# Computational Genomics

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

## Gene Finding and HMM (cont'd)

**Eric Xing**

**Lecture 7, February 4, 2009**

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

# 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
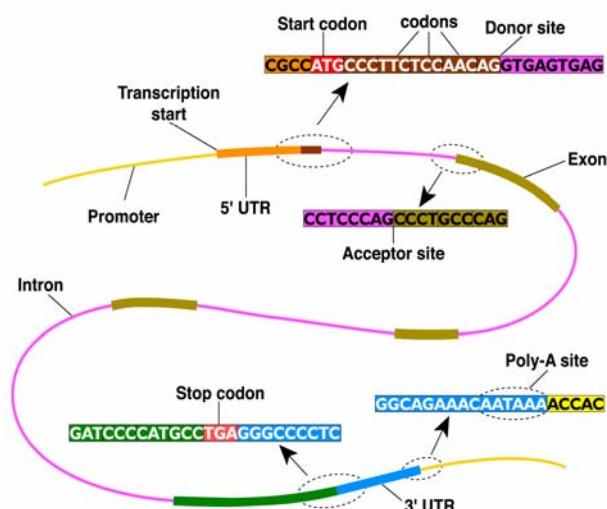    GIVEN: 10,000 rolls of the casino player, but we don't see when he changes dice

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

# Typical structure of a gene

## Some Facts About Human Genes

- Comprise about 3% of the genome
- Average gene length: ~ 8,000 bp
- Average of 5-6 exons/gene
- Average exon length: ~200 bp
- Average intron length: ~2,000 bp
- ~8% genes have a single exon

- Some exons can be as small as 1 or 3 bp.
  - HUMFMR1S is not atypical: 17 exons 40-60 bp long, comprising 3% of a 67,000 bp gene

5

## Supervised ML estimation

6

3

# 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 $\left\{ (x_{n,t}, y_{n,t}) : t = 1:T, n = 1:N \right\}$ as $N \times T$ observations of, e.g., a Gaussian, and apply learning rules for Gaussian …**

---

# 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
    ```
    x = 2, 1, 5, 6, 1, 2, 3, 6, 2, 3
    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 $i \mapsto j$ 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? $A_{ij} = \sum_{n,t} \langle y_{n,t-1}^i y_{n,t}^j \rangle$    $B_{ik} = \sum_{n,t} \langle y_{n,t}^i \rangle x_{n,t}^k$ ,
  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?

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

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

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

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

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

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

11

---

# The Baum Welch algorithm

- The complete log likelihood

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

- The expected complete log likelihood

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

- EM
  - The E step

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

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

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

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

12

6

# The Baum-Welch algorithm -- comments

Time Complexity:
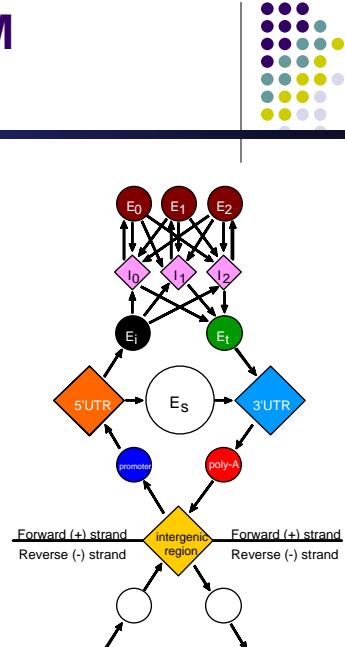
$$\text{\# iterations} \times O(K^2 N)$$

- 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

13

---

# The Idea Behind a GHMM GeneFinder

- **States** represent standard gene features: intergenic region, exon, intron, perhaps more (promotor, 5'UTR, 3'UTR, Poly-A,..).

- **Observations** embody state-dependent base composition, dependence, and signal features.

- In a GHMM, **duration** must be included as well.

