

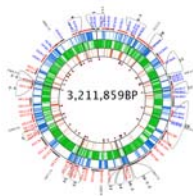
Computational Genomics

10-810/02-710, Spring 2009

Gene Finding and HMM

Eric Xing

Lecture 4, January 26, 2009



Reading: Durbin chap 3.

© Eric Xing @ CMU, 2005-2009

1

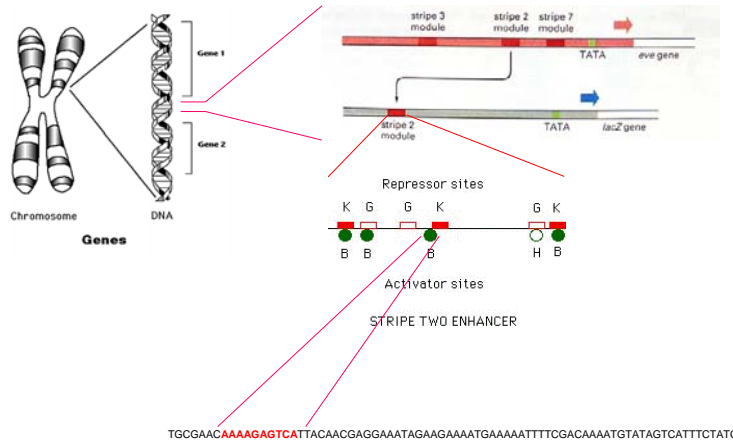
Please correct

- 10-810-09s-instr@cs should be 10810-09s-instr@cs
- 10-810-09s-announce@cs should be 10810-09s-announce@cs.

© Eric Xing @ CMU, 2005-2009

2

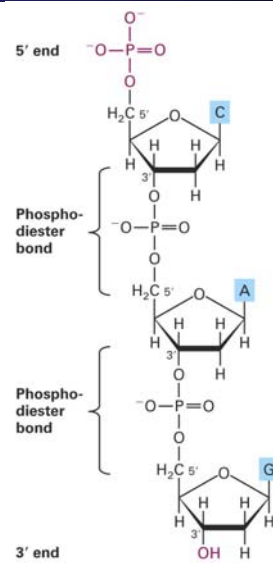
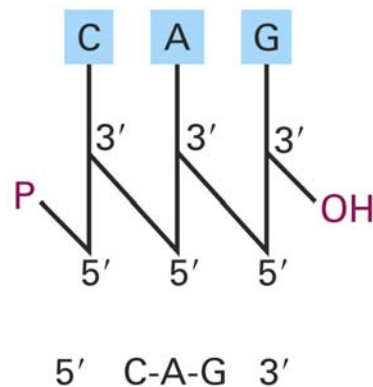
Hierarchical structure of the genome



© Eric Xing @ CMU, 2005-2009

3

The DNA strand has a chemical polarity



© Eric Xing @ CMU, 2005-2009

Writing DNA sequence



- One strand is written by listing its bases in 5' to 3' order

5' ACCGTTACT 3'

- Each strand uniquely determines the complementary strand, which runs in the opposite direction:

5' ACCGTTACT 3'

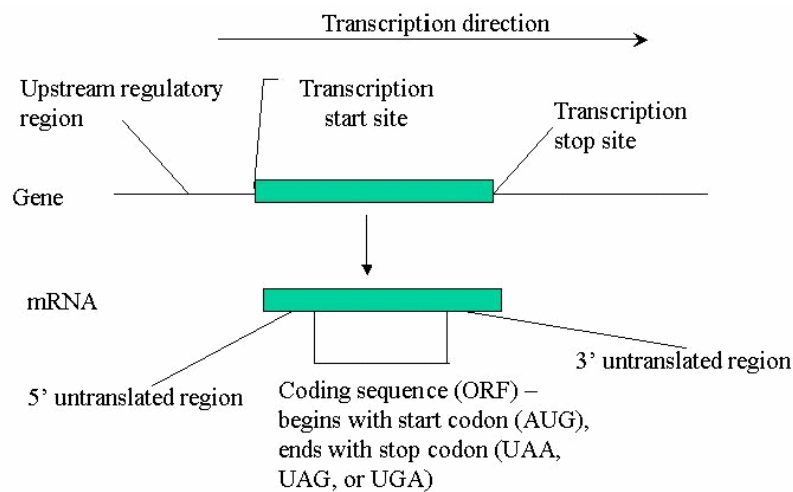
3' TGGCAATGA 5'

- So the reverse complement of ACCGTTACT is written TGGCAATGA
- In general people write one strand and in 5' to 3' order
 - This is the ordering that a polymerase or a ribosome scan the sequence
 - Establishes a common standard for genome nomenclatures

