

# Lecture 10

August 23, 2016

## Formal grammars



Reminder: start the recording

# Announcements

- Project 4: DNA
  - Due at 11:45 p.m. on 9/1 (one week from Thursday)
- Project 5: Naïve Bayesian language classifier
  - Due at 11:45 p.m. on Thursday Sept. 8th  
<http://courses.washington.edu/ling473/Project5.pdf>
- Reminder: Writing Assignment  
<http://courses.washington.edu/ling473/writing-assignment.html>
  - Due at **11:45 p.m.** on Tuesday, Sept. 6<sup>th</sup>

# Self-quiz answers

Two fair coins are tossed. One of them shows heads and the other rolls under the couch. What is that chance that the hidden coin is showing tails?

There is much controversy surrounding the frequentist vs. Bayesian interpretation of this type of problem

[http://en.wikipedia.org/wiki/Marilyn\\_vos\\_Savant](http://en.wikipedia.org/wiki/Marilyn_vos_Savant)

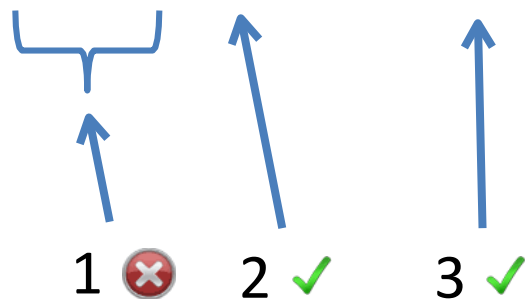
[http://en.wikipedia.org/wiki/Boy\\_or\\_Girl\\_paradox](http://en.wikipedia.org/wiki/Boy_or_Girl_paradox)

In the non-frequentist approach, the answer would be  $\frac{2}{3}$

(H,H) (H,T) (T,H) (T,T)

One head is showing already

(H,H) (H,T) (T,H) ~~(T,T)~~



2 out of 3 remaining possibilities have a hidden 'tails'

In[1]:= `□ □ Tuples` , `T` ,

Out[1]= `H, H , H, T , T, H , T, T`

Size of the sample space:

In[2]:= `n` `Length` `□`

Out[2]= 4

Assuming heads and tails are equally likely, the unconditional probability of one head and one tail:

In[4]:= `Pr H` `□ T` `□ 1` `□` 
$$\frac{\text{Length Select } \square, \text{MemberQ } \square 1, H \quad \text{MemberQ } \square 1, T \&}{n}$$

Out[4]=  $\frac{1}{2}$

In[7]:= `Pr T` `□ 1` `□` 
$$\frac{\text{Length Select } \square, \text{MemberQ } \square 1, T \&}{n}$$

Out[7]=  $\frac{3}{4}$

Probability of exactly one head and one tail given that one is head:

In[8]:= `Pr H` `□ T` `□ 1` `H` `□ 1` 
$$\frac{\text{Pr } H \square T \square 1}{\text{Pr } T \square 1}$$

Out[8]=  $\frac{2}{3}$

There are two boxes. One contains two black marbles and two white marbles. The other contains four black marbles and one white marble. In this experiment, a box will be selected randomly, and then a marble will be drawn randomly from the selected box.

- What is the probability that the marble will be black?
- We run the experiment and it turns out that the selected marble is white. What is the probability that the first box was selected?

```
In[19]:= B1 = {b, b, w, w}
```

```
In[20]:= B2 = {b, b, b, b, w}
```

Assuming the selection of either of the two boxes is equally likely, we record the prior probability that the first box was selected (and complement) from the problem statement:

```
In[21]:= Pr FS = Pr FS^C = 1/2
```

```
Out[21]=
```

$$\frac{1}{2}$$

Conditional probabilities of drawing black given the box selection:

```
In[22]:= Pr DB FS = Count B1, b / Length B1
```

```
Out[22]=
```

$$\frac{1}{2}$$

```
In[23]:= Pr DB FS^C = Count B2, b / Length B2
```

```
Out[23]=
```

$$\frac{4}{5}$$

Overall prior probability of drawing black:

```
In[33]:= Pr DB = Pr FS Pr DB FS + Pr FS^C Pr DB FS^C
```

```
Out[33]=
```

$$\frac{13}{20}$$

```
In[34]:= N
```

```
Out[34]=
```

$$0.65$$

For part (b.), use Bayes' theorem

In[35]:  $\Pr DW \square 1 \square \Pr DB$

Out[35]=

$$\frac{7}{20}$$

In[36]:  $\Pr DW FS \square \frac{\text{Count B1, } w}{\text{Length}}$

Out[36]=

$$\frac{1}{2}$$

Now it is a simple matter to calculate the probability that the first was previously selected, given the observation of a white marble:

In[37]:  $\Pr FS DW \square \frac{\Pr DW FS \Pr FS}{\Pr DW}$

Out[37]=

$$\frac{5}{7} = 0.714286$$



The high school Shakespeare club has 5 freshman boys, 7 freshman girls, and 6 sophomore boys. How many sophomore girls must be in the club in order for gender and class to be independent when a student is chosen at random from the club?

Let  $n$  be the number of sophomore girls

In[5]:= Pr boy, freshman  $\square \frac{5}{n \square 18}$

In[6]:= Pr boy  $\square \frac{11}{n \square 18}$

In[7]:= Pr freshman  $\square \frac{12}{n \square 18}$

In[8]:= Solve Pr boy, freshman  $\square$  Pr boy Pr freshman ,  $n$

Out[8]=  $n \square \frac{42}{5} = 8.4$

It is *not possible* for gender and class to be independent



The linguistics section at the library has three books on Austronesian languages. We choose two linguistics books at random. The probability of them both being on Austronesian language is  $\frac{1}{1650} = 0.000606061$ . How many linguistics books are there?

```
In[185]:= Replace n Solve  $\left(\frac{3}{n}\right)\left(\frac{2}{n-1}\right) \square .000606061, n$ 
```

```
Out[185]=  $\square 99., 100.$ 
```

Discard a negative solution:

```
In[186]:= Select  $\square, \square \square 0 \&$ 
```

```
Out[186]= 100.
```

A multiple choice exam is given. A problem has four possible answers, and exactly one answer is correct. The student is allowed to select as many answers as he likes. If his chosen subset contains the correct answer, the student receives three points, but he loses one point for each wrong answer in his chosen subset. What is the expected score?

zero.

- A die is rolled  $TT$  times until **6** is shown. What is the probability distribution for  $TT$ ?
- *Geometric* with  $pp = \frac{1}{6}$
- Expected value: 6

Note: this would be the same answer for 1, 2, 3, 4, 5...

A die is rolled  $TT$  times until **6** is shown. Let  $EE$  denote the event that  $TT > 3$ . Let  $FF$  denote the event that  $TT > 6$ . Calculate  $PP(EE)$  and  $PP(FF|EE)$ .

- from last week's lecture: for a geometric distribution,  
 $PP(XX > xx) = (1 - pp)^{xx}$

$$\left(1 - \frac{1}{6}\right)^3 = \frac{125}{216} = .5787$$

## 7c. alternate



$$\begin{aligned} & PR(TT > 3) \\ &= 1 - PP(TT \leq 3) \\ &= 1 - \left( PR(TT = 1) + PP(TT = 2) + PP(TT = 3) \right) \\ &\quad \quad \quad 3 \\ &= 1 - \blacklozenge (1 - pp)^{xx-1} pp \\ &\quad \quad \quad xx=1 \\ &= \frac{125}{216} \\ &= .5787 \end{aligned}$$

7d.

$$\begin{aligned} PR(FE|EE) &= \\ &= PR(TT > 6 | TT > 3) \\ &= \frac{PR(TT > 6, TT > 3)}{P(T > 3)} \\ &= \frac{PR(TT > 6)}{PR(TT > 3)} \\ &= \frac{(1 - pp)^6}{(1 - pp)^3} \\ &= \frac{125}{216} \end{aligned}$$

In other words, given that you've rolled the die 3 times and you have not gotten a **6** yet, what is the probability that you won't get a **6** within the next three rolls?

Notice: the probability hasn't changed, suggesting that the trials (each roll) are independent.

# Tagging objective function

Predict a sequence of tags  $tt$  based on the probability of tags and words  $PP(tt_{ii} | ww_{ii})$ . Given sentence

$$S = (ww_0, ww_1, \dots, ww_n)$$
$$tt = \underset{tt_{ii}}{\text{argmax}} PP(tt_{ii} | ww_{ii}).$$

“t-hat”



“ $tt$  is the best sequence of tags that match a tag  $tt_{ii}$  to its word  $ww_{ii}$ .”

This material is also covered in section 5.5 (p.139) of Jurafsky & Martin, 2<sup>nd</sup> ed.



## Simplistic tagger

$$S = (w_0, w_1, \dots, w_n)$$
$$t = \operatorname{argmax}_{t_i} PP(t_i | w_i)$$

repeated from last slide

This is surely the function we want to maximize, but it's not clear how to calculate the probabilities  $PP(t|w)$ .

Simplistic tagger: Why don't we use probabilities calculated from a corpus ?

