

A Practical Part-of-Speech Tagger

Doug Cutting
Julian Kupiec
Jan Pedersen
Penelope Sibun

Introduce

- What is a part-of-speech tagger system?
 - A system that uses context to assign parts of speech to words.
- What is the requirements of POS system?
 - Robust
 - Efficient
 - Accurate
 - Tunable
 - Reusable

Introduce

- What is the advantage of the POS system proposed in this paper?
 - Few resource (only a lexicon and some unlabeled training text)
 - Facilitate higher level analysis
 - Running in linear complexity

Brief of Methodology

- Approaches used in building text taggers
 - Rule-based approach (TAGGIT)
 - Rule-based approach with finite-state machines
 - Statistical methods (Markov model)

Brief of Methodology

- Types used in training
 - Make use of a tagged training corpus
 - Manually tag a small amount of text
 - Train a partially accurate model
 - Correct the tag in the text
 - Retrain the model
 - Without using tagged training corpus
 - Forward-Backward algorithm (Hidden Markov Model)

Hidden Markov Modeling

- Definition
 - HMM = (S, X, A, B, Π)
 - State Sequence: $S = \{S_1, S_2, \dots, S_T\}$ ($S_i \in W, 1 \leq i \leq T$)
 - Observation Sequence: $X = \{X_1, X_2, \dots, X_T\}$ ($1 \leq i \leq T$)
 - Transfer Probability: $A = \{a_{ij}\}$ ($1 \leq i, j \leq N$)
where a_{ij} is the probability of moving from state i to state j
 - Output Probability: $B = \{b_{jk}\}$ ($1 \leq j \leq N, 1 \leq k \leq M, M=|W|$)
where b_{jk} is the probability of generating S_k at state j
 - $\Pi = \{\pi_i\}$ where π_i is the probability of starting in state i

Hidden Markov Modeling

- Assumption for HMM
 - Assumption 1: Markov assumption-- probability of the occurrence of word S_i at time t depends only on occurrence of word S_{i-1} at time $t-1$. $P(S_1, S_2, \dots, S_n)$ can be estimated as $\prod_{i=1}^n P(S_i | S_{i-1})$
 - Assumption 2: Output independent assumption--for all i , the S_i, X_i (observation sequence) are independent of all S_j and X_j , if $i \neq j$.
 $P(X_1, \dots, X_n | S_1, \dots, S_n)$ can be estimated as $\prod_{i=1}^n P(X_i | S_i)$

Hidden Markov Modeling



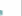



- Before going on, let's see a small example
 - Let's say there are three types of weather: sunny ☀️, rainy 🌧️, and foggy 🌫️.
 - Precondition: A guy is locked in a room, the only piece of evidence is whether the person who comes in carries an umbrella 🌂 or not.
 - Task: ask about the weather outside based on the observations.

Hidden Markov Modeling

- Example continue
 - Finding the probability of a certain weather $S \in \{\text{sunny, Rainy, foggy}\}$ can only be based on the observation X_i (whether umbrella is carried). This conditional probability $P(S_i|X_i)$ can be rewritten according to Bayes's rule:
 - $P(S_i|X_i) = P(X_i|S_i)P(S_i)/P(X_i)$
 - Because X_i is a known observation sequence, $P(X_i) = 1$, then $P(S_i|X_i) = P(X_i|S_i)P(S_i) = \prod_{i=1}^n P(S_i|S_{i-1}) \prod_{i=1}^n P(X_i|S_i)$

Hidden Markov Modeling

- Example continue
 - Task: Calculate the likelihood for the weather on these three days to have been $\{S_1=\text{sunny}, S_2=\text{foggy}, S_3=\text{sunny}\}$ and the person comes in without an umbrella (F).
 - $P(S_1=\text{sunny}, S_2=\text{foggy}, S_3=\text{sunny} | X_1=F, X_2=F, X_3=F)$

Today's weather	Tomorrow's weather		
			
	0.8	0.05	0.15
	0.2	0.6	0.2
	0.2	0.3	0.5

Weather	Probability of umbrella
Sunny	0.1
Rainy	0.8
Foggy	0.3

Hidden Markov Modeling

- Example continue

$P(S_1=\text{sunny}, S_2=\text{foggy}, S_3=\text{sunny} | X_1=F, X_2=F, X_3=F)$
 $= P(X_1=F | S_1=\text{sunny}) * P(X_2=F | S_2=\text{foggy}) * P(X_3=F | S_3=\text{sunny}) * P(S_3=\text{sunny} | S_2=\text{foggy}) * P(S_2=\text{foggy} | S_1=\text{sunny}) * P(S_1=\text{sunny})$
 $= 0.9 * 0.7 * 0.9 * 1/3 * 0.15 * 0.2$
 $= 0.0057$

So, The probability of {sunny, foggy, sunny} is 0.0057

Hidden Markov Modeling

- HMM in speech recognition
 - Task
 - Estimate the model parameters A , B and Π from a training set
 - Find the most likely sequence of words W given some acoustic input.

