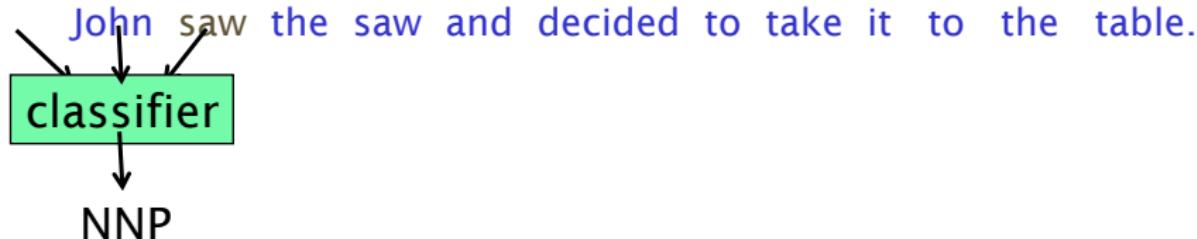




Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



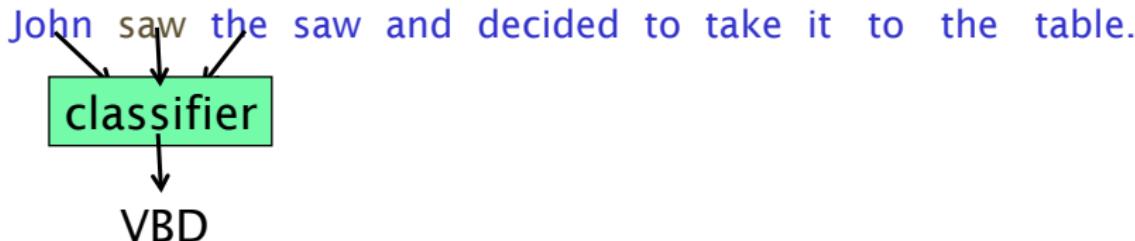
Slide from Ray Mooney

Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



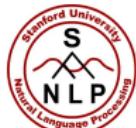
Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).



Slide from Ray Mooney

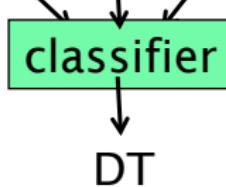
Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



Slide from Ray Mooney

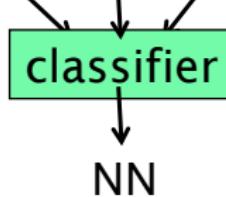
Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Sequence Labeling as Classification

- Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table.



Slide from Ray Mooney

Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



5. Word-Window classification

- **Idea:** classify a word in its context window of neighboring words.
- For example, **Named Entity Classification** of a word in context:
 - Person, Location, Organization, None
- A simple way to classify a word in context might be to **average** the word vectors in a window and to classify the average vector
 - Problem: that would **lose position information**



Window classification: Softmax

- Train softmax classifier to classify a center word by taking concatenation of word vectors surrounding it in a window
- Example: Classify “Paris” in the context of this sentence with window length 2:

... museums in Paris are amazing ...

● ● ● ●
● ● ● ●
● ● ● ●
● ● ● ●
● ● ● ●

$$\mathbf{X}_{\text{window}} = [\mathbf{x}_{\text{museums}} \quad \mathbf{x}_{\text{in}} \quad \mathbf{x}_{\text{Paris}} \quad \mathbf{x}_{\text{are}} \quad \mathbf{x}_{\text{amazing}}]^T$$

- Resulting vector $\mathbf{x}_{\text{window}} = \boxed{\mathbf{x} \in \mathbb{R}^{5d}}$, a column vector!



Simplest window classifier: Softmax

- With $x = x_{window}$ we can use the same softmax classifier as before

*predicted model
output
probability*

$$\hat{y}_y = p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$$

- With cross entropy error as before:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right)$$

- How do you update the word vectors?
- Answer: Just take derivatives like last week and optimize



Slightly more complex: Multilayer Perceptron

- Introduce an additional layer in our softmax classifier with a non-linearity.
- MLPs are fundamental building blocks of more complex neural systems!
- Assume we want to classify whether the center word is a Location
- Similar to word2vec, we will go over all positions in a corpus. But this time, it will be supervised s.t. positions that are true NER Locations should assign high probability to that class, and others should assign low probability.



Neural Network Feed-forward Computation

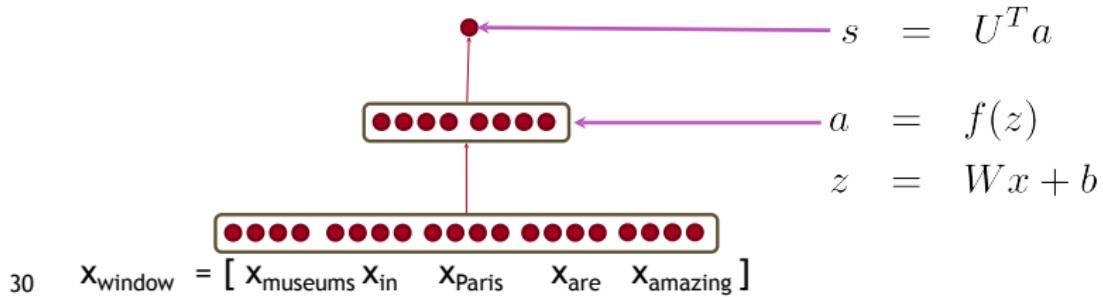
$$\text{score}(x) = U^T a \in \mathbb{R}$$

We compute a window's score with a 3-layer neural net:

- $s = \text{score}(\text{"museums in Paris are amazing"})$

$$s = U^T f(Wx + b)$$

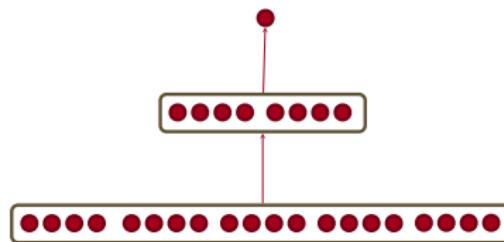
$$x \in \mathbb{R}^{20 \times 1}, W \in \mathbb{R}^{8 \times 20}, U \in \mathbb{R}^{8 \times 2}$$





Main intuition for extra layer

The middle layer learns **non-linear interactions** between the input word vectors.



$$\begin{aligned} \mathbf{x}_{\text{window}} = [& \mathbf{x}_{\text{museums}} \quad \mathbf{x}_{\text{in}} \quad \mathbf{x}_{\text{Paris}} \quad \mathbf{x}_{\text{are}} \\ & \mathbf{x}_{\text{amazing}}] \end{aligned}$$

Example: only if “*museums*” is first vector should it matter that “*in*” is in the second position



Content

- 1 Sequence labeling problems
- 2 Word window classification
- 3 Hidden Markov models (HMMs)
- 4 Graphical models for sequence labeling



Weakness of word classification

- A weakness of the word classification method for sequence labeling is that it is not capable to make use of the dependencies between POS tags in prediction.
- An extreme example:

It is true for all that that that that that that that refers to is not the same that that that refers to.

