

# Natural Language Processing

## CSE 325/425



Sihong Xie

### Lecture 7:

- Forward and Backward algorithm

Recall : Bigram  $\Pr(O) | \Pr(o_{t-1} \rightarrow o_t)$  ①

In the future we will talk about PCFG for  $\Pr(O) | \text{PCFG}$  ②

## Likelihood of a sentence

Given: observed sentence  $O = [o_1, \dots, o_T], o_t \in V$   
a fixed HMM  $(A, B, \pi)$

Find:  $\Pr(O|A, B, \pi)$   $\triangleq$  Likelihood

Useful for evaluating the validity of a sentence.

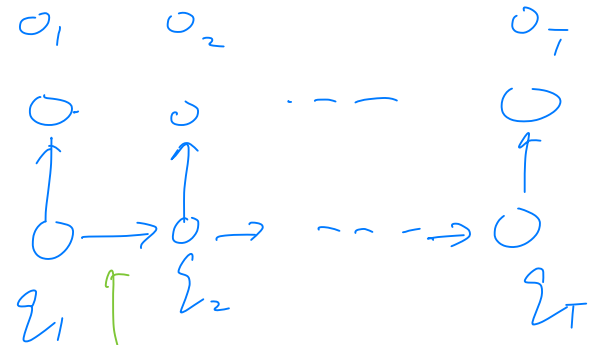
By total probability,  $\Pr(O|A, B, \pi) = \sum_Q \Pr(O, Q|A, B, \pi)$

The summation is over all possible hidden state sequences  $Q = [q_1, \dots, q_T], q_t \in S$

A brute-force enumeration will have time complexity

$o_t$  is dependent of  $\xi_t$

$\xi_t$  is dependent of  $\xi_{t-1}$



$$A = N \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} A_{ij} = \Pr(\xi_t = j | \xi_{t-1} = i)$$

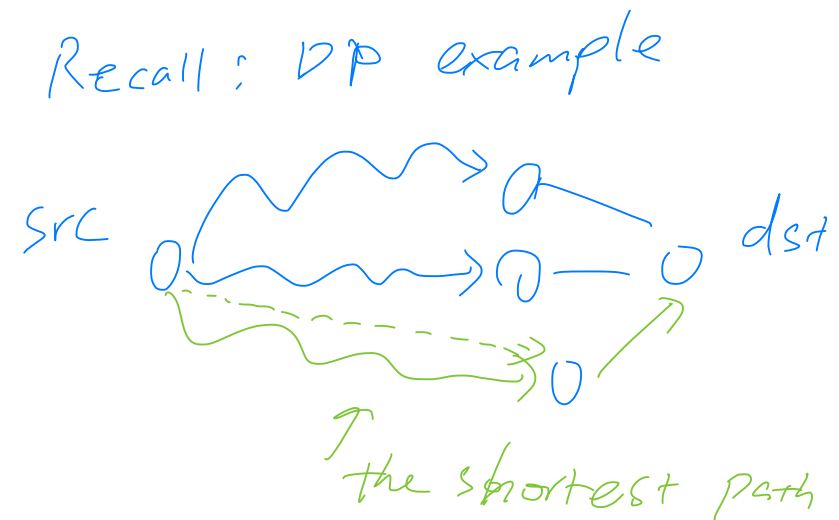
$$B = N \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

$$\# \text{ of } Q = N^T$$

# Likelihood of a sentence

Forward algorithm:  $\Rightarrow P(c|A, B, \pi)$

- It is an instance of dynamic programming algorithm:
  - solves a large problem using cached solutions to sub-problems;
  - overlapping sub-problems: can re-use solutions to sub-problems;
  - optimal structures: optimal sub-problem solutions  $\Rightarrow$  optimal larger problem solutions.
- It is a specific example of belief propagation or message passing algorithm for probabilistic graphical models (PGM)
  - more general PGM includes Bayesian networks and Markov random fields;
  - here we focus on a linear chain (sequence);
  - later on will discuss trees in syntactic parsing;
  - see CSE326 note for more details;



# Likelihood of a sentence

Two simple identities:

- Order of product and summation can be exchanged

Sum-product

$$ax + ay + bx + by = (a + b)(x + y)$$

→ forward / backward

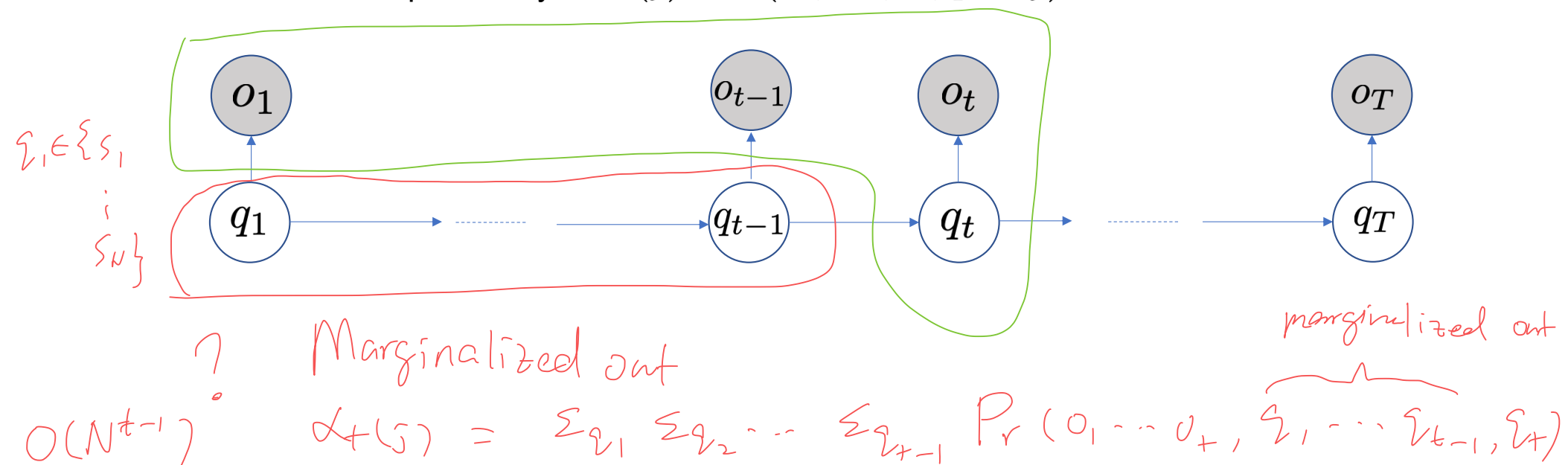
product-sum

- Order of product and max can be exchanged

— Viterbi;

$$\max(ax, ay, bx, by) = \max(a, b) \times \max(x, y)$$

Define the *forward* probability  $\alpha_t(j) = P(o_1, \dots, o_t, q_t = j)$



# Likelihood of a sentence

Forward algorithm: a gentle start. Compute the following

Base  
case  
of the  
recursion

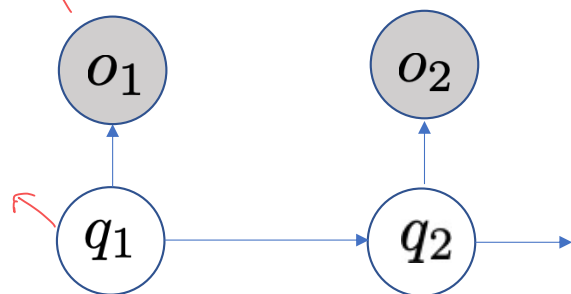
$$\alpha_1(i) = \Pr(o_1, z_1 = i)$$

$$= \Pr(z_1 = i) \Pr(o_1 | z_1 = i)$$

$$= \pi(i) \times b_{i o_1} \quad o_1 \in \{1, \dots, V\}$$

word index in  $V$

$i$ -th row  $o_1$ -th col



$\Pr(o_1, z_1 = i)$

$i = 1, \dots, N$

$N = \# \text{ of pos-tags}$

$$\begin{bmatrix} \alpha_1(1) \\ \alpha_1(2) \\ \vdots \\ \alpha_1(N) \end{bmatrix} \begin{matrix} \xrightarrow{\text{green lines}} \end{matrix} \begin{bmatrix} \alpha_2(1) \\ \alpha_2(2) \\ \vdots \\ \alpha_2(N) \end{bmatrix}$$

