

# HMM and Part of Speech Tagging

Adam Meyers  
New York University



# Outline

- Parts of Speech Tagsets
- Rule-based POS Tagging
- HMM POS Tagging
- Transformation-based POS Tagging



# Part of Speech Tags Standards

- There is no standard set of parts of speech that is used by all researchers for all languages.
- The most commonly used English tagset is that of the Penn Treebank at the University of Pennsylvania:
  - pdf of POS annotation guidelines at class website
    - [http://cs.nyu.edu/courses/spring12/CSCI-GA.2590-001/ptb\\_tagguide.pdf](http://cs.nyu.edu/courses/spring12/CSCI-GA.2590-001/ptb_tagguide.pdf)
- To map several POS tagsets to each other, see Table 1 in:
  - <http://nlp.cs.nyu.edu/meyers/Annotation%20Compatibility%20Working%20Group%20Report%202006.html>
- POS tagsets:
  - Assume Particular Tokenizations, e.g., *Mary's* → *Mary* + *'s*
  - Distinguish inflectional morphology: e.g., *eat*/**VB**, *eats*/**VBZ**, *ate*/**VBD**
  - Are exceptionless – there are tags for all possible tokens



# The Penn Treebank II POS tagset

- Verbs: VB, VBP, VBZ, VBD, VBG, VBN
  - base, present-non-3rd, present-3rd, past, -ing, -en
- Nouns: NNP, NNPS, NN, NNS
  - proper/common, singular/plural (singular includes mass + generic)
- Adjectives: JJ, JJR, JJS (base, comparative, superlative)
- Adverbs: RB, RBR, RBS, RP (base, comparative, superlative, particle)
- Pronouns: PRP, PP\$ (personal, possessive)
- Interrogatives: WP, WP\$, WDT, WRB (compare to: PRP, PP\$, DT, RB)
- Other Closed Class: CC, CD, DT, PDT, IN, MD
- Punctuation: # \$ . , : ( ) “ ” ’ `
- Weird Cases: FW(*deja vu*), SYM (@), LS (*1, 2, a, b*), TO (*to*), POS(*'s, '*), UH (*no, OK, well*), EX (*it/there*)
- Newer tags: HYPH, PU



# Part of Speech Tagging

- POS taggers assign 1 POS tag to each input token
  - *The/DT silly/JJ man/NN is/VBZ a/DT professor/NN ./PU*
- Different ways of breaking down POS tagging:
  - Use separate “tokenizer”, program that divides string into list of tokens – POS tagger processes output
  - Incorporate tokenizer into POS tagger
- Different ways of breaking down parsing:
  - Use separate POS tagger – output of tagger is input to parser
  - Assign POS tags as part of parsing (assumed previously)
- Accurate POS tagging is “easier” than accurate parsing
  - POS tags may be sufficient information for some tasks



# Some Tokenization Rules for English

- 1) Divide at spaces and hyphens.
- 2) Divide before punctuation that is followed by: a space or the end of the line
  - Define punctuation as any non-letter/non-number:
    - `!@#\$%^&\*()-\_+={[]\|:;'"<>./
  - Punctuation followed by a space, other punctuation, or at the end of line should be separated from words:
    - ...*and he left.*") → *and* *he* *left* . " )
- 3) Break off the following as separate tokens when followed by a space or end of line:
  - 's, n't, 'd, 've, 'm, 'll, 're, ... (a short list)
- 4) Abbreviations are exceptions to rule 2:
  - Period after abbreviations should not be separate from words
    - Most cases covered by list of 100 items (or if sentence end is known)
  - Final periods are not duplicated after abbreviations (consistency issues)



# Rule-based POS Tagger

- Method
  - Assign lists of potential POS tags to each word based on dictionary
  - Manual rules for Out of Vocabulary (OOV) words
    - Ex: Non-initial capital → NNP; ends in S → VBZ|NNS; default → NN|JJ; etc.
  - Apply hand-written constraints until each word has only one possible POS
- Sample Constraints:
  - 1) DT cannot immediately precede a verb
  - 2) No verb can immediately precede a tensed verb (VBZ, VBP, VBD)
- Example:
  - The/DT book/{NN|VB|VBP} is/VBZ on/IN the/DT table{NN|VB|VBP}
  - The/DT book/NN is/VBZ on/IN the/DT table/NN
    - DT cannot precede VB or VBP
    - VBZ cannot be preceded by VB or VBP



# Introduction to Probability

- Estimate of probability of future event based on past observations

$$P(event) = \frac{\text{num of events}}{\text{num of trials}}$$

- Conditional Probability: probability of X given Y

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

- Examples:
  - If out of 1000 words, *the* occurs 200 times, we predict:
    - a random word is 20% likely to be *the*
  - If the word after *the* is a noun 120 times and an adjective 60 times:
    - A word following *the* is
      - $120/200 = 60\%$  likely to be a noun
      - $60/200 = 30\%$  likely to be an adjective





