

## Markov Models

- Can be used for continuous or discrete sequences
- continuous data can be turned into discrete by vector quantization

Markov models are generative models of sequences!

$$W_t = \text{weather over time } \in \{r, s\}$$

rain	sun	1	2	3	4	5	...
		r	r	s	s	r	

ex:  $P(W_7 = s | W_1, \dots, W_6) \rightarrow$  dependence on previous states

$P(W_t | W_{t-1}) \rightarrow$  1st order Markov

$$\text{ex: } n = \text{sunny} \quad r = \text{rainy} \quad \pi \rightarrow \text{starting probabilities}$$

	n	r	
w_{t-1}	w_t	n	r

$$\pi = [0.9 \quad 0.1]$$

$$\begin{matrix} n & \begin{bmatrix} 0.9 & 0.1 \\ 0.5 & 0.5 \end{bmatrix} \\ r & \end{matrix} = a \quad S = \{n, r\}$$

$$a_{ij} = P(W_t = s_j | W_{t-1} = s_i)$$

$\nearrow$   
State  $s_j$

## Generative Markov Sequence

1. select  $s \in S$  according to starting probability  $\pi$

$$\pi = [0.9 \quad 0.1]$$

2. Select  $s_t$  from row corresponding to  $s_{t-1}$

$$a = \begin{bmatrix} n & r \\ 0.9 & 0.1 \\ 0.5 & 0.5 \end{bmatrix}_r^n w_{t-1}$$

ex:  $\boxed{t} \quad 1 \quad 2 \quad 3$

$\boxed{s} \quad n \quad n \quad r$

$$S = \{n, r\}$$

\* code and simulate!

In reality ...

- you calculate the probabilities over all possible transition states

ex:

$$\begin{aligned} P(W_0 = n) &= \pi_1 \\ P(W_1 = n) &= \pi_1 \cdot d_{nn} \\ &= 0.9 \cdot 0.9 \\ &= 0.81 \end{aligned}$$

two possible outcomes  
each have a probability

$$\begin{aligned} P(W_1 = r) &= \pi_1 \cdot d_{nr} \\ &= 0.9 \cdot 0.1 \\ &= 0.09 \end{aligned}$$

$$d = \begin{bmatrix} n & r \\ 0.9 & 0.1 \\ 0.5 & 0.5 \end{bmatrix} \quad \begin{matrix} w_t \\ n \\ r \\ w_{t-1} \end{matrix}$$
$$\pi = \begin{bmatrix} n & r \\ 0.9 & 0.1 \end{bmatrix}$$

in the general case ... (first order)

$$\begin{aligned} P(W_0) &= \pi \\ P(W_1) &= P(W_0) \cdot d \\ &\quad \begin{matrix} 1 \times 2 \\ 1 \times 2 \end{matrix} \quad \begin{matrix} 2 \times 2 \end{matrix} \end{aligned}$$

$$P(W_t) = P(W_{t-1}) \cdot d = P(W_{t-2}) \cdot d^2 = \underbrace{P(W_0)}_{\pi} \cdot d^t$$

in python: (try it out)

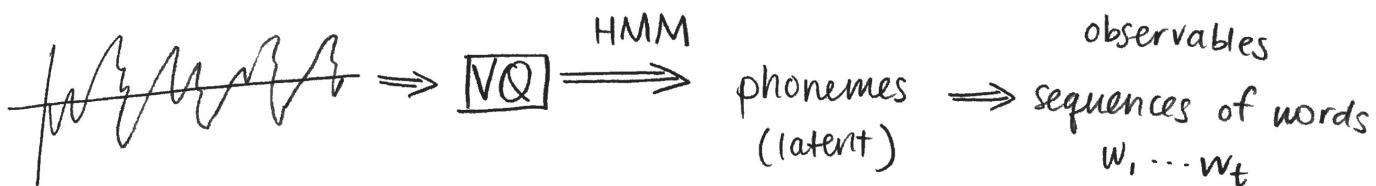
```
pi = np.array([1.0, 0.0])
d = np.array([[0.9, 0.1], [0.5, 0.5]])
P_Wt = np.dot(pi, np.linalg.matrix_power(d, t))
t is the # time points
```

