# Natural Language Processing

## Sequence labeling

Yulia Tsvetkov

yuliats@cs.washington.edu

PAUL G. ALLEN SCHOOL
OF COMPUTER SCIENCE & ENGINEERING

# Announcements

- HW 2 released
  - Start early

- HW2 content will be covered in lectures by the end of the next week

# Readings

- J&M SLP3 https://web.stanford.edu/~jurafsky/slp3/8.pdf
- Collins (2011) http://www.cs.columbia.edu/~mcollins/hmms-spring2013.pdf
- Eis 7.1-7.4, 8.1

# Levels of linguistic knowledge

| speech | text |
| phonetics | |
| | orthography |
| phonology | |

| morphology |
| lexemes |
| syntax |
| semantics |
| pragmatics |
| discourse |

"shallower"

"deeper"

| Phonetics | The study of the sounds of human language |
|---|---|
| Phonology | The study of sound systems in human language |
| Morphology | The study of the formation and internal structure of words |
| Syntax | The study of the formation and internal structure of sentences |
| Semantics | The study of the meaning of sentences |
| Pragmatics | The study of the way sentences with their semantic meanings are used for particular communicative goals |

# Factorizing solutions for linguistic analysis

- Formalism
  - map text to some representation
- Theoretical grounding from linguistics
  - why this representation?
- An algorithmic solution
  - how to solve the mapping problem?
    - Rule based
    - Supervised learning: symbolic or neural solutions
    - Unsupervised learning

# Supervised algorithms for **sequence labeling** problems

Map a sequence of words to a sequence of labels

- Part-of-speech tagging (Church, 1988; Brants, 2000)
- Named entity recognition (Bikel et al., 1999)
- Text chunking and shallow parsing (Ramshaw and Marcus, 1995)
- Word alignment of parallel text (Vogel et al., 1996)
- Compression (Conroy and O'Leary, 2001)
- Acoustic models, discourse segmentation, etc.

# Part of speech tagging

**PART OF SPEECH**  DT  VBZ  DT  JJ  NN

**WORDS**  This  is  a  simple  sentence

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

# Example of POS tagging

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**

Preliminary/JJ findings/NNS were/VBD **reported/VBN** in/IN today/NN **'s/POS** New/NNP England/NNP Journal/NNP of/IN Medicine/NNP ./.

# 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 |
|---|---|---|
| ADJ | ADP | PUNCT |
| ADV | AUX | SYM |
| INTJ | CCONJ | X |
| NOUN | DET | |
| PROPN | NUM | |
| VERB | PART | |
| | PRON | |
| | SCONJ | |

# Why POS tagging

- Goal: resolve ambiguities
- Text-to-speech
  - record, lead, protest
- Lemmatization
  - saw/V → see, saw/N → saw
- Preprocessing for harder disambiguation problems
  - syntactic parsing
  - semantic parsing

# Ambiguities in POS tags

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

# Ambiguities in POS tags

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

# Most frequent class baseline

- Assigning each token to **the class it occurred in most often** in the training set

- 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

$$\hat{y}_i = \underset{y \in \mathcal{L}}{\mathrm{argmax}}\, s(\boldsymbol{x}, i, y)$$

# Generative sequence labeling: Hidden Markov Models

# Markov Chain: weather



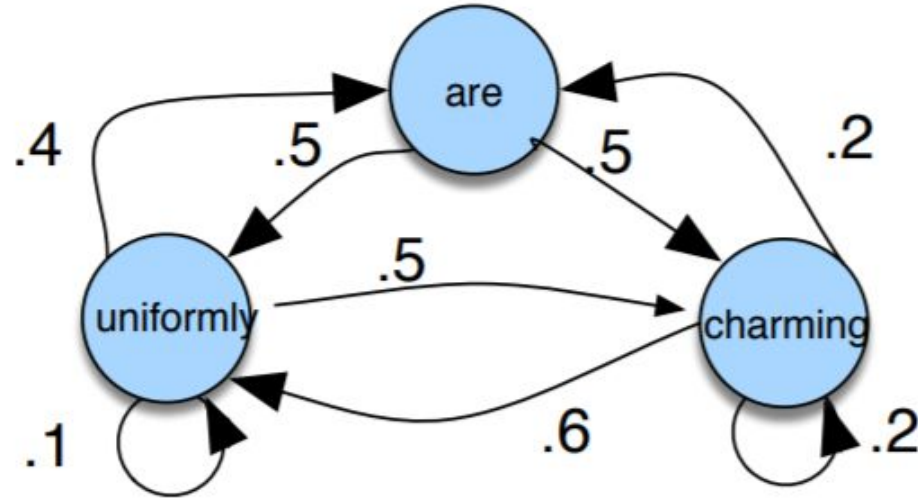**Markov Assumption:** $P(q_i = a | q_1 \ldots q_{i-1}) = P(q_i = a | q_{i-1})$

the future is independent of the past given the present

# Markov chain

Formally, a Markov chain is specified by the following components:

$Q = q_1 q_2 \ldots q_N$ — a set of $N$ **states**

$A = a_{11} a_{12} \ldots a_{n1} \ldots a_{nn}$ — a **transition probability matrix** $A$, each $a_{ij}$ representing the probability of moving from state $i$ to state $j$, s.t. $\sum_{j=1}^{n} a_{ij} = 1 \quad \forall i$

$\pi = \pi_1, \pi_2, \ldots, \pi_N$ — 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, $\sum_{i=1}^{n} \pi_i = 1$
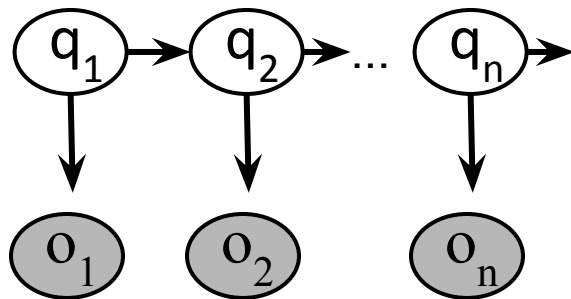
# Markov Chain: words



$$\pi = [0.1, 0.7, 0.2]$$

the future is independent of the past given the present

# Hidden Markov Models

- In real world many events are not observable
- Speech recognition: we observe acoustic features but not the phones
- POS tagging: we observe words but not the POS tags



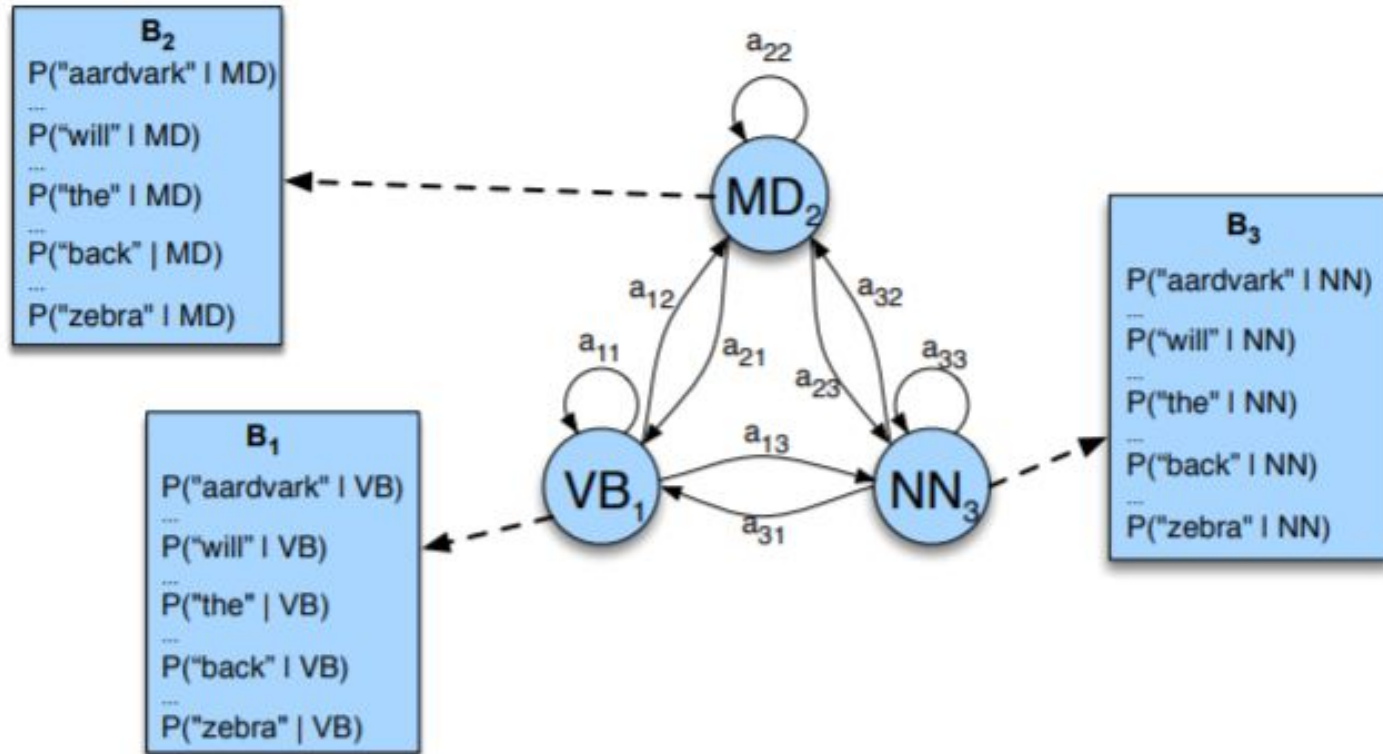**Markov Assumption:** $P(q_i|q_1 \ldots q_{i-1}) = P(q_i|q_{i-1})$