It	is	true	for	all	that	that	that	that	that	that	refers	to	is	not	the	same	that	that	that	that	refers	to	.
					pron.	conj.	det.	noun	rel.pron.	det.	noun						noun	rel.pron.	det.	noun			
					(adj.)	"that"	which	(adj.)	"that"							"that"	which	(adj.)	"that"				

- The word windows for some of the occurrences of “that” in this sentence are the same, if the window size is not greater than 2.
- The word window classification method will not be able to distinguish those “that”.
- The POS tags of the previous words may be helpful!



Sequence Labeling as Classification Using Outputs as Inputs

- Better input features are usually the **categories** of the surrounding tokens, but these are not available yet.
- Can use category of either the preceding or succeeding tokens by going forward or back and using previous output.

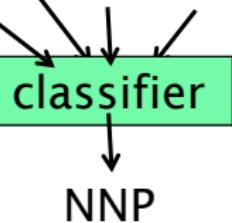
Slide from Ray Mooney

Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Forward Classification

John saw the saw and decided to take it to the table.

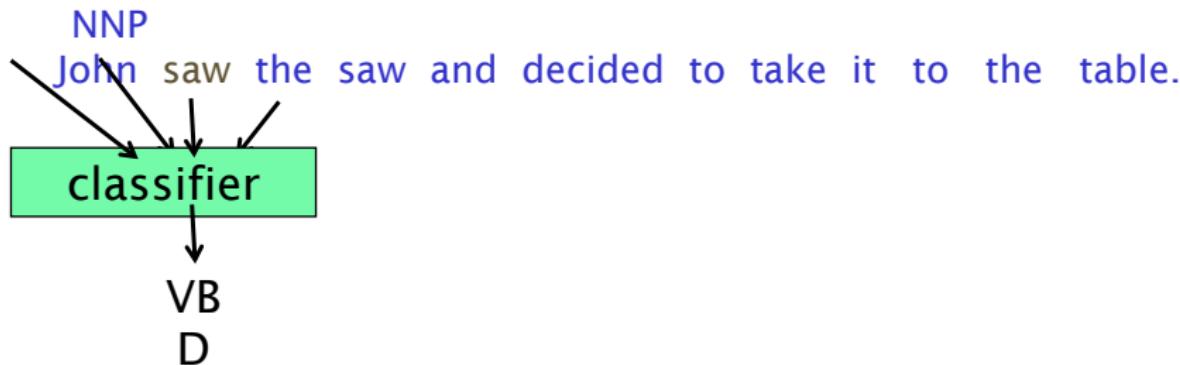


Slide from Ray Mooney

Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Forward Classification

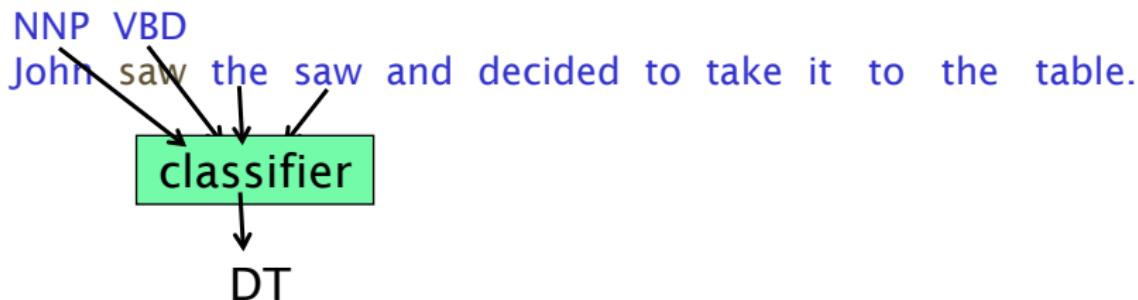


Slide from Ray Mooney

Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Forward Classification



Slide from Ray Mooney

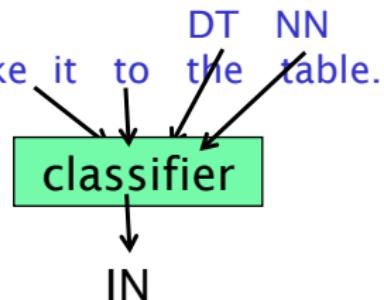
Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Backward Classification

- Disambiguating “to” in this case would be even easier backward.

John saw the saw and decided to take it to the table.



Slide from Ray Mooney

Daniel Jurafsky and James H. Martin, Part-of-speech tagging (slides)



Content

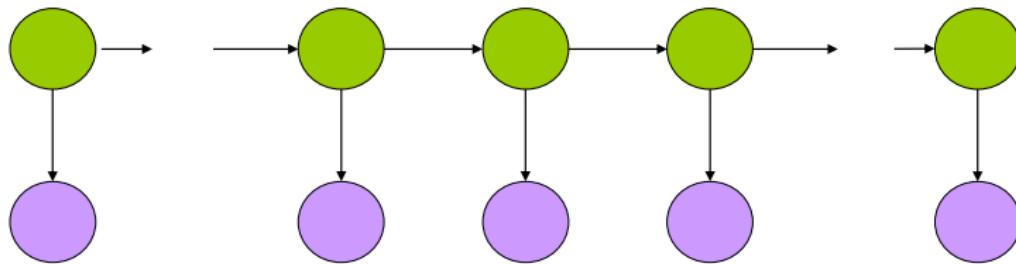
3

Hidden Markov models (HMMs)

- Model definition
- Inference in an HMM
- Decoding
- Forward Procedure
- Backward Procedure
- Decoding Solution
- Viterbi Algorithm
- Parameter Estimation
- HMM Applications



What is an HMM?

