

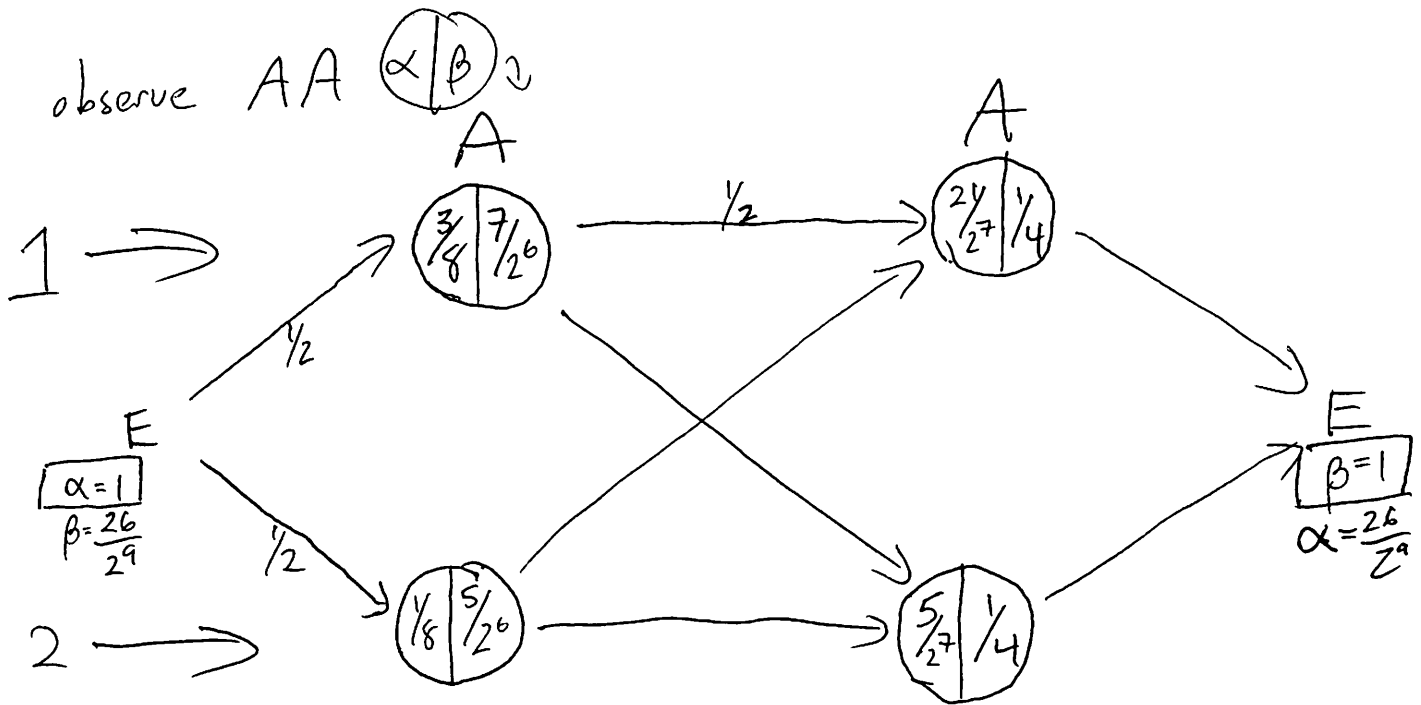
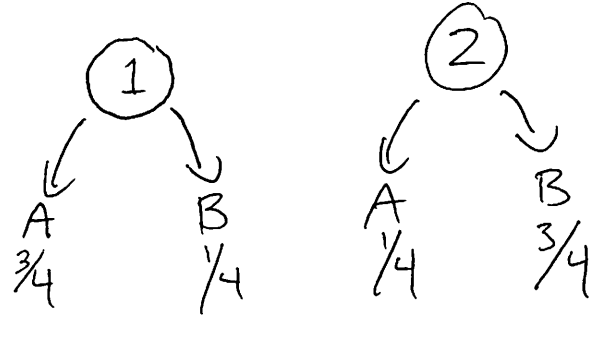
Hidden Markov Models

