

# Hidden Markov Models

#### Instructor: Jessica Wu -- Harvey Mudd College

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

Based on tutorial by Lawrence Rabiner

# **Basic Problems for HMMs**

Use the compact notation  $\lambda = (A, B, \pi)$ .

	single path	all paths		
scoring / evaluation	scoring $O$ , one path $q$ $P(q, O \mid \lambda)$ probability of a path & observations	scoring $O$ , all paths $P(O) = \Sigma_q P(q, O \mid \lambda)$ probability of observations over all paths [forward-backward algorithm]		
decoding	$\begin{array}{l} \text{most likely path } q^{*} \\ q^{*} = \operatorname{argmax}_{q} P(q, \ O \mid \lambda) \\ \text{[Viterbi decoding]} \end{array}$	$\begin{array}{l} \text{path containing most likely state} \\ \text{at any time point} \\ \hat{q_{i}} = \{q_{t} \mid q_{t} = \operatorname{argmax}_{S_{i}} P(q_{t} = S_{i} \mid O, \lambda)\} \\ \text{[posterior decoding]} \end{array}$		
learning	$\begin{array}{l} \text{supervised learning of } \lambda^{*} \\ \lambda^{*} = \operatorname{argmax}_{\lambda} P(q, \ O \mid \lambda) \\ \text{unsupervised learning of } \lambda^{*} \\ \lambda^{*} = \operatorname{argmax}_{\lambda} \max_{q} P(q, \ O \mid \lambda) \\ \text{[Viterbi training]} \end{array}$	$\begin{aligned} & \text{unsupervised learning} \\ \lambda^* = \operatorname{argmax}_{\lambda} \Sigma_q \ P(q, \ O \mid \lambda) \\ & \text{[Baum-Welch training]} \end{aligned}$		

Based on slides by Manolis Kellis

states	$S = \{S_1, ,  S_N\}$	states $q_t \in S$	
observations	$V = \{v_1,  \dots,  v_M\}$	observations $O_t \in V$	
initial state distribution	$ \begin{aligned} \pi &= \{\pi_i\} \\ \pi_i &= P(q_1 = S_i) \end{aligned} $	$1 \leq i \leq N$	
state transition probability distribution	$egin{aligned} A &= \{a_{ij}\}\ a_{ij} &= P(q_t = S_j \mid q_{t-1} = S_i) \end{aligned}$	$1 \leq i, j \leq N$	
observation symbol probability distribution	$B = \{ b_j(k) \} \\ b_j(k) = P(v_k \text{ at } t \mid q_t = S_j \}$	$1 \le j \le N,  1 \le k \le M$	
Forward-Backward Algorithm (Scoring)			
forward variable	$\alpha_t(i) = P(O_1 \ O_2 \ \dots \ O_t, \ q_t = S_i \mid \lambda)$	probability of partial observations $O_1 O_2 \ldots O_t$ (until time $t$ ) and state $S_i$ at time $t$ , given model $\lambda$	
backward variable	$\beta_t(i) = P(O_{t+1} \ O_{t+2} \ \dots \ O_T \mid q_t = S_i, \lambda)$	probability of partial observations $O_{t+1}O_T$ given state $S_i$ at time $t$ and model $\lambda$	
Posterior Decoding Algorithm			
	$\begin{aligned} \gamma_t(i) &= P(q_t = S_i \mid O, \lambda) \\ \gamma_t(i) &= \alpha_t(i)\beta_t(i) \ / \ \Sigma_i \ \alpha_t(i)\beta_t(i) \end{aligned}$	probability of being in state $S_i$ at time $t_{\!\!\!\!\!\!\!\!}$ given observations $O$ and model $\lambda$	
Viterbi Algorithm (Decoding)			
	$\delta_t(i) = \max_{q_1 \dots q_{t-1}} P(q_1 \dots q_t = S_i, O_1 \dots O_t \mid \lambda)$	best score (highest probability) along single path, at time $t$ , which accounts for first $t$ observations and ends in state $S_i$	

# Markov Models and Markov Chains Learning Goals • Describe the properties of a Markov chain • Describe the relationship between a Markov Chain and a Markov Model • Describe the elements of a Markov Chain

# Markov Chains and Markov Models

A Markov chain is a stochastic process with the Markov property.

