

X represents data

θ_i represents i^{th} of N parameters ($i \in \{1, \dots, N\}$).

$\theta_i^{(j)}$ represents j^{th} sampled value of i^{th} parameter

Gibbs Sampler:

Initialize with parameter values $\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_N^{(0)}$.

We require posterior probability of initial state to be positive

$$\Pr(\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_N^{(0)} \mid X) > 0$$

Set $k = 0$.

Step 1: Sample $\theta_1^{(k+1)}$ from $\Pr(\theta_1 \mid X, \theta_2^{(k)}, \dots, \theta_N^{(k)})$

Step 2: Sample $\theta_2^{(k+1)}$ from $\Pr(\theta_2 \mid X, \theta_1^{(k+1)}, \theta_3^{(k)}, \dots, \theta_N^{(k)})$

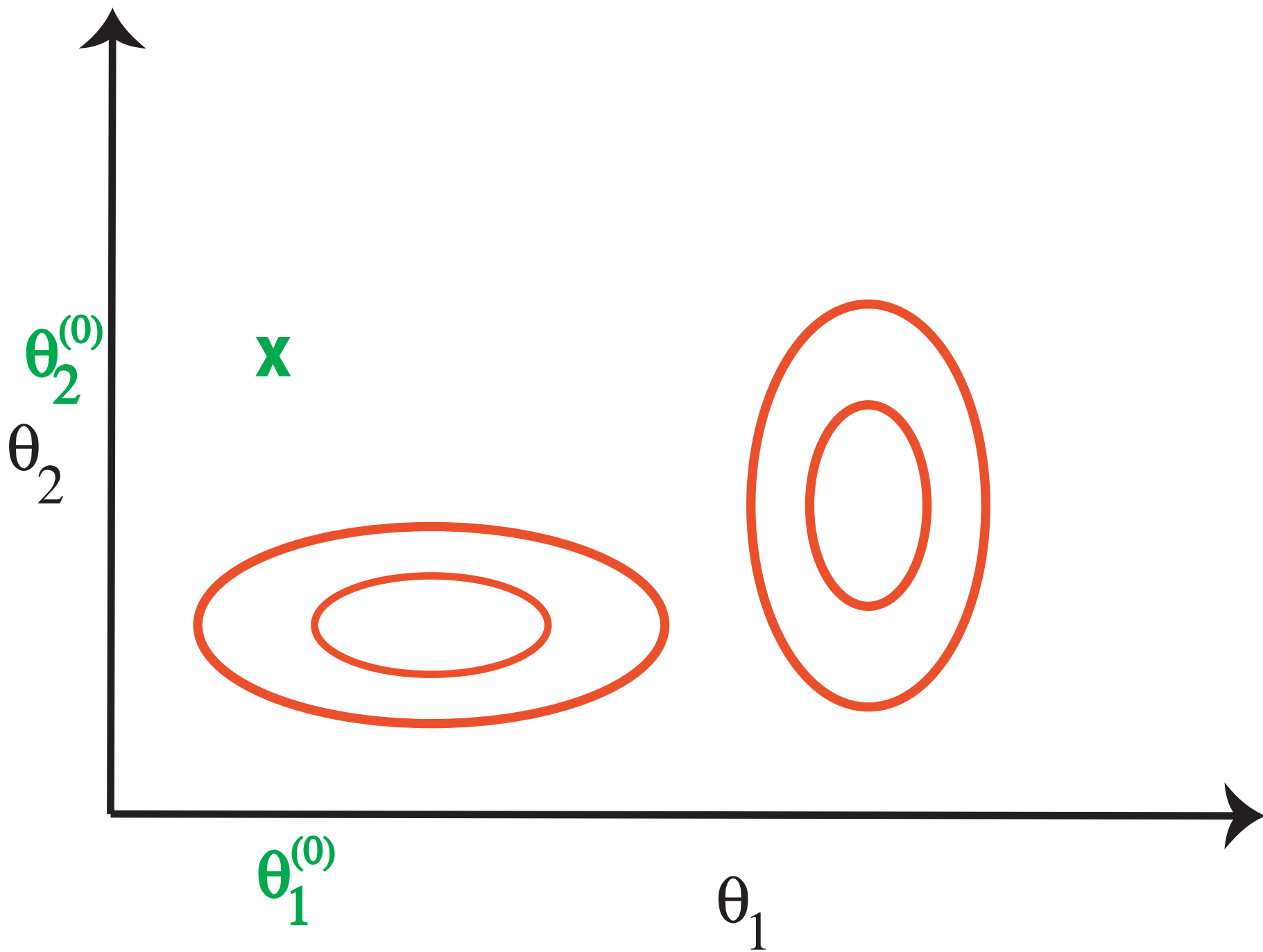
Step 3:

