

Bayesian Phylogenetic Inference using Sequential Monte Carlo Algorithms

Alexandre Bouchard-Côté^{*}, Sriram Sankararaman^{*}, and Michael I. Jordan^{*,†}

^{*} Computer Science Division, University of California Berkeley

⁺ Department of Statistics, University of California Berkeley



Phylogenetic tree inference

Topic of this talk: **integration** over the space of trees using Sequential Monte Carlo (SMC)

Motivation: Bayesian approach to phylogenetic inference

- + Put a prior on trees, use the posterior for reconstruction
- Heavy use of integrals over the space of trees: e.g. for handling nuisance parameters, computing minimum risk estimators, Bayes factors, etc.

Prelude: a parallel with the simpler problem of **maximization** over the space of trees

Maximization over phylogenies

Two strategies: Local and sequential search

Key difference: representation



Maximization: *local* strategy

Meta-algorithm:

1. Start at arbitrary state



Maximization: *local* strategy



Meta-algorithm:

- 1. Start at arbitrary state
- 2. Iterate:
 - i. Evaluate neighbors
 - ii. Move to a nearby tree

Maximization: *local* strategy



Meta-algorithm:

- 1. Start at arbitrary state
- 2. Iterate:
 - i. Evaluate neighbors
 - ii. Move to a nearby tree
- 3. Return best state visited

Example: stochastic annealing

Maximization over phylogenies

Two strategies: Local and sequential search

Key difference: representation



Maximization over phylogenies

Two strategies: Local and sequential search

Key difference: representation



Maximization: sequential strategy

Meta-algorithm:

1. Start at the initial, unconstrained *partial state*

Maximization: sequential strategy



Meta-algorithm:

- 1. Start at the initial,
 - unconstrained partial state
- 2. Iterate:
 - i. Extend partial state
 - ii. Estimate best successor

Maximization: sequential strategy



Meta-algorithm:

- 1. Start at the initial,
 - unconstrained partial state
- 2. Iterate:
 - i. Extend partial state
 - ii. Estimate best successor
- 3. Return best final state

Example: neighbor joining

Parallel

Classification of phylogenetic algorithms

	Local strategy	Sequential strategy
Maximization	Stochastic annealing,	Neighbor-joining,
Integration		

Parallel

Classification of phylogenetic algorithms

	Local strategy	Sequential strategy	
Maximization	Stochastic annealing,	Neighbor-joining,	
Integration	MCMC algorithms	???	

Parallel

Classification of phylogenetic algorithms



Outline

Background: Importance sampling and Sequential Monte Carlo

SMC for phylogenetic inference

Framework for designing proposals

Experiments: comparisons with MCMC

Preview: Comparative advantages

SMC

- + Trivial to parallelize
- + Easier to get data likelihood estimate
- + No burn-in

MCMC

- + Easier to resample hyper-parameters
- + Easier to design proposal distribution

Preview: Comparative advantages

SMC

- + Trivial to parallelize
- + Easier to get data likelihood estimate
- + No burn-in

MCMC

- + Easier to resample hyper-parameters
- + Easier to design proposal distribution

Not exclusive: the two approaches can be combined







Farget distribution:
$$T|Y \stackrel{d}{=} \pi$$

with density: $\frac{\gamma(t)}{Z}$





Background: Importance Sampling (IS)



Background: Importance Sampling (IS)



IS : Approximation for π

- 1. Sample trees from a proposal $q: t_i \sim q$
- 2. Compute weights

$$w_i = \gamma(t_i) / q(t_i)$$

3. Normalize weights

Background: Importance Sampling (IS)



Background:

Problem with importance sampling: π is high-dimensional



Background: Importance Sampling



Background: Importance Sampling

$$\int \frac{q}{\pi} \approx \frac{\gamma}{Z}$$

SMC: a sequence of proposals



Background: Importance Sampling

$$\int \frac{q}{\pi} \approx \frac{\gamma}{Z}$$

SMC: a sequence of proposals



SMC for phylogenies: π_r are distributions over partial states (forest)



Monday, July 5, 2010



SMC : Approximation for $\,\pi\,$



```
2. Iterate :
```

i. Sample partial states







 π_2







Basic result: SMC is asymptotically consistent



Basic result: SMC is asymptotically consistent



Basic result: SMC is asymptotically consistent



Compare:

Weights along a SMC path

$$w \cdot w' \cdot w'' = \frac{\gamma(p'')}{q(\bot \to p'')}$$

Importance sampling

$$w = \frac{\gamma(t)}{q(t)}$$

Issue: Over-counting



Useful abstraction: *q* induce a partial order (poset) *P*



Useful abstraction: q induce a partial order (poset) P



Useful abstraction: *q* induce a partial order (poset) *P*

Poset's Hesse diagram:



Useful abstraction: *q* induce a partial order (poset) *P*

Poset's Hesse diagram:



Example: a proposal that has a tree-shaped Hesse diagram.

- 1. Pick a pair of trees to merge uniformly at random
- 2. Pick a height for the new tree such that



Experiments: setup

	Synthetic-small	Synthetic-med	Real data
Source	Generated from the model		Subset of HGDP
Likelihood model	Brownian motion on frequencies		
Number of sites	100		11,511
Number of nodes	25	51	25
Number of leaves	13	26	13

Synthetic experiments

Goal: comparison against MCMC

Competitor: standard MCMC sampler, 4 tempering chains, shared sum-product implementation

Metric: symmetric clade difference of the Minimum Bayes Risk reconstructed tree to the generating tree

Datapoints computed by increasing the number of particles (for SMC) and the number of sampling steps (for MCMC)

Comparison with MCMC

Synthetic-small



Comparison with MCMC

Synthetic-medium



Experiments on real data

Goal: show that the method scales to large number of sites

Number of particle (10,000) determined using synthetic experiments, timing experiments with different numbers of cores:





- **SMC can be applied to a wide range of phylogenetic models;** previous work limited to Coalescent priors [Teh et al. 07]
- Order theoretic framework for designing proposals
- **Experiments:** There are regimes where SMC outperforms MCMC
- **Promising applications** of SMC in phylogenetic inference:
 - 1. Quickly analyze large datasets
 - 2. Initialization and large step proposal for MCMC chains