- like you did for Assignment 3

# Simplistic tagger

DT	NN	VBD	RB		IN	DT	NN	,	CC	DT	VBG		NNS	VBD		DT	NN	VBD
the	cold	passed	reluctantly		from	the	earth	,	and	the	retiring	fogs	revealed	an	army	stretched		

IN	IN	DT	NNS	,	VBG	.	IN	DT	NN		VBN	IN	JJ	TO	VB	,	DT	NN	VBN
out	on	the	hills	,	resting	.	as	the	landscape	changed	from	brown	to	green	,	the	army	awakened	

,	CC	VBD	TO	VB		IN	NN		IN	DT	NN	IN	NNS	.
,	and	began	to	tremble	with	eagerness	at	the	noise	of	rumors	.		

$$\operatorname{argmax}_{tt}^{PP}(tt|\text{the}) = \text{DT} \quad \checkmark$$

$$\operatorname{argmax}_{tt}^{PP}(tt|\text{cold}) = \text{JJ} \quad \times$$

# How well does the simplistic tagger work?

- Such a POS tagger is not really usable

Most probable POS tag: JJ  
Correct tag: NN

Most probable POS tag: VBG  
Correct tag: JJ

DT		VBD	RB		IN	DT	NN	,	CC	DT		NNS	VBD		DT	NN	VBD
the	cold	passed	reluctantly		from	the	earth	,	and	the	retiring	fogs	revealed		an	army	stretched

IN	IN	DT	NNS	,	VBG	.	IN	DT	NN		VBN	IN	JJ	TO	JJ	,	DT	NN	VBN
out	on	the	hills	,	resting	.	as	the	landscape		changed	from	brown	to	green	,	the	army	awakened

,	CC	VBD	TO	VB		IN	NN		IN	DT	NN	IN	NNS	.
,	and	began	to	tremble		with	eagerness		at	the	noise	of	rumors	.

# Use Bayes Theorem



Of course, you have  
this memorized

$$PP(AA|BB) = \frac{PP(BB|AA)PP(AA)}{PP(BB)}$$

Remember, this was  
our objective function

$$tt = \operatorname{argmax}_{tt} PP(tt_{ii} | ww_{ii})$$

$$tt = \operatorname{argmax}_{tt} \frac{PP(ww_{ii} | tt_{ii})PP(tt_{ii})}{PP(ww_{ii})}$$



This is one of the most important slides of this entire class

For each evaluated value of  $i$ ,  $PP(w_{ii})$  will be the same. We can cancel it.

$$tt = \operatorname{argmax}_{tt} \frac{PP(w_{ii} | tt_{ii}) PP(tt_{ii})}{\cancel{PP(w_{ii})}}$$
$$tt = \operatorname{argmax}_{tt} PP(w_{ii} | tt_{ii}) PP(tt_{ii})$$

The best sequence of tags is determined by the probability of each word given its tag and also the probability of that tag.

I repeat...

$$tt = \operatorname{argmax}_{tt} PP(w_{ii} | tt_{ii}) PP(tt_{ii})$$

“likelihood”

“prior”



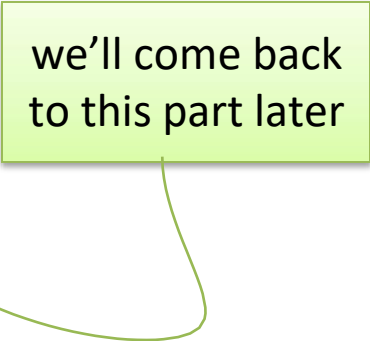
“We compute the most probable tag sequence... by multiplying the **likelihood** and the **prior probability** for each tag sequence and choosing the tag sequence for which this product is greatest.

“Unfortunately, this is still too hard to compute directly...”

Jurafsky & Martin (paraphrase) p.140

We still need to make some assumptions.

we'll come back  
to this part later



$$tt = \operatorname{argmax}_{tt} PP(w_{ii} | tt_{ii}) PP(tt_{ii})$$

**Assumption 1:** If we want to use corpus probabilities to estimate  $PP(w_{ii} | tt_{ii})$ , we need to formally note that we're assuming

$$PP'(w_{ii} | tt_{ii}) \approx \underset{ii}{\blacklozenge ?} PP(w_{ii})$$

“The only POS tag a word depends on is its own.”

# Any progress?

- So wait: if we're assuming the only *POS tag* a *word* depends on is its own, how is this going to be better than the *simplistic tagger* from before, which assumed that the only *word* a *POS tag* depends on is its own?
- In other words, Why is  $PP(w|t)$  going to work better than  $PP(t|w)$ ?
- Hint:  $|\Omega|$
- Hint:  $|T| \ll |W|$

Answer: because there are a lot more distinct words than tags, conditioning on *tags* rather than *words* increases the resolution of the corpus measurements



## example

$$PP(\text{cold}|\text{NN}) = .00002$$

$$PP(\text{cold}|\text{JJ}) = .00040$$



$$PR(\text{JJ}|\text{cold}) = .97$$

$$PR(\text{NN}|\text{cold}) = .03$$

This value will drown out our calculation and we'd never tag "cold" as a noun!

$$t_i = \operatorname{argmax}_{t_i} P(t_i | w_{i-1} \dots w_1, t_{i-1} \dots t_1)$$

**Assumption 2:** The only tags that a tag  $t_i$  depends on are the  $n$  previous tags,  $t_{i-n} \dots t_{i-1}$ . For example, in a POS bigram model:

$$P(t_i) \approx \sum_{t_{i-1}} P(t_i, t_{i-1})$$

This is known as the **bigram assumption**: “The only POS tag(s) a POS tag depends on are the ones immediately preceding it.”

# Putting it together

$$tt = \operatorname{argmax}_{tt} PP (ww_i | tt_{ii}) PP (tt_{ii})$$

$$PP' (ww_{ii} | tt_{ii}) \approx \blacklozenge ? \left( \begin{matrix} PP & | & \end{matrix} \right)$$

$ww_{ii} \quad tt_{ii}$   
 $ii$

$$PP' (tt_{ii}) \approx \blacklozenge ? \left( \begin{matrix} PP & | & \end{matrix} \right) \begin{matrix} tt_{ii} \\ tt_{ii-1} \end{matrix}$$

$ii$

$$tt = \operatorname{argmax}_{tt} \blacklozenge ? \left( \begin{matrix} PP & | & \end{matrix} \right) \begin{matrix} ww_{ii} \\ tt_{ii} \end{matrix} \blacklozenge ? \left( \begin{matrix} PP & | & \end{matrix} \right) \begin{matrix} tt_{ii} \\ tt_{ii-1} \end{matrix}$$

$ii \qquad ii$

$$tt = \operatorname{argmax}_{tt} \blacklozenge ? \left( \begin{matrix} PP & | & \end{matrix} \right) \begin{matrix} ww_{ii} \\ tt_{ii} \end{matrix} \left( \begin{matrix} PP & | & \end{matrix} \right) \begin{matrix} tt_{ii} \\ tt_{ii-1} \end{matrix}$$

$ii$



Reminder: estimating  $PP(w_{ii}|t_{ii})$  from a corpus

Definition of  
conditional probability

→  $PR(AA|BB) = \frac{PP(AA, BB)}{PP(BB)}$

$$PR(AA|BB) = \frac{\frac{\text{count}(AA, BB)}{|\Omega|}}{\frac{\text{count}(BB)}{|\Omega|}}$$

word likelihood

$$PR(w_{ii}|t_{ii}) = \frac{\text{count}(w_{ii}, t_{ii})}{\text{count}(t_{ii})}$$

Reminder: estimating  $PP(tt_{ii} | tt_{ii-1})$  from a corpus

Definition of  
conditional probability

$$\rightarrow PR(AA | BB) = \frac{PP(AA, BB)}{PP(BB)}$$

$$PR(AA | BB) = \frac{\frac{\text{count}(AA, BB)}{\cancel{|\Omega|}}}{\frac{\text{count}(BB)}{\cancel{|\Omega|}}}$$

$$PP(tt_{ii} | tt_{ii-1}) = \frac{\text{count}(tt_{ii-1}, tt_{ii})}{\text{count}(tt_{ii-1})}$$

# POS tagging objective function

$$tt = \underset{tt}{\operatorname{argmax}} \left( \frac{\text{count}(ww_{ii}, tt_{ii})}{\text{count}(tt_{ii})} \times \frac{\text{count}(tt_{ii-1}, tt_{ii})}{\text{count}(tt_{ii-1})} \right)$$

Best POS tag  
sequence

How often does word  
 $ww_{ii}$  occur with tag  $tt_{ii}$   
in the corpus?

How often does  $tt_{ii}$   
follow  $tt_{ii-1}$  in the  
corpus?

This might seem a little backwards (especially if you aren't familiar with Bayes' theorem). We're trying to find the best *tag sequence*, but we're using  $PP(ww|tt)$ , which seems to be predicting *words*.

This compares: “If we are expecting an **adjective** (based on the tag sequence), how likely is it that the adjective will be ‘cold?’” **versus** “If we are expecting a **noun**, how likely is it that the noun will be ‘cold?’”

DT		VBD	RB		IN	DT	NN	,
the	cold	passed	reluctantly		from	the	earth	,

“If we are expecting an **adjective**, how likely is it that the adjective will be ‘cold?’” (high) **WEIGHTED BY** our chance of seeing the sequence **DT JJ** (medium)

*versus*

“If we are expecting a **noun**, how likely is it that the noun will be ‘cold?’” (medium) **WEIGHTED BY** our chance of seeing the sequence **DT NN** (very high)

**THE WINNER: NN**





# Multiplying probabilities

- We're multiplying a whole lot of probabilities together
- What do we know about probability values?
$$0 \leq p \leq 1$$
- What happens when you multiply a lot of these together?
- This is an important consideration in computational linguistics. We need to worry about **underflow**.

