

+

Natural Language Processing 2

CS151
David Kauchak

+



Admin

- Project reports
- Keep up with the homework
- NLP spring course?

+

Language translation

Yo quiero Taco Bell

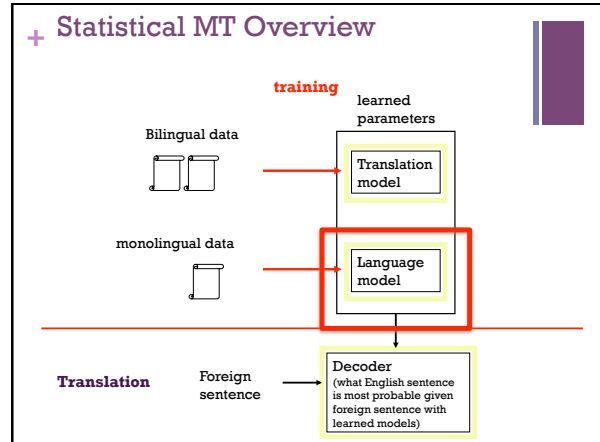


+ Noisy channel model

model $p(e | f) \propto p(f | e)p(e)$

translation model language model

how do foreign sentences get translated to English sentences? what do English sentences look like?



+ Language modeling

Step 1: decompose the probability

$P(\text{I think today is a good day to be me}) =$

- $P(\text{I}) \times$
- $P(\text{think} | \text{I}) \times$
- $P(\text{today} | \text{I think}) \times$
- $P(\text{is} | \text{I think today}) \times$
- $P(\text{a} | \text{I think today is}) \times$
- $P(\text{good} | \text{I think today is a}) \times$
- ...

+ The n-gram Approximation

Assume each word depends only on the previous n-1 words (e.g. trigram: three words total)

- $P(\text{is} | \text{I think today}) \approx P(\text{is} | \text{think today})$
- $P(\text{a} | \text{I think today is}) \approx P(\text{a} | \text{today is})$
- $P(\text{good} | \text{I think today is a}) \approx P(\text{good} | \text{is a})$

+ Estimating probabilities

- How do we find probabilities? $P(\text{is} | \text{think today})$
- Get real text, and start counting!

$$P(\text{is} | \text{think today}) = \frac{\text{count}(\text{think today is})}{\text{count}(\text{think today})}$$

+ Smoothing

Is this sentence reasonable?

The Singing Toadstools are a new band.

$$P(\text{are} | \text{Singing Toadstools}) = \frac{\text{count}(\text{Singing Toadstools are})}{\text{count}(\text{Singing Toadstools})}$$

"singing toadstools are" Advanced search

⚠ No results found for "singing toadstools are".

+ Smoothing: Add One (Laplacian)

- Add one smoothing: $P(c | ab) \approx \frac{C(abc) + 1}{C(ab) + V}$
- Works very badly. DO NOT DO THIS
- Add delta smoothing: $P(c | ab) \approx \frac{C(abc) + \delta}{C(ab) + \delta V}$
- Still very bad. DO NOT DO THIS

+ A better idea

Trigram: $P(\text{is} | \text{think today})$

Bigram: $P(\text{is} | \text{today})$

Unigram: $P(\text{is})$

How are these probabilities related?

Is one more likely to be 0 than another?

Could we combine them somehow?

+ Two general approaches

■ Interpolation

- $p^*(z | x, y) = \lambda p(z | x, y) + \mu p(z | y) + (1 - \lambda - \mu) p(z)$
- Combine the probabilities with some linear combination

■ Backoff

$$P(z | xy) = \begin{cases} \frac{C^*(xyz)}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P(z | y) & \text{otherwise} \end{cases}$$

- Combine the probabilities by "backing off" to lower models only when we don't have any information

+ Smoothing: Simple Interpolation

$$P(z | xy) \approx \lambda \frac{C(xyz)}{C(xy)} + \mu \frac{C(yz)}{C(y)} + (1 - \lambda - \mu) \frac{C(z)}{C(\bullet)}$$

- Trigram is very context specific, very noisy
- Unigram is context-independent, smooth
- Interpolate Trigram, Bigram, Unigram for best combination
- **How should we determine λ and μ ?**

+ Smoothing: Finding parameter values

- Split data into training, "heldout", test
- Try lots of different values for λ , μ on heldout data, pick best
- Two approaches for finding these efficiently
 - EM (expectation maximization)
 - "Powell search" – see Numerical Recipes in C

+ Smoothing: Jelinek-Mercer

- Simple interpolation:

$$P_{smooth}(z | xy) = \lambda \frac{C(xyz)}{C(xy)} + (1 - \lambda) P_{smooth}(z | y)$$

- Should all bigrams be smoothed equally?