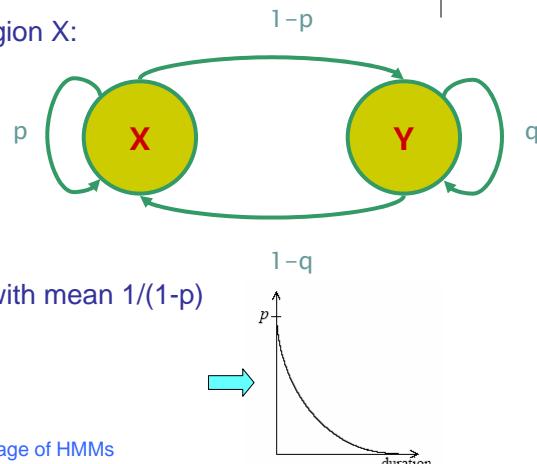- Finally, **reading frames** and **both strands** must be dealt with.

14

7

# Modeling the Duration of States

- Length distribution of region X:

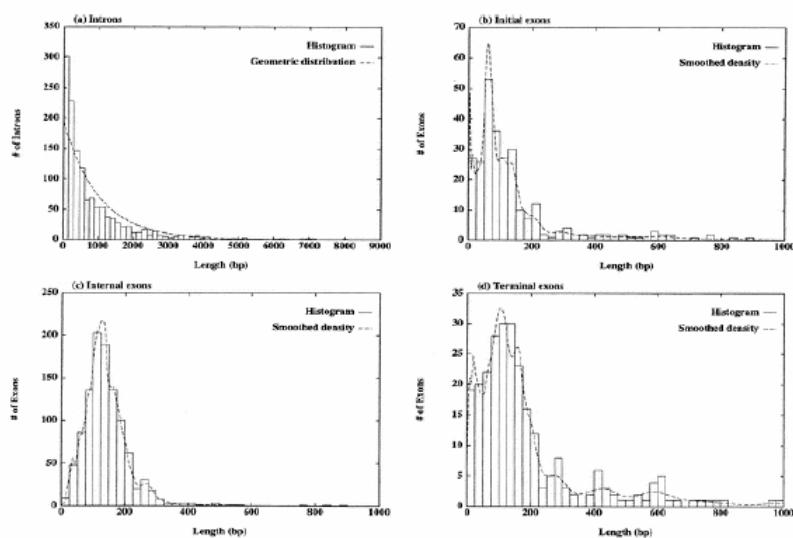  $E[l_X] = 1/(1-p)$



- 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

15

# Observed Duration Time



16

8

# 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 $\lambda\tau$:

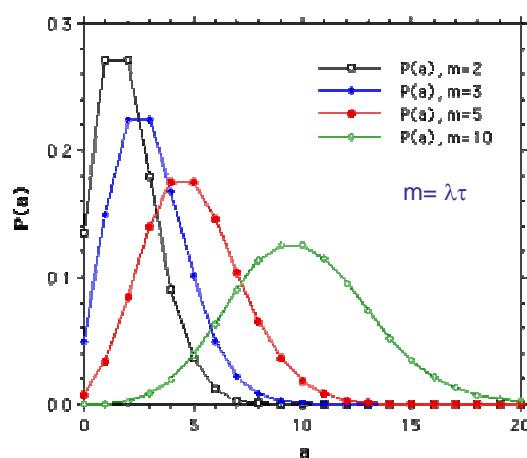$$p(A(t+\tau) - A(n) = n) = e^{-\lambda\tau}\frac{(\lambda\tau)^n}{n!}$$

  - 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!
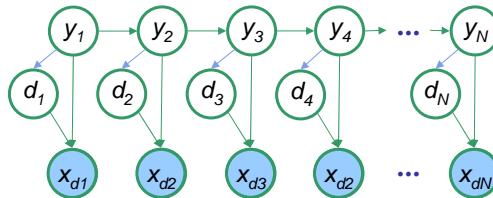
# Generalized HMM

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

19

---

# Higher-order HMMs

- **The Genetic Code**

  - 3 nucleotides make 1 amino acid

  - Statistical dependencies in triplets

- **Question:**

  - Recognize protein-coding segments with an HMM



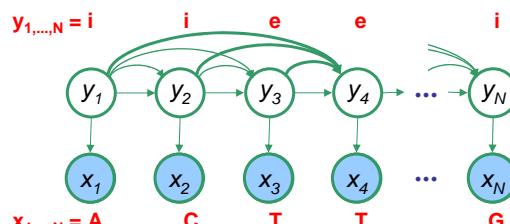|   | U | C | A | G |
|---|---|---|---|---|
| U | UUU / UUC phe<br>UUA / UUG leu | UCU / UCC / UCA / UCG ser | UAU / UAC tyr<br>UAA *Stop*<br>UAG *Stop* | UGU / UGC cys<br>UGA *Stop*<br>UGG *Stop* |
| C | CUU / CUC / CUA / CUG leu | CCU / CCC / CCA / CCG pro | CAU / CAC his<br>CAA / CAG gln | CGU / CGC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

