

CS4740/5740
Introduction to Natural Language Processing
Final Exam Solutions
May 13, 2014

1 Parsing

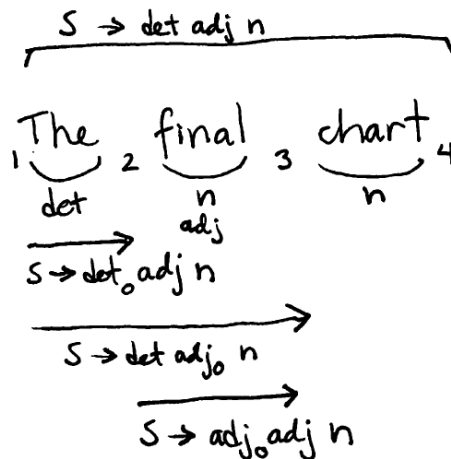


Figure 1: Final bottom-up chart.

1. See Figure 1.
2. See Figure 2.
3. The Earley algorithm incorporates (limited) top-down predictions to prevent the consideration of any complete edge (i.e., a part-of-speech for a token, a constituent type for a phrase) unless it was expected to occur at that position in the sentence (according to the top-down edges introduced into the chart).
4. (a) The specific probabilities don't matter; for each constituent type, they just need to add to 1.

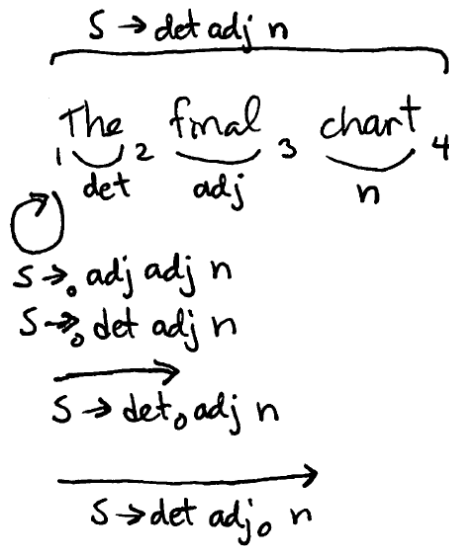


Figure 2: Final Earley chart.

production rule	probability
$S \rightarrow VP$	1.0
$VP \rightarrow \text{Verb NP}$	0.7
$VP \rightarrow \text{Verb NP PP}$	0.3
$NP \rightarrow \text{NP PP}$	0.3
$NP \rightarrow \text{Det Noun}$	0.7
$PP \rightarrow \text{Prep Noun}$	1.0

(b) We need a different PP attachment in each sentence, but the grammar and parser force the PP attachment to be the same in each case. Semantically, “with the binoculars” should attach to (i.e., modify) the noun phrase “the man” in the second sentence (since you can’t use binoculars to call someone); in the first sentence, the PP attachment is ambiguous — the sentence’s interpretation would depend on the context — but most readers prefer the reading that attaches “with the binoculars” to the verb (i.e., saw with the binoculars).

(c) Lexicalization of production rules can capture lexical specific preference of certain rule expansions. In order to mitigate the sparse data problem, we will lexicalize with respect to the head word of the left hand side of each production rule, instead of all nonterminals in each

production rule. In particular, the rules expending from VP should be modified as

production rule	probability
VP (x) \rightarrow Verb NP	(p_x)
VP (x) \rightarrow Verb NP PP	(q_x)

where

$x \in \{ \text{Cut, Ask, Find, ... } \},$

$p_x \stackrel{\text{def}}{=} P(\text{VP} (x) \rightarrow \text{Verb NP} \mid \text{VP}, x),$

$q_x \stackrel{\text{def}}{=} P(\text{VP} (x) \rightarrow \text{Verb NP PP} \mid \text{VP}, x),$

and $p_x + q_x = 1.$

Comment

Because we didn't restrict lexicalization to head words of the right hand side of rules, it is okay to propose lexicalized PCFGs in many different ways; in particular, you don't have to condition on the head word, you can condition on the entire combination of words for all non-terminals, as long as you made it clear what you are conditioning on, although it would be much less practical.

