

BIOS 731

Advanced Statistical Computing

Fall 2020

Lecture 13

Applications of MCMC and SMC

Steve Qin

Review

- Gibbs sampler
- Grouping and collapsing
- Convergence check
- Sequential Monte Carlo
 - Acceptance rejection method
 - Importance sampling

Importance sampling

- *Importance sampling:*

to evaluate $E_f[h(X)] = \int h(x)f(x)dx$

based on generating a sample X_1, \dots, X_n from a given distribution g and approximating

$$E_f[h(X)] \approx \frac{1}{m} \sum_{j=1}^m \frac{f(X_j)}{g(X_j)} h(X_j)$$

which is based on

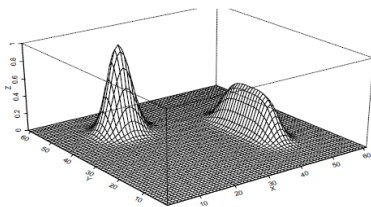
$$E_f[h(X)] = \int h(x) \frac{f(x)}{g(x)} g(x) dx$$

3

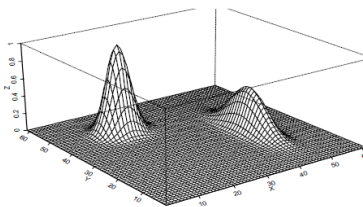
Another example

$$f(x, y) = 0.5e^{-90(x-0.5)^2 - 45(y+0.1)^4} + e^{-45(x+0.4)^2 - 60(y-0.5)^2}$$

(a)



(b)



- Both grid-point method and vanilla Monte Carlo methods wasted resources on “boring” desert area.

4

Another example

- Use proposal function

$$g(x, y) \propto 0.5e^{-90(x-0.5)^2-10(y+0.1)^2} + e^{-45(x+0.4)^2-60(y-0.5)^2},$$

with $(x, y) \in [-1, 1] \times [-1, 1]$, a truncated mixture of bivariate Gaussian

$$0.46\mathcal{N}\left[\begin{pmatrix} 0.5 \\ -0.1 \end{pmatrix}, \begin{pmatrix} \frac{1}{180} & 0 \\ 0 & \frac{1}{20} \end{pmatrix}\right] + 0.54\mathcal{N}\left[\begin{pmatrix} -0.4 \\ 0.5 \end{pmatrix}, \begin{pmatrix} \frac{1}{90} & 0 \\ 0 & \frac{1}{120} \end{pmatrix}\right]$$

Vanilla Monte Carlo

$$\hat{\mu} = 0.1307$$

$$\text{std}(\hat{\mu}) = 0.009$$

Importance Sampling

$$\hat{\mu} = 0.1259$$

$$\text{std}(\hat{\mu}) = 0.0005$$

5

Sequential importance sampling

- For high dimensional problem, how to design trial distribution is challenging.
- Suppose the target density of $\mathbf{x} = (x_1, x_2, \dots, x_d)$ can be decomposed as

$$\pi(\mathbf{x}) = \pi(x_1)\pi(x_2 | x_1) \cdots \pi(x_d | x_1, \dots, x_{d-1})$$

then constructed trial density as

$$g(\mathbf{x}) = g_1(x_1)g_2(x_2 | x_1) \cdots g_d(x_d | x_1, \dots, x_{d-1})$$

6

Sequential importance sampling

$$w(\mathbf{x}) = \frac{\pi(x_1)\pi(x_2 | x_1) \cdots \pi(x_d | x_1, \dots, x_{d-1})}{g_1(x_1)g_2(x_2 | x_1) \cdots g_d(x_d | x_1, \dots, x_{d-1})}$$

Suggest a recursive way of computing and monitoring importance weight. Denote

$$\mathbf{x}_t = (x_1, x_2, \dots, x_t)$$

then we have

$$w_t(\mathbf{x}_t) = w_{t-1}(\mathbf{x}_{t-1}) \frac{\pi(x_t | \mathbf{x}_{t-1})}{g_t(x_t | \mathbf{x}_{t-1})}$$

7

Sequential importance sampling

- Advantages of the recursion scheme
 - Can stop generating further components of \mathbf{x} if the partial weight is too small.
 - Can take advantage of $\pi(x_t | \mathbf{x}_{t-1})$ in designing $g_t(x_t | \mathbf{x}_{t-1})$
- However, the scheme is impractical since requires the knowledge of marginal distribution $\pi(\mathbf{x}_t)$.

8

Sequential importance sampling

- Add another layer of complexity:
- Introduce a sequence of “auxiliary distributions” $\pi_1(x_1)\pi_2(\mathbf{x}_2)\pi_d(\mathbf{x})$ such that $\pi_t(\mathbf{x}_t)$ is a reasonable approximation of the marginal distribution $\pi(\mathbf{x}_t)$, for $t = 1, \dots, d-1$ and $\pi_d = \pi$.
- Note the π_d are only required to be known up to a normalizing constant.

9

The SIS procedure

For $t = 2, \dots, d$,

- Draw $X_t = x_t$ from $g_t(x_t | x_{t-1})$, and let $\mathbf{x}_t = (\mathbf{x}_{t-1}, x_t)$
- Compute $u_t = \frac{\pi_t(\mathbf{x}_t)}{\pi_{t-1}(\mathbf{x}_{t-1}) g_t(x_t | \mathbf{x}_{t-1})}$ and let $w_t = w_{t-1} u_t$
- u_t : incremental weight.
- The key idea is to break a difficult task into manageable pieces.
- If w_t is getting too small, reject.

10

An application example of SIS

- Assume
 - Constant population size N ,
 - Evolve in non-overlapping generation,
 - The chromosomal region is sufficiently small,
 - No recombination,
 - “haplotype”: each chromosome only has one parent.