①

Transition Matrix

	1	2	E
1	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$
2	$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
E	$\frac{1}{2}$	$\frac{1}{2}$	0

special E boundary state



$\alpha(i, t)$ = total mass of all paths that get to state i @ time t and emit symbol w_t , up to and including the emission

$\beta(i, t)$ = total mass of all paths that leave state i @ time t , not including the emission @ time t

$$\alpha(1, 2) = \left(\frac{3}{8} \cdot \frac{1}{2} + \frac{1}{8} \cdot \frac{1}{4}\right) \cdot \frac{3}{4} = \frac{21}{2^7}$$

$$\beta(1, 1) = \frac{1}{2} \cdot \frac{3}{4} \cdot \frac{1}{4} + \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{7}{2^6}$$

$$\alpha(2, 2) = \left(\frac{3}{8} \cdot \frac{1}{4} + \frac{1}{8} \cdot \frac{1}{2}\right) \cdot \frac{1}{4} = \frac{5}{2^7}$$

$$\beta(2, 1) = \frac{1}{4} \cdot \frac{5}{4} \cdot \frac{1}{4} + \frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{5}{2^6}$$

$$\alpha(E, 0) = 1$$

$$\alpha(i, t+1) = \left(\sum_j \alpha_{jt} P(j \rightarrow i) \right) P(i \Rightarrow w_{t+1})$$

$$\beta(E, L_{t+1}) = 1$$

$$\beta(i, t-1) = \sum_j P(i \rightarrow j) P(j \Rightarrow w_t) \beta_{jt}$$

NOTE

$\alpha(E, L_{t+1}) = \beta(E, 0) = Z$ total sequence probability
also for any t , $\sum_{i \in S} \alpha_{it} \beta_{it} = Z$, b/c all paths must go through exactly one state at time t
(great for debugging EM!)

updates

$$\mathbb{E} \text{ count of } 1 \rightarrow 1 \text{ @ time } 1 = \frac{\alpha(1,1) P(1 \rightarrow 1) P(1 \rightarrow A) \beta(1,2)}{Z}$$

$$\mathbb{E} \text{ count of } 2 \rightarrow H \text{ @ time } 1 = \frac{\alpha(2,1) \beta(2,1)}{Z}$$

for EM just sum all these \mathbb{E} counts over all times in all sequences, normalize, and that's it.

NOTE

$\mathbb{E}(1 \rightarrow 1) = \frac{9/2^8}{2^6/2^8} = 9/13 > 1/2$ so the transition $1 \rightarrow 1$ is "encouraged" to increase

Decoding

③

Viterbi

need for Dyn
Prog table
not needed
for argmax

$V(i, t) =$ best path through i at time t

$$= \underset{j}{\text{best path from}} \operatorname{argmax} V(j, t-1) \cdot \text{prob } P(j \rightarrow i) P(i \rightarrow w_t)$$

Max Prob Tag

for time t pick $\operatorname{argmax}_i \alpha_i t \beta_i t$

can lead to different results from Viterbi, often better fit to real world needs (# of correct tags)

A Rainbow of HMMs

(4)

Continuous Outputs

Speech Recognition \rightarrow FFT snippets

OCR \rightarrow Like MNIST, but w/ a Language Model
framework

Stocks \rightarrow use generative/predictive power

Hidden Semi-Markov Models

- HMM assumes #self transition RV is geometrically distributed

- HSMM uses arbitrary different distribution, but can't use normal FB algorithm

Sticky HMMs encourage self transitions through bias in the prior - important for HDP-HMMs, which have potentially ∞ # of states

used in Unsupervised POS tagging

consider

TOK	"	That's	what	she	said	"	she	said	.	
POS	"	DT	VBZ	WP	PRP	VPD	"	PRP	VPD	.

the 2 "said"s will be confused b/c of their different contexts (followed by " vs .)

UPOS Tagging Cont

(3)

also, we get 2 classes of "the",
because of different noun types

Tricks to improve performance

Condition on more nodes (Tri/Bi/Uni-
something)

Use Morphology to figure out rare words
i.e. if running, swimming, etc are
all in a class, then spelunking
should be too, even without such
context info.

posterior Regularization — how do you
constrain one word to be generated by
a small # of states?