Hidden Markov Modeling

- HMM in speech recognition

- Probability of the entire sequence

- Forward probabilities: $\alpha_{t+1}(j) = \left[\sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(S_{t+1}) \quad 1 \leq t \leq T-1,$

- Backward probabilities: $\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(S_{t+1}) \beta_{t+1}(j) \quad (1 \leq t \leq T-1, !)$

- $P = \sum_i \alpha_t(i) \beta_t(i) = \sum_j \alpha_{t+1}(j) \beta_{t+1}(j)$
 $= \sum_i \alpha_t(i) (\sum_j a_{ij} b_j(S_{t+1}) \beta_{t+1}(j)) = \sum_j (\sum_i \alpha_t(i) a_{ij} b_j(S_{t+1})) \beta_{t+1}(j)$
 $= \sum_{i=1}^N \sum_{j=1}^N \alpha_t(i) a_{ij} b_j(S_{t+1}) \beta_{t+1}(j)$

Hidden Markov Modeling

- HMM in speech recognition

Hence we can estimate a_{ij} by:

$$\hat{a}_{ij} = \frac{\gamma_{ij}}{\sum_{j=1}^N \gamma_{ij}} = \frac{\sum_{t=1}^{T-1} \alpha_t(i) a_{ij} b_j(S_{t+1}) \beta_{t+1}(j)}{\sum_{t=1}^{T-1} \alpha_t(i) \beta_t(i)}. \quad (3)$$

Similarly, b_{jk} and π_i can be estimated as follows:

$$\hat{b}_{jk} = \frac{\sum_{t: S_t = k} \alpha_t(j) \beta_t(j)}{\sum_{t=1}^T \alpha_t(j) \beta_t(j)} \quad (4)$$

and

$$\hat{\pi}_i = \frac{1}{P} \alpha_1(i) \beta_1(i). \quad (5)$$

Hidden Markov Modeling

- Viterbi Algorithm

To find the most probable such sequence we start by defining $\phi_1(i) = \pi_i b_i(S_1)$ for $1 \leq i \leq N$ and then perform the recursion

$$\phi_t(j) = \max_{1 \leq i \leq N} [\phi_{t-1}(i) a_{ij}] b_j(S_t) \quad (6)$$

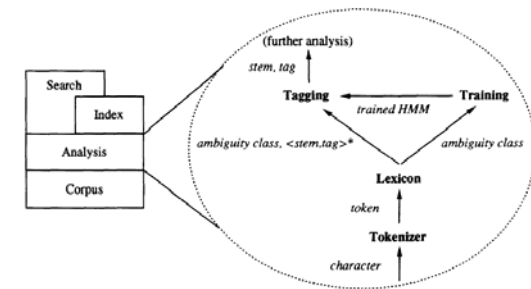
and

$$\psi_t(j) = \max_{1 \leq i \leq N} [\phi_{t-1}(i) a_{ij}]$$

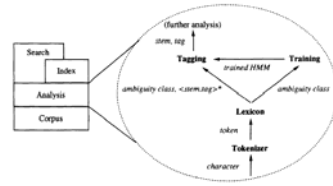
Notice: According to other reference book, $\Psi_t(j) = \max_{i,j} [\phi_{t-1}(i) a_{ij}]$ rather than $\Psi_t(j) = \max_i \phi_{t-1}(i)$

Architecture

- Tagger Modules in System Context



Architecture



- Five components
 - Term: a word stem annotated with part-of-speech
 - Corpus: provides text in a generic manner
 - Analysis: extracts terms from the text
 - Index: stores term occurrence statistics
 - Search: utilizes statistics to resolve queries

Architecture

- Tokenizer vs. Lexicon Implementation
 - Two basic classes of Tokenizer
 - Sentence boundary
 - Word
 - Other classes: numbers, paragraph boundaries
 - Three Stage of Lexicon
 - Construct lexicon manually
 - Guess ambiguity classes automatically (suffix)
 - Use default ambiguity classes

Evaluation

- Time Complexity: $O(kTN)$
- Space Complexity: $T+2N(T+N+M+1)$
- Accuracy: 96%
- Robust: fragments and ungrammaticalities
- Tunable: support tuning itself
(various tagset, lexicon / empirical and prior information)

Application

- Phrase Recognition
- Word Sense Disambiguation
- Grammatical Function Assignment

Reference

- Barbara Resch, Hidden Markov Models--
A Tutorial for the Courses Computational
Intelligence, <http://www.igi.tugraz.at/lehre/CI>