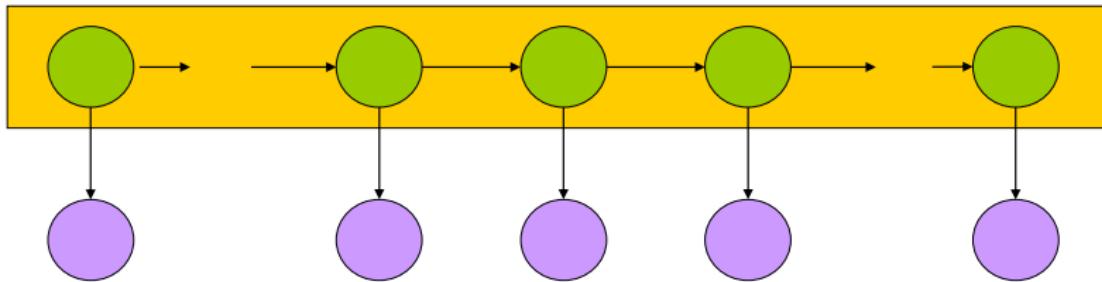


- Graphical Model
- Circles indicate states
- Arrows indicate probabilistic dependencies between states

David Meir Blei, Hidden Markov Models, 1999 (Slides)



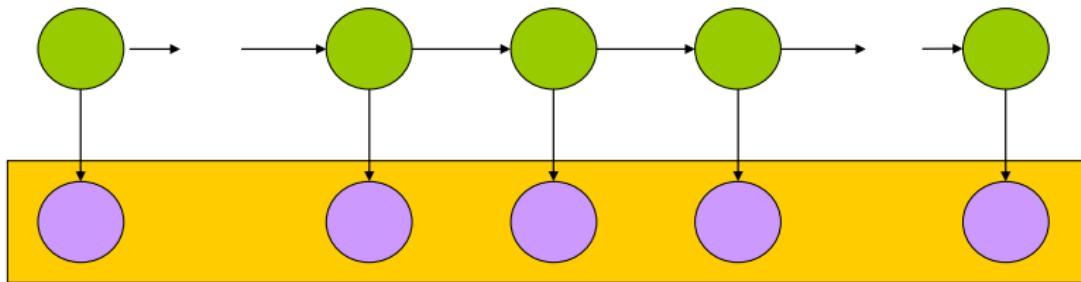
What is an HMM?



- Green circles are *hidden states*
- Dependent only on the previous state
- “The past is independent of the future given the present.”



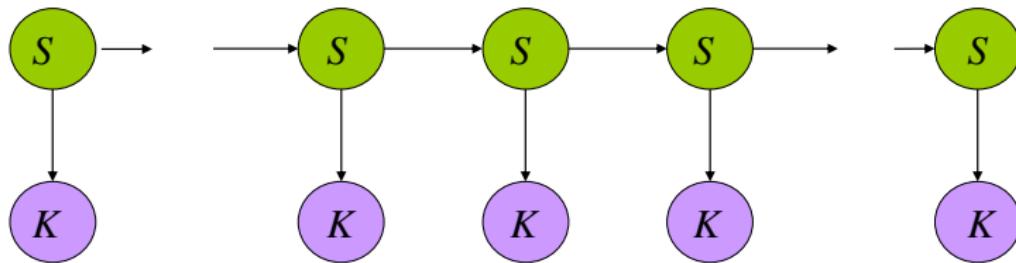
What is an HMM?



- Purple nodes are *observed states*
- Dependent only on their corresponding hidden state



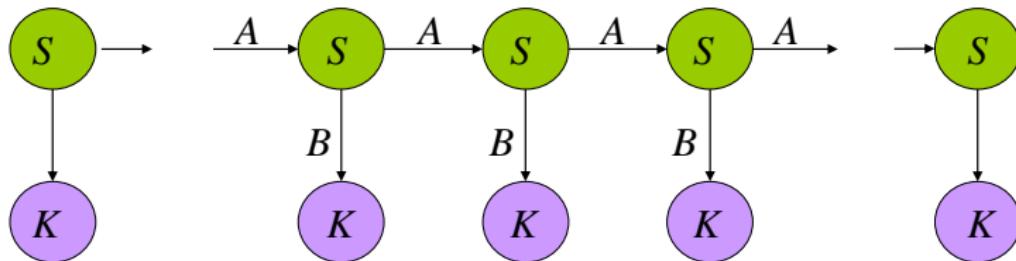
HMM Formalism



- $\{S, K, \Pi, A, B\}$
- $S : \{s_1 \dots s_N\}$ are the values for the hidden states
- $K : \{k_1 \dots k_M\}$ are the values for the observations



HMM Formalism

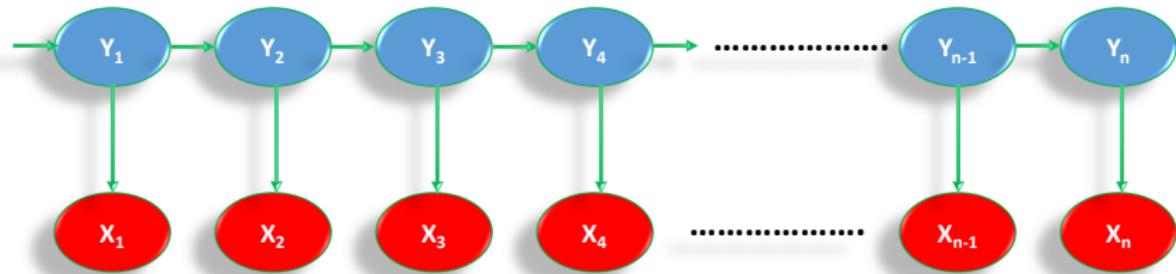


- $\{S, K, \Pi, A, B\}$
- $\Pi = \{\pi_i\}$ are the initial state probabilities
- $A = \{a_{ij}\}$ are the state transition probabilities
- $B = \{b_{ik}\}$ are the observation state probabilities