11

Population genetics example

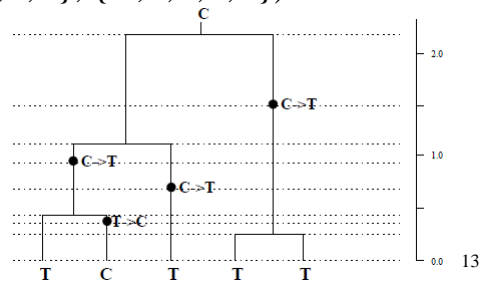
- Notation:
 - E : set of all possible genetic types,
 - μ : mutation rate per chromosome per generation,
 - $P = (P_{\alpha\beta})$: the mutation transition matrix,
 - If a parental segment of type $\alpha \in E$,

its progeny is $\begin{cases} \alpha & \text{with prob. } 1 - \mu, \\ \beta & \text{with prob. } \mu P_{\alpha\beta}. \end{cases}$

12

Example data

- From Stephens and Donnelly (2000)
- $E = \{C, T\}$
- The history $H = (H_{-k}, H_{-(k-1)}, \dots, H_{-1}, H_0)$
 $= (\{C\}, \{C, C\}, \{C, T\}, \{C, C, T\}, \{C, T, T\}, \{T, T, T\},$
 $\{T, T, T, T\}, \{C, T, T, T, T\}, \{C, T, T, T, T\})$

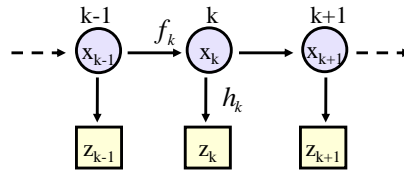


Particle filter

Michael Rubinstein

IDC

Dynamic System



State equation: $x_k = f_k(x_{k-1}, v_k)$

x_k state vector at time instant k
 f_k state transition function, $f_k : R^{N_x} \times R^{N_v} \rightarrow R^{N_x}$
 v_k i.i.d process noise

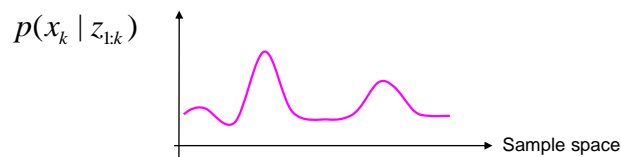
Stochastic diffusion

Observation equation: $z_k = h_k(x_k, w_k)$

z_k observations at time instant k
 h_k observation function, $h_k : R^{N_x} \times R^{N_w} \rightarrow R^{N_z}$
 w_k i.i.d measurement noise

© Michael Rubinstein

Recursive Bayes filter



- Prediction:

$$p(x_k | z_{1:k-1}) = \int p(x_k | x_{k-1}) p(x_{k-1} | z_{1:k-1}) dx_{k-1}$$

(1)

- Update:

$$p(x_k | z_{1:k}) = \frac{p(z_k | x_k) p(x_k | z_{1:k-1})}{p(z_k | z_{1:k-1})}$$

(2)

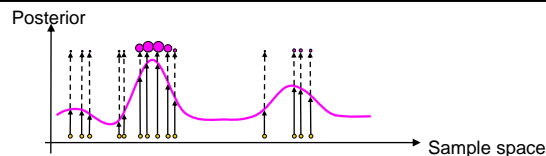
$$p(z_k | z_{1:k-1}) = \int p(z_k | x_k) p(x_k | z_{1:k-1}) dx_k$$

© Michael Rubinstein

Particle filtering

- Many variations, one general concept:

Represent the posterior pdf by a set of randomly chosen weighted samples (particles)



- Randomly Chosen = Monte Carlo (MC)
- As the number of samples become very large – the characterization becomes an equivalent representation of the true pdf

© Michael Rubinstein

Particle filtering

- Compared to methods we've mentioned last time
 - Can represent any arbitrary distribution
 - multimodal support
 - Keep track of many hypotheses as there are particles
 - **Approximate representation of complex model rather than exact representation of simplified model**
- The basic building-block: *Importance Sampling*

© Michael Rubinstein

Monte Carlo integration

- Evaluate complex integrals using probabilistic techniques
- Assume we are trying to estimate a complicated integral of a function f over some domain D :

$$F = \int_D f(\vec{x}) d\vec{x}$$

- Also assume there exists some PDF p defined over D

© Michael Rubinstein

Monte Carlo integration

- Then

$$F = \int_D f(\vec{x}) d\vec{x} = \int_D \frac{f(\vec{x})}{p(\vec{x})} p(\vec{x}) d\vec{x}$$

- But

$$\int_D \frac{f(\vec{x})}{p(\vec{x})} p(\vec{x}) d\vec{x} = E \left[\frac{f(\vec{x})}{p(\vec{x})} \right], x \sim p$$

- This is true for any PDF p over D !

© Michael Rubinstein

Monte Carlo integration

- Now, if we have i.i.d random samples $\vec{x}_1, \dots, \vec{x}_N$ sampled from p , then we can approximate

$E\left[\frac{f(\vec{x})}{p(\vec{x})}\right]$ by

$$F_N = \frac{1}{N} \sum_{i=1}^N \frac{f(\vec{x}_i)}{p(\vec{x}_i)}$$

- Guaranteed by law of large numbers:

$$N \rightarrow \infty, F_N \xrightarrow{a.s.} E\left[\frac{f(\vec{x})}{p(\vec{x})}\right] = F$$

© Michael Rubinstein

Importance Sampling (IS)

- What about $p(\vec{x})=0$?
- If p is very small, f/p can be arbitrarily large, ‘damaging’ the average
 - Design p such that f/p is bounded
 - Rule of thumb: take p similar to f as possible
- The effect: get more samples in ‘important’ areas of f , i.e. where f is large

© Michael Rubinstein

IS for Bayesian estimation

- We draw the samples from the importance density $q(x_{0:k} | z_{1:k})$ with importance weights

$$w_k^i \propto \frac{p(x_{0:k} | z_{1:k})}{q(x_{0:k} | z_{1:k})}$$

- Sequential update (after some calculation...)