About 4,370,000 results (0.11 seconds) [Advanced search](#)

About 47,200 results (0.11 seconds) [Advanced search](#)

+ Smoothing: Jelinek-Mercer

- Simple interpolation:

$$P_{smooth}(z | xy) = \lambda \frac{C(xyz)}{C(xy)} + (1 - \lambda)P_{smooth}(z | y)$$

- Smooth a little after "The Dow", more after "Adobe acquired"

$$P_{smooth}(z | xy) = \lambda(C(xy)) \frac{C(xyz)}{C(xy)} + (1 - \lambda(C(xy)))P_{smooth}(z | y)$$

+ Smoothing: Jelinek-Mercer continued

$$P_{smooth}(z | xy) = \lambda(C(xy)) \frac{C(xyz)}{C(xy)} + (1 - \lambda(C(xy)))P_{smooth}(z | y)$$

- Bin counts by frequency and assign a λ s for each bin
- Find λ s by cross-validation on held-out data

+ Backoff models: absolute discounting

$$P_{absolute}(z | xy) = \begin{cases} \frac{C(xyz) - D}{C(xy)} & \text{if } C(xyz) > 0 \\ \alpha(xy)P_{absolute}(z | y) & \text{otherwise} \end{cases}$$

- Subtract some absolute number from each of the counts (e.g. 0.75)
 - will have a large effect on low counts
 - will have a small effect on large counts

+ Kneser-Ney

- Idea: not all counts should be discounted with the same value

P(Francisco | eggplant) vs
P(stew | eggplant)

If we've never seen either, which should be more likely? why?

What would an interpolated/backoff model say?

What is the problem?

+ Kneser-Ney

- Idea: not all counts should be discounted with the same value
- “Francisco” is common, so backoff/interpolated methods say it is likely
 - But it only occurs in context of “San”
- “Stew” is common, and in many contexts
- **Weight backoff by number of contexts word occurs in**

$P(\text{Francisco} \mid \text{eggplant})$ **low**
 $P(\text{stew} \mid \text{eggplant})$ **higher**

+ Smoothing

- Lots of other approaches...
 - Good Turing: estimate unseen events based on 1-count events
 - Katz, Witten-Bell,...
- In practice, Kneser-Ney works very well (or minor modifications of it)

+ Other language model ideas?

- Skipping models: rather than just the previous 2 words, condition on the previous word and the 3rd word back, etc.
- Caching models: phrases seen are more likely to be seen again (helps deal with new domains)
- Clustering:
 - some words fall into categories (e.g. Monday, Tuesday, Wednesday...)
 - smooth probabilities with category probabilities
- Domain adaptation:
 - interpolate between a general model and a domain specific model

+ Language model evaluation

- We have two different language models (i.e. two different probability distributions over English)
- How can we determine which is better?
 - Idea 1: use it in our MT system and see which works better
 - Idea 2: should predict actual English sentences with high probability

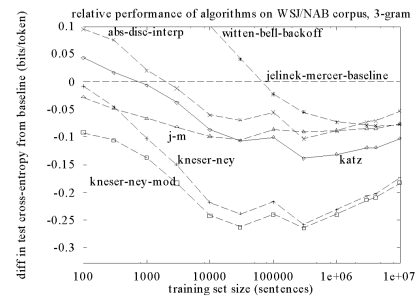
+ Perplexity

- Ask the two models to predict the likelihood of some test data
- The one with the higher probability is better
- Perplexity standardizes this idea by averaging over the probability of all words:

$$\max_q \sqrt[n]{\prod_{i=1}^n P(w_i | w_{1..i-1})} \cong \min \frac{\sum_{i=1}^n \log p(w_i | w_{1..i-1})}{n}$$