Trellis Diagram

- An HMM can be graphically depicted by Trellis diagram



Anantharaman Narayana Iyer, Natural Language Processing, Unit 2 – Tagging Problems and HMM (Slides)



Another example: Coin tossing

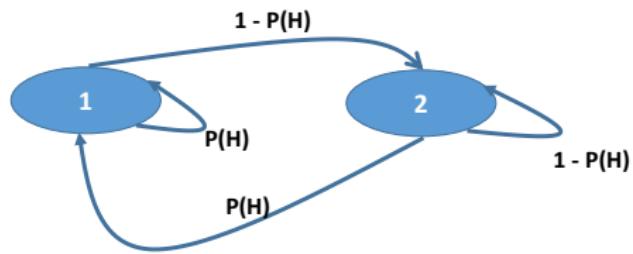


Fig (a)

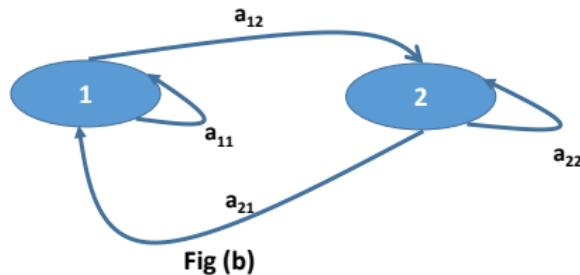


Fig (b)

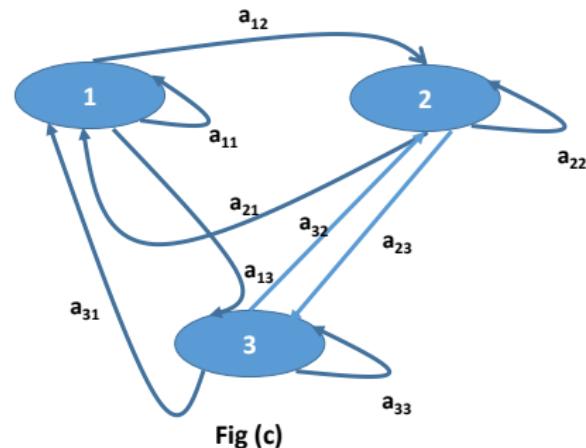


Fig (c)

	State 1	State 2	State 3
P(H)	0.5	0.75	0.25
P(T)	0.5	0.25	0.75

Anantharaman Narayana Iyer, Natural Language Processing, Unit 2 – Tagging Problems and HMM (Slides)



Content

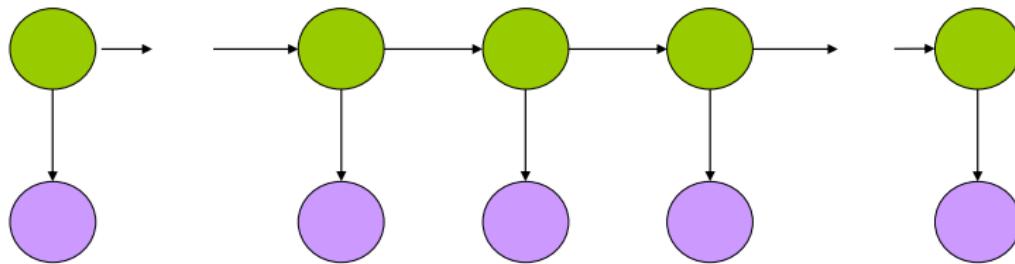
3

Hidden Markov models (HMMs)

- Model definition
- Inference in an HMM
- Decoding
- Forward Procedure
- Backward Procedure
- Decoding Solution
- Viterbi Algorithm
- Parameter Estimation
- HMM Applications



Inference in an HMM



- Compute the probability of a given observation sequence
- Given an observation sequence, compute the most likely hidden state sequence
- Given an observation sequence and set of possible models, which model most closely fits the data?

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Content

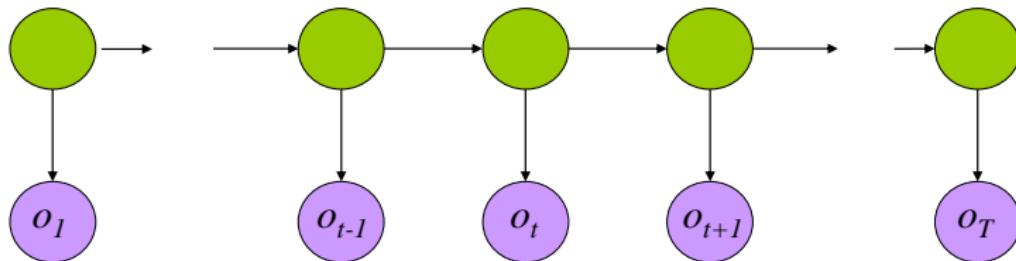
3

Hidden Markov models (HMMs)

- Model definition
- Inference in an HMM
- **Decoding**
- Forward Procedure
- Backward Procedure
- Decoding Solution
- Viterbi Algorithm
- Parameter Estimation
- HMM Applications



Decoding



Given an observation sequence and a model,
compute the probability of the observation sequence

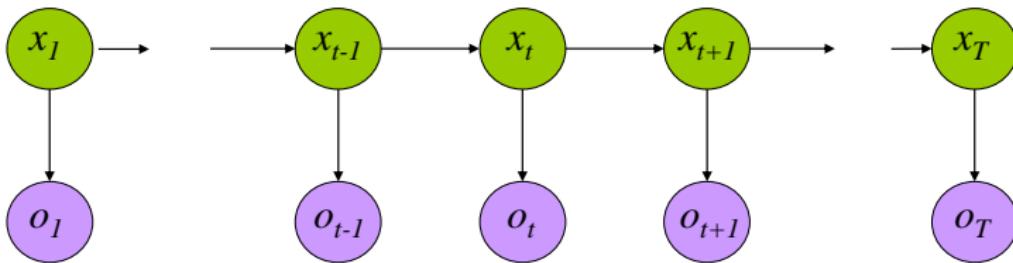
$$O = (o_1 \dots o_T), \mu = (A, B, \Pi)$$