#### Stochastic process

- probabilistic counterpart to a deterministic process
- a collection of r.v.'s that evolve over time

#### Markov property

• memoryless: conditional probability distribution of future states depends only on present state

		system state is	
		fully observable	partially observable
	autonomous	Markov chain	Hidden Markov Models
system is	system is controlled	Markov decision process	partially observable Markov decision process





# $\begin{aligned} & \text{Solution to Q1} \\ O &= \{S_3, S_3, S_3, S_1, S_1, S_3, S_2, S_3\} \\ P(O \mid \text{Model}) \\ &= P(S_3, S_3, S_3, S_1, S_1, S_3, S_2, S_3 \mid \text{Model}) \\ &= P(S_3) \ P(S_3 \mid S_3) \ P(S_3 \mid S_3) \ P(S_1 \mid S_3) \\ P(S_1 \mid S_1) \ P(S_3 \mid S_1) \ P(S_2 \mid S_3) \ P(S_3 \mid S_2) \\ &= \pi_3 \cdot a_{33} \cdot a_{33} \cdot a_{31} \cdot a_{11} \cdot a_{13} \cdot a_{32} \cdot a_{23} \\ &= (1)(0.8)(0.8)(0.1)(0.4)(0.3)(0.1)(0.2) \\ &= 1.536 \times 10^{-4} \end{aligned}$



## Solution to Q2

 $O = \{ S_i, S_i, S_i, \dots, S_i, S_j \neq S_i \}$ 

 $P(O \mid \text{Model}, q_1 = S_i) = (a_{ii})^{d-1}(1 - a_{ii}) = p_i(d)$ where  $p_i(d)$  is the (discrete) PDF of duration d in state i.

Notice that  $D_i \sim \text{geometric}(p)$ , where  $p = 1 - a_{ii}$  is the probability of success (exiting state *i*) and there are d - 1 failures before the first success.

Then 
$$\overline{D_i} = \frac{1}{p} = \frac{1}{1 - a_{ii}}$$

Intuition: Consider a fair die. If the probability of success (a "1") is p=1/6, it will take 1/p=6 rolls until a success.

 $\begin{array}{l} \text{the "math" way: } X \sim \operatorname{geom}(p) \\ \overline{X} = \sum_{k=1}^{\infty} k(1-p)^{k-1}p \\ = p \sum_{k=1}^{\infty} k(1-p)^{k-1} \\ = p \frac{1}{(1-(1-p))^2} = \frac{p}{p^2} = \frac{1}{p} \end{array} \text{for } x \in \mathbb{R}, \ |x| \leq 1, \\ \sum_{n=1}^{\infty} nx^{n-1} = \frac{1}{(1-x)^2} \end{aligned}$ 

For example, the expected number of consecutive days of rainy weather is  $1/a_{11} = 1/0.6 = 1.67$ ; for cloudy, 2.5; for sunny, 5.



# **Hidden Markov Models**

Learning Goals

- Describe the difference between Markov Chains and Hidden Markov Models
- Describe applications of HMMs
- Describe the elements of a HMM
- Describe the basic problems for HMMs



Now we would like to model pairs of sequences.

- There exists an underlying stochastic process that is **hidden** (not observable directly).
- But it affects observations (that we can collect directly).



# HMMs are Everywhere

application	states	observations
weather inference	seasons	
dishonest casino (casino has fair die and loaded die, casino switches between dice on average once every 20 turns)	dice used	
missile tracking	position	
speech recognition	phoneme	
NLP part-of-speech tagging	part of speech	
computational biology	protein structure	
medicine	disease (state of progression)	

# Elements of an HMMMA S-tuple (S, V, $\pi$ , A, B), where• S: finite set of states $\{S_1, \ldots, S_N\}$ • V: finite set of observations per state $\{v_1, \ldots, v_M\}$ • $\pi$ : initial state distribution $\{\pi_i\}$ • A: state transition probability distribution $\{a_{ij}\}$ • B: observation symbol probability distribution $\{b_j(k)\}$ $b_j(k) = P(v_k \text{ at } t \mid q_t = S_j\}, 1 \le j \le N$ $1 \le k \le M$

A HMM outputs only emitted symbols  $O = \{O_1, ..., O_T\}$ , where  $O_t \in V$ . Both the underlying states and random walk between states are hidden.

