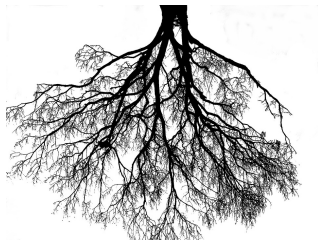


Probabilistic Context-free Grammars and Other Syntactic Language Models



Jon Dehdari

January 4, 2016

Kelsey and other Grammers

- A **grammar** here is another word for a language model
- They consist of four sets $G = \langle \Sigma, N, S, P \rangle$
 - terminals** – word types; lowest nodes in syntax trees
Examples: *dog, the, eats*
 - non-terminals** – phrasal types; middle nodes in syntax trees
Examples: *VP, DET, NP*
 - start symbol** – “S”; the top node in syntax trees

Kelsey and other Grammars

- A **grammar** here is another word for a language model
- They consist of four sets $G = \langle \Sigma, N, S, P \rangle$
 - terminals** – word types; lowest nodes in syntax trees
Examples: *dog, the, eats*
 - non-terminals** – phrasal types; middle nodes in syntax trees
Examples: *VP, DET, NP*
 - start symbol** – “S”; the top node in syntax trees
 - production rules** – recursive symbol substitutions

Examples:

$S \rightarrow NP VP$

$NP \rightarrow DET N$

$NP \rightarrow ADJ N$

$VP \rightarrow V NP$

$VP \rightarrow V$

$N \rightarrow dog$

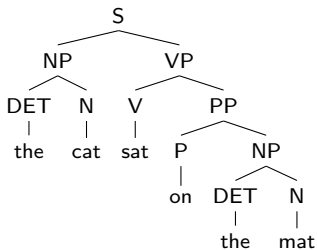
$N \rightarrow cat$

$V \rightarrow barks$

$DET \rightarrow the$

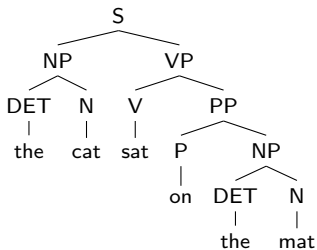
Visualization

- Sentences are often visualized using **derivation trees**, also known as **parse trees** or **syntax trees**
- Example:



Visualization

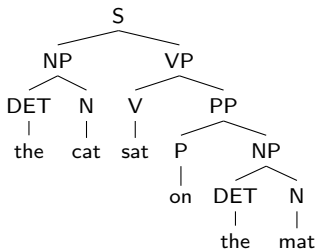
- Sentences are often visualized using **derivation trees**, also known as **parse trees** or **syntax trees**
- Example:



<i>S</i>	→	<i>NP VP</i>
<i>NP</i>	→	<i>DET N</i>
<i>DET</i>	→	<i>the</i>
<i>N</i>	→	<i>cat</i>
<i>VP</i>	→	<i>V PP</i>
<i>V</i>	→	<i>sat</i>
<i>PP</i>	→	<i>P NP</i>
<i>N</i>	→	<i>mat</i>

Visualization

- Sentences are often visualized using **derivation trees**, also known as **parse trees** or **syntax trees**
- Example:

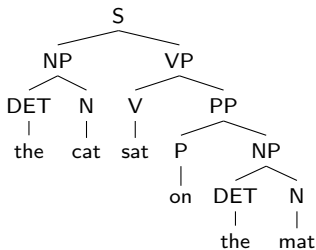


<i>S</i>	→	<i>NP VP</i>
<i>NP</i>	→	<i>DET N</i>
<i>DET</i>	→	<i>the</i>
<i>N</i>	→	<i>cat</i>
<i>VP</i>	→	<i>V PP</i>
<i>V</i>	→	<i>sat</i>
<i>PP</i>	→	<i>P NP</i>
<i>N</i>	→	<i>mat</i>

- Originally these trees were **mere visualizations** of how you could generate a grammatical sentence, given a grammar

Visualization

- Sentences are often visualized using **derivation trees**, also known as **parse trees** or **syntax trees**
- Example:

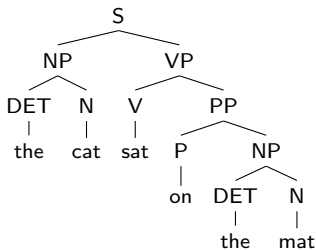


<i>S</i>	→	<i>NP VP</i>
<i>NP</i>	→	<i>DET N</i>
<i>DET</i>	→	<i>the</i>
<i>N</i>	→	<i>cat</i>
<i>VP</i>	→	<i>V PP</i>
<i>V</i>	→	<i>sat</i>
<i>PP</i>	→	<i>P NP</i>
<i>N</i>	→	<i>mat</i>