# Underflow

- When multiplying many probability terms together, we need to prevent underflow
  - Due to limitations in the computer's internal representation of floating point numbers, the product quickly becomes zero
- We usually work with the logarithm of the probability values
- This is known as the “log-prob”  
$$= \log_{10} pp$$

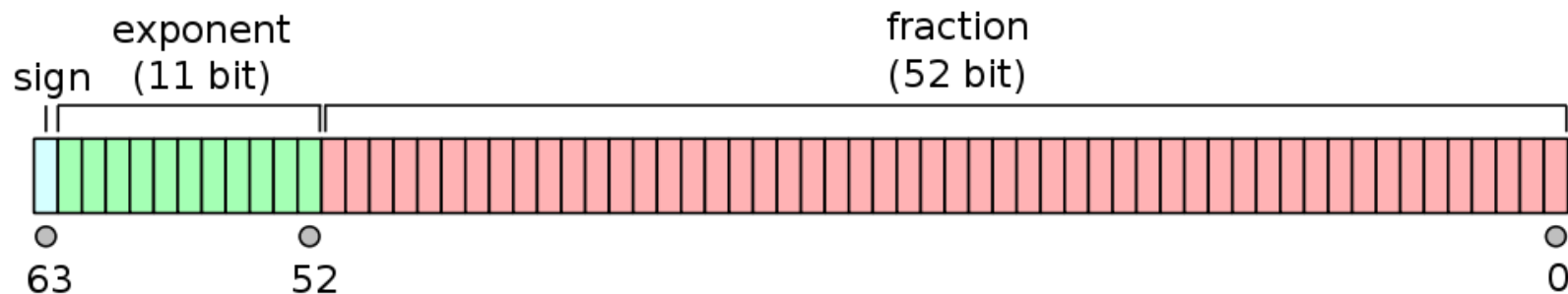
- 32-bit “single” “float”

Diagram illustrating the IEEE 754 single-precision floating-point format. The 32-bit word is divided into three fields:

- sign (1 bit)
- exponent (8 bits)
- fraction (23 bits)

The bit pattern shown is 00111100 followed by 23 zeros, which equals 0.15625. The bit indices 31, 30, 23, 22, and 0 are marked.

- $$\approx \pm 1.8 \times 10^{308}$$



# logarithms refresher

definition:

$$\log_{bb} xx = yy: xx = bb^{yy}$$

$$bb^{xx} \times bb^{yy} = bb^{xx+yy}$$

$$\log xx yy = \log xx + \log yy$$

$$\log \underset{\textit{ii}}{\blacklozenge ?} xx_{ii} = \underset{\textit{ii}}{\blacklozenge ?} \log xx_{ii}$$

$$\frac{bb^{xx}}{bb^{yy}} = bb^{xx-yy}$$

$$\log \frac{xx}{yy} = \log xx - \log yy$$



Write an expression for Bayes' theorem as log-probabilities

# Bayes' theorem as log-prob

$$PP(AA|BB) = \frac{PP(BB|AA)PP(AA)}{PP(BB)}$$

$$\log PP(AA|BB) = \log PP(BB|AA) + \log PP(AA) - \log PP(BB)$$

## Remember this?

$$tt = \operatorname{argmax}_{tt} \underset{ii}{\blacklozenge ?} \left( \frac{\operatorname{count}(ww_{ii}, tt_{ii})}{\operatorname{count}(tt_{ii})} \times \frac{\operatorname{count}(tt_{ii-1}, tt_{ii})}{\operatorname{count}(tt_{ii-1})} \right)$$

$$tt = \operatorname{argmax}_{tt} \underset{ii}{\blacklozenge ?} \left( \log \frac{\operatorname{count}(ww_{ii}, tt_{ii})}{\operatorname{count}(tt_{ii})} + \log \frac{\operatorname{count}(tt_{ii-1}, tt_{ii})}{\operatorname{count}(tt_{ii-1})} \right)$$

Wait, how can you do that, there was no “log”  
outside of the  $\Pi$  !

# argmax magic

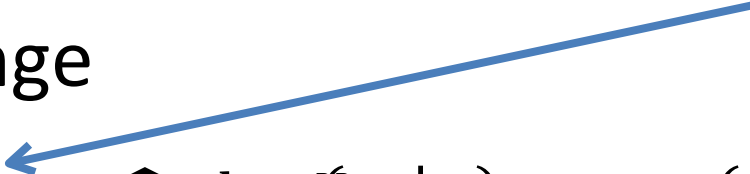
- Doesn't matter. Since argmax doesn't care about the actual answer, but rather just the *sequence that gives it*, we can drop the overall log
  - this is valid so long as  $\log xx$  is a monotonically increasing function
- argmax will find the same "best" tag sequence when looking at either **probabilities** or **log-probs** because both functions will peak at the same point

$$tt = \underset{ii}{\operatorname{argmax}} \log P \mid ) ww_{ii} tt_{ii} (+|\log P) \checkmark tt_{ii} tt_{ii-1}$$



# Hidden Markov Model

- This is the foundation for the **Hidden Markov Model** (HMM) for POS tagging
- To proceed further and solve the argmax is still a challenge

$$t_t = \underset{t_t}{\operatorname{argmax}} \log P(w_{t_t} | t_{t-1}) + \log P(t_t | t_{t-1})$$


- Computing this naively is still  $O(|T|^n)$

# Dynamic programming

- The **Viterbi algorithm** is typically used to decode Hidden Markov Models
  - You might get to implement it in Ling 570
- It is a **dynamic programming** technique
  - We maintain a trellis of partial computations
- This approach reduces the problem to  $O(|T|^2nn)$  time

# POS Trigram model

Recall the bigram assumption:

$$PP'(tt_{ii}) \approx \underset{ii}{\blacklozenge ?}^{PP}(tt_{ii} | tt_{ii-1})$$

We can improve the tagging accuracy by extending to a trigram (or larger) model

$$PP'(tt_{ii}) \approx \underset{ii}{\blacklozenge ?}^{PP}(tt_{ii} | tt_{ii-2}, tt_{ii-1})$$

## Data sparsity

- However, we might start having a problem if we try to get a value for  $PP(tt_i | tt_{i-2}, tt_{i-1})$  by counting in the corpus

$$\frac{\text{count}(tt_{i-2}, tt_{i-1}, tt_i)}{\text{count}(tt_{i-2}, tt_{i-1})}$$

...it was a butterfly in distress that she...

The count of this in our training set is likely to be zero

# Unseens

- Our model will predict zero probability for something that we actually encounter
  - This counts as a failure of the model
- This is a pervasive problem in corpus linguistics
  - At runtime, how do you deal with observations that you never encountered during training (**unseen data**)?

# Smoothing

- We don't want our model to have a discontinuity between something infrequent and something unseen
- Various techniques address this problem:
  - add-one smoothing
  - Good-Turing method
  - Assume unseens have probability of the rarest observation
  - Ideally, smoothing preserves the validity of your probability space

# Formal grammars

## Parsing

This lecture adapts some slides from:

- Andrew McCallum, (UMass)
- Chris Manning
- Jason Eisner
- Norman Landis (Fairleigh Dickinson Univ.)

# Constituents

- In lecture 1, we talked about constituents
- Constituents help us organize language into structures

[**Thing** The dog] is [**Place** in the garden]

[**Thing** The dog] is [**Property** fierce]

[**Action** [**Thing** The dog] is chasing [**Thing** the cat]]

[**State** [**Thing** The dog] was sitting [**Place** in the garden] [**Time** yesterday]]

[**Action** [**Thing** We] ran [**Path** out into the water]]

[**Action** [**Thing** The dog] barked [**Property/Manner** loudly]]

[**Action** [**Thing** The dog] barked [**Property/Amount** nonstop for five hours]]



# Word categories

- Traditional parts of speech

Noun	Names of things	boy, cat, truth
Verb	Action or state	become, hit
Pronoun	Used for noun	I, you, we
Adverb	Modifies V, Adj, Adv	sadly, very
Adjective	Modifies noun	happy, clever
Conjunction	Joins things	and, but, while
Preposition	Relation of N	to, from, into
Interjection	An outcry	ouch, oh, alas, psst

# Substitution test

- Adjective:

The {sad, intelligent, green, fat, ...} one is in the corner.

- Noun:

The {cat, mouse, dog} ate the bug.

- Verb:

Kim {loves, eats, makes, buys, moves} potato chips.

# Constituency

- The idea: Groups of words may behave as a single unit or phrase, called a **constituent**
- Sentences have parts, some of which appear to have subparts, which have subparts...
- These groupings of words that go together we will call constituents.
- e.g. Noun Phrase
  - Kermit the frog
  - they
  - December twenty-sixth
  - the reason he is running for president

# Constituent Phrases

For constituents, we usually name them as phrases based on the word that **heads** the constituent

<i>the man from Amherst</i>	is a Noun Phrase (NP) because the head <i>man</i> is a noun
<i>extremely clever</i>	is an Adjective Phrase (AP) because the head <i>clever</i> is an adjective
<i>down the river</i>	is a Prepositional Phrase (PP) because the head <i>down</i> is a preposition
<i>killed the rabbit</i>	is a Verb Phrase (VP) because the head <i>killed</i> is a verb

Note that a word is a constituent (a little one). Sometimes words also act as phrases. In:

*Joe grew potatoes.*

*Joe* and *potatoes* are both nouns and noun phrases.

