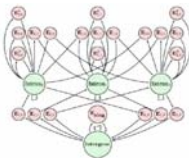# Computational Genomics

**10-810/02-710, Spring 2009**

## Gene Finding and HMM

**Eric Xing**

**Lecture 5, January 28, 2009**

**Reading:** Durbin Chap 3, class assignment

1

---

# The HMM Algorithms

**Questions:**

- **Decoding**: What is the most likely DNA parsing?  Viterbi
- **Evaluation**: What is the probability of the observed sequence?  Forward
- **Decoding**: What is the probability that the state of the 3rd position is Bk or gene, given the observed sequence? Forward-Backward
- **Learning**: Under what parameterization are the observed sequences most probable? Baum-Welch (EM)

2

# The likelihood of a sequence

- We want to calculate $P(\mathbf{x})$, the likelihood of $\mathbf{x}$, given the HMM
  - Sum over all possible ways of generating $\mathbf{x}$:

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} \pi_{y_1} \prod_{t=2}^{T} a_{y_{t-1}, y_t} \prod_{t=1}^{T} p(x_t \mid y_t)$$

- Complexity?

- Why useful?

3

---

# The Forward Algorithm

- We want to calculate $P(\mathbf{x})$, the likelihood of $\mathbf{x}$, given the HMM
  - Sum over all possible ways of generating $\mathbf{x}$:

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \cdots \sum_{y_N} \pi_{y_1} \prod_{t=2}^{T} a_{y_{t-1}, y_t} \prod_{t=1}^{T} p(x_t \mid y_t)$$

  - To avoid summing over an exponential number of paths $\mathbf{y}$, define

$$\alpha(y_t^k = 1) = \alpha_t^k \stackrel{\text{def}}{=} P(x_1, \ldots, x_t, y_t^k = 1) \qquad \text{(the forward probability)}$$
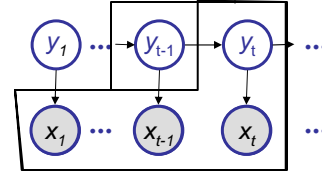
  - The recursion:

$$\alpha_t^k = p(x_t \mid y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

$$P(\mathbf{x}) = \sum_k \alpha_T^k$$

4

2

# The Forward Algorithm – derivation

- Compute the forward probability:



$$\alpha_t^k = P(x_1, ..., x_{t-1}, x_t, y_t^k = 1)$$

$$= \sum_{y_{t-1}} P(x_1, ..., x_{t-1}, x_t, y_{t-1}, y_t^k = 1)$$

$$= \sum_{y_{t-1}} P(x_1, ..., x_{t-1}, y_{t-1}) P(y_t^k = 1 \mid y_{t-1}, x_1, ..., x_{t-1}) P(x_t \mid y_t^k = 1, x_1, ..., x_{t-1}, y_{t-1})$$

$$= \sum_{y_{t-1}} P(x_1, ..., x_{t-1}, y_{t-1}) P(y_t^k = 1 \mid y_{t-1}) P(x_t \mid y_t^k = 1)$$

$$= P(x_t \mid y_t^k = 1) \sum_i P(x_1, ..., x_{t-1}, y_{t-1}^i = 1) P(y_t^k = 1 \mid y_{t-1}^i = 1)$$

$$= P(x_t \mid y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

Chain rule : $P(A, B, C) = P(A) P(B \mid C) P(C \mid A, B)$

---

# The Forward Algorithm

- We can compute $\alpha_t^k$ for all $k$, $t$, using dynamic programming!

