

# Natural Language Processing

## Sequence labeling

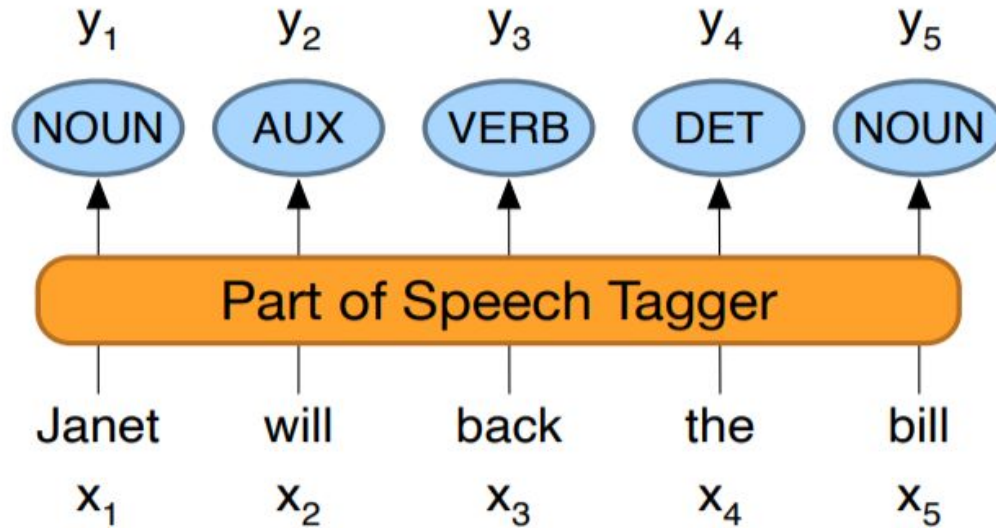
Yulia Tsvetkov

[yuliats@cs.washington.edu](mailto:yuliats@cs.washington.edu)

# Announcements

- HW 2 review

# Part of speech tagging



# Why POS tagging

- Goal: resolve ambiguities
- Text-to-speech
  - record, lead, protest
- Lemmatization
  - saw/V → see, saw/N → saw
- Machine translation – the meaning of a word depends on its POS tag
  - “what is your address” vs. “we should address this comment”
  - “grand challenge” vs. “challenge the team”
- Sentiment analysis – adjectives are the major opinion holders
  - “good” vs. “bad”, “excellent” vs. “terrible”
- Preprocessing for harder disambiguation problems
  - syntactic parsing
  - semantic parsing

# Parts of speech

- **Open classes**

- nouns
- verbs
- adjectives
- adverbs

- **Closed classes**

- prepositions
- determiners
- pronouns
- conjunctions
- auxiliary verbs

# Parts of speech, more fine-grained classes

- **Open classes**

- nouns
  - proper
  - common
    - count
    - mass
- verbs
- adjectives
- adverbs
  - directional
  - degree
  - manner
  - temporal

*Actually, I ran home extremely quickly yesterday*

# Parts of speech, closed classes

**prepositions:** on, under, over, near, by, at, from, to, with

**particles:** up, down, on, off, in, out, at, by

**determiners:** a, an, the

**conjunctions:** and, but, or, as, if, when

**pronouns:** she, who, I, others

**auxiliary verbs:** can, may, should, are

**numerals:** one, two, three, first, second, third

# Part of speech tagsets

- Penn treebank tagset (Marcus et al., 1993)

| Tag  | Description                   | Example             | Tag  | Description           | Example            | Tag | Description          | Example              |
|------|-------------------------------|---------------------|------|-----------------------|--------------------|-----|----------------------|----------------------|
| CC   | coordinating conjunction      | <i>and, but, or</i> | PDT  | predeterminer         | <i>all, both</i>   | VBP | verb non-3sg present | <i>eat</i>           |
| CD   | cardinal number               | <i>one, two</i>     | POS  | possessive ending     | <i>'s</i>          | VBZ | verb 3sg pres        | <i>eats</i>          |
| DT   | determiner                    | <i>a, the</i>       | PRP  | personal pronoun      | <i>I, you, he</i>  | WDT | wh-determ.           | <i>which, that</i>   |
| EX   | existential 'there'           | <i>there</i>        | PRPS | possess. pronoun      | <i>your, one's</i> | WP  | wh-pronoun           | <i>what, who</i>     |
| FW   | foreign word                  | <i>mea culpa</i>    | RB   | adverb                | <i>quickly</i>     | WPS | wh-possess.          | <i>whose</i>         |
| IN   | preposition/<br>subordin-conj | <i>of, in, by</i>   | RBR  | comparative<br>adverb | <i>faster</i>      | WRB | wh-adverb            | <i>how, where</i>    |
| JJ   | adjective                     | <i>yellow</i>       | RBS  | superlatv. adverb     | <i>fastest</i>     | \$  | dollar sign          | <i>\$</i>            |
| JJR  | comparative adj               | <i>bigger</i>       | RP   | particle              | <i>up, off</i>     | #   | pound sign           | <i>#</i>             |
| JJS  | superlative adj               | <i>wildest</i>      | SYM  | symbol                | <i>+, %, &amp;</i> | "   | left quote           | <i>' or "</i>        |
| LS   | list item marker              | <i>1, 2, One</i>    | TO   | "to"                  | <i>to</i>          | "   | right quote          | <i>' or "</i>        |
| MD   | modal                         | <i>can, should</i>  | UH   | interjection          | <i>ah, oops</i>    | (   | left paren           | <i>[, (, {, &lt;</i> |
| NN   | sing or mass noun             | <i>llama</i>        | VB   | verb base form        | <i>eat</i>         | )   | right paren          | <i>], ), }, &gt;</i> |
| NNS  | noun, plural                  | <i>llamas</i>       | VBD  | verb past tense       | <i>ate</i>         | ,   | comma                | <i>,</i>             |
| NNP  | proper noun, sing.            | <i>IBM</i>          | VBG  | verb gerund           | <i>eating</i>      | .   | sent-end punc        | <i>. ! ?</i>         |
| NNPS | proper noun, plu.             | <i>Carolinas</i>    | VBN  | verb past part.       | <i>eaten</i>       | :   | sent-mid punc        | <i>: ; ... --</i>    |



# Example of POS tagging

- There are 7 children there.

There/**PRON** are/**AUX** 70/**NUM** children/**NOUN** there/**ADV** ./**PUNC**

- Preliminary findings were reported in today's NY TIMES.

Preliminary/**ADJ** findings/**NOUN** were/**AUX** reported/**VERB** in/**ADP**  
today/**NOUN** 's/**PART** NY/**PROPN** TIMES/**PROPN** ./**PUNC**

# The Universal Dependencies

## Universal Dependencies

Universal Dependencies (UD) is a framework for consistent annotation of grammar (parts of speech, morphological features, and syntactic dependencies) across different human languages. UD is an open community effort with over 300 contributors producing more than 150 treebanks in 90 languages. If you're new to UD, you should start by reading the first part of the Short Introduction and then browsing the annotation guidelines.

