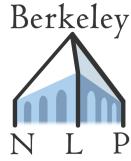


# Statistical NLP

## Spring 2011

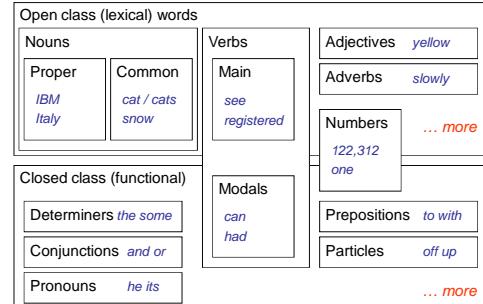


### Lecture 6: POS / Phrase MT

Dan Klein – UC Berkeley

## Parts-of-Speech (English)

- One basic kind of linguistic structure: syntactic word classes



CC	conjunction, coordinating
CD	numeral, cardinal
DT	determiner
EX	existential there
FW	foreign word
IN	preposition or conjunction, subordinating
JJ	adjective, ordinal
JJR	adjective, comparative
JJS	adjective, superlative
MD	modality
NN	noun, common, singular or mass
NNP	noun, proper, singular
NNPS	noun, proper, plural
NNS	noun, common, plural
POS	genitive marker
PRP	pronoun, personal
PRP\$	pronoun, possessive
RB	adverb
RBR	adverb, comparative
RBS	adverb, superlative
RP	particle
TO	"to" as preposition or infinitive marker
UH	interjection
VB	verb, base form
VBD	verb, past tense
VBG	verb, present participle or gerund
VBN	verb, past participle
VBP	verb, present tense, not 3rd person singular
VBZ	verb, present tense, 3rd person singular
WDT	WH-detector
WP	WH-pronoun
WP\$	WH-pronoun, possessive
WRB	Wh-adverb

and both but either or
mid-1890 nine-nine-thy-0.5 one
a all an every no that the
there
gemeinschaft hund ich jeux
among whether cut on by if
third ill-mannered regrettable
braver cheaper taller
best biggest easiest greatest
can might will would
cabbage thermosist investment
Motown Cougar Yvette Liverpool
Americans Materials States
undergraduates bric-a-brac averages
's
hers himself it we them
her his mine my our ours their thy your
occasionally maddeningly adventurously
further gloomier heavier less-perfectly
best biggest easiest worst
aboard away back by on open through
to
huh hokey uh whammo shucks heck
ask bring fire see take
pledged swiped registered saw
stirring focusing approaching erasing
dilapidated imitated reunified unsettled
twist appear comprise molt postpone
bases reconsists marks uses
that what whatever which whichever
that what whatever which who whom
whose
however wherever where why

## Why POS Tagging?

- Useful in and of itself (more than you'd think)
  - Text-to-speech: record, lead
  - Lemmatization: *saw[v]* → *see*, *saw[n]* → *saw*
  - Quick-and-dirty NP-chunk detection: *grep (JJ | NN)\* {NN | NNS}*
- Useful as a pre-processing step for parsing
  - Less tag ambiguity means fewer parses
  - However, some tag choices are better decided by parsers

```

IN
DT NNP NN VBD VBN RP NN NNS
The Georgia branch had taken on loan commitments ...

VBN
DT NN IN NN VBD NNS VBD
The average of interbank offered rates plummeted ...
  
```

## Part-of-Speech Ambiguity

- Words can have multiple parts of speech

VBD	VB
VBN	VBZ
VBD	VBP
NNP	NN
NNS	NNS
CD	CD
NN	NN

Fed raises interest rates 0.5 percent

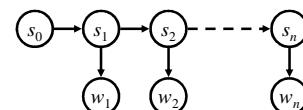
Mrs./NNP Shafer/NNP never/RB got/VBD around/RP to/TO joining/VBG  
 All/DT we/WP gotta/VBN do/VB is/VBZ go/VB around/IN the/DT corner/NN  
 Chateau/NNP Petrus/NNP costs/VBZ around/RB 250/CD

- Two basic sources of constraint:
  - Grammatical environment
  - Identity of the current word

- Many more possible features:
  - Suffixes, capitalization, name databases (gazetteers), etc...

## Classic Solution: HMMs

- We want a model of sequences  $s$  and observations  $w$



$$P(s, w) = \prod_i P(s_i | s_{i-1}) P(w_i | s_i)$$

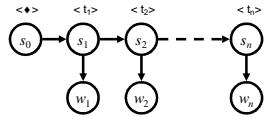
- Assumptions:

- States are tag n-grams
- Usually a dedicated start and end state / word
- Tag/state sequence is generated by a markov model
- Words are chosen independently, conditioned only on the tag/state
- These are totally broken assumptions: why?