$$\alpha_2(i) \triangleq \Pr(o_1, o_2, z_2 = i)$$

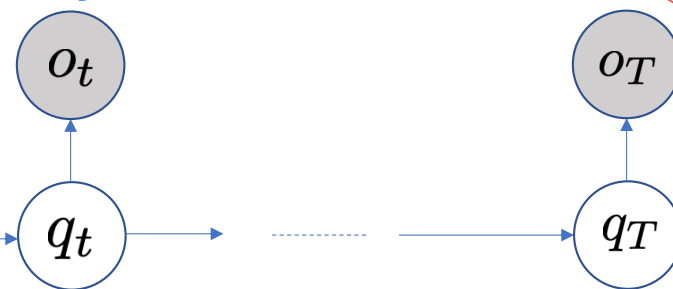
R.V. marginalized out

$$= \sum_j \Pr(o_1, o_2, z_1 = j, z_2 = i)$$

$$= \sum_j \Pr(o_1, z_1 = j) \times \Pr(z_2 = i | z_1 = j) \times \Pr(o_2 | z_2 = i)$$

$$= \sum_j \alpha_1(j) \times a_{ji} \times b_{i o_2}$$

$o_2 \in \{1, \dots, V\}$



$$\begin{bmatrix} \alpha_t(1) \\ \vdots \\ \alpha_t(N) \end{bmatrix}$$

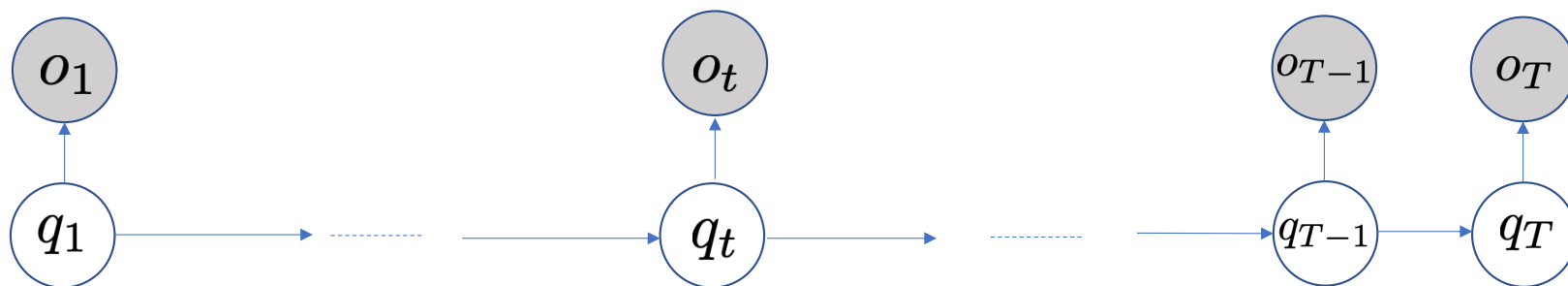
General Recursion.

$$\alpha_t(i) = \sum_j \alpha_{t-1}(j) \times a_{ji} \times b_{i o_t}$$

# Likelihood of a sentence

Forward algorithm: compute the likelihood of an entire sequence

$$\begin{aligned}
 \Pr(O) &= \sum_j P(O, q_T = j) = \sum_j \alpha_T(j) \\
 &= \sum_j \sum_i \alpha_{T-1}(i) \times a_{ij} \times b_{jO_{T-1}} \\
 &= \sum_j \sum_i \left[ \sum_k \alpha_{T-2}(k) \times a_{ki} \times b_{iO_{T-2}} \right] \times a_{ij} \times b_{jO_{T-1}} \\
 &\vdots
 \end{aligned}$$



Sum-product:

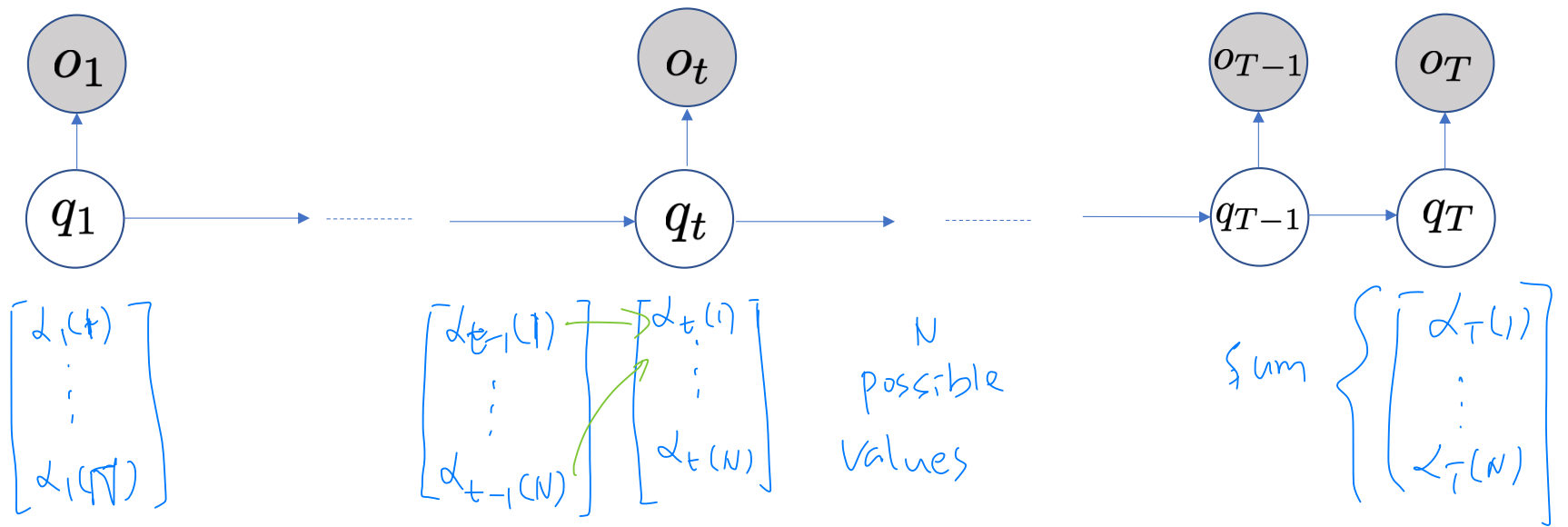
$O(N^T)$

$$= \sum_{\tilde{v}_T} \sum_{\tilde{v}_{T-1}} \sum_{\tilde{v}_{T-2}} \dots \sum_{\tilde{v}_1} \left[ \pi(\tilde{v}_1) b_{\tilde{v}_1 o_1} \times a_{\tilde{v}_1 \tilde{v}_2} \times b_{\tilde{v}_2 o_2} \times a_{\tilde{v}_2 \tilde{v}_3} \times \dots \times a_{\tilde{v}_{T-1} \tilde{v}_T} b_{\tilde{v}_T o_T} \right] \begin{bmatrix} \alpha_T(1) \\ \vdots \\ \alpha_T(N) \end{bmatrix}$$

# Likelihood of a sentence

Forward algorithm: compute the likelihood of an entire sequence

- $O(N)$  1. Initialize  $\alpha_1(i)$  for each value  $i$  of the first hidden state  $q_1$ .  $\mathcal{L}_1(i) = \pi(i) b_{i|o_1}$
- $O(T)$  2. for  $t = 2, \dots, T$  # go through word locations
- $O(N)$  for  $j = 1, \dots, N$  # go through pos-tags of the  $t$ -th location
- $O(NT)$  compute  $\alpha_t(j) = \sum_k \alpha_{t-1}(k) a_{kj} b_j(o_t)$  ← B matrix
- A matrix
- previous  $\alpha$  array
- } All are given or calculated.
- $O(N)$  3. Return  $\Pr(O) = \sum_j P(O, q_T = j) = \sum_j \alpha_T(j)$