- [Short introduction to UD](#)
- [UD annotation guidelines](#)
- More information on UD:
  - [How to contribute to UD](#)
  - [Tools for working with UD](#)
  - [Discussion on UD](#)
  - [UD-related events](#)
- Query UD treebanks online:
  - [SETS treebank search](#) maintained by the University of Turku
  - [PML Tree Query](#) maintained by the Charles University in Prague
  - [Kontext](#) maintained by the Charles University in Prague
  - [Grew-match](#) maintained by Inria in Nancy
  - [INESS](#) maintained by the University of Bergen
- [Download UD treebanks](#)

| Open class words      | Closed class words    | Other                 |
|-----------------------|-----------------------|-----------------------|
| <a href="#">ADJ</a>   | <a href="#">ADP</a>   | <a href="#">PUNCT</a> |
| <a href="#">ADV</a>   | <a href="#">AUX</a>   | <a href="#">SYM</a>   |
| <a href="#">INTJ</a>  | <a href="#">CCONJ</a> | <a href="#">X</a>     |
| <a href="#">NOUN</a>  | <a href="#">DET</a>   |                       |
| <a href="#">PROPN</a> | <a href="#">NUM</a>   |                       |
| <a href="#">VERB</a>  | <a href="#">PART</a>  |                       |
|                       | <a href="#">PRON</a>  |                       |
|                       | <a href="#">SCONJ</a> |                       |

# The Universal Dependencies

|              | tag   | tag name                  | description   | examples  |
|--------------|-------|---------------------------|---|---|
| open class   | ADJ   | Adjective                 | noun modifiers describing properties  | red, young, awesome                             |
|              | ADV   | Adverb                    | verb modifiers of time, place, manner   | very, slowly, home (as in "go home"), yesterday |
|              | NOUN  | Noun                      | words for persons, places, things, etc.   | algorithm, cat, mango, beauty                   |
|              | VERB  | Verb                      | words for actions and processes   | draw, provide, run                              |
|              | PROPN | Proper noun               | name of a person, organization, location, etc.                                  | Regina, IBM, Vietnam                            |
|              | INTJ  | Interjection              | exclamation, greeting, yes/no response, etc.                                    | oh, um, yes, hello                              |
| closed class | ADP   | Adposition                | (preposition/postposition): marks a noun's spatial, temporal, or other relation | in, on, by, under                               |
|              | AUX   | Auxiliary                 | helping a verb making tense, aspect, mood, etc.                                 | can, may, should, are                           |
|              | CCONJ | Coordinating conjunction  | joins two phrases or clauses  | and, or, but                                    |
|              | DET   | Determiner                | marks noun phrase properties  | a, an, the, this                                |
|              | NUM   | Numeral                   |   | one, two, first, second                         |
|              | PART  | Particle                  | a preposition-like form used together with a verb                               | up, down, on, off, in, out, at, by              |
|              | PRON  | Pronoun                   | a shorthand for referring to an entity or event                                 | you, she, who, I others                         |
|              | SCONJ | Subordinating conjunction | joins a main clause with a subordinate clause such as sentential complement     | that, which                                     |
| other        | PUNCT | Punctuation               |   | ; ! ? ()  |
|              | SYM   | Symbol                    | various symbols   | \$ %  |
|              | X     | Other                     |   | all the rest...                                 |

Universal Dependencies tagset (Nivre et al., 2016) (17 tags)

# Ambiguities in POS tags

| <b>Types:</b>      |           | <b>WSJ</b>            | <b>Brown</b>          |
|--------------------|-----------|-----------------------|-----------------------|
| <b>Unambiguous</b> | (1 tag)   | 44,432 ( <b>86%</b> ) | 45,799 ( <b>85%</b> ) |
| <b>Ambiguous</b>   | (2+ tags) | 7,025 ( <b>14%</b> )  | 8,050 ( <b>15%</b> )  |

# Ambiguities in POS tags

| <b>Types:</b>      |           | <b>WSJ</b> |              | <b>Brown</b> |              |
|--------------------|-----------|------------|--------------|--------------|--------------|
| <b>Unambiguous</b> | (1 tag)   | 44,432     | <b>(86%)</b> | 45,799       | <b>(85%)</b> |
| <b>Ambiguous</b>   | (2+ tags) | 7,025      | <b>(14%)</b> | 8,050        | <b>(15%)</b> |
| <b>Tokens:</b>     |           |            |              |              |              |
| <b>Unambiguous</b> | (1 tag)   | 577,421    | <b>(45%)</b> | 384,349      | <b>(33%)</b> |
| <b>Ambiguous</b>   | (2+ tags) | 711,780    | <b>(55%)</b> | 786,646      | <b>(67%)</b> |

# POS-tagging is a disambiguation task

- look at this small building in the **back/NN**
  - the main player in the team took a **back/JJ** seat
  - a majority of senators will **back/VERB** the decision
  - we should enable the country to buy **back/PART** its debt
  - I was twenty-one **back/ADV** then
- 
- try yourself with an online POS tagger, e.g., <https://parts-of-speech.info/>

# POS tagging algorithms

- We will introduce classic and state-of-the-art sequence labeling approaches
  - Hidden Markov Model (HMM) + Viterbi, generative
  - RNNs + Conditional Random Field (CRF), discriminative

# Most frequent class baseline

- A baseline: assign each word its **most probable** (frequent) tag
- Always compare a classifier against a baseline at least as good as the most frequent class baseline
- The WSJ training corpus and test on sections 22-24 of the same corpus the most-frequent-tag baseline achieves an accuracy of 92.34%.
- 97% tag accuracy achievable by most algorithms (HMMs, MEMMs, neural networks, rule-based algorithms)

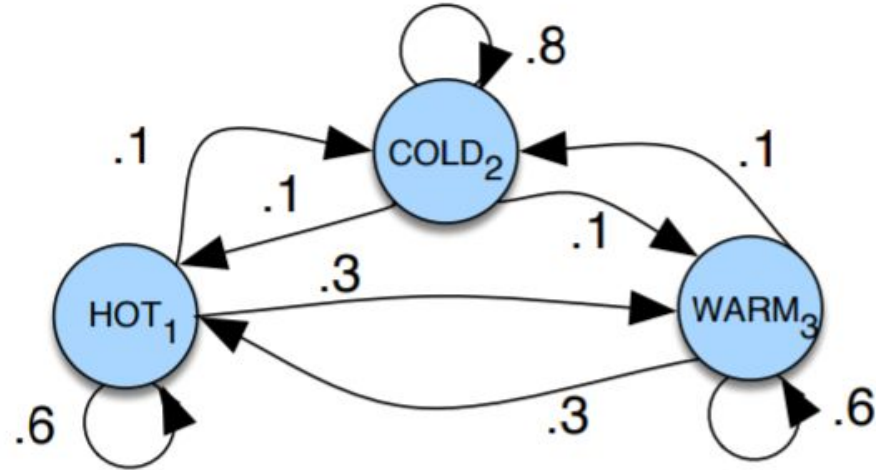


# Sequence labeling as text classification

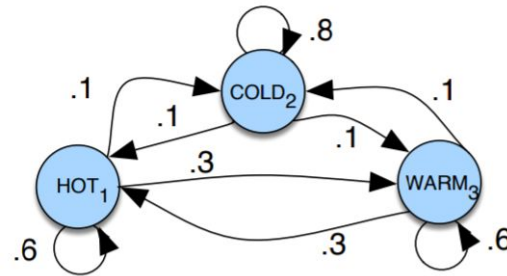
$$\hat{y}_i = \operatorname{argmax}_{y \in \mathcal{L}} s(\mathbf{x}, i, y)$$

# Generative sequence labeling: Hidden Markov Models

# Markov Chain: weather



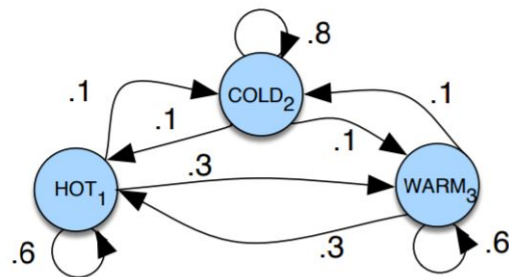
# Markov Chain: weather



- A model that assigns probabilities to sequences of random variables
  - Each variable can take a value from some finite set
- Markov assumption: only the current state matters for predicting the next state, all states before the current have no impact
- For a set of variables:  $q_1, q_2, \dots, q_i$

$$P(q_i = a | q_1, \dots, q_{i-1}) = P(q_i = a | q_{i-1})$$

# Markov Chain: weather



- Given initial probabilities of  $[0.1, 0.7, 0.2]$  compute the probability for:
  - HOT WARM WARM HOT
  - HOT COLD HOT COLD