# More Math Terminology

- N instances of a variable looked at individually:

$X^n$  is the same as  $\{X_1, X_2, X_3, \dots, X_n\}$  in sequence

- For example,

$$\prod_{i=1}^n P(X_i)$$

means the product of instances of X from 1 to n

- Max = the maximum number in a set
- Argmax = the formula that maximizes a particular argument of the formula



# Probabilistic Models of POS tagging

- For tokens  $w_1, \dots, w_n$ , find the most probable corresponding sequence of possible tags  $t_1, \dots, t_n$ 
  - We assume that *probable* means something like “most frequently observed in some manually tagged corpus of words”.
- Penn Treebank II (a common training corpus)
  - 1 million words from the Wall Street Journal
  - Tagged for POS (and other attributes)
- The specific sequence (sentence) is not in the training corpus
  - Therefore the actual “probability” is 0
  - Common practice: estimate probability given assumptions, e.g.,
    - Choose the most frequent possible sequence of tags independently of how frequently each tag is assigned to each token



# Probabilistic Assumptions of HMM Tagging

- $\hat{t} = \underset{t_1^n}{\operatorname{argmax}} P(t_1^n | w_1^n)$ 
  - Choose the tag sequence of length n that is most probable given the input token sequence
- Bayes Rule:
  - $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$
  - Way to derive the probability of x given y when you know the probability of y given x
- Applying Bayes Rule to Tag Probability
  - $\hat{t} = \underset{t_1^n}{\operatorname{argmax}} \frac{P(w_1^n | t_1^n) P(t_1^n)}{P(w_1^n)}$



# Simplifying Assumptions for HMMs

- Simplification: Drop the denominator
  - Denominator is same for all the tag sequences
  - $\hat{t} = \underset{t_1^n}{\operatorname{argmax}} P(w_1^n | t_1^n) P(t_1^n)$
  - For each tag sequence calculate the product of:
    - The probability of the word sequence given the tag sequence (**likelihood**)
    - The probability of the tag sequence (**prior probability**)
  - Still too hard
- 2 simplifying assumptions make it possible to estimate the probability of tag sequences given word sequences:
  - 1) If the probability of a word is only dependent on its own POS tag,
    - $P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i)$
  - 2) If the probability of a POS tag is only dependent on the previous POS tag,
    - $P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$
- The result of these assumptions:  $\hat{t} \approx \underset{t_1^n}{\operatorname{argmax}} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$
- HMM taggers are fast and achieve precision/recall scores of about 93-95%



# Estimating Probability of $\hat{t}$

- We assume that:  $\hat{t} \approx \underset{t_1^n}{argmax} \prod_{i=1}^n P(w_i|t_i)P(t_i|t_{i-1})$
- Acquire frequencies from a training corpus:
  - Word Frequency with given POS
    - suppose **book** occurs 14 times in a corpus: 10 times (.001) as **NN** (there are 10000 instances of **NN** in the corpus); 3 times (.003) as **VBP** (the corpus has 1000 **VBPs**), and 1 instance of book (.005) as **VB** (the corpus has 500 **VBs**).
  - Given the previous tag, how often does each tag occur
    - suppose **DT is** followed by **NN** 80,000 times (.53), **JJ** 30,000 times (.2), **NNS** 20,000 times (.13), **VCN** 3,000 (.02) times, ... out of a total of 150,000 occurrences of **DT**
- All possible tags for sequence:
  - *The/DT book/{NN|VB|VBP} is/VBZ on/IN the/DT table/{NN|VB|VBP}*
- Hypothetical probabilities for highest scoring tag sequence:
  - *The/DT book/NN is/VBZ on/IN the/DT table/NN*
  - *The/DT=.4, book/NN=.001, is/VBZ=.02, on/IN=.1, the/DT=.4, table/NN=.0005,*
  - **B DT = .61, DT NN = .53, NN VBZ = .44, VBZ IN = .12, IN DT = .05, DT NN = .53 NN E .31**
  - $\prod_{i=1}^n P(w_i|t_i)P(t_i|t_{i-1}) = (.4 \times .61)(.001 \times .53)(.02 \times .44)(.1 \times .12)(.4 \times .05)(.0005 \times .53)(1 \times .31) \approx 2.4 \times 10^{-13}$

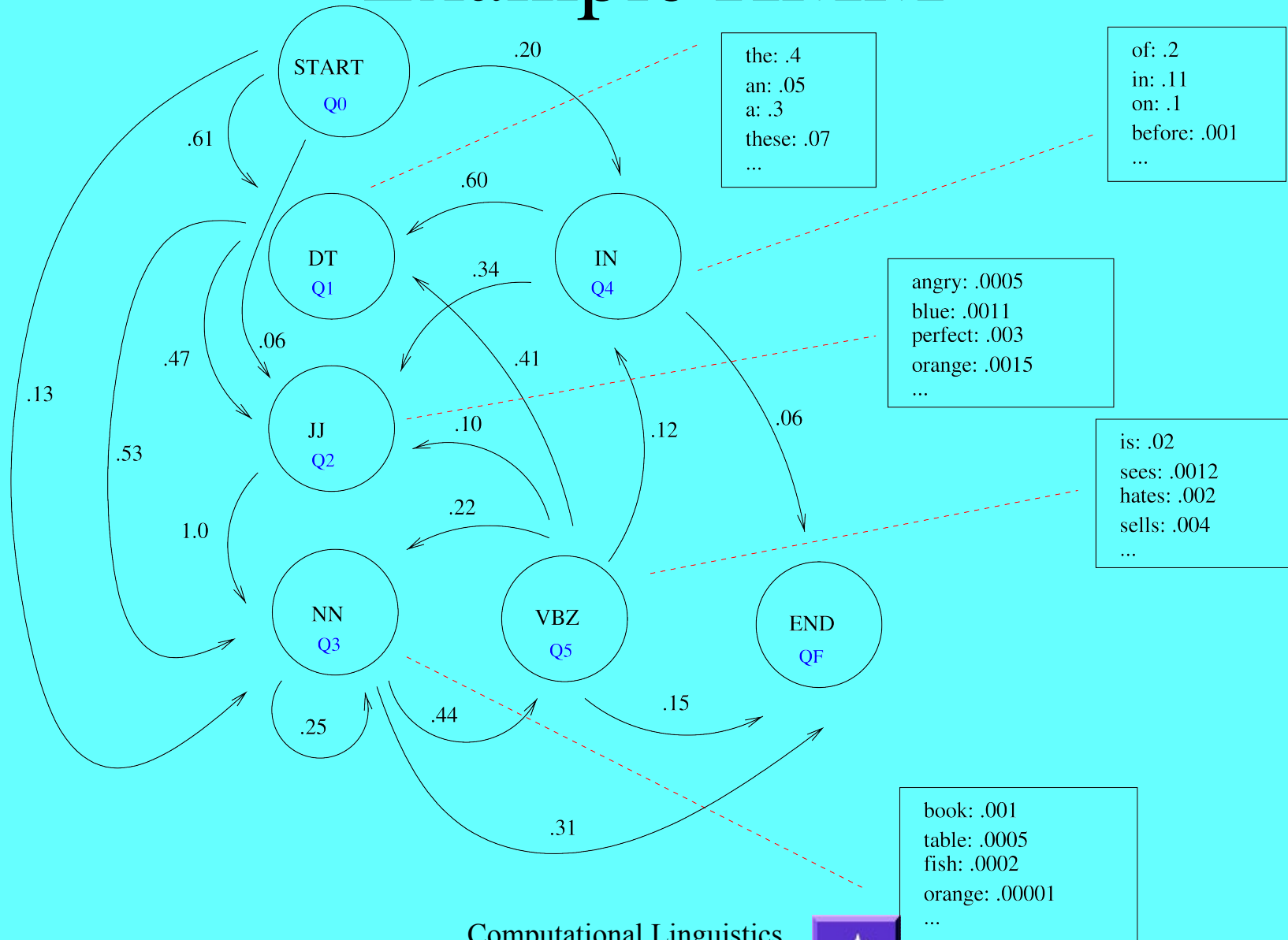


# Defining an HMM

- A Weighted Finite-state Automaton (WFSA)
  - Each transition arc is associated with a probability
  - The sum of all arcs outgoing from a single node is 1
- Markov chain is a WFSA in which an input string uniquely determine path through the Automaton
- Hidden Markov Model (HMM) is a slightly different case because some information (previous POS tags) is unknown (or hidden)
- HMM consists of the following:
  - $\mathbf{Q}$  = set of states:  $\mathbf{q}_0$  (start state), ...,  $\mathbf{q}_F$  (final state)
  - $\mathbf{A}$  = transition probability matrix of  $\mathbf{n} \times \mathbf{n}$  probabilities of transitioning between any pair of  $\mathbf{n}$  states ( $\mathbf{n} = \mathbf{F}+1$ ). Called: *prior probability* or *transition probability* of a tag sequence
  - $\mathbf{O}$  = sequence of  $\mathbf{T}$  observations (**words**) from a vocabulary  $\mathbf{V}$
  - $\mathbf{B}$  = sequence of observation likelihoods (probability of observation generated at state) – Called *likelihood* (of word sequence given tag sequence), aka *emission probability*



# Example HMM



# Viterbi Algorithm for HMM

Observed\_Words =  $w_1 \dots w_T$

- States =  $q_0, q_1 \dots q_N, q_F$

$A = N \times N$  matrix such that  $a_{i,j}$  is the probability of the transition from  $q_i$  to  $q_j$

$B$  = lookup table such that  $b_i(w_t)$  is the probability that POS  $i$  is realized as word  $t$

**viterbi** =  $(N+2) \times T$  matrix # columns are states, rows are words

**backpointer** =  $(N+2) \times T$  matrix # highest scoring previous cells for viterbi

for states  $q$  from 1 to  $N$ :

initialize **viterbi**[ $q,1$ ] to  $a_{0,q} * b_q(w_1)$  # score transition  $0 \rightarrow q$  given  $w_1$

initialize **backpointer**[ $q,1$ ] to 0 (start state)

for word  $w$  from 2 to  $T$ :

for state  $q$  from 1 to  $N$ :

# for  $T-1 \times N$  ( $w,q$ ) pairs

**viterbi**[ $q,w$ ]  $\leftarrow \max_{q'=1}^N \text{viterbi}[q',t-1] * a_{q',q} * b_q(w_t)$  # score = maximum previous \* prior \* likelihood

**backpointer**[ $q,w$ ]  $\leftarrow \operatorname{argmax}_{q'=1}^N \text{viterbi}[q',t-1] * a_{q',q}$  # backpointer = maximum previous

**viterbi**[ $q_F,T$ ]  $\leftarrow \max_{q=1}^N \text{viterbi}[q,T] * a_{q,q_F}$  # score = maximum previous \* prior \* likelihood

**backpointer**[ $q_F,T$ ]  $\leftarrow \operatorname{argmax}_{q=1}^N \text{viterbi}[q,T] * a_{q,q_F}$  # backpointer = maximum previous

- return(best\_path) # derive by following backpointers from ( $q_F,T$ ) to  $q_0$





# *Walk Through: The book fell off the table.* (ignoring period)

$$1 * .4 * .61 * .00001 * .53 * .02 * .33 * .1 * .12 * .4 * .6 * .54 * .0005 * .33 * 1 = 2.19 * 10^{-15}$$

	0	1: <i>The</i>	2: <i>orange</i>	3: <i>is</i>	4: <i>on</i>	5: <i>the</i>	6: <i>table</i>	7
Start	1							
DT		.4 * .61				4*.6		
IN					.1*.12			
JJ			.0015 * .47					
VBZ								
NN			.00001 * .53	.02*.44			.53*.0005	
End								.33 * 1



# Comments on Viterbi Trace

- Initialize scores for first column: transitions from 0 to each possible state given: *the*
  - The probability of reaching Q1 matching the first item on the tape (*the*) will be  $.4 \times .61 = .244$  (this is also the only possibility)
- The adjective sense of *orange* is more likely locally, but leads to a dead end
- The transitions from B and the transition to E are necessary parts of the process.



# Go to Ralph's Viterbi Demo for *Fish Sleep*



# N-grams

- Prior Probability in our HMM was estimated as:
  - $\text{freq}(\text{currentPOS}, \text{previousPOS}) / \text{freq}(\text{previousPOS})$
  - This is a bigram
- Trigram:  $\text{freq}(\text{POS-2}, \text{POS-1}, \text{POS}) / \text{freq}(\text{POS-2}, \text{POS-1})$
- N-gram:  $\text{freq}(\text{POS-N} \dots \text{POS}) / \text{freq}(\text{POS-N} \dots \text{POS-1})$
- N-grams
  - Used a lot in computational linguistics
  - N-grams can be calculated for: words, POS, other attributes
  - Chapter 4 has more info on N-grams
- J & M show how to extend HMM to tri-grams by calculating the prior based on two previous POS tags instead of one
  - They also show a way of combining unigram, bigram and trigrams



# Unknown (OOV) Words

- Possibility 1
  - Assume they are equally ambiguous between all POS
    - Refinements:
      - Assume distribution of unknown = overall POS distribution
      - Assume distribution of unknown = distribution of 1 time words
    - Use N-grams to predict the correct tag
- Possibility 2
  - Use morphology (prefixes, suffixes), orthography (uppercase/lowercase), hyphenation
- Possibility 3: Some combination



# Homework

- <http://cs.nyu.edu/courses/spring16/CSCI-UA.0480-011/homework4.html>