© Eric Xing @ CMU, 2005-2009

5

Gene structure in prokaryotes



© Eric Xing @ CMU, 2005-2009

6

Gene structure in prokaryotes

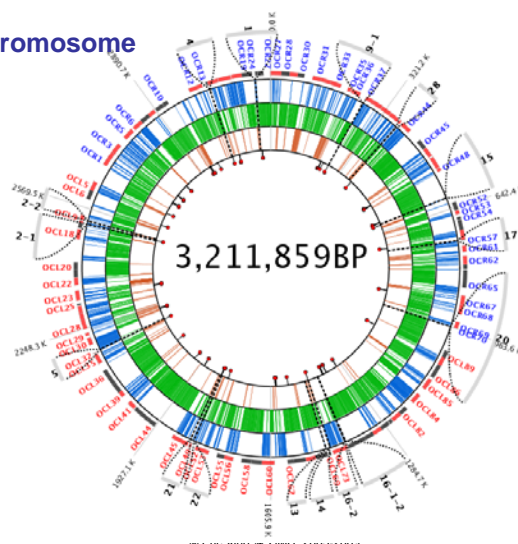
- A protein-coding gene consists of the following, in 5' to 3' order
 - An **upstream regulatory region**, generally < 50 bp, which turns transcription on and off.
 - A **transcription start site** where RNA polymerase incorporates 1st nucleotide into nascent mRNA.
 - A 5' **untranslated region**, generally < 30bp, that is transcribed into mRNA but not translated.
 - The **translation start site** marking the start of the coding region. Consists of a **start codon**, which causes the start of translation
 - The **coding region** of the gene (typically=1000bp), consisting of a sequence of codons.
 - The **translation stop site** marking the end of coding region. Consists of a **stop codon**, which causes the release of the polypeptide at conclusion of translation.
 - A 3' **untranslated region**, transcribed into RNA but not translated.
 - The **transcription stop site** marking where the RNA polymerase concludes transcription.

© Eric Xing @ CMU, 2005-2009

7

The bacterial genome

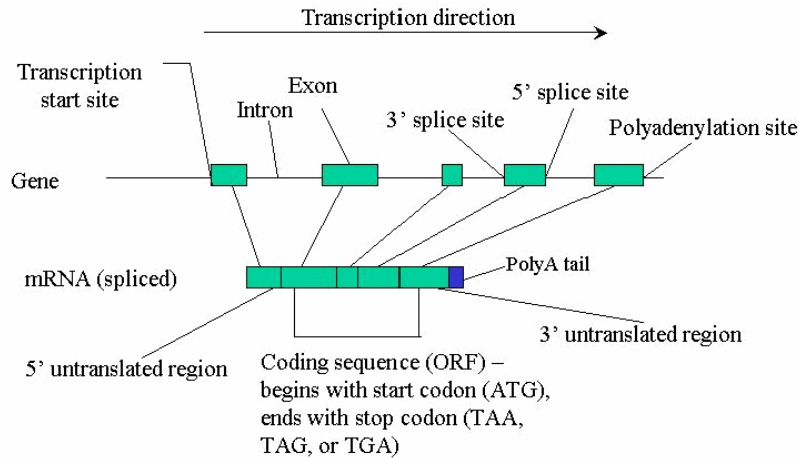
The E. coli chromosome



© Eric Xing @ CMU, 2005-2009

8

Gene structure in eukaryotes



© Eric Xing @ CMU, 2005-2009

9

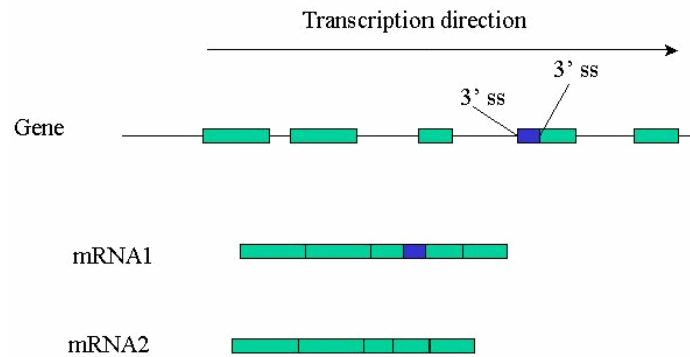
Gene structure in eukaryotes

- A typical gene consists of the following, in 5' to 3' order
 - An **upstream regulatory region**, often larger and more complex than in prokaryotes, parts of which may be several thousand bases or more upstream of transcription start site.
 - A **transcription start site**.
 - A **5' untranslated region**, often larger than in prokaryotes, and which may include sequences playing a role in translation regulation.
 - The **coding sequence**, which unlike the case with prokaryotes, may be interrupted by non-coding regions called introns. These are spliced out of the transcript to form the mature mRNA (and sometimes the splicing can occur in more than one way).
 - The **translation stop site**.
 - A **3' untranslated region**, which may contain sequences involved in translational regulation.
 - A **polyadenylation (polyA) signal**, which indicates to the cell's RNA processing machinery that the RNA transcript is to be cleaved and a poly-adenine sequence (AAAAAA...) tail appended to it.
 - The **transcription stop site**.

