---

**1**

# BTRY 4840/6840, CS 4775
# Computational Genetics and Genomics
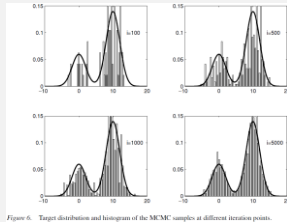


*Figure 6. Target distribution and histogram of the MCMC samples at different iteration points.*

October 25, 2018

Andrieu *et al.* (2003)

---

**2**

# Announcements

- Problem set 4 due today

- Problem set 5 out today

- Final project proposals due Tuesday (Oct 30)

- Reading on MCMC posted on website
  – A bit advanced: no need to understand it, but a good resource for learning about this topic

---

**3**

# Today's lecture

- More on MCMC
  – General properties, reversibility
  – Gibbs sampling and reversibility
  – Gibbs sampling for motif finding
  – STRUCTURE
  – Convergence

---

**4**

# MCMC background: Sampling

- Suppose we can draw samples $x^{(1)}, x^{(2)}, \dots, x^{(N)}$ from a pdf of interest $p(x)$

- Can then estimate several properties such as:

$$E[f(x)] = \int f(x)p(x)dx \approx \frac{1}{N}\sum_{i=1}^{N} f(x^{(i)}) \qquad [f \text{ any function}]$$

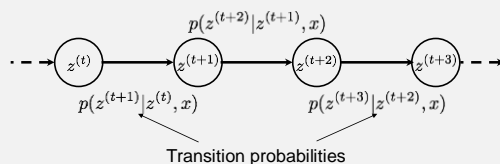$$p(\tilde{x}_j, \dots, \tilde{x}_k) = \int p(x_j, \dots, x_k) I(x_j = \tilde{x}_j, \dots, x_k = \tilde{x}_k) dx$$

$$\approx \frac{1}{N}\sum_{i=1}^{N} I\left(x_j^{(i)} = \tilde{x}_j, \dots, x_k^{(i)} = \tilde{x}_k\right)$$

- Key idea: each sample $x^{(i)}$ obtained in proportion to its probability, approximations implicitly include $p(x)$

---

**5**

# MCMC typically used to infer hidden variables

- Often with MCMC we:
  – Have observed data $x$ (often multivariate)
  – Want to infer unobserved variables $z$ (or parameters) related to $x$ via a model
  – That is, we wish to sample from $p(z|x)$, even if this is complex
  – We do so using a Markov chain



Transition probabilities

---

**6**

# Markov chain Monte Carlo (MCMC) basics

- Want to sample from some very complex $p(x)$ that is:
  – Hard to sample from directly
  – Quick to compute $p(x)$ for given $x$

- MCMC in a nutshell:
  – Aim: construct a Markov chain with each state $x^{(i)}$ a sample from the distribution of interest $p(x)$
  – How? Define transition probabilities between states such that the stationary distribution = $p(x)$
    • Initial samples often *not* from stationary distribution: called <u>burn-in</u>
    • Informally: after huge number of iterations, samples are from $p(x)$
  – When $x^{(i)} \sim p(x)$, get unbiased sample of distribution

- <u>Questions: How do we define transition probabilities? Will the chain converge to the stationary distribution?</u>

## Ensuring stationary distribution of chain is distribution we want to sample from [7]

- Markov chain is <u>reversible</u> wrt distribution $p(z^{(t)}|x)$ if:
  $$p(z^{(t)}|x)p(z^{(t+1)}|z^{(t)},x) = p(z^{(t+1)}|x)p(z^{(t)}|z^{(t+1)},x)$$
  - Also called <u>detailed balance</u>
  - ➢ Transition $z^{(t)} \rightarrow z^{(t+1)}$ or reverse with equal probability
- Let $\pi(z^{(t)}|x)$ be stationary distribution of chain, then:
  - ↪ $\pi(z^{(t+1)}|x) = \sum_{z^{(t)}} \pi(z^{(t)}|x)p(z^{(t+1)}|z^{(t)},x)$ [by def of stationarity]
  - ➢ Multiplying sample from stationary distribution by transition probabilities produces a value from the stationary distribution
  - Synonyms: <u>stationary</u>, <u>invariant</u>, <u>equilibrium</u> distribution
- In fact, reversibility wrt $p(z^{(t)}|x)$ implies it is stationary dist.
  $$\sum_{z^{(t)}} p(z^{(t)}|x)p(z^{(t+1)}|z^{(t)},x) = \sum_{z^{(t)}} p(z^{(t+1)}|x)p(z^{(t)}|z^{(t+1)},x)$$
  $$= p(z^{(t+1)}|x)\sum_{z^{(t)}} p(z^{(t)}|z^{(t+1)},x) = p(z^{(t+1)}|x)$$

