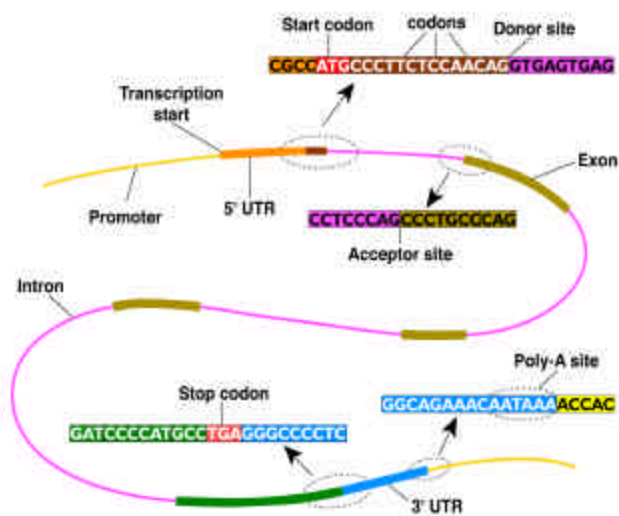


The Hidden Markov Models for sequence parsing



Gene structure in eukaryotes



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



Question # 1 – Evaluation



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

Question # 2 – Decoding



GIVEN

A sequence of rolls by the casino player

1245526462146146136136661664661636616366163616515615115146123562344

QUESTION

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

Question # 3 – Learning



GIVEN

A sequence of rolls by the casino player

124552646214614613613666166466163661636616361651561511514612356234

QUESTION

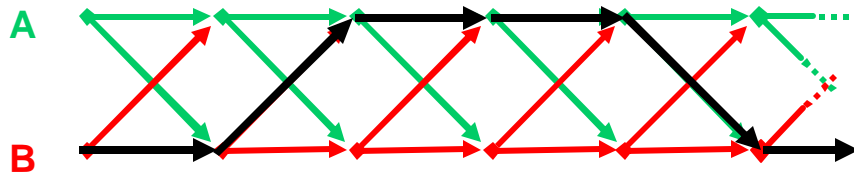
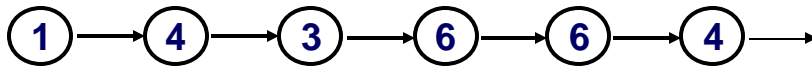
How “loaded” is the loaded die? How “fair” is the fair die? How often does the casino player change from fair to loaded, and back?

This is the **LEARNING** question in HMMs

A stochastic generative model



Observed sequence:



Hidden sequence (a parse or segmentation):



Definition (of HMM)



Definition: A hidden Markov model (HMM)

- **Observation alphabet** $\Sigma = \{ b_1, b_2, \dots, b_M \}$
- **Set of hidden states** $Q = \{ 1, \dots, K \}$
- **Transition probabilities** between any two states

$$a_{ij} = P(y_t=j|y_{t-1}=i)$$

$$a_{i,1} + \dots + a_{i,K} = 1, \text{ for all states } i = 1 \dots K$$

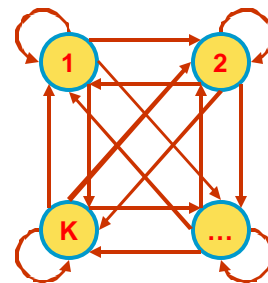
- **Start probabilities** $a_{0,i}$

$$a_{0,1} + \dots + a_{0,K} = 1$$

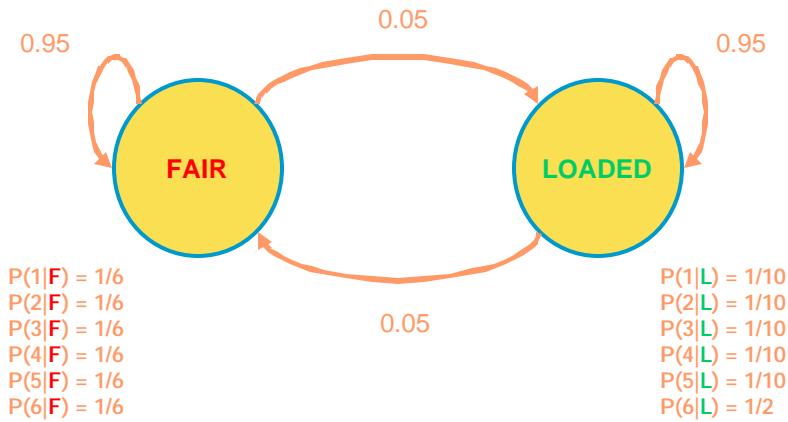
- **Emission probabilities** associated with each state

