

Last Class:

1. EM

Today: Part-of-Speech Tagging

1. Background
2. HMM Tagger

Slide CS674–2

POS tags

“There are 10 parts of speech, and they are all troublesome.”

-*Mark Twain*

- POS tags are also known as word classes, morphological classes, or lexical tags.
- Typically much larger than Twain’s 10:
 - Penn Treebank: 45
 - Brown corpus: 87
 - C7 tagset: 146

Slide CS674–3

Penn Treebank Tagset

| Tag | Description | Example | Tag | Description | Example |
|------|-----------------------|------------------------|------|-----------------------|------------------------|
| CC | Coordin. Conjunction | <i>and, but, or</i> | SYM | Symbol | <i>+, %, &</i> |
| CD | Cardinal number | <i>one, two, three</i> | TO | “to” | <i>to</i> |
| DT | Determiner | <i>a, the</i> | UH | Interjection | <i>ah, oops</i> |
| EX | Existential ‘there’ | <i>there</i> | VB | Verb, base form | <i>eat</i> |
| FW | Foreign word | <i>mea culpa</i> | VBD | Verb, past tense | <i>ate</i> |
| IN | Preposition/sub-conj | <i>of, in, by</i> | VBG | Verb, gerund | <i>eating</i> |
| JJ | Adjective | <i>yellow</i> | VBN | Verb, past participle | <i>eaten</i> |
| JJR | Adj., comparative | <i>bigger</i> | VBP | Verb, non-3sg pres | <i>eat</i> |
| JJS | Adj., superlative | <i>wildest</i> | VBZ | Verb, 3sg pres | <i>eats</i> |
| LS | List item marker | <i>1, 2, One</i> | WDT | Wh-determiner | <i>which, that</i> |
| MD | Modal | <i>can, should</i> | WP | Wh-pronoun | <i>what, who</i> |
| NN | Noun, sing. or mass | <i>llama</i> | WP\$ | Possessive wh- | <i>whose</i> |
| NNS | Noun, plural | <i>llamas</i> | WRB | Wh-adverb | <i>how, where</i> |
| NNP | Proper noun, singular | <i>IBM</i> | \$ | Dollar sign | <i>\$</i> |
| NNPS | Proper noun, plural | <i>Carolinas</i> | # | Pound sign | <i>#</i> |
| PDT | Predeterminer | <i>all, both</i> | “ | Left quote | <i>(‘ or “)</i> |
| POS | Possessive ending | <i>’s</i> | ” | Right quote | <i>(’ or ”)</i> |
| PP | Personal pronoun | <i>I, you, he</i> | (| Left parenthesis | <i>([, { , <</i> |
| PP\$ | Possessive pronoun | <i>your, one’s</i> |) | Right parenthesis | <i>([, } , ></i> |
| RB | Adverb | <i>quickly, never</i> | , | Comma | <i>,</i> |
| RBR | Adverb, comparative | <i>faster</i> | . | Sentence-final punc | <i>(. ! ?)</i> |
| RBS | Adverb, superlative | <i>fastest</i> | : | Mid-sentence punc | <i>(: ; ... - -)</i> |
| RP | Particle | <i>up, off</i> | | | |

Slide CS674–4

Why do POS tagging?

1. Provides a lot of information about the word and its neighbors. Useful for speech recognition.
2. Can tell us something about how the word is pronounced. Useful for speech synthesis systems.
3. Can be used in IR systems...to aid stemming algorithms, to select nouns.
4. Can aid WSD algorithms.
5. Used in ASR language models, e.g. in class-based N-gram language models.
6. Critical for partial parsing algorithms.

Slide CS674–5

Part-of-Speech Tagging Baseline

How hard is p-o-s tagging?

Given word w , find the most likely tag t , i.e. find the tag that maximizes: $P(t|w)$

Maximum Likelihood Estimator: 90% accuracy rate.

