

Statistical NLP

Spring 2009

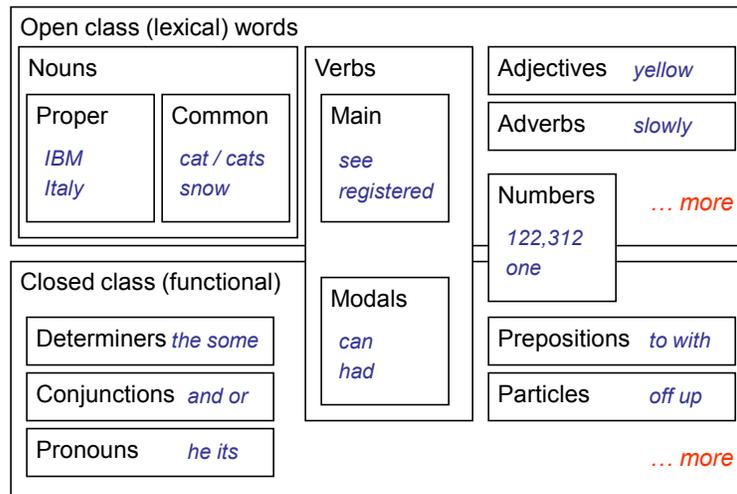


Lecture 6: Parts-of-Speech

Dan Klein – UC Berkeley

Parts-of-Speech (English)

- One basic kind of linguistic structure: syntactic word classes



CC	conjunction, coordinating	and both but either or
CD	numeral, cardinal	mid-1890 nine-thirty 0.5 one
DT	determiner	a all an every no that the
EX	existential there	there
FW	foreign word	gemeinschaft hund ich jeux
IN	preposition or conjunction, subordinating	among whether out on by if
JJ	adjective or numeral, ordinal	third ill-mannered regrettable
JJR	adjective, comparative	braver cheaper taller
JJS	adjective, superlative	bravest cheapest tallest
MD	modal auxiliary	can may might will would
NN	noun, common, singular or mass	cabbage thermostat investment subhumanity
NNP	noun, proper, singular	Motown Cougar Yvette Liverpool
NNPS	noun, proper, plural	Americans Materials States
NNS	noun, common, plural	undergraduates bric-a-brac averages
POS	genitive marker	's
PRP	pronoun, personal	hers himself it we them
PRP\$	pronoun, possessive	her his mine my our ours their thy your
RB	adverb	occasionally maddeningly adventurously
RBR	adverb, comparative	further gloomier heavier less-perfectly
RBS	adverb, superlative	best biggest nearest worst
RP	particle	aboard away back by on open through
TO	"to" as preposition or infinitive marker	to
UH	interjection	huh howdy uh whammo shucks heck
VB	verb, base form	ask bring fire see take
VBD	verb, past tense	pleaded swiped registered saw
VBG	verb, present participle or gerund	stirring focusing approaching erasing
VBN	verb, past participle	dilapidated imitated reunified unsettled
VBP	verb, present tense, not 3rd person singular	twist appear comprise mold postpone
VBZ	verb, present tense, 3rd person singular	bases reconstructs marks uses
WDT	WH-determiner	that what whatever which whichever
WP	WH-pronoun	that what whatever which who whom
WP\$	WH-pronoun, possessive	whose
WRB	Wh-adverb	however whenever where why

Part-of-Speech Ambiguity

- Example

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

Fed raises interest rates 0.5 percent

Mrs./NNP Shaefer/NNP never/RB got/VBD **around/RP** to/TO joining/VBG
 All/DT we/PRP 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:
 - ... but we won't be able to use them for a while

Part-of-Speech Tagging

Republicans warned Sunday that the Obama administration 's \$ 800 billion economic stimulus effort will lead to what one called a " financial disaster . "

The administration is also readying a second phase of the financial bailout program launched by the Bush administration last fall.

Why POS Tagging?

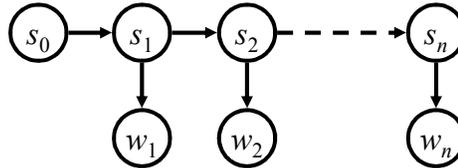
- Useful in and of itself
 - Text-to-speech: record, lead
 - Lemmatization: saw[v] → see, saw[n] → saw
 - Quick-and-dirty NP-chunk detection: `grep {JJ | NN}* {NN | NNS}`
- Useful as a pre-processing step for parsing?
 - Less tag ambiguity means fewer parses
 - However, some tag choices are better decided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

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

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

HMMs

- We want a model of sequences s and observations w



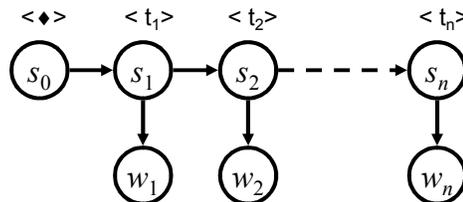
$$P(\mathbf{s}, \mathbf{w}) = \prod_i P(s_i | s_{i-1}) P(w_i | s_i)$$

- 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?

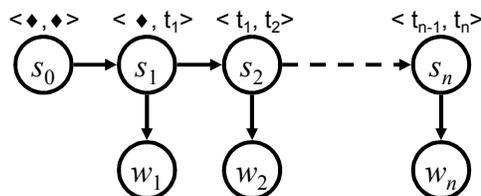
Transitions

- Transitions $P(s|s')$ encode well-formed tag sequences

- 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)

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-Turning smoothing
 - Issue: unknown