Compare with:

*The man from Amherst* *grew* *beautiful russet potatoes*.

We say *Joe* counts as a noun phrase because it appears in a place that a larger noun phrase could have been.

# Evidence for constituency

- They appear in similar environments (before a verb)
  - *Kermit the frog comes on stage*
  - *They come to Massachusetts every summer*
  - *December twenty-sixth comes after Christmas*
  - *The reason he is running for president comes out only now.*
- But not each individual word in the constituent
  - \*The comes out... \*is comes out... \*for comes out...
- The constituent can be placed in a number of different locations
  - *On December twenty-sixth* I'd like to fly to Florida.
  - I'd like to fly *on December twenty-sixth* to Florida.
  - I'd like to fly to Florida *on December twenty-sixth*.
- But not split apart:
  - \**On December* I'd like to fly *twenty-sixth* to Florida.
  - \**On* I'd like to fly *December twenty-sixth* to Florida.

# Context-free grammar (CFG)

- or, “Phrase structure grammar”
- or, “Backus-Naur Form” (BNF)
- The most common way of modeling constituency
- The idea of basing a grammar on constituent structure dates back to Wilhem Wundt (1890), but not formalized until Chomsky (1956), and, independently, by Backus (1959).
- This is a particular type of **formal grammar**.
- Before we look at CFGs grammar, let’s put them in “context...”

# Formal Languages

- The set of all possible strings for an **alphabet**  $\Sigma$  is written as  $\Sigma^*$
- A **language** is a prescribed subset of  $\Sigma^*$

$$L \subset \Sigma^*$$

- Example:

$$\Sigma = \{a, b, c\}$$

$$\Sigma^* = \{\epsilon, a, b, c, ab, ac, ba, bc, ca, aaa, \dots\}$$

one language might be the set of strings of length less than or equal to 2:

$$L = \{\epsilon, a, b, c, aa, ab, ac, ba, bb, bc, ca, cb, cc\}$$

# “Strings”

- Note, when we talk about formal grammars, a **string** can be either a string of characters or a string of words (or symbols, etc.)
- Don't get confused by the more computer-science use/definition of the term
  - i.e. “a sequence of characters in contiguous memory, possibly zero-terminated”



# Rules

- It is useful to constrain the set of strings in a language to be other than  $\Sigma^*$
- Let's accept the existence of constituents
- Let's assume that constituents are defined by a set of **rules**,  $RR$
- Here's how we'll constrain  $\Sigma^*$  : constituents can only be *juxtaposed* as permitted by  $RR$

This introduces an ordering constraint which is prevalent in many—but not all—natural languages. But aside from  $LL$ 's constraints on ordering, what practical purpose(s) might ordering serve in a language?

# Modeling language

- Working towards a theoretical model of language
- So far we have:

$\Sigma$  – a set of symbols (words, letters, ...)

$RR$  – a set of rules

- We'll need to keep track of the particular system we're defining. Such a system is called a **grammar**:

$$GG = \langle \Sigma, RR, \dots \rangle$$

## Parsing v. Generation

- Given a configuration of rules in grammar  $GG$ , find the string(s)  $SS$  that can be formed. This is called **generation**

A language  $LL$  is the set of all strings that can be *generated* by grammar  $GG$

- Given a string  $SS$  in a language  $LL$ , find a rule configuration that generates  $SS$ . This is called **parsing**

## Formal grammar - Formal language

- A **formal grammar** is a (constrained) set of rules for forming strings in a language
  - The rules describe the **syntax** of a language: how to form strings from the symbols in  $\Sigma$
  - $GG$  restricts  $LL$  to some subset of  $\Sigma^*$
- A **formal language** is the set of all strings produced by a formal grammar

$$LL \subset \Sigma^*$$

- We can classify formal languages according to the *constraints* placed on the grammar rules

# Grammaticality

- Grammars define formal languages
- i.e., the set of all sentences (strings of words) that can be derived by the grammar.
- Sentences in this set said to be (prescriptively) **grammatical** or **felicitous**.
- Sentences outside this set said to be (prescriptively) **ungrammatical** or **infelicitous**.

## Defining grammars

- A *grammar* is defined by the tuple

$$G = \langle V, \Sigma, S, R \rangle$$

- $V$  is a finite set of *variables* or *preterminals*
- $\Sigma$  is a finite set of symbols, called *terminals*
- $S$  is in  $V$  and is called the *start symbol*
- $R$  is a finite set of *productions*, which are *rules* of the form

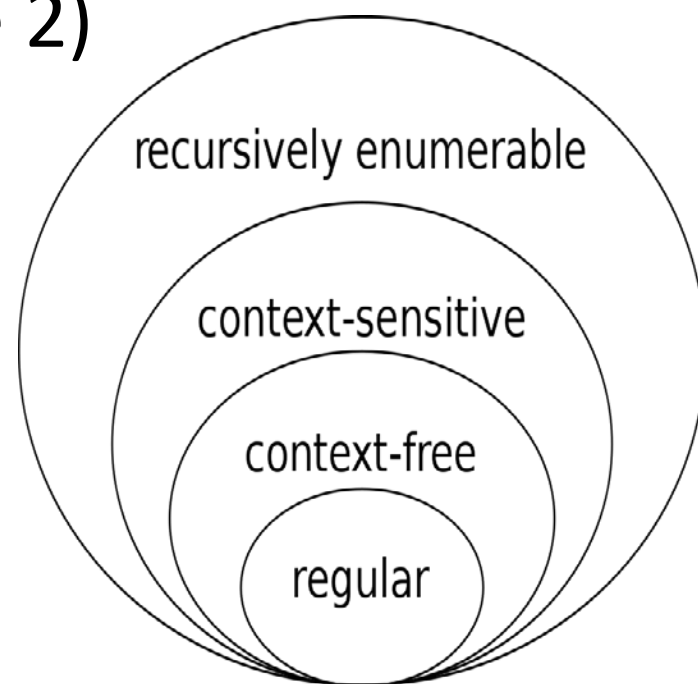
$$\alpha \rightarrow \beta$$

where  $\alpha$  and  $\beta$  are strings consisting of terminals and variables.

# Types of formal grammars

- Unrestricted, recursively enumerable (type 0)
- Context-sensitive (type 1)
- Context-free grammars (type 2)
- Regular grammars (type 3)

The Chomsky hierarchy



# Types of formal grammars

- Unrestricted, recursively enumerable (type 0)

$$\alpha\alpha \rightarrow \beta\beta$$

$\alpha\alpha$  and  $\beta\beta$  are any string of terminals and non-terminals

- Context-sensitive (type 1)

$$\alpha\alpha XX\beta\beta \rightarrow \alpha\alpha\alpha\alpha\beta\beta$$

$XX$  is a nonterminal;  $\alpha\alpha, \beta\beta, \alpha\alpha$  are any string of terminals and non-terminals;  $\alpha\alpha$  may not be empty.

- Context-free grammars (type 2)

$$XX \rightarrow \alpha\alpha$$

$XX$  is a nonterminal;  $\alpha\alpha$  are any string of terminals and non-terminals

- Regular grammars (type 3)

$$XX \rightarrow \alpha\alpha YY$$

(i.e. a *right* regular grammar)

$XX, YY$  are nonterminals;  $\alpha\alpha$  is any string of terminals;  $YY$  may be absent.

increasing restriction





# Equivalence with automata

- Formal grammar classes are defined according to the type of automaton that can accept the language
- There are various types of automata, and many correspond to certain types of formal grammars

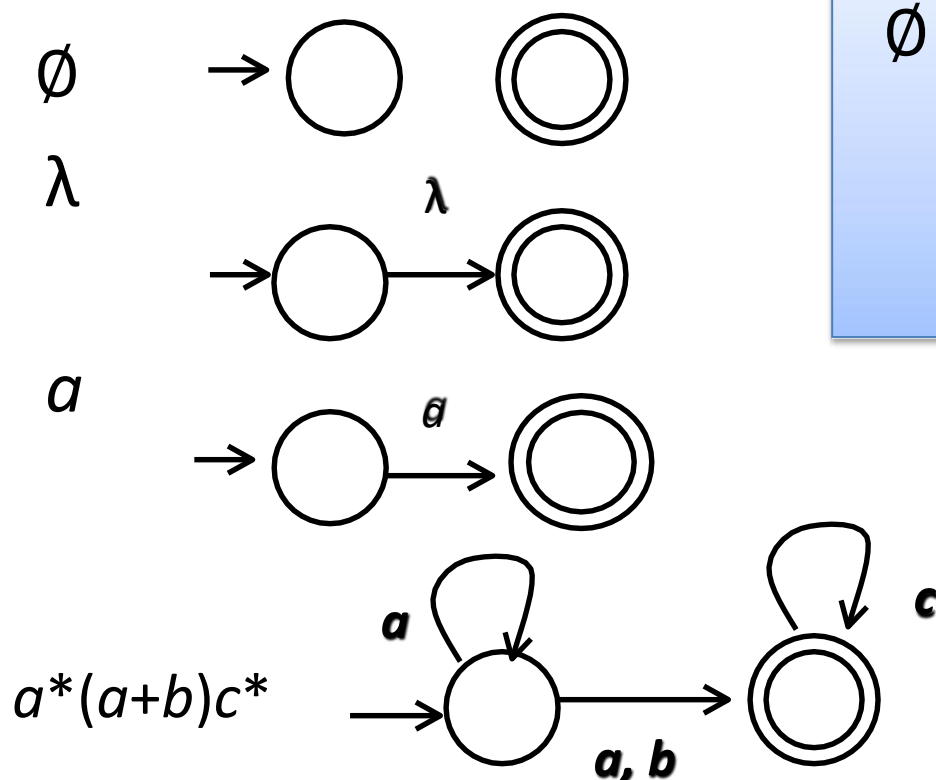
Grammar type	Accepted by
Recursively Enumerable	Turing machine
Context sensitive	non-deterministic linear bounded automaton (LBA)
Context Free	non-deterministic pushdown automaton (NDPA)
Regular	Finite State Automaton

# Regular languages

- Let's look at the most restricted case
- A **regular language** (over an *alphabet*  $\Sigma$ ) is any language for which there exists a finite state machine (finite automaton) that recognizes (accepts) it
- This class is equivalent to regular expressions (but no capture groups)

# Regular Expressions and Regular Languages

- There are simple finite automata corresponding to the simple regular expressions:



$\emptyset$  - the empty set  
 $\lambda$  - the empty string  
 $aa \in \Sigma$

Each of these has an initial state and one accepting state.