# Markov chain

- Formally, a Markov chain specified by the following components

$$Q = q_1 q_2 \dots q_N$$

a set of  $N$  **states**

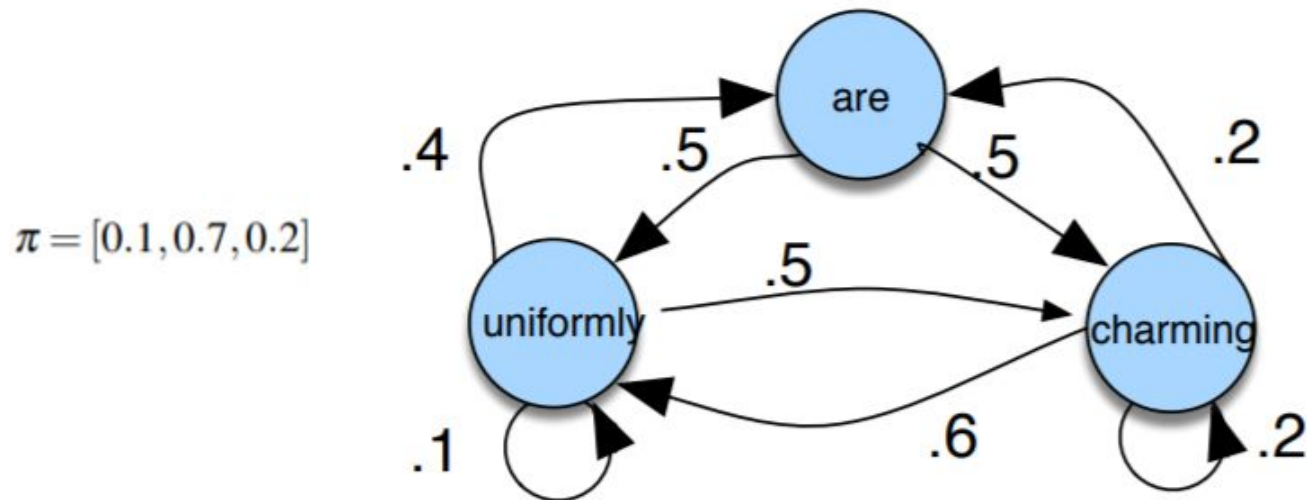
$$A = a_{11} a_{12} \dots a_{n1} \dots a_{nn}$$

a **transition probability matrix**  $A$ , each  $a_{ij}$  representing the probability of moving from state  $i$  to state  $j$ , s.t.  $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$

$$\pi = \pi_1, \pi_2, \dots, \pi_N$$

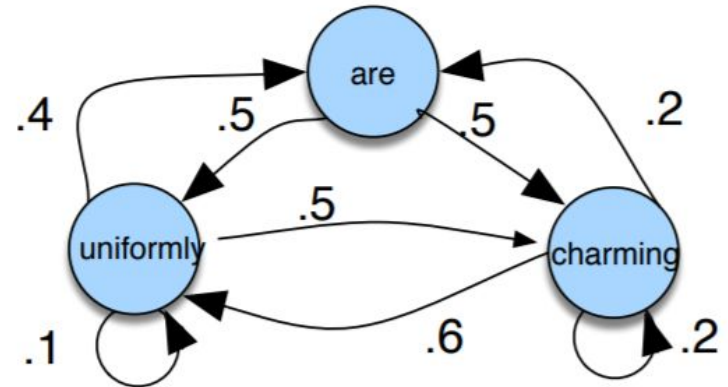
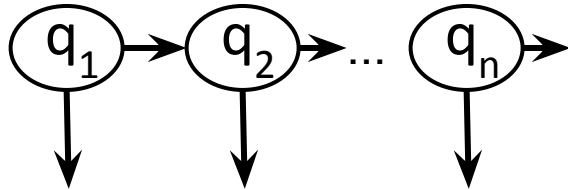
an **initial probability distribution** over states.  $\pi_i$  is the probability that the Markov chain will start in state  $i$ . Some states  $j$  may have  $\pi_j = 0$ , meaning that they cannot be initial states. Also,  $\sum_{i=1}^n \pi_i = 1$

# Markov Chain: words



- only the current state matters for predicting the next state, all states before the current have no impact

# Markov Chain: words

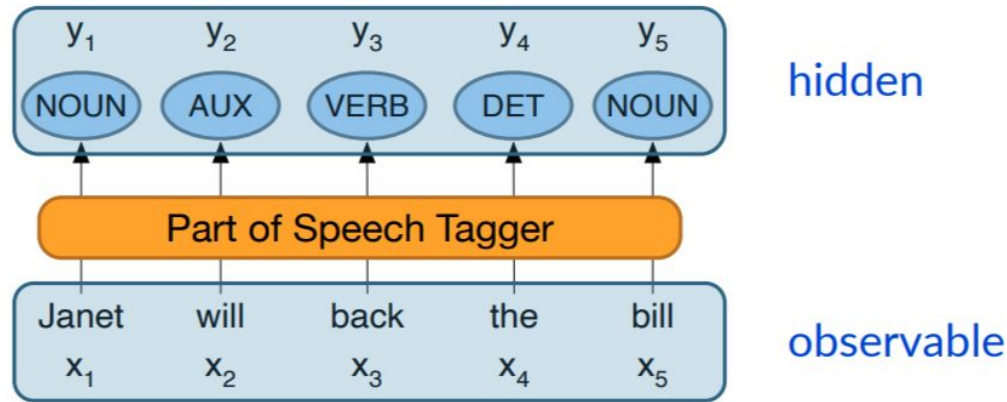


- only the current state matters for predicting the next state, all states before the current have no impact

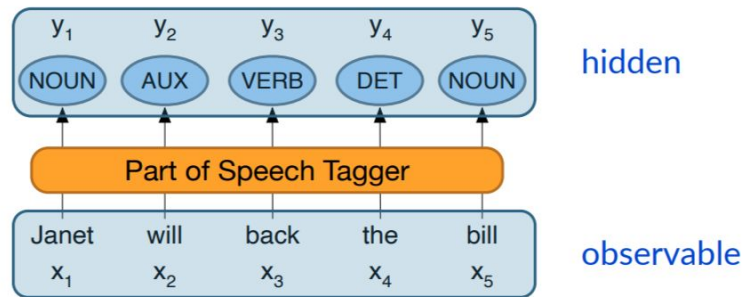
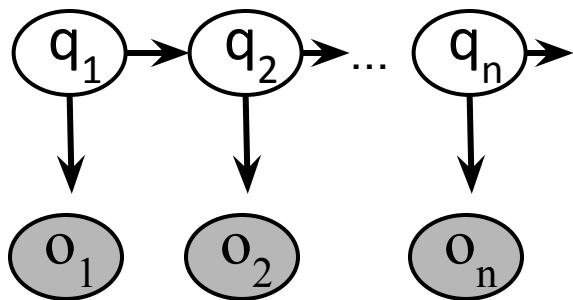


# Hidden Markov Models

- We use a Markov chain for computing  $P$  for a sequence of observable events
- In many cases the events we are interested in are hidden
  - e.g., we don't observe POS tags in a text