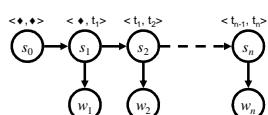
## States

- States encode what is relevant about the past
  - Transitions  $P(s'|s)$  encode well-formed tag sequences leading to next state  $s'$

- In a bigram tagger, states = tags



- In a trigram tagger, states = tag pairs



## Estimating Transitions

- Use standard smoothing methods to estimate transitions:  
$$P(t_i | t_{i-1}, t_{i-2}) = \lambda_2 \hat{P}(t_i | t_{i-1}, t_{i-2}) + \lambda_1 \hat{P}(t_i | t_{i-1}) + (1 - \lambda_1 - \lambda_2) \hat{P}(t_i)$$
  - Can get a lot fancier (e.g. KN smoothing) or use higher orders, but in this case it doesn't buy much
  - One option: encode more into the state, e.g. whether the previous word was capitalized (Brants 00)
  - **BIG IDEA:** The basic approach of state-splitting turns out to be very important in a range of tasks

# Estimating Emissions

$$P(\mathbf{s}, \mathbf{w}) = \prod_i P(s_i|s_{i-1})P(w_i|s_i)$$

- Emissions are trickier:
    - Words we've never seen before
    - Words which occur with tags we've never seen them with
    - One option: break out the Good-Turing smoothing
    - Issue: unknown words aren't black boxes:

343,127,23      11-year      Minteria      reintroducibly
  - Basic solution: unknown words classes (affixes or shapes)

D<sup>+</sup>, D<sup>-</sup>, D<sup>+</sup>      D<sup>+</sup>-x<sup>+</sup>      Xx<sup>+</sup>      x<sup>+</sup>-"ly"
  - [Brants 00] used a suffix trie as its emission model

## Disambiguation (Inference)

- Problem: find the most likely (Viterbi) sequence under the model  

$$t^* = \arg \max_t P(t|w)$$

▪ Given model parameters, we can score any tag sequence
<♦,♦>      <♦,NNP>    <NNP, VBZ>    <VBZ, NN>    <NN, NNS>    <NNS, CD>    <CD, NN>    <STOP,>
NNP      VBZ      NN      NNS      CD      NN      .
Fed    raises    interest    rates    0.5    percent    .

- In principle, we're done – list all possible tag sequences, score each one, pick the best one (the Viterbi state sequence)

NNP VBZ NN NNS CD NN → logP = -23  
NNP NNS NN NNS CD NN → logP = -29  
NNP VBZ VB NNS CD NN → logP = -27

## Finding the Best Trajectory

- Too many trajectories (state sequences) to list
  - Option 1: Beam Search
 

```

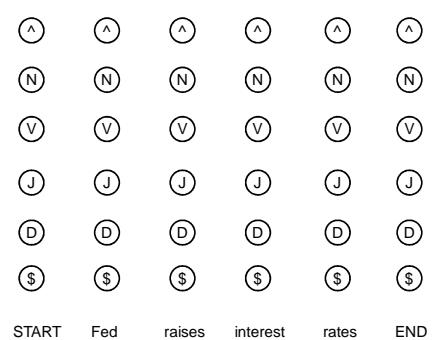
graph LR
    Root[\"Fed:NNP\"] --> VBN[\"Fed:VBN\"] 
    Root --> VBD[\"Fed:VBD\"]

    VBN --> NNS1[\"Fed:NNP raises:NNS\"] 
    VBN --> VBZ1[\"Fed:NNP raises:VBZ\"] 

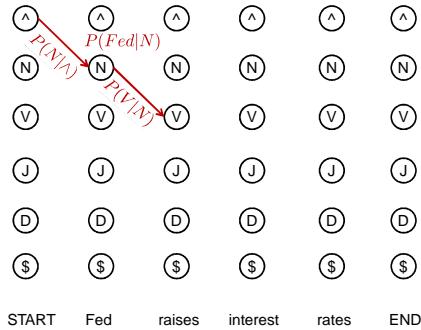
    VBD --> NNS2[\"Fed:NNP raises:NNS\"] 
    VBD --> VBZ2[\"Fed:NNP raises:VBZ\"] 

    VBZ1 --- Cross1(( ))
    NNS1 --- Cross2(( ))
    NNS2 --- Cross3(( ))
  
```
  - A beam is a set of partial hypotheses
  - Start with just the single empty trajectory
  - At each derivation step:
    - Consider all continuations of previous hypotheses
    - Discard most, keep top k, or those within a factor of the best
  - Beam search works ok in practice
    - ... but sometimes you want the optimal answer
    - ... and you need optimal answers to validate your beam search
    - ... and there's usually a better option than naïve beams

