Hidden Markov Models

part 1

CBB 261 / CPS 261

B. Majoros

What is an HMM?

The short answer: it's an <u>opaque</u>, <u>nondeterministic</u> machine that <u>emits variable-length</u> sequences of <u>discrete¹</u> <u>symbols</u>.



What we mean by "Hidden" in "Hidden Markov Model": We can take the machine apart and look inside to see how many states there are and how they are connected, but we can't look inside the machine while it's running. (However, we can make inferences about what probably happened inside the machine, based on its output.)

¹...usually, though not always...



What is an HMM?

An HMM is a *stochastic machine* $M = (Q, \alpha, P_t, P_e)$ consisting of the following:

- a finite set of states, $Q = \{q_0, q_1, \dots, q_m\}$
- a finite <u>alphabet</u> $\alpha = \{s_0, s_1, \dots, s_n\}$
- a <u>transition</u> distribution $P_t: Q \times Q \mapsto \mathbb{R}$
- an <u>emission</u> distribution $P_e: Q \times \alpha \mapsto \mathbb{R}$





Probability of a Sequence

S=YRYRY $\phi = (q0, q1, q2, q1, q2, q1, q0)$



$P(\text{YRYRY}|M_1) =$

 $a_{0 \rightarrow 1} \times b_{1,Y} \times a_{1 \rightarrow 2} \times b_{2,R} \times a_{2 \rightarrow 1} \times b_{1,Y} \times a_{1 \rightarrow 2} \times b_{2,R} \times a_{2 \rightarrow 1} \times b_{1,Y} \times a_{1 \rightarrow 0}$ =1 × 1 × 0.15 × 1 × 0.3 × 1 × 0.15 × 1 × 0.3 × 1 × 0.05 =0.00010125 Duke

Another Example

 $M_{2} = (Q, \alpha, P_{t}, P_{e})$ $Q = \{q_{0}, q_{1}, q_{2}, q_{3}, q_{4}\}$ $\alpha = \{A, C, G, T\}$



Finding the Most Probable Path



Decoding with an HMM

$$\phi_{\max} = \operatorname{argmax}_{\phi} P(\phi|S) = \operatorname{argmax}_{\phi} \frac{P(\phi,S)}{P(S)}$$

$$= \operatorname{argmax}_{\phi} P(\phi,S)$$

$$= \operatorname{argmax}_{\phi} P(S|\phi) P(\phi)$$

$$F(S|\phi) = \prod_{i=0}^{L-1} P_e(x_i | y_{i+1})$$

$$F(\phi) = \prod_{i=0}^{L} P_i(y_{i+1} | y_i)$$

$$F(\phi) = \prod_{i=0}^{L} P_i(y_{i+1} | y_i)$$

$$F(\phi) = \operatorname{argmax}_{\phi} P(\phi) P(\phi)$$

$$F(y_{i+1} | y_i)$$

$$F(\phi) = \operatorname{argmax}_{\phi} P(\phi) P(\phi)$$

$$F(y_{i+1} | y_i)$$

$$F(\phi) = \operatorname{argmax}_{\phi} P(\phi) P(\phi)$$

The Viterbi Algorithm

$$V(i,k) = \begin{cases} \max_{j} V(j,k-1)P_t(q_i | q_j)P_e(x_k | q_i) & \text{if } k > 0, \\ P_t(q_i | q_0)P_e(x_0 | q_i) & \text{if } k = 0. \end{cases}$$



In the final column:

$$P(\phi_{\max}) = \max_{i} V(i, L-1) P_t(q_0 | q_i)$$



Viterbi: Traceback

$$V(i,k) = \begin{cases} \max_{j} V(j,k-1)P_{t}(q_{i}|q_{j})P_{e}(x_{k}|q_{i}) & \text{if } k > 0, \\ P_{t}(q_{i}|q_{0})P_{e}(x_{0}|q_{i}) & \text{if } k = 0. \end{cases}$$
$$T(i,k) = \begin{cases} \arg_{j} XV(j,k-1)P_{t}(q_{i}|q_{j})P_{e}(x_{k}|q_{i}) & \text{if } k > 0, \\ 0 & \text{if } k = 0. \end{cases}$$

T(T(T(...,T(T(i,L-1),L-2)...,2),1),0) = 0



(Log) Viterbi Algorithm in Pseudocode





The Forward Algorithm : Probability of a Sequence

$$F(i,k) = \begin{cases} 1 & \text{for } k = 0, i = 0 \\ 0 & \text{for } k > 0, i = 0 \\ 0 & \text{for } k = 0, i > 0 \end{cases}$$

$$F(i,k) = \begin{cases} 1 & \text{for } k = 0, i = 0 \\ 0 & \text{for } k = 0, i > 0 \\ 0 & \text{for } k = 0, i > 0 \end{cases}$$

$$F(i,k) = \begin{cases} 2^{|-1|} \sum_{j=0}^{|-1|} F(j,k-1)P_t(q_i \mid q_j)P_e(x_{k-1} \mid q_i) & \text{for } 1 \le k \le |S|, \\ 1 \le i < |Q| \end{cases}$$