Particle update

$$x_k^i \sim q(x_k | x_{k-1}^i, z_k)$$

Weight update

$$w_k^i = w_{k-1}^i \frac{p(z_k | x_k^i) p(x_k^i | x_{k-1}^i)}{q(x_k^i | x_{k-1}^i, z_k)}$$

© Michael Rubinstein

Sequential Importance Sampling (SIS)

$$\left[\{x_k^i, w_k^i\}_{i=1}^N \right] = \text{SIS} \left[\{x_{k-1}^i, w_{k-1}^i\}_{i=1}^N, z_k \right]$$

- FOR $i=1:N$

– Draw $x_k^i \sim q(x_k | x_{k-1}^i, z_k)$

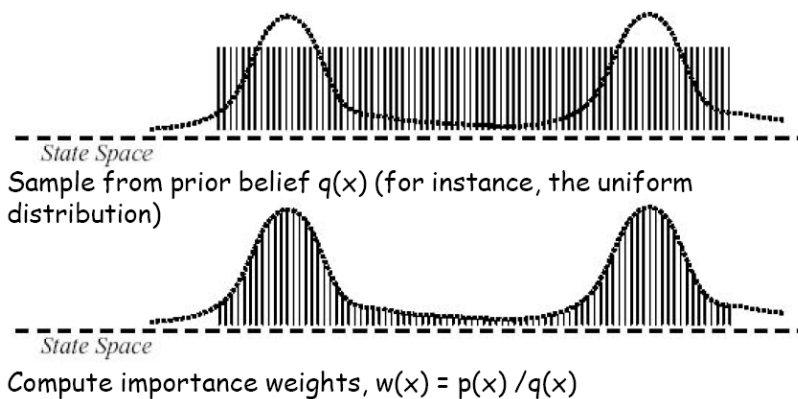
– Update weights $w_k^i = w_{k-1}^i \frac{p(z_k | x_k^i) p(x_k^i | x_{k-1}^i)}{q(x_k^i | x_{k-1}^i, z_k)}$

- END

- Normalize weights

© Michael Rubinstein

Choice of importance density



Hsiao et al.

© Michael Rubinstein

Choice of importance density

- Most common (suboptimal): the transitional prior

$$q(x_k | x_{k-1}^i, z_k) = p(x_k | x_{k-1}^i)$$

$$\Rightarrow w_k^i = w_{k-1}^i \frac{p(z_k | x_k^i) p(x_k^i | x_{k-1}^i)}{q(x_k^i | x_{k-1}^i, z_k)} = w_{k-1}^i p(z_k | x_k^i)$$

Grid filter weight update:

$$w_{k|k}^i = \frac{w_{k|k-1}^i p(z_k | x_k^i)}{\sum_{j=1}^{N_s} w_{k|k-1}^j p(z_k | x_k^j)}$$

© Michael Rubinstein

The degeneracy phenomenon

- Unavoidable problem with SIS: after a few iterations most particles have negligible weights
 - Large computational effort for updating particles with very small contribution to $p(x_k | z_{1:k})$
- Measure of degeneracy - the effective sample size:

$$N_{eff} = \frac{1}{\sum_{i=1}^N (w_k^i)^2}$$

- Uniform: $N_{eff} = N$, severe degeneracy: $N_{eff} = 1$

© Michael Rubinstein

Resampling

- The idea: when degeneracy is above some threshold, eliminate particles with low importance weights and multiply particles with high importance weights

$$\{x_k^i, w_k^i\}_{i=1}^N \rightarrow \{x_k^{i*}, \frac{1}{N}\}_{i=1}^N$$

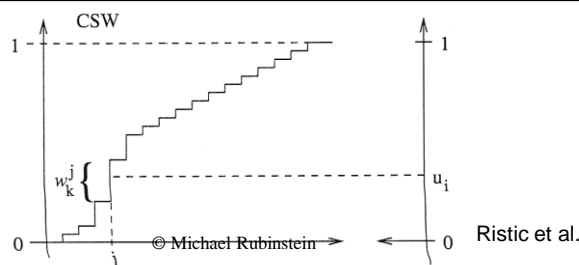
- The new set is generated by sampling with replacement from the discrete representation of $p(x_k | z_{1:k})$ such that $\Pr\{x_k^{i*} = x_k^j\} = w_k^j$

© Michael Rubinstein

Resampling

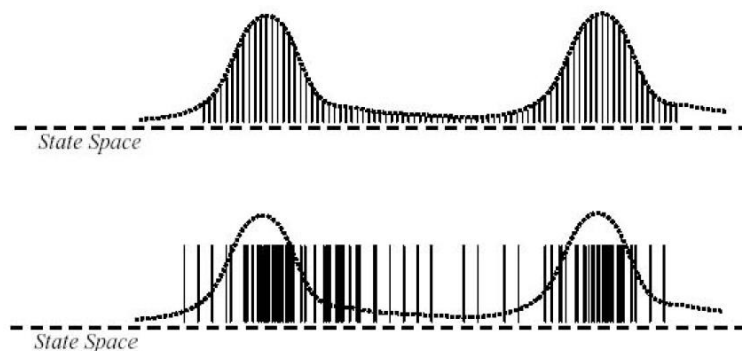
$$[\{x_k^{i*}, w_k^i\}_{i=1}^N] = \text{RESAMPLE}[\{x_k^i, w_k^i\}_{i=1}^N]$$

- Generate N i.i.d variables $u_i \sim U[0,1]$
- Sort them in ascending order
- Compare them with the cumulative sum of normalized weights



Resampling

- Complexity: $O(N \log N)$
 - $O(N)$ sampling algorithms exist



© Michael Rubinstein

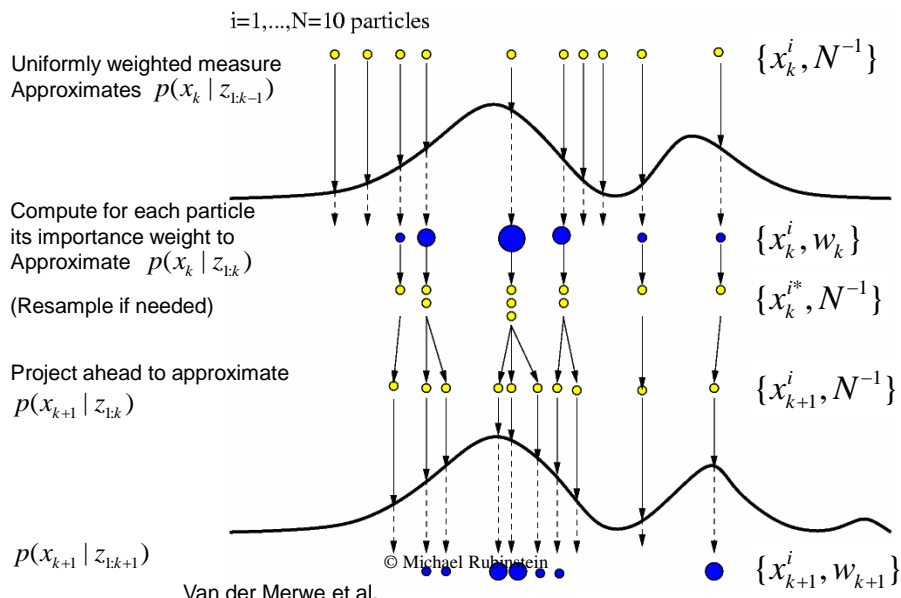
Hsiao et al.