\* Markov models converge onto a fixed point depending on the trans kernel  $d$  (regardless of  $\pi$ ) yay eigenvalues and dynamical systems

## Hidden Markov Models

first used in domain of speech recognition (Rabiner)

continuous, acoustic signal



- calculate  $P(\text{phoneme sequence} | \text{acoustic sequence})$

ex:  $P(\text{sequence}) = 0.1 \times 0.5 \times 0.5 \times 0.9 \times 0.9 \dots \text{etc}$

## HMM continued...

def:  $\theta = [s, \pi, a]$

you can find  $P(\text{seq}; \theta)$  using NLL function:  
 $L(\theta; \text{seq})$

$$\pi_t = P(X_0 = s_i)$$

can also estimate params from data

$$a_{ij} = P(X_{t+1} = S_d | X_t = S_i)$$

$$\hat{\pi}_i = \frac{\# \text{ times } s_i \text{ occurs at start of sequence}}{m_i \text{ states}}$$

$$\hat{a}_{ij} = \frac{\# \text{ times } ij \text{ transition occurs}}{\# \text{ times in state } i}$$

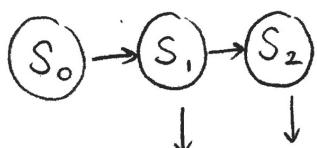
$$\frac{\partial L(D; \theta)}{\partial \pi} = 0, \quad \frac{\partial L(D; \theta)}{\partial a} = 0$$

HMM setup has 2 more parameters:

$$\theta = [s^{1 \dots n}, o^{1 \dots k}, \pi^{1 \times n}, a^{n \times n}, b^{n \times k}]$$

↓              ↓              ↓              ↓              ↓  
 n hidden states    k observed states    initial prob state of s    transition kernel    emission probabilities  
 ↑              ↑              ↑              ↑              ↑  
 \* S and O also go from  $1 \leq t \leq T$  time steps!!

=  
 $P(\text{observe a state} | \text{hidden state})$



$$a_{ij} = P(S_{t+1} = s_j | S_t = s_i)$$

$$o_1, o_2$$

$$b_{ik} = P(o_t = o_k | S_t = s_i)$$

(probably the most important page)

## Inference in HMMs

and

## Estimation in HMMs

Inference problems:

- 1) Filtering : inferring hidden state  $S$  from observations  $O$
- 2) Smoothing : using future observations  $O$  to smooth predictions of past hidden states
- 3) Decoding : find choice of hidden variable assignments that has highest probability given known observations

## Filtering $P(S_t | O_1, \dots, O_t)$

Dishonest Casino

$$S = \{F, L\} \quad O = \{h, t\} \quad \pi = [1, 0]$$

$$b = F \begin{bmatrix} h & t \\ 0.5 & 0.5 \\ 0.1 & 0.9 \end{bmatrix} \quad (n \times k)$$

$$\text{ex: } \left. \begin{array}{l} P(S_0 = F) = 1 \\ P(S_0 = L) = 0 \end{array} \right\} \pi$$

$$a = F \begin{bmatrix} F & L \\ 0.95 & 0.05 \\ 0.1 & 0.9 \end{bmatrix} \quad s_{t-1}$$

$$\text{find } P(S_1 = F) = \sum_{S_0} P(S_1 = F | S_0) P(S_0) \Rightarrow P(S_1 = F | O_1 = h) = \frac{P(O_1 = h | S_1 = F) P(S_1 = F)}{P(O_1 = h)}$$

we use forward algorithm to estimate all the probabilities:

Small example:  $P(S_1 = F | O_1 = h) = 0.5 \times 0.95 = 0.475 \Rightarrow \text{normalize!!}$

$$P(S_1 = L | O_1 = h) = 0.1 \times 0.05 = 0.0005 \Rightarrow \sum_i p_i = 0.48$$

$$P(S_1 | O_1 = h) = [0.475 \quad 0.0005] ./ 0.48$$

$$= [0.99 \quad 0.01]$$

for our example:  
 $n=2$  states

$T=2$  states

	F	L
$O$	$\alpha_0$	1
heads	$\alpha_1$	0.99 0.01
tails	$\alpha_2$	0.058 0.942

$$P(S_1 = F | O_1 = t) = 0.5 \times 0.1 = 0.05$$

$$P(S_1 = L | O_1 = t) = 0.9 \times 0.9 = 0.81$$

$$\sum_i p_i = 0.86$$

$$P(S_1 | O_1 = t) = [0.058 \quad 0.942] \Rightarrow \text{these are } \alpha_1(i)$$