# HMMs as a Generative Model

Given *S*, *V*,  $\pi$ , *A*, *B*, the HMM can be used as a generator to give an observation sequence

$$O = O_1 O_2 \dots O_T.$$

- 1) Choose initial state  $q_1 = S_i \operatorname{according}$  to initial state distribution  $\pi.$
- 2) Set t = 1.
- 3) Choose  $O_t = v_k$  according to symbol probability distribution in state  $S_i$ , i.e.,  $b_i(k)$ .
- 4) Transit to new state  $q_{t+1} = S_j$  according to state transition probability distribution for state  $S_i$ , i.e.,  $a_{ij}$ .
- 5) Set t = t + 1. Return to step 3 if t < T. Otherwise stop.

# **Scoring HMMs**

Learning Goals

- Describe how to score an observation over a single path and over multiple paths
- Describe the forward algorithm





## **The Forward Algorithm**

1) Initialization  $\alpha_1(i) = \pi_i b_i(O_1), \ 1 \le i \le N$ 

2) Induction

$$\alpha_{t+1}(j) = \left[\sum_{i=1}^{N} \alpha_t(i)a_{ij}\right] b_j(O_{t+1}) \quad \begin{array}{l} 1 \le t \le T-1 \\ 1 \le j \le N \\ \text{Perform for all states for given } t, \end{array}$$

3) Termination

$$P(O|\lambda) = \sum_{i=1}^{N} \alpha_T(i)$$

AT

[Proofs for Initialization and Termination Steps]

$$lpha_1(i) = P(O_1, q_1 = S_i | \lambda)$$
  
=  $P(O_1 | q_1 = S_i, \lambda) P(q_1 = S_i | \lambda)$   
=  $b_i(O_1) \pi_i$ 

 $P(O|\lambda) = P(O_1 \dots O_T|\lambda)$  $= \sum_{i=1}^{N} P(O_1 \dots O_T, q_T = S_i | \lambda)$  $= \sum_{i=1}^{N} \alpha_T(i)$ 

 $1 \le t \le T - 1$ 

 $1 \leq j \leq N$ 

then advance t.



# The Forward Variable

We showed the induction step for  $\alpha_{t+1}(j)$  through intuition. Can we prove it?  $\begin{aligned} \alpha_{t+1}(j) &= P(O_1 \dots O_{t+1}, q_{t+1} = S_j | \lambda) \\ &= \sum_{q_1 \dots q_t} P(O_1 \dots O_t, O_{t+1}, q_1 \dots q_t, q_{t+1} = S_j | \lambda) \\ &= \left[ \sum_{q_1 \dots q_t} P(O_1 \dots O_t, q_1 \dots q_t, q_{t+1} = S_j | \lambda) \right] \frac{P(O_{t+1} | q_{t+1} = S_j, \lambda)}{P(O_{t+1} | q_{t+1} = S_j, \lambda)} \\ &= \left[ \sum_{i=1}^N \sum_{q_1 \dots q_{t-1}} P(O_1 \dots O_t, q_1 \dots q_{t-1}, q_t = S_i, q_{t+1} = S_j | \lambda) \right] \frac{b_j(O_{t+1})}{b_j(O_{t+1})} \\ &= \left[ \sum_{i=1}^N \sum_{q_1 \dots q_{t-1}} P(O_1 \dots O_t, q_1 \dots q_{t-1}, q_t = S_i | \lambda) \frac{P(q_{t+1} = S_j | q_t = S_i, \lambda)}{D(q_{t+1} = S_j | q_t = S_i, \lambda)} \right] \frac{b_j(O_{t+1})}{b_j(O_{t+1})} \\ &= \left[ \sum_{i=1}^N P(O_1 \dots O_t, q_t = S_i | \lambda) a_{ij} \right] \frac{b_j(O_{t+1})}{D(q_{t+1})} \\ &= \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] \frac{b_j(O_{t+1})}{D(q_{t+1})} \right] \frac{b_j(O_{t+1})}{D(q_{t+1})} \\ &= \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] \frac{b_j(O_{t+1})}{D(q_{t+1})} \end{bmatrix}$ 



# **Decoding HMMs**

Learning Goals

- Describe how to decode the state sequence
- Describe the Viterbi algorithm











# **Posterior Decoding**

We found the individually most likely state  $q_t$  at time t.