Compute $P(O | \mu)$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



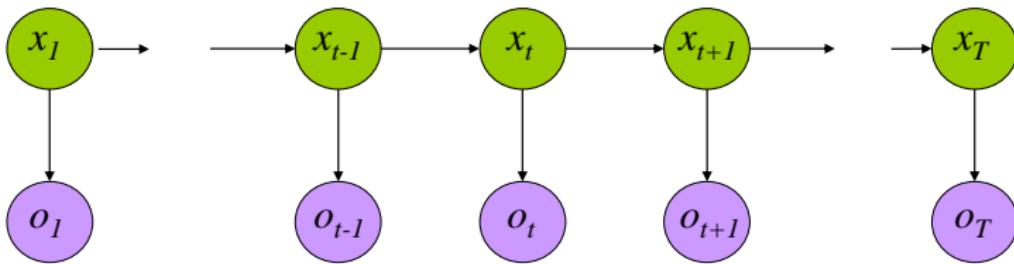
Decoding



$$P(O | X, \mu) = b_{x_1 o_1} b_{x_2 o_2} \dots b_{x_T o_T}$$



Decoding

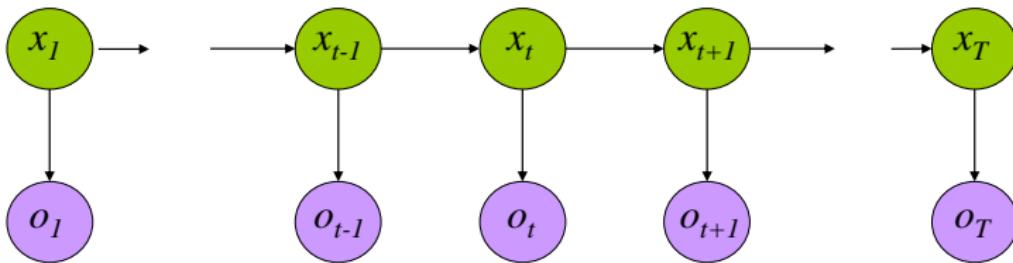


$$P(O | X, \mu) = b_{x_1 o_1} b_{x_2 o_2} \dots b_{x_T o_T}$$

$$P(X | \mu) = \pi_{x_1} a_{x_1 x_2} a_{x_2 x_3} \dots a_{x_{T-1} x_T}$$



Decoding



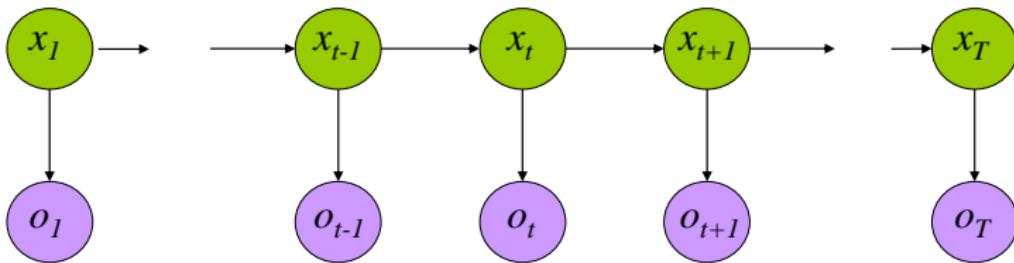
$$P(O | X, \mu) = b_{x_1 o_1} b_{x_2 o_2} \dots b_{x_T o_T}$$

$$P(X | \mu) = \pi_{x_1} a_{x_1 x_2} a_{x_2 x_3} \dots a_{x_{T-1} x_T}$$

$$P(O, X | \mu) = P(O | X, \mu)P(X | \mu)$$



Decoding



$$P(O | X, \mu) = b_{x_1 o_1} b_{x_2 o_2} \dots b_{x_T o_T}$$

$$P(X | \mu) = \pi_{x_1} a_{x_1 x_2} a_{x_2 x_3} \dots a_{x_{T-1} x_T}$$

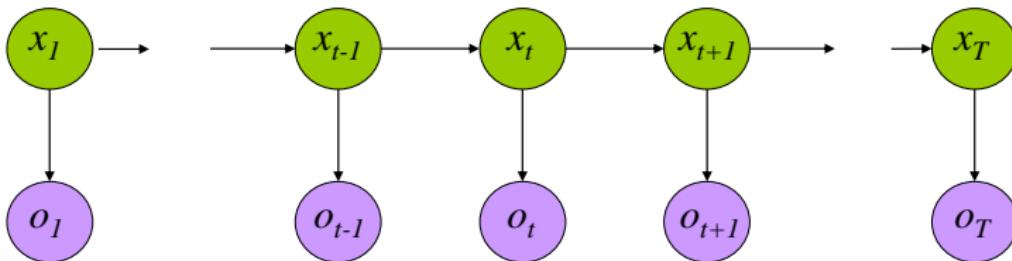
$$P(O, X | \mu) = P(O | X, \mu) P(X | \mu)$$

$$P(O | \mu) = \sum_X P(O | X, \mu) P(X | \mu)$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Decoding



$$P(O | \mu) = \sum_{\{x_1 \dots x_T\}} \pi_{x_1} b_{x_1 o_1} \prod_{t=1}^{T-1} a_{x_t x_{t+1}} b_{x_{t+1} o_{t+1}}$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Content

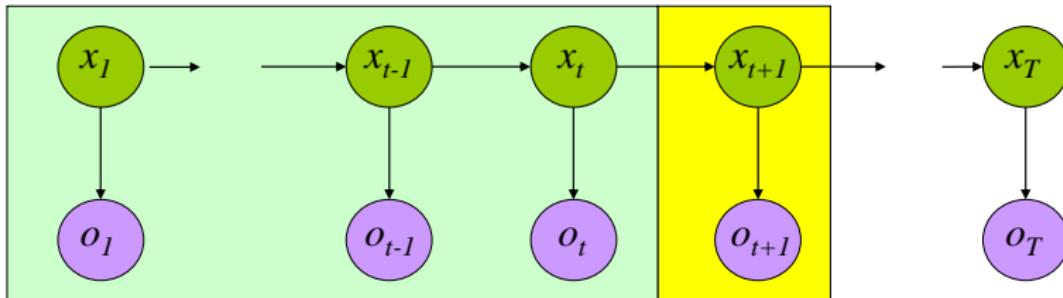
3

Hidden Markov models (HMMs)