- Or... can be seen as the average log of the probability

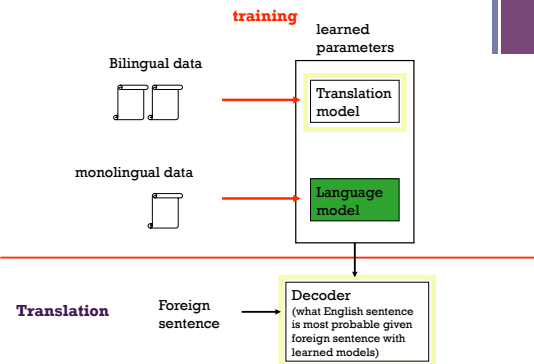
+ Smoothing results



+ Language Modeling Toolkits

- SRI
 - <http://www.speech.sri.com/projects/srilm/>
- CMU
 - http://www.speech.cs.cmu.edu/SLM_info.html

+ Statistical MT Overview



+ The Classic Translation Model

Word Substitution/Permutation
[IBM Model 3, Brown et al., 1993] $p(f|e)$

Generative story:

Mary did not slap the green witch

How do we get from English to Spanish?

Maria no dió una bofetada a la bruja verde

+ The Classic Translation Model

Word Substitution/Permutation
[IBM Model 3, Brown et al., 1993] $p(f|e)$

Generative story:

Mary did not slap the green witch

Mary not slap slap slap the green witch n(3|slap)

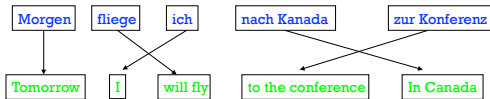
Mary not slap slap slap NULL the green witch P-Null

Maria no dió una bofetada a la verde bruja t(la|the)

Maria no dió una bofetada a la bruja verde d(j|i)

Probabilities can be learned from raw bilingual text.

+ Phrase-based translation model



- Phrase-based machine translation:
 - break English into phrases
 - probabilistically translate those phrases, $p(f_{phrase} | e_{phrase})$
 - with some probability, phrases can be reordered

How can we learn these?

+ How to Learn the Phrase Translation Table?

- Start with word alignment, build phrases from that

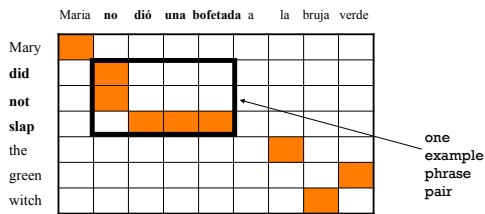
	Maria	no	dió	una	bofetada	a	la	bruja	verde
Mary									
did									
not									
slap									
the									
green									
witch									

This word-to-word alignment is a by-product of training a translation model like IBM-Model-3.

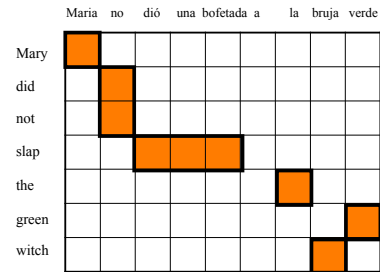
This is the best (or "Viterbi") alignment.

+ How to Learn the Phrase Translation Table?

- Collect all phrase pairs that are consistent with the word alignment

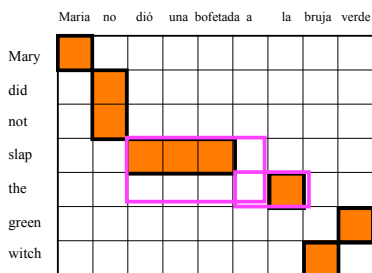


+ Word Alignment Induced Phrases



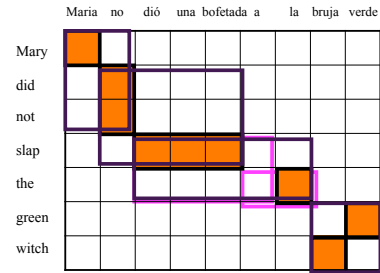
(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)