2 HMMs and Part-of-Speech Tagging

- Given that w_i is the i th word, and t_i is the i th part-of-speech tag, what we want here is $P(t_i|t_{i-1})$, estimated as $\frac{\text{count-of}(t_i, t_{i-1})}{\text{count-of}(t_{i-1})}$

$$P(V|N) = \frac{3}{9}$$

$$P(N|ADJ) = \frac{4}{4}$$

$$P(,|N) = \frac{2}{9}$$

$$P(.|N) = \frac{3}{9}$$

$$P(N|V) = \frac{1}{3}$$

$$P(N|N) = \frac{1}{9}$$

$$P(N|DET) = \frac{1}{3}$$

If you don't include sentence-boundary tokens, then $P(N|.) = \frac{1}{3}$

If you do, then $P(N| < s >) = \frac{2}{3}$

Many students gave bigram probabilities for words, either as the joint $P(\text{current} - \text{word} \wedge \text{previous} - \text{word})$ or as the conditional $P(\text{current} - \text{word} | \text{previous} - \text{word})$. This might be reasonable in language modeling, but for HMMs we only want bigram probabilities for part-of-speech tags. Otherwise the probabilities cannot be combined into a single probability for the sentence and tags. Also note that joint probabilities cannot be combined at all - e.g. $P(w_0, w_1) * P(w_1, w_2)$ isn't a sensible probability.

- What we want here is $P(w_i | t_i)$, estimated as $\frac{\text{count-of}(w_i, t_i)}{\text{count-of}(t_i)}$

$$P(\text{Mary} | N) = \frac{2}{9}$$

$$P(\text{lamb} | N) = \frac{4}{9}$$

$$P(\text{fleece} | N) = \frac{1}{9}$$

$$P(\text{snow} | N) = \frac{1}{9}$$

$$P(\text{white} | N) = \frac{1}{9}$$

You can optionally include 0/9 for all other words.

1. One VERY simple baseline is to select the most frequent part-of-speech for a particular word type.
2. $10 \times 40^2 = 16,000$
3. See Figure 3.

3 Information Extraction

1. See Figure 4.
2. See Figure 4.
3. See Figure 5 for one possible extraction pattern.
4. Noun phrase coreference resolution is important for IE because it allows information to be extracted w.r.t. a particular entity in sentences or clauses where it is not mentioned explicitly. It is also important for allowing all of the information extracted w.r.t. an entity or event to be appropriately merged even when it is spread throughout a document.

$$\begin{aligned}
S_N(w_3) &= \max \left\{ \begin{array}{l} S_N(w_2) \times P(N|N) \\ S_V(w_2) \times P(N|V) \\ S_D(w_2) \times P(N|D) \end{array} \right\} \times P(\text{eat} | N) \\
b_N(w_3) &= \text{argmax}_{t \in \{N, V, D\}} S_N(w_3) \\
S_V(w_3) &= \max \left\{ \begin{array}{l} S_N(w_2) \times P(V|N) \\ S_V(w_2) \times P(V|V) \\ S_D(w_2) \times P(V|D) \end{array} \right\} \times P(\text{eat} | V) \\
b_V(w_3) &= \text{argmax}_{t \in \{N, V, D\}} S_V(w_3) \\
S_D(w_3) &= \max \left\{ \begin{array}{l} S_N(w_2) \times P(D|N) \\ S_V(w_2) \times P(D|V) \\ S_D(w_2) \times P(D|D) \end{array} \right\} \times P(\text{eat} | D) \\
b_D(w_3) &= \text{argmax}_{t \in \{N, V, D\}} S_D(w_3)
\end{aligned}$$

Figure 3: Viterbi Calculations

^{COMPANY}
AT&T is in talks to buy DirectTV for at least
 \$50 billion, and the two sides are actively
 working toward an announcement. If completed,
 a deal would give AT&T the country's second-largest
wireless carrier, control of the country's largest
satellite television provider. AT&T has grown interested
 in DirectTV in recent months because of its 20 million
^{LOCATION}
U.S. subscribers.

optional

Figure 4: NEs and NP coref chains.