© Eric Xing @ CMU, 2005-2009

10

Alternative splicing

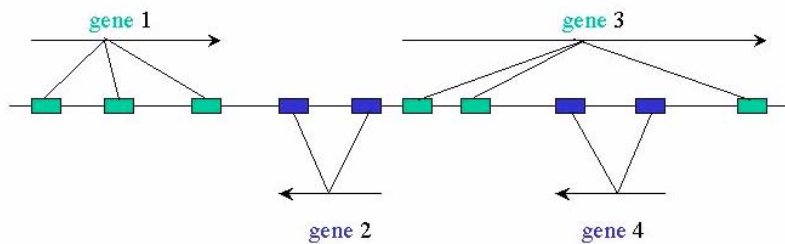


© Eric Xing @ CMU, 2005-2009

11

Eukaryotic genome structure

- Genes may be transcribed in either direction, and can overlap



© Eric Xing @ CMU, 2005-2009

12

Probabilities on Sequences



- Let \mathcal{S} be the space of DNA or protein sequences of a given length n . Here are some simple assumptions for assigning probabilities to sequences:
 - Equal frequency assumption:** All residues are equally probable at any position; i.e., $P(X_{i,r}) = P(X_{i,q})$ for any two residues r and q , for all i .
 - this implies that $P(X_{i,r}) = \theta_r = 1/|\mathcal{A}|$, where \mathcal{A} is the residue alphabet (1/20 for proteins, 1/4 for DNA)
 - Independence assumption:** whether or not a residue occurs at a position is independent of what residues are present at other positions.
 - probability of a sequence

$$P(X_1, X_2, \dots, X_N) = \theta_{r_1} \cdot \theta_{r_2} \cdot \dots \cdot \theta_{r_N} = \theta_r^N$$

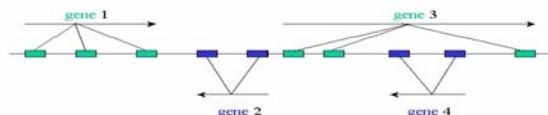
© Eric Xing @ CMU, 2005-2009

13

Failure of Equal Frequency Assumption for (real) DNA



- For most organisms, the nucleotides composition is significantly different from 0.25 for each nucleotide, e.g.,
 - H. influenza* .31 A, .19 C, .19 G, .31 T
 - P. aeruginosa* .17 A, .33 C, .33 G, .17 T
 - M. janaschii* .34 A, .16 C, .16 G, .34 T
 - S. cerevisiae* .31 A, .19 C, .19 G, .31 T
 - C. elegans* .32 A, .18 C, .18 G, .32 T
 - H. sapiens* .30 A, .20 C, .20 G, .30 T
- Note symmetry: $A \approx T$, $C \approx G$, even though we are counting nucleotides on just one strand. Explanation:
 - although individual biological features may have non-symmetric composition, usually features are distributed ~ randomly w.r.t. strand, so get symmetry.



© Eric Xing @ CMU, 2005-2009

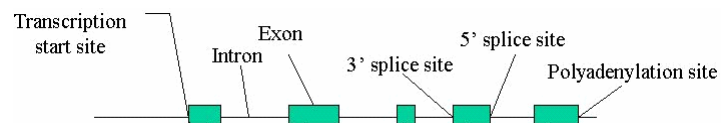
14

General Hypothesis Regarding Unequal Frequency



- Neutralist hypothesis: mutation bias (e.g., due to nucleotide pool composition)
- Selectionist hypothesis: natural selection bias

Probabilistic segmentation



Models for Homogeneous Sequence Entities



- Probabilities models for long "homogeneous" sequence entities, such as:
 - exons (ORFs)
 - introns
 - inter-genetic background
 - protein coiled-coil (other other structural) regions
- Assumptions:
 - no consensus, no recurring string patterns
 - have distinct but uniform residue-composition (i.e., same for all sites)
 - every site in the entity are iid samples from the same model
- The model:
 - a single multinomial: $X \sim \text{Mul}(1, \theta)$

© Eric Xing @ CMU, 2005-2009

17

The Multinomial Model for Sequence



- For a site i , define its residue identity to be a multinomial random vector:

$$X_i = \begin{bmatrix} X_{i,A} \\ X_{i,C} \\ X_{i,G} \\ X_{i,T} \end{bmatrix}, \quad \text{where} \quad \begin{array}{l} X_{i,j} \in [0,1], \text{ and } \sum_{j \in \{A,C,G,T\}} X_{i,j} = 1 \\ X_{i,j} = 1 \text{ w.p. } \theta_j, \quad \sum_{j \in \{A,C,G,T\}} \theta_j = 1 \end{array}$$

- The probability of an observation $s_i=A$ (i.e, $x_{i,A}=1$) at site i :

$$\begin{aligned} p(x_i = (\text{say}, A)) &= P(X_{i,j} = 1, \text{ where } j = A \text{ index the observed nucleotide}) \\ &= \theta_A = \theta_A^{x_A} \times \theta_C^{x_C} \times \theta_G^{x_G} \times \theta_T^{x_T} = \prod_k \theta_k^{x_k} = \theta^x \end{aligned}$$