Generic PF

- $$\left[\{x_k^i, w_k^i\}_{i=1}^N \right] = \text{PF} \left[\{x_{k-1}^i, w_{k-1}^i\}_{i=1}^N, z_k \right]$$
- Apply SIS filtering $\left[\{x_k^i, w_k^i\}_{i=1}^N \right] = \text{SIS} \left[\{x_{k-1}^i, w_{k-1}^i\}_{i=1}^N, z_k \right]$
 - Calculate N_{eff}
 - IF $N_{\text{eff}} < N_{\text{thr}}$
 - $\left[\{x_k^i, w_k^i\}_{i=1}^N \right] = \text{RESAMPLE} \left[\{x_k^i, w_k^i\}_{i=1}^N \right]$
 - END

© Michael Rubinstein

Generic PF



PF variants

- Sampling Importance Resampling (SIR)
- Auxiliary Sampling Importance Resampling (ASIR)
- Regularized Particle Filter (RPF)
- Local-linearization particle filters
- Multiple models particle filters (maneuvering targets)
- ...

© Michael Rubinstein

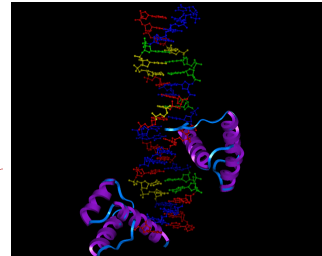
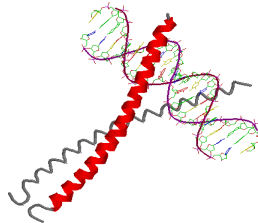
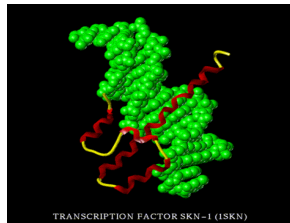
Sampling Importance Resampling (SIR)

- A.K.A Bootstrap filter, Condensation

- **Initialize** $\{x_0^i, w_0^i\}_{i=1}^N$ from prior distribution X_0
- For $k > 0$ do
 - **Resample** $\{x_{k-1}^i, w_{k-1}^i\}_{i=1}^N$ into $\{x_{k-1}^{i*}, \frac{1}{N}\}_{i=1}^N$
 - **Predict** $x_k^i \sim p(x_k | x_{k-1} = x_{k-1}^{i*})$
 - **Reweight** $w_k^i = p(z_k | x_k = x_k^i)$
 - **Normalize** weights
 - **Estimate** \hat{x}_k (for display)

© Michael Rubinstein

Appliation: Transcription Factor Binding Sites Discovery



35

Example: cyclic receptor protein (CRP)

cole1	t a a t g t t t g t g c t g g t t t t t g t g g c a t c g g g c g a g a a t a g c g c g t g g t g t g a a a g a c t g t t t t t g a t c g t t t t c a c a a a a t g g a g t c c a c a g t c t t g a c a g
ecoarabop	g a c a a a a a c g c g t a a c a a a a g t g t c t a t a a t c a c g g c a g a a a a g t c c a c a t t g a t t a t t g c a c g g c g t c a c a c t t g c t a t g c c a t a g c a t t t t a t c a t a a g
ecobglr1	a c a a a t c c c a a t a a c t t a a t t a t t g g g a t t g t t a t a t a a c t t t a t a a t t c c t a a a a t t a c a c a a g t t a a t a a c t g t g a g c a t g g t c a t a t t t t a t c a a t
ecocrp	c a c a a a g c g a a a g c t a t g c t a a a c a g t c a g g a t g c t a c a g t a a t a c a t t g a t g t a c t g c a t g t a t g c a a a g g a c g t c a c a t t a c c t g c a g t a c a g t g a t a g c
ecocya	a c g g t g c t a c a c t t g t a g t a g c g c a t c t t t c t t a c g g t c a a t c a g c a t g g t g t t a a a t t g a t c a c g t t t t a g a c c a t t t t t t c g t c g t g a a c t a a a a a a c c
ecodecop	a g t g a a t t a t t t g a a c c a g a t c g c a t t a c a g t g a t g c a a a c t t g t a a g t a g a t t t c c t t a a t t g t g a t g t a t c g a a g t g t g t g c g g g a g t a g a t g t t a g a a t a
ecogale	g c g c a t a a a a a c g g c t a a a t t c t t g t g t a a a c g a t t c c a c t a a t t t a t t c c a t g t c a c a c t t t t c g c a t c t t t g t t a t g c t a t g g t t a t t t c a c c a t a a g c c
ecoilvbpr	g c t c c g g g c g g g g t t t t t t g t t a t c t g c a a t t c a g t c a c a a a a c g t g a t c a a c c c c t c a a t t t t c c t t t g t c g a a a a a t t t c c a t t g t c t c c c t g t a a a g c t g t
ecolac	a a c g c a a t t a a t g t g a g t a g c t c a c t c a t t a g g c a c c c c a g g c t t t a c a c t t t a t g t t c c g g c t c g t a t g t t g t g a a t t g t g a g c g g a t a c a a t t t c a c
ecomale	a c a t t a c c g c c a a t t c t g t a a c a g a g a t c a c a c a a a g c a c g g t g g g g c g t a g g g g c a a g g a g a t g g a a a g a g t t g c c g t a t a a a a a c a t a g a g t c c g t t a
ecomalk	g g a g g a g c g g g g g g a t g a g a a c c g g c t t c t g t a a c t a a a c c g a g g t c a t g t a a g g a a t t t c g t a g t g t g t g c a a a a a t c g t g g c g a t t t a t g t g c a
ecomalt	g a t c a g c g t c g t t t t a g g t g a g t g t t a a t a a g a t t t g g a a t t g t g a c a c a g t g c a a t t c a g a c a t a a a a a a a a a a g t c a t c g t g c a t t a g a a a g t t t c
ecoempa	g c t g a c a a a a a a g a t t a a a c a t a c c t t a t a c a a g a c t t t t t t t c a t a t g c c t g a c g a g t t c a c a c t t g t a a g t t t t c a a c t a c g t t g t a g a c t t t a c a t c g c
ecotnaa	t t t t t t a a c a t a a a a t c t t a c g t a a t t t a a t c t t t a a a a a a g c a t t a a t t g t c c c c g a a c g a t t g t g a t t c g a t t c a c a t t t a a c a a t t t c a
ecouxu1	c c c a t g a g a g a g a a t g t t t g t g a t g g t t a c c c a a t t a g a a t t c g g g a t t g a c a t g t t t a c c a a a a g g t a g a a c t a t a c g c c a t c t a t c c g a t g c a a g c
pbr-p4	c t g g c t a a c t a t g c g c a t c a g a g c a g a t t t a c t g a g a g t g c a c c a t a t g c g g t g t g a a t a c c g c a g a t g c g t a g g a g a a a a c c g c a t c a g g c g c t c
tru9cat	c t g t g a c g g a a g a c a c t t c g c a g a a t a a a t c c t g g t g t c c c t g t t g a t a c c g g g a a g c c t g g g c c a a c t t t g g c g a a a t g a g a c g t t g a t c g g c a c g
(tdc)	g a t t t t t a c t a t t a a c t t g t g t a t t t a a a g g t a t t a a t t g t a a t a a c g a t a c t c g g a a g t a t t g a a a g t a a t t t t g a g t g g t c g c a c a t c c t g a t t