# Regular grammars

(type 3)

- Regular grammars can be generated by FSMs
- Equivalence with RegEx
- Definition of a *right-regular grammar*:
  - Every rule in  $R$  is of the form
    - $A \rightarrow aB$  or
    - $A \rightarrow a$  or
    - $S \rightarrow \lambda$  (to allow  $\lambda$  to be in the language)
    - where  $A$  and  $B$  are variables (perhaps the same, but  $B$  can't be  $S$ ) in  $V$
    - and  $a$  is any terminal symbol

# Regular grammars

(type 3)

- Example:

$$VV = \{SS, AA\}$$

$$\Sigma = \{aa, bb, cc\}$$

$$RR = \{SS \rightarrow aaSS, SS \rightarrow bbAA, AA \rightarrow \lambda\lambda, AA \rightarrow ccAA\}$$

$$SS = SS$$

RegEx:  $a^*bc^*$

- Cannot express  $aa^{nn}bb^{nn}$

# Parsing regular grammars

(type 3)

- Parsingspace:  $O(1)$
- Parsing time:  $O(n)$

## Pumping lemma

- Two slides ago, we asserted that a regular language cannot express  $aa^{nn}bb^{nn}, nn \geq 0$ . How would you prove this?
- The **pumping lemma for regular languages** describes a property of all regular languages
  - Some other language classes have pumping lemmata too
- One way to prove that a particular language is not regular is to demonstrate that a string of the language does *not* satisfy this pumping lemma

# Pumping lemma for regular languages

- Specifically, the pumping lemma says that any regular language has a value  $p \geq 1$  such that any string in the regular language  $L$  can be decomposed into

$$xyy^i, \quad |y| \geq 1, \quad |xy| \leq p, \quad i \geq 0$$

such that

$$xyy^i \in L,$$

- Since there's no way to satisfy this with the string  $aaaaabbbbbbb$ , the language  $a^n b^n$  is not regular.

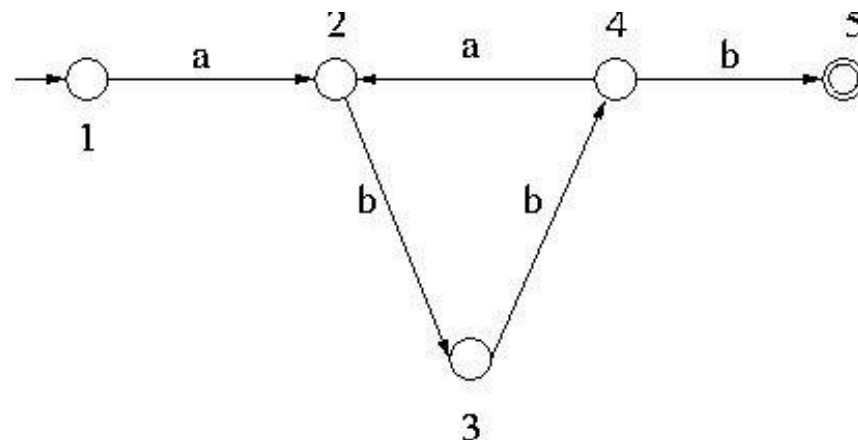


## Pumping lemma

- Every string must have a non-empty “middle section” which can be repeated an arbitrary number of times, giving new strings which are all in the language
- Since there’s no way to satisfy this with the string  $aaaaabbabbbb$ , the language  $aa^{nn}bb^{nn}$  is not regular.

# Pumping lemma

- For any string in a regular language, there should be a part somewhere within the first  $nn$  characters that can be pumped
- Informally, this means that, if there is a loop in the automaton, you can keep going around it as many times as you like and still be generating acceptable strings



NFA accepting  $a(bba)^*bbb$

# Pumping lemma

- Show that the language  $aa^{nn}bb^{nn}$  is not a regular language

aaabbb

try  $pp = 2$

a a abbb

a aa abbb



try  $pp = 3$

aa a bbb

aa aa bbb



try  $pp = 4$

aaa b bb

aaa bb bb



$\forall$  regular languages  $LL, \exists pp :$

$ww = xxyyxx, ww \in LL$

$|y| \geq 1, |xxy| \leq pp$

$\forall ii \geq 0: xxyy^{ii}xx \in LL$

# Linear grammars

(between types 3 and 2)

- Relax regular grammars slightly to defeat that pumping lemma
- Proper superset of regular grammars (type 3)
- Proper subset of context-free grammars (type 2)

$$V = \{S\}$$

$$\Sigma = \{aa, bb\}$$

$$R = \{SS \rightarrow aaSSbb, SS \rightarrow \varepsilon\}$$

$$S = S$$

ab, aabb, aaabbb, aaaabbbb, ...

# Context-free grammar

(type 2)

- What if we allow our production rule in a linear grammar to have more than one nonterminal?

...we get the most important type of grammar studied in linguistics: the **context-free grammar** (CFG)

# Context free grammar

(type 2)

$$GG = (VV, \Sigma, RR, SS)$$

$VV = \{ \text{pre-terminals } vv_0, vv_1, \dots \}$

$\Sigma = \{ \text{terminals } ww_0, ww_1, \dots \}$

$RR = \{ \text{rules } r_i: vv \rightarrow \alpha\alpha \}$

$SS = \text{start symbol, } SS \in VV$

$\alpha\alpha$ : sequence of terminals and pre-terminals (or  $\emptyset$ )

$LL$ : the language generated by  $GG$

## Context-sensitive grammars

(type 1)

- No rule can make a string shorter
- Can be accepted by a 'linear bounded automaton' (LBA), a nondeterministic Turing machine which uses space  $O(n)$

## Unrestricted grammar

(type 0)

- The most general class
- Recognized by Turing machine
- Generate recursively-enumerable languages

# Example CFG

$GG = (VV, \Sigma, RR, \$S$   
 $VV = \{ SS, NNPP, NNOONN, VVPP, DDDDtt, NNttNNnn, VVDDrrbb, AANNxx \}$   
 $\Sigma = \{ tt\textcolor{red}{t}aatt, tt\textcolor{red}{t}iitt, aa, tt\textcolor{red}{t}DD, mmaann, bbttttbb, ffffiiff\textcolor{red}{t}tt, mmDDaaff, iinnccffNN\textcolor{red}{i}iDD, rrDDaaii, iittDDtt \}$   
 $SS = SS$

$R = \{$   

$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a \mid the$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid man$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid read$
$NP \rightarrow Det NOM$	$Aux \rightarrow does$
$NOM \rightarrow Noun$	
$NOM \rightarrow Noun NOM$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	

  
 $\}$



# Grammar rewrite rules

S --> NP VP

--> Det NOM VP

--> The NOM VP

--> The Noun VP

--> The man VP

--> The man Verb NP

--> The man read NP

--> The man read Det NOM

--> The man read this NOM

--> The man read this Noun

--> The man read this book

S → NP VP

S → Aux NP VP

S → VP

NP → Det NOM

NOM → Noun

NOM → Noun NOM

VP → Verb

VP → Verb NP

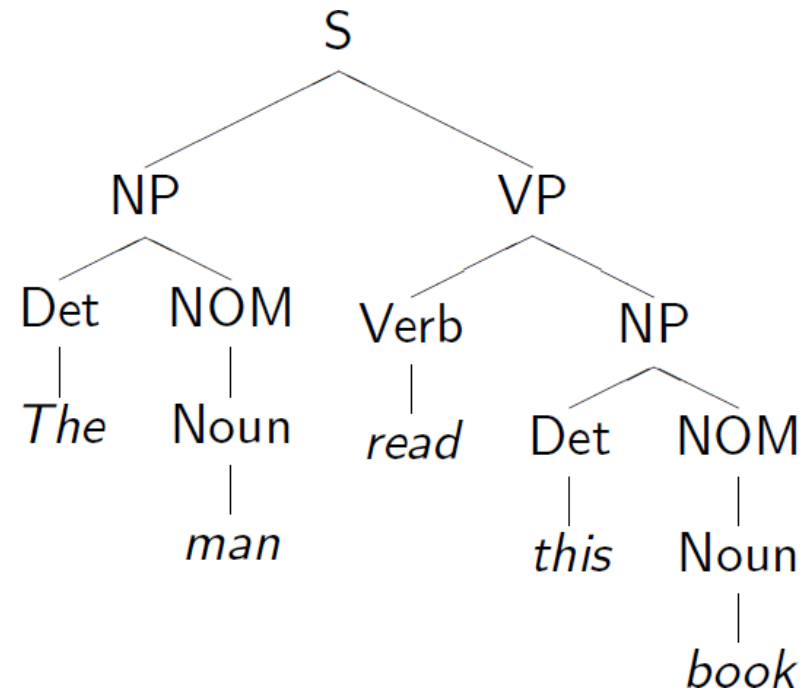
Det → *that* | *this* | *a* | *the*

Noun → *book* | *flight* | *meal* | *man*

Verb → *book* | *include* | *read*

Aux → *does*

# Parse tree



# PCFG

- Probabilistic context-free grammar
- Adds probabilities to each rule
- Each distinct left-hand-side gets a probability mass 1.0
- Rule weights can be estimated from corpora