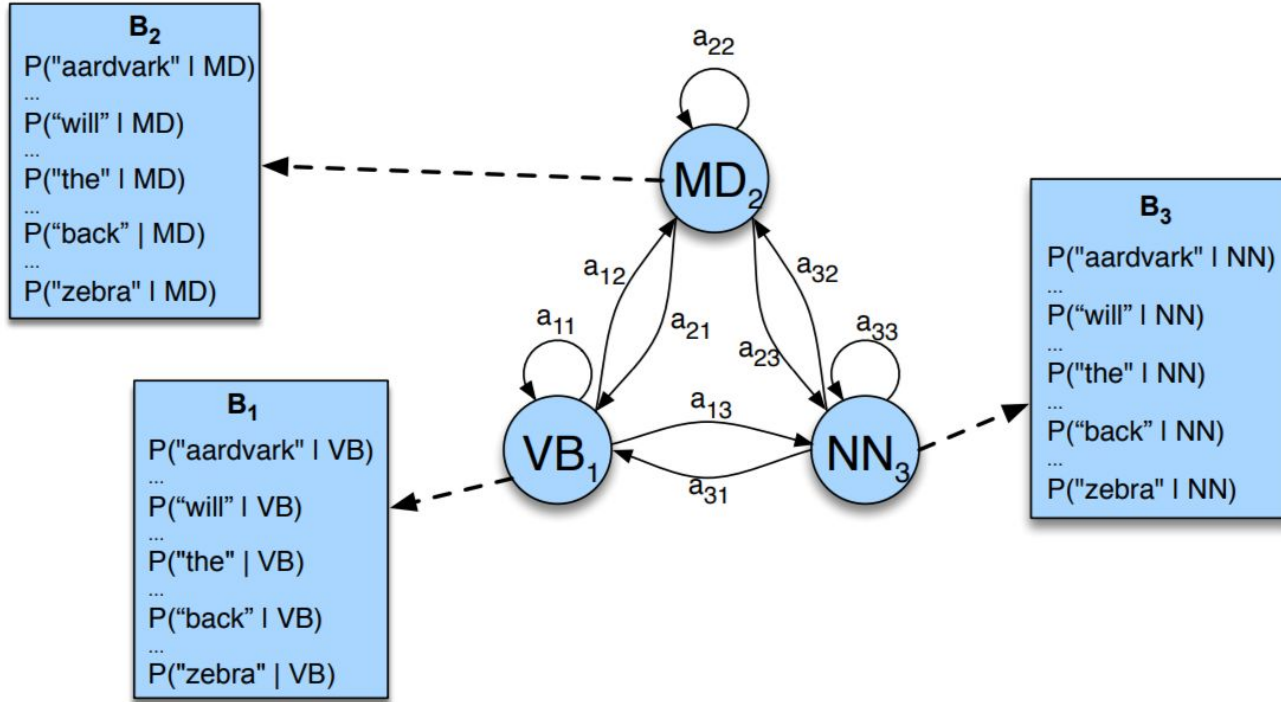
# Hidden Markov Models



**Markov Assumption:**  $P(q_i | q_1 \dots q_{i-1}) = P(q_i | q_{i-1})$

**Output Independence:**  $P(o_i | q_1 \dots q_i, \dots, q_T, o_1, \dots, o_i, \dots, o_T) = P(o_i | q_i)$

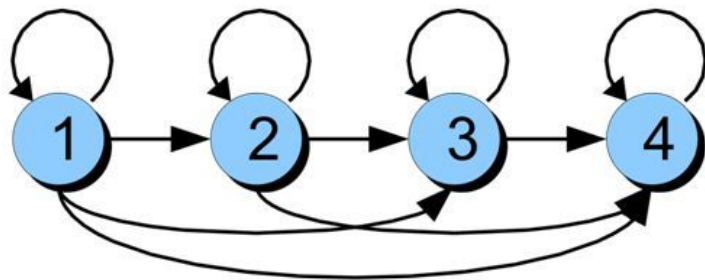
# Hidden Markov Models



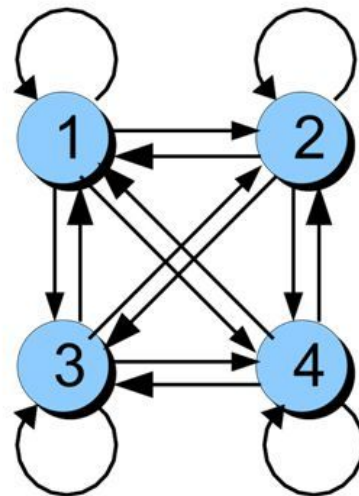
# Hidden Markov Models (HMMs)

|  |  |
|--|--|
| $Q = q_1 q_2 \dots q_N$                | a set of $N$ <b>states</b>   |
| $A = a_{11} \dots a_{ij} \dots a_{NN}$ | a <b>transition probability matrix</b> $A$ , each $a_{ij}$ representing the probability of moving from state $i$ to state $j$ , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$   |
| $O = o_1 o_2 \dots o_T$                | a sequence of $T$ <b>observations</b> , each one drawn from a vocabulary $V = \{v_1, v_2, \dots, v_V\}$  |
| $B = b_i(o_t)$                         | a sequence of <b>observation likelihoods</b> , also called <b>emission probabilities</b> , each expressing the probability of an observation $o_t$ being generated from a state $q_i$  |
| $\pi = \pi_1, \pi_2, \dots, \pi_N$     | an <b>initial probability distribution</b> over states. $\pi_i$ is the probability that the Markov chain will start in state $i$ . Some states $j$ may have $\pi_j = 0$ , meaning that they cannot be initial states. Also, $\sum_{i=1}^N \pi_i = 1$ |

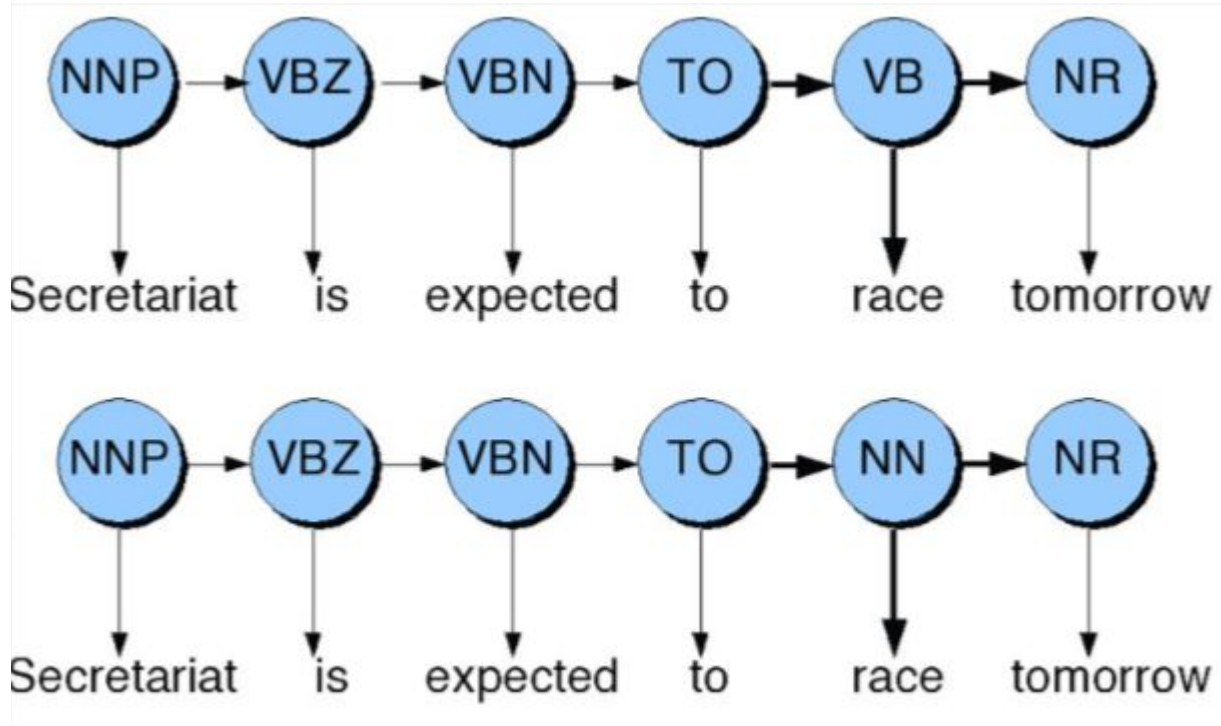
# Types of HMMs



Bakis = left-to-right



Ergodic =  
fully-connected



# HMMs in language technologies

- Part-of-speech tagging (Church, 1988; Brants, 2000)
- Named entity recognition (Bikel et al., 1999) and other information extraction tasks
- Text chunking and shallow parsing (Ramshaw and Marcus, 1995)
- Word alignment of parallel text (Vogel et al., 1996)
- Acoustic models in speech recognition (emissions are continuous)
- Discourse segmentation (labeling parts of a document)

# HMM parameters

$$Q = q_1 q_2 \dots q_N$$

a set of  $N$  **states**

$$A = a_{11} \dots a_{ij} \dots a_{NN}$$

a **transition probability matrix**  $A$ , each  $a_{ij}$  representing the probability of moving from state  $i$  to state  $j$ , s.t.  $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$

$$O = o_1 o_2 \dots o_T$$

a sequence of  $T$  **observations**, each one drawn from a vocabulary  $V = v_1, v_2, \dots, v_V$

$$B = b_i(o_t)$$

a sequence of **observation likelihoods**, also called **emission probabilities**, each expressing the probability of an observation  $o_t$  being generated from a state  $q_i$

$$\pi = \pi_1, \pi_2, \dots, \pi_N$$

an **initial probability distribution** over states.  $\pi_i$  is the probability that the Markov chain will start in state  $i$ . Some states  $j$  may have  $\pi_j = 0$ , meaning that they cannot be initial states. Also,  $\sum_{i=1}^n \pi_i = 1$



# HMMs: algorithms

Forward

**Problem 1 (Likelihood):** Given an HMM  $\lambda = (A, B)$  and an observation sequence  $O$ , determine the likelihood  $P(O|\lambda)$ .

