# Natural Language Processing
## CSE 325/425



Sihong Xie

Lecture 6:
- Part-of-Speech (POS)
- POS tagging
- Hidden Markov Models (HMM)

# Part-of-Speech

English word classes: cover a few common classes that will be used in tagging.

1. Nouns: pronouns (she, he, I, who, others), proper nouns (Russia), countable nouns (desk), mass noun (air)

2. Verbs: participles (paced), gerund (pacing), auxiliaries (be, do, have, can, may, should)

3. Adjectives: comparative, superlative.

4. Adverbs: I went *Church yesterday*

5. Prepositions: in, on, over, ...

6. Particles: phrasal verb like "go *over*".
   - Easy to be confused with prepositions.
   - Combination of verb and particle does not have their meanings combined simply.

7. Determiners: a, the, an

8. Conjunctions: and, but, that, when

9. Other smaller classes.

# Part-of-Speech

Syntatic information: how words are ordered in a sentence.

- noun-verb

- determiner-noun

- adjective-noun

- verb-adverb

- preposition-noun

Useful for grammar checking: go to (a?the?) hospital

Semantic information: meaning of a word in a context. Useful for:

- machine translation (building a building): building -> (建 vs. 楼)

- question-answering, need to understand the semantics of what a person is asking.

- relation extraction (Bill Gates founded MS): gates (verb vs. noun)

- event extraction (They went to a concert): concert (verb vs. noun)

- entity extraction (I will visit DC): DC?

Syntactic info ⇐⇒ Semantic info

# POS tagging

The Penn Tree Bank tagsets.

- JJ (adjective), JJR (comparative), JJS (superlative)
- NN, NNS (plural), NNP (pronoun), NNPS (pronoun plural)
- RB (adverb), RBR (comparative), RBS (superlative)
- VB, VBD (past tense), VBG (present participle, gerund), *going*
  VBN (past participle), VBZ (third-person single) *goes*
- RP: particle (as in *go on*)
- IN: preposition (as in *go to church*)

*go*
*gone*

No free

→ large-scale training corpus
for Language models,
POS tagger

Syntactic Parsing

See SLP3 for a more
comprehensive List & examples.

Example: → Determiner
- The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.
- There/EX are/VBP 70/CD children/NNS there/RB

Def of Pos tagging task:

Given a sentence = $[w_1, w_2, \cdots, w_n]$

output a single Pos tag for each word in $[w_1 \cdots w_n]$

*corpus with sentences POS-tagged* → *NLP model (POS-tagger)* → *make predictions on future sentences.*

# POS tagging

## Why POS-tagging is difficult

*VBZ*

- *Does __that__ flight serve dinner (that is used as a determiner)*
  - NN   VB   NN

- *I thought __that__ your flight was earlier (that is used as a complementizer)*
  - *clause*

*VB get   RP around*

main sentence   $q_1$   $q_2$   $q_3$

## Difficult even for humans

- Mrs./NNP Shaefer/NNP never/RB got/VBD around/RP to/TO joining/VBG
- All/DT we/PRP gotta/VBN do/VB is/VBZ go/VB around/IN the/DT corner/NN
- Chateau/NNP Petrus/NNP costs/VBZ around/RB 250/CD

*Adverb*

## Naïve solution?

- always predict the most frequent tag for each word.
- contexts are useful

*good performance (90%) on easy words,*

Most **common** words are ambiguous; most words are not ambiguous.   *Such as "difficult"/JJ*

# Hidden Markov Model ( HMM )

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

- Vocabulary $V$

- Set of N POS tags $S = \{s_1, \ldots, s_N\}$

$T$ = length of the sentence

$T = 9$

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

| All | we | gotta do | is | go around the corner |
|-----|-----|----------|-----|----------------------|

| DT | PRP | VBN | VB | VBZ | VB | IN | DT | NN |
|----|-----|-----|----|-----|----|----|----|----|

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

- MAP (maximum a posterior) prediction

$$Q^* = \arg \max_Q \boxed{\Pr(Q|O)}$$

posterior probability of seq Q given seq O.

- Using the Bayes theorem:

$Pr(0) =$ Prob. of seeing the sentence $O$

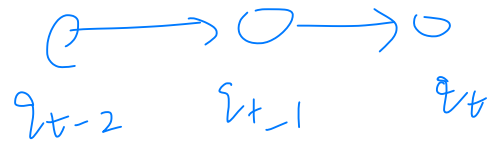$$\Pr(Q|O) = \frac{\Pr(O|Q)\Pr(Q)}{\Pr(O)}$$

$\sum_Q P(0,Q) =$

$Pr(O|Q)$ : likelihood of observing seq $O$, conditioned on seq $Q$.