#### The Good

maximizes expected number of correct states

#### The Bad

- may result in invalid path (not all  $S_i \rightarrow S_j$  transitions may be possible)
- ⇒ most probable state is most likely to be correct at any instant, but sequence of individually probable states is not likely to be most probable sequence

# **Viterbi Decoding**

Goal: Find single best state sequence.

 $q^* = \operatorname{argmax}_{q} P(q \mid O, \lambda) = \operatorname{argmax}_{q} P(q, O \mid \lambda)$ 

Define  $\delta_t(i) = \max_{q_1,\dots,q_{t-1}} P(q_1\dots q_t = S_i, O_1\dots O_t | \lambda)$ 

i.e. the best score (highest probability) along a single path, at time t, which accounts for the first t observations and ends in state  $S_i$ .

Compare to algorithm for  $\alpha_t(i) = P(O_1 \dots O_t, q_t = S_i \mid \lambda)$ . To determine best path to  $q_{t+1} = S_{i}$  compute  $\delta_{t+1}(j)$ .

- $\Rightarrow \delta_t(i)$ • best path to  $q_t = S_i$
- transition from  $q_t = S_i$  to  $q_{t+1} = S_j$   $\Rightarrow a_{ij}$  emission of  $O_{t+1}$  from  $q_{t+1} = S_j$   $\Rightarrow b_j(O_{t+1})$
- to retrieve state sequence, also need traceback pointer  $\psi_t(i) =$  state  $S_i$  that maximizes  $\delta_t(i)$

# **The Viterbi Algorithm**

- 1) Initialization
- 2) Induction

3) Termination

4) Path (state sequence) backtracking



# **The Viterbi Algorithm**

 $1 \le t \le T - 1$ 

 $1 \le t \le T - 1$ 

 $1 \le j \le N$ 



2) Induction

$$\delta_{t+1}(j) = \left[ \max_{1 \le i \le N} \delta_t(i) a_{ij} \right] b_j(O_{t+1})$$
  
$$\psi_{t+1}(j) = \arg\max_{1 \le i \le N} \delta_t(i) a_{ij}$$

 $1 \le i \le N$ 3) Termination  $p^* = \max_{1 \le i \le N} \delta_T(i)$   $q^*_T = \arg \max_{1 \le i \le N} \delta_T(i)$   $q^*_T = \arg \max_{1 \le i \le N} \delta_T(i)$ 4) Path (state sequence) backtracking  $q^*_t = \psi_{t+1}(q^*_{t+1}), \ t = T - 1, T - 2, \dots, 1$ 

# The Viterbi Algorithm

- similar to forward algorithm (use max instead of sum)
- use DP table to compute
- same complexity as forward algorithm

#### **Practical Issues**

- underflow issues  $\rightarrow$  use log probabilities for model
- for logs of products of probabilities, use sum of logs

# Learning HMMs

Learning Goals

- Describe how to learn HMM parameters
- Describe the Baum-Welch algorithm

# Learning

#### Goal

Adjust the model parameters  $\lambda = (A, B, \pi)$  to maximize  $P(O \mid \lambda)$ , i.e. the probability of the observation sequence(s) given the model.

#### Supervised Approach

Assume we have complete data (we know the underlying states). Use MLE.

# Supervised Learning Example

 $\begin{array}{ll} \text{state space} & S = \{1, 2\} \\ \text{observation space} & V = \{e, f, g, h\} \\ \text{training set} & \begin{bmatrix} 1 & 2 \\ e & g \end{bmatrix} & \begin{bmatrix} 1 & 2 \\ e & h \end{bmatrix} & \begin{bmatrix} 1 & 2 \\ f & h \end{bmatrix} & \begin{bmatrix} 1 & 2 \\ f & g \end{bmatrix} \\ \end{array}$ 

What are the optimal model parameters?



# Pseudocounts

For small training set, the parameters may overfit.

- $P(O \mid \lambda)$  is maximized but  $\lambda$  is unreasonable
- probabilities of 0 are problematic

Add pseudocounts to represent our prior belief.