# Likelihood of a sentence

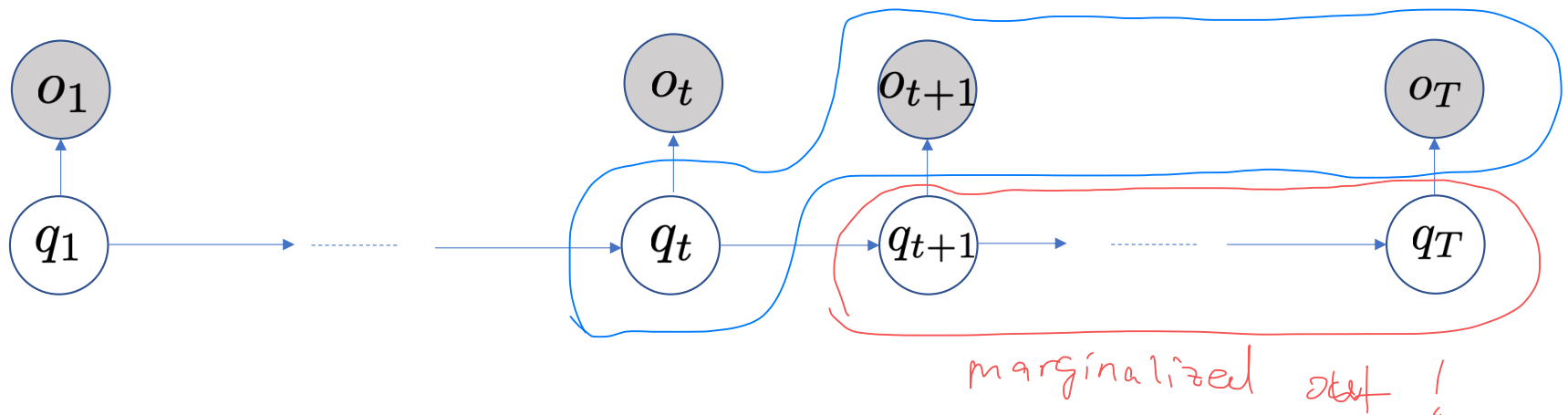
Given: observed sentence  $O = [o_1, \dots, o_T], o_t \in V$

a fixed HMM  $(A, B, \pi)$

Find:  $\Pr(O|A, B, \pi)$

An alternative way to evaluate  $\Pr(O|A, B, \pi) = \sum_Q \Pr(O, Q|A, B, \pi)$

Define the *backward* probability  $\beta_t(j) = P(o_{t+1}, \dots, o_T | q_t = j)$   $j \in \{1, \dots, N\}$





# Likelihood of a sentence

Backward algorithm:

Base case:

$$\beta_T(i) = \Pr(o_{T+1} \dots o_T | \hat{v}_T = i)$$

$$= \Pr(\phi | \hat{v}_T = i)$$

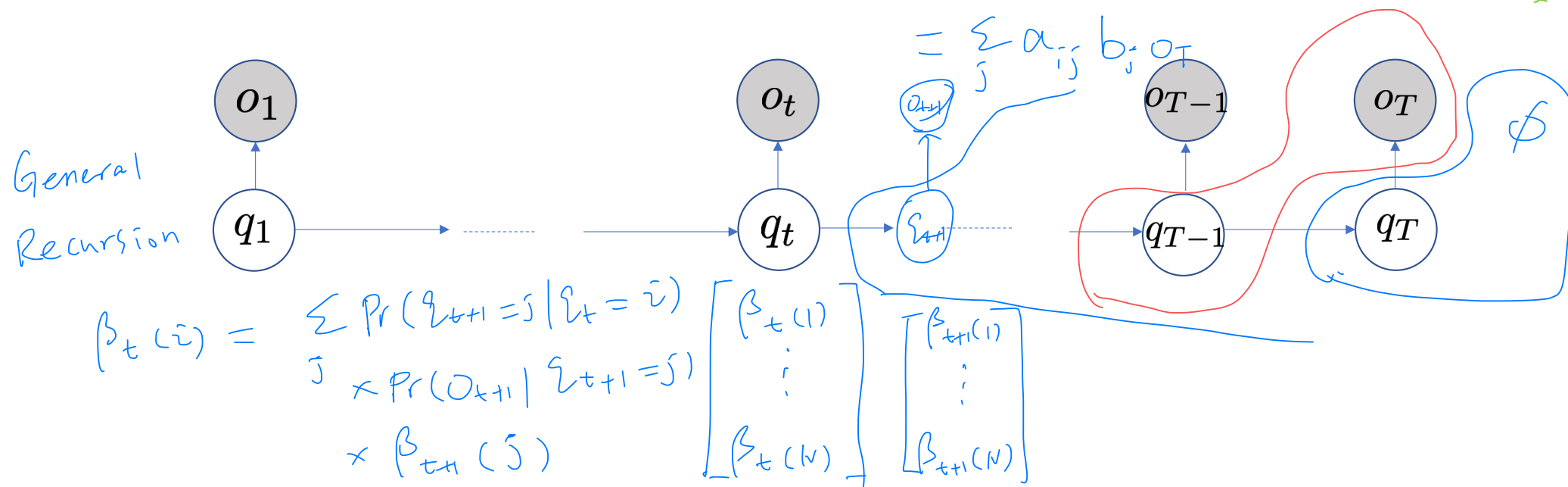
$$= 1$$

$$\beta_{T-1}(i) \triangleq \Pr(o_T | \hat{v}_{T-1} = i)$$

$$= \sum_j \Pr(\hat{v}_T = j | \hat{v}_{T-1} = i)$$

$$\times \Pr(o_T | \hat{v}_T = j)$$

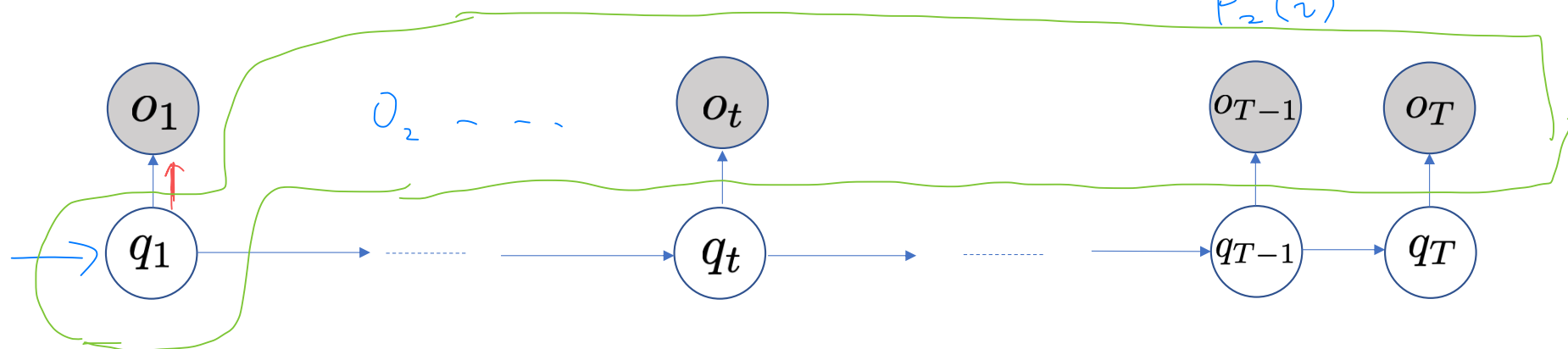
$$\times \Pr(o_{T+1} \dots o_T | \hat{v}_T = j)$$



# Likelihood of a sentence

Backward algorithm: compute the likelihood of an entire sequence

$$\begin{aligned}
 \Pr(O) &= \sum_j P(O, q_1 = j) = \sum_j \pi(j) \boxed{\Pr(o_1 | q_1 = j)} \boxed{\Pr(o_2 \dots o_T | q_1 = j)} \\
 &= \sum_j \pi(j) b_{jo_1} \boxed{\beta_1(j)} \\
 &= \sum_j \pi(j) b_{jo_1} \left[ \sum_{\bar{v}} \Pr(q_2 = \bar{v} | q_1 = j) \Pr(o_2 | q_2 = \bar{v}) \right] \beta_2(\bar{v})
 \end{aligned}$$



$$\begin{aligned}
 O(N^T) &= \sum_{\bar{v}_1} \sum_{\bar{v}_2} \dots \sum_{\bar{v}_T} \pi(\bar{v}_1) b_{\bar{v}_1 o_1} a_{\bar{v}_1 \bar{v}_2} b_{\bar{v}_2 o_2} \dots a_{\bar{v}_{T-1} \bar{v}_T} b_{\bar{v}_T o_T}
 \end{aligned}$$