**Initialization:**

$$\alpha_1^k = P(x_1 \mid y_1^k = 1) \pi_k$$

$$\alpha_1^k = P(x_1, y_1^k = 1)$$
$$= P(x_1 \mid y_1^k = 1) P(y_1^k = 1)$$
$$= P(x_1 \mid y_1^k = 1) \pi_k$$

**Iteration:**

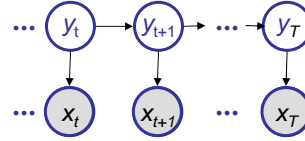$$\alpha_t^k = P(x_t \mid y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

**Termination:**

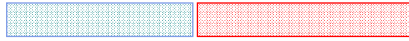$$P(\mathbf{x}) = \sum_k \alpha_T^k$$

3

# The Backward Algorithm

- We want to compute $P(y_t^k = 1 \mid \mathbf{x})$ ,

  the posterior probability distribution on the $t$ th position, given $\mathbf{x}$

  - We start by computing

$$P(y_t^k = 1, \mathbf{x}) = P(x_1,...,x_t, y_t^k = 1, x_{t+1},...,x_T)$$
$$= P(x_1,...,x_t, y_t^k = 1)P(x_{t+1},...,x_T \mid x_1,...,x_t, y_t^k = 1)$$
$$= P(x_1...x_t, y_t^k = 1)P(x_{t+1}...x_T \mid y_t^k = 1)$$

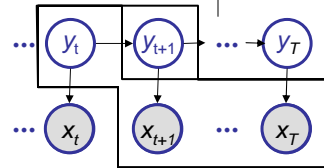**Forward, $\alpha_t^k$**     **Backward,** $\beta_t^k = P(x_{t+1},...,x_T \mid y_t^k = 1)$

  - The recursion:

$$\beta_t^k = \sum_i a_{k,i} p(x_{t+1} \mid y_{t+1}^i = 1)\beta_{t+1}^i$$

---

# The Backward Algorithm – derivation

- Define the backward probability:

$$\beta_t^k = P(x_{t+1},...,x_T \mid y_t^k = 1)$$
$$= \sum_{y_{t+1}} P(x_{t+1},...,x_T, y_{t+1} \mid y_t^k = 1)$$
$$= \sum_i P(y_{t+1}^i = 1 \mid y_t^k = 1) p(x_{t+1} \mid y_{t+1}^i = 1, y_t^k = 1) P(x_{t+2},...,x_T \mid x_{t+1}, y_{t+1}^i = 1, y_t^k = 1)$$
$$= \sum_i P(y_{t+1}^i = 1 \mid y_t^k = 1) p(x_{t+1} \mid y_{t+1}^i = 1) P(x_{t+2},...,x_T \mid y_{t+1}^i = 1)$$
$$= \sum_i a_{k,i} p(x_{t+1} \mid y_{t+1}^i = 1)\beta_{t+1}^i$$

Chain rule: $P(A,B,C \mid \alpha) = P(A,\alpha)P(B \mid C,\alpha)P(C \mid A,B,\alpha)$

# The Backward Algorithm

- We can compute $\beta_t^k$ for all $k$, $t$, using dynamic programming!

**Initialization:**

$$\beta_T^k = 1, \ \forall k$$

**Iteration:**

$$\beta_t^k = \sum_i a_{k,i} P(x_{t+1} \mid y_{t+1}^i = 1) \beta_{t+1}^i$$

**Termination:**

$$P(\mathbf{x}) = \sum_k \alpha_1^k \beta_1^k$$

# Example:

# Posterior decoding

- We can now calculate

$$P(y_t^k = 1 \mid \mathbf{x}) = \frac{P(y_t^k = 1, \mathbf{x})}{P(\mathbf{x})} = \frac{\alpha_t^k \beta_t^k}{P(\mathbf{x})}$$

- Then, we can ask
  - What is the most likely state at position $t$ of sequence $\mathbf{x}$:

$$k_t^* = \arg\max_k P(y_t^k = 1 \mid \mathbf{x})$$

  - Note that this is an MPA of a single hidden state, what if we want to a MPA of a whole hidden state sequence?
  - Posterior Decoding: $\left\{ y_t^{k_t^*} = 1 : t = 1 \cdots T \right\}$
  - This is different from MPA of a whole sequence of hidden states
  - This can be understood as *bit error rate* vs. *word error rate*