# CFGs can express recursion

- Example of seemingly endless recursion of embedded prepositional phrases:

PP  $\rightarrow$  Prep NP

NP  $\rightarrow$  Noun PP

[S The mailman ate his [NP lunch [PP with his friend [PP from the cleaning staff [PP of the building [PP at the intersection [PP on the north end [PP of town]]]]]]]].

Most programming languages are type-2 grammars (CFGs)

# Chomsky Normal Form

- Any CFG can be converted into a form where all rules are of the form:

$$XX \rightarrow YYY$$

$$XX \rightarrow aa$$

$$SS \rightarrow \lambda$$

(S is the only terminal that can go to the empty string)

Convert  $WW \rightarrow XX YYaa YY$  to Chomsky Normal Form

# CNF conversion

- Steps:
  1. Make  $S$  non-recursive
  2. Eliminate  $\lambda\lambda$  (except  $SS \rightarrow \lambda\lambda$ )
  3. Eliminate all chain rules
  4. Remove unused symbols

Convert the following  
grammar to Chomsky  
Normal Form:

$$SS \rightarrow AASSAA$$
$$SS \rightarrow aaBB$$
$$AA \rightarrow BB$$
$$AA \rightarrow SS$$
$$BB \rightarrow bb$$
$$BB \rightarrow \lambda\lambda$$

# CNF conversion

$SS \rightarrow AASSAA$   
 $SS \rightarrow aaBB$   
 $AA \rightarrow BB$   
 $AA \rightarrow SS$   
 $BB \rightarrow bb$   
 $BB \rightarrow \lambda\lambda$

$SS \rightarrow AASSAA$   
 $SS \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow SS$   
 $BB \rightarrow bb$   
 $BB \rightarrow \lambda\lambda$   
 $UU_{aa} \rightarrow aa$

$SS \rightarrow AAXX$   
 $SS \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow SS$   
 $BB \rightarrow bb$   
 $BB \rightarrow \lambda\lambda$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow SSAA$

$SS \rightarrow SS'$   
 $SS' \rightarrow AAXX$   
 $SS' \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow SS'$   
 $BB \rightarrow bb$   
 $BB \rightarrow \lambda\lambda$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow SS'AA$

$SS \rightarrow SS'$   
 $SS' \rightarrow AAXX$   
 $SS' \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow SS'$   
 $BB \rightarrow bb$   
  
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow SS'AA$   
 $XX \rightarrow SS'$   
 $SS' \rightarrow XX$

$SS \rightarrow SS'$   
 $SS' \rightarrow AAXX$   
 $SS' \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow SS'$   
 $BB \rightarrow bb$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow SS'AA$   
 $XX \rightarrow SS'$   
 $SS' \rightarrow XX$

$SS \rightarrow XX$   
 $XX \rightarrow AAXX$   
 $XX \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow XX$   
 $BB \rightarrow bb$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow XXAA$   
 ~~$XX \rightarrow XX$~~   
 ~~$SS' \rightarrow XX$~~

~~$SS \rightarrow XX$~~   
 $XX \rightarrow AAXX$   
 $XX \rightarrow UU_{aa}BB$   
 $AA \rightarrow BB$   
 $AA \rightarrow XX$   
 $BB \rightarrow bb$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow XXAA$   
 $SS \rightarrow AAXX$   
 $SS \rightarrow UU_{aa}BB$   
 $SS \rightarrow XXAA$

$XX \rightarrow AAXX$   
 $XX \rightarrow UU_{aa}BB$   
 $AA \rightarrow bb$   
 $AA \rightarrow XX$   
 $BB \rightarrow bb$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow XXAA$   
 $SS \rightarrow AAXX$   
 $SS \rightarrow UU_{aa}BB$   
 $SS \rightarrow XXAA$

$XX \rightarrow AAXX$   
 $XX \rightarrow UU_{aa}BB$   
 $AA \rightarrow bb$   
 ~~$AA \rightarrow XX$~~   
 $BB \rightarrow bb$   
 $UU_{aa} \rightarrow aa$   
 $XX \rightarrow XXAA$   
 $SS \rightarrow AAXX$   
 $SS \rightarrow UU_{aa}BB$   
 $SS \rightarrow XXAA$   
 $AA \rightarrow AAXX$   
 $AA \rightarrow UU_{aa}BB$   
 $AA \rightarrow XXAA$

all done




# Parsing context-free grammars (type 2)

- CFGs are widely used to represent surface syntax in natural languages
- Space complexity:
  - you'll need at least one stack
  - space use will depend on the amount of recursion in the input
- Time complexity
  - Generally  $O(n^3)$

# Parsing

- The opposite of generation: find the structure from a string
- Essentially a search problem
- Find all structure that match an input string
- Two approaches
  - Bottom-up
  - Top-down



This refers to the parse tree, not the search tree. So either method can be done breadth-first or depth-first

## Recognizer v. parser

- **Recognizer** (acceptor) is a program that determines whether a sentence is accepted by the grammar or not
- A **parser** determines this as well, and if the sentence is accepted, it also returns the structural configuration(s) of grammar rules for the sentence
- Some parsing systems may also produce compositional semantics

# Soundness and completeness

- Correctness: a parser is **sound** if every parse it returns is correct
- A parser **terminates** if it is guaranteed not to enter an infinite loop
- A parser is **complete** for grammar  $GG$  and sentence  $SS$  if it is sound, produces every possible parse for  $SS$ , and terminates
- Often, we settle for sound but incomplete parsers
  - probabilistic parsers may be able to return the *bb*-best parses

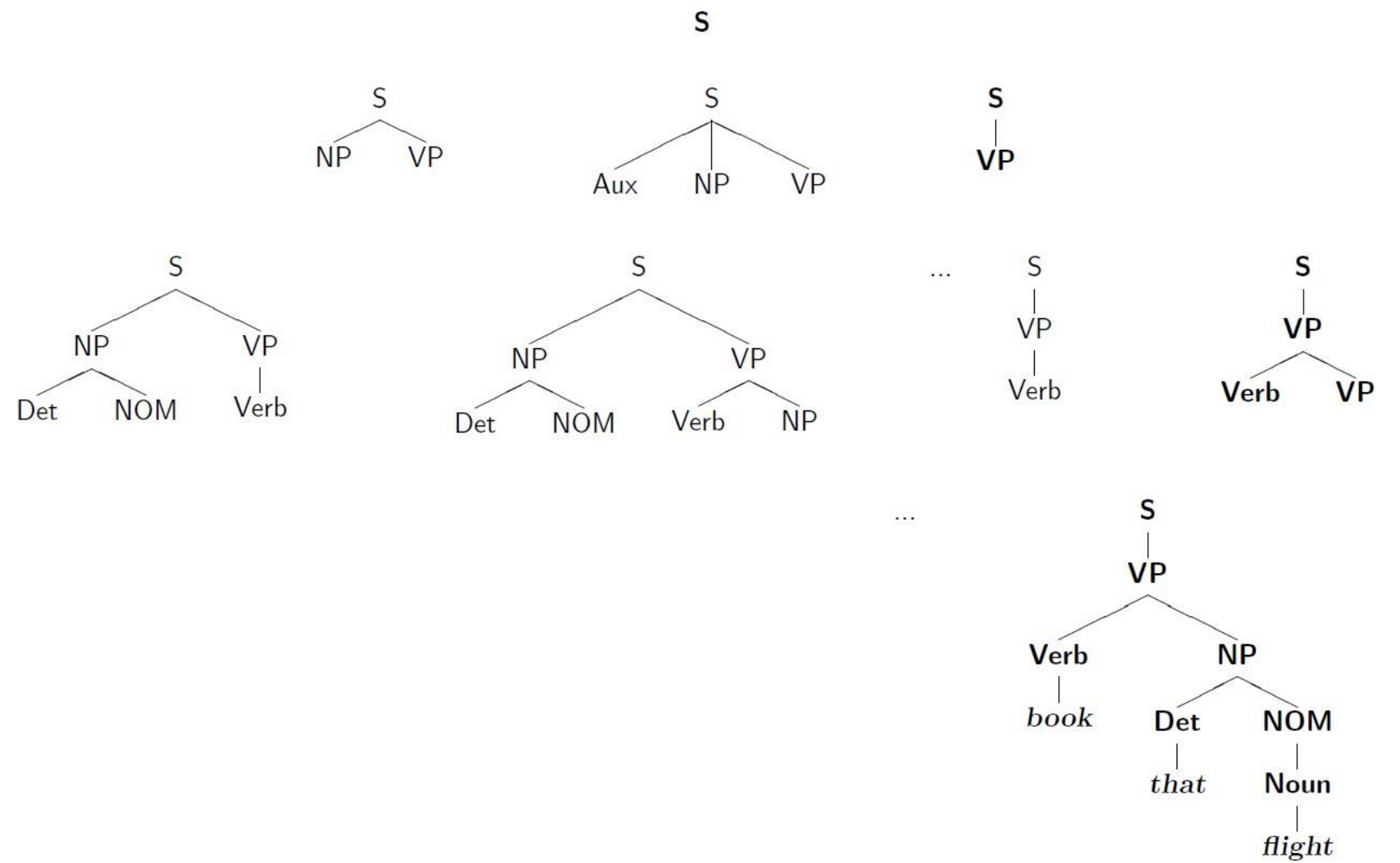
# Top-down parsing

- Create a list of goal constituents
- Rewrite goals by matching a goal on the list with the left-hand-side of a rule
- Replace with the right-hand-side
- If there is more than one match to the LHS, try different rules (breadth-first or depth-first)