Viterbi

**Problem 2 (Decoding):** Given an observation sequence  $O$  and an HMM  $\lambda = (A, B)$ , discover the best hidden state sequence  $Q$ .

Forward-backward;  
Baum-Welch

**Problem 3 (Learning):** Given an observation sequence  $O$  and the set of states in the HMM, learn the HMM parameters  $A$  and  $B$ .

Forward

**Problem 1 (Likelihood):** Given an HMM  $\lambda = (A, B)$  and an observation sequence  $O$ , determine the likelihood  $P(O|\lambda)$ .

Viterbi

**Problem 2 (Decoding):** Given an observation sequence  $O$  and an HMM  $\lambda = (A, B)$ , discover the best hidden state sequence  $Q$ .

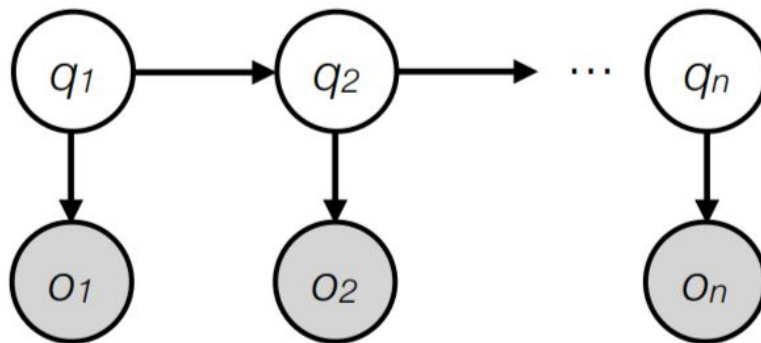
Forward-backward;  
Baum-Welch

**Problem 3 (Learning):** Given an observation sequence  $O$  and the set of states in the HMM, learn the HMM parameters  $A$  and  $B$ .

# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n \mid w_1^n)$$



# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\begin{aligned}\hat{t}_1^n &= \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \\ &= \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n)P(t_1^n)}{P(w_1^n)}\end{aligned}$$

# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\begin{aligned}\hat{t}_1^n &= \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \\ &= \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n)P(t_1^n)}{P(w_1^n)} \\ &= \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n)P(t_1^n)\end{aligned}$$

# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\begin{aligned}\hat{t}_1^n &= \operatorname{argmax}_{t_1^n} P(t_1^n \mid w_1^n) \\ &= \operatorname{argmax}_{t_1^n} \frac{P(w_1^n \mid t_1^n)P(t_1^n)}{P(w_1^n)} \\ &= \operatorname{argmax}_{t_1^n} P(w_1^n \mid t_1^n)P(t_1^n)\end{aligned}$$

simplifying assumptions:

# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\begin{aligned}\hat{t}_1^n &= \operatorname{argmax}_{t_1^n} P(t_1^n \mid w_1^n) \\ &= \operatorname{argmax}_{t_1^n} \frac{P(w_1^n \mid t_1^n)P(t_1^n)}{P(w_1^n)} \\ &= \operatorname{argmax}_{t_1^n} P(w_1^n \mid t_1^n)P(t_1^n)\end{aligned}$$

simplifying assumptions:

$$P(w_1^n \mid t_1^n) \approx \prod_{i=1}^n P(w_i \mid t_i)$$

# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\begin{aligned}\hat{t}_1^n &= \operatorname{argmax}_{t_1^n} P(t_1^n \mid w_1^n) \\ &= \operatorname{argmax}_{t_1^n} \frac{P(w_1^n \mid t_1^n)P(t_1^n)}{P(w_1^n)} \\ &= \operatorname{argmax}_{t_1^n} P(w_1^n \mid t_1^n)P(t_1^n)\end{aligned}$$

simplifying assumptions:

$$P(w_1^n \mid t_1^n) \approx \prod_{i=1}^n P(w_i \mid t_i)$$

$$P(t_1^n) \approx \prod_{i=1}^n P(t_i \mid t_{i-1})$$



# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n \overset{\text{emission, } B}{P(w_i | t_i)} \overset{\text{transition, } A}{P(t_i | t_{i-1})}$$

# HMM tagging as decoding

**Decoding:** Given as input an HMM  $\lambda = (A, B)$  and sequence of observations  $O = o_1, o_2, \dots, o_n$ , find the most probable sequence of states  $Q = q_1, q_2, \dots, q_n$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n \overset{\text{emission, } B}{P(w_i | t_i)} \overset{\text{transition, } A}{P(t_i | t_{i-1})}$$

How many possible choices?

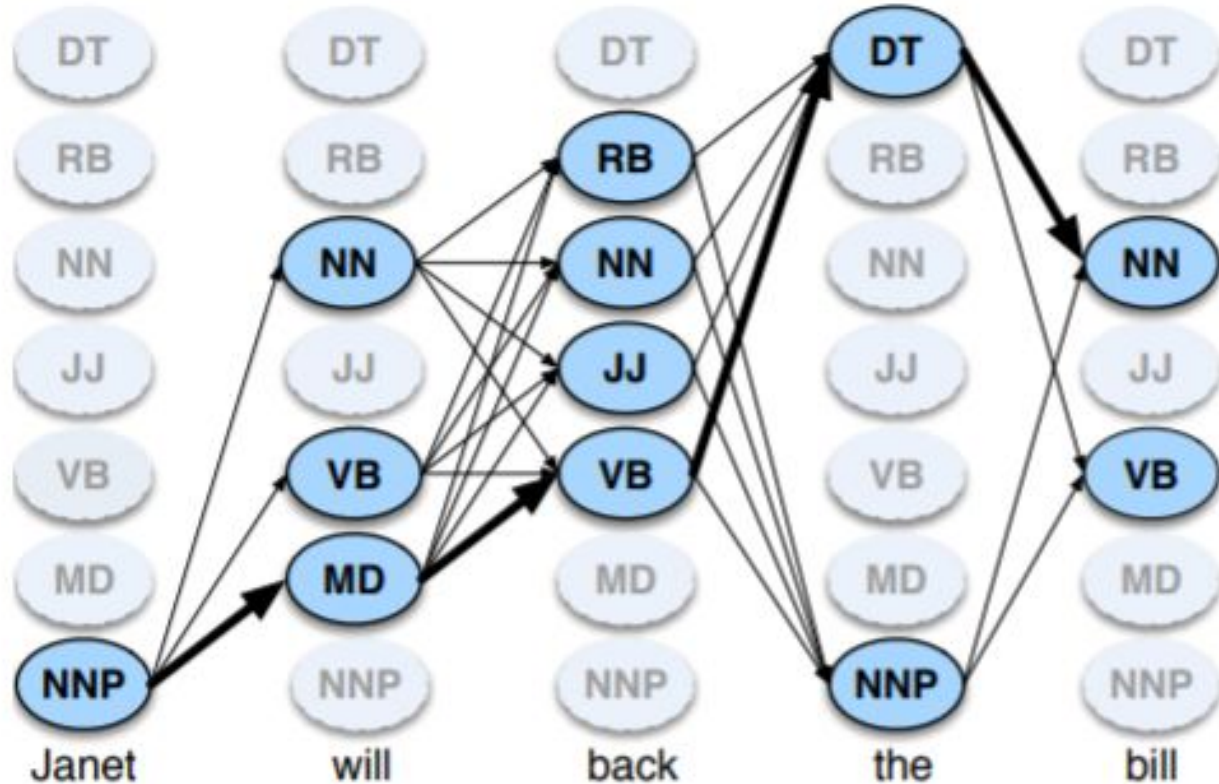
# Part of speech tagging example

|       | I | suspect | the | present | forecast | is | pessimistic | . |
|-------|---|---------|-----|---------|----------|----|-------------|---|
| noun  | • | •       | •   | •       | •        | •  |             |   |
| adj.  |   | •       |     | •       | •        |    | •           |   |
| adv.  |   |         |     | •       |          |    |             |   |
| verb  |   | •       |     | •       | •        | •  |             |   |
| num.  | • |         |     |         |          |    |             |   |
| det.  |   |         | •   |         |          |    |             |   |
| punc. |   |         |     |         |          |    |             | • |

With this very simple tag set,  $7^8 = 5.7$  million labelings.  
(Even restricting to the possibilities above, 288 labelings.)