- Model definition
- Inference in an HMM
- Decoding
- **Forward Procedure**
- Backward Procedure
- Decoding Solution
- Viterbi Algorithm
- Parameter Estimation
- HMM Applications



Forward Procedure



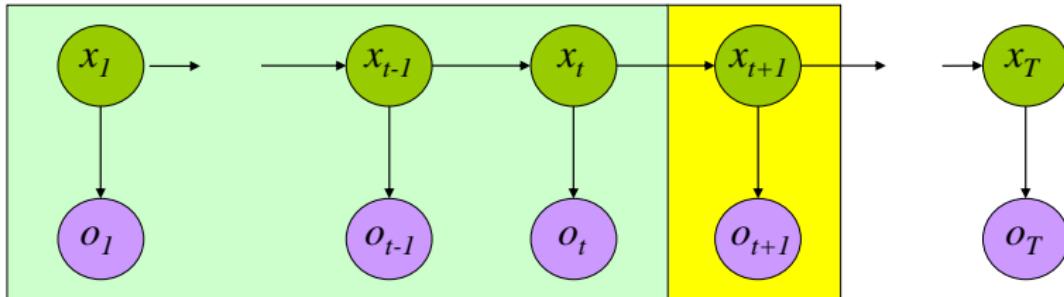
- Special structure gives us an efficient solution using *dynamic programming*.
- **Intuition:** Probability of the first t observations is the same for all possible $t+1$ length state sequences.

• **Define:** $\alpha_i(t) = P(o_1 \dots o_t, x_t = i \mid \mu)$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



$$\alpha_j(t+1)$$

$$= P(o_1 \dots o_{t+1}, x_{t+1} = j)$$

$$= P(o_1 \dots o_{t+1} | x_{t+1} = j) P(x_{t+1} = j)$$

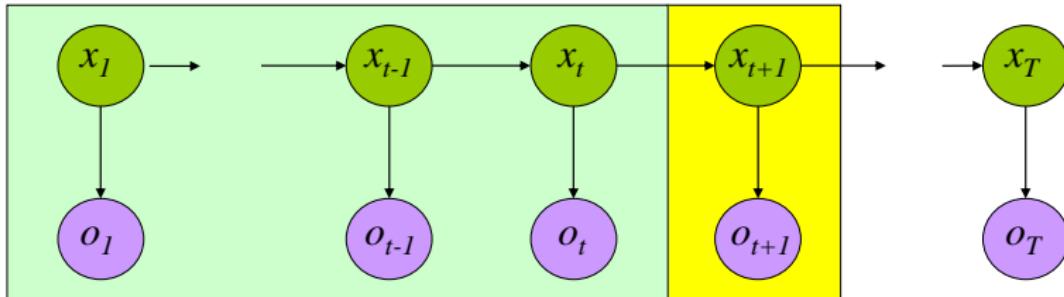
$$= P(o_1 \dots o_t | x_{t+1} = j) P(o_{t+1} | x_{t+1} = j) P(x_{t+1} = j)$$

$$= P(o_1 \dots o_t, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j)$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



$$\alpha_j(t+1)$$

$$= P(o_1 \dots o_{t+1}, x_{t+1} = j)$$

$$= P(o_1 \dots o_{t+1} | x_{t+1} = j)P(x_{t+1} = j)$$

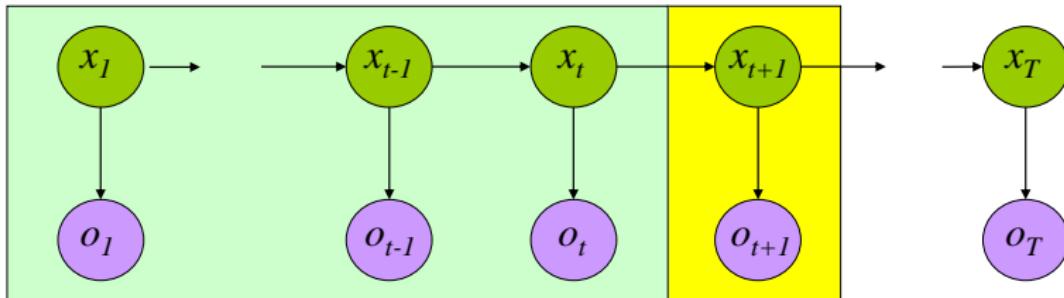
$$= P(o_1 \dots o_t | x_{t+1} = j)P(o_{t+1} | x_{t+1} = j)P(x_{t+1} = j)$$

$$= P(o_1 \dots o_t, x_{t+1} = j)P(o_{t+1} | x_{t+1} = j)$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



$$\alpha_j(t+1)$$

$$= P(o_1 \dots o_{t+1}, x_{t+1} = j)$$

$$= P(o_1 \dots o_{t+1} | x_{t+1} = j)P(x_{t+1} = j)$$

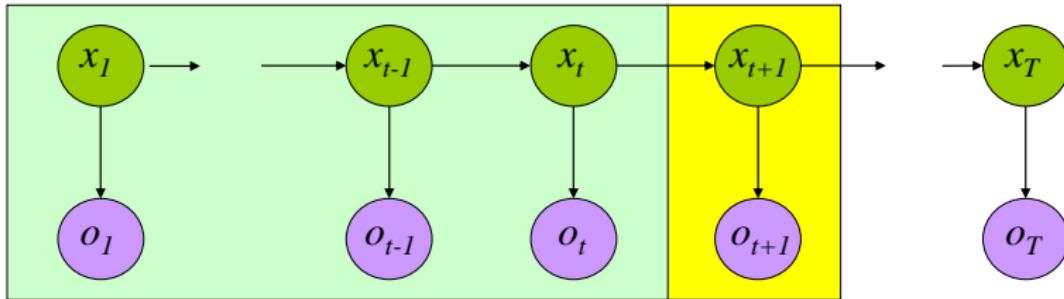
$$= P(o_1 \dots o_t | x_{t+1} = j)P(o_{t+1} | x_{t+1} = j)P(x_{t+1} = j)$$

$$= P(o_1 \dots o_t, x_{t+1} = j)P(o_{t+1} | x_{t+1} = j)$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



