}\}} p(H_k | H_{k-1})p(H_{k-1})$. it is much better to try to simulate from the backward transition probabilities $Pr(H_{k-1} | H_k)$. Those probabilities are unknown but they may be approximated

$$p(H_k) = \sum_{\{H_{k-1}\}} \frac{p(H_k \mid H_{k-1})}{\hat{p}(H_{k-1} \mid H_k)} p(H_{k-1}) \hat{p}(H_{k-1} \mid H_k)$$
$$= \mathbb{E}_{\hat{p}} \sum_{\{H_{k-1}\}} \left[\frac{p(H_k \mid H_{k-1})}{\hat{p}(H_{k-1} \mid H_k)} \mid H_k \right].$$

the new importance sampling of Delorio & Griffiths 2004 :

= the recurrence between ancestral samples $p(H_k) = \sum_{\{H_{k-1}\}} p(H_k \mid H_{k-1}) p(H_{k-1})$. it is much better to try to simulate from the backward transition probabilities $Pr(H_{k-1}|H_k)$. Those probabilities are unknown but they may be approximated

$$p(H_k) = \sum_{\substack{\{H_{k-1}\}\\ \{H_{k-1}\}}} \frac{p(H_k \mid H_{k-1})}{\hat{p}(H_{k-1} \mid H_k)} p(H_{k-1})\hat{p}(H_{k-1} \mid H_k)$$

$$= \mathbb{E}_{\hat{p}} \sum_{\substack{\{H_{k-1}\}\\ \{H_{k-1}\}}} \left[\frac{p(H_k \mid H_{k-1})}{\hat{p}(H_{k-1} \mid H_k)} \mid H_k \right].$$

$$= 8$$

- the new importance sampling of Delorio & Griffiths 2004 :
 - = the recurrence between ancestral samples $p(H_k) = \sum_{\{H_{k-1}\}} p(H_k \mid H_{k-1})p(H_{k-1})$ it is much better to simulate genealogies from approximated backward transition probabilities $\hat{p}(H_{k-1} \mid H_k)$

$$p(H_0) = \mathbb{E}_{\hat{p}} \left[\frac{p(H_0 \mid H_{-1})}{\hat{p}(H_{-1} \mid H_0)} \cdots \frac{p(H_{-m+1} \mid H_{-m})}{\hat{p}(H_{-m} \mid H_{-m+1})} \times p(H_{-m}) \right],$$

- the new importance sampling of Delorio & Griffiths 2004 :
 - = the recurrence between ancestral samples $p(H_k) = \sum_{\{H_{k-1}\}} p(H_k \mid H_{k-1})p(H_{k-1})$. iii is much better to simulate genealogies from approximated backward transition probabilities $\hat{p}(H_{k-1} \mid H_k)$

$$p(H_0) = \mathbb{E}_{\hat{p}} \begin{bmatrix} p(H_0 \mid H_{-1}) \\ \hat{p}(H_{-1} \mid H_0) \\ \cdot p(H_{-m} \mid H_{-m+1}) \end{bmatrix}$$

$$\cdot \frac{p(H_{-m+1} \mid H_{-m})}{\hat{p}(H_{-m} \mid H_{-m+1})}$$

$$\times p(H_{-m})],$$

the probability of the data for a given genealogy *z* is the product of all transition importance weights $w_{IS}(H_k, H_{k-1})$ between ancestral states from H_0 to H_{-m} = the MRCA

- the new importance sampling of Delorio & Griffiths 2004 :
 - = the recurrence between ancestral samples $p(H_k) = \sum_{\{H_{k-1}\}} p(H_k \mid H_{k-1})p(H_{k-1})$. iii is much better to simulate genealogies from approximated backward transition probabilities $\hat{p}(H_{k-1} \mid H_k)$

$$p(H_0) = \mathbb{E}_{\hat{p}} \left[\frac{p(H_0 \mid H_{-1})}{\hat{p}(H_{-1} \mid H_0)} \cdots \frac{p(H_{-m+1} \mid H_{-m})}{\hat{p}(H_{-m} \mid H_{-m+1})} \times p(H_{-m}) \right],$$

The approach of Delorio & Griffiths 2004 is much more efficient, 1,000 times less trees to explore!