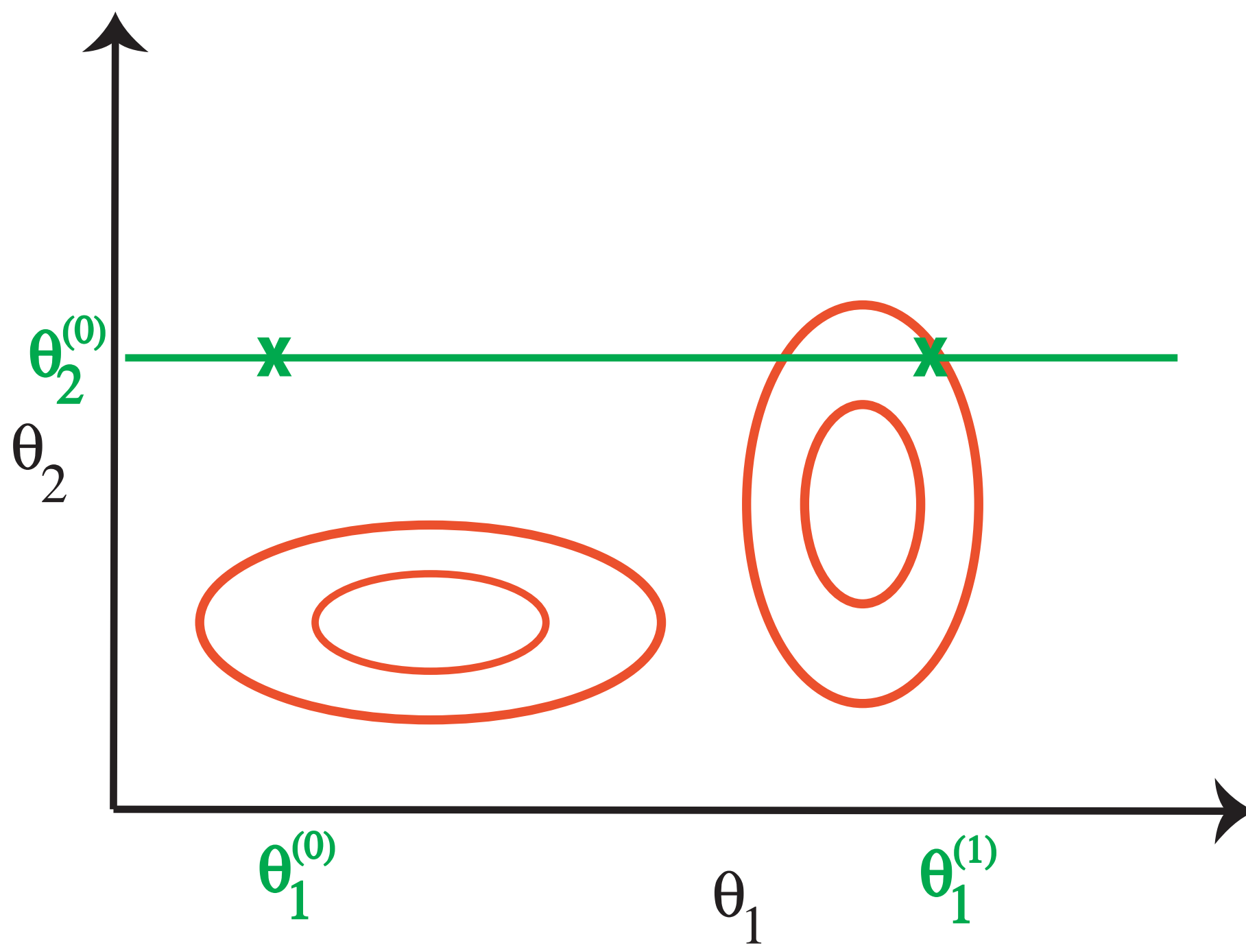
Sample $\theta_3^{(k+1)}$ from $\Pr(\theta_3 \mid X, \theta_1^{(k+1)}, \theta_2^{(k+1)}, \theta_4^{(k)}, \dots, \theta_N^{(k)})$