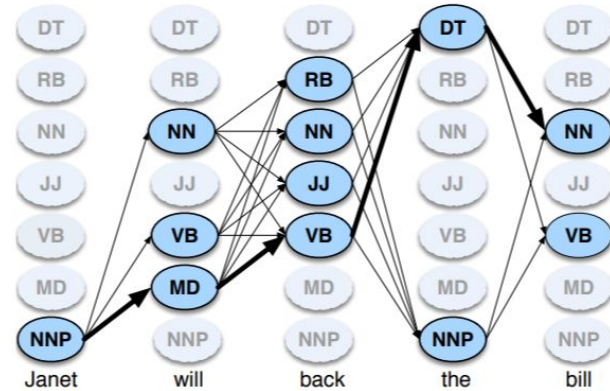
# The Viterbi algorithm

# The Viterbi algorithm



# The Viterbi algorithm

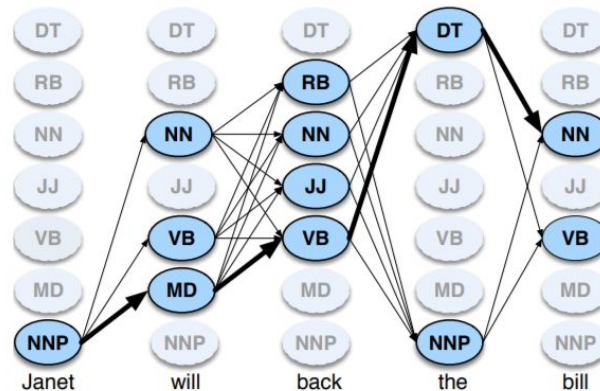
$v_{t-1}(i)$  the **previous Viterbi path probability** from the previous time step  
 $a_{ij}$  the **transition probability** from previous state  $q_i$  to current state  $q_j$   
 $b_j(o_t)$  the **state observation likelihood** of the observation symbol  $o_t$  given the current state  $j$



$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

# The Viterbi algorithm

$v_{t-1}(i)$  the **previous Viterbi path probability** from the previous time step  
 $a_{ij}$  the **transition probability** from previous state  $q_i$  to current state  $q_j$   
 $b_j(o_t)$  the **state observation likelihood** of the observation symbol  $o_t$  given the current state  $j$

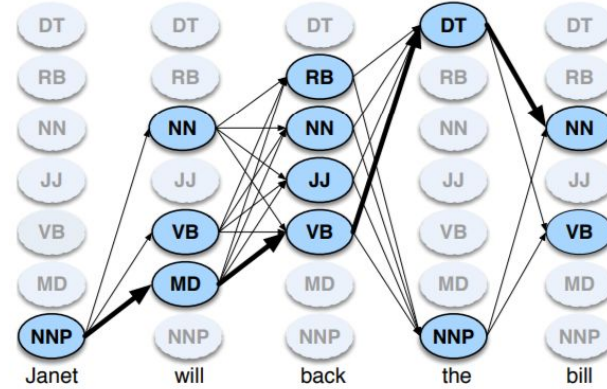


$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

previous  
 Viterbi path  
 probability

# The Viterbi algorithm

$v_{t-1}(i)$  the **previous Viterbi path probability** from the previous time step  
 $a_{ij}$  the **transition probability** from previous state  $q_i$  to current state  $q_j$   
 $b_j(o_t)$  the **state observation likelihood** of the observation symbol  $o_t$  given the current state  $j$



transition  
probability

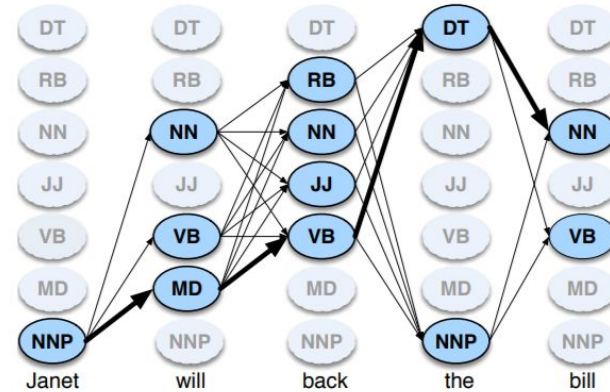
$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

previous  
Viterbi path  
probability



# The Viterbi algorithm

$v_{t-1}(i)$  the **previous Viterbi path probability** from the previous time step  
 $a_{ij}$  the **transition probability** from previous state  $q_i$  to current state  $q_j$   
 $b_j(o_t)$  the **state observation likelihood** of the observation symbol  $o_t$  given the current state  $j$



$$v_t(j) = \max_{i=1}^N \underbrace{v_{t-1}(i)}_{\text{previous Viterbi path probability}} \underbrace{a_{ij}}_{\text{transition probability}} \underbrace{b_j(o_t)}_{\text{state observation likelihood}}$$

# The Viterbi algorithm

**function** VITERBI(*observations* of len  $T$ , *state-graph* of len  $N$ ) **returns** *best-path*, *path-prob*

create a path probability matrix  $viterbi[N, T]$

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$$

$$backpointer[s, 1] \leftarrow 0$$

**for** each time step  $t$  **from** 2 **to**  $T$  **do**

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

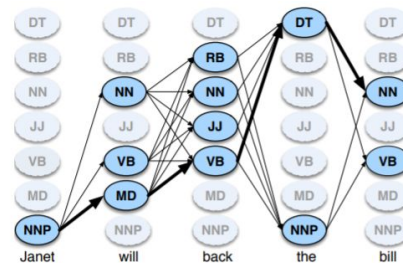
$$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$$

$$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$$

$bestpath \leftarrow$  the path starting at state  $bestpathpointer$ , that follows  $backpointer[]$  to states back in time

**return**  $bestpath$ ,  $bestpathprob$



# The Viterbi algorithm

**function** VITERBI(*observations* of len  $T$ , *state-graph* of len  $N$ ) **returns** *best-path*, *path-prob*

create a path probability matrix  $viterbi[N, T]$

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$$

$$backpointer[s, 1] \leftarrow 0$$

initialization

**for** each time step  $t$  **from** 2 **to**  $T$  **do**

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

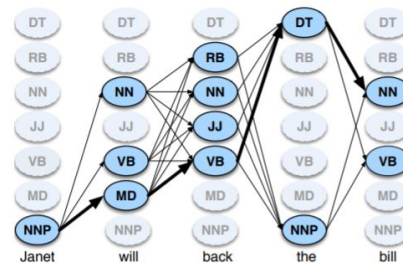
$$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$$

$$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$$

$bestpath \leftarrow$  the path starting at state  $bestpathpointer$ , that follows  $backpointer[]$  to states back in time

**return**  $bestpath$ ,  $bestpathprob$



# The Viterbi algorithm

**function** VITERBI(*observations* of len  $T$ , *state-graph* of len  $N$ ) **returns** *best-path*, *path-prob*

create a path probability matrix  $viterbi[N, T]$

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$$

$$backpointer[s, 1] \leftarrow 0$$

initialization

**for** each time step  $t$  **from** 2 **to**  $T$  **do**

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

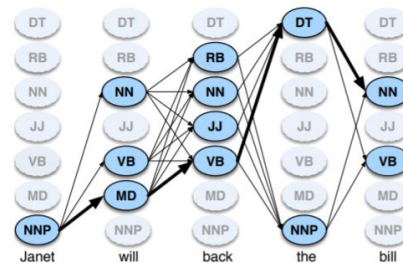
recursion

$$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$$

$$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$$

*bestpath*  $\leftarrow$  the path starting at state *bestpathpointer*, that follows *backpointer*[ ] to states back in time