# The State Lattice / Trellis



## The State Lattice / Trellis



## The Viterbi Algorithm

- Dynamic program for computing

$$\delta_i(s) = \max_{s_0 \dots s_{i-1}, w} P(s_0 \dots s_i, w_1 \dots w_i)$$

- The score of a best path up to position i ending in state s

$$\delta_0(s) = \begin{cases} 1 & \text{if } s = <\bullet, \bullet> \\ 0 & \text{otherwise} \end{cases}$$

$$\delta_i(s) = \max_{s'} P(s | s') P(w | s') \delta_{i-1}(s')$$

- Also can store a backtrace (but no one does)

$$\psi_i(s) = \arg \max_{s'} P(s | s') P(w | s') \delta_{i-1}(s')$$

- Memoized solution

- Iterative solution

## So How Well Does It Work?

- Choose the most common tag
  - 90.3% with a bad unknown word model
  - 93.7% with a good one

- TnT (Brants, 2000):
  - A carefully smoothed trigram tagger
  - Suffix trees for emissions
  - 96.7% on WSJ text (SOA is ~97.5%)

- Noise in the data
  - Many errors in the training and test corpora

DT NN IN NN VBD NNS VBD  
The average of interbank offered rates plummeted ...

- Probably about 2% guaranteed error from noise (on this data)

JJ	JJ	NN
chief executive officer		
NN	JJ	NN
chief executive officer		
JJ	NN	NN
chief executive officer		
NN	NN	NN
chief executive officer		

## Overview: Accuracies

- Roadmap of (known / unknown) accuracies:

- Most freq tag: ~90% / ~50%

- Trigram HMM: ~95% / ~55%

- TnT (HMM++): 96.2% / 86.0%

Most errors on unknown words

- Maxent P(t|w): 93.7% / 82.6%

- MEMM tagger: 96.9% / 86.9%

- Cyclic tagger: 97.2% / 89.0%

- Upper bound: ~98%

## Common Errors

- Common errors [from Toutanova & Manning 00]

	JJ	NN	NNP	NNPS	RB	RP	IN	VB	VBD	VBN	VBP	Total
JJ	0	177	56	0	61	2	5	10	15	108	0	488
NN	244	0	103	0	12	1	1	29	5	6	19	525
NNP	107	106	0	132	5	0	7	5	1	2	0	427
NNPS	1	0	110	0	0	0	0	0	0	0	0	142
RB	72	21	7	0	0	16	138	1	0	0	0	295
RP	0	0	0	0	39	0	65	0	0	0	0	104
IN	11	0	1	0	169	103	0	1	0	0	0	323
VB	17	64	9	0	2	0	1	4	7	85	189	
VBD	10	5	3	0	0	0	3	0	143	2	166	
VBN	101	3	3	0	0	0	0	3	108	0	1	221
VBP	5	34	3	1	1	0	2	49	6	3	0	104
Total	626	536	348	144	317	122	279	102	140	269	108	3651

NN/JJ NN  
official knowledge

VBD RP//IN DT NN  
made up the story

RB VBD/VBN NNS  
recently sold shares

## Corpus-Based MT

Modeling correspondences between languages

Sentence-aligned parallel corpus:

Yo lo haré mañana  
I will do it tomorrow

Hasta pronto  
See you soon

Hasta pronto  
See you around

Machine translation system:

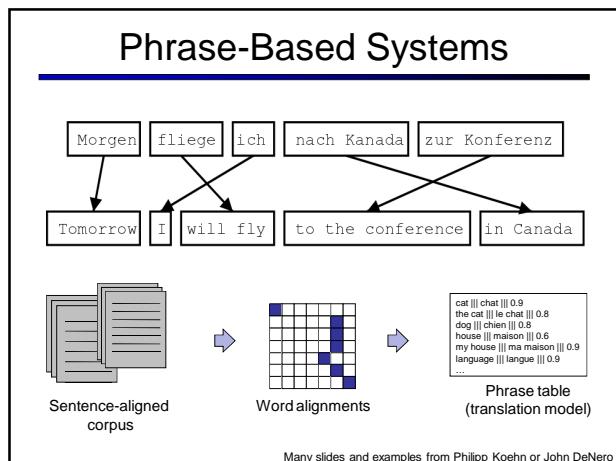
Yo lo haré pronto  
Novel Sentence

Model of translation

I will do it soon  
I will do it around  
See you tomorrow

# Phrase-Based Translation Overview

		<i>The decoder...</i>
<b>Input:</b>	lo haré   rápidamente   .	tries different segmentations,
<b>Translations:</b>	I'll do it   quickly   .	translates phrase by phrase,
	quickly   I'll do it   .	and considers reorderings.

