

# Natural Language Processing

## CSE 325/425



Sihong Xie

### Lecture 8:

- HMM prediction and training

# Predict optimal hidden states

Want to find the best sequence of POS tags for a given sentence.

- Vocabulary  $V$
- Set of N POS tags  $S = \{s_1, \dots, s_n\}$

- Observed sentence  $O = [o_1, \dots, o_T], o_t \in V$

All we gotta do is go around the corner

- Hidden states:  $Q = [q_1, \dots, q_T], q_t \in S$

DT PRP VBN VB VBZ VB IN DT NN

- MAP (maximum a posterior) prediction

$$Q^* = \arg \max_Q \Pr(Q|O) \Leftrightarrow \arg \max_Q \Pr(O|Q) \Pr(Q)$$

- Using the Bayes theorem:

$$\Pr(Q|O) = \frac{\Pr(O|Q) \Pr(Q)}{\Pr(O)}$$

← Constant w.r.t.  $Q$

- Difficulty: there are exponentially many possible sequences.

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



# Predict optimal hidden states

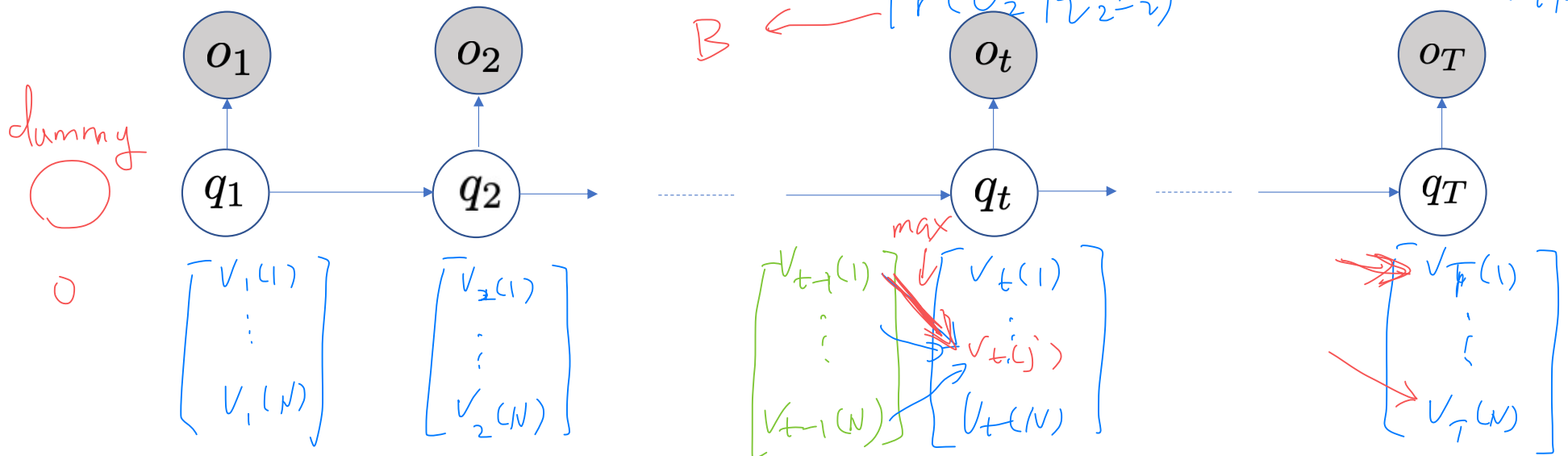
Viterbi Algorithm: Compute the base case, the next one, and the general case