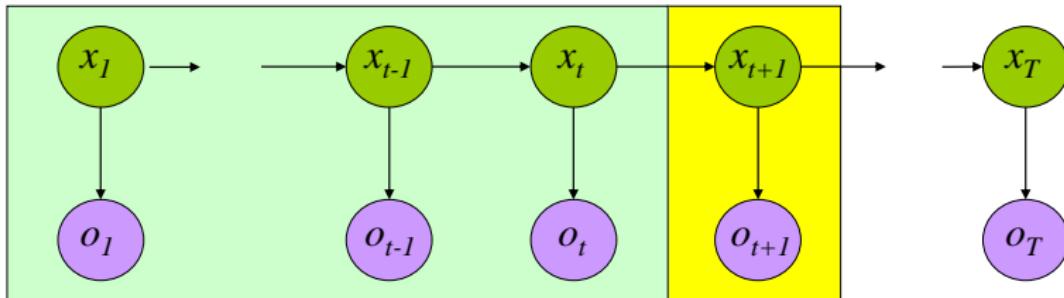
$$\alpha_j(t+1)$$

$$\begin{aligned}&= P(o_1 \dots o_{t+1}, x_{t+1} = j) \\&= P(o_1 \dots o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j) \\&= P(o_1 \dots o_t \mid x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)P(x_{t+1} = j) \\&= P(o_1 \dots o_t, x_{t+1} = j)P(o_{t+1} \mid x_{t+1} = j)\end{aligned}$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j)$$

$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_{t+1} = j | x_t = i) P(x_t = i) P(o_{t+1} | x_{t+1} = j)$$

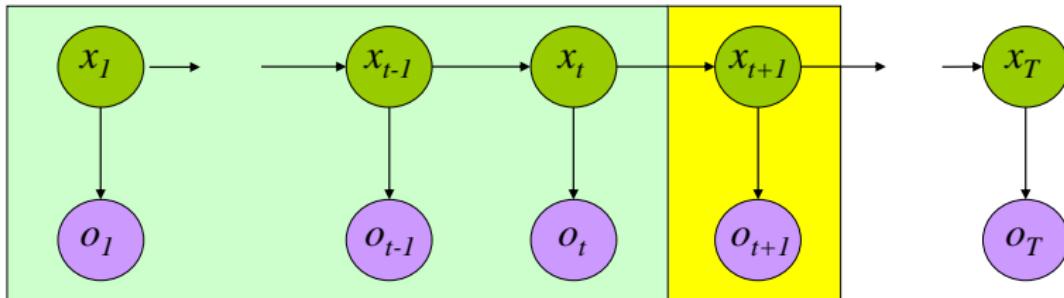
$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i) P(x_{t+1} = j | x_t = i) P(o_{t+1} | x_{t+1} = j)$$

$$= \sum_{i=1 \dots N} \alpha_i(t) a_{ij} b_{j o_{t+1}}$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j)$$

$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_{t+1} = j | x_t = i) P(x_t = i) P(o_{t+1} | x_{t+1} = j)$$

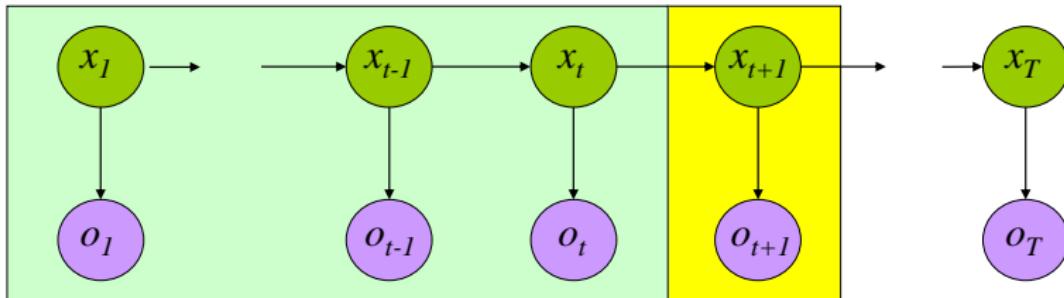
$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i) P(x_{t+1} = j | x_t = i) P(o_{t+1} | x_{t+1} = j)$$

$$= \sum_{i=1 \dots N} \alpha_i(t) a_{ij} b_{j o_{t+1}}$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)



Forward Procedure



$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j)$$

$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_{t+1} = j | x_t = i) P(x_t = i) P(o_{t+1} | x_{t+1} = j)$$

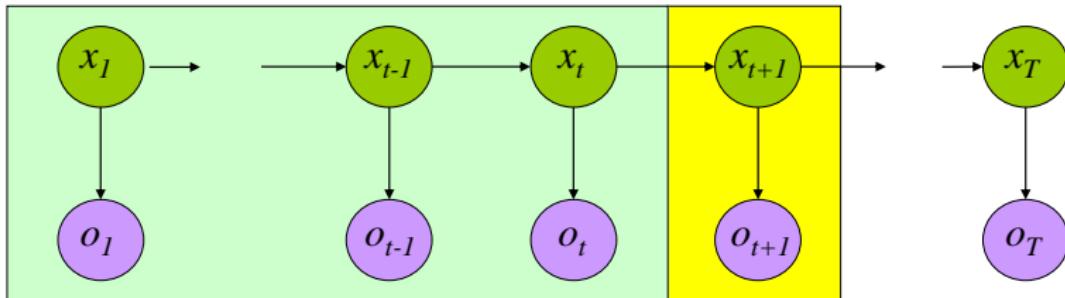
$$= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i) P(x_{t+1} = j | x_t = i) P(o_{t+1} | x_{t+1} = j)$$

$$= \sum_{i=1 \dots N} \alpha_i(t) a_{ij} b_{j o_{t+1}}$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)