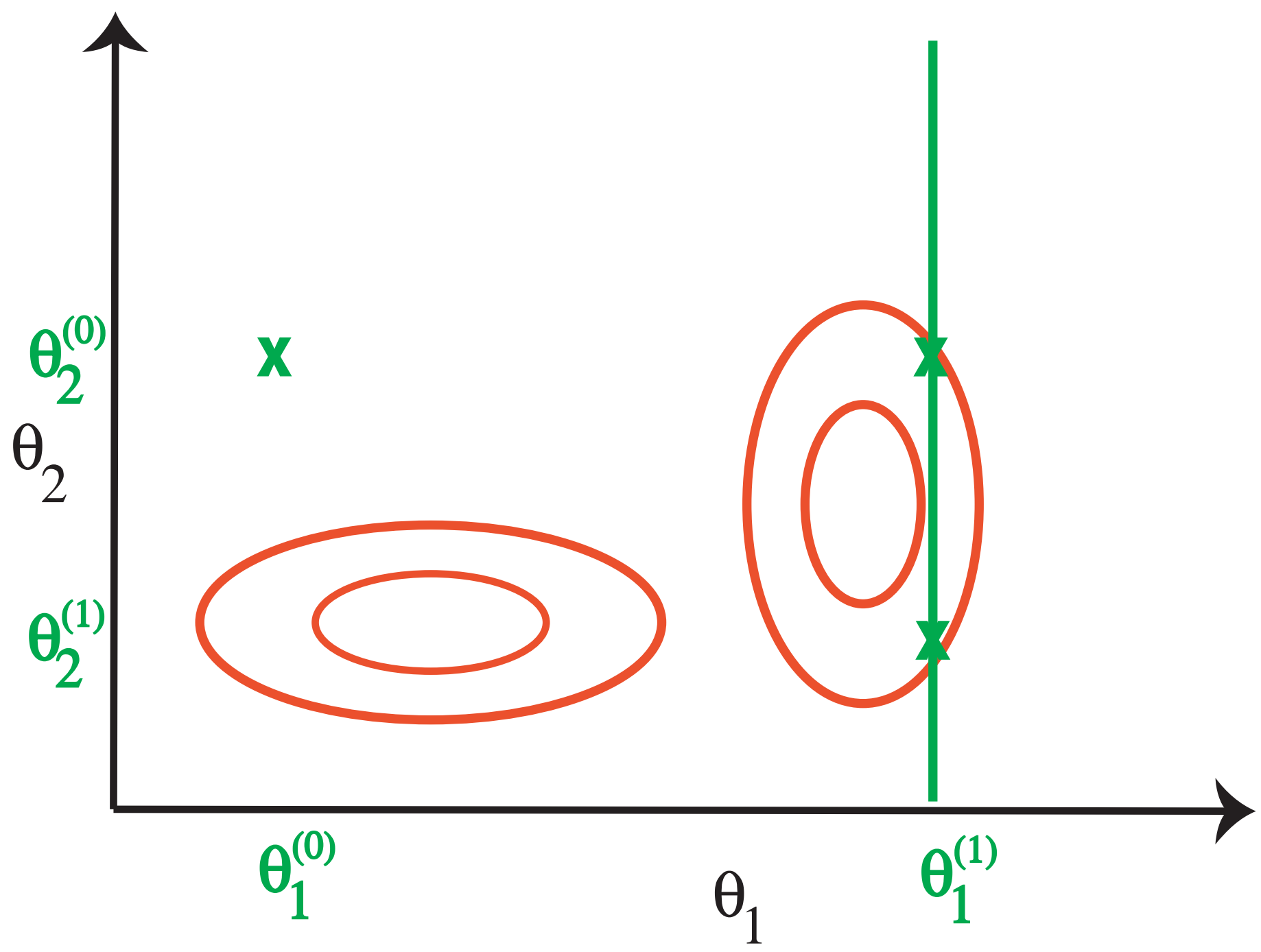
...