## example

$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a \mid the$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid man$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid read$
$NP \rightarrow Det NOM$	$Aux \rightarrow does$
$NOM \rightarrow Noun$	
$NOM \rightarrow Noun NOM$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	

*Book that flight.*



# Problems with top-down parsing

- Left-recursive rules lead to infinite recursion  
$$NNPP \rightarrow NNPP PPPP$$
- Poor performance when there are many matches for an LHS
  - If there are many rules for S, there's no way to eliminate irrelevant ones. In other words, it does the useless work of expanding things that there is no evidence for
- Doesn't work at the terminals (lexemes)
- Can't make use of common substructure



## Bottom-up parsing

- Bottom-up parsing is data-directed
- Start with the string to be parsed
- Match right-hand-sides, condense to LHS
  - Still need to choose when there are multiple possible matches for the RHS
  - Can use breadth-first or depth-first search
- Parsing is complete when all you have left is the start symbol

## example

$S \rightarrow NP VP$	$Det \rightarrow that \mid this \mid a \mid the$
$S \rightarrow Aux NP VP$	$Noun \rightarrow book \mid flight \mid meal \mid man$
$S \rightarrow VP$	$Verb \rightarrow book \mid include \mid read$
$NP \rightarrow Det NOM$	$Aux \rightarrow does$
$NOM \rightarrow Noun$	
$NOM \rightarrow Noun NOM$	
$VP \rightarrow Verb$	
$VP \rightarrow Verb NP$	

*Book that flight.*

## Shift-reduce parsing

Stack	Input remaining	Action
( )	Book that flight	shift
(Book)	that flight	reduce, Verb $\rightarrow$ book, (Choice #1 of 2)
(Verb)	that flight	shift
(Verb that)	flight	reduce, Det $\rightarrow$ that
(Verb Det)	flight	shift
(Verb Det flight)		reduce, Noun $\rightarrow$ flight
(Verb Det Noun)		reduce, NOM $\rightarrow$ Noun
(Verb Det NOM)		reduce, NP $\rightarrow$ Det NOM
(Verb NP)		reduce, VP $\rightarrow$ Verb NP
(Verb)		reduce, S $\rightarrow$ V
(S)		SUCCESS!

Ambiguity may lead to the need for backtracking.

# Shift-reduce parser

- Start with the sentence in an input buffer
  - Shift: push the next input symbol onto the stack
  - Reduce: if a RHS matches the top elements of the stack, pop those elements off and push the LHS
- If either shift or reduce are possible, choose arbitrarily
- If you end up with only the start symbol on the stack, you have a parse
- Otherwise, you can backtrack

## Shift-reduce parser

- In the top-down parser, the main decision was which production rule to pick
  - In a bottom-up shift-reduce parser, the decisions are:
    - Should we shift, or reduce
    - If we reduce, then by which rule
- Both of these decisions can be revisited when backtracking

# Problems with bottom-up parsing

- No obvious way to generate structures that generate empty surface
- Lexical ambiguity can explode the search space
- Useless constituents can be built locally

Top-down and bottom-up parsers can both be extremely inefficient on real-world NLP parsing problems. Complexity may approach  $OO(bb^{nn})$  in the sentence length.

# Parsing is hard

- Left-recursive structures must be found, not predicted
- Empty categories must be predicted, not found
- When backtracking, don't redo any work
- Get linguists on board to think about computability from the start

# Next time

- Clustering
- Classifiers
- Overview of Information Theory



# C# Tutorial (continued...)

# IEnumerable, yield, and deferred execution

- Before describing the trie data structure, let's look at iterators which enumerate a sequence of elements

Examples in C#. If you use another language, it will be instructive to think about how to adapt the solutions to your language

- Enumeration is obvious when the data is at hand and you want to use it all:

```
String[] data = { "able", "bodied", "cows", "don't", "eat", "fish" };  
  
foreach (String s in data)  
    Console.WriteLine(s);
```

We can pass (a reference to) the array around too, no problem

```
String[] data = { "able", "bodied", "cows", "don't", "eat", "fish" };  
// ...  
ProcessSomeStrings(data);  
// ...  
  
void ProcessSomeStrings(String[] the_strings)  
{  
    foreach (String s in the_strings)  
        Console.WriteLine(s);  
}
```

## What if we only want to “process” the four-letter words?

```
String[] data = { "able", "bodied", "cows", "don't", "eat", "fish" };  
// ...
```

```
List<String> filtered = new List<String>();  
foreach (String s in data)  
    if (s.Length == 4)  
        filtered.Add(s);  
ProcessSomeStrings(filtered);  
// ...
```

This doesn't seem very nice. For one thing, we have to use more memory and waste time copying the elements we care about to a new list.

```
void ProcessSomeStrings(List<String> the_strings)  
{  
    foreach (String s in the_strings)  
        Console.WriteLine(s);  
}
```

Is there a way to pass this function enough information to filter the *original list* itself, where it lies?

- Remember the non-filtered example for a second

```
void ProcessSomeStrings(String[] the_strings)
{
    foreach (String s in the_strings)
        Console.WriteLine(s);
}
```

- The processing function doesn't really care about the fact that the data is in an array
- This violates an important programming maxim:

A flexible interface *demands the least* and *provides the most*:

- Inputs* are as general as possible (allowing clients to supply any level of specificity, i.e. be lazy)
- Outputs* are as specific as possible (allowing clients to capitalize on work products, i.e. be lazy).

```
void ProcessSomeStrings(String[] the_strings)
{
    foreach (String s in the_strings)
        Console.WriteLine(s);
}
```

The extra (unused) demands this function is making by asking for String[]:

- That the strings all be in memory at the same time
  - That the strings be randomly accessible by an index
  - That the number of strings be known and fixed before the function starts
- 
- To modify this to comply with the maxim, we first ask:
  - Q: What is the absolute minimum that this function actually needs to accomplish it's work?
  - Answer: a way to iterate strings

# Interfaces

- `IEnumerable<T>` is one of many system-defined **interfaces** that a class can elect to implement

An **interface** is a named set of zero or more function signatures with no implementation(s)

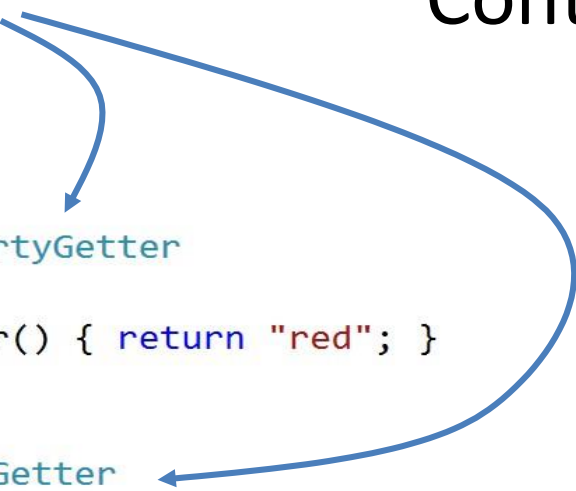
- To implement an interface, a class defines a matching implementation for every function in the interface
- Interfaces are sometimes described as contracts
- You can define and use a reference to an interface just like any other object reference

## Contrived Example

```
interface IPropertyGetter
{
    String GetColor();
}

class Strawberry : IPropertyGetter
{
    public String GetColor() { return "red"; }
}

class Ferrari : IPropertyGetter
{
    public String GetColor() { return "yellow"; }
}
```



- This looks like C++ class inheritance
  - yes, but it's more ad-hoc
  - C# classes can have **single inheritance** of other classes, and **multiple inheritance** of interfaces
  - Interfaces can inherit from other interfaces (not shown)



# IEnumerable<T>

- This is one of the simplest interfaces defined in the BCL (base class libraries)
- This interface provides just one thing: a way to iterate over elements of type T
- All of the system arrays, collections, dictionaries, hash sets, etc. implement IEnumerable<T>
  - Implementing IEnumerable<T> on your own classes can be very useful, but you don't need to worry about that
  - For now, what's important is that you get to use it, because it's available on all of the system collections

# IEnumerator<T>

- **IEnumerable<T>** has only one function, which allows a caller or caller(s) to obtain an *enumerator* object which is able to iterate over elements
  - The actual enumerator object is an object that implements a different interface, called **IEnumerator<T>**
  - This “factory” design allows a caller to initiate and maintain several simultaneous iterations if needed
  - The enumerator object, **IEnumerator<T>** can only:
    - Get the current element
    - Move to the next element
    - Tell you if you’ve reached the end
  - Note: There’s no count
    - ICollection inherits from IEnumerable to provide this

# Interfaces as function arguments

- Using interfaces as function arguments allows you to require the absolute minimum functionality the function actually needs
- In this way, the ad-hoc nature of interfaces allows us to comply with the maxim

```
void ProcessSomeStrings(IEnumerable<String> the_strings)
{
    foreach (String s in the_strings)
        Console.WriteLine(s);
}
```

Now, this function is exposing the **weakest (most general) requirement** possible for the processing it has to do. This provides more flexibility to callers since they can choose whatever level of specificity is convenient. The function can be used in the widest possible variety of situations.