- The probability of a sequence (x_1, x_2, \dots, x_N) :

$$\begin{aligned} P(x_1, x_2, \dots, x_N) &= \prod_{i=1}^N P(x_i) = \prod_{i=1}^N \prod_k \theta_k^{x_{i,k}} \\ &= \prod_k \theta_k^{\sum_{i=1}^N x_{i,k}} = \prod_k \theta_k^{n_k} \end{aligned}$$

© Eric Xing @ CMU, 2005-2009

18

Parameter Estimation



- Maximum likelihood estimation: $\theta = \arg \max_{\theta} P(D | \theta)$

- multinomial parameters:

$$\{\theta_1, \theta_2, \dots\} = \arg \max_{\theta} \prod_k \theta_k^{n_k}, \quad \text{s.t. } \sum_k \theta_k = 1$$

It can be shown that: $\theta_k^{\text{ML}} = n_k / N$

- Bayesian estimation:

- Dirichlet distribution:

$$P(\theta) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1} = C(\alpha) \prod_k \theta_k^{\alpha_k - 1}$$

- Posterior distribution of θ under the Dirichlet prior:

$$P(\theta | x_1, \dots, x_N) \propto \prod_k \theta_k^{\alpha_k - 1} \prod_k \theta_k^{n_k} = \prod_k \theta_k^{\alpha_k - 1 + n_k}$$

- Posterior mean estimation:

$$\theta_k = \int \theta_k P(\theta | D) d\theta = \int \theta_k \prod_k \theta_k^{\alpha_k - 1 + n_k} d\theta = \frac{n_k + \alpha_k}{N + |\alpha|}$$

© Eric Xing @ CMU, 2005-2009

19

Models for Homogeneous Sequence Entities, ctd



- Limitations

- non-uniform residue composition (e.g., CG rich regions)
- non-coding structural regions (MAR, centromere, telomere)
- di- or tri- nucleotide couplings
- estimation bias
- evolutionary constraints

© Eric Xing @ CMU, 2005-2009

20

Site Models



- Probabilities models for short sequences, such as:
 - splice sites
 - translation start sites
 - promoter elements
 - protein "motifs"
- Assumptions:
 - different examples of sites can be aligned without indels (insertions/deletions) such that tend to have similar residues in same positions
 - drop equal frequency assumption; instead have position-specific frequencies
 - retain independence assumption (for now)

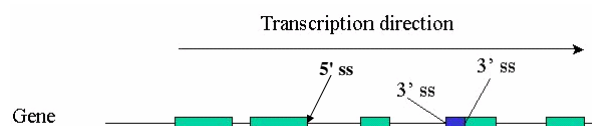
© Eric Xing @ CMU, 2005-2009

21

Site Models ctd.



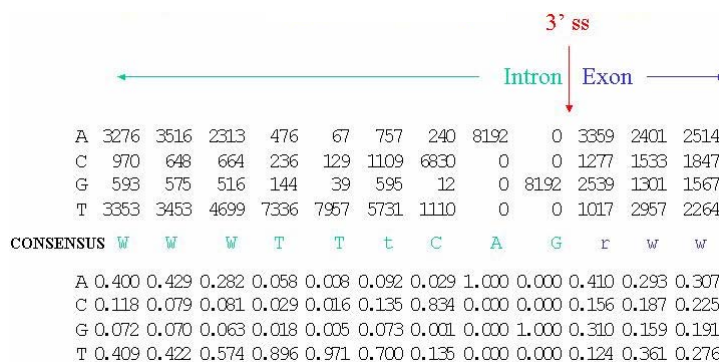
- Applies to short segments (<30 residues) where precise residue spacing is structurally or functionally important, and certain positions are highly conserved
 - DNA/RNA sequence binding sites for a single protein or RNA molecule
 - Protein internal regions structurally constrained due to folding requirements; or surface regions functionally constrained because bind certain ligands
- Example: *C. elegans* splice sites



© Eric Xing @ CMU, 2005-2009

22

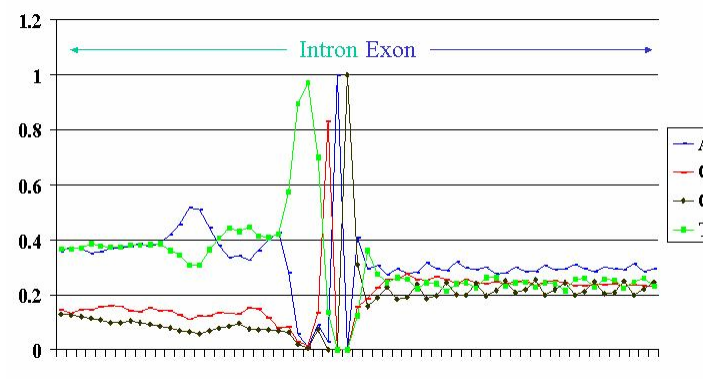
Nucleotide Counts for 8192 *C. elegans* 3' Splice Sites