$\xi_0$  is dummy

$$v_1(i) = \max_{\xi_0} \Pr(\xi_0, \xi_1=i, o_1) = \Pr(\xi_1=i) \Pr(o_1|\xi_1=i)$$

$$v_2(i) = \max_{\xi_1} \Pr(\xi_1, \xi_2=i, o_1, o_2) = \max_{\xi_1} \Pr(\xi_1=j) \Pr(o_1|\xi_1=j) \Pr(\xi_2=i|\xi_1=j) \Pr(o_2|\xi_2=i)$$

$$v_t(i) = \max_{\xi_1 \dots \xi_{t-1}} \Pr(\xi_1 \dots \xi_{t-1}, \xi_t=i, o_1 \dots o_t) = \max_{\xi_{t-1}} \max_{\xi_1 \dots \xi_{t-2}} \Pr(\xi_1 \dots \xi_{t-2}, \xi_{t-1}=j) \times \Pr(\xi_t=i|\xi_{t-1}=j) \times \Pr(o_t|\xi_t=i)$$



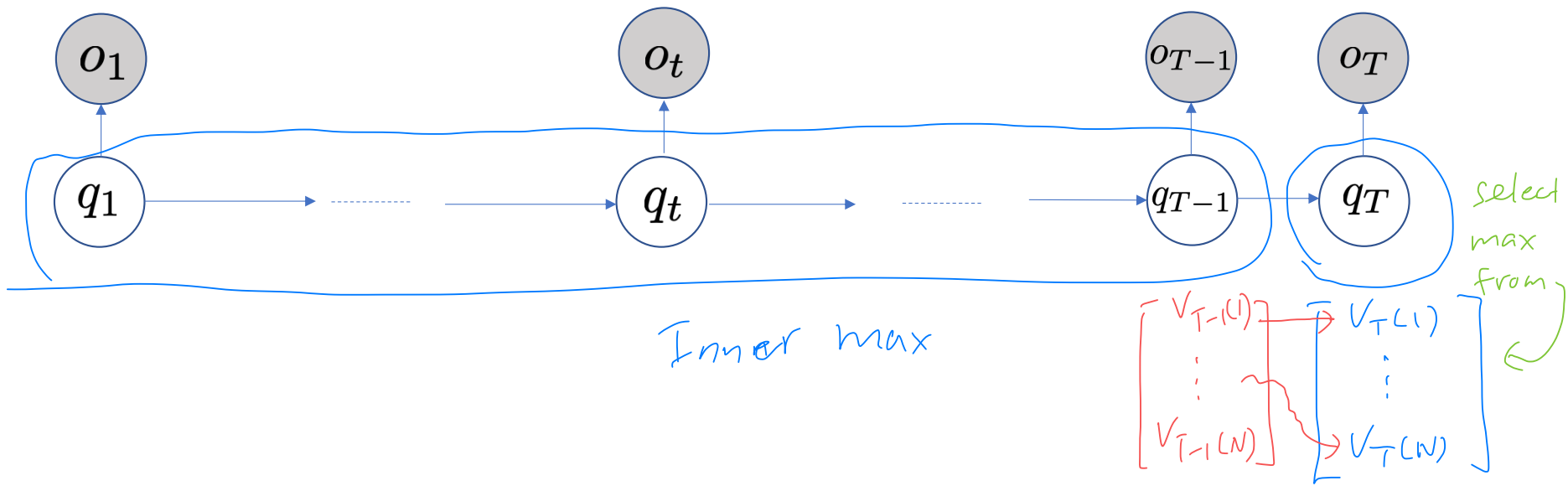
$$= \max_{v_T(\hat{v})} \left[ \max_{v_{T-1}(\hat{j})} \left[ \max_{v_1 \dots v_{T-2}} \Pr(v_1 \dots v_{T-2}, o_1 \dots o_{T-2}) \Pr(v_{T-1} | v_{T-2}) \Pr(o_{T-1} | v_{T-1} = \hat{j}) \right] \times \Pr(v_T | v_{T-1}) \Pr(o_T | v_T = \hat{v}) \right]$$

# Predict optimal hidden states

Viterbi algorithm: compute maximum probability

$$\Pr(Q^* | O) = \max_Q \Pr(Q | O) = \max_Q \Pr(v_1 \dots v_{T-1}, v_T, o_1 \dots o_{T-1}, o_T)$$

$$= \max_{v_T} \max_{v_1 \dots v_{T-1}} \left[ \Pr(v_1 \dots v_{T-1}, o_1 \dots o_{T-1}) \times \Pr(v_T = \hat{v} | v_{T-1} = \hat{j}) \times \Pr(o_T | v_T = \hat{v}) \right]$$



$$v_1(i) = \pi(i) \times b_{i|o_1}$$

# Predict optimal hidden states

Viterbi algorithm: compute maximum probability and the optimal sequence (decoding)