**Example:**
**MPA of X ?**
**MPA of (X, Y) ?**

| x | y | P(x,y) |
|---|---|--------|
| 0 | 0 | 0.35 |
| 0 | 1 | 0.05 |
| 1 | 0 | 0.3 |
| 1 | 1 | 0.3 |

---

# Computational Complexity and implementation details

- What is the running time, and space required, for Forward, and Backward?

$$\alpha_t^k = p(x_t \mid y_t^k = 1) \sum_i \alpha_{t-1}^i a_{i,k}$$

$$\beta_t^k = \sum_i a_{k,i} p(x_{t+1} \mid y_{t+1}^i = 1) \beta_{t+1}^i$$

$$V_t^k = p(x_t \mid y_t^k = 1) \max_i a_{i,k} V_{t-1}^i$$

Time: $O(K^2 N)$;          Space: $O(KN)$.

- Useful implementation technique to avoid underflows
  - Viterbi:               sum of logs
  - Forward/Backward:   rescaling at each position by multiplying by a constant

# Learning HMM: two scenarios

- **Supervised learning**: estimation when the "right answer" is known
  - **Examples:**
    GIVEN: a genomic region x = $x_1 \ldots x_{1,000,000}$ where we have good (experimental) annotations of the CpG islands
    GIVEN: the casino player allows us to observe him one evening, as he changes dice and produces 10,000 rolls

- **Unsupervised learning**: estimation when the "right answer" is unknown
  - **Examples:**
    GIVEN: the porcupine genome; we don't know how frequent are the CpG islands there, neither do we know their composition
    GIVEN: 10,000 rolls of the casino player, but we don't see when he changes dice

- **QUESTION:** Update the parameters $\theta$ of the model to maximize $P(x|\theta)$ --- Maximal likelihood (ML) estimation

© Eric Xing @ CMU, 2005-2009

© Eric Xing @ CMU, 2005-2009

---

# Supervised ML estimation

- Given $x = x_1 \ldots x_N$ for which the true state path $y = y_1 \ldots y_N$ is known,
  - **Define:**
    $A_{ij}$ = # times state transition $i \rightarrow j$ occurs in **y**
    $B_{ik}$ = # times state $i$ in **y** emits $k$ in **x**

  - We can show that the **maximum likelihood** parameters $\theta$ are:

$$a_{ij}^{ML} = \frac{\#(i \rightarrow j)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=2}^{T} y_{n,t-1}^i y_{n,t}^j}{\sum_n \sum_{t=2}^{T} y_{n,t-1}^i} = \frac{A_{ij}}{\sum_{j'} A_{ij'}}$$

$$b_{ik}^{ML} = \frac{\#(i \rightarrow k)}{\#(i \rightarrow \bullet)} = \frac{\sum_n \sum_{t=1}^{T} y_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^{T} y_{n,t}^i} = \frac{B_{ik}}{\sum_{k'} B_{ik'}}$$

  - What if y is continuous? We can treat $\{(x_{n,t}, y_{n,t}) : t = 1:T, n = 1:N\}$ as $N \times T$ observations of, e.g., a Gaussian, and apply learning rules for Gaussian ...

7

# Supervised ML estimation, ctd.

- **Intuition:**
  - When we know the underlying states, the best estimate of $\theta$ is the average frequency of transitions & emissions that occur in the training data

- **Drawback:**
  - Given little data, there may be **overfitting**:
    - $P(x|\theta)$ is maximized, but $\theta$ is unreasonable
      - **0 probabilities – VERY BAD**