- Originally these trees were **mere visualizations** of how you could generate a grammatical sentence, given a grammar
- Then people started to think of these trees as the actual **structure** of a sentence

Visualization

- Sentences are often visualized using **derivation trees**, also known as **parse trees** or **syntax trees**
- Example:



<i>S</i>	→	<i>NP VP</i>
<i>NP</i>	→	<i>DET N</i>
<i>DET</i>	→	<i>the</i>
<i>N</i>	→	<i>cat</i>
<i>VP</i>	→	<i>V PP</i>
<i>V</i>	→	<i>sat</i>
<i>PP</i>	→	<i>P NP</i>
<i>N</i>	→	<i>mat</i>

- Originally these trees were **mere visualizations** of how you could generate a grammatical sentence, given a grammar
- Then people started to think of these trees as the actual **structure** of a sentence
- Confusion ensued

Context-free Grammars

- A **context-free grammar** (CFG) is a generative model that can generate context-free languages, which are somewhere in the middle of the formal language hierarchy
- Many, but not all, phenomena in natural languages can be generated by CFGs

Context-free Grammars

- A **context-free grammar** (CFG) is a generative model that can generate context-free languages, which are somewhere in the middle of the formal language hierarchy
- Many, but not all, phenomena in natural languages can be generated by CFGs
- Context-free production rules have the general form of a non-terminal rewriting to a sequence (string) of terminals and/or non-terminals ($A \rightarrow \alpha$)

Context-free Grammars

- A **context-free grammar** (CFG) is a generative model that can generate context-free languages, which are somewhere in the middle of the formal language hierarchy
- Many, but not all, phenomena in natural languages can be generated by CFGs
- Context-free production rules have the general form of a non-terminal rewriting to a sequence (string) of terminals and/or non-terminals ($A \rightarrow \alpha$)
- CFGs can generate and recognize **center embedding**, but not more complex word order phenomena, so effectively CFG parse trees have **no crossing lines**

Context-free Grammars

- A **context-free grammar** (CFG) is a generative model that can generate context-free languages, which are somewhere in the middle of the formal language hierarchy
- Many, but not all, phenomena in natural languages can be generated by CFGs
- Context-free production rules have the general form of a non-terminal rewriting to a sequence (string) of terminals and/or non-terminals ($A \rightarrow \alpha$)
- CFGs can generate and recognize **center embedding**, but not more complex word order phenomena, so effectively CFG parse trees have **no crossing lines**
- Non-projective dependency grammars are more or less equivalent to CFGs (they have the same weak generative capacity)

Treebanks

- It's a lot of work to define a language model by hand (including context-free grammars), so another way is to annotate treebanks
- Example: (S (NP (DET the) (N cat))(VP (V sat)(PP (P on)(NP (DET the) (N mat))))))

Trebanks

- It's a lot of work to define a language model by hand (including context-free grammars), so another way is to annotate treebanks
- Example: (S (NP (DET the) (N cat))(VP (V sat)(PP (P on)(NP (DET the) (N mat))))))
- There are treebanks for about 10–20 languages, the Penn Treebank being the most well-known for English

Treebanks

- It's a lot of work to define a language model by hand (including context-free grammars), so another way is to annotate treebanks
- Example: (S (NP (DET the) (N cat))(VP (V sat)(PP (P on)(NP (DET the) (N mat))))))
- There are treebanks for about 10–20 languages, the Penn Treebank being the most well-known for English
- Treebanks can be annotated with various grammatical annotations, like **constituency** / **phrase-structure** (as above), **dependency grammar** (as we saw last class), categorial grammar, HPSG, etc.
- Most of these annotation styles can be approximately mapped to other styles

Treebanks

- It's a lot of work to define a language model by hand (including context-free grammars), so another way is to annotate treebanks
- Example: (S (NP (DET the) (N cat))(VP (V sat)(PP (P on)(NP (DET the) (N mat))))))
- There are treebanks for about 10–20 languages, the Penn Treebank being the most well-known for English
- Treebanks can be annotated with various grammatical annotations, like **constituency** / **phrase-structure** (as above), **dependency grammar** (as we saw last class), categorial grammar, HPSG, etc.
- Most of these annotation styles can be approximately mapped to other styles
- Here is a link to a list of syntactic treebanks

PCFGs

- We can induce a **probabilistic context-free grammar** (PCFG) from the treebank
- With multiple annotated sentences, we can get probabilities for production rules. Example:

1.0	<i>S</i>	→ <i>NP VP</i>
0.6	<i>NP</i>	→ <i>DET N</i>
0.4	<i>NP</i>	→ <i>ADJ N</i>
0.7	<i>VP</i>	→ <i>V NP</i>
0.3	<i>VP</i>	→ <i>V</i>
0.8	<i>N</i>	→ <i>dog</i>
0.2	<i>N</i>	→ <i>cat</i>
1.0	<i>V</i>	→ <i>barks</i>
1.0	<i>DET</i>	→ <i>the</i>

PCFGs

- We can induce a **probabilistic context-free grammar** (PCFG) from the treebank
- With multiple annotated sentences, we can get probabilities for production rules. Example:

1.0	<i>S</i>	→ <i>NP VP</i>
0.6	<i>NP</i>	→ <i>DET N</i>
0.4	<i>NP</i>	→ <i>ADJ N</i>
0.7	<i>VP</i>	→ <i>V NP</i>
0.3	<i>VP</i>	→ <i>V</i>
0.8	<i>N</i>	→ <i>dog</i>
0.2	<i>N</i>	→ <i>cat</i>
1.0	<i>V</i>	→ <i>barks</i>
1.0	<i>DET</i>	→ <i>the</i>

- Notice that the probabilities for each left-hand side must sum to one (unity)

Parameter Estimation

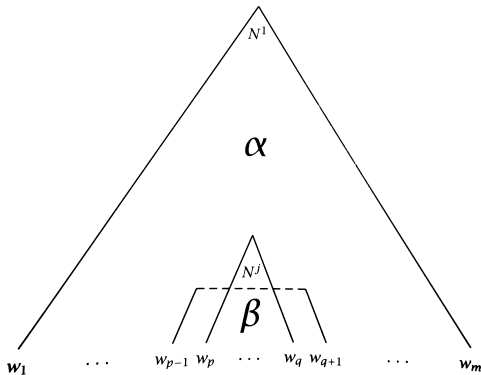
- So how do we get these probabilities?
- If we have a treebank, we can start with just counting how often productions occur (maximum likelihood estimation)

Parameter Estimation

- So how do we get these probabilities?
- If we have a treebank, we can start with just counting how often productions occur (maximum likelihood estimation)
- If we don't have a treebank, we can still use unannotated text, and apply the **inside-outside algorithm**
- The inside–outside algorithm is just the expectation–maximization (EM) algorithm applied to trees
- We start by randomly initializing probabilities to all possible rule productions, then use EM to search for good rule probabilities that maximize the likelihood of the training set

Inside-Outside Algorithm

- The inside-outside algorithm uses inside- and outside-probabilities:
 - Inside probability: $\beta_j(p, q) = P(w_{pq} | N_{pq}^j, G)$
 - Outside probability: $\alpha_j(p, q) = P(w_{1(p-1)}, N_{pq}^j, w_{(q+1)m} | G)$



String Probabilities

- We use the inside probability of the entire sentence to get the probability of that sentence:

$$P(w_{1:m}|G) = P(N^1 \xRightarrow{*} w_{1:m}|G) = \beta(1, m)$$

- Inside probabilities are calculated recursively & compositionally for each rule production
- We can do this because rule productions in context-free grammars are, well, context-free!

String Probabilities

- We use the inside probability of the entire sentence to get the probability of that sentence:

$$P(w_{1:m}|G) = P(N^1 \xRightarrow{*} w_{1:m}|G) = \beta(1, m)$$

- Inside probabilities are calculated recursively & compositionally for each rule production
- We can do this because rule productions in context-free grammars are, well, context-free!
- Probabilities of ambiguous parses at a given non-terminal are summed, since either parse could have produced the final substring

PCFGs vs. n -gram Language Models (Lexicalized Probabilistic Regular Grammars)

- PCFGs can better handle long-distance dependencies like subject-verb agreement and filler-gap dependencies
- PCFGs usually give worse perplexity than n -gram LMs. Why?

PCFGs vs. n -gram Language Models (Lexicalized Probabilistic Regular Grammars)

- PCFGs can better handle long-distance dependencies like subject-verb agreement and filler-gap dependencies
- PCFGs usually give worse perplexity than n -gram LMs. Why? Mostly because PCFGs are unlexicalized – they use pre-terminals (word classes / POS tags). Thus they fail to account for local co-occurrences like multiword expressions and proper names.

PCFGs vs. n -gram Language Models (Lexicalized Probabilistic Regular Grammars)

- PCFGs can better handle long-distance dependencies like subject-verb agreement and filler-gap dependencies
- PCFGs usually give worse perplexity than n -gram LMs. Why? Mostly because PCFGs are unlexicalized – they use pre-terminals (word classes / POS tags). Thus they fail to account for local co-occurrences like multiword expressions and proper names.
- PCFGs take longer to train
- PCFGs need manually-annotated treebanks to give decent results
- PCFG parsers (eg. CKY) are usually not incremental