

GRAMMARS

David Kauchak
CS159 – Fall 2014

some slides adapted from
Ray Mooney

Admin

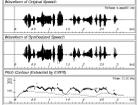
Assignment 3

Quiz #1

How was the lab last Thursday?

Simplified View of Linguistics

Phonetics



→ /waddiyasai/

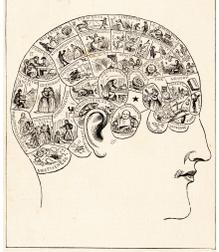
Nikolai Trubetzkoy in *Grundzüge der Phonologie* (1939) defines **phonology** as "the study of sound pertaining to the system of language," as opposed to **phonetics**, which is "the study of sound pertaining to the act of speech."

<http://en.wikipedia.org/wiki/Phonology>

Phonetics: "The study of the pronunciation of words"
Phonology: "The areas of linguistics that describes the systematic way that sounds are differently realized in different environments"

--The book

Not to be confused with...



phrenology

Context free grammar

$$S \rightarrow NP VP$$

left hand side (single symbol) right hand side (one or more symbols)

Formally...

$$G = (NT, T, P, S)$$

NT: finite set of nonterminal symbols

T: finite set of terminal symbols, NT and T are disjoint

P: finite set of productions of the form
 $A \rightarrow \alpha, A \in NT \text{ and } \alpha \in (T \cup NT)^*$

$S \in NT$: start symbol

CFG: Example

Many possible CFGs for English, here is an example (fragment):

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow \text{boy} \mid \text{girl}$
 $V \rightarrow \text{sees} \mid \text{likes}$
 $Adj \rightarrow \text{big} \mid \text{small}$
 $Adv \rightarrow \text{very}$
 $DetP \rightarrow a \mid \text{the}$

Grammar questions

Can we determine if a sentence is grammatical?

Given a sentence, can we determine the syntactic structure?

Can we determine how likely a sentence is to be grammatical? to be an English sentence?

Can we generate candidate, grammatical sentences?

Which of these can we answer with a CFG? How?

Grammar questions

Can we determine if a sentence is grammatical?

- Is it accepted/recognized by the grammar
- Applying rules right to left, do we get the start symbol?

Given a sentence, can we determine the syntactic structure?

- Keep track of the rules applied...

Can we determine how likely a sentence is to be grammatical? to be an English sentence?

- Not yet... no notion of "likelihood" (probability)

Can we generate candidate, grammatical sentences?

- Start from the start symbol, randomly pick rules that apply (i.e. left hand side matches)

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

S

What can we do?

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

S

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

NP VP

What can we do?

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

NP VP

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

DetP N VP

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

DetP N VP

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

the boy VP

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

the boy likes NP

Derivations in a CFG

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

the boy likes a girl

Derivations in a CFG; Order of Derivation Irrelevant

$S \rightarrow NP VP$
 $VP \rightarrow V NP$
 $NP \rightarrow DetP N \mid AdjP NP$
 $AdjP \rightarrow Adj \mid Adv AdjP$
 $N \rightarrow boy \mid girl$
 $V \rightarrow sees \mid likes$
 $Adj \rightarrow big \mid small$
 $Adv \rightarrow very$
 $DetP \rightarrow a \mid the$

the boy likes a girl

Derivations of CFGs

String rewriting system: we derive a string

Derivation history shows constituent tree:

the boy likes a girl

Parsing

Parsing is the field of NLP interested in automatically determining the syntactic structure of a sentence

parsing can be thought of as determining what sentences are "valid" English sentences

As a by product, we often can get the structure

Parsing

Given a CFG and a sentence, determine the possible parse tree(s)

S -> NP VP
 NP -> N
 NP -> PRP
 NP -> N PP
 VP -> V NP
 VP -> V NP PP
 PP -> IN N
 PRP -> I
 V -> eat
 N -> sushi
 N -> tuna
 IN -> with

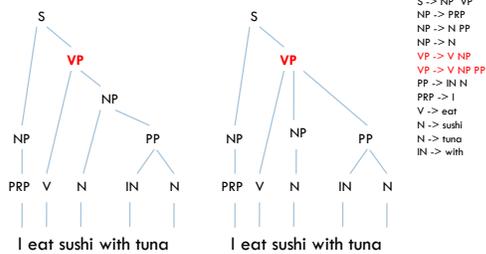
I eat sushi with tuna

What parse trees are possible for this sentence?