To improve reliability: *need to use some of the local context.*

Slide CS674–6

Approaches

1. **rule-based**: involve a large database of hand-written disambiguation rules, e.g. that specify that an ambiguous word is a noun rather than a verb if it follows a determiner.
2. **stochastic**: resolve tagging ambiguities by using a training corpus to compute the probability of a given word having a given tag in a given context.
HMM tagger, Maximum Likelihood Tagger, Markov model tagger
3. **hybrid**: E.g. transformation-based tagger (Brill tagger); learns symbolic rules based on a corpus.
4. **ensemble methods**: combine the results of multiple taggers.

Slide CS674–7

HMM Tagger

Given $W = w_1, \dots, w_n$, find $T = t_1, \dots, t_n$ that maximizes

$$P(t_1, \dots, t_n | w_1, \dots, w_n)$$

Restate using Bayes' rule:

$$(P(t_1, \dots, t_n) * P(w_1, \dots, w_n | t_1, \dots, t_n)) / P(w_1, \dots, w_n)$$

Ignore denominator...

Make independence assumptions...

Slide CS674–8

Independence Assumptions (factor 1)

$P(t_1, \dots, t_n)$: approximate using **n-gram model**

bigram $\prod_{i=1,n} P(t_i | t_{i-1})$

trigram $\prod_{i=1,n} P(t_i | t_{i-2} t_{i-1})$

Slide CS674–9

Independence Assumptions (factor 2)

$P(w_1, \dots, w_n | t_1, \dots, t_n)$: approximate by assuming that a word appears in a category independent of its neighbors

$$\prod_{i=1,n} P(w_i | t_i)$$

Assuming bigram model:

$$P(t_1, \dots, t_n) * P(w_1, \dots, w_n | t_1, \dots, t_n) \approx$$

$$\prod_{i=1,n} P(t_i | t_{i-1}) * P(w_i | t_i)$$

Slide CS674–10

Hidden Markov Models

Equation can be modeled by an HMM.

- **states**: represent a possible lexical category
- **transition probabilities**: bigram probabilities
- **observation probabilities, lexical generation probabilities**: indicate, for each word, how likely that word is to be selected if we randomly select the category associated with the node.

Slide CS674–11

Viterbi Algorithm

c: number of lexical categories

$P(w_t | t_i)$: lexical generation probabilities

$P(t_i | t_j)$: bigram probabilities

Find most likely sequence of lexical categories T_1, \dots, T_n for word sequence.

Initialization

For i = 1 to c do

$$\text{SCORE}(i,1) = P(t_i | \phi) * P(w_1 | t_i)$$

$$\text{BPTR}(i,1) = 0$$

Slide CS674–12

Iteration

For t = 2 to n

For i = 1 to c

$$\text{SCORE}(i,t) =$$

$$\text{MAX}_{j=1,c} (\text{SCORE}(j, t-1) * P(t_i | t_j)) * P(w_t | t_i)$$

$$\text{BPTR}(i,t) = \text{index of } j \text{ that gave max}$$

Identify Sequence

$$T(n) = i \text{ that maximizes } \text{SCORE}(i,n)$$

For i = n-1 to 1 do

$$T(i) = \text{BPTR}(T(i+1), i+1)$$

Slide CS674–13

Results

- Effective if probability estimates are computed from a large corpus
- Effective if corpus is of the same style as the input to be classified
- Consistently achieve accuracies of 96% or better using trigram model
- Cuts error rate in half vs. naive algorithm (90% accuracy rate)
- Can be smoothed using backoff or deleted interpolation...

Slide CS674–14

Extensions

- Can train HMM tagger on unlabeled data using the EM algorithm, starting with a dictionary that lists which tags can be assigned to which words.
- EM then learns the word likelihood function for each tag, and the tag transition probabilities.
- Merialdo (1994) showed, however, that a tagger trained on even a small amount hand-tagged data works better than one trained via EM.

Slide CS674–15