## Example

```
String[] d1 = { "able", "bodied", "cows", "don't", "eat", "fish" };  
ProcessSomeStrings(d1);
```


```
List<String> d2 = new List<String> { "clifford", "the", "big", "red", "dog" };  
ProcessSomeStrings(d2);
```

```
HashSet<String> d3 = new HashSet<String> { "these", "must", "be", "distinct" };  
ProcessSomeStrings(d3);
```

```
Dictionary<String,int> d4 =  
    new Dictionary<String, int> { { "the", 334596 }, { "in", 153024 } };  
ProcessSomeStrings(d4.Keys);
```

```
void ProcessSomeStrings(IEnumerable<String> the_strings)  
{  
    foreach (String s in the_strings)  
        Console.WriteLine(s);  
}
```

Python users might not be impressed, but the difference is that this is all 100% strongly typed



# Iteration is efficient

- That's cool, `IEnumerable<T>` lets a function **not care** about where a sequence of elements is coming from
  - We don't copy the elements around
  - Iterators let us access elements right from their source
- All of those examples iterate over elements that **already exist** somewhere
- Is there a way to iterate over data that's generated on-the-fly, doesn't exist yet, or is never persisted at all?
- Yes!

# Iterating over on-the-fly data

```
IEnumerable<String> GetNewsStories(int desired_count)
{
    for (int i = 0; i < desired_count; i++)
        yield return RealtimeNewswireSource.GetLatestStory();
}
```

see next slide

```
// ...
IEnumerable<String> d5 = GetNewsStories(7);
ProcessSomeStrings(d5);
// ...
```

This is exactly the same as before, but this time there's no "collection" of elements sitting anywhere

```
void ProcessSomeStrings(IEnumerable<String> the_strings)
{
    foreach (String s in the_strings)
        Console.WriteLine(s);
}
```

This function doesn't care. In fact, it can't even tell.

# yield keyword

- The **yield** keyword makes it easy to define your own custom iterator functions
- Any function that contains the `yield` keyword becomes special
  - It must be declared as returning an `IEnumerable<T>`
  - Deferred execution means that the function's body is not necessarily invoked when you "call" it
  - It must deliver zero or more elements of type `T` using:  
`yield return t;`
  - Sometime later, control may continue immediately after this statement to allow you to yield additional elements
  - It may signal the end of the sequence by using:  
`yield break;`



# Custom iterator function example

```
IEnumerable<String> GetNewsStories(int desired_count)
{
    for (int i = 0; i < desired_count; i++)
        yield return RealtimeNewswireSource.GetLatestStory();
}
```

code from this custom iterator function is *not* executed at this point.

```
// ...
IEnumerable<String> d5 = GetNewsStories(7);
ProcessSomeStrings(d5);
// ...
```

d5 refers to an iterator that “knows how” to get a certain sequence of strings when asked

```
void ProcessSomeStrings(IEnumerable<String> the_strings)
{
    foreach (String s in the_strings)
        Console.WriteLine(s);
}
```

This finally demands the strings, causing our custom iterator function to execute—interleaved with this loop!



# Closures

- Lambda expressions automatically capture local variables that they reference, which are then passed around as part of the lambda variable
  - Caution: languages do this differently with respect to reference (the lambda expression will modify the original value) versus value (the lambda expression has a snapshot of the value)
- This can lead to interesting scoping issues

# Lambda expressions

```
// recall Select(ch => ('a' <= ch && ch <= 'z') || ch == '\\' ? ch : ' ');

String s = "Al's 20 fat-ish oxen.";
Func<Char, Char> myfunc = (ch) => ('a' <= ch && ch <= 'z') || ch == '\\' ? ch : ' ';
IEnumerable<Char> iech = s.Select(myfunc);
// iech is now a deferred enumerator for the characters in: " l's    fat ish oxen "

Func<Char, Char> myfunc = (ch) =>
{
    if ('a' <= ch && ch <= 'z')
        return ch;
    if (ch == '\\')
        return ch;
    return ' ';
};

Func<String, int, bool> string_is_longer_than = (s, i) => s.Length > i;

bool b = string_is_longer_than("hello", 3);    // true
```

# Closure example

```
int x = 3;  
Action a = () =>  
    {  
        Console.WriteLine(x);  
    };  
a();           // prints 3  
x = 5;  
a();           // prints 5
```

# State machine example

```
using System;
using System.Collections.Generic;
using System.Linq;

static class Program
{
    static class MainClass
    {
        enum State { Zero, One, Two };

        static Dictionary<State, Func<Char, State>> machine = new Dictionary<State, Func<Char, State>>
        {
            {
                State.Zero, (ch) => { return State.Two; }
            },
            {
                State.One, (ch) => { return State.One; }
            },
        };

        static void Main(String[] args)
        {
            String s = "the string to parse";
            int i = 0;

            State state = State.Zero;
            while (i < s.Length)
                state = machine[state](s[i++]);

        }
    }
}
```

A dictionary  
of lambda  
functions

The state machine

## LINQ in C#

- Sequences: `IEnumerable<T>`
- Deferred execution
- Type inference
- Strong typing
  - despite the 'var' keyword
  - (C# 4.0 now has the 'dynamic' keyword, which allows true runtime typing where desired)

# LINQ operators

- Filter/Quantify:  
Where, ElementAt, First, Last, OfType
- Aggregate:  
Count, Any, All, Sum, Min, Max
- Partition/Concatenate:  
Take, Skip, Concat
- Project/Generate:  
Select, SelectMany, Empty, Range, Repeat
- Set:  
Union, Intersect, Except, Distinct
- Sort/Ordering:  
OrderBy, ThenBy, OrderByDescending, Reverse
- Convert/Render:  
Cast, ToArray, ToList, ToDictionary

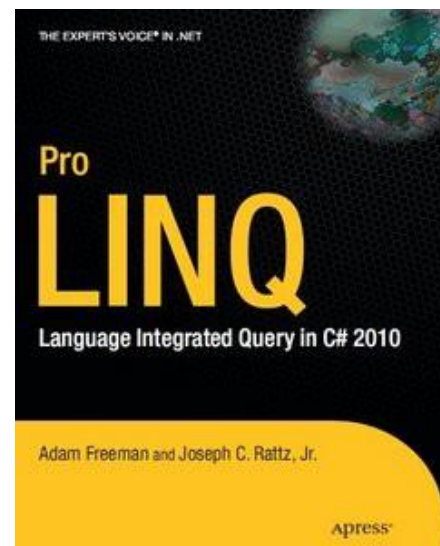
# Example

```
Dictionary<String, int> pet_sym_map =  
    File.ReadAllLines("/programming/analytical-grammar/erg-funcs/pet-symbol-key.txt")  
    .Select(s => s.Split(sc7, StringSplitOptions.RemoveEmptyEntries))  
    .Where(rgs => rgs.Length == 3)  
    .Select(rgs => new { id = int.Parse(rgs[0]), sym = rgs[1].Trim().ToLower() })  
    .GroupBy(a => a.sym)  
    .Select(grp => grp.ArgMin(a => a.id))  
    .ToDictionary(a => a.sym, a => a.id);
```

# C# LINQ (Language Integrated Query)

- Declarative operations on sequences
- Recommendation:

Joseph C. Rattz, Jr. (2007) *Pro LINQ: Language Integrated Query in C# 2008*. Apress.






# Programming Paradigms: Procedural

- Procedural (“imperative”) programming
  - FORTRAN (1954) grew out of hardware assembly languages, which are necessarily procedural
  - We explicitly specify the (synchronous) steps for doing something (i.e. an algorithm)

```
int Factorial(int n)
{
    int f = 1;
    for (int i = n; i > 1; i--)
        f *= i;
    return f;
}
```

# Functional Programming

- A type of declarative programming
  -  Constraint-based syntax formalisms such as unification grammars (LFG, HPSG, ...) are also declarative
- Like function definitions in math, the “program” describes asynchronous relationships
- Functions are stateless and should have no side-effects
- Immutable values
  - As a bonus, this really facilitates concurrent programming
- Scheme, Haskell, F#

# F# Example

- F# interactive on patas:

```
$ mono /opt/fsharp/bin/fsi.exe --gui-  
      (you must use an ANSI terminal)
```

```
> let rec factorial = function  
    | 0 -> 1  
    | n -> n * factorial(n -  
1);; val factorial : int -> int  
> factorial 5;;  
val it : int =  
120  
> #quit;;
```