BIOS 731 Advanced Statistical Computing Fall 2022

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

Appliation: Transcription Factor Binding Sites Discovery



Example: cyclic receptor protein (CRP)

cole1	taatgtttgtgctggtttttgtggcatcgggcgagaatagcgcgtggtgtgaaagactgtttttttgatcgttttcacaaaaatggaagtccacagtcttgacag
ecoarabop	gacaaaaaacgcgtaacaaaagtgtctataatcacggcagaaaagtccacattgattatttgcacggcgtcacactttgctatgccatagcatttttatccataagtgccatagcatttttatccataagtgccatagcatttttatccataagtgccatagcatttttatccataagtgccatagcatttttatccataagtgccatagcatttttatccataagtgccatagcatttttatccataagtgccatagcatttttatccatagtgccatagcattttttatccataagtgccatagcattttttatccatagtgcgtcatagcattttttatccatagtgcgtcatagcattttttatccatagtgcgtcatagcattttttatccatagtgcgtcatagtgcgtcatagtgcgtcatagcatttttttt
ecobglr1	$a cas a {\tt c} cca {\tt a} a {\tt c} t {\tt t} a {\tt t} t {\tt$
ecocrp	cacaa agcga a agct at gct a a a a cagt cagg at gct a cagt a a tacatt g at gt a ct g cat gt at g ca a a g a c g t cacatt a c g t g cagt a c agt t g at a g c g t c a c at t a c g t g c agt a c agt t g at a g c g t c a c at t a c g t g c agt a c agt t g at a g c g c a g c g c a g t a c agt t g at a g c g c g c a g t a c a g t g at g c a g c g c g c g c g c g c g c g c g
ecocya	acggtgctacacttgtatgtagcgcatctttctttacggtcaatcagcatggtgttaaattgatcacgttttagaccattttttcgtcgtgaaactaaaaaaaa
ecodecop	agtgaattatttgaaccagatcgcattacagtgatgcaaacttgtaagtagatttccttaattgtgatgtgtgtg
ecogale	gcgcataaaaaacggctaaattcttgtgtaaacgattccactaatttattccatgtcacacttttcgcatctttgttatgctatggttatttcataccataagccacaatttattccatgtcacacttttcgcatctttgttatgctatggttatttcataccataagccacaatttattccatgtcacacttttcgcatctttgttatgctatggttatttcataccataagccacatttttcgcatcttttgttatgctatggttatttcataccataagccacattttttgtgttatgctatggttatttcatgctatggttatttcataccataagccacatttttttt
ecoilvbpr	gctccggcggggttttttgttatctgcaattcagtacaaaacgtgatcaacccctcaattttccctttgctgaaaaattttccattgtctcccctgtaaagctgt
ecolac	a a cgc a a tta a tgt g a g tt a g ct cact ca
ecomale	a cattaccg ccaattctg taacagagatcacacaaag cg cg tg gg gg cg tagg gg caagg ag gg at gg aa gg gg tg ccg ta taa ag aa actag ag t ccg t t ta cab a caa ag a cab
ecomalk	ggaggaggaggaggaggaggaggaggaggacacggcttctgtgaactaaaccgaggtcatgtaaggaatttcgtgatgttgcttgc
ecomalt	gatcagcgtcgtttttaggtgagttgttaataaagatttgggaattgtgacacagtgcaaattcagacacataaaaaaacgtcatcgcttgcattagaaaggtttct
ecoompa	$g {\tt ctgacaaaaaagattaaacataccttatacaagactttttttt$
ecotnaa	ttttttaaa cattaaa attcttacgtaatttataatctttaaa aagcatttaatattgctccccgaacgattgtgattcgattcacatttaaa caatttcaga
ecouxu1	eq:ccatgagagtgaaattgttgtggttgacccaattggattgacatgtcttaccaaaaggtagaacttatacgccatctcatccgatgcaagccatgtcttaccaaaggtagaacttatacgccatctcatccgatgcaagccatgtcttaccaaaggtagaacttatacgccatgtgttgacgtgatgaacttatacgccatgtcttaccaaggtagaacttatgtgtgttaacccaaggtagaacttatacgccatgtgtgttaacccaaggtagaacttatacgccatgtgtgtg
pbr-p4	${\tt ctggcttaactatgcggcatcagagcagattgtactgagagtgcaccatatgcggtgtgaaataccgcacagatgcgtaaggagaaaataccgcatcaggcgctcatgcgctaaggagaaaataccgcatcaggcgctcatgcgctgtgaaataccgcatgcgtaaggagaaaataccgcatcaggcgctcatgcgctgcatcaggcgctcatgcgctgcatcaggagaaaataccgcatcaggcgctcatgcgctgcatgcgctgcatcaggagaaaataccgcatcaggcgctcatgcgctgcatcaggcgctgcatcaggagaaaataccgcatcaggcgctgcatggcgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgctgcatcaggcgcgcatcaggcgctgcatcaggcgcatcaggcgcatcaggcgctgcatcaggcgcgctgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcgcatcaggcgcatcaggcgcgcatcaggcgcatcaggcgcgcatcaggcgcgcatcaggcgcatcaggcgcatcaggcgcatcaggcgcgcgc$
${ m trn9cat}$	${\tt ctgtgacggaagatcacttcgcagaataaataaatcctggtgtccctgttgataccgggaagccctgggccaacttttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatcactttggcgaaaatgagacgttgatcggcacggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagacgttgatcggaagatgagaaggatgagacgtgagaaatgagacgttgatcggaagatgagacgtgagaagatgagacgttgatcggaagacgtgagaagatgagacgtgagaaggatgagacggagaagatgagacggaagatgagacggaagatgagacggaagatgagaaggatgagacggaagatgagacggaagatgagacggaagatgagacggagaaggatgaggaggagagaggaggagaggaggaggag$
(tdc)	gatttttatactttaacttgttgatatttaaaggtatttaattgtaataacgatactctggaaagtattgaaagttaatttgtgagtggtcgcacatatcctgtt