36

Stormo and Hartzell, 1989

Motif identification model

a_1
 aaaggtcga ^{a_1} gtagctactcga ^{a_2} tcgatactagcaatcgttaccctagctcgatcgaaa
 acgtgagatcagctatgaccga ^{a_2} tagctactcga ^{a_3} tataaccg
 gaa ^{a_3} tagctactcga ^{a_4} tcgatactagcaatcgttaccctagctcgatcgagatggaaa
 ...
 acgtgagatcagctatcgatcgattga ^{a_l} taactactcgtacgtat

Alignment variable $A = \{a_1, a_2, \dots, a_J\}$

39

Posterior distributions

- The posterior conditional distribution for alignment variable A

$$p(a_j = l \mid \theta_0, \boldsymbol{\theta}, \mathbf{R}_j, A_{-j}) \propto \prod_{k=1}^4 \theta_{0k}^{h_k(\mathbf{R}_j)} \prod_{i=1}^w \prod_{k=1}^4 \left(\frac{\theta_{ik}}{\theta_{0k}} \right)^{h_k(r_{j,l+i-1})} \propto \prod_{i=1}^w \prod_{k=1}^4 \left(\frac{\theta_{ik}}{\theta_{0k}} \right)^{h_k(r_{j,l+i-1})}$$

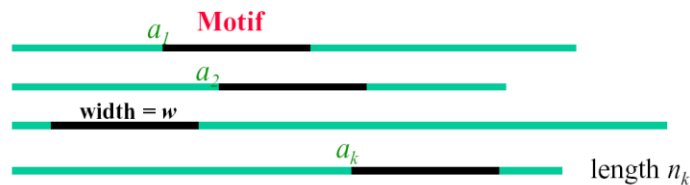
DNA sequence data

$$\mathbf{R} = (\mathbf{R}_1, \dots, \mathbf{R}_J)$$

Lawrence *et al.* *Science* 1993, Liu *et al.* *JASA* 1995

40

Motif Alignment Model



The missing data: Alignment variable: $A = \{a_1, a_2, \dots, a_k\}$

- Every **non-site positions** follows a common multinomial with $\mathbf{p}_0 = (p_{0,1}, \dots, p_{0,20})$
- Every position i in the motif element follows probability distribution $\mathbf{p}_i = (p_{i,1}, \dots, p_{i,20})$

41

Statistical Model

- **Objects:**
 - Seq: sequence data to search for motif
 - θ_0 : non-motif (genome background) probability
 - θ : motif probability matrix parameter
 - π : site locations
- **Problem:** $P(\theta, \pi \mid \text{seq}, \theta_0)$
- **Approach:** alternately estimate
 - π by $P(\pi \mid \theta, \text{seq}, \theta_0)$
 - θ by $P(\theta \mid \pi, \text{seq}, \theta_0)$

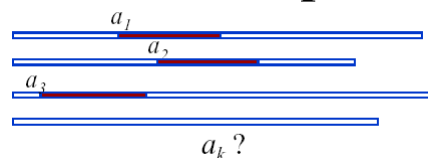
42

The Algorithm

- Initialize by choosing random starting positions
- Iterate the following steps many times;
 - Randomly or systematically choose a sequence to exclude
 - Carry out the predictive-updating step to update the starting position
 - Stop when no more observable changes in likelihood.

43

The Predictive Updating Step



- Compute predictive frequencies of each position i in motif
 - c_{ij} = count of amino acid type j at position i .
 - c_{0j} = count of amino acid type j in all non-site positions.
 - $q_{ij} = (c_{ij} + b_j) / (K - I + B)$, $B = b_1 + \dots + b_K$ "pseudo-counts"
- Sample from the predictive distribution of a_k

$$P(a_k = l + 1) \propto \prod_{i=1}^w \frac{q_{l, R_k(l+i)}}{q_{0, R_k(l+i)}} \quad 44$$

References

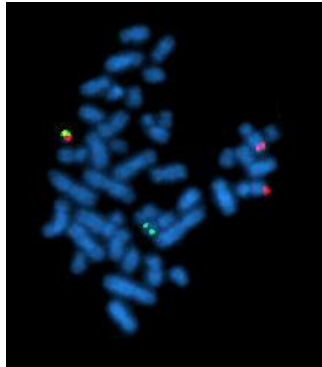
- Lawrence et al. (1993) *Science*.
- Liu, Neuwald and Lawrence (1995) *JASA*.
- Liu and Lawrence (1999) *Bioinformatics*.

45

Infer the 3D shape of
chromosomes

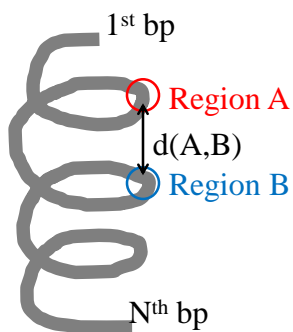
Microscopic Methods

- Fluorescent *in situ* hybridization (**FISH**)



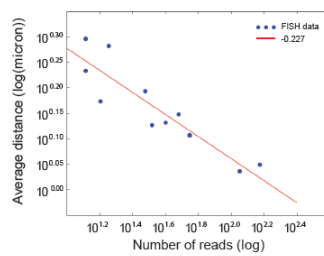
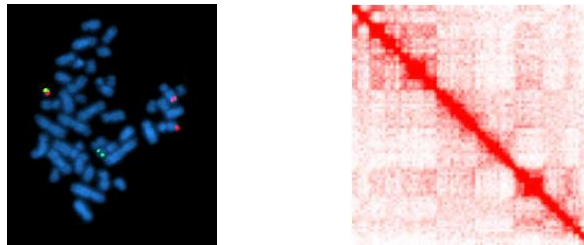
<http://en.wikipedia.org/wiki/Cytogenetics>⁴⁷

FISH Data Representation



3D chromosomal
structure

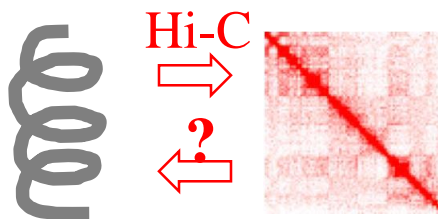
Contact Frequency vs. Spatial Distance



49

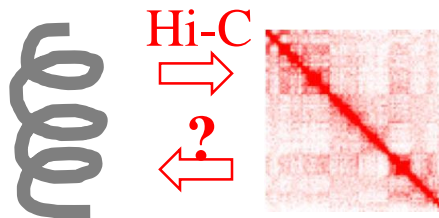
Lieberman-Aiden, et al, 2009

Problem setting



50

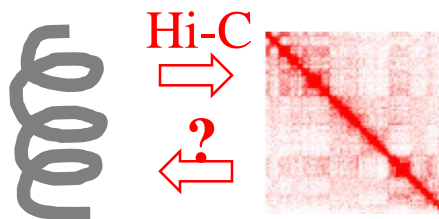
Problem setting



- Challenges:
 - Sequencing uncertainties
 - Biases: enzyme, GC content, mappability

51

Problem setting



- Challenges:
 - Sequencing uncertainties
 - Biases: enzyme, GC content, mappability

52

Yaffe and Tanay, 2011

Beads-on-a-string Representation

ACGTAGCTAGATACTGTAGTGTAGTTTGGAACTGAGGG

53

Beads-on-a-string Representation

ACGTAGCTAGATACTGTAGTGTAGTTTGGAACTGAGGG

54

Beads-on-a-string Representation

ACGTAGCTAG ATACTGTAGT GTAGTTTGGA ACCTGAGGG

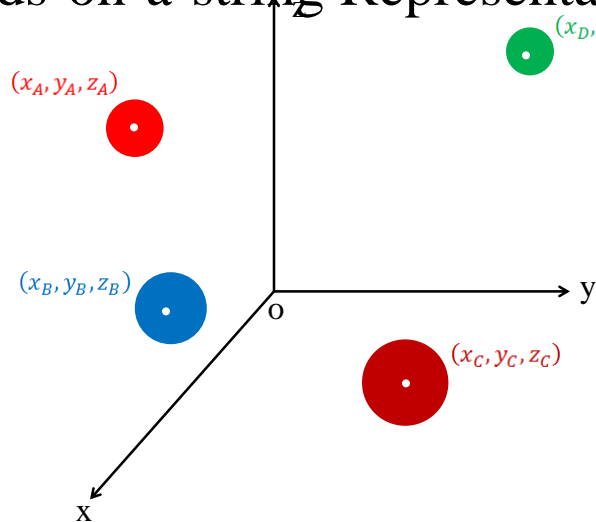
55

Beads-on-a-string Representation



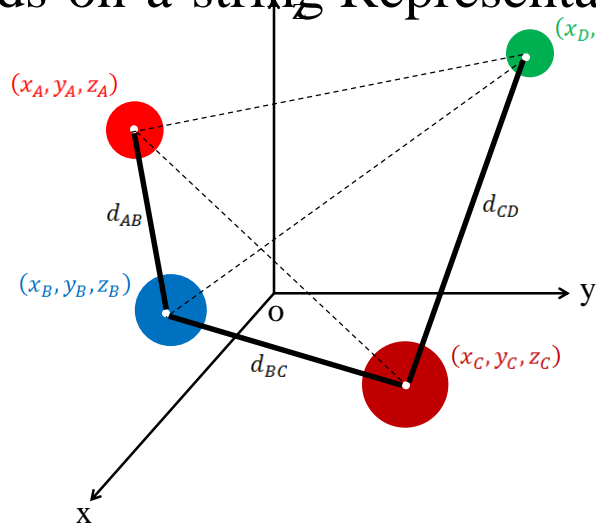
56

Beads-on-a-string Representation



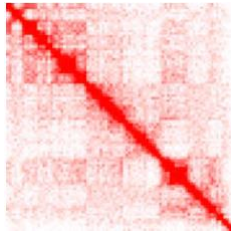
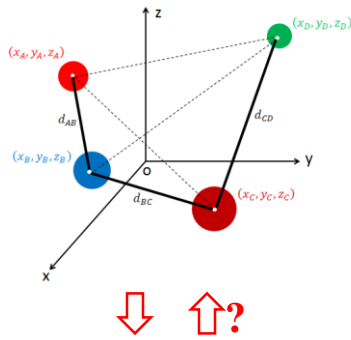
57

Beads-on-a-string Representation



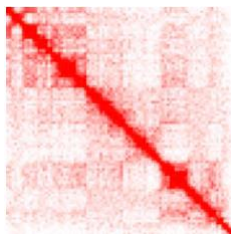
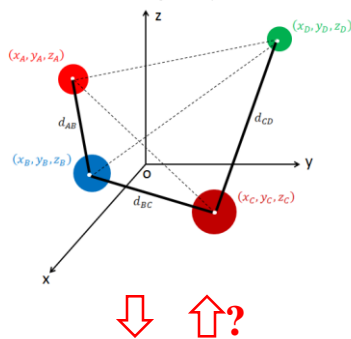
58

Bayesian Statistical Model



59

Bayesian Statistical Model

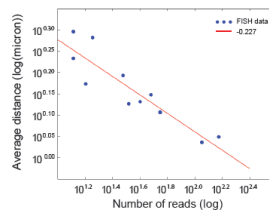


- u_{ij} : # of reads between loci i and j
- (x_i, y_i, z_i) : Euclidian coordinates of locus i
- d_{ij} : **spatial distance** between loci i and j
- e_i : # of enzyme cut site in locus i
- g_i : GC content of locus i
- m_i : mappability of locus i

Hi-C read counts: population summation

$$u_{ij} \sim \text{Poisson}(\theta_{ij})$$

Hi-C read counts vs. spatial distance: log-log linear



$$\log(\theta_{ij}) = \beta_0 + \beta_1 \log(d_{ij}) + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$$

60

Lieberman-Aiden, et al, 2009

Bayesian Statistical Model

- Likelihood:

$$L(u_{ij}, 1 \leq i < j \leq N | x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) = \prod_{1 \leq i < j \leq N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}}}{u_{ij}!}$$

$$\begin{aligned} \log(\theta_{ij}) = & \beta_0 + \beta_1 \log \left(\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \right) \\ & + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j) \end{aligned}$$

61

Bayesian Statistical Model

- Likelihood: $\binom{N}{2}$ data points, $3N + 5$ parameters

$$L(u_{ij}, 1 \leq i < j \leq N | x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) = \prod_{1 \leq i < j \leq N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}}}{u_{ij}!}$$

$$\begin{aligned} \log(\theta_{ij}) = & \beta_0 + \beta_1 \log \left(\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \right) \\ & + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j) \end{aligned}$$

62

Bayesian Statistical Model

- Likelihood: $\binom{N}{2}$ data points, $3N + 5$ parameters

$$L(u_{ij}, 1 \leq i < j \leq N | x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) = \prod_{1 \leq i < j \leq N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}}}{u_{ij}!}$$

$$\begin{aligned} \log(\theta_{ij}) = & \beta_0 + \beta_1 \log \left(\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2} \right) \\ & + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j) \end{aligned}$$

- Posterior distribution

$$\begin{aligned} & \pi(x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m | u_{ij}, 1 \leq i < j \leq N) \\ & \propto L(u_{ij}, 1 \leq i < j \leq N | x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) \text{prior} \end{aligned}$$

63

Statistical Inference

- Algorithm: **B**ayesian 3D **c**onstructor for **H**i-C data (**BACH**)

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m | u_{ij}, 1 \leq i < j \leq N)$$

64

Statistical Inference

- Algorithm: Bayesian 3D constructor for Hi-C data (BACH)

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m | u_{ij}, 1 \leq i < j \leq N)$$

- Initialization 1: use Poisson regression to obtain the initial values for

$\beta_0, \beta_e, \beta_g, \beta_m$. Set $\beta_1 = -1$.

$$u_{ij} \sim \text{Poisson}(\theta_{ij}) \quad \log(\theta_{ij}) = \beta_0 + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$$

65

Statistical Inference

- Algorithm: Bayesian 3D constructor for Hi-C data (BACH)

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m | u_{ij}, 1 \leq i < j \leq N)$$

- Initialization 1: use Poisson regression to obtain the initial values for

$\beta_0, \beta_e, \beta_g, \beta_m$. Set $\beta_1 = -1$.

$$u_{ij} \sim \text{Poisson}(\theta_{ij}) \quad \log(\theta_{ij}) = \beta_0 + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$$

- Initialization 2: use sequential important sampling to get the initial 3D chromosomal structure $\{x_i, y_i, z_i, 1 \leq i \leq N\}$.

66

Statistical Inference

- Algorithm: Bayesian 3D constructor for Hi-C data (BACH)

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m | u_{ij}, 1 \leq i < j \leq N)$$

- Initialization 1: use Poisson regression to obtain the initial values for $\beta_0, \beta_e, \beta_g, \beta_m$. Set $\beta_1 = -1$.

$$u_{ij} \sim \text{Poisson}(\theta_{ij}) \quad \log(\theta_{ij}) = \beta_0 + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$$

- Initialization 2: use sequential important sampling to get the initial 3D chromosomal structure $\{x_i, y_i, z_i, 1 \leq i \leq N\}$.
- Refinement: use Gibbs sampler with hybrid Monte Carlo to refine the initial values for parameters.

67

SIS in BACH: Outline

- Goal: use sequential importance sampling to sequentially put N loci into 3D space, i.e. sample from:

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N | u_{ij}, 1 \leq i < j \leq N)$$

68

SIS in BACH: Outline

- Goal: use sequential importance sampling to **sequentially** put N loci into 3D space, i.e. sample from:

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N | u_{ij}, 1 \leq i < j \leq N)$$

- Bridging distributions:

$$\pi_t(x_i, y_i, z_i, 1 \leq i \leq t | u_{ij}, 1 \leq i < j \leq t)$$

69

SIS in BACH: Outline

- Goal: use sequential importance sampling to **sequentially** put N loci into 3D space, i.e. sample from:

$$\pi(x_i, y_i, z_i, 1 \leq i \leq N | u_{ij}, 1 \leq i < j \leq N)$$

- Bridging distributions:

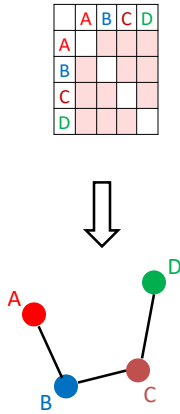
$$\pi_t(x_i, y_i, z_i, 1 \leq i \leq t | u_{ij}, 1 \leq i < j \leq t)$$

- Proposal distributions (given the first $t-1$ loci, put the t th locus in to 3D space):

$$g_t(x_t, y_t, z_t | x_i, y_i, z_i, 1 \leq i \leq t-1, u_{ij}, 1 \leq i < j \leq t)$$

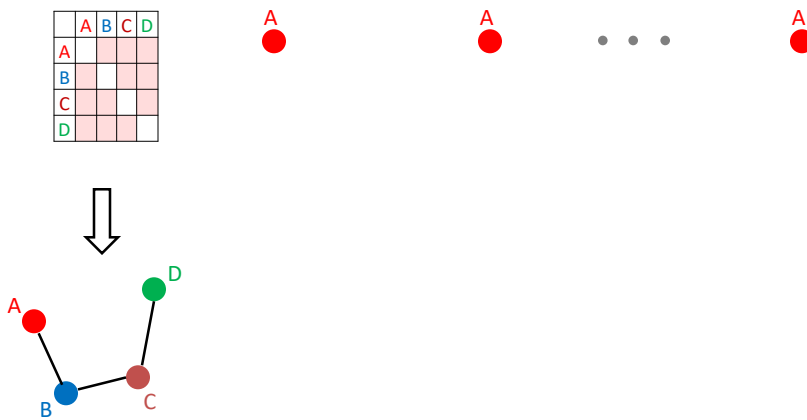
70

SIS in BACH: Illustration



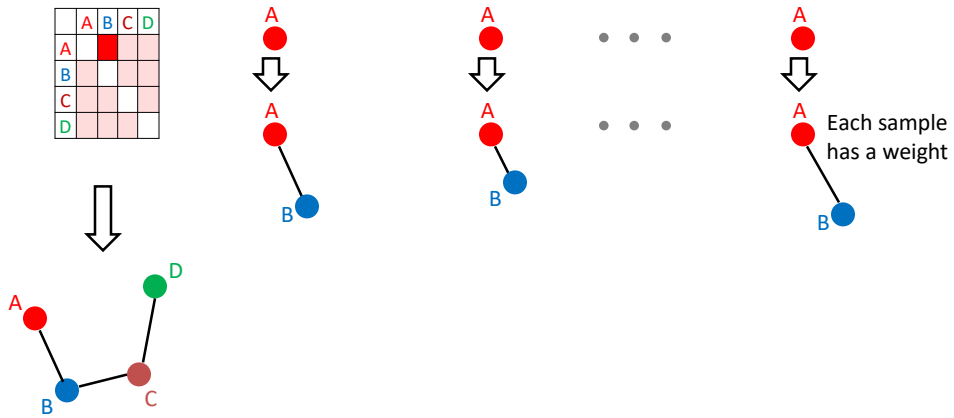
71

SIS in BACH: Illustration



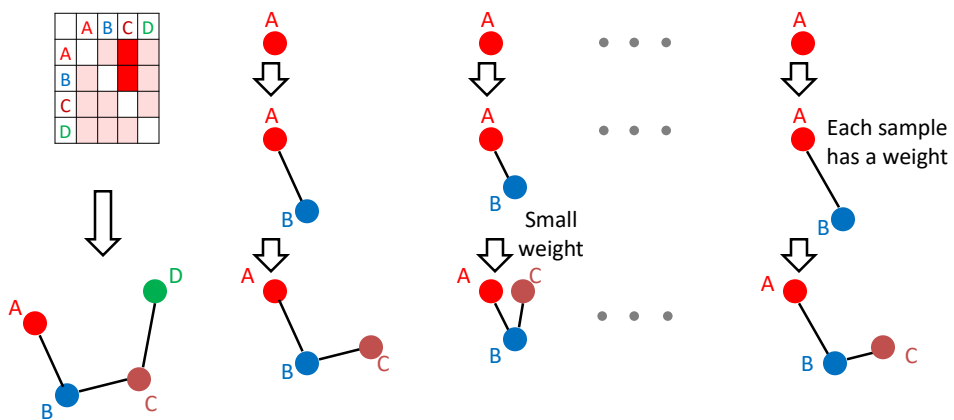
72

SIS in BACH: Illustration



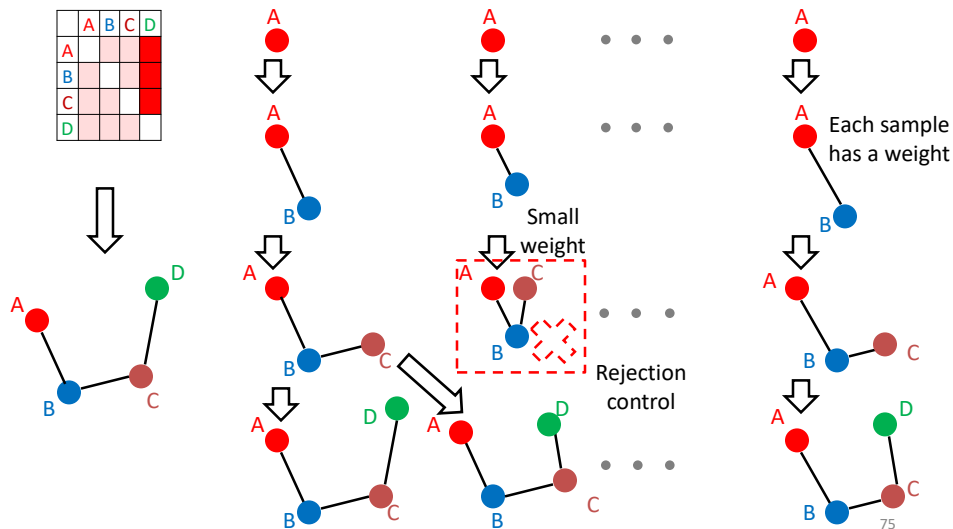
73

SIS in BACH: Illustration



74

SIS in BACH: Illustration



Hybrid Monte Carlo

- Goal: do efficient group move to refine initial 3D chromosomal structure, since local 3D coordinates are highly correlated.
- Combine molecular dynamics with Metropolis acceptance-rejection rule.

Hybrid Monte Carlo in BACH

- Goal: sampling from

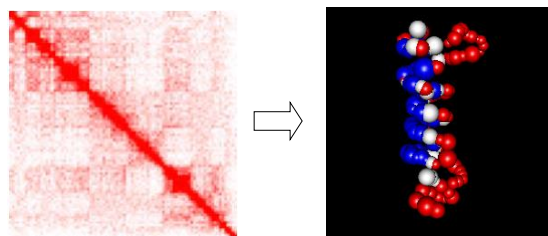
$$\pi(x_i, y_i, z_i, 1 \leq i \leq N | u_{ij}, 1 \leq i < j \leq N)$$

- Take partial derivate of log likelihood over 3D coordinates $(x_i, y_i, z_i, 1 \leq i \leq N)$.
- Run the leap-frog algorithm, adaptively tune the time interval to achieve acceptance rate $\sim 90\%$.

77

Conclusions

- BACH: reconstruct chromosome 3D structures from Hi-C data
- Remove systematic biases
- Predicted spatial distances are consistent with FISH data
- Elongation of chromatin is highly associated with genetic/epigenetic features.
- Separation of compartments of A and B can be visualized.



78

References

- **Hu M**, Deng K, Qin ZS, Dixon J, Selvaraj S, Fang J, Ren B, Liu JS. (2013) Bayesian inference of three-dimensional chromosomal organization. *PLoS Comput Biol.* **9** e1002893.
<http://www.people.fas.harvard.edu/~junliu/BACH/>
- Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, **Hu M**, Liu JS and Ren B. (2012) Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature* , 485, 376-380.