

Maximum Entropy Models

Pawan Goyal

CSE, IIT Kharagpur

Week 4, Lecture 4

Practice Question

Consider the maximum entropy model for POS tagging, where you want to estimate $P(\text{tag}|\text{word})$. In a hypothetical setting, assume that *tag* can take the values *D*, *N* and *V* (short forms for Determiner, Noun and Verb). The variable *word* could be any member of a set *V* of possible words, where *V* contains the words *a*, *man*, *sleeps*, as well as additional words. The distribution should give the following probabilities

- $P(D|a) = 0.9$
- $P(N|man) = 0.9$
- $P(V|sleeps) = 0.9$
- $P(D|\text{word}) = 0.6$ for any word other than *a*, *man* or *sleeps*
- $P(N|\text{word}) = 0.3$ for any word other than *a*, *man* or *sleeps*
- $P(V|\text{word}) = 0.1$ for any word other than *a*, *man* or *sleeps*

It is assumed that all other probabilities, not defined above could take any values such that

$\sum_{\text{tag}} P(\text{tag}|\text{word}) = 1$ is satisfied for any word in *V*.

- Define the features of your maximum entropy model that can model this distribution. Mark your features as f_1, f_2 and so on. Each feature should have the same format as explained in the class. **[Hint: 6 Features should make the analysis easier]**
- For each feature f_i , assume a weight λ_i . Now, write expression for the following probabilities in terms of your model parameters
 - ▶ $P(D|cat)$
 - ▶ $P(N|laughs)$
 - ▶ $P(D|man)$
- What value do the parameters in your model take to give the distribution as described above. (i.e. $P(D|a) = 0.9$ and so on. *You may leave the final answer in terms of equations*)

Features for POS Tagging (Ratnaparakhi, 1996)

The specific word and tag context available to a feature is

$$h_i = \{w_i, w_{i+1}, w_{i+2}, w_{i-1}, w_{i-2}, t_{i-1}, t_{i-2}\}$$

Features for POS Tagging (Ratnaparakhi, 1996)

The specific word and tag context available to a feature is

$$h_i = \{w_i, w_{i+1}, w_{i+2}, w_{i-1}, w_{i-2}, t_{i-1}, t_{i-2}\}$$

Example: $f_j(h_i, t_i) = 1$ if $\text{suffix}(w_i) = \text{"ing"} \& t_i = \text{VBG}$

Example Features

Condition	Features
w_i is not rare	$w_i = X \quad \& \quad t_i = T$
w_i is rare	X is prefix of w_i , $ X \leq 4 \quad \& \quad t_i = T$
	X is suffix of w_i , $ X \leq 4 \quad \& \quad t_i = T$
	w_i contains number $\quad \& \quad t_i = T$
	w_i contains uppercase character $\quad \& \quad t_i = T$
	w_i contains hyphen $\quad \& \quad t_i = T$
$\forall w_i$	$t_{i-1} = X \quad \& \quad t_i = T$
	$t_{i-2}t_{i-1} = XY \quad \& \quad t_i = T$
	$w_{i-1} = X \quad \& \quad t_i = T$
	$w_{i-2} = X \quad \& \quad t_i = T$
	$w_{i+1} = X \quad \& \quad t_i = T$
	$w_{i+2} = X \quad \& \quad t_i = T$

Example Features

<i>Word:</i>	the	stories	about	well-heeled	communities	and	developers
<i>Tag:</i>	DT	NNS	IN	JJ	NNS	CC	NNS
<i>Position:</i>	1	2	3	4	5	6	7

Example Features

Word:	the	stories	about	well-heeled	communities	and	developers
Tag:	DT	NNS	IN	JJ	NNS	CC	NNS
Position:	1	2	3	4	5	6	7

$w_i = \text{about}$ & $t_i = \text{IN}$
 $w_{i-1} = \text{stories}$ & $t_i = \text{IN}$
 $w_{i-2} = \text{the}$ & $t_i = \text{IN}$
 $w_{i+1} = \text{well-heeled}$ & $t_i = \text{IN}$
 $w_{i+2} = \text{communities}$ & $t_i = \text{IN}$
 $t_{i-1} = \text{NNS}$ & $t_i = \text{IN}$
 $t_{i-2}t_{i-1} = \text{DT NNS}$ & $t_i = \text{IN}$

