# Learning Sequence Motif Models Using Gibbs Sampling

BMI/CS 776 www.biostat.wisc.edu/bmi776/ Spring 2011 Mark Craven craven@biostat.wisc.edu

#### **Goals for Lecture**

the key concepts to understand are the following

- Markov Chain Monte Carlo (MCMC) and Gibbs sampling
- Gibbs sampling applied to the motif-finding task
- parameter tying
- incorporating prior knowledge using Dirichlets and Dirichlet mixtures

# Gibbs Sampling: An Alternative to EM

- EM can get trapped in local minima
- one approach to alleviate this limitation: try different (perhaps random) initial parameters
- Gibbs sampling exploits randomized search to a much greater degree
- can view it as a stochastic analog of EM for this task
- in theory, Gibbs sampling is less susceptible to local minima than EM
- [Lawrence et al., Science 1993]

# Gibbs Sampling Approach

- in the EM approach we maintained a distribution  $Z_i$  over the possible motif starting points for each sequence
- in the Gibbs sampling approach, we'll maintain a specific starting point for each sequence  $a_i$  but we'll keep randomly resampling these

# Gibbs Sampling Algorithm for Motif Finding

given: length parameter *W*, training set of sequences choose random positions for *a* do pick a sequence  $X_i$ estimate *p* given current motif positions *a* (using all sequences but  $X_i$ ) (predictive update step) sample a new motif position  $a_i$  for  $X_i$  (sampling step) until convergence return: *p*, *a* 

# Markov Chain Monte Carlo (MCMC)

 Consider a Markov chain in which, on each time step, a grasshopper randomly chooses to stay in its current state, jump one state left or jump one state right.



- •
- let  $P^{(t)}(u)$  represent the probability of being in state u at time t in the random walk  $P^{(0)}(0) = 1 \qquad P^{(0)}(+1) = 0 \qquad P^{(0)}(+2) = 0$

 $P^{(1)}(0) = 0.5 \qquad P^{(1)}(+1) = 0 \qquad P^{(1)}(+2) = 0$   $P^{(1)}(0) = 0.5 \qquad P^{(1)}(+1) = 0.25 \qquad P^{(1)}(+2) = 0$   $P^{(2)}(0) = 0.375 \qquad P^{(2)}(+1) = 0.25 \qquad P^{(2)}(+2) = 0.0625$   $\vdots \qquad \vdots \qquad \vdots$   $P^{(100)}(0) \approx 0.11 \qquad P^{(100)}(+1) \approx 0.11 \qquad P^{(100)}(+2) \approx 0.11$ 

#### The Stationary Distribution

• let *P*(*u*) represent the probability of being in state *u* at any given time in a random walk on the chain

$$P^{(t)}(u) \approx P^{(t+1)}(u)$$

$$P^{(t+1)}(u) = \sum_{v} P^{(t)}(v)\tau(u \mid v)$$
probability of
state v
probability of
transition  $v \rightarrow u$ 

 the stationary distribution is the set of such probabilities for all states

# Markov Chain Monte Carlo (MCMC)

- we can view the motif finding approach in terms of a Markov chain
- each state represents a configuration of the starting positions (*a<sub>i</sub>* values for a set of random variables A<sub>1</sub> ... A<sub>n</sub>)
- transitions correspond to changing selected starting positions (and hence moving to a new state)

ACAT <mark>CCG</mark>		AC <mark>ATC</mark> CG
CGACTAC		CGACTAC
ATTGAGC		ATTGAGC
CGTTGAC		CGTTGAC
GAGTGAT		GAGTGAT
TCGTTGG	$\tau(y \mid u)$	TCGTTGG
ACAGGAT	u(v + u)	ACAGGAT
TAGCTAT		TAGCTAT
GCTACCG		GCTACCG
GGCCTCA		GGCCTCA
state u		state v

# Markov Chain Monte Carlo

- for the motif-finding task, the number of states is enormous •
- key idea: construct Markov chain with stationary • distribution equal to distribution of interest; use sampling to find most probable states
- detailed balance:

state *u* 

$$P(u)\tau(v \mid u) = P(v)\tau(u \mid v)$$
probability of
state u
probability of
transition  $u \rightarrow v$ 

when detailed balance holds:

$$\frac{1}{N}\lim_{N\to\infty}count(u) = P(u)$$

# MCMC with Gibbs Sampling

Gibbs sampling is a special case of MCMC in which

- Markov chain transitions involve changing one variable at a time
- transition probability is conditional probability of the changed variable given all others
- i.e. we sample the joint distribution of a set of random variables  $P(A_1...A_n)$  by iteratively sampling from  $P(A_i | A_1 ... A_{i-1}, A_{i+1} ... A_n)$