© Eric Xing @ CMU, 2005-2009

23

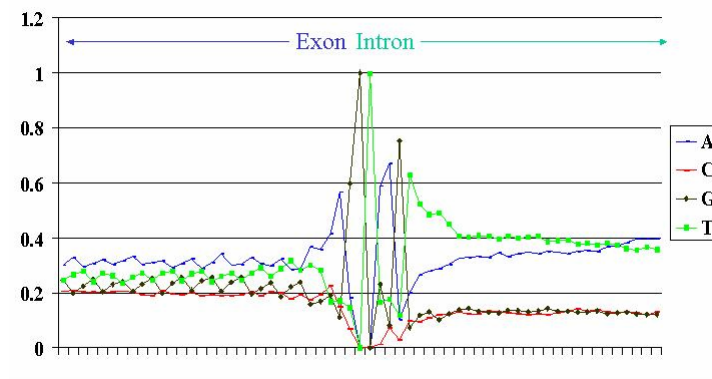
3' Splice Site - *C. elegans*



© Eric Xing @ CMU, 2005-2009

24

5' Splice Sites - *C. elegans*



© Eric Xing @ CMU, 2005-2009

25

Limitation of Homogeneous Site Models

- Failure to allow indels means variably spaced subelements are "smeared", e.g.:
 - branch site, for 3' splice sites;
 - coding sequences, for both 3' and 5' sites
- Independence assumption
 - usually OK for protein sequences (after correcting for evolutionary relatedness)
 - often fails for nucleotide sequences; examples:
 - 5' sites (Burge-Karlin observation);
 - background (dinucleotide correlation, e.g., GC repeats).

© Eric Xing @ CMU, 2005-2009

26

Why Correlation?



- Splicing involves pairing of a small RNA with the transcription at the 5' splice site.
- The RNA is complementary to the 5' srRNA consensus sequence.
- A mismatch at position -1 tends to destabilize the pairing, and makes it more important for other positions to be correctly paired.
- Analogy can be easily drew for other DNA and protein motifs.

Comparing Alternative Probability Models



- We will want to consider more than one model at a time, in the following situations:
 - To differentiate between two or more hypothesis about a sequence
 - To generate increasingly refined probability models that are progressively more accurate
- First situation arises in testing biological assertion, e.g., "is this a coding sequence?" Would compare two models:
 1. one associated with a hypothesis H_{coding} which attaches to a sequence the probability of observing it under experiment of drawing a random coding sequence from the genome
 2. one associate with a hypothesis $H_{noncoding}$ which attaches to a sequence the probability of observing it under experiment of drawing a random non-coding sequence from the genome.

Likelihood Ratio Test

- The posterior probability of a model given data is:

$$P(M|D) = P(D|M)P(M)/P(D)$$

- Given that all models are equally probable *a priori*, the posterior probability ratio of two models given the same data reduce to a *likelihood ratio*:

$$LR(M_a, M_b | D) = \frac{P(D | M_a)}{P(D | M_b)}$$

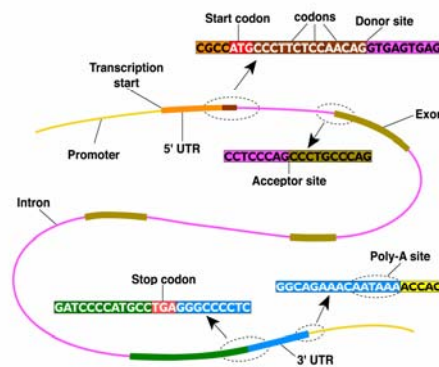
- the numerator and the denominator may both be very small!
- The log likelihood ratio (LLR) is the logarithm of the likelihood ratio:

$$LLR(M_a, M_b | D) = \log P(D | M_a) - \log P(D | M_b)$$

© Eric Xing @ CMU, 2005-2009

29

The Hidden Markov Models for sequence parsing



© E.

Gene Finding

- Given un-annotated sequences,
- delineate:
 - transcription initiation site,
 - exon-intron boundaries,
 - transcription termination site,
 - a variety of other motifs: promoters, polyA sites, branching sites, etc.
- The hidden Markov model (HMM)



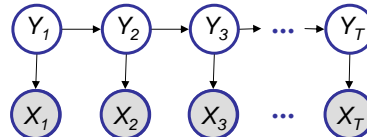
© Eric Xing @ CMU, 2005-2009

31

Hidden Markov Models

The underlying source:
genomic entities,
dice,

The sequence:
Ploy NT,
sequence of rolls,



© Eric Xing @ CMU, 2005-2009

32

Example: The Dishonest Casino



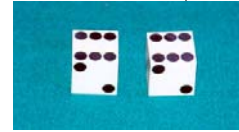
A casino has two dice:

- Fair die
 $P(1) = P(2) = P(3) = P(5) = P(6) = 1/6$
- Loaded die
 $P(1) = P(2) = P(3) = P(5) = 1/10$
 $P(6) = 1/2$

Casino player switches back-&-forth
between fair and loaded die once every
20 turns

Game:

1. You bet \$1
2. You roll (always with a fair die)
3. Casino player rolls (maybe with fair die,
maybe with loaded die)
4. Highest number wins \$2