- **Example:**
  - Given 10 casino rolls, we observe
    $$\text{x = 2, 1, 5, 6, 1, 2, 3, 6, 2, 3}$$
    $$\text{y = F, F, F, F, F, F, F, F, F, F}$$
  - Then: $a_{FF} = 1;$ $a_{FL} = 0$
    $b_{F1} = b_{F3} = .2;$
    $b_{F2} = .3; b_{F4} = 0; b_{F5} = b_{F6} = .1$

# Pseudocounts

- Solution for small training sets:
  - Add pseudocounts
    $A_{ij}$ = # times state transition $i \rightarrow j$ occurs in $\mathbf{y}$ + $R_{ij}$
    $B_{ik}$ = # times state $i$ in $\mathbf{y}$ emits $k$ in $\mathbf{x}$ + $S_{ik}$

  - $R_{ij}, S_{ij}$ are pseudocounts representing our prior belief
  - Total pseudocounts: $R_i = \Sigma_j R_{ij}$ , $S_i = \Sigma_k S_{ik}$ ,
    - --- "strength" of prior belief,
    - --- total number of imaginary instances in the prior

- Larger total pseudocounts $\Rightarrow$ strong prior belief

- Small total pseudocounts: just to avoid 0 probabilities --- smoothing

# Unsupervised ML estimation

- Given $x = x_1 \ldots x_N$ for which the true state path $y = y_1 \ldots y_N$ is **unknown**,

  - **EXPECTATION MAXIMIZATION**

  0. Starting with our best guess of a model $M$, parameters $\theta$:
  1. Estimate $A_{ij}$, $B_{ik}$ in the training data
     - How? $A_{ij} = \sum_{n,t} \left\langle y_{n,t-1}^i y_{n,t}^j \right\rangle$ $\quad B_{ik} = \sum_{n,t} \left\langle y_{n,t}^i \right\rangle x_{n,t}^k$,
  2. Update $\theta$ according to $A_{ij}$, $B_{ik}$
     - Now a "supervised learning" problem
  3. Repeat 1 & 2, until convergence

  **This is called the Baum-Welch Algorithm**

  We can get to a provably more (or equally) likely parameter set $\theta$ each iteration

---

# How to compute expected count?

$$B_{ik} = \sum_{n,t} \left\langle y_{n,t}^i \right\rangle x_{n,t}^k$$

$$\langle y_{n,t}^i \rangle = P(Y_{n,t}^i = 1 | \mathbf{x}_n)$$

$$= \frac{\alpha_{n,t}^i \beta_{n,t}^i}{P(\mathbf{x}_n)}$$

$$A_{ij} = \sum_{n,t} \left\langle y_{n,t-1}^i y_{n,t}^j \right\rangle$$

$$\langle y_{n,t-1}^i y_{n,t}^j \rangle = P(Y_{n,t-1}^i = 1, Y_{n,t}^j = 1 | \mathbf{x}_n)$$

$$= \frac{\alpha_{n,t-1}^i a_{i,j} x_{n,t}^j \beta_{n,t}^j}{P(\mathbf{x}_n)}$$

# The Baum Welch algorithm

- The complete log likelihood

$$\ell_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y}) = \log p(\mathbf{x}, \mathbf{y}) = \log \prod_n \left( p(y_{n,1}) \prod_{t=2}^{T} p(y_{n,t} \mid y_{n,t-1}) \prod_{t=1}^{T} p(x_{n,t} \mid x_{n,t}) \right)$$

- The expected complete log likelihood

$$\langle \ell_c(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y}) \rangle = \sum_n \left( \langle y_{n,1}^i \rangle_{p(y_{n,1}|\mathbf{x}_n)} \log \pi_i \right) + \sum_n \sum_{t=2}^{T} \left( \langle y_{n,t-1}^i y_{n,t}^j \rangle_{p(y_{n,t-1}, y_{n,t}|\mathbf{x}_n)} \log a_{i,j} \right) + \sum_n \sum_{t=1}^{T} \left( x_{n,t}^k \langle y_{n,t}^i \rangle_{p(y_{n,t}|\mathbf{x}_n)} \log b_{i,k} \right)$$