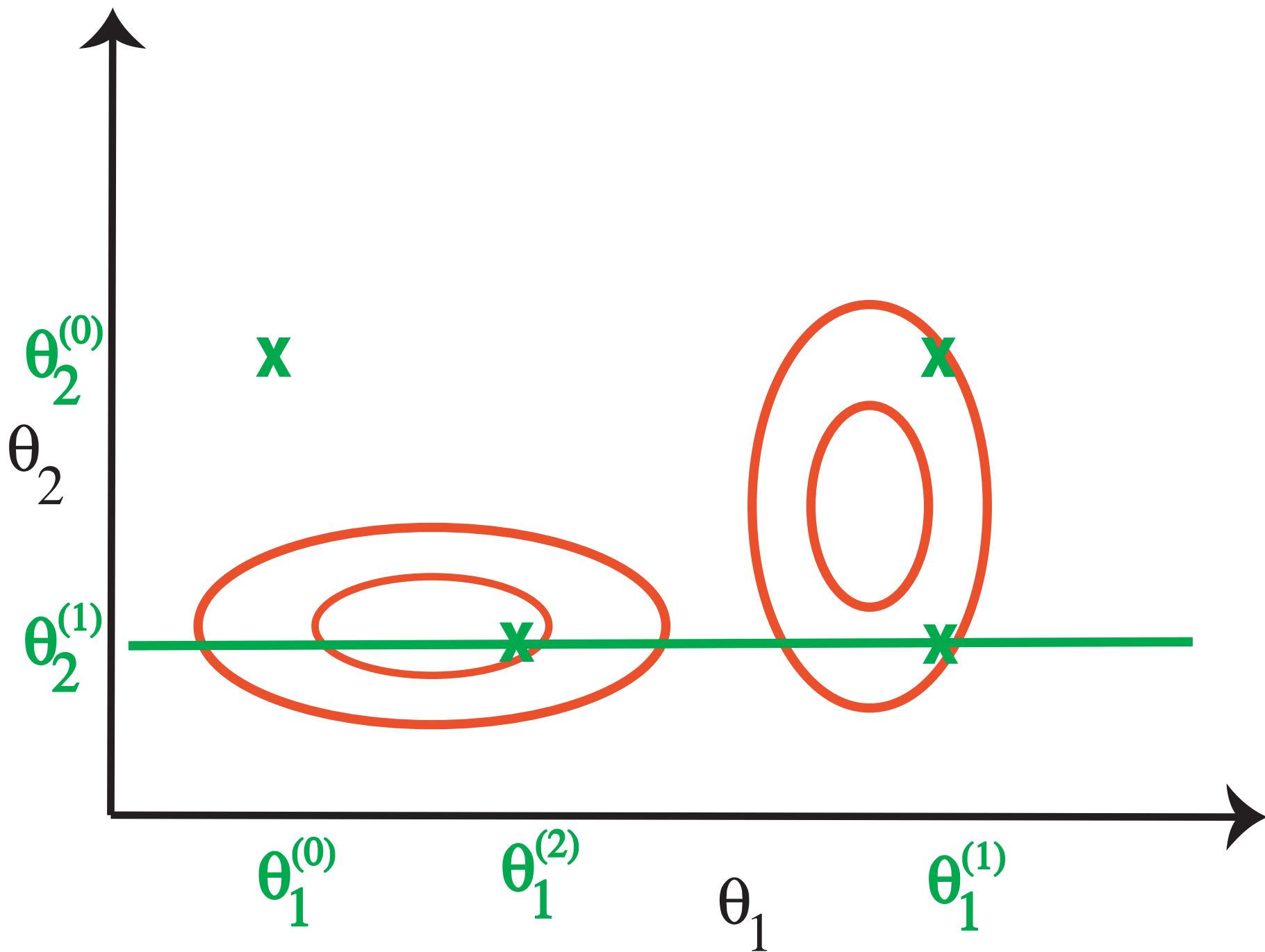
Step N: Sample $\theta_N^{(k+1)}$ from $\Pr(\theta_N \mid X, \theta_1^{(k+1)}, \dots, \theta_{N-1}^{(k+1)})$

Step N+1: Set $k = k + 1$. Goto Step 1.