$$e_{ib} = P(x_t = b | y_t = i)$$

$$e_{i,b_1} + \dots + e_{i,b_M} = 1, \text{ for all states } i = 1 \dots K$$



The dishonest casino model

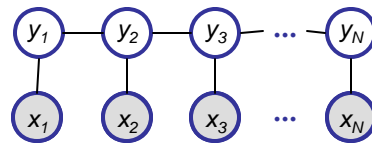


Likelihood of a parse



Given a sequence $x = x_1, \dots, x_N$
and a parse $y = y_1, \dots, y_N$,

To find how likely is the parse:
(given our HMM and the sequence)



$$\begin{aligned}
 P(x, y) &= P(x_1, \dots, x_N, y_1, \dots, y_N) && \text{(Joint probability)} \\
 &= P(y_1) P(x_1 | y_1) P(y_2 | y_1) P(x_2 | y_2) \dots P(y_N | y_{N-1}) P(x_N | y_N) \\
 &= P(y_1) P(y_2 | y_1) \dots P(y_N | y_{N-1}) \times P(x_1 | y_1) P(x_2 | y_2) \dots P(x_N | y_N) \\
 &= P(y_1, \dots, y_N) P(x_1, \dots, x_N | y_1, \dots, y_N) \\
 &= a_{0,y_1} a_{y_1,y_2} \dots a_{y_{N-1},y_N} e_{y_1,x_1} \dots e_{y_N,x_N}
 \end{aligned}$$

$$P(x) = \sum_y P(x, y) \quad \text{(Marginal probability)}$$

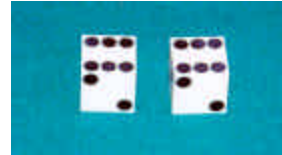
$$P(y|x) = P(x, y)/P(x) \quad \text{(Posterior probability)}$$

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 = Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair, Fair?

(say initial probs $a_{0\text{Fair}} = \frac{1}{2}$, $a_{0\text{Loaded}} = \frac{1}{2}$)

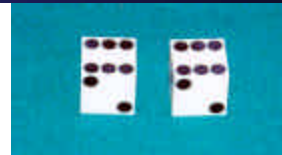
$$\frac{1}{2} \times P(1 \mid \text{Fair}) P(\text{Fair} \mid \text{Fair}) P(2 \mid \text{Fair}) P(\text{Fair} \mid \text{Fair}) \dots P(4 \mid \text{Fair}) =$$

$$\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = .00000000521158647211 = 0.5 \times 10^{-9}$$

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

p = Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded, Loaded?

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

$$\frac{1}{2} \times (1/10)^8 \times (1/2)^2 (0.95)^9 = .0000000078781176215 = 0.79 \times 10^{-9}$$

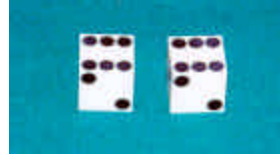
Therefore, it is after all 6.59 times more likely that the die is fair all the way, than that it is loaded all the way

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

$y = L, L, \dots, L$?