How did you do it?

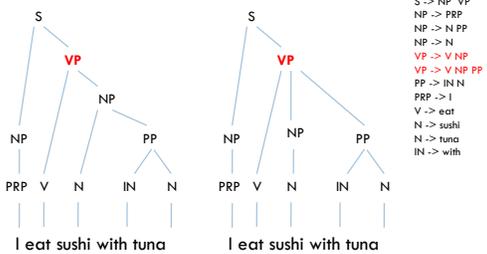
What if the grammar is much larger?

Parsing



What is the difference between these parses?

Parsing ambiguity

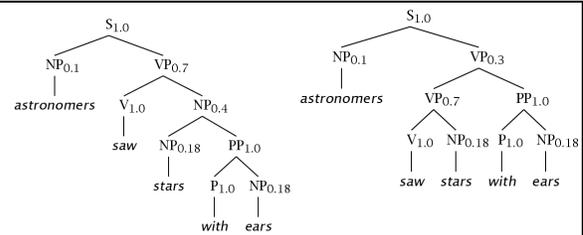


How can we decide between these?

A Simple PCFG

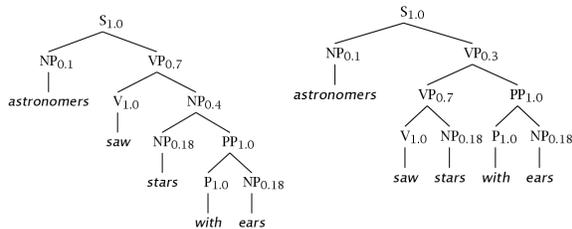
Probabilities!

S	→	NP VP	1.0	NP	→	NP PP	0.4
VP	→	V NP	0.7	NP	→	astronomers	0.1
VP	→	VP PP	0.3	NP	→	ears	0.18
PP	→	P NP	1.0	NP	→	saw	0.04
P	→	with	1.0	NP	→	stars	0.18
V	→	saw	1.0	NP	→	telescope	0.1



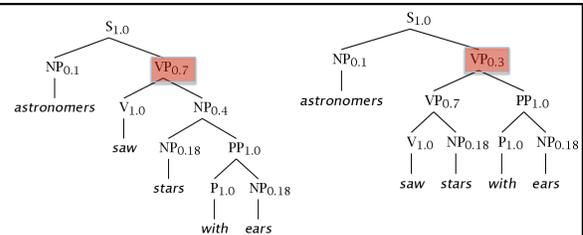
Just like *n*-gram language modeling, PCFGs break the sentence generation process into smaller steps/probabilities

The probability of a parse is the product of the PCFG rules



What are the different interpretations here?

Which do you think is more likely?



$$= 1.0 * 0.1 * 0.7 * 1.0 * 0.4 * 0.18$$

$$* 1.0 * 1.0 * 0.18$$

$$= 0.0009072$$

$$= 1.0 * 0.1 * 0.3 * 0.7 * 1.0 * 0.18$$

$$* 1.0 * 1.0 * 0.18$$

$$= 0.0006804$$

Parsing problems

- Pick a model
 - e.g. CFG, PCFG, ...
- Train (or learn) a model
 - What CFG/PCFG rules should I use?
 - Parameters (e.g. PCFG probabilities)?
 - What kind of data do we have?
- Parsing
 - Determine the parse tree(s) given a sentence

PCFG: Training

If we have example parsed sentences, how can we learn a set of PCFGs?

S → NP VP	0.9
S → VP	0.1
NP → Det A N	0.5
NP → NP PP	0.3
NP → Prop N	0.2
A → ε	0.6
A → Adj A	0.4
PP → Prep NP	1.0
VP → V NP	0.7
VP → VP PP	0.3

English

Extracting the rules

S	→	NP VP
NP	→	PRP
PRP	→	I
VP	→	V NP
V	→	eat
NP	→	N PP
N	→	sushi
PP	→	IN N
IN	→	with
N	→	tuna

What CFG rules occur in this tree?

Estimating PCFG Probabilities

We can extract the rules from the trees

S	→	NP VP	1.0
NP	→	PRP	0.7
PRP	→	I	0.3
VP	→	V NP	1.0
V	→	eat	1.0
NP	→	N PP	1.0
N	→	sushi	1.0
PP	→	P NP	1.0
P	→	with	1.0
N	→	saw	1.0
...			