Stormo and Hartzell, 1989

Example: cyclic receptor protein (CRP)

cole1	$taatgtttgtgctggtttttgtggcatcgggcgagaatagcgcgtggtgtgaaagactgtttt{tttgatcgttttcacaaaaatggaagtccacagtcttgacagagtccacagtcttgacagagtccacagtcttgacagagtccacagtcttgacagagtccacagtcttgacagagtccacagtcttgacaggtgtggaagtccacagtcttgacaggtgtggaagagtccacagtcttgacaggtgtggaaggtgggggg$
ecoarabop	gacaaaaacgcgtaacaaaagtgtctataatcacggcagaaaagtccacattgattattgcacggcgtcacactttgctatgccatagcatttttatccataagcatttttatccataagcatttttatccataagcatttttatccataagcatttttatccataagcatttttatccataagcattgctatagcattgccatagcatttttatccataagcattgccatagcatttttatccataagcattgccatagcatttttatccataagcattgccatagcattgccatagcatttttatccataagcattgccatagcatgcat
ecobglr1	a caa a to ccaa ta a ctta a tta ttggg atttg tta ta ta ta a cttta ta a a tto ccaa a a tta ccaa a g tta a ta ttg ttgg catgg to a ta tta tcaa tta ccaa a sta ccaa a g tta a ta ttg ttgg catgg to a ta tta tcaa tta ccaa a sta
ecocrp	cacaa a a g c f a t g c t a 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 c a a g g a c g t c a c a g t t g c a g t g c a g t a c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t t g a t a g c a g t g c a g t a g c a g t t g a t a g c a g t t g c a g t a c a g t t g a t a g c a g t t g c a g t a c a g t t a c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t g c a g t a c a g t t a c a g t t a c a g t t a c a g t t a c a g t t a c a g t t a c a g t t a c a g t t a c a g t
ecocya	$acggtgctacacttgtatgtagcgcatctttctttacggtcaatcagcatgg \\ \begin{tabular}{lllllllllllllllllllllllllllllllllll$
ecodecop	$agtgaatta {\tt tttgaaccagatcgcatta} cagtgatgcaaacttgtaagtagatttccttaattgtgatgtgtatcgaagtgtgttgcggagtagatgttagaata$
ecogale	$gcgcataaaaaacggctaaattcttgtgtaaacgattccactaa \\ \underline{ttattccatgtcacactt} \\ ttcgcatctttgttatgctatggttatttcataccataagccataagccataggttatttcataccataagccataggttatttcataccataagccataggttatttcataccataagccataggttatttcataccataggttatttcataccataggttattttcataggttatttcataggttattttatttcataggttatttt$
ecoilvbpr	gctccggcggggttttttgttatctgcaattcagtacaaaa cgtgatcaacccctcaattttccctttgctgaaaaattttccattgtctcccctgtaaagctgt
ecolac	$a acgcaattaa { \tt gtgagttagctcactcat} aggcaccccaggctttacactttatgcttccggctcgtatgttgtgtgggaattgtgagcggataacaatttcactttatgcttccggctcgtatgttgtgtgtg$
ecomale	acattaccgccaatte <mark>tgtaacagagatcacaca</mark> agogacggtgggggggtaggggaggatgggaaggatggaaagaggttgccgtataaagaaactagagtccgttta
ecomalk	ggaggaggcgggaggatgagaacacggcttctgtgaactaaaccgaggtcatgtaaggaatt cgtgatgttgcttgcaaaaatcgtggcgattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattgtgcgaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggatgtgcggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattttatgtgcgcaacaaggattgtgcggatgtgcgaacaaggattttatgtgcgcaacaaggattatgtgcgcgattttatgtgcgcaacaaggattgtgcgaacaaggattgtgcgaacaaggattttatgtgcgcaacaaggattgtgcgcaacaaggatgtggcgattttatgtgcgcaacaaggattgtgcgaacaaggatgatgtgcgaacaaggattgtgcgaacaaggatgtggggatggat
$\operatorname{ecomalt}$	gatcagcgtcgtttttaggtgagttgttaataaagatttggaat $ttggaaattgtgacacagtgcaaattgagacacataaaaaaagtcatcgcttgcattagaaaggtttct$
ecoompa	gctgacaaaaaagattaaacataccttatacaagactttttttt
ecotnaa	ttttttaaa cattaaaattcttacgtaatttataatctttaaaaaaagcatttaatattgctccccgaacga ttgtgattcgattc
ecouxu1	cccatgagagtgaaattgt <mark>igtgatgtggttaacccaa</mark> ttagaattcgggattgacatgtcttaccaaaaggtagaacttatacgccatctcatccgatgcaagc
pbr-p4	ctggcttaactatgcggcatcagagcagattgtactgagagtgcaccatatgcg <mark>t gtgaaataccgcacaga</mark> igcgtaaggagaaaataccgcatcaggcgctc
${ m trn9cat}$	${\tt ctgtgacggaagatcacttcgcagaataaataaatcctggtgtccctgttgataccgggaagccctgggccaacttttggcgaaaatgagacgttgatcggcacgggatgatgggcacgttgatcgggatgatgggcacgttgatcgggatggggatgggggggg$
(tdc)	$gatttttatactttaacttgttgatatttaaaggtatttaattgtaataacgatactctggaaagtattgaaagttaat { $