© Eric Xing @ CMU, 2005-2009

33

Puzzles Regarding the Dishonest Casino



GIVEN: A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

QUESTION

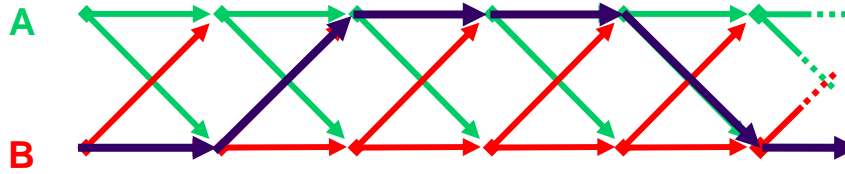
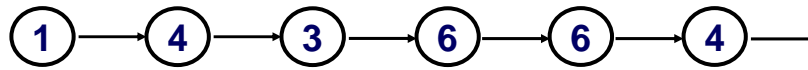
- How likely is this sequence, given our model of how the casino works?
 - This is the **EVALUATION** problem in HMMs
- What portion of the sequence was generated with the fair die, and what portion with the loaded die?
 - This is the **DECODING** question in HMMs
- How “loaded” is the loaded die? How “fair” is the fair die? How often does the casino player change from fair to loaded, and back?
 - This is the **LEARNING** question in HMMs

© Eric Xing @ CMU, 2005-2009

34

A Stochastic Generative Model

- Observed sequence:



- Hidden sequence (a parse or segmentation):



© Eric Xing @ CMU, 2005-2009

35

Definition (of HMM)

- Observation space

Alphabetic set: $C = \{c_1, c_2, \dots, c_K\}$
Euclidean space: \mathbb{R}^d

- Index set of hidden states

$I = \{1, 2, \dots, M\}$

- Transition probabilities between any two states

$$p(y_t^j = 1 | y_{t-1}^i = 1) = a_{i,j},$$

or $p(y_t | y_{t-1}^i = 1) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in I.$

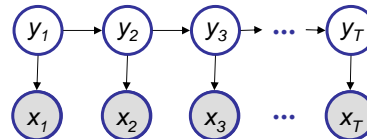
- Start probabilities

$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

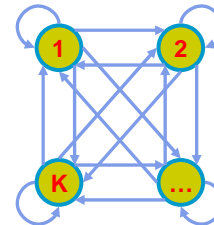
- Emission probabilities associated with each state

$p(x_t | y_t^i = 1) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in I.$
or in general:

$$p(x_t | y_t^i = 1) \sim f(\cdot | \theta_i), \forall i \in I.$$



Graphical model



State automata

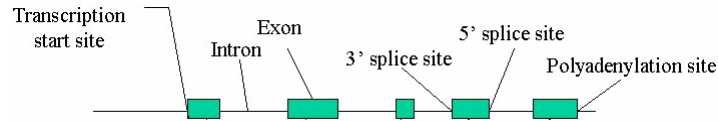
© Eric Xing @ CMU, 2005-2009

36

Probability of a Parse

- What is a parse?

1245526462146146136136661664661636616366163616515615115146123562344



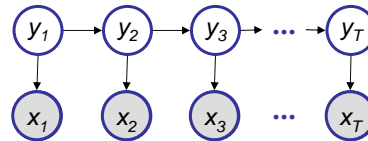
- How to score a parse?

© Eric Xing @ CMU, 2005-2009

37

Probability of a Parse

- Given a sequence $\mathbf{x} = x_1, \dots, x_T$ and a parse $\mathbf{y} = y_1, \dots, y_T$,
- To find how likely is the parse: (given our HMM and the sequence)



$$\begin{aligned}
 p(\mathbf{x}, \mathbf{y}) &= p(x_1, \dots, x_T, y_1, \dots, y_T) && \text{(Joint probability)} \\
 &= p(y_1) p(x_1 | y_1) p(y_2 | y_1) p(x_2 | y_2) \dots p(y_T | y_{T-1}) p(x_T | y_T) \\
 &= p(y_1) P(y_2 | y_1) \dots p(y_T | y_{T-1}) \times p(x_1 | y_1) p(x_2 | y_2) \dots p(x_T | y_T) \\
 &= p(y_1, \dots, y_T) p(x_1, \dots, x_T | y_1, \dots, y_T)
 \end{aligned}$$