## Will obtain stationary distribution under specific conditions [8]

- We must define transition equations such that the stationary distribution they imply is the distribution we wish to sample
  - In practice, we leverage reversibility to accomplish this
- To ensure convergence to the stationary distribution, transition equations must also be:
  - <u>Aperiodic</u>: the chain doesn't revisit states in periodic fashion
  - <u>Irreducible:</u> for any state, there is a positive probability of visiting all other states. (That is, the transition graph is connected.)
- Given these conditions, the MCMC will converge to stationary distribution at some point
  - Initially it will generally not sample from stationary distribution, but will eventually

## Gibbs sampling overview [9]

- Gibbs sampling: applies to multivariate distributions
  $$P(z_1, z_2, \ldots, z_n)$$
  - One of many MCMC approaches
  - Examples:
    - Motif finding: variables are start positions in $t$ sequences
    - STRUCTURE method (will discuss more later)
- General Gibbs sampling algorithm:
  1. Randomly initialize starting point $z_j^{(0)}$, for $1 \le j \le n$
  2. For $i = 0$ to $N - 1$:
     - Sample $z_1^{(i+1)} \sim p\left(z_1 \middle| z_2^{(i)}, z_3^{(i)}, \ldots, z_n^{(i)}\right)$
     - Sample $z_2^{(i+1)} \sim p\left(z_2 \middle| z_1^{(i+1)}, z_3^{(i)}, \ldots, z_n^{(i)}\right)$
     - Sample $z_n^{(i+1)} \sim p\left(z_n \middle| z_1^{(i+1)}, z_2^{(i+1)}, \ldots, z_{n-1}^{(i+1)}\right)$

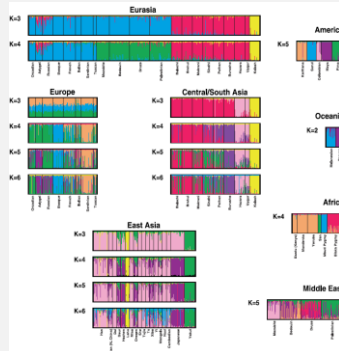## Gibbs samplers are reversible by construction [10]

- Gibbs sampling idea:
  - Successively update single variable, treating all others as known
    - "Known" values are most recent samples of other variables
  - Variable to be updated is sampled conditional on others
- Must show transition probabilities imply reversibility wrt $P(z_1, \ldots, z_n)$
  - Dropping super scripts for simplicity:
  - Want to show that:
    $$P(z_1, \ldots, z_{j-1}, z_j, z_{j+1}, \ldots, z_n)P(\tilde{z}_j|z_1, \ldots, z_{j-1}, z_{j+1}, \ldots, z_n) =$$
    $$P(z_1, \ldots, z_{j-1}, \tilde{z}_j, z_{j+1}, \ldots, z_n)P(z_j|z_1, \ldots, z_{j-1}, z_{j+1}, \ldots, z_n)$$
  - By definition of conditional probability, both sides equal to
    $\frac{P(z_1, \ldots, z_{j-1}, z_j, z_{j+1}, z_n)P(z_1, \ldots, z_{j-1}, \tilde{z}_j, z_{j+1}, \ldots, z_n)}{P(z_1, \ldots, z_{j-1}, z_{j+1}, \ldots, z_n)}$, so reversibility does hold
- Irreducibility and aperiodicity also necessary:
  - Is specific to each problem; can be simple to demonstrate

## STRUCTURE uses Gibbs sampling



Rosenberg *et al.* (2002)

## STRUCTURE algorithm overview [12]

- Given $X$: genotypes for $n$ samples; $K$: # populations
- Want to estimate:
  - Matrix $Z$ : ancestral population $1, \ldots, K$ of each allele
    - $z_l^{(i,a)}$: ancestral population of allele $a \in \{1,2\}$ at locus $l$ in sample $i$
  - Matrix $P$ containing allele frequencies in each population
    - $p_{klj}$: frequency of allele $j$ at locus $l$ in population $k$
  - Matrix Q containing admixture proportions for each sample
    - $q_k^{(i)}$: proportion of genome of sample $i$ that originated in $k$
1. Initialize $Z^{(0)}$ by random sample from uniform prior
2. Iterate:
   1. Sample $P^{(m)}, Q^{(m)}$ from $P(P, Q|X, Z^{(m-1)})$
   2. Sample $Z^{(m)}$ from $P(Z|X, P^{(m)}, Q^{(m)})$