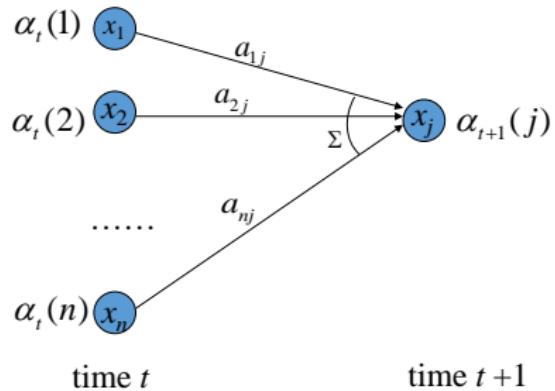


Forward Procedure



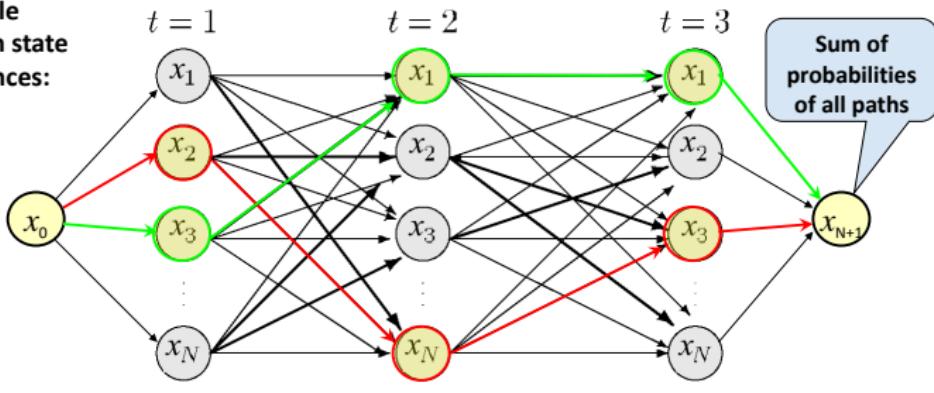
$$\begin{aligned}
 &= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i, x_{t+1} = j) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_{t+1} = j | x_t = i) P(x_t = i) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} P(o_1 \dots o_t, x_t = i) P(x_{t+1} = j | x_t = i) P(o_{t+1} | x_{t+1} = j) \\
 &= \sum_{i=1 \dots N} \alpha_i(t) a_{ij} b_{j o_{t+1}}
 \end{aligned}$$

David Meir Blei, Hidden Markov Models, 1999 (Slides)

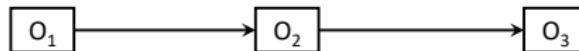




possible hidden state sequences:



observation sequence:



of possible hidden state sequences: N^T

Time complexity of Forward algorithm: $O(N^2T)$



Forward Algorithm Summary

1.

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

Multiplications

= N

2. for $t = 1, 2, \dots, T-1, 1 \leq j \leq N$

$$\alpha_{t+1}(j) = [\sum_{i=1 \text{ to } N} \alpha_t(i) * a_{ij}] * b_j(O_{t+1})$$

= (N+1)N(T-1)

3. Finally we have:

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

= 0

Total:

N+(N+1)N(T-1)

Anantharaman Narayana Iyer, Natural Language Processing, Unit 2 – Tagging Problems and HMM (Slides)



Content

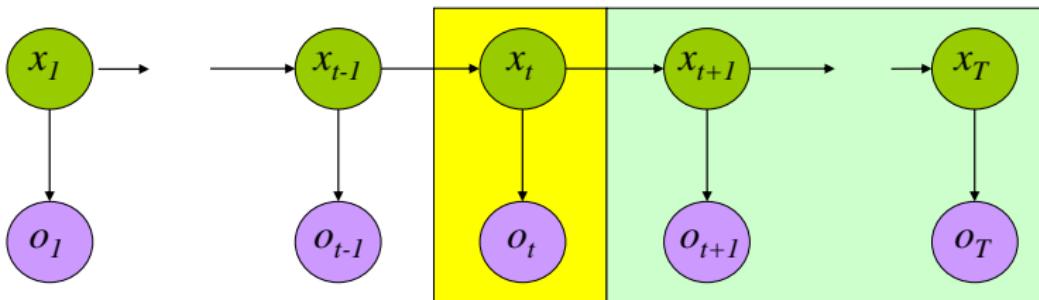
3

Hidden Markov models (HMMs)

- Model definition
- Inference in an HMM
- Decoding
- Forward Procedure
- **Backward Procedure**
- Decoding Solution
- Viterbi Algorithm
- Parameter Estimation
- HMM Applications



Backward Procedure



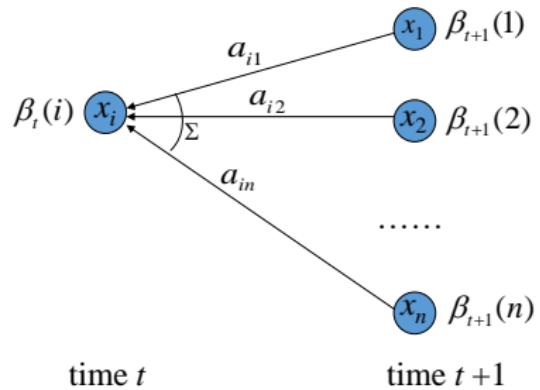
$$\beta_i(T + 1) = 1$$

$$\beta_i(t) = P(o_t \dots o_T \mid x_t = i)$$

$$\beta_i(t) = \sum_{j=1 \dots N} a_{ij} b_{io_t} \beta_j(t + 1)$$

Probability of the rest
of the states given the
first state

David Meir Blei, Hidden Markov Models, 1999 (Slides)





Backward Algorithm: Summary

1.

$$\beta_T(i) = 1, 1 \leq i \leq N$$

$= 0$

Multiplications

2. for $t = T-1, T-2, \dots, 1, 1 \leq i \leq N$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} * b_j(O_{t+1}) * \beta_{t+1}(j)$$

$= (2N)N(T-1)$

3. Finally we have:

$$P(O|\lambda) = \sum_{i=1}^N \pi_i * b_i(O_1) * \beta_1(i)$$

$= 2N$

Total:

$$2N + (2N)N(T-1)$$

Anantharaman Narayana Iyer, Natural Language Processing, Unit 2 – Tagging Problems and HMM (Slides)