20

10

# 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

$y_{1,\dots,N} = i \quad i \quad e \quad e \quad i$

$x_{1,\dots,N} = A \quad C \quad T \quad T \quad G$

21

# Inference on Higher-order HMMs

- Building 1$^{st}$-order HMM on "mega" state

- Use FB algorithm as usual
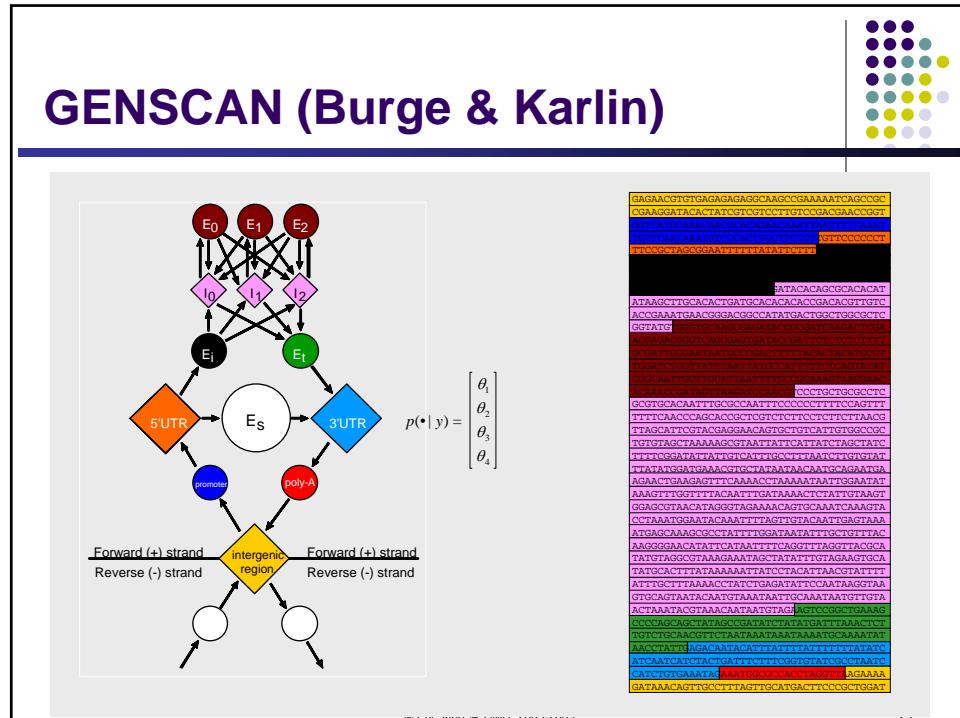
  - $P(Q_2 \mid R)$

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

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

$y_{1,\dots,N} = i \quad i \quad e \quad e \quad i$

$x_{1,\dots,N} = A \quad C \quad T \quad T \quad G$

$Q_1 \quad Q_2 \quad Q_3$

$y_1, y_2, y_3 \quad y_2, y_3, y_4 \quad y_3, y_4, y_5 \quad \dots$

$X_1, X_2, X_3 \quad X_2, X_3, X_4 \quad X_3, X_4, X_5 \quad \dots$

$R_1 \quad R_2 \quad R_3$

22

11

# GENSCAN (Burge & Karlin)



© Eric Xing @ CMU, 2005-2009

# Gene Finding



© Eric Xing @ CMU, 2005-2009                                    24

# Generalized HMM Gene finder

TAAT ATGTCCACGG GTATTGAG CATTGTACACGGG GTATTGAG CATGTAA TGAA

25

# Allowing for inserted exons

26

# Summary

**The HMM Algorithm:**

- **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)

27

# 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

28