#### Gibbs Sampling Approach possible state transitions when first sequence is selected • ACATCCG CGACTAC ACATCCG ACATCCG ATTGAGC CGACTAC CGACTAC CGTTGAC ATTGAGC • ACATCCG ATTGAGC GAGTGAT CGTTGAC CGACTAC CGTTGAC TCGTTGG ACATCCG GAGTGAT ATTGAGC GAGTGAT ACAGGAT CGACTAC TCGTTGG CGTTGAC TCGTTGG TAGCTAT ACATCCG ATTGAGC ACAGGAT GAGTGAT ACAGGAT **GCTACCG** CGACTAC CGTTGAC TAGCTAT TCGTTGG TAGCTAT GGCCTCA ATTGAGC GAGTGAT **GCTACCG** ACAGGAT **GCTACCG** CGTTGAC TCGTTGG GGCCTCA TAGCTAT GGCCTCA GAGTGAT ACAGGAT **GCTACCG** TCGTTGG TAGCTAT GGCCTCA ACAGGAT **GCTACCG** TAGCTAT GGCCTCA **GCTACCG** GGCCTCA

# **Gibbs Sampling Approach**

• How do we get the transition probabilities when we don't know what the motif looks like?



# **Sampling New Motif Positions**

• for each possible starting position,  $A_i = j$ , compute the likelihood ratio (leaving sequence *i* out of estimates of *p*)

$$LR(j) = \frac{\prod_{k=j}^{j+W-1} p_{c_k, k-j+1}}{\prod_{k=j}^{j+W-1} p_{c_k, 0}}$$

• randomly select a new starting position  $A_i = j$  with probability LR(j)



#### The Phase Shift Problem

- Gibbs sampler can get stuck in a local maximum that corresponds to the correct solution shifted by a few bases
- solution: add a special step to shift the *a* values by the same amount for all sequences. Try different shift amounts and pick one in proportion to its probability score

#### **Convergence of Gibbs**



# Using Background Knowledge to Bias the Parameters

let's consider two ways in which background knowledge can be exploited in the motif finding process

- 1. accounting for palindromes that are common in DNA binding sites
- 2. using Dirichlet mixture priors to account for biochemical similarity of amino acids

# Using Background Knowledge to Bias the Parameters

 Many DNA motifs have a palindromic pattern because they are bound by a protein *homodimer*: a complex consisting of two identical proteins



#### **Representing Palindromes**

 parameters in probabilistic models can be "tied" or "shared"



 during motif search, try tying parameters according to palindromic constraint; accept if it increases likelihood test (half as many parameters)



$$\begin{bmatrix} p_{a,0} & p_{a,1} & \cdots & p_{a,W} \\ p_{c,0} & p_{c,1} & \cdots & p_{c,W} \\ p_{g,0} & p_{g,1} & \cdots & p_{g,W} \\ p_{t,0} & p_{t,1} & \cdots & p_{t,W} \end{bmatrix}$$

$$p_{a,1} \equiv p_{t,W} = \frac{n_{a,1} + n_{t,W} + d_{a,1} + d_{t,W}}{\sum_{b} (n_{b,1} + d_{b,1}) + \sum_{b} (n_{b,W} + d_{b,W})}$$

#### **Using Dirichlet Mixture Priors**

recall that the EM/Gibbs update the parameters by:

$$p_{c,k} = \frac{n_{c,k} + d_{c,k}}{\sum_{b} (n_{b,k} + d_{b,k})}$$

- Can we use background knowledge to guide our choice of pseudocounts ( d<sub>c,k</sub>)?
- suppose we're modeling protein sequences...

#### **Amino Acids**

- Can we encode prior knowledge about amino acid properties into the motif finding process?
- there are classes of amino acids that share similar properties

			FOLAR, UNCHARGED		
Alanine Ala A MW = 89	- оос <sub>Н<sub>3</sub>№</sub> >сн	г - СН <sub>3</sub>	OUPS H-	сн <sup>- соо-</sup> № Н <sub>3</sub>	Glycine Gly G MW = 75
Valine Val V MW = 117	- оос <sub>Н<sub>3</sub>№</sub> >сн	- сң <sup>сн</sup> з снз	но-сн <sub>2</sub> -	сн <sup>соо-</sup>	Serine Ser S MW = 105
Leucine Leu L MW = 131	- оос <sub>Н<sub>3</sub>№</sub> >сн	і - сн <sub>2</sub> - сң <sup>сн</sup> 3 сн <sub>3</sub>	<sup>ОН</sup> >сн - сн <sub>3</sub> -сн -	сн <sup>&lt;соо-</sup>	Threonine Thr T MWV = 119
Isoleucine Ile I MW = 131	- оос <sub>Н<sub>3</sub>№</sub> >сн	н-сң <sup>сн</sup> 3 сн <sub>2</sub> -сн <sub>3</sub>	HS - CH <sub>2</sub>	- сн < <sup>СОО<sup>-</sup> <sup>№</sup> Н<sub>3</sub></sup>	Cysteine Cys C MW = 121
Phenylalanine Phe F MW = 131	- оос <sub>Н<sub>3</sub>№</sub> >сн	I-СН <sub>2</sub>	но - 🖉 <b>-</b> сн <sub>2</sub>	- сңС <sup>соо-</sup>	Tyrosine Tyr Y MW = 181
Tryptophan Trp W MW = 204	- оос н <sub>з</sub> ү	- сн <sub>2</sub> - с	0 C - CH2	-сн <sup>соо-</sup>	Asparagine Asp N MW = 132
Methionine Met M MW = 149	- оос <sub>Н<sub>3</sub>№</sub> _сн	- CH <sub>2</sub> - CH <sub>2</sub> - S - CH <sub>3</sub>	NH <sub>2</sub> 0 С - СН <sub>2</sub> - СН <sub>2</sub>	- сн < <sup>соо-</sup>	Glutamine Gln Q MWV = 146
Proline Pro P MW = 115	-000 C		<sup>+</sup> NH <sub>3</sub> – CH <sub>2</sub> – (СН	POLAR BASIC	Lysine Lys K MW = 146
Aspartic acid Asp D MW = 133		с I - СН <sub>2</sub> - С <sup>0</sup>	NH <sub>2</sub> NH <sub>2</sub> C - NH - (CH	<sub>2</sub> )3-сн <sup>соо</sup>	Arginine Arg R MW = 174
Glutamine acid Glu E MW = 147	- 00C H <sub>3</sub> <sup>N</sup> >CH	- сн <sub>2</sub> - сн <sub>2</sub> - с	/=Ç-CH₂- HN≫NH +	сн <sup>соо-</sup> <sup>№</sup> <sup>№</sup> <sup>№</sup>	Histidine His H MVV = 155