- EM
  - The E step

  $$\gamma_{n,t}^i = \langle y_{n,t}^i \rangle = p(y_{n,t}^i = 1 \mid \mathbf{x}_n)$$
  $$\xi_{n,t}^{i,j} = \langle y_{n,t-1}^i y_{n,t}^j \rangle = p(y_{n,t-1}^i = 1, y_{n,t}^j = 1 \mid \mathbf{x}_n)$$

  - The M step ("symbolically" identical to MLE)

  $$\pi_i^{ML} = \frac{\sum_n \gamma_{n,1}^i}{N} \qquad a_{ij}^{ML} = \frac{\sum_n \sum_{t=2}^{T} \xi_{n,t}^{i,j}}{\sum_n \sum_{t=1}^{T-1} \gamma_{n,t}^i} \qquad b_{ik}^{ML} = \frac{\sum_n \sum_{t=1}^{T} \gamma_{n,t}^i x_{n,t}^k}{\sum_n \sum_{t=1}^{T-1} \gamma_{n,t}^i}$$

19

---

# The Baum-Welch algorithm -- comments

Time Complexity:

  # iterations $\times$ O(K$^2$N)

- Guaranteed to increase the log likelihood of the model

- Not guaranteed to find globally best parameters

- Converges to local optimum, depending on initial conditions

- Too many parameters / too large model:  Overt-fitting

20

10

# Higher-order HMMs

- **The Genetic Code**

  - 3 nucleotides make 1 amino acid

  - Statistical dependencies in triplets

- **Question:**

  - Recognize protein-coding segments with an HMM

| | U | C | A | G |
|---|---|---|---|---|
| U | UUU phe<br>UUC<br>UUA leu<br>UUG | UCU<br>UCC ser<br>UCA<br>UCG | UAU tyr<br>UAC<br>UAA *Stop*<br>UAG *Stop* | UGU cys<br>UGC<br>UGA *Stop*<br>UGG *Stop* |
| C | CUU<br>CUC leu<br>CUA<br>CUG | CCU<br>CCC pro<br>CCA<br>CCG | CAU his<br>CAC<br>CAA gln<br>CAG | CGU<br>CGC arg<br>CGA<br>CGG |
| A | AUU<br>AUC ile<br>AUA<br>AUG met | ACU<br>ACC thr<br>ACA<br>ACG | AAU asn<br>AAC<br>AAA lys<br>AAG | AGU ser<br>AGC<br>AGA arg<br>AGG |
| G | GUU<br>GUC val<br>GUA<br>GUG | GCU<br>GCC ala<br>GCA<br>GCG | GAU asp<br>GAC<br>GAA glu<br>GAG | GGU<br>GGC gly<br>GGA<br>GGG |

---

# Higher-order HMMs

- Every state of the HMM emits 1 nucleotide

- Transition probabilities:

  Probability of a state at one position, given those of 3 previous positions (triplets):
  $$P(y_i \mid y_{i-1}, y_{i-2}, y_{i-3})$$

- Emission probabilities:
  $$P(x_i \mid y_i)$$

- Algorithms extend with small modifications

$y_{1,...,N} = i$     i     e     e     i

$y_1 \rightarrow y_2 \rightarrow y_3 \rightarrow y_4 \cdots y_N$

$x_1 \quad x_2 \quad x_3 \quad x_4 \cdots x_N$

$x_{1,...,N} = A$    C    T    T    G

# Inference on Higher-order HMMs

- Building 1st-order HMM on "mega" state

- Use FB algorithm as usual

  - $P(Q_2|R)$

  $\rightarrow P(Y_2, Y_3, Y_4|X)$

  $\rightarrow P(Y_3|X) = \Sigma_{y2,y4} P(Y_2, Y_3, Y_4|X)$

$y_{1,\ldots,N} =$    **i**    **i**    **e**    **e**    **i**

$y_1 \rightarrow y_2 \rightarrow y_3 \rightarrow y_4 \rightarrow \ldots \rightarrow y_N$