5. Some examples: (1) "the country's largest satellite television provider" = DirecTV. This would allow the system to know that DirecTV is the

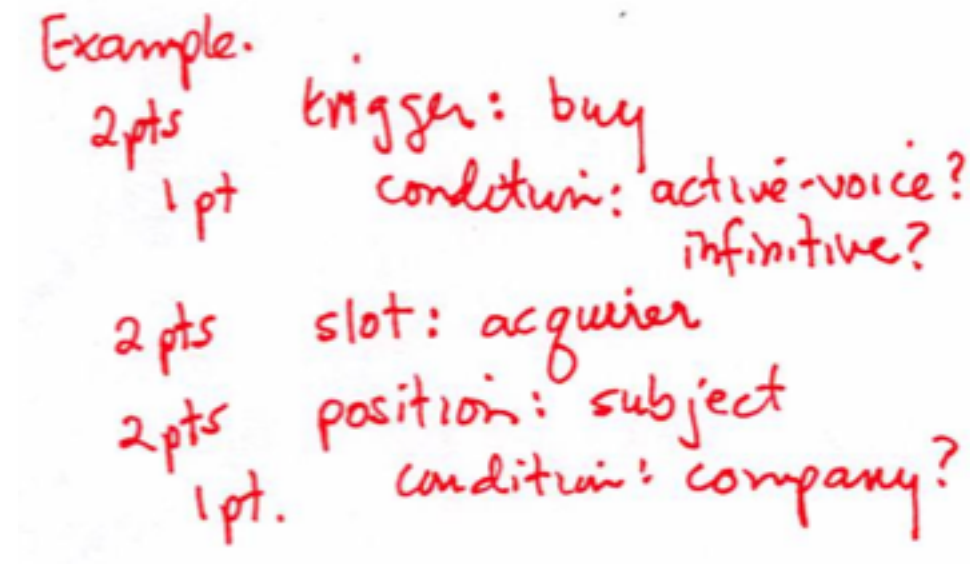


Figure 5: Extraction pattern.

entity being controlled in: “control of the country’s largest satellite television provider”.

(2) “its” in “its 20 million U.S. subscribers” refers to DirecTV. This would potentially allow the IE system to identify the reasons for the acquisition.

4 Grab Bag

1. False. Zipf’s law relates a term’s frequency to its rank in a (large) corpus, but does not state anything about the difference in term frequencies in different corpora.
2. False. WordNet is fundamentally arranged around the synset, not the lexeme. A lexeme is one surface form with several related senses, whereas one synset is one sense with several possible surface forms.
3. Language modeling can be used in language identification by cre-

ating a language model for each of the target languages. Remember that an n-gram language model gives the probability of a word given the n-1 preceding words – $P(x_i|x_{i-1}, x_{i-2}, \dots, x_{i-(n-1)})$. A language model can be used to compute the probability of a sequence of words by combining the probabilities of individual words by multiplying, adding, or using another form of combination. This way we can compute the probability of a sequence of words (text) as

$$P(x_0, x_1, \dots, x_k) = P(x_0|x_{-1}, x_{-2}, \dots, x_{-(n-1)}) \cdot P(x_1|x_0, x_{-1}, \dots, x_{-n}) \dots P(x_n|x_{k-1}, x_{k-2}, \dots, x_{k-(n-1)})$$

for a specific language model.

Using this formula, and given text and a language model for each of the target languages we can compute the probability that each of the language models assigns for the text and decide that the text is from the language associated with the language model that gives the highest probability.

For training we will need a set of (untagged) texts from each of the target languages.

Note that if we decide to use unigram language models, this is equivalent to counting the relative frequency of words in each of the target languages and using this relative frequency as the probability of each word.