# **Using Dirichlet Mixture Priors**

- since we're estimating multinomial distributions (frequencies of amino acids at each motif position), a natural way to encode prior knowledge is using Dirichlet distributions
- let's consider
  - the Beta distribution
  - the Dirichlet distribution
  - mixtures of Dirichlets

#### The Beta Distribution

- suppose we're taking a Bayesian approach to estimating the parameter  $\theta$  of a weighted coin
- the Beta distribution provides an appropriate prior

$$P(\theta) = \frac{\Gamma(\alpha_h + \alpha_t)}{\Gamma(\alpha_h)\Gamma(\alpha_t)} \theta^{\alpha_h - 1} (1 - \theta)^{\alpha_t - 1}$$

where

- $\alpha_h$  # of "imaginary" heads we have seen already
- $\alpha_t$  # of "imaginary" tails we have seen already.







Beta(19,39)

#### The Beta Distribution

suppose now we're given a data set *D* in which we observe *D<sub>h</sub>* heads and *D<sub>t</sub>* tails

$$P(\theta \mid D) = \frac{\Gamma(\alpha + D_h + D_t)}{\Gamma(\alpha_h + D_h)\Gamma(\alpha_t + D_t)} \theta^{\alpha_h + D_h - 1} (1 - \theta)^{\alpha_t + D_t - 1}$$

= Beta(
$$\alpha_h$$
 +  $D_h$ , $\alpha_t$  +  $D_t$ )

 the posterior distribution is also Beta: we say that the set of Beta distributions is a *conjugate* family for binomial sampling

#### The Dirichlet Distribution

- for discrete variables with more than two possible values, we can use *Dirichlet* priors
- Dirichlet priors are a *conjugate* family for multinomial data

$$P(\theta) = \frac{\Gamma\left(\sum_{i=1}^{K} \alpha_{i}\right)}{\prod_{i=1}^{K} \Gamma(\alpha_{i})} \prod_{i=1}^{K} \theta_{i}^{\alpha_{i}-1}$$

• if  $P(\theta)$  is Dirichlet $(\alpha_1, \ldots, \alpha_K)$ , then  $P(\theta|D)$  is Dirichlet $(\alpha_1+D_1, \ldots, \alpha_K+D_K)$ , where  $D_i$  is the # occurrences of the *i*<sup>th</sup> value



# Mixture of Dirichlets

- we'd like to have Dirichlet distributions characterizing amino acids that tend to be used in certain "roles"
- Brown et al. [ISMB '95] induced a set of Dirichlets from "trusted" protein alignments
  - "large, charged and polar"
  - "polar and mostly negatively charged"
  - "hydrophobic, uncharged, nonpolar"
  - etc.

# <section-header>Construction of the optimiser of the opt

# Using Dirichlet Mixture Priors

• recall that the EM/Gibbs update the parameters by:

$$p_{c,k} = \frac{n_{c,k} + d_{c,k}}{\sum_{b} (n_{b,k} + d_{b,k})}$$

 we can set the pseudocounts using a *mixture* of Dirichlets:

$$d_{c,k} = \sum_{j} P(\alpha^{(j)} | \mathbf{n}_{k}) \alpha_{c}^{(j)}$$

• where  $\alpha^{(j)}$  is the *j*<sup>th</sup> Dirichlet component



# Motif Finding: EM and Gibbs

- these methods compute local, multiple alignments
- both methods try to optimize the likelihood of the sequences
- EM converges to a local maximum
- Gibbs will converge to a global maximum, *in the limit;* in a reasonable amount of time, probably not
- can take advantage of background knowledge by
  - tying parameters
  - Dirichlet priors
- there are many other methods for motif finding
- in practice, motif finders often fail
  - motif "signal" may be weak
  - large search space, many local minima