$w_{i-1} = \text{about}$ & $t_i = \text{JJ}$
 $w_{i-2} = \text{stories}$ & $t_i = \text{JJ}$
 $w_{i+1} = \text{communities}$ & $t_i = \text{JJ}$
 $w_{i+2} = \text{and}$ & $t_i = \text{JJ}$
 $t_{i-1} = \text{IN}$ & $t_i = \text{JJ}$
 $t_{i-2}t_{i-1} = \text{NNS IN}$ & $t_i = \text{JJ}$
 $\text{prefix}(w_i) = \text{w}$ & $t_i = \text{JJ}$
 $\text{prefix}(w_i) = \text{we}$ & $t_i = \text{JJ}$
 $\text{prefix}(w_i) = \text{wel}$ & $t_i = \text{JJ}$
 $\text{prefix}(w_i) = \text{well}$ & $t_i = \text{JJ}$
 $\text{suffix}(w_i) = \text{d}$ & $t_i = \text{JJ}$
 $\text{suffix}(w_i) = \text{ed}$ & $t_i = \text{JJ}$
 $\text{suffix}(w_i) = \text{led}$ & $t_i = \text{JJ}$
 $\text{suffix}(w_i) = \text{eled}$ & $t_i = \text{JJ}$
 $w_i \text{ contains hyphen}$ & $t_i = \text{JJ}$

Conditional Probability

Given a sentence $\{w_1, \dots, w_n\}$, a tag sequence candidate $\{t_1, \dots, t_n\}$ has conditional probability:

$$P(t_1, \dots, t_n | w_1 \dots, w_n) = \prod_{i=1}^n p(t_i | x_i)$$

Conditional Probability

Given a sentence $\{w_1, \dots, w_n\}$, a tag sequence candidate $\{t_1, \dots, t_n\}$ has conditional probability:

$$P(t_1, \dots, t_n | w_1 \dots, w_n) = \prod_{i=1}^n p(t_i | x_i)$$

A *Tag Dictionary* is used, which, for each known word, lists the tags that it has appeared with in the training set.

Search Algorithm

Let $W = \{w_1, \dots, w_n\}$ be a test sentence, s_{ij} be the j th highest probability tag sequence up to and including word w_i .

Search Algorithm

Let $W = \{w_1, \dots, w_n\}$ be a test sentence, s_{ij} be the j th highest probability tag sequence up to and including word w_i .

Search description

- Generate tags for w_1 , find top N , set $s_{1j}, 1 \leq j \leq N$, accordingly.

Search Algorithm

Let $W = \{w_1, \dots, w_n\}$ be a test sentence, s_{ij} be the j th highest probability tag sequence up to and including word w_i .

Search description

- Generate tags for w_1 , find top N , set s_{1j} , $1 \leq j \leq N$, accordingly.
- Initialize $i = 2$
 - ▶ Initialize $j = 1$
 - ▶ Generate tags for w_i , given $s_{(i-1)j}$ as previous tag context, and append each tag to $s_{(i-1)j}$ to make a new sequence
 - ▶ $j = j + 1$, repeat if $j \leq N$

Search Algorithm

Let $W = \{w_1, \dots, w_n\}$ be a test sentence, s_{ij} be the j th highest probability tag sequence up to and including word w_i .

Search description

- Generate tags for w_1 , find top N , set s_{1j} , $1 \leq j \leq N$, accordingly.
- Initialize $i = 2$
 - ▶ Initialize $j = 1$
 - ▶ Generate tags for w_i , given $s_{(i-1)j}$ as previous tag context, and append each tag to $s_{(i-1)j}$ to make a new sequence
 - ▶ $j = j + 1$, repeat if $j \leq N$
- Find N highest probability sequences generated by above loop, set s_{ij} accordingly

Search Algorithm

Let $W = \{w_1, \dots, w_n\}$ be a test sentence, s_{ij} be the j th highest probability tag sequence up to and including word w_i .

Search description

- Generate tags for w_1 , find top N , set $s_{1j}, 1 \leq j \leq N$, accordingly.
- Initialize $i = 2$
 - ▶ Initialize $j = 1$
 - ▶ Generate tags for w_i , given $s_{(i-1)j}$ as previous tag context, and append each tag to $s_{(i-1)j}$ to make a new sequence
 - ▶ $j = j + 1$, repeat if $j \leq N$
- Find N highest probability sequences generated by above loop, set s_{ij} accordingly
- $i = i + 1$, repeat if $i \leq n$

Search Algorithm

Let $W = \{w_1, \dots, w_n\}$ be a test sentence, s_{ij} be the j th highest probability tag sequence up to and including word w_i .

Search description

- Generate tags for w_1 , find top N , set $s_{1j}, 1 \leq j \leq N$, accordingly.
- Initialize $i = 2$
 - ▶ Initialize $j = 1$
 - ▶ Generate tags for w_i , given $s_{(i-1)j}$ as previous tag context, and append each tag to $s_{(i-1)j}$ to make a new sequence
 - ▶ $j = j + 1$, repeat if $j \leq N$
- Find N highest probability sequences generated by above loop, set s_{ij} accordingly
- $i = i + 1$, repeat if $i \leq n$
- Return highest probability sequence s_{n1}

A Good Reference

Berger et al., *A Maximum Entropy Approach to Natural Language Processing*, Computational Linguistics, Vol. 22, No. 1.