**Output Independence:** $P(o_i|q_1 \ldots q_i, \ldots, q_T, o_1, \ldots, o_i, \ldots, o_T) = P(o_i|q_i)$
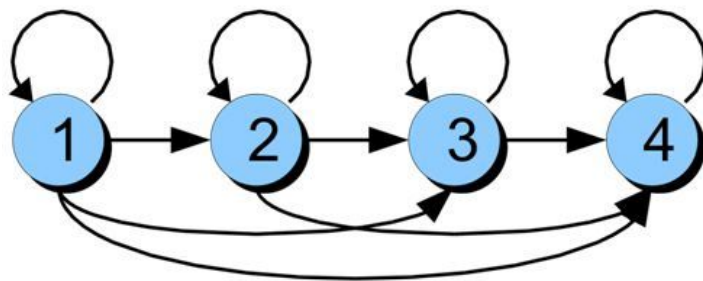
# Hidden Markov Models (HMMs)

$Q = q_1 q_2 \dots q_N$     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. $\sum_{j=1}^{N} a_{ij} = 1 \;\; \forall i$

$O = o_1 o_2 \dots o_T$     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, $\sum_{i=1}^{n} \pi_i = 1$
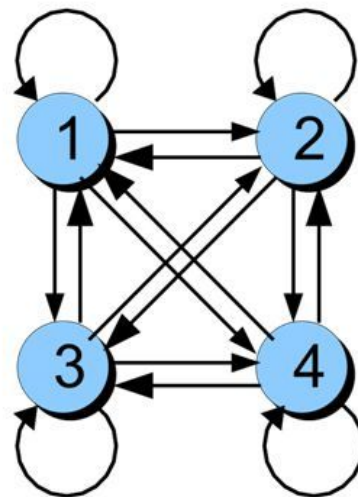
# HMM example

# 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

$Q = q_1 q_2 \ldots q_N$ — a set of $N$ **states**

→ $A = a_{11} a_{12} \ldots a_{n1} \ldots a_{nn}$ — a **transition probability matrix** $A$, each $a_{ij}$ representing the probability of moving from state $i$ to state $j$, s.t. $\sum_{j=1}^{n} a_{ij} = 1 \quad \forall i$

$O = o_1 o_2 \ldots o_T$ — a sequence of $T$ **observations**, each one drawn from a vocabulary $V = v_1, v_2, \ldots, v_V$

→ $B = b_i(o_t)$ — a sequence of **observation likelihoods**, also called **emission probabilities**, each expressing the probability of an observation $o_t$ being generated from a state $i$

→ $q_0, q_F$ — a special **start state** and **end (final) state** that are not associated with observations, together with transition probabilities $a_{01} a_{02} \ldots a_{0n}$ out of the start state and $a_{1F} a_{2F} \ldots a_{nF}$ into the end state

# HMMs: algorithms

Forward

Viterbi

Forward–
Backward;
Baum–Welch

| **Problem 1 (Likelihood):** | Given an HMM $\lambda = (A, B)$ and an observation sequence $O$, determine the likelihood $P(O\|\lambda)$. |
|---|---|
| **Problem 2 (Decoding):** | Given an observation sequence $O$ and an HMM $\lambda = (A, B)$, discover the best hidden state sequence $Q$. |
| **Problem 3 (Learning):** | Given an observation sequence $O$ and the set of states in the HMM, learn the HMM parameters $A$ and $B$. |