1. Initialize  $v_1(i)$  for each value  $i$  of the first hidden state  $q_1$ .

2. for  $t = 2, \dots, T$

for  $j = 1, \dots, N$

compute  $v_t(j) = \max_k v_{t-1}(k) a_{kj} b_j(o_t)$

record back-pointers  $p_t(j) = \arg \max_k v_{t-1}(k) a_{kj} b_j(o_t) = k^*$

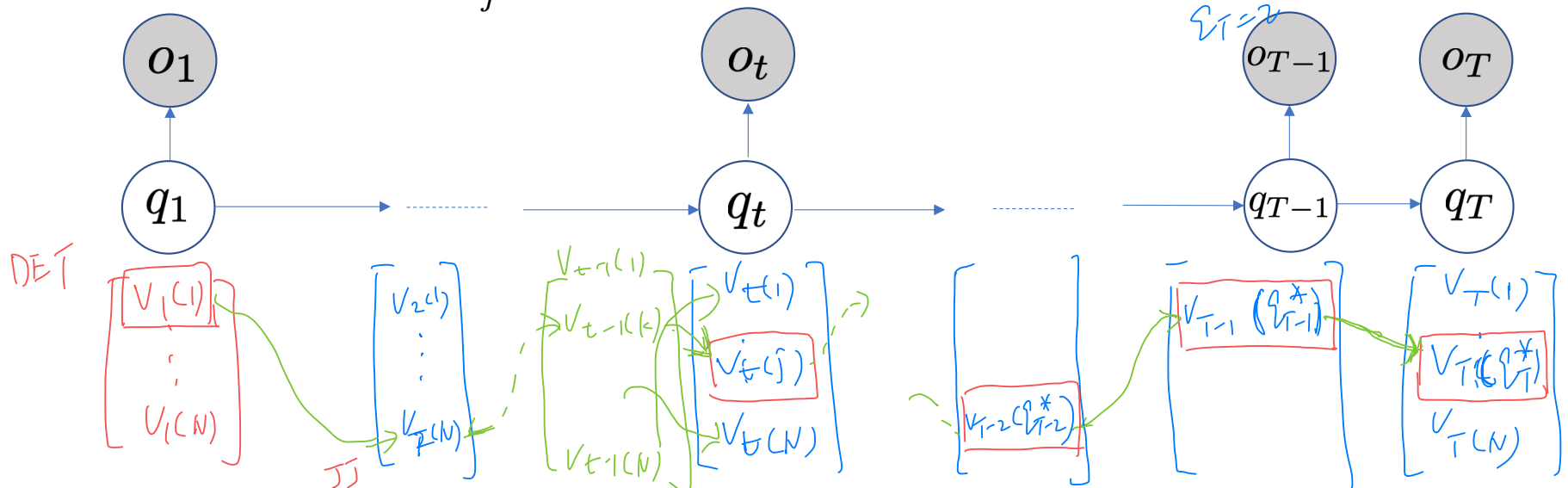
3. Backtracking to find  $Q^*$

4. Return  $\Pr(Q^*|O) = \max_j v_T(j)$

$$a_{kj} = \Pr(q_t = j | q_{t-1} = k)$$

$$\Pr(Q^*|O) = \Pr(Q^*|o)$$

$$q_T^* = \arg \max_{\tilde{j}} v_T(\tilde{j})$$



# Predict optimal hidden states

A running example of Viterbi.

- Suppose you are in a casino with a cheating dealer X.
- X flips a fair coin if X decides not to cheat and a biased coin otherwise.
- X follows a Markov chain to change between cheating and no cheating.
- You only observe a sequence of flips.
- Guess which flips are from a biased coin?

*If cheating*

$$P_Y(H) = 0.9$$