Stormo and Hartzell[§] 1989

6

Transcription factor binding site (TFBS)



Motif identification model

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

Posterior distributions

• The posterior conditional distribution for alignment variable *A*

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

DNA sequence data

$$R = (R_1, ..., R_J)$$

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

Motif Alignment Model



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

- Every non-site positions follows a common multinomial with p₀=(p_{0,1},..., p_{0,20})
- Every position *i* in the motif element follows probability distribution $p_i = (p_{i,1}, ..., p_{i,20})$

9

Statistical Model

- Objects:
 - Seq: sequence data to search for motif
 - $-\theta_0$: non-motif (genome background) probability
 - $-\theta$: motif probability matrix parameter
 - $-\pi$: site locations
- Problem: $P(\theta, \pi | \text{seq}, \theta_0)$
- Approach: alternately estimate
 - $-\pi$ by P($\pi \mid \theta$, seq, θ_0)
 - $-\theta$ by P($\theta \mid \pi$, seq, θ_0)
 - 10

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.

11





• Compute predictive frequencies of each position *i* in motif

 C_{ii} = count of amino acid type *j* at position *i*.

 $c_{0j}^{'}$ = count of amino acid type j in all non-site positions. $q_{ij} = (c_{ij} + b_j)/(K - 1 + B), \quad B = b_1 + \dots + b_K \text{ "pseudo-counts"}$

• Sample from the predictive distribution of a_k

$$P(a_{k} = l+1) \propto \prod_{i=1}^{w} \frac{q_{i,R_{k}(l+i)}}{q_{0,R_{k}(l+i)}}$$

