### Nebraska Nebraska Introduction CSCE 478/878 Lecture 9: Hidden Markov Models Stephen Sco Useful for modeling/making predictions on sequential Introduction ntroduction data Dutline Outline • E.g., biological sequences, text, series of larkov hains Stephen Scott larkov hains sounds/spoken words lidden • Will return to graphical models that are generative sscott@cse.unl.edu ・ロン・西と・川川・山・山・

#### Markov Chains Nebraska Nebraska Outline Bioinformatics Example: CpG Islands 178/87 Hidder larko Model • Focus on nucleotide sequences: Sequences of tephen Sc Markov chains symbols from alphabet {A, C, G, T} ntroductior ntroduction Hidden Markov models (HMMs) • The sequence "CG" (written "CpG") tends to appear Outline Dutline Formal definition more frequently in some places than in others Markov Chains larkov hains • Finding most probable state path (Viterbi algorithm) • Such **CpG islands** are usually 10<sup>2</sup>-10<sup>3</sup> bases long • Forward and backward algorithms Questions: Specifying an HMM Given a short segment, is it from a CpG island? Oiven a long segment, where are its islands? ・



# Markov Chains Nebraska Modeling CpG Islands (cont'd) enhen So

Dutline

P(A | T)

### 

## Nebraska

tephen Sco

ntroduction

Nebraska

tephen Sco

troduction Dutline

Hidden

/odels

What's Hidden?

### Markov Chains The Markov Property

- A first-order Markov model (what we study) has the property that observing symbol  $\mathbf{x}_i$  while in state  $\pi_i$ depends **only** on the previous state  $\pi_{i-1}$  (which generated  $\mathbf{x}_{i-1}$ )
- Standard model has 1-1 correspondence between symbols and states, thus

 $P(\mathbf{x}_i \mid \mathbf{x}_{i-1}, \dots, \mathbf{x}_1) = P(\mathbf{x}_i \mid \mathbf{x}_{i-1})$ 

and

$$P(\mathbf{x}_1,\ldots,\mathbf{x}_L) = P(\mathbf{x}_1) \prod_{i=2}^L P(\mathbf{x}_i \mid \mathbf{x}_{i-1})$$

#### Markov Chains Nebraska Begin and End States

- For convenience, can add special "begin" (B) and "end" (E) states to clarify equations and define a distribution over sequence lengths
- Emit empty (null) symbols  $\mathbf{x}_0$  and  $\mathbf{x}_{L+1}$  to mark ends of sequence



#### Markov Chains Nebraska Markov Chains for Discrimination

- How do we use this to differentiate islands from non-islands?
- Define two Markov models: islands ("+") and non-islands ("-")
  - Each model gets 4 states (A, C, G, T)
  - Take training set of known islands and non-islands
  - Let  $c_{st}^+$  = number of times symbol *t* followed symbol *s* in an island:

$$\hat{P}^+(t \mid s) = rac{c_{st}^+}{\sum_{t'} c_{st'}^+}$$

• Now score a sequence  $X = \langle \mathbf{x}_1, \dots, \mathbf{x}_L \rangle$  by summing the log-odds ratios:

$$\log\left(\frac{\hat{P}(X\mid+)}{\hat{P}(X\mid-)}\right) = \sum_{i=1}^{L+1} \log\left(\frac{\hat{P}^+(\mathbf{x}_i\mid\mathbf{x}_{i-1})}{\hat{P}^-(\mathbf{x}_i\mid\mathbf{x}_{i-1})}\right)$$

#### Nebraska Hidden Markov Models

478/87

/larko

ntroduction

Hidden arkov odels

Second CpG question: Given a long sequence, where are its islands?

- Could use tools just presented by passing a fixed-width window over the sequence and computing scores
- Trouble if islands' lengths vary
- Prefer single, unified model for islands vs. non-islands



 Within the + group, transition probabilities similar to those for the separate + model, but there is a small chance of switching to a state in the - group







#### Nebřaška Linde Hidden Markov Models Example: The Occasionally Dishonest Casino

Stephen Sco

Dutline



Given a sequence of rolls, what's hidden?

#### Nebraska Linden Markov Models The Viterbi Algorithm

Outline

• Probability of seeing symbol sequence X and state sequence  $\pi$  is

$$P(X, \pi) = P(\pi_1 \mid 0) \prod_{i=1}^{L} P(\mathbf{x}_i \mid \pi_i) P(\pi_{i+1} \mid \pi_i)$$

• Can use this to find most likely path:

$$\pi^* = \operatorname{argmax} P(X, \pi)$$

and trace it to identify islands (paths through "+" states)

・ヨン・ヨン・ヨン

• There are an exponential number of paths through chain, so how do we find the most likely one?