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad \ldots \quad x_N$

$x_{1,\ldots,N} =$    **A**    **C**    **T**    **T**    **G**

$Q_1 \qquad Q_2 \qquad Q_3$

$(y_1, y_2, y_3) \rightarrow (y_2, y_3, y_4) \rightarrow (y_3, y_4, y_5) \ldots$

$(x_1, x_2, x_3) \quad (x_2, x_3, x_4) \quad (x_3, x_4, x_5) \ldots$

$R_1 \qquad R_2 \qquad R_3$

23

---

# Modeling the Duration of States

- Length distribution of region X:

  $E[l_X] = 1/(1-p)$

$1-p$

$p \quad$ **X** $\qquad$ **Y** $\quad q$

$1-q$

- Geometric distribution, with mean 1/(1-p)
  - (homework: derive this)

  - This is a significant disadvantage of HMMs
  - Several solutions exist for modeling different length distributions

$p$

duration

24

12

# Observed Duration Time

# Poisson Point Process

- A counting process that represents the total number of occurrences of discrete events during a temporal/spatial interval

  - the number of occurrences in any internal of length $\tau$ is Poisson distributed with parameter $\lambda\tau$:

$$p(A(t+\tau) - A(n) = n) = e^{-\lambda\tau}\frac{(\lambda\tau)^n}{n!}$$



  - the number of occurrences in disjoint intervals are independent

  - the duration of the interval between two consecutive occurrences has the following distribution:

$$p(\tau < s) = 1 - e^{-\lambda s}$$

# Poisson point process



$m = \lambda \tau$

Truncation is needed at both ends!

---

# Generalized HMM

Upon entering a state:

1. Choose duration d, according to probability distribution
2. Generate d letters according to emission probs
3. Take a transition to next state according to transition probs



Disadvantage: Increase in complexity:

Time: $O(D^2)$
Space: $O(D)$

where D = maximum duration of state

## Comparative Genomics

Human

Mouse

29

## A pairwise comparison between human and mouse genome

30

15

## Aligning One Locus



Exon 1  Exon 2  Exon 3  Exon 4
Intron 1  Intron 2  Intron 3

Splice site
GGTGAG

Splice site
CAG

Translation
Initiation
ATG

Branchpoint
CTG**A**C

Stop codon
TAG/TGA/TAA

31

## Three Pairwise Alignments

32

16

## Example: a human/mouse ortholog



**Human Locus**

Alignment:

CDS

**Mouse Locus**

- coding exons
- noncoding exons
- introns
- intergenic regions
- intergenic regions
- strong alignment
- weak alignment

33

## Paired HMM

Alignments correspond 1-to-1 with sequences of states M, I, J



M (+1, +1)

I (+1, 0)

J (0, +1)

```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---
TAG-CTATCAC--GACCGC-GGTCGATTTGCCCGACC
IMMJMMMMMMMJJMMMMMMMJMMMMMMMIIMMMMMIII
```

34

17

# Let's score the transitions

Alignments correspond 1-to-1 with sequences of states M, I, J



```
-AGGCTATCACCTGACCTCCAGGCCGA--TGCCC---
TAG-CTATCAC--GACCGC-GGTCGATTTGCCCGACC
IMMJMMMMMMMJJMMMMMMMJMMMMMMMIIMMMMMIII
```

# A Pair HMM for alignments

# Gene Finding

# Generalized HMM Gene finder

19

# Generalized Pair-HMM gene finder

# Hierarchical state transition in pHMM

## Allowing for inserted exons

## Acknowledgments

- **Serafim Batzoglou**: for some of the slides adapted or modified from his lecture slides at Stanford University
- **Lior Pachter'**: for some of the slides modified from his lectures at UC Berkeley