- Difficulty: there are exponentially many possible sequences.

How many possible Q sequences ?

$Pr(Q)$ : prob. of observing seq $Q$

$O(N^T)$ : exponential time complex.

# The Markov assumption

Motivations:

• Linguistically, any two POS tags do not occur independently.

   o   P(NN|DT) and P(JJ | DT) should be higher than P(DT | JJ)

   o   e.g., *a yellow cab* vs. *yellow a cab* vs. *yellow cab a*

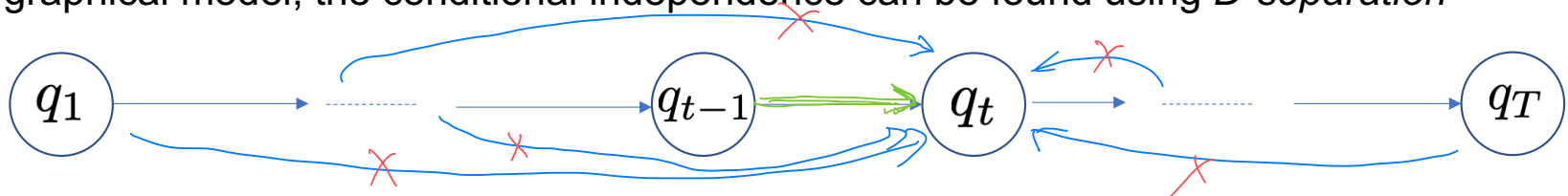• Computationally, reduce the complexity of computing the prior $\Pr(Q)$

Conditioned on the **previous** POS tag, the probability of the **current** tag is

• Transition probability

$$\Pr(q_t | q_1, \ldots, q_{t-1}) = \Pr(q_t | q_{t-1})$$

Using a graphical model, the conditional independence can be found using *D-separation*

$q_1$      $q_{t-1}$      $q_t$      $q_T$

# The Markov assumption

$O(N^2) = \begin{matrix} N \\ N \\ \vdots \\ N \end{matrix}$

N pos tags

exponential ?!

fully connected
directed
graph

P(N|N)

P(N|V)

P(D|A)

N = Noun
V = Verb
D = Determiner
A = Adjective

Each row of A is a prob dist.

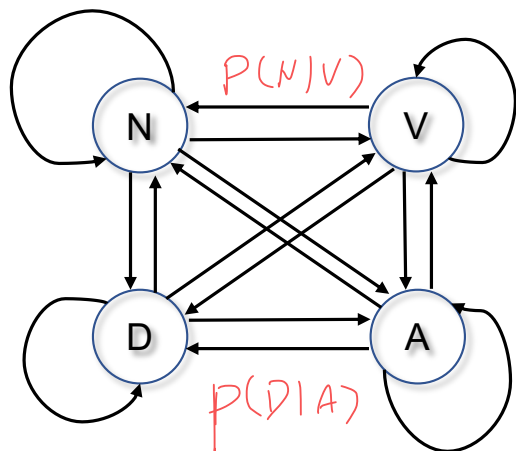Transition probability matrix $A$

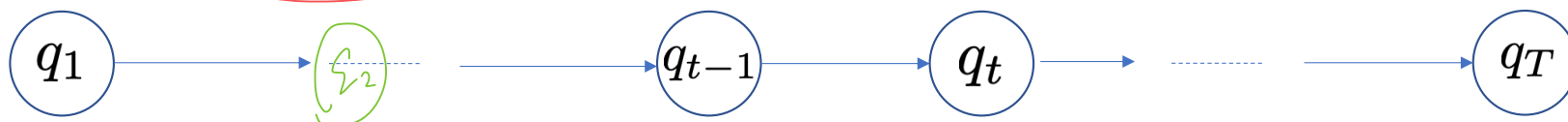$A_{ij}$ = Probability of going from hidden state i to hidden state j,

i-th row

j-th col

$\begin{array}{c} & N \quad V \quad D \quad A \\ \begin{array}{c} N \\ V \\ D \\ A \end{array} & \left[ \begin{array}{cccc} \oslash & \oslash & \oslash & \oslash \end{array} \right] \end{array}$

$q_1 \rightarrow q_2 \cdots \rightarrow q_{t-1} \rightarrow q_t \rightarrow \cdots \rightarrow q_T$