+ Word Alignment Induced Phrases



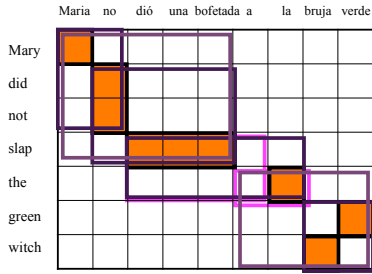
(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)
 (a la, the) (dió una bofetada a, slap the)

+ Word Alignment Induced Phrases



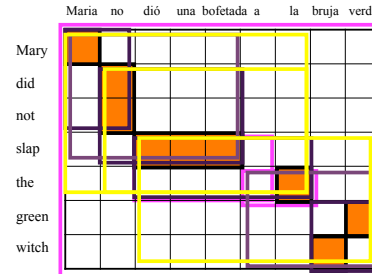
(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)
 (a la, the) (dió una bofetada a, slap the)
 (Maria no, Mary did not) (no dió una bofetada, did not slap), (dió una bofetada a la, slap the)
 (bruja verde, green witch)

+ Word Alignment Induced Phrases



(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)
 (a la, the) (dió una bofetada a, slap the)
 (Maria no, Mary did not) (no dió una bofetada, did not slap), (dió una bofetada a la, slap the)
 (bruja verde, green witch) (Maria no dió una bofetada, Mary did not slap)
 (a la bruja verde, the green witch) ...

+ Word Alignment Induced Phrases

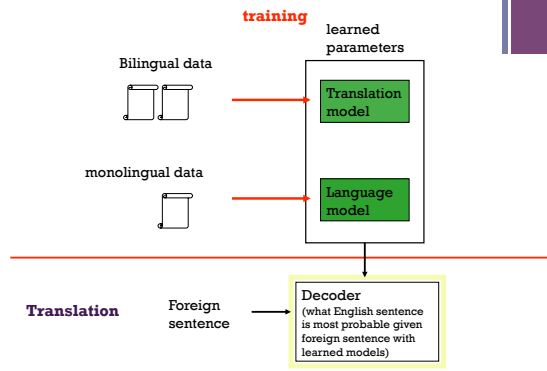


(Maria, Mary) (no, did not) (slap, dió una bofetada) (la, the) (bruja, witch) (verde, green)
 (a la, the) (dió una bofetada a, slap the)
 (Maria no, Mary did not) (no dió una bofetada, did not slap), (dió una bofetada a la, slap the)
 (bruja verde, green witch) (Maria no dió una bofetada, Mary did not slap)
 (a la bruja verde, the green witch) ...
 (Maria no dió una bofetada a la bruja verde, Mary did not slap the green witch)

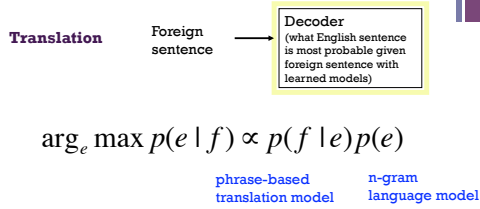
+ Learning phrase probabilities

- Extract all phrase pairs from our bilingual corpus
- Estimate $p(f_{\text{phrase}} | e_{\text{phrase}})$ using MLE
- Prune the rules
 - remove rules based on low-frequency events
 - remove rules with very low probabilities

+ Statistical MT Overview



+ Translation (aka decoding)



Ideas?

+ Decoding

- Of all conceivable English word strings, find the one maximizing $P(e) \times P(f | e)$
- Decoding is an NP-complete problem (for many translation models)
 - (Knight, 1999)
- AI to the rescue: pose as a search problem