- large pseudocounts  $\rightarrow$  large regularization
- small pseudocounts  $\rightarrow$  small regularization (just to avoid P = 0)

# Learning

#### **Unsupervised Approach**

- we do not know the underlying states
- no known way to analytically solve for optimal model

#### Ideas

- use iterative algorithm to locally maximize  $P(O \mid \lambda)$
- either gradient descent or EM work
- Baum-Welch algorithm based on EM is most popular

# **Unsupervised Learning**

#### Goal

 $\hat{\pi}_i =$  expected frequency (number of times) in state  $S_i$  at time (t = 1)

 $\hat{a}_{ij} = \frac{\text{expected number of transitions from state } S_i \text{ to state } S_j}{\text{expected number of transitions from state } S_i}$ 

 $\hat{b}_j(k) = \frac{\text{expected number of times in state } S_j \text{ and observing symbol } v_k}{\text{expected number of times in state } S_j}$ 

Recall  $\gamma_t(i) = P(q_t = S_i \mid O, \lambda)$ , i.e. the probability of being in state  $S_i$  at time t, given observation sequence O and model  $\lambda$ . Can we use this to solve for any of the above terms?



# **Expected Number of Transitions**

Define

$$\xi_t(i, j) = P(q_t = S_i, q_{t+1} = S_j \mid O, \lambda)$$

i.e. the probability of being in state  $S_i$  at time t, and state  $S_j$  at time t+1, given the model and the observation sequence.



# **Unsupervised Learning**

#### Goal

$$\hat{\pi}_{i} = \text{expected frequency (number of times) in state } S_{i} \text{ at time } (t = 1)$$

$$= \gamma_{1}(i) = \sum_{t=1}^{T-1} \xi_{t}(i, j)$$

$$\hat{a}_{ij} = \frac{\text{expected number of transitions from state } S_{i} \text{ to state } S_{j}}{\text{expected number of transitions from state } S_{i}} = \sum_{t=1}^{T-1} \gamma_{t}(i)$$

$$= \sum_{t=1}^{T} \text{ s.t. } O_{t} = v_{k} \gamma_{t}(j)$$

$$\hat{b}_{j}(k) = \frac{\text{expected number of times in state } S_{j} \text{ and observing symbol } v_{k}}{\text{expected number of times in state } S_{j}} = \sum_{t=1}^{T} \gamma_{t}(j)$$



# **Baum-Welch Algorithm** Initialization • Set $\lambda = (A, B, \pi)$ to random initial conditions (or using prior information) Iteration (repeat until convergence) • Compute $\alpha_t(i)$ and $\beta_t(i)$ using forward-backward algo $\Rightarrow$ Compute $P(O \mid \lambda)$ [E-step] • Compute $\gamma_t(i)$ and $\xi_t(i,j)$ $\Rightarrow$ Update model parameters [M-step]

# **Baum-Welch Algorithm**

• Time complexity:  $O(N^2 T) \cdot (\# \text{ iterations})$ 

• Guaranteed to increase likelihood  $P(O \mid \lambda)$  via EM but *not* guaranteed to find globally optimal  $\lambda^*$ 

#### **Practical Issues**

- Use multiple training sequences (sum over them)
- Apply smoothing to avoid zero counts and improve generalization (add pseudocounts)

# HMMs and Protein Structure

One biological application of HMMs is to determine the secondary structure (i.e. the general three-dimensional shape) of a protein. This general shape is made up of alpha helices, beta sheets, and other structures. In this problem, we will assume that the amino acid composition of these regions is governed by an HMM.

To keep this problem relatively simple, we do not use actual transition values or emission probabilities. The start state is always "other". We will use the state transition probabilities and emission probabilities below.

		alaba	hota	othor
		aipiia	Dela	other
	alpha	0.7	0.1	0.2
	beta	0.2	0.6	0.2
	other	0.3	0.3	0.4

e.g.  $P(Alpha Helix \rightarrow Beta Sheet) = 0.1$ 

Based on exercise by Manolis Kellis

amino acid	alpha	beta	other
Μ	0.35	0.10	0.05
L	0.30	0.05	0.15
N	0.15	0.30	0.20
E	0.10	0.40	0.15
А	0.05	0.00	0.20
G	0.05	0.15	0.25

# **Protein Structure Questions**

1) What is the probability  $P(q = O\alpha, O = ML)$ ?

- 2) How many paths could give rise to the sequence O = MLN? What is the total probability P(O)?
- 3) Give the most likely state transition path  $q^*$  for the amino acid sequence MLN using the Viterbi algorithm. What is  $P(q^*, O)$ ?

Compare this to P(O) above. What does this say about the reliability of the Viterbi path?

Based on exercise by Manolis Kellis