   | Forms a Gibbs sampler |

## STRUCTURE Gibbs sampling

- $p_{klj}$ values are sampled from Dirichlet distribution at each locus, parameters for each allele set to $\lambda_j + n_{klj}$, where
  - $n_{klj}$ = counts of allele $j$ in population $k$, locus $l$ (from $Z$)
- $q_k^{(i)}$ values are sampled from Dirichlet for each sample, parameters for each population set to $\alpha + m_k^{(i)}$, where
  - $m_k^{(i)}$ = counts of alleles in sample $i$ from population $k$ (from $Z$)
- $z_l^{(i,a)}$ are sampled from: [$x_l^{(i,a)}$ is allele in sample $i$, locus $l$]

$$P\left(z_l^{(i,a)} = k \middle| X, P, Q\right) = \frac{q_k^{(i)} P\left(x_l^{(i,a)} \middle| P, z_l^{(i,a)} = k\right)}{\sum_{k'} q_{k'}^{(i)} P\left(x_l^{(i,a)} \middle| P, z_l^{(i,a)} = k'\right)}$$

where $P\left(x_l^{(i,a)} \middle| P, z_l^{(i,a)} = k\right) = p_{klx_l^{(i,a)}}$

## Dirichlet distribution

- Multivariate distribution where values of each variant sum to 1
  - Example: frequencies of $J$ alleles at a given locus
  - (Generalization of Beta distribution [which is univariate])
- Uses <u>concentration parameters</u> $\boldsymbol{\alpha}$ corresponding to each variable

$$f(x_1, \dots, x_J; \alpha_1, \dots, \alpha_J) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^{J} x_i^{\alpha_i - 1}$$

$$\boxed{\mathrm{E}[X_i] = \frac{\alpha_i}{\sum_j \alpha_j}}$$

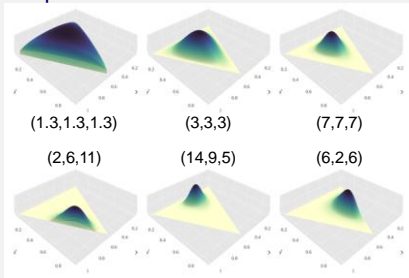Intuition: on average gives fraction of $\alpha$s contributed by $i$

$$\mathrm{Var}[X_i] = \frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}, \qquad \alpha_0 = \sum_{i=1}^{J} \alpha_i$$

## Concentration parameters and variance

- Distribution more strongly peaked for larger concentration parameters

Example 3-variant Dirichlet, $\boldsymbol{\alpha}$ shown

(1.3,1.3,1.3)    (3,3,3)    (7,7,7)

(2,6,11)    (14,9,5)    (6,2,6)

Plot: Empetrisor (wikipedia)

## Gibbs sampling for motif finding

Algorithm:
1. Choose <u>random start sites $s^{(i)}$ for all $i = 1 \dots t$</u>
2. For $i = 1 \dots t$, do:
   1. Compute position weight matrix using other <u>starts $s^{(j)}, j \neq i$</u>
      - Should include pseudocounts if using (good idea to do so)
   2. Compute likelihood of motif start at all positions $1 \dots n - k$ in $x^{(i)}$
      - For a given start position $r$, likelihood is $\sum_{j=1}^{k} w_{j, x_{r+j-1}^{(i)}}$
   3. Sample new start position $s^{(i)}$ in proportion to likelihood
      - Note: normalizing likelihoods (summing over all possible starts) gives the probability of motif starting at each site
   - Repeat for $N$ (large) iterations OR
     Until $w_{j,\cdot}$ doesn't change much for a pass over $t$

## Likelihood of motif positions

- Given sequences $X = \boldsymbol{x}_1, \dots, \boldsymbol{x}_t$ of length $n$, start positions $\boldsymbol{s} = (s_1, \dots, s_t)$, <u>background model $\boldsymbol{\theta}$</u>, and <u>motif model $\boldsymbol{\pi}$</u>, motif length $k$
- Likelihood is:

$$p(X, \boldsymbol{s} | \boldsymbol{\theta}, \boldsymbol{\pi}) = \prod_{i=1}^{t} p(s_i) \prod_{j=1}^{n} p(x_{i,j} | s_i, \boldsymbol{\theta}, \boldsymbol{\pi})$$

$$= \prod_{i=1}^{t} \frac{1}{n - k + 1} \left[ \prod_{j=1}^{s_i - 1} p(x_{i,j} | \boldsymbol{\theta}) \right] \left[ \prod_{j=s_i}^{s_i + k - 1} p(x_{i,j} | \boldsymbol{\pi}_{j-s_i+1}) \right] \left[ \prod_{j=s_i+k}^{n} p(x_{i,j} | \boldsymbol{\theta}) \right]$$

$$= \prod_{i=1}^{t} \frac{1}{n - k + 1} \left[ \prod_{j=1}^{n} p(x_{i,j} | \boldsymbol{\theta}) \right] \left[ \prod_{j=s_i}^{s_i + k - 1} \frac{p(x_{i,j} | \boldsymbol{\pi}_{j-s_i+1})}{p(x_{i,j} | \boldsymbol{\theta})} \right]$$