Given

$$\pi, [ ] \quad A [ ] \quad B [ ]$$

# Likelihood of a sentence

Backward algorithm: compute the likelihood of an entire sequence

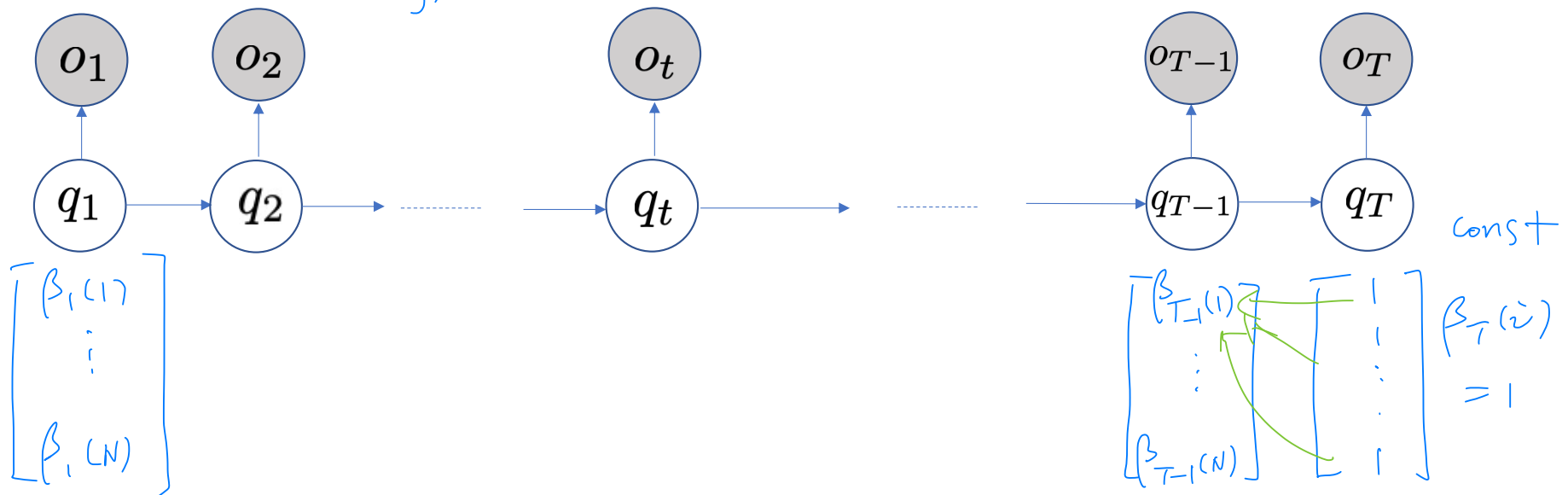
1. Initialize  $\beta_T(i)$  for each value  $i$  of the last hidden state  $q_T$ .
2. for  $t = T - 1, \dots, 1$   
for  $j = 1, \dots, N$   
compute  $\beta_t(j) = \sum_k \beta_{t+1}(k) a_{jk} b_j(o_t)$

MLE Alg  
to estimate  $\pi, A, B$

$$Pr(\text{corpus} | A, B, \pi)$$

$j$ -th row &  $0_t$ -th column.  
B matrix

3. Return  $Pr(O) = \sum_j \underbrace{P(q_1 = j)}_{\pi(j)} b_j(o_1) \beta_1(j)$



# Summary

Given a sentence  $O = [o_1, \dots, o_T], o_t \in V$

and HMM parameters  $(A, B, \pi)$

Compute  $\Pr(O|A, B, \pi)$

Forward algorithm

- start from the first POS tag and go forward.
- end with the last POS tag and find sentence likelihood.

Backward algorithm

- Start from the last POS tag and go backward.
- end with the first POS tag and find sentence likelihood.