References

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

13

Infer the 3D shape of chromosomes

Slides from Ming Hu

Microscopic Methods

• Fluorescent *in situ* hybridization (FISH)



http://en.wikipedia.org/wiki/CJ5genetics

FISH Data Representation



3D chromosomal structure

Contact Frequency vs. Spatial Distance



Problem setting



Problem setting



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

Problem setting



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

20 Yaffe and Tanay, 2011

Beads-on-a-string Representation

ACGTAGCTAGATACTGTAGTGTAGTTTGGAACCTGAGGG

21

Beads-on-a-string Representation

ACGTAGCTAGATACTGTAGTGTAGTTTGGAACCTGAGGG

Beads-on-a-string Representation

ACGTAGCTAG ATACTGTAGT GTAGTTTGGA ACCTGAGGG

23

Beads-on-a-string Representation







Bayesian Statistical Model



Bayesian Statistical Model



Bayesian Statistical Model

• Likelihood:

$$\begin{split} L(u_{ij}, 1 \le i < j \le N | x_i, y_i, z_i, 1 \le i \le N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) &= \prod_{1 \le i < j \le N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}!}}{u_{ij}!} \\ \log(\theta_{ij}) &= \beta_0 + \beta_1 \log\left(\sqrt{\left(x_i - x_j\right)^2 + \left(y_i - y_j\right)^2 + \left(z_i - z_j\right)^2}\right) \\ &+ \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j) \end{split}$$

Bayesian Statistical Model

• Likelihood: $\binom{N}{2}$ data points, 3N + 5 parameters $L(u_{ij}, 1 \le i < j \le N | x_i, y_i, z_i, 1 \le i \le N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) = \prod_{1 \le i < j \le N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}}}{u_{ij}!}$ $\log(\theta_{ij}) = \beta_0 + \beta_1 \log\left(\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}\right)$ $+\beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$

Bayesian Statistical Model

• Likelihood: $\binom{N}{2}$ data points, 3N + 5 parameters $L(u_{ij}, 1 \le i < j \le N | x_i, y_i, z_i, 1 \le i \le N, \beta_0, \beta_1, \beta_e, \beta_g, \beta_m) = \prod_{1 \le i < j \le N} \frac{e^{-\theta_{ij}} \theta_{ij}^{u_{ij}}}{u_{ij}!}$ $\log(\theta_{ij}) = \beta_0 + \beta_1 \log\left(\sqrt{(x_i - x_j)^2 + (y_i - y_j)^2 + (z_i - z_j)^2}\right)$ $+\beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$

Posterior distribution

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

31

Statistical Inference

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

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

Statistical Inference

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

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

► Initialization 1: use Poisson regression to obtain the initial values for $\beta_0, \beta_e, \beta_g, \beta_m$. Set $\beta_1 = -1$. $u_{ij} \sim Poisson(\theta_{ij}) \log(\theta_{ij}) = \beta_0 + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$

33

Statistical Inference

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

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

- ➤ Initialization 1: use Poisson regression to obtain the initial values for $\beta_0, \beta_e, \beta_g, \beta_m$. Set $\beta_1 = -1$. $u_{ij} \sim Poisson(\theta_{ij}) \log(\theta_{ij}) = \beta_0 + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$
- ▶ Initialization 2: use sequential important sampling to get the initial 3D chromosomal structure $\{x_i, y_i, z_i, 1 \le i \le N\}$.

Statistical Inference

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

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

- ➤ Initialization 1: use Poisson regression to obtain the initial values for $\beta_0, \beta_e, \beta_g, \beta_m$. Set $\beta_1 = -1$. $u_{ij} \sim Poisson(\theta_{ij}) \log(\theta_{ij}) = \beta_0 + \beta_e \log(e_i e_j) + \beta_g \log(g_i g_j) + \beta_m \log(m_i m_j)$
- ▶ Initialization 2: use sequential important sampling to get the initial 3D chromosomal structure $\{x_i, y_i, z_i, 1 \le i \le N\}$.
- Refinement: use Gibbs sampler with hybrid Monte Carlo to refine the initial values for parameters.

Sequential Importance Sampling





















Sequential Importance Sampling (SIS) Algorithm:

- (1) Design bridging distributions $\pi_t(\vec{x}_t)$ and proposal distributions $g_t(x_t|\vec{x}_{t-1})$
- (2) Sequentially draw weighted samples $x_t \sim g_t(x_t | \vec{x}_{t-1})$, and update weight

$$w_{t} = \frac{w_{t-1}\pi_{t}(\vec{x}_{t})}{\pi_{t-1}(\vec{x}_{t-1})g_{t}(x_{t}|\vec{x}_{t-1})}$$
⁴²

SIS in BACH: Outline

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

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

SIS in BACH: Outline

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

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

• Bridging distributions:

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

44

SIS in BACH: Outline

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

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

• Bridging distributions:

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

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

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

SIS in BACH: Illustration



SIS in BACH: Illustration



SIS in BACH: Illustration



48



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.

51 Duane, et al, 1987

Hybrid Monte Carlo in BACH

• Goal: sampling from

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

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

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.



53

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/