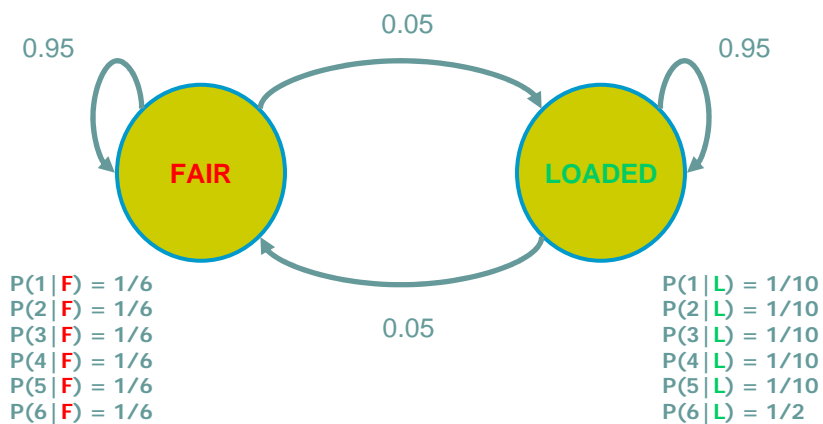
$$\begin{aligned}
 \text{Let } \pi_{y_1} &\stackrel{\text{def}}{=} \prod_{i=1}^M [\pi_i]^{y'_1}, \quad a_{y_t, y_{t+1}} \stackrel{\text{def}}{=} \prod_{i,j=1}^M [a_{ij}]^{y'_t y'_{t+1}}, \quad \text{and } b_{y_t, x_t} \stackrel{\text{def}}{=} \prod_{i=1}^M \prod_{k=1}^K [b_{ik}]^{y'_t x'_t}, \\
 &= \pi_{y_1} a_{y_1, y_2} \dots a_{y_{T-1}, y_T} b_{y_1, x_1} \dots b_{y_T, x_T}
 \end{aligned}$$

- Marginal probability: $p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) = \sum_{y_1} \sum_{y_2} \dots \sum_{y_N} \pi_{y_1} \prod_{t=2}^T a_{y_{t-1}, y_t} \prod_{t=1}^T p(x_t | y_t)$
- Posterior probability: $p(\mathbf{y} | \mathbf{x}) = p(\mathbf{x}, \mathbf{y}) / p(\mathbf{x})$

© Eric Xing @ CMU, 2005-2009

38

The Dishonest Casino Model



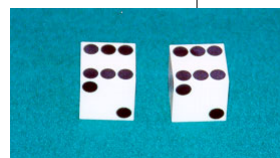
© Eric Xing @ CMU, 2005-2009

39

Example: the Dishonest Casino



- Let the sequence of rolls be:
 - $x = 1, 2, 1, 5, 6, 2, 1, 6, 2, 4$
- Then, what is the likelihood of
 - $y = \text{Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair?}$
(say initial probs $a_{0\text{Fair}} = 1/2, a_{0\text{Loaded}} = 1/2$)



$$\frac{1}{2} \times P(1 | \text{Fair}) P(\text{Fair} | \text{Fair}) P(2 | \text{Fair}) P(\text{Fair} | \text{Fair}) \dots P(4 | \text{Fair}) =$$

$$\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = .00000000521158647211 = 5.21 \times 10^{-9}$$

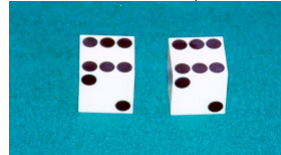
© Eric Xing @ CMU, 2005-2009

40

Example: the Dishonest Casino



- So, the likelihood the die is fair in all this run is just 5.21×10^{-9}



- OK, but what is the likelihood of
 - π = Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded?

$$\frac{1}{2} \times P(1 \mid \text{Loaded}) P(\text{Loaded} \mid \text{Loaded}) \dots P(4 \mid \text{Loaded}) =$$

$$\frac{1}{2} \times (1/10)^8 \times (1/2)^2 (0.95)^9 = .00000000078781176215 = 0.79 \times 10^{-9}$$

- Therefore, it is after all 6.59 times more likely that the die is fair all the way, than that it is loaded all the way

© Eric Xing @ CMU, 2005-2009

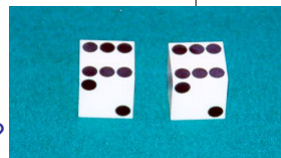
41

Example: the Dishonest Casino



- Let the sequence of rolls be:

- $x = 1, 6, 6, 5, 6, 2, 6, 6, 3, 6$



- Now, what is the likelihood $\pi = F, F, \dots, F$?

- $\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = 0.5 \times 10^{-9}$, same as before

- What is the likelihood $y = L, L, \dots, L$?

$$\frac{1}{2} \times (1/10)^4 \times (1/2)^6 (0.95)^9 = .00000049238235134735 = 5 \times 10^{-7}$$

- So, it is 100 times more likely the die is loaded

© Eric Xing @ CMU, 2005-2009

42

Applications of HMMs

- **Some early applications of HMMs**

- finance, but we never saw them
- speech recognition
- modelling ion channels

- **In the mid-late 1980s HMMs entered genetics and molecular biology, and they are now firmly entrenched.**

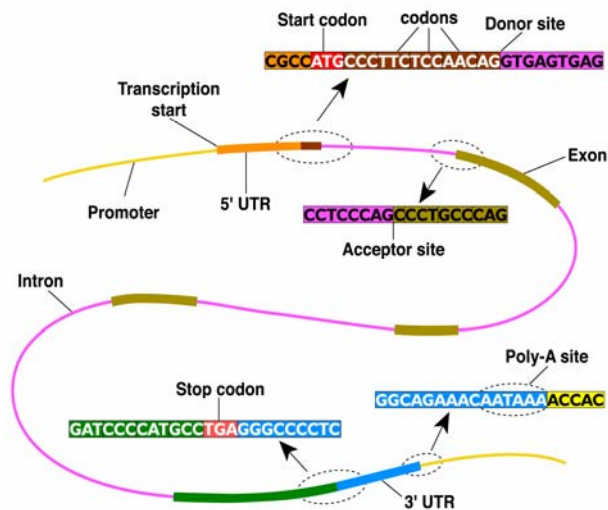
- **Some current applications of HMMs to biology**