$\frac{1}{2} \times (1/10)^4 \times (1/2)^6 (0.95)^9 = .00000049238235134735 = 0.5 \times 10^{-7}$

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

The three main questions on HMMs



1. Evaluation

GIVEN an HMM M , and a sequence x ,
FIND Prob (x | M)

2. Decoding

GIVEN an HMM M , and a sequence x ,
FIND the sequence π of states that maximizes, e.g., $P(x | \pi, M)$

3. Learning

GIVEN an HMM M , with unspecified transition/emission probs.,
and a sequence x ,

FIND parameters $\theta = (e_{ik}, a_{ij})$ that maximize $P(x | \theta)$

Applications of HMMs



Some early applications of HMMs

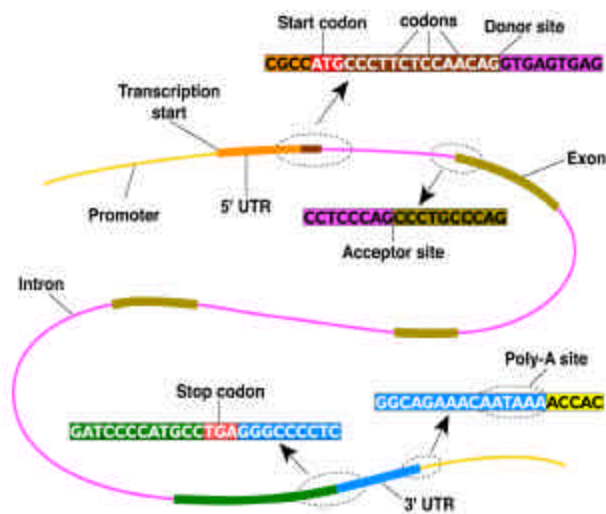
- finance, but we never saw them
- speech recognition
- modelling ion channels

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

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

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.

The Forward Algorithm



We want to calculate

$P(x)$ = probability of x , given the HMM

Sum over all possible ways of generating x :

$$P(x) = \sum_{\pi} P(x, \pi) = \sum_{\pi} P(x | \pi) P(\pi)$$

To avoid summing over an exponential number of paths π ,
define

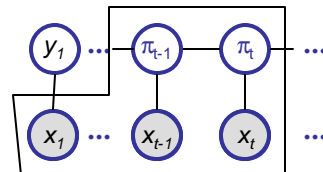
$$f_k(t) = P(x_1 \dots x_t, \pi_t = k) \quad (\text{the forward probability})$$

The Forward Algorithm – derivation



Compute the forward probability:

$$f_k(t) = P(x_1 \dots x_{t-1}, x_t, \pi_t = k)$$



$$= \sum_{\pi_{t-1}} P(x_1 \dots x_{t-1}, x_t, \pi_{t-1}, \pi_t = k)$$

$$= \sum_{\pi_{t-1}} P(x_1 \dots x_{t-1}, \pi_{t-1}) P(\pi_t = k | \pi_{t-1}, x_1 \dots x_{t-1}) P(x_t | \pi_t = k, \pi_{t-1}, x_1 \dots x_{t-1})$$

$$= \sum_{\pi_{t-1}} P(x_1 \dots x_{t-1}, \pi_{t-1}) P(\pi_t = k | \pi_{t-1}) P(x_t | \pi_t = k)$$

$$= \sum_i P(x_1 \dots x_{t-1}, \pi_{t-1} = i) P(\pi_t = k | \pi_{t-1} = i) P(x_t | \pi_t = k)$$

$$= e_k(x_t) \sum_i f_i(t-1) a_{ik}$$

The Forward Algorithm



We can compute $f_k(t)$ for all k, t , using dynamic programming!

Initialization:

$$f_0(0) = 1$$
$$f_k(0) = 0, \text{ for all } k > 0$$

Iteration:

$$f_i(t) = e_i(x_t) \sum_k f_k(t-1) a_{ik} \quad (a_{0k} \text{ is a vector of initial probability})$$

Termination:

$$P(x) = \sum_k f_k(T)$$

The Backward Algorithm



We want to compute $P(\pi_t = k | x)$,

the probability distribution on the t^{th} position, given x

We start by computing

$$P(\pi_t = k, x) = P(x_1 \dots x_t, \pi_t = k, x_{t+1} \dots x_N)$$
$$= P(x_1 \dots x_t, \pi_t = k) P(x_{t+1} \dots x_N | x_1 \dots x_t, \pi_t = k)$$
$$= P(x_1 \dots x_t, \pi_t = k) P(x_{t+1} \dots x_N | \pi_t = k)$$



Forward, $f_k(i)$

Backward, $b_k(i)$

The Backward Algorithm – derivation



Define the backward probability:

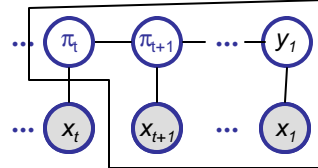
$$b_k(t) = P(x_{t+1} \dots x_N | \pi_t = k)$$

$$= \sum_{\pi_{t+1}} P(x_{t+1}, x_{t+2}, \dots, x_N, \pi_{t+1} | \pi_t = k)$$

$$= \sum_i P(\pi_{t+1} = i | \pi_t = k) P(x_{t+1} | \pi_t = k) P(x_{t+2}, x_{t+2}, \dots, x_N | \pi_{t+1} = i)$$

$$= \sum_i a_{k,i} e_i(x_{t+1}) P(x_{t+2}, \dots, x_N | \pi_{t+1} = i)$$

$$= \sum_i a_{k,i} e_i(x_{t+1}) b_i(t+1)$$



The Backward Algorithm



We can compute $b_k(t)$ for all k, t , using dynamic programming

Initialization:

$$b_k(T) = 1, \text{ for all } k$$

Iteration:

$$b_k(t) = \sum_i e_i(x_{t+1}) a_{k,i} b_i(t+1)$$

Termination:

$$P(x) = \sum_i a_{0,i} e_i(x_1) b_i(1)$$

Posterior Decoding



We can now calculate

$$P(\pi_t = k | x) = \frac{f_k(t) b_k(t)}{P(x)}$$

Then, we can ask

What is the most likely state at position t of sequence x :

Define π^* by Posterior Decoding:

$$\pi_t^* = \operatorname{argmax}_k P(\pi_t = k | x)$$

Decoding



GIVEN $x = x_1 x_2 \dots x_T$

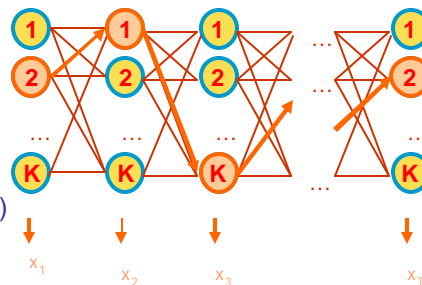
We want to find $\pi = \pi_1, \dots, \pi_T$,
such that $P(\pi|x)$ is maximized

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

We can use dynamic programming!

$$\text{Let } V_k(t) = \max_{\{\pi_1, \dots, \pi_{t-1}\}} P(x_1 \dots x_{t-1}, \pi_1, \dots, \pi_{t-1}, x_t, \pi_t = k)$$

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



Decoding – main idea



Given that for each possible state k ,
and for a fixed position t ,

$$V_k(t) = \max_{\{\pi_1, \dots, \pi_{t-1}\}} P(x_1 \dots x_{t-1}, \pi_1, \dots, \pi_{t-1}, x_t, \pi_t = k)$$

What is $V_i(t+1)$?

From definition,

$$\begin{aligned} V_i(t+1) &= \max_{\{\pi_1, \dots, \pi_t\}} P(x_1 \dots x_t, \pi_1, \dots, \pi_t, x_{t+1}, \pi_{t+1} = i) \\ &= \max_{\{\pi_1, \dots, \pi_t\}} P(x_{t+1}, \pi_{t+1} = i \mid x_1 \dots x_t, \pi_1, \dots, \pi_t) P(x_1 \dots x_t, \pi_1, \dots, \pi_t) \\ &= \max_{\{\pi_1, \dots, \pi_t\}} P(x_{t+1}, \pi_{t+1} = i \mid \pi_t) P(x_1 \dots x_{t-1}, \pi_1, \dots, \pi_{t-1}, x_t, \pi_t) \\ &= \max_k P(x_{t+1}, \pi_{t+1} = i \mid \pi_t = k) \max_{\{\pi_1, \dots, \pi_{t-1}\}} P(x_1 \dots x_{t-1}, \pi_1, \dots, \pi_{t-1}, x_t, \pi_t = k) \\ &= e_i(x_{t+1}) \max_k a_{ki} V_k(t) \end{aligned}$$