Gibbs Sampling - another Markov chain Monte Carlo technique for sampling from posterior distributions (developed by Geman and Geman 1984)

Skeleton of an Example:

$\alpha^{(t)}$ represents t^{th} sampled pairwise alignment

$\mathcal{S}^{(t)}$ represents t^{th} set of sampled parameters controlling nucleotide substitution (or amino acid replacement)

$\delta^{(t)}$ represents t^{th} set of sampled parameters controlling insertion and deletion

Data are two non-aligned sequences A and B

To sample from joint posterior density $P(\alpha, S, \delta|A, B)$:

1. Randomly pick initial values $\alpha^{(0)}, S^{(0)}, \delta^{(0)}$
2. Set $t = 1$
3. Sample $\alpha^{(t)}$ from $P(\alpha|S^{(t-1)}, \delta^{(t-1)}, A, B)$
4. Sample $S^{(t)}$ from $P(S|\alpha^{(t)}, \delta^{(t-1)}, A, B)$
5. Sample $\delta^{(t)}$ from $P(\delta|\alpha^{(t)}, S^{(t)}, A, B)$
6. Set $t = t + 1$, Go to Step 3

Chip Lawrence, Jun Liu, and collaborators have done much work on adapting Gibbs Sampling and related Markov chain Monte Carlo techniques to the detection and characterization of subtle sequence patterns.

Most basic case:

N non-aligned sequences

Let $S = \{S_1, S_2, \dots, S_N\}$ represent these sequences

A priori knowledge is that each sequence will have exactly one ungapped motif of length W

a_i represents the position in sequence i where the ungapped motif of length W begins

Model:

1. All sequences are independent realizations of some process (Specifically, the process generates sequences according to some model but the length of the sequence is assumed given)
2. Amino acid residues at a site in a sequence are determined independently of residues at all other positions
3. Residue types of all sequence positions that are not in the motif are sampled from some background model probability distribution
4. Each of the W positions in the motif is associated with its own set of 20 amino acid frequencies

(Dirichlet priors or a mixture of Dirichlet priors can be used for both the background and motif-position amino acid frequencies.)

5. All possible starting positions for the motif within a sequence can be considered equally likely *a priori*

More Notation: Let π be the vector that contains information about residue frequencies at all background positions and at all motif positions

Goal:

1. Find the set of starting positions of the motif in Sequences 1, 2, ..., N.
2. Estimate the frequencies of all amino acids in position 1 of the motif, position 2 of the motif, ..., position W of the motif.

Outline of Gibbs Procedure (details follow):

Step 1. Make random initial guess as to values of all parameters in models and denote these guesses by $\pi^{(0)}, a_1^{(0)}, a_2^{(0)}, \dots, a_N^{(0)}$.

Step 2. Set $t = 1$

Step 3. Sample $\pi^{(t)}$ from $P(\pi | a_1^{(t-1)}, a_2^{(t-1)}, \dots, a_N^{(t-1)}, S)$

Step 4.

Sample $a_1^{(t)}$ from $P(a_1 | \pi^{(t)}, a_2^{(t-1)}, \dots, a_N^{(t-1)}, S) = P(a_1 | \pi^{(t)}, S_1)$

$a_2^{(t)}$ from $P(a_2 | \pi^{(t)}, a_1^{(t)}, a_3^{(t-1)}, \dots, a_N^{(t-1)}, S) = P(a_2 | \pi^{(t)}, S_2)$

\dots
 $a_N^{(t)}$ from $P(a_N | \pi^{(t)}, a_1^{(t)}, a_2^{(t)}, \dots, a_{N-1}^{(t)}, S) = P(a_N | \pi^{(t)}, S_N)$

Step 5. Set $t=t+1$

Step 6. Go to Step 3

Gibbs Motif Implementation (an “almost” Gibbs Sampler)

0. Start by randomly picking a position in each sequence where the motif begins. This defines the “motif alignment” and the background positions.