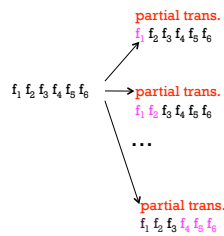
+ Search

$$\arg_e \max p(f | e)p(e)$$

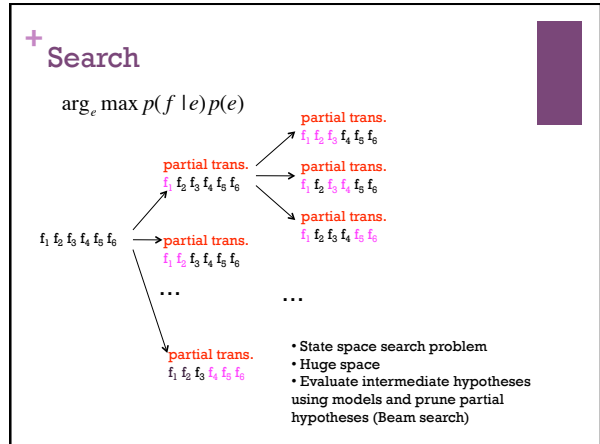
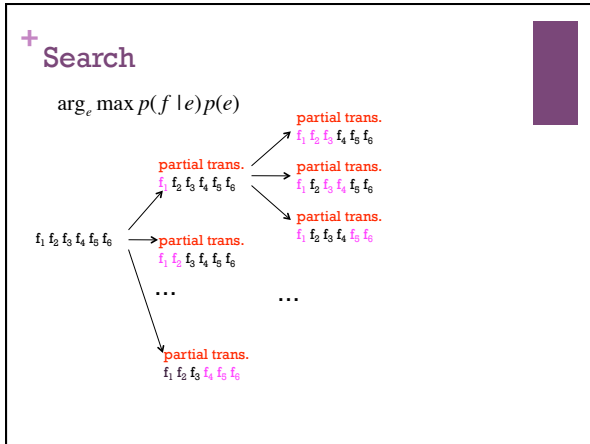
$f_1 f_2 f_3 f_4 f_5 f_6$

+ Search

$$\arg_e \max p(f | e)p(e)$$



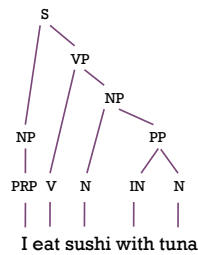
Key: partial translations represent *inorder* translations on the English side



- + Phrased-based approach**
- Works very well!
 - Still more or less state of the art
 - Google translate currently still uses this approach
 - much larger n-gram language model (5-grams >)
 - but... we can do better

- + Some things are hard to model with phrases...**
- Word order within constituents:
 - English NPs: art adj n *the big boy*
 - Hebrew NPs: art n art adj *ha yeled ha gadol*
 - Spanish NPs: art n adj *el niño grande*
 - Constituent structure:
 - English is SVO: Subj Verb Obj *I saw the man*
 - Modern Arabic is VSO: Verb Subj Obj
 - Different verb syntax:
 - Verb complexes in English vs. in German
 - I can eat the apple Ich kann den apfel essen*

+ Syntax!



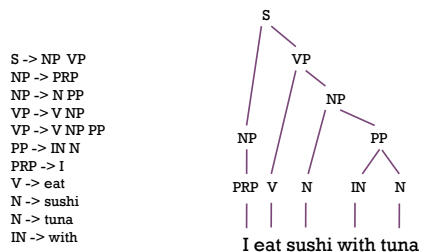
All of the examples before were motivated by the syntax of the languages

+ 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
- What types of grammars have you seen in CS?

+ Context free grammars

- Specify a language (i.e. set of valid sentences)
- If we look at the derivation it forms a tree



+ Parsing

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

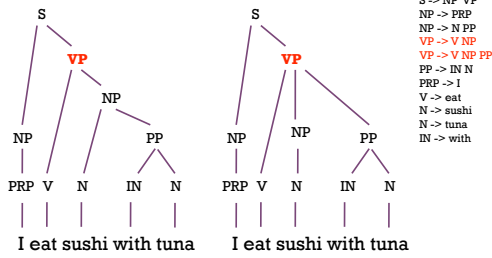
I eat sushi with tuna

What parse trees are possible for this sentence?

How can you algorithmically determine this?

What if the grammar is much larger?