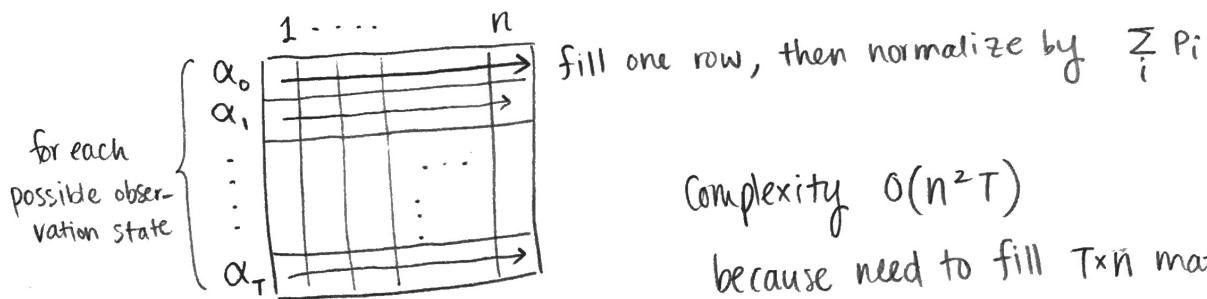
(defined on next page)

## generalization of forward algorithm

Define  $\alpha_t(i) = P(S_t = s_i; o_1 \dots o_T)$

$$\alpha_{t+1}(j) = b_j(o_{t+1}) \sum_{i=1}^{n_{\text{states}}} a_{ij} \alpha_t(i) \quad \text{with } \alpha_0(i) = \pi_i$$

$n$  states (hidden)  $1 \leq i, j \leq n$



Complexity  $O(n^2 T)$

because need to fill  $T \times n$  matrix  $n$  times!

$$0 \leq j \leq n \\ 1 \leq t \leq T$$

$$S_0 \rightarrow \circled{S_1} \rightarrow \circled{S_2} \rightarrow \dots \rightarrow S_n$$

$$\downarrow$$

$$o_1 \quad o_2$$

$$o_K$$

$$\text{complexity } T \times n \quad T \times n$$

$$T \times n \Rightarrow O(n^2 T)$$

## Smoothing

$$(R < t)$$

$$P(S_R | o_1 \dots o_t) = \frac{P(o_{R+1} \dots o_t | S_R, o_1 \dots o_R) P(S_R; o_1 \dots o_R)}{P(o_1 \dots o_t)}$$

$$\underbrace{\alpha R}_{\text{backward algorithm}} \underbrace{P(o_{R+1} \dots o_t | S_R) P(S_R; o_1 \dots o_R)}_{\alpha_R \text{ (forward algorithm/filtering!!)}}$$

backward algorithm

$\alpha_R$  (forward algorithm/filtering!!)

## backward algorithm

Define  $\beta_R = P(o_{R+1} \dots o_t | S_R) \quad \beta_T(i) = 1 \quad 1 \leq i \leq n$

$$\beta(t) = \sum_{j=1}^n a_{ij} \beta_{t+1}(j) b_j(o_{t+1})$$

$$O(n^2 T)$$

## Posterior Decoding

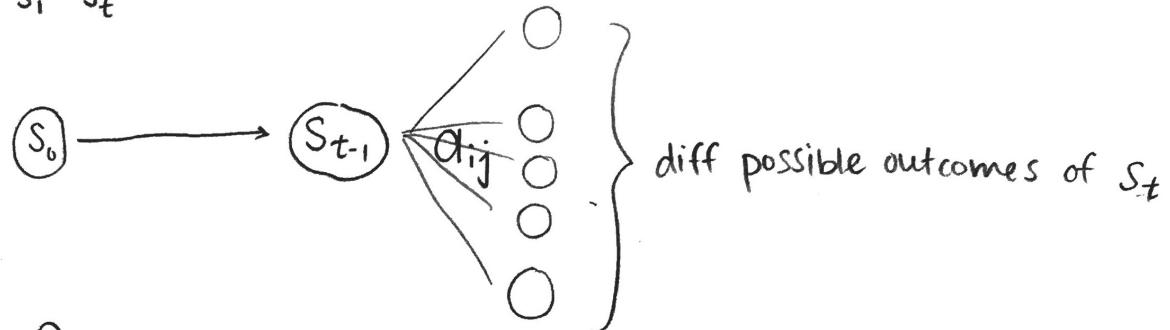
$$\alpha_t(i) \beta_t(i) \propto P(S_t = s_i | o_1, \dots, o_T) \quad \text{combine forward/backward}$$