1. Randomly pick a sequence and temporarily remove it from remaining sequences (Note: this would not be done in a “proper” Gibbs Sampler). Estimate amino acid frequencies at site j in motif with $(c_{t,j,r} + b_{t,j,r}) / (c_t + b_t)$ where: $c_{t,j,r}$ is # residues of type r at position j for motif t , $b_{t,j,r}$ is # residues that are pseudocounts for type r at position j for motif t , c_t is total # residues per position, and b_t is total number # pseudocounts per position (Note: we are effectively using estimates of posterior means rather than sampling from posterior of residue frequencies).

Estimate background frequencies of amino acid type r with $(c_r + q_r) / (c + q)$ where: c_r is # background counts of type r , q_r is # pseudocounts for background of type r , c is total # background counts, and q is total # background pseudocounts.

2. Add the removed sequence back to motif alignment. Do this by randomly picking starting position for motif according to its posterior probability.

A. Posterior probability for each starting position is joint probability of data and starting position divided by probability of data.

B. Probability of data is simply sum over all possible starting positions of the joint probability of data and starting position

C. Joint probability of data and starting position is product of probability of data given starting position and probability of starting position (probability of starting position could be treated uniform over all possible starting positions)

D. Prob. of data given starting position is product of appropriate frequencies (either motif or background) for all residues in sequence

① Example: $W=2$ (motif length), $N=3$ (# of sequences)

π contains parameters for nucleotide frequencies

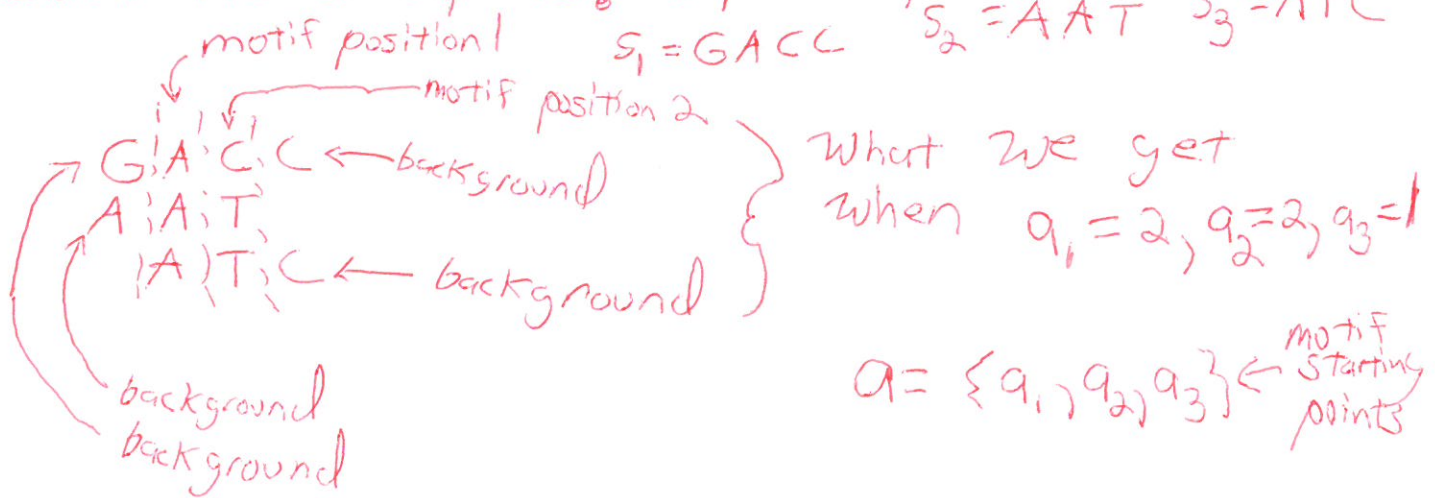
$$\pi = \{ \pi_B, \pi_1, \pi_2 \} \text{ where}$$

$$\pi_B = \{ \pi_{BA}, \pi_{BC}, \pi_{BG}, \pi_{BT} \} \sim \text{background frequencies}$$

$$\pi_1 = \{ \pi_{1A}, \pi_{1C}, \pi_{1G}, \pi_{1T} \} \sim \text{frequencies at motif position 1}$$

$$\pi_2 = \{ \pi_{2A}, \pi_{2C}, \pi_{2G}, \pi_{2T} \} \sim \text{frequencies at motif position 2}$$