+ Parsing ambiguity

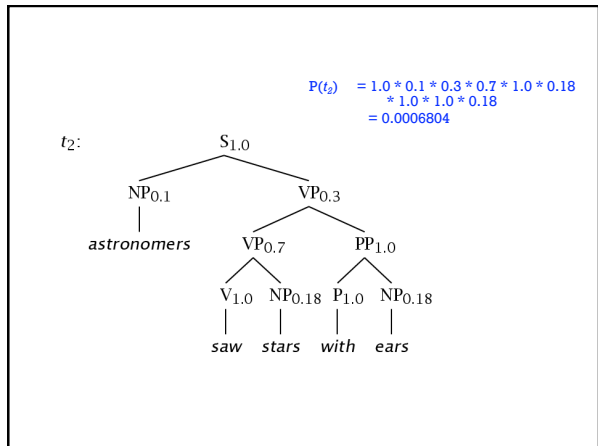
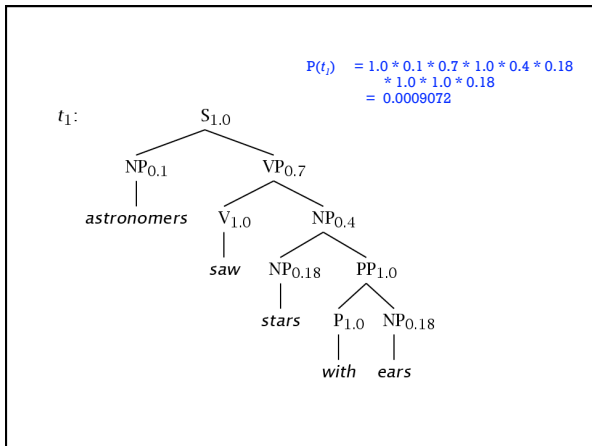


How can we decide between these?

+ A Simple PCFG

Probabilities!

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

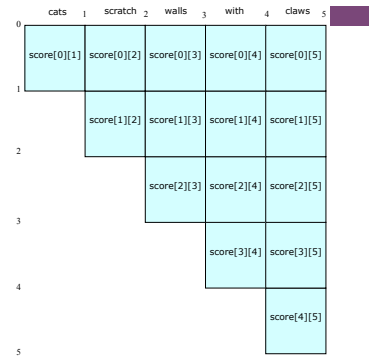


+ Parsing

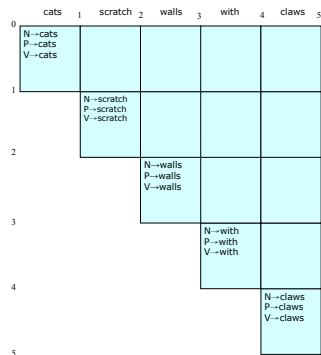
- Top-down parsing
 - ends up doing a lot of repeated work
 - doesn't take into account the words in the sentence until the end!
- Bottom-up parsing
 - constrain based on the words
 - avoids repeated work (dynamic programming)
 - CKY parser