ntroductior

Dutline

Nebraska

### Hidden Markov Models The Forward Algorithm

Given a sequence *X*, find  $P(X) = \sum_{\pi} P(X, \pi)$ 

Use dynamic programming like Viterbi, replacing max with sum, and  $v_k(i)$  with  $f_k(i) = P(\mathbf{x}_1, \dots, \mathbf{x}_i, \pi_i = k)$  (= prob. of observed sequence through  $\mathbf{x}_i$ , stopping in state k)

f<sub>0</sub>(0) = 1, f<sub>k</sub>(0) = 0 for k > 0
For i = 1 to L; for l = 1 to M (# states)
f<sub>l</sub>(i) = P(x<sub>i</sub> | l) ∑<sub>k</sub> f<sub>k</sub>(i − 1)P(l | k)

• 
$$P(X) = \sum_{k} f_k(L) P(0 \mid k)$$

To avoid underflow, can again use logs, though exactness of results compromised

< ロ > < 個 > < 三 > < 三 > < 三 > の へ の



### Nebraska

Dutline

larkov hains

### Hidden Markov Models Example Use of Forward/Backward Algorithm

- Define g(k) = 1 if k ∈ {A<sub>+</sub>, C<sub>+</sub>, G<sub>+</sub>, T<sub>+</sub>} and 0 otherwise
- Then  $G(i \mid X) = \sum_{k} P(\pi_i = k \mid X) g(k)$  = probability that  $\mathbf{x}_i$  is in an island
- For each state k, compute  $P(\pi_i = k \mid X)$  with forward/backward algorithm
- Technique applicable to any HMM where set of states is partitioned into classes

• Use to label individual parts of a sequence

#### Nebraska Linden Markov Models Specifying an HMM

- Two problems: defining structure (set of states) and parameters (transition and emission probabilities)
- Start with latter problem, i.e., given a training set  $X_1, \ldots, X_N$  of independently generated sequences, learn a good set of parameters  $\theta$
- Goal is to maximize the (log) likelihood of seeing the training set given that θ is the set of parameters for the HMM generating them:





Nebraska

Dutline

arkov hains

Nebraska
----------

Dutline Aarkov Chains

### Hidden Markov Models Specifying an HMM: The Baum-Welch Algorithm

- Used for estimating params when state seq unknown
- Special case of expectation maximization (EM)
- Start with arbitrary P(l | k) and P(b | k), and use to estimate A<sub>kl</sub> and E<sub>k</sub>(b) as expected number of occurrences given the training set<sup>1</sup>:

$$A_{k\ell} = \sum_{j=1}^{N} \frac{1}{P(X_j)} \sum_{i=1}^{L} f_k^j(i) P(\ell \mid k) P(\mathbf{x}_{i+1}^j \mid \ell) b_\ell^j(i+1)$$

(Prob. of transition from k to  $\ell$  at position i of sequence j, summed over all positions of all sequences)

Specifying an HMM: The Baum-Welch Algorithm (cont'd)

Hidden Markov Models

$$E_k(b) = \sum_{j=1}^N \sum_{i:\mathbf{x}_i^j = b} P(\pi_i = k \mid X_j) = \sum_{j=1}^N \frac{1}{P(X_j)} \sum_{i:\mathbf{x}_i^j = b} f_k^j(i) b_k^j(i)$$

- Use these (& pseudocounts) to recompute  $P(\ell \mid k)$  and  $P(b \mid k)$
- After each iteration, compute log likelihood and halt if no improvement

### Nebraska

Stephen Sco

Outline

/larkov Chains

idden

### Hidden Markov Models Specifying an HMM: Structure

- How to specify HMM states and connections?
- States come from background knowledge on problem, e.g., size-4 alphabet, +/-, ⇒ 8 states
- Connections:
  - Tempting to specify complete connectivity and let Baum-Welch sort it out
  - Problem: Huge number of parameters could lead to local max
  - Better to use background knowledge to invalidate some connections by initializing  $P(\ell \mid k) = 0$ 
    - Baum-Welch will respect this

#### Nebraska Lindon Markov Models Specifying an HMM: Silent States

ture 9: idden

tenhen Sr

Outline

arkov hains

- May want to allow model to generate sequences with certain parts deleted
  - E.g., when aligning DNA or protein sequences against a fixed model or matching a sequence of spoken words against a fixed model, some parts of the input might be omitted



• **Problem:** Huge number of connections, slow training, local maxima

・ロット 4回マスポット 中マットロッ

Nebraska

Hidden I	Markov	Models	
Specifying an HMM: Silent States (cont'd)			



• Silent states (like begin and end states) don't emit symbols, so they can "bypass" a regular state



- If there are no purely silent loops, can update Viterbi, forward, and backward algorithms to work with silent states
- Used extensively in **profile HMMs** for modeling sequences of protein families (aka **multiple alignments**)

・ロト・(四ト・(三ト・(三)・)へ(の)

・ロト・西ト・モー・モー ひゃう