Computation time and model complexity with Griffiths & Tavaré (1984) algorithm

Number of genealogies (= iterations) and time to correctly infer the likelihood of a sample at a single parameter point (one vector Θ of parameter values)





• the new importance sampling of Delorio & Griffiths 2004 :

possible genealogies (= coalescent trees) with mutations are build backward in time event by event (i.e. H_k , each time the sample configuration changes) until the MRCA is found.

- Those coalescent tree simulations (absorbing Markov chains) are used to explore the genealogy space
- The importance sampling fonction $\hat{p}(H_{k-1} | H_k)$ is used to more efficiently explore the genealogical space (i.e. more likely genealogies, with the more likely events)

the parameter space is explored independently

• the new importance sampling of Delorio & Griffiths 2004 :

$$p(H_{k}) = \sum_{H_{k-1}} w_{\hat{p}}(H_{k}, H_{k-1}) \times \hat{p}(H_{k-1} | H_{k}) \times p(H_{k-1})$$
$$p(H_{0}) = E_{IS} \begin{bmatrix} k = -m+1 = MRCA \\ \prod_{k=0}^{m} w_{IS}(H_{k}, H_{k-1}) \times p(H_{-m}) \end{bmatrix}$$

• the recurrence : $p(H_k) = \sum_{H_{k-1}} w_{IS}(H_k, H_{k-1}) \bullet \hat{p}(H_{k-1} | H_k) \bullet p(H_{k-1})$ $p(H_0) = E_{IS} \begin{bmatrix} k = -m + 1 = MRCA \\ \prod_{k=0}^{k} w_{IS}(H_k, H_{k-1}) \bullet p(H_{-m}) \end{bmatrix}$



91

- Coalescent tree building
- 1. Start with the sample configuration H_0
- Draw randomly an event among all possible events (=coa ou mig ou mut) from the IS transition probabilities
 - → new ancestral configuration H_{k-1}
- 3. compute and store the IS transition weight $w_{IS}(H_{k-1}, H_k)$
- 4. Go back to 2 until the MRCA is found



- Coalescent tree building
- 1. Start with the sample configuration H_0
- Draw randomly an event among all possible events (=coa ou mig ou mut) from the IS transition probabilities
 - → new ancestral configuration H_{k-1}
- 3. compute and store the IS transition weight $w_{IS}(H_{k-1}, H_k)$
- 4. Go back to 2 until the MRCA is found



- Coalescent tree building
- 1. Start with the sample configuration H_0
- Draw randomly an event among all possible events (=coa ou mig ou mut) from the IS transition probabilities
 - → new ancestral configuration H_{k-1}
- 3. compute and store the IS transition weight $w_{IS}(H_{k-1}, H_k)$
- 4. Go back to 2 until the MRCA is found



- Coalescent tree building
- 1. Start with the sample configuration H_0
- Draw randomly an event among all possible events (=coa ou mig ou mut) from the IS transition probabilities
 - → new ancestral configuration H_{k-1}
- 3. compute and store the IS transition weight $w_{IS}(H_{k-1}, H_k)$
- 4. Go back to 2 until the MRCA is found



- Coalescent tree building
- 1. Start with the sample configuration H_0
- Draw randomly an event among all possible events (=coa ou mig ou mut) from the IS transition probabilities
 - \rightarrow new ancestral configuration H_{k-1}
- 3. compute and store the IS transition weight $w_{IS}(H_{k-1}, H_k)$
- 4. Go back to 2 until the MRCA is found

probability of the MRCA = probability of the allelic state of the MRCA in the stationnary distribution of the mutation model

for most model it is equal to 1/K, with K the number of possible allelic state



• Probability of the sample for a given coalescent tree:

All transition weight $w_{IS}(H_{k-1}, H_k)$ were computed and stored

$$p(H_0 | G_z) = \prod_{k=0}^{k=-m+1=MRCA} w_{IS}(H_k, H_{k-1}) \times p(H_{-m})$$

$$p(H_0 | G_z) = \prod_{k=0}^{k=-m+1=MRCA} w_{IS}(H_k, H_{k-1}) \times 1/K$$

• Probability of the sample using Monte Carlo integration over a large number (*Z*) of coalescent trees:

$$p(H_0 | G_z) = \prod_{k=0}^{k=-m+1=MRCA} w_{IS}(H_k, H_{k-1}) \times 1/K$$

$$p(H_0) = E_{IS} \Big[p(H_0 | G_z) \Big]$$

$$p(H_0) \approx \frac{1}{Z} \sum_{z} p(H_0 | G_z)$$

$$L(P|D) = \frac{1}{K} \sum_{k=1}^{K} \Pr(D|G_k; P)$$

 the likelihood of the sample L(P|D)=p(H₀) is computed for many points (random or on a grid) over the parameter space and the likelihood surface is interpolated using Kriging







- ML point estimate and Confidence intervals are determined from this interpolated likelihood surface
- no convergence required!

IBD and coalescence based maximum likelihood inferences

In theory, Maximum Likelihood methods (ML) should be more powerful than moment based methods (F_{ST}) because :

- > Use all the information present in the genetic data
- Powerful maximum likelihood statistical framework
- \succ Possible to make inference on parameters other than $D\sigma^2$
 - ✓ Migration rates (Nm)
 - \checkmark Shape of the distribution
 - \checkmark Total population size
 - ✓ Mutation rate

IBD and maximum likelihood inference

Two main approaches for maximum likelihood demographic inference under the coalescent :

- <u>"MCMC"</u>: Felsenstein et al. (e.g. MsVar, IM, LAMARCK, MIGRATE)
 High computation times → difficult to test
- <u>"IS"</u>: Griffiths et al. (GENETREE, MIGRAINE) Much less developed than MCMC, simpler models Only recently : 1D and 2D IBD, One pop with variable size ("MsVar like" model)
IBD and maximum likelihood inference Griffiths et al. (IS, software MIGRAINE) IBD 1D, recent development for 2D IBD (in prep)

Likelihood and Approximate Likelihood Analyses of Genetic Structure in a Linear Habitat: Performance and Robustness to Model Mis-Specification

François Rousset* and Raphaël Leblois†

*Université, Montpellier 2, CNRS, Institut des Sciences de l'Évolution, France; and †Unité Origine, Structure et Évolution de la Biodiversité, Museum National d'Histoire Naturelle, Paris, France

Mol. Biol. Evol. 24(12):2730-2745. 2007

Demic model of IBD on a circle or on a line with absorbing boundaries

IS much faster than MCMC (10x + easy parallel computing)

Number of parameters reduced by consideration of homogeneous IBD model

1- First results under **stepping stone migration** (i.e. no middle/long distance migrants):

very good precision and robustness on Nm inference :

Rel biais = [0.04-0.12] and Rel RMSE=[0.15-0.5]

relatively good precision for Nµ

Rel biais =[0.04-0.40] and Rel RMSE=[0.25-0.8])

1- First results under **stepping stone migration** (i.e. no middle/long distance migrants):

 $N\mu$ slightly influenced by the total number of sub-populations considered in the analysis ("Ghost populations")



2- geometric dispersal distance migrants): $\frac{m}{2}(1-g)g^{|k-1|}$, (i.e. with middle/long

Large g more long distance, large $D\sigma^2 = \sigma^2 = m(1+g)/(1-g)^2$.

 $D\sigma^2$ and Nm inferences much more precise and robust than for g

large *m* and *g* (i.e. more migrants, at larger distances)

→ more influence of the ghost/unsampled pops and of the mutation process

Stronger effect for $N\mu$ and g than Nm,

not much effect for $D\sigma^2$ (compensation of different bias)

2- geometric dispersal distance migrants): $\frac{m}{2}(1-g)g^{|k-1|}$, (i.e. with middle/long

Large $g \implies$ more long distance, large $D\sigma^2$ $\sigma^2 = m(1+g)/(1-g)^2$.

 $D\sigma^2$ and Nm inferences much more precise and robust than for g

ML more accurate than moment based regression method when analyzed under the good model (i.e. nb of sub-pops and mutation processes well specified)

Hopefully the results are also very accurate for most cases with misspecifications

3- test on a real data set : the 1D damselflies data set



Not much information on *g*, because of a strong correlation with *Nm*

Lines of equal $4D\sigma^2$ values

IBD and IS inference (MIGRAINE)

3- test on a real data set : the 1D damselflies data set



 $Nb = 4D\sigma^2 \rightarrow 108 [50-220] 2N\mu = 0.04 [0.02-0.07]$ More information about Nb than Nm and g separetely More information about Nm than g.