score[i][j] has all of the partial parses spanning words i to j-1

build tree from the bottom up

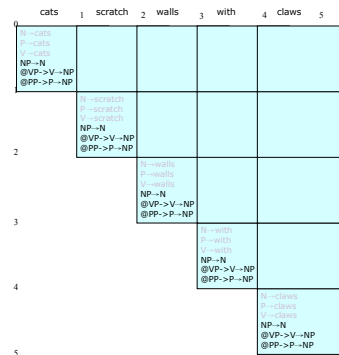


+

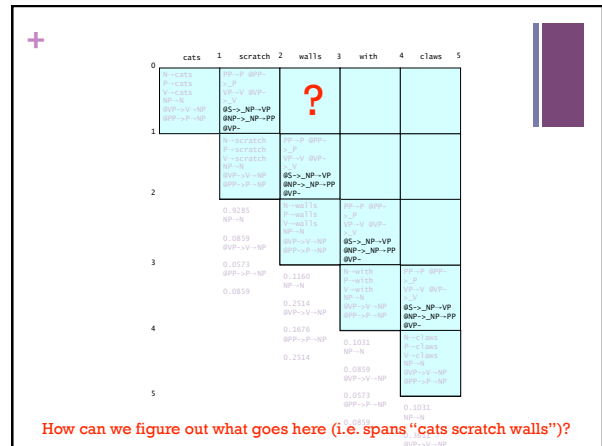
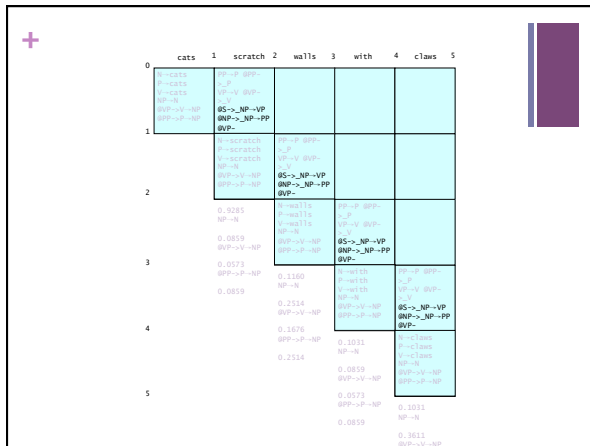
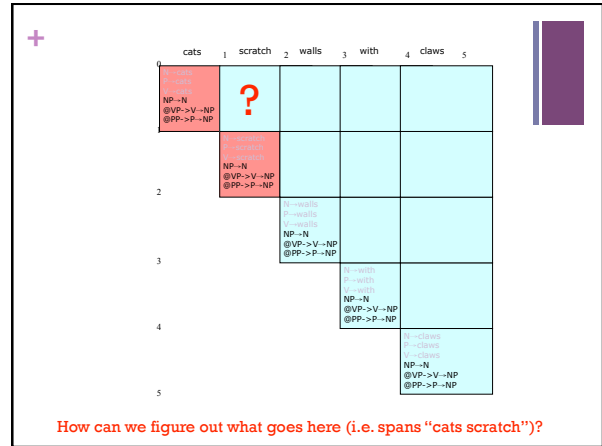
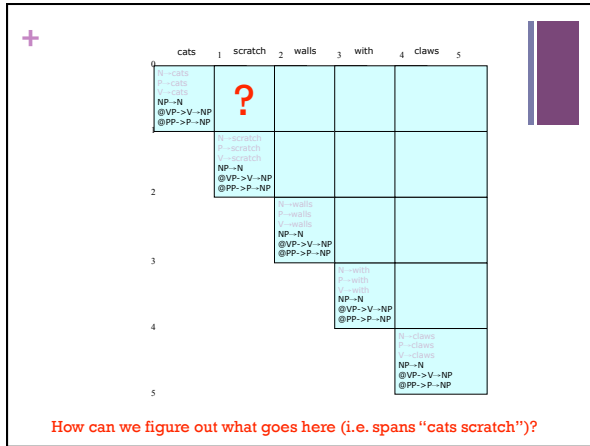


Start at the bottom level

+



Be careful about unary rules (i.e. rules with 1 element on the RHS)



+

	cats	1	scratch ₂	walls	3	with	4	claws	5
1	S → NP VP 0.4275 NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475	NP → N 0.9992 VP → V NP 0.4475
2									
3									
4									
5									

- +
- ### Parsing
- Works well in practice
 - many publicly available
 - Stanford, CMU, Berkeley all have good ones available
 - Many with online interfaces:
 - <http://nlp.stanford.edu:8080/parser/>
 - How do we obtain the grammars?
 - Tree banks: contain examples sentences with human annotated parse trees
 - learn rules and probabilities from these tree banks
 - Improvements
 - Lexicalize the rules
 - S → NP VP becomes S(eat) → NP(I) VP(eat)
 - Machine learning approaches

- +
- ### Other types of parsing
- Part of speech tagging
 - I/PRP eat/VBP sushi/NN with/IN tuna/NN ./.
 - Dependency parsing
 - nsubj(eat-2, I-1) – who's doing the eating?
 - dobj(eat-2, sushi-3) – what are they eating?
 - prep_with(sushi-3, tuna-5) – what is the sushi with?

+

How do people do it?

The horse raced past the barn fell.

The horse *that was* raced past the barn fell.

+ How do people do it?

The old man the boat.

The old *people* man the boat.

+ How do people do it?

The man whistling tunes pianos

The man *who is* whistling tunes pianos

+ How do people do it?

The government plans to raise taxes were defeated.

The plans *of the government* to raise taxes were defeated.

+ Garden path effect

- People tend to parse this a sentence as they read it
 - not bottoms-up
- For garden path sentences, the initial parse is incorrect
 - have to go back and reparse the sentence