How do we go from the extracted CFG rules to PCFG rules?

Estimating PCFG Probabilities

Extract the rules from the trees

Calculate the probabilities using MLE

$$\alpha \rightarrow \beta \quad \Rightarrow \quad p(\alpha \rightarrow \beta | \alpha)$$

$$P(\alpha \rightarrow \beta | \alpha) = \frac{\text{count}(\alpha \rightarrow \beta)}{\sum_{\gamma} \text{count}(\alpha \rightarrow \gamma)} = \frac{\text{count}(\alpha \rightarrow \beta)}{\text{count}(\alpha)}$$

Estimating PCFG Probabilities

	Occurrences
S → NP VP	10
S → V NP	3
S → VP PP	2
NP → N	7
NP → N PP	3
NP → DT N	6

$$P(S \rightarrow V NP) = ?$$

$$P(S \rightarrow V NP) = P(S \rightarrow V NP | S) = \frac{\text{count}(S \rightarrow V NP)}{\text{count}(S)} = \frac{3}{15}$$

Grammar Equivalence

What does it mean for two grammars to be equal?

Grammar Equivalence

Weak equivalence: grammars generate same set of strings

- Grammar 1: NP → DetP N and DetP → a | the
- Grammar 2: NP → a N | the N

Strong equivalence: grammars have same set of derivation trees

- With CFGs, possible only with useless rules
- Grammar 2: NP → a N | the N
- Grammar 3: NP → a N | the N, DetP → many

Normal Forms

There are weakly equivalent **normal forms** (Chomsky Normal Form, Greibach Normal Form)

A CFG is in Chomsky Normal Form (CNF) if all productions are of one of two forms:

- $A \rightarrow BC$ with A, B, C nonterminals
- $A \rightarrow \alpha$, with A a nonterminal and α a terminal

Every CFG has a weakly equivalent CFG in CNF

CNF Grammar

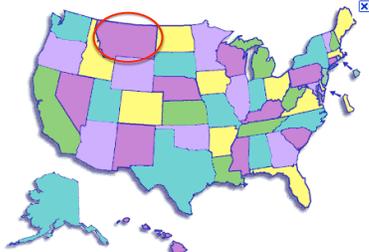
<p>S -> VP VP -> VB NP VP -> VB NP PP NP -> DT NN NP -> NN NP -> NP PP PP -> IN NP DT -> the IN -> with VB -> film VB -> trust NN -> man NN -> film NN -> trust</p>	<p>S -> VP VP -> VB NP VP -> VP2 PP VP2 -> VB NP NP -> DT NN NP -> NN NP -> NP PP PP -> IN NP DT -> the IN -> with VB -> film VB -> trust NN -> man NN -> film NN -> trust</p>
---	---

Probabilistic Grammar Conversion

Original Grammar	Chomsky Normal Form
S → NP VP 0.8	S → NP VP 0.8
S → Aux NP VP 0.1	S → XI VP 0.1
	XI → Aux NP 1.0
S → VP 0.1	S → book include prefer 0.01 0.004 0.006
	S → Verb NP 0.05
	S → VP PP 0.03
NP → Pronoun 0.2	NP → I he she me 0.1 0.02 0.02 0.06
NP → Proper-Noun 0.2	NP → Houston NVA 0.16 .04
NP → Det Nominal 0.6	NP → Det Nominal 0.6
Nominal → Noun 0.3	Nominal → book flight meal money 0.03 0.15 0.06 0.06
Nominal → Nominal Noun 0.2	Nominal → Nominal Noun 0.2
Nominal → Nominal PP 0.5	Nominal → Nominal PP 0.5
VP → Verb 0.2	VP → book include prefer 0.1 0.04 0.06
VP → Verb NP 0.5	VP → Verb NP 0.5
VP → VP PP 0.3	VP → VP PP 0.3
PP → Prep NP 1.0	PP → Prep NP 1.0

States





What is the capitol of this state? Helena (Montana)

Grammar questions

Can we determine if a sentence is grammatical?

Given a sentence, can we determine the syntactic structure?

Can we determine how likely a sentence is to be grammatical? to be an English sentence?

Can we generate candidate, grammatical sentences?

Next time:
parsing