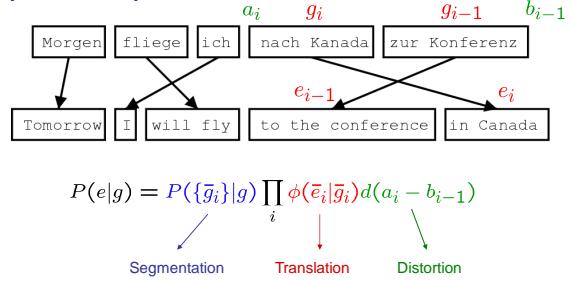


Many slides and examples from Philipp Koehn or John DeNero

# Phrase-Based Decoding

## The Pharaoh “Model”

[Koehn et al, 2003]



$$P(e|g) = P(\{\bar{g}_i\}|g) \prod \phi(\bar{e}_i|\bar{g}_i) d(a_i - b_{i-1})$$

## Segmentation      Translation      Distortion

## Distortion

# The Pharaoh “Model”

## Phrase Weights

How the MT community estimates  $P(\bar{f}|\bar{e})$

### *Parallel training sentences*

provide phrase pair counts.

Gracias , lo haré de muy buen gusto.  
Thank you . I shall do so gladly.

lo haré ↲ I shall do so  
44 times in the corpus

All phrase pairs are counted

*and counts are normalized*

Gracias [lo haré de muy buen grado]  
Thank you [I shall do so gladly]

$$P(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{f}, \bar{e})}{\text{count}(\bar{e})}$$

## Phrase-Based Decoding

Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
did not	give	a	slap	by	the	green	witch	green
no	slap	to	the					

slap the witch

## Monotonic Word Translation

Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
did not	give	a	slap	by	the	green	witch	green
no	slap	to	the					

no slap the green witch

- Cost is LM \* TM

- It's an HMM?

  - P(e|e<sub>1</sub>, e<sub>2</sub>)
  - P(f|e)

- State includes

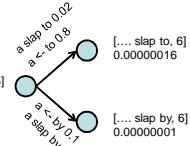
  - Exposed English
  - Position in foreign

- Dynamic program loop?

```
for (fPosition in 1...|f|)
    for (eContext in bestEContexts[fPosition])
        for (eOption in translations[fPosition])
            score = scores[fPosition-1][eContext] * LM(eContext) * TM(eOption, fWord[fPosition])
            bestEContexts.maybeAdd(eContext[2]+eOption, score)
```

score = scores[fPosition-1][eContext] \* LM(eContext) \* TM(eOption, fWord[fPosition])

bestEContexts[eContext[2]+eOption] =<sub>max</sub> score



0.00000016

0.000000016

0.000000001

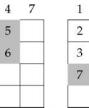
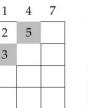
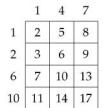
0.0000000001

## Beam Decoding

- For real MT models, this kind of dynamic program is a disaster (why?)
- Standard solution is beam search: for each position, keep track of only the best k hypotheses

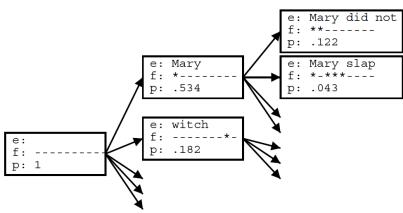
```
for (fPosition in 1...|f|)
    for (eContext in bestEContexts[fPosition])
        for (eOption in translations[fPosition])
            score = scores[fPosition-1][eContext] * LM(eContext) * TM(eOption, fWord[fPosition])
            bestEContexts.maybeAdd(eContext[2]+eOption, score)
```

- Still pretty slow... why?
- Useful trick: cube pruning (Chiang 2005)



Example from David Chiang

## Non-Monotonic Phrasal MT



## Pruning: Beams + Forward Costs

Maria no dio una bofetada a la bruja verde

e: Mary did not  
f: \*-----  
p: 0.154

better  
partial  
translation

e: the  
f: -----  
p: 0.354

covers  
easier part  
--> lower cost

- Problem: easy partial analyses are cheaper

  - Solution 1: use beams per foreign subset

  - Solution 2: estimate forward costs (A\*-like)

## The Pharaoh Decoder

Maria	no	dio	una	bofetada	a	la	bruja	verde
Mary	not	give	a	slap	to	the	witch	green
	did not						green	witch
	no		slap		by			
					to			
					the			
				slap			the	witch

Maria	no	dio	una	bofetada	a la	bruja	verde
Mary	did not	give	a	slap	to	the	witch
	no						green
			slap		by		witch
					to		
					the		
				slap			the
							witch

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

## Hypothesis Lattices