Consider these 3 sequences: Sequence 1, Sequence 2, Sequence 3
 $s_1 = GACC$ $s_2 = AAT$ $s_3 = ATC$



Dirichlet Priors [Independent Across motif positions & background frequencies]
 (α 's below are called pseudocounts)

$$P(\pi) = P(\pi_B, \pi_1, \pi_2) = P(\pi_B) P(\pi_1) P(\pi_2)$$

$$P(\pi_B) = C_B \pi_{BA}^{\alpha_{BA}-1} \pi_{BC}^{\alpha_{BC}-1} \pi_{BG}^{\alpha_{BG}-1} \pi_{BT}^{\alpha_{BT}-1}$$

$$P(\pi_1) = C_1 \pi_{1A}^{\alpha_{1A}-1} \pi_{1C}^{\alpha_{1C}-1} \pi_{1G}^{\alpha_{1G}-1} \pi_{1T}^{\alpha_{1T}-1}$$

$$P(\pi_2) = C_2 \pi_{2A}^{\alpha_{2A}-1} \pi_{2C}^{\alpha_{2C}-1} \pi_{2G}^{\alpha_{2G}-1} \pi_{2T}^{\alpha_{2T}-1}$$

② Likelihoods $P(S|a, \pi) = P(s_1|a_1, \pi) P(s_2|a_2, \pi) P(s_3|a_3, \pi)$

$$P(s_1 = GACC | a_1, \pi) = P(s_1 = GACC | a_1 = 2, \pi)$$

$$= \pi_{BG} \pi_{IA} \pi_{2C} \pi_{BC}$$

$$P(s_2 = AAT | a_2 = 2, \pi) = \pi_{BA} \pi_{IA} \pi_{2A}$$

$$P(s_3 = ATC | a_3 = 1, \pi) = \pi_{IA} \pi_{2T} \pi_{BC}$$

$$P(\pi | a_1 = 2, a_2 = 2, a_3 = 1, S) = \frac{P(\pi | a_1, a_2, a_3) P(S | \pi, a_1, a_2, a_3)}{\cancel{P(a_1, a_2, a_3)}} P(S | a_1, a_2, a_3)$$

$$= \frac{P(\pi) P(S | \pi, a_1, a_2, a_3)}{\cancel{P(a_1, a_2, a_3)}} P(S | a_1, a_2, a_3)$$

$$= C P(\pi_B) P(\pi_1) P(\pi_2) P(s_1 | a_1, \pi) P(s_2 | a_2, \pi) P(s_3 | a_3, \pi)$$

$$= C \left[\pi_{BA}^{\alpha_{BA}-1+1} \pi_{BC}^{\alpha_{BC}-1+2} \pi_{BG}^{\alpha_{BG}-1+1} \pi_{BT}^{\alpha_{BT}-1} \right]$$

$$\times \left[\pi_{IA}^{\alpha_{IA}-1+3} \pi_{IG}^{\alpha_{IG}-1} \pi_{IC}^{\alpha_{IC}-1} \pi_{IT}^{\alpha_{IT}-1} \right]$$

← All Dirichlet Distributions!

$$\times \left[\pi_{2A}^{\alpha_{2A}-1} \pi_{2C}^{\alpha_{2C}-1+1} \pi_{2G}^{\alpha_{2G}-1} \pi_{2T}^{\alpha_{2T}-1+2} \right]$$

$$\textcircled{3} \quad P(a_i=2 | \pi, s_i) = \frac{P(s_i, a_i=2 | \pi)}{P(s_i | \pi)} = \frac{P(s_i | a_i=2, \pi) P(a_i=2 | \pi)}{\sum_{j=1}^3 P(s_i | a_i=j, \pi) P(a_i=j | \pi)}$$

$$= \frac{P(s_i | a_i=2, \pi) P(a_i=2)}{\sum_{j=1}^3 P(s_i | a_i=j, \pi) P(a_i=j)} = \frac{P(s_i | a_i=2, \pi) \cdot \frac{1}{3}}{\sum_{j=1}^3 P(s_i | a_i=j, \pi) \cdot \frac{1}{3}}$$

Note: $P(s_i | a_i=2, \pi)$ is explained at top of previous page