$q_1 \in \{N, V, D, A\}$

$q_t \in \{N, V, D, A\}$

$q_T \in \{N, V, D, A\}$

$P(q_2 | q_1 = V)$

Sample $q_2$ using the probabilities $P(N|V), P(V|V), P(D|V), P(A|V)$
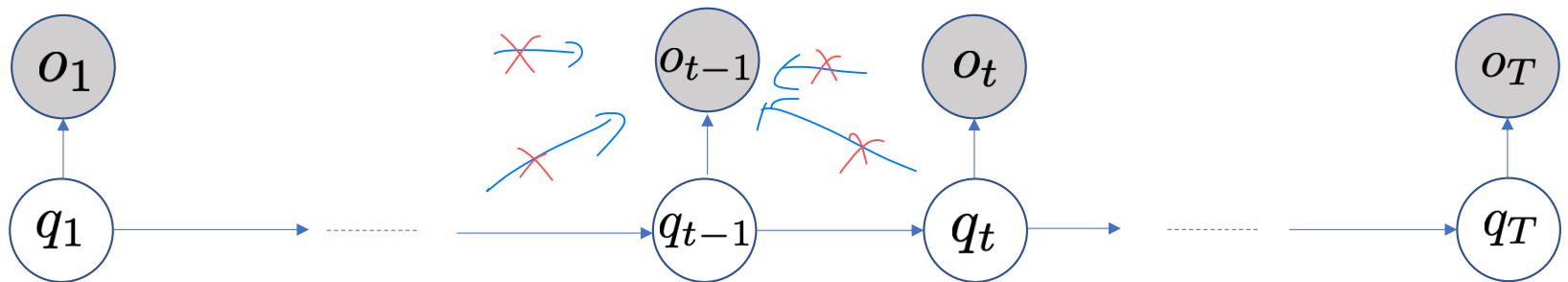
# Emission probability

Motivations:

- Linguistically, a word is selected (partially) based on the current POS tag.

- Computationally, reduce the complexity of computing the likelihood $\Pr(O|Q)$

Conditioned on the **current** POS tag, the probability of the **current** word is

- Emission probability

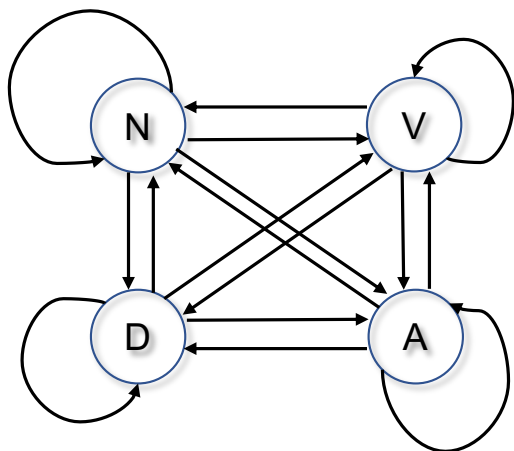$$\Pr(o_t|q_1, \ldots, q_T, o_1, \ldots, o_{t-1}, o_{t+1}, \ldots, o_T) = \Pr(o_t|q_t)$$

Using a graphical model

Each row of $B$ is a prob. dist.

$$B = |S| \left\{ \begin{bmatrix} \, . & . & . & - & - & . \end{bmatrix} \right.$$