Viterbi Algorithm: find most possible sequence of hidden variable assignments

recall there are  $2^T$  possible sequences for a 2 state sequence

\* but there is a special property of most likely sequences

$$\underset{s_1, \dots, s_t}{\operatorname{argmax}} \quad P(s_1, \dots, s_t | o_1, \dots, o_t)$$



$$\text{Define } \delta_t(i) = \underset{s_0, \dots, s_{t-1}}{\operatorname{argmax}} \quad P(s_t = s_i, s_0, \dots, s_{t-1}, o_1, \dots, o_t)$$

$$\delta_0 = \pi_i \quad 1 \leq i \leq n \quad O(n^2 T)$$

$$\delta_{t+1}(j) = b_j(o_{t+1}) \max_i a_{ij} \delta_t(i)$$

	1	2	...	-	n
$\delta_0$					
$\delta_1$					
$\vdots$					
$\delta_T$					

## Estimating HMMs from data

- given observation sequences and corresponding hidden states, find

O

S

$$\Theta = [\pi, a, b]$$

$$\hat{\pi}_i = \frac{\# \text{ times } S_0 = i}{\# \text{ sequences}}$$

$$\hat{a}_{ij} = \frac{\# \text{ times } i \rightarrow j \text{ transitions}}{\# \text{ times in state } i}$$

$$\hat{b}_k(i) = \frac{\# \text{ times } i \in O_k \text{ occurs}}{\# \text{ times in state } i}$$

Using EM to guess an algorithm:

1. Guess  $\Theta = [\pi, a, b]$

2. Use sequence  $O, \dots, O_T$  and  $\Theta$  to infer  $\hat{a}_{ij}$ ,  $\hat{b}_k(i)$ , and  $\hat{\pi}_i$

3. Update  $\Theta$

until  
no  
change  
in  $\Theta$

## Baum-Welch (EM) for HMMs

1. Guess  $\hat{\Theta} = [\hat{\pi}, \hat{a}, \hat{b}]$  E-step

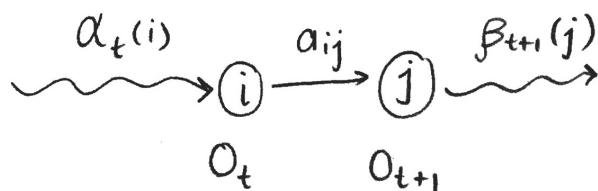
→ What is probability of being in state  $i$  at time  $t$  and moving to state  $j$   
 estimate of  $\hat{a}_{ij}$  →  $\hat{\xi}_t(i, j) = P(S_t = i, S_{t+1} = j | O, \hat{\Theta})$

$$= \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^n \sum_{j=1}^n \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}$$

→  $\gamma_t(i) = \text{expected times in } i$

$$= P(S_t = s_i | O, \hat{\Theta})$$

$$= \sum_{j=1}^n \hat{\xi}_t(i, j)$$



M-Step Given  $\hat{\gamma}_t$  and  $\tilde{z}_t$

$$\text{re-estimate } \hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \tilde{z}_t(i, j)}{\sum_{t=1}^{T-1} \hat{\gamma}_t(i)} \quad \hat{\pi}_i = \hat{\gamma}_1(i)$$

$$\hat{b}_i(k) = \frac{\sum_{t=1}^T \hat{\gamma}_t(i) \underset{o_t=k}{\text{(or)}}}{\sum_{t=1}^T \hat{\gamma}_t(i)} \quad \begin{matrix} (\# \text{ times in state } i \text{ and observed symbol } k) \\ (\# \text{ times in state } i) \end{matrix}$$