## Probability of start positions

- <u>Posterior probability of start position for sequence $i$ is</u>

$$p(s_i = p | \boldsymbol{x}_i, \boldsymbol{\theta}, \boldsymbol{\pi}) = \frac{p(s_i = p, \boldsymbol{x}_i | \boldsymbol{\theta}, \boldsymbol{\pi})}{p(\boldsymbol{x}_i | \boldsymbol{\theta}, \boldsymbol{\pi})}$$

$$= \frac{\frac{1}{L - k + 1} \left[ \prod_{j=1}^{n} p(x_{i,j} | \boldsymbol{\theta}) \right] \left[ \prod_{j=p}^{p+k-1} \frac{p(x_{i,j} | \boldsymbol{\pi}_{j-p+1})}{p(x_{i,j} | \boldsymbol{\theta})} \right]}{\frac{1}{L - k + 1} \left[ \prod_{j=1}^{n} p(x_{i,j} | \boldsymbol{\theta}) \right] \left[ \sum_{s_i=1}^{n-k+1} \prod_{j=s_i}^{s_i+k-1} \frac{p(x_{i,j} | \boldsymbol{\pi}_{j-s_i+1})}{p(x_{i,j} | \boldsymbol{\theta})} \right]}$$

$$= \frac{\prod_{j=p}^{p+k-1} \frac{p(x_{i,j} | \boldsymbol{\pi}_{j-p+1})}{p(x_{i,j} | \boldsymbol{\theta})}}{\sum_{s_i=1}^{n-k+1} \prod_{j=s_i}^{s_i+k-1} \frac{p(x_{i,j} | \boldsymbol{\pi}_{j-s_i+1})}{p(x_{i,j} | \boldsymbol{\theta})}}$$

---

**20**

## Thinning in MCMC

- After convergence, MCMC samples from stationary distribution, however…

- Sample $z^{(i)}$ is conditioned on previous sample $z^{(i-1)}$
  - Note: in Gibbs sampling, $z$ is multivariate

- Successive samples $z^{(i-1)}, z^{(i)}$ are *not* independent
  - Samples separated by some number of iterations are (effectively) independent
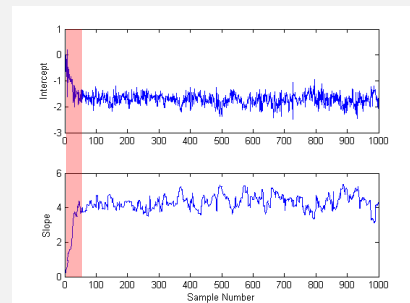  - Can address this by <u>thinning</u>: keep only every $k$th sample

---

**21**

## Convergence of MCMC (1)

- How do we know if the MCMC has converged?

- Simple approach: plot log likelihood in each iteration, visually inspect
  - At stationarity and with good mixing, likelihood will have roughly constant mean and variance

- Plots for parameters being estimated will also have roughly constant mean and variance

- Can make informed guess about burn-in length $B$ from these plots

---

**22**

## Visually inspection: trace plots



Plot: iteration number vs. value of a parameter being estimated
Top plot has higher variance: "mixes" better than bottom
Relatively constant mean and variance *suggests* stationarity

---

**23**

## Convergence of MCMC (2)

- Best to run more than one chain at different places and ensure that they converge to the same distribution

- More principled:
  - Run ≥ 2 chains from high variance (overdispersed) starting values
  - Compare within-chain and between-chain variance of estimated parameters
  - Expect similar variance for the comparison (i.e., within-chain and between-chain should be roughly the same)
    - If not, should run longer
  - [Above based on Gelman-Rubin diagnostic; this diagnostic gives recommendation for when to run chain longer]

- Other principled approaches exist (beyond scope here)

---

**24**

## Gibbs sampling algorithm with burn-in, thinning

- … with $B$ burn-in iterations, $S$ sampling iterations, thinning interval $k$

1. Initialize $z^{(0)}$ such that $p(z^{(0)}|x) > 0$
2. For $i = 1 \dots (B + S)$:
    1. For $j = 1 \dots |z|$:
        - Sample $z_j^{(i)}$ from $p\left(z_j^{(i)}\big|z_1^{(i+1)}, \dots, z_{j-1}^{(i+1)}, z_{j+1}^{(i)}, \dots, z_n^{(i)}, x\right)$
    2. If $(i > B$ and $(i \bmod k) = 0)$ retain sample $z^{(i)}$

# Final notes

- More on MCMC
  - Can use reversibility to ensure desired stationary distribution
  - Gibbs samplers are reversible by construction
  - Gibbs sampling for motif finding
  - STRUCTURE
  - Convergence