- mapping chromosomes
- aligning biological sequences
- predicting sequence structure
- inferring evolutionary relationships
- finding genes in DNA sequence

© Eric Xing @ CMU, 2005-2009

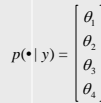
43

Typical structure of a gene



© Eric Xing @ CMU, 2005-2009

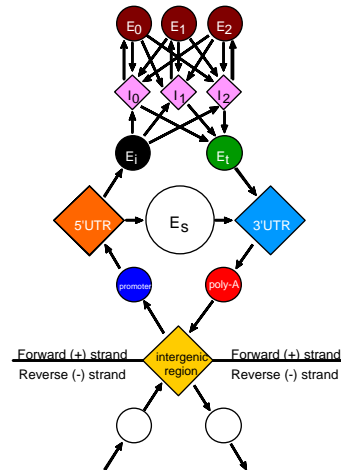
44

[illegible]

- 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

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.



© Eric Xing @ CMU, 2005-2009

47

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 B_k or gene, given the observed sequence? **Forward-Backward**
- **Learning:** Under what parameterization are the observed sequences most probable? **Baum-Welch (EM)**

© Eric Xing @ CMU, 2005-2009

48

Decoding

- GIVEN $\mathbf{x} = x_1, \dots, x_T$, we want to find $\mathbf{y} = y_1, \dots, y_T$, such that $P(\mathbf{y}|\mathbf{x})$ is maximized:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \operatorname{argmax}_{\pi} P(\mathbf{y}, \mathbf{x})$$

Viterbi decoding

- GIVEN $\mathbf{x} = x_1, \dots, x_T$, we want to find $\mathbf{y} = y_1, \dots, y_T$, such that $P(\mathbf{y}|\mathbf{x})$ is maximized:

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}) = \operatorname{argmax}_{\pi} P(\mathbf{y}, \mathbf{x})$$

- Let

$$V_t^k = \max_{\{y_1, \dots, y_{t-1}\}} P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t^k = 1)$$

= Probability of most likely sequence of states ending at state $y_t = k$

- The recursion:

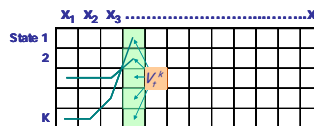
$$V_t^k = p(x_t | y_t^k = 1) \max_i a_{i,k} V_{t-1}^i$$

- Underflows are a significant problem

$$p(x_1, \dots, x_t, y_1, \dots, y_t) = \pi_{y_1} a_{y_1, y_2} \cdots a_{y_{t-1}, y_t} b_{y_1, x_1} \cdots b_{y_t, x_t}$$

- These numbers become extremely small – underflow

- Solution: Take the logs of all values: $V_t^k = \log p(x_t | y_t^k = 1) + \max_i (\log(a_{i,k}) + V_{t-1}^i)$



The Viterbi Algorithm – derivation



- Define the viterbi probability:

$$\begin{aligned}
 V_{t+1}^k &= \max_{\{y_1, \dots, y_t\}} P(x_1, \dots, x_t, y_1, \dots, y_t, x_{t+1}, y_{t+1}^k = 1) \\
 &= \max_{\{y_1, \dots, y_t\}} P(x_{t+1}, y_{t+1}^k = 1 | x_1, \dots, x_t, y_1, \dots, y_t) P(x_1, \dots, x_t, y_1, \dots, y_t) \\
 &= \max_{\{y_1, \dots, y_t\}} P(x_{t+1}, y_{t+1}^k = 1 | y_t) P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t) \\
 &= \max_i P(x_{t+1}, y_{t+1}^k = 1 | y_t^i = 1) \max_{\{y_1, \dots, y_{t-1}\}} P(x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t, y_t^i = 1) \\
 &= \max_i P(x_{t+1}, y_{t+1}^k = 1) a_{i,k} V_t^i \\
 &= P(x_{t+1}, y_{t+1}^k = 1) \max_i a_{i,k} V_t^i
 \end{aligned}$$

© Eric Xing @ CMU, 2005-2009

51

The Viterbi Algorithm



- Input: $\mathbf{x} = x_1, \dots, x_T$

Initialization:

$$V_1^k = P(x_1 | y_1^k = 1) \pi_k$$

Iteration:

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

$$\text{Ptr}(k, t) = \arg \max_i a_{i,k} V_{t-1}^i$$

Termination:

$$P(\mathbf{x}, \mathbf{y}^*) = \max_k V_T^k$$

TraceBack:

$$y_T^* = \arg \max_k V_T^k$$

$$y_{t-1}^* = \text{Ptr}(y_t^*, t)$$

© Eric Xing @ CMU, 2005-2009

52

Time complexity of Viterbi