$$|V|$$

# Generating a sentence using HMM
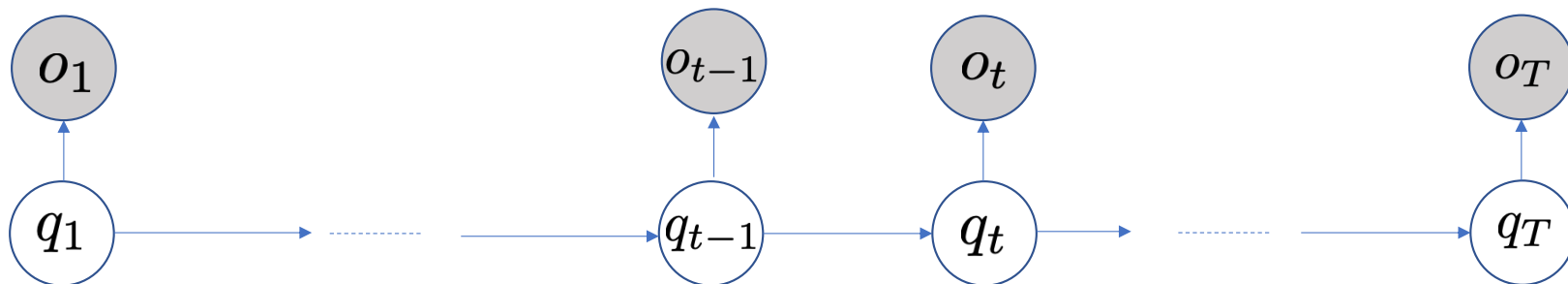


N = Noun
V = Verb
D = Determiner
A = Adjective

Emission probability matrix $B$

$B_{io}$ = emission probability of
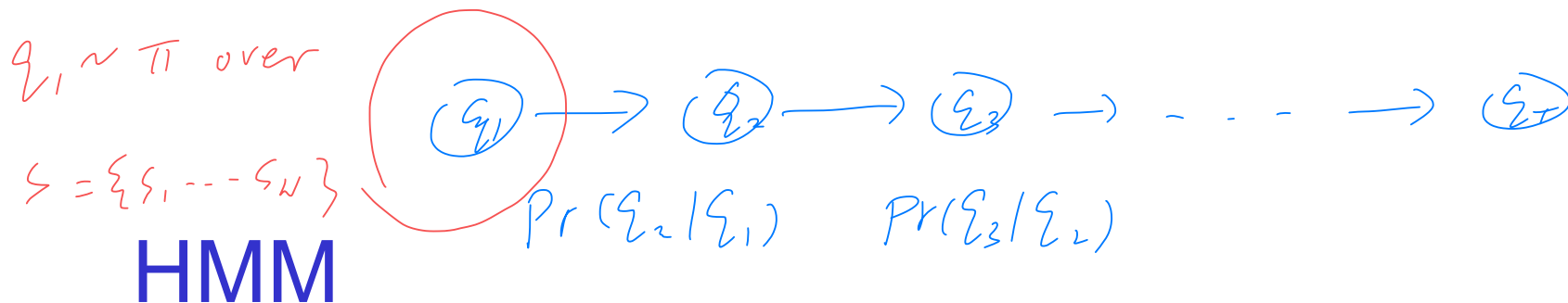
word $o$ from hidden state $i$.

$i$-th row     $o$-th col

$o_t \in V$

$q_t \in \{N, V, D, A\} = S$

$\xi_i \in \{N, V, D, A\}$

$\pi = [Pr(N), Pr(V), Pr(D), Pr(A)]$

High       low       High       higher

$q_1 \sim \pi$ over

$S = \{S_1, \dots S_N\}$

$\boxed{q_1} \rightarrow \boxed{q_2} \rightarrow \boxed{q_3} \rightarrow \cdots \rightarrow \boxed{q_T}$

$Pr(q_2 | q_1) \qquad Pr(q_3 | q_2)$

# HMM

Now we can simplify the following probabilities

- Prior

$= Pr(q_1 | \emptyset)$

$$\Pr(Q) = \boxed{\Pr(q_1)}\Pr(q_2|q_1)\dots\Pr(q_T|q_{T-1}) = \Pr(q_1)\prod_{t=2}^{T}\Pr(q_t|q_{t-1})$$

with the starting probability $\Pr(q_1) = \pi_{q_1}$ = probability of using $q_1$ as the first POS tag.

- Likelihood

$q_1 \in S$

$$\Pr(O|Q) = \prod_{t=1}^{T}\Pr(o_t|q_t) = \prod_{t=1}^{T} B_{q_t, o_t}$$

$\uparrow$ emission prob.

$\pi = [Pr(S_1) \cdots Pr(S_N)]$

$\underbrace{\qquad\qquad}_{N}$

- Posterior

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

MAP:

$Q^* = \arg\max_Q Pr(Q|O) = \dfrac{Pr(O|Q)\,Pr(Q)}{\sum_Q Pr(O,Q)} = \dfrac{Pr(O|Q)\,Pr(Q)}{\sum_Q Pr(O|Q)Pr(Q)}$

# HMM

Similar to n-gram model

- The transitions between tags is a bi-gram model  *over pos – tags*

*Bi–gram*

$$Pr(w_t \mid w_{t-1})$$

Different from n-gram model

- The transitions happen on the hidden state (POS-tag) level;
- Words are independent given the POS sequence.

*vs.* $Pr(y_t \mid q_{t-1})$

Three tasks with HMM:

- find the likelihood of a POS-tag sequence (forward algorithm);   $P(Q \mid O)$

$max\, P(Q\mid O)$

- find the most likely POS-tag sequence given an observed sentence (Viterbi algorithm);  $Q$
- find the start, transition, and emission probabilities (MLE with backward algorithm).

*train HMM from training corpus (pos – tagged sentences)*

$\overset{\Delta}{=}$ *Estimating*

*A · transition prob*
*B · emission Prob.*
*$T_1$ : start prob*

*EM – alg w. forward backward*