The Viterbi Algorithm



Input: $x = x_1 \dots x_N$

Initialization:

$$\begin{aligned} V_0(0) &= 1 && (0 \text{ is the imaginary first position}) \\ V_k(0) &= 0, \text{ for all } k > 0 \end{aligned}$$

Iteration:

$$\begin{aligned} V_j(t) &= e_j(x_t) \times \max_k a_{kj} V_k(t-1) \\ \text{Ptr}_j(t) &= \text{argmax}_k a_{kj} V_k(t-1) \end{aligned}$$

Termination:

$$P(x, \pi^*) = \max_k V_k(T)$$

Traceback:

$$\begin{aligned} \pi_T^* &= \text{argmax}_k V_k(T) \\ \pi_{t-1}^* &= \text{Ptr}_{\pi_t^*}(t) \end{aligned}$$

Viterbi Algorithm – a practical detail



Underflows are a significant problem

$$P(x_1, \dots, x_t, \pi_1, \dots, \pi_t) = a_{0\pi_1} a_{\pi_1\pi_2} \dots a_{\pi_{t-1}\pi_t} e_{\pi_1}(x_1) \dots e_{\pi_t}(x_t)$$

These numbers become extremely small – underflow

Solution: Take the logs of all values

$$V_i(t) = \log e_k(x_t) + \max_k [V_k(t-1) + \log a_{ki}]$$

Computational Complexity



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

Time: $O(K^2N)$

Space: $O(KN)$

Useful implementation technique to avoid underflows

Viterbi: sum of logs

Forward/Backward: rescaling at each position by multiplying by a constant

Model learning: two scenarios



1. **Supervised learning:** estimation when the “right answer” is known

Examples:

GIVEN: a genomic region $x = x_1 \dots x_{1,000,000}$ 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

2. **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 θ of the model to maximize $P(x|\theta)$

Supervised ML estimation



Given $x = x_1 \dots x_N$

for which the true $\pi = \pi_1 \dots \pi_N$ is known,

Define:

A_{kl} = # times $k \rightarrow l$ transition occurs in π
 $E_k(b)$ = # times state k in π emits b in x

We can show that the **maximum likelihood** parameters θ are:

$$a_{kl} = \frac{A_{kl}}{\sum_i A_{ki}} \quad e_k(b) = \frac{E_k(b)}{\sum_c E_k(c)} \quad \text{(Homework!)}$$

Supervised ML estimation



Intuition: When we know the underlying states,
Best estimate 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 θ is unreasonable
0 probabilities – VERY BAD

Example:
Given 10 casino rolls, we observe
 $x = 2, 1, 5, 6, 1, 2, 3, 6, 2, 3$
 $\pi = F, F, F, F, F, F, F, F, F, F$
Then: $a_{FF} = 1$; $a_{FL} = 0$
 $e_F(1) = e_F(3) = .2$;
 $e_F(2) = .3$; $e_F(4) = 0$; $e_F(5) = e_F(6) = .1$

Pseudocounts



Solution for small training sets:

Add pseudocounts

A_{kl} = # times $k \rightarrow l$ transition occurs in $\pi + r_{kl}$
 $E_k(b)$ = # times state k in π emits b in $x + r_k(b)$

r_{kl} , $r_k(b)$ are pseudocounts representing our prior belief

Larger pseudocounts \Rightarrow **strong prior belief**

Small pseudocounts ($\epsilon < 1$): just to avoid 0 probabilities --- **smoothing**

Unsupervised ML estimation



Given $x = x_1 \dots x_N$
for which the true $\pi = \pi_1 \dots \pi_N$ is unknown,

EXPECTATION MAXIMIZATION

0. Starting with our best guess of a model M , parameters θ :
1. Estimate A_{kl} , $E_k(b)$ in the training data
2. Update θ according to A_{kl} , $E_k(b)$
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 θ each iteration

The Baum-Welch algorithm -- comments



Time Complexity:

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

- Guaranteed to increase the log likelihood of the model

$$P(\theta | x) = P(x, \theta) / P(x) = P(x | \theta) / (P(x) P(\theta))$$

- Not guaranteed to find globally best parameters

Converges to local optimum, depending on initial conditions

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