**return** *bestpath*, *bestpathprob*



# The Viterbi algorithm

**function** VITERBI(*observations* of len  $T$ , *state-graph* of len  $N$ ) **returns** *best-path*, *path-prob*

create a path probability matrix  $viterbi[N, T]$

**for each state  $s$  from 1 to  $N$  do**

$$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$$

$$backpointer[s, 1] \leftarrow 0$$

**for each time step  $t$  from 2 to  $T$  do**

**for each state  $s$  from 1 to  $N$  do**

$$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t) \quad \leftarrow v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

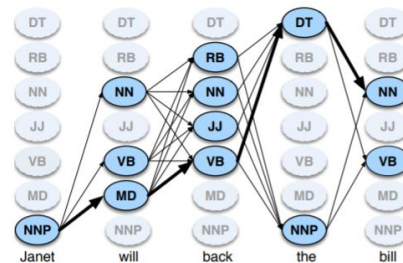
$$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$$

$$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$$

$bestpath \leftarrow$  the path starting at state  $bestpathpointer$ , that follows  $backpointer[]$  to states back in time

**return**  $bestpath$ ,  $bestpathprob$



initialization  
recursion

# The Viterbi algorithm

**function** VITERBI(*observations* of len  $T$ , *state-graph* of len  $N$ ) **returns** *best-path*, *path-prob*

create a path probability matrix  $viterbi[N, T]$

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$$

$$backpointer[s, 1] \leftarrow 0$$

**for** each time step  $t$  **from** 2 **to**  $T$  **do**

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t) \quad \leftarrow v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

$$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$$

$$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$$

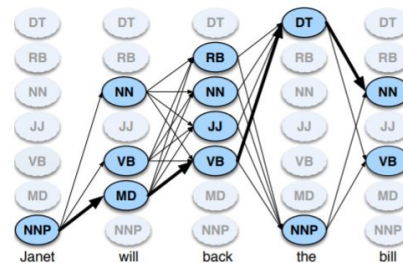
$bestpath \leftarrow$  the path starting at state  $bestpathpointer$ , that follows  $backpointer[]$  to states back in time

**return**  $bestpath$ ,  $bestpathprob$

initialization

recursion

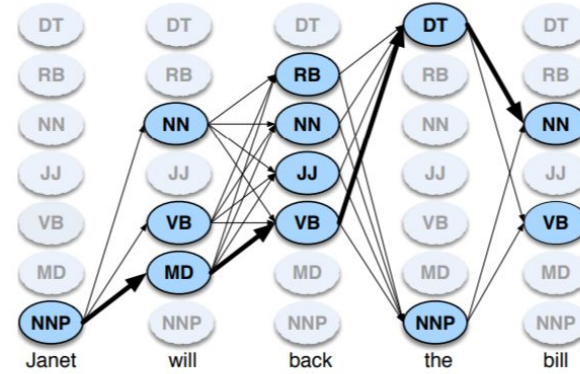
termination



# The Viterbi algorithm

|     | NNP    | MD     | VB     | JJ     | NN     | RB     | DT     |
|-----|--------|--------|--------|--------|--------|--------|--------|
| <s> | 0.2767 | 0.0006 | 0.0031 | 0.0453 | 0.0449 | 0.0510 | 0.2026 |
| NNP | 0.3777 | 0.0110 | 0.0009 | 0.0084 | 0.0584 | 0.0090 | 0.0025 |
| MD  | 0.0008 | 0.0002 | 0.7968 | 0.0005 | 0.0008 | 0.1698 | 0.0041 |
| VB  | 0.0322 | 0.0005 | 0.0050 | 0.0837 | 0.0615 | 0.0514 | 0.2231 |
| JJ  | 0.0366 | 0.0004 | 0.0001 | 0.0733 | 0.4509 | 0.0036 | 0.0036 |
| NN  | 0.0096 | 0.0176 | 0.0014 | 0.0086 | 0.1216 | 0.0177 | 0.0068 |
| RB  | 0.0068 | 0.0102 | 0.1011 | 0.1012 | 0.0120 | 0.0728 | 0.0479 |
| DT  | 0.1147 | 0.0021 | 0.0002 | 0.2157 | 0.4744 | 0.0102 | 0.0017 |

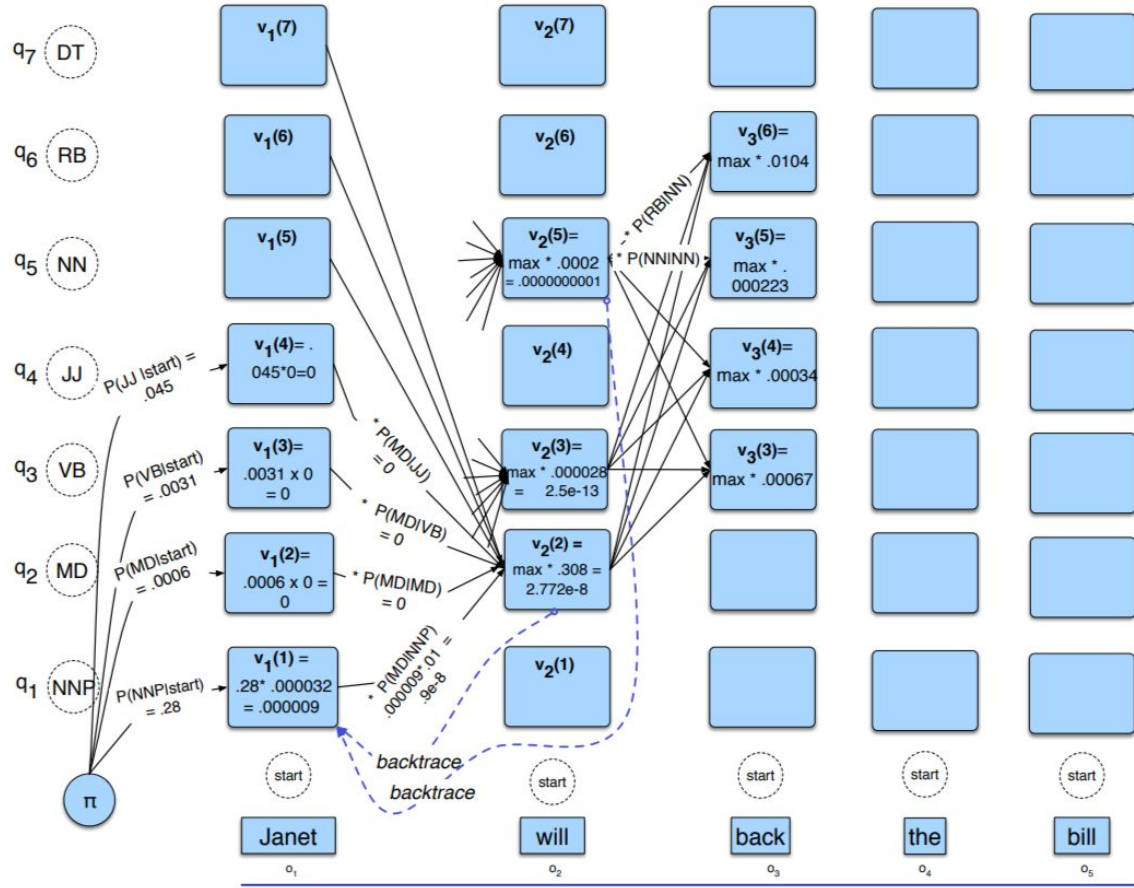
|     | Janet    | will     | back     | the      | bill     |
|-----|----------|----------|----------|----------|----------|
| NNP | 0.000032 | 0        | 0        | 0.000048 | 0        |
| MD  | 0        | 0.308431 | 0        | 0        | 0        |
| VB  | 0        | 0.000028 | 0.000672 | 0        | 0.000028 |
| JJ  | 0        | 0        | 0.000340 | 0        | 0        |
| NN  | 0        | 0.000200 | 0.000223 | 0        | 0.002337 |
| RB  | 0        | 0        | 0.010446 | 0        | 0        |
| DT  | 0        | 0        | 0        | 0.506099 | 0        |



A

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

B





# The Viterbi algorithm

**function** VITERBI(*observations* of len  $T$ , *state-graph* of len  $N$ ) **returns** *best-path*, *path-prob*

create a path probability matrix  $viterbi[N, T]$

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, 1] \leftarrow \pi_s * b_s(o_1)$$

$$backpointer[s, 1] \leftarrow 0$$

**for** each time step  $t$  **from** 2 **to**  $T$  **do**

**for** each state  $s$  **from** 1 **to**  $N$  **do**

$$viterbi[s, t] \leftarrow \max_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$backpointer[s, t] \leftarrow \operatorname{argmax}_{s'=1}^N viterbi[s', t-1] * a_{s', s} * b_s(o_t)$$

$$bestpathprob \leftarrow \max_{s=1}^N viterbi[s, T]$$

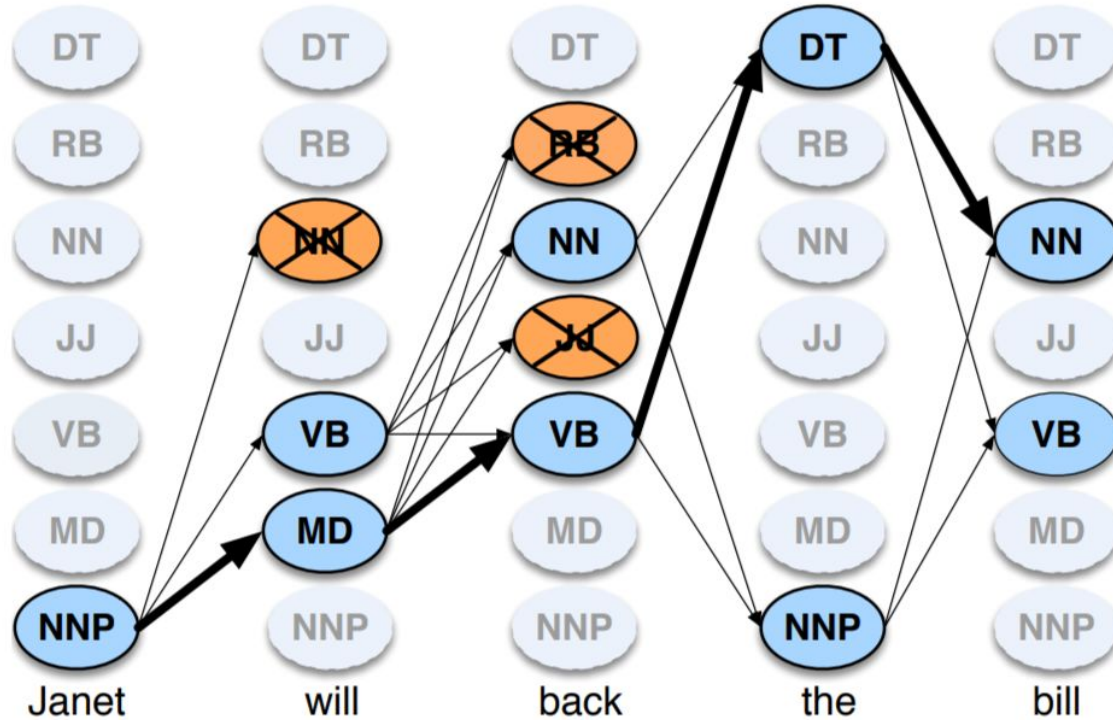
$$bestpathpointer \leftarrow \operatorname{argmax}_{s=1}^N viterbi[s, T]$$

$bestpath \leftarrow$  the path starting at state  $bestpathpointer$ , that follows  $backpointer[]$  to states back in time

**return**  $bestpath$ ,  $bestpathprob$

Computational complexity in  $N$  and  $T$ ?

# Beam search



# HMMs: algorithms

|                                 |                                |   |
|---------------------------------|--------------------------------|---|
| Forward                         | <b>Problem 1 (Likelihood):</b> | Given an HMM $\lambda = (A, B)$ and an observation sequence $O$ , determine the likelihood $P(O \lambda)$ .     |
| Viterbi                         | <b>Problem 2 (Decoding):</b>   | Given an observation sequence $O$ and an HMM $\lambda = (A, B)$ , discover the best hidden state sequence $Q$ . |
| Forward-backward;<br>Baum-Welch | <b>Problem 3 (Learning):</b>   | Given an observation sequence $O$ and the set of states in the HMM, learn the HMM parameters $A$ and $B$ .      |

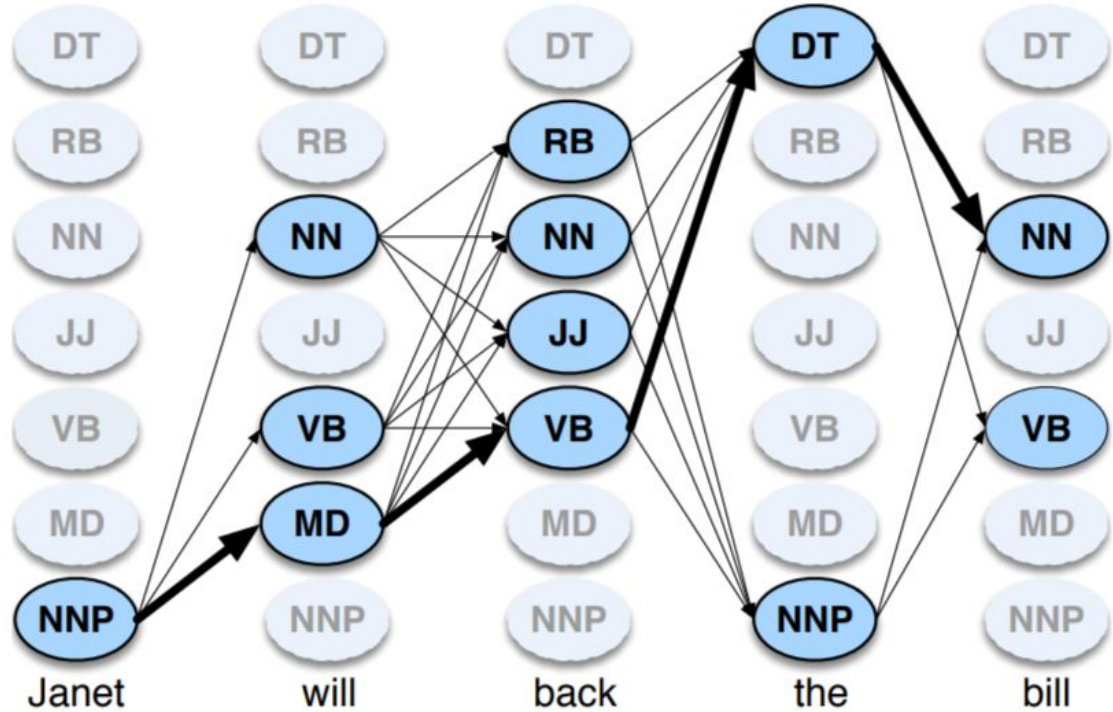
# HMMs: algorithms

|                                 |                                |   |
|---------------------------------|--------------------------------|---|
| Forward                         | <b>Problem 1 (Likelihood):</b> | Given an HMM $\lambda = (A, B)$ and an observation sequence $O$ , determine the likelihood $P(O \lambda)$ .     |
| Viterbi                         | <b>Problem 2 (Decoding):</b>   | Given an observation sequence $O$ and an HMM $\lambda = (A, B)$ , discover the best hidden state sequence $Q$ . |
| Forward-backward;<br>Baum-Welch | <b>Problem 3 (Learning):</b>   | Given an observation sequence $O$ and the set of states in the HMM, learn the HMM parameters $A$ and $B$ .      |

# The Forward algorithm

- Just sum instead of max!

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t)$$



# Viterbi

- n-best decoding
- relationship to sequence alignment

| Citation                    | Field              |
|-----------------------------|--------------------|
| Viterbi (1967)              | information theory |
| Vintsyuk (1968)             | speech processing  |
| Needleman and Wunsch (1970) | molecular biology  |
| Sakoe and Chiba (1971)      | speech processing  |
| Sankoff (1972)              | molecular biology  |
| Reichert et al. (1973)      | molecular biology  |
| Wagner and Fischer (1974)   | computer science   |