F(i,k) represents the probability $P(S_{0..k-1}|q_i)$ that the machine emits the subsequence $x_0...x_{k-1}$ by any path ending in state q_i —i.e., so that symbol x_{k-1} is emitted by state q_i .

$$P(S \mid M) = \sum_{i=0}^{|Q|-1} F(i, |S|) P_t(q_0 \mid q_i)$$



The Forward Algorithm in Pseudocode

```
procedure unscaledForwardAlg(Q, P<sub>t</sub>, P<sub>e</sub>, S, \lambda_{trans}, \lambda_{emit})
          \forall_{i \in [0, |Q|-1]} \forall_{k \in [0, |S|]} F[i][k] \leftarrow 0;
1.
    F[0][0]←1;
2.
3.
    for k \leftarrow 1 up to |S| do
4.
    s←S[k-1];
              foreach q_i \in \lambda_{\text{emit}}[s] do
                                                                                               fill out the
5.
                                                                                               DP matrix
6.
                  sum \leftarrow 0;
7.
                  foreach q_i \in \lambda_{trans}[i] do
                     sum \leftarrow sum + P_t (q_i | q_j) * F[j][k-1];
8.
9.
                 F[i][k] \leftarrow sum^* P_e(s|q_i);
10.
          return F;
procedure unscaledForwardProb (Q, P<sub>t</sub>, P<sub>e</sub>, S, \lambda_{trans}, \lambda_{emit})
          F \leftarrow unscaled Forward Alg (Q, P<sub>t</sub>, P<sub>e</sub>, S, \lambda_{trans}, \lambda_{emit});
1.
2. sum←0;
                                                                                               sum over
3.
    len←|S|;
                                                                                               the final
    s←S[len-1]
4.
                                                                                               column to
5.
          foreach q_i \in \lambda_{emit}[s] do
                                                                                               get P(S)
6.
              sum \leftarrow sum + F[i][len] * P_t(q_0|q_i);
7.
           return sum;
                                                                                                        Duke
```

Training an HMM from Labeled Sequences

CGATATTCGATTCTACGCGCGTATACTAGCTTATCTGATC 01111112222222111112222111111222211110

| S | | | | to state | | | | | | | | | |
|--------|-------|---|--------|----------|----------|--|--|--|--|--|--|--|--|
| sition | | | 0 | 1 | 2 | | | | | | | | |
| rar | from | 0 | 0 (0%) | 1 (100%) | 0 (0%) | | | | | | | | |
| 1 | state | 1 | 1 (4%) | 21 (84%) | 3 (12%) | | | | | | | | |
| | | 2 | 0 (0%) | 3 (20%) | 12 (80%) | | | | | | | | |

| | | | symbol | | | | | | | | | |
|-------------|---|------------|------------|------------|------------|--|--|--|--|--|--|--|
| SUC | | Α | С | G | Т | | | | | | | |
| in state | 1 | 6 (24%) | 7 (28%) | 5 (20%) | 7 (28%) | | | | | | | |
| 0 | 2 | 3 (20%) | 3 (20%) | 2 (13%) | 7 (47%) | | | | | | | |

 $a_{i,j} = \frac{A_{i,j}}{\sum_{h=0}^{|Q|-1} A_{i,h}}$

 $e_{i,k} = \frac{E_{i,k}}{\sum_{h=0}^{|\Sigma|-1} E_{i,k}}$



Recall: Eukaryotic Gene Structure



Using an HMM for Gene Prediction



the input sequence: AGCTAGCAGTATGT the most probable path: ______ the gene prediction: _____







Recall: Sensitivity and Specificity

$$Sn = \frac{TP}{TP + FN}$$

$$Sp = \frac{TP}{TP + FP}$$

$$F = \frac{2 \times Sn \times Sp}{Sn + Sp}$$



HOMER, version H_3



| | nuc | cleoti | des | splice sites | | start cod | /stop lons | е | genes | | | |
|----------------|-----|--------|-----|-----------------|----|--------------|---------------|----|-------|---|----|---|
| | Sn | Sp | F | Sn | Sp | Sn | Sp | Sn | Sp | F | Sn | # |
| baseline | 100 | 28 | 44 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H ₃ | 53 | 88 | 66 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

HOMER, version H_5



| | nuc | cleoti | des | splice sites | | start/stop codons | | exons | | | genes | | |
|-----------------------|-----|--------|-----|-----------------|----|----------------------|----|-------|----|---|-------|---|------|
| | Sn | Sp | F | Sn | Sp | Sn | Sp | Sn | Sp | F | Sn | # | |
| H ₃ | 53 | 88 | 66 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| H ₅ | 65 | 91 | 76 | 1 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | Juke |



| | nucleotides | | | SITES | | coaons | | exons | | | genes | | |
|------------------------|-------------|----|----|-------|--------------------|--------|----|-------|----|----|-------|-------|--------|
| | Sn | Sp | F | Sn | Sp | Sn | Sp | Sn | Sp | F | Sn | # | |
| H ₅ | 65 | 91 | 76 | 1 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 | |
| H ₁₇ | 81 | 93 | 87 | 34 | 48 | 43 | 37 | 19 | 24 | 21 | 7 | 35 | ke |
| | | | | | ALC: NOT THE OWNER | | | | | | | UNIVI | ERSITY |

Maintaining Phase Across an Intron





Recall: Weight Matrices







Summary of HOMER Results



Figure 6.10: Nucleotide F-score (y-axis) on a test set of 500 A. thaliana genes, as a function of number of states (x-axis) in an HMM for a simple gene finder.

| | nue | cleoti | des | splice sites | | start coa | t/stop lons | e | exons | genes | | |
|-----------------|-----|--------|-----|-----------------|----|--------------|----------------|----|-------|-------|----|----|
| | Sn | Sp | F | Sn | Sp | Sn | Sp | Sn | Sp | F | Sn | # |
| baseline | 100 | 28 | 44 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H ₃ | 53 | 88 | 66 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| H ₅ | 65 | 91 | 76 | 1 | 3 | 3 | 3 | 0 | 0 | 0 | 0 | 0 |
| H ₁₇ | 81 | 93 | 87 | 34 | 48 | 43 | 37 | 19 | 24 | 21 | 7 | 35 |
| H ₂₇ | 83 | 93 | 88 | 40 | 49 | 41 | 36 | 23 | 27 | 25 | 8 | 38 |
| H ₇₇ | 88 | 96 | 92 | 66 | 67 | 51 | 46 | 47 | 46 | 46 | 13 | 65 |
| H ₉₅ | 92 | 97 | 94 | 79 | 76 | 57 | 53 | 62 | 59 | 60 | 19 | 93 |



Higher-order Markov Models

| | P(G) |
|------------------------|--------|
| 0 th order: | ACGCTA |

 $1^{st} \text{ order:} \qquad ACGCTA$

2nd order:

P(G|AC) ACGCTA

$$P_{e}(g_{n} | g_{0}...g_{n-1},q_{j}) \approx \frac{C(g_{0}...g_{n},q_{j})}{\sum_{s \in \alpha} C(g_{0}...g_{n-1}s,q_{j})}$$



Higher-order Markov Models

| | order | order nucleotides | | splice sites | | starts/ stops | | exons | | | genes | | |
|--------------|-------|--------------------------|----|-----------------|----|------------------|----|-------|----|----|-------|----|-----|
| | | Sn | Sp | F | Sn | Sp | Sn | Sp | Sn | Sp | F | Sn | # |
| H_{95}^0 | 0 | 92 | 97 | 94 | 79 | 76 | 57 | 53 | 62 | 59 | 60 | 19 | 93 |
| H_{95}^{l} | 1 | 95 | 98 | 97 | 87 | 81 | 64 | 61 | 72 | 68 | 70 | 25 | 127 |
| H_{95}^2 | 2 | 98 | 98 | 98 | 91 | 82 | 65 | 62 | 76 | 69 | 72 | 27 | 136 |
| H_{95}^3 | 3 | 98 | 98 | 98 | 91 | 82 | 67 | 63 | 76 | 69 | 72 | 28 | 140 |
| H_{95}^{4} | 4 | 98 | 97 | 98 | 90 | 81 | 69 | 64 | 76 | 68 | 72 | 29 | 143 |
| H_{95}^{5} | 5 | 98 | 97 | 98 | 90 | 81 | 66 | 62 | 74 | 67 | 70 | 27 | 137 |



Variable-Order Markov Models

$$P_e^{IMM} (g_n | g_0 \dots g_{n-1}) = \begin{cases} \lambda_n^G P_e(g_n | g_0 \dots g_{n-1}) + (1 - \lambda_n^G) P_e^{IMM}(g_n | g_1 \dots g_{n-1}) & \text{if } n > 0 \\ P_e(g_n) & \text{if } n = 0 \end{cases}$$

$$\lambda_n^G = \begin{cases} 1 & \text{if } m \ge 400 \\ 0 & \text{if } m < 400 \text{ and } c < 0.5 \\ \frac{c}{400} \sum_{x \in \alpha} C(g_0 \dots g_{n-1}x) & \text{otherwise} \end{cases}$$

Interpolation Results



Figure 7.4: Relative accuracy of Markov chains (MC), IMC's, three-periodic MC's (3PMC), and three-periodic IMC's (3PIMC) on a particular task involving classification of human DNA sequences as coding versus noncoding.





• An HMM is a *stochastic generative model* which emits sequences

 Parsing with an HMM can be accomplished using a *decoding* algorithm (such as *Viterbi*) to find the most probable state-path generating the input sequence

 When state-labeled sequences are available, training of HMMs can be accomplished using *labeled sequence training*

• Otherwise, training of HMMs can be accomplished using *Expectation-Maximization (EM) (next lesson...)*

•*Posterior decoding* can be used to estimate the probability that a given symbol or substring was generate by a particular state (*next lesson...*)