6. IBD and IS inference (MIGRAINE)

4 - Comparison with demographic estimates and the moment based regression method on the damselflies example

	Inference of $D\sigma^2$			
	Direct	Indirect	Indirect	
	(demo)	(Regression)	(MIGRAINE)	
Site 1	277	222	108	
(1D)		[66-392]	[50-220]	



"Effective" demographic estimates are probably overestimated (not corrected for temporal variations in density)

CI obtained by the regression method overlaps widely with the one given by MLE.

6. IBD and IS inference (MIGRAINE)

4 - Comparison with demographic estimates and the moment based regression method on the damselflies example

	Inference of $D\sigma^2$			
	Direct (demo)	Indirect (Regression)	Indirect (MIGRAINE)	
Site 1 (1D)	277	222 [66-392]	108 [50-220]	



Both genetic methods may estimate (with different small-sample biases) the same effective $D\sigma^2$ and the demographic estimate may be slightly overestimated.

Further comparisons necessary to demonstrate systematic differences of this magnitude.

6. IBD and IS inference (MIGRAINE)

4 - Comparison with demographic estimates and the moment based regression method on the damselflies example

	Inference of $D\sigma^2$			
	Direct	Indirect	Indirect	
	(demo)	(Regression)	(MIGRAINE)	
Site 1	277	222	108	
(1D)		[66-392]	[50-220]	



Other possible explanations for the observed differences:

- Shape of the dispersal distribution (i.e. not geometric in reality)
- Influence of past demographic processes/fluctuations
- Mutation processes, edge effects, number of sub-populations, binning (but showed only moderate effects on simulations)

ML and IBD : Conclusions

+ Very good performances, even when the model is mis-specified

- Very slow for large network of populations (>100)
- some problems for large migration rates, long distance migration, and small population sizes (due to the coalescent approximations)
- impossible to model continuous populations (ABC methods??)
- => geographic data binning needed to deal with continuous samples
- need to test robustness to past demographic fluctuations

- not much information in "classical" samples on the shape of the dispersal distribution (i.e. inference of param. others than $D\sigma^2$ and Nm)

+ may be used for other developments (e.g. IBD between habitats, landscape genetics)

6. IBD and maximum likelihood inference??

Two main approaches for maximum likelihood demographic inference under the coalescent :

- "MCMC": Felsenstein et al. (MIGRATE) High computation times → difficult to test Bad results (simulations and comparison with demography) under IBD but MIGRATE is not especially designed for IBD
- 3. "IS": Griffiths et al. (MIGRAINE) Much less developed than MCMC, simpler models Only recently: Linear IBD (i.e. one dimension)

Preliminary results

First test of MIGRATE : comparison with demographic data

Human data : villages of New Guinea





Limited dispersal : few kilometers per generation

Demographic data : Wood et al. *Am. Nat.* 1985 Genetic data (allozymes) : Long et al. *Am. J. Phys. Anth.*1986



moment based regression method (a_r) \rightarrow inference of σ^2 : 1.4 km²/generation



Complementary tests using simulations

11 samples (•) of 20
individuals evolving on a
lattice of 40 000
(200x200) subpopulations

5 loci KAM 10 alleles Mutation rate of 5.10⁻⁴

stepping stone migration







Possible explanations...(1)

Inherent Bias of the method?

Yes

Observed on simulations by Beerli et Felsenstein (2001) : Expected bias when low number of migrants

Possible explanations...(2)

Wrong mutation model?

Number of populations

Slow convergence of MCMC?

Inherent bias of the method?

No major effects



15

10

20

25

30

Possible explanations...(3)



Not easy to solve in practice

Many possible explanations...

Inherent bias to method? Yes

Slow convergence of MCMC?

Total nb of sub-populations VS nb of sampled subpopulations? No major effects

? - not easy to solve in practice

Too many parameters to infer?

Expected to have an important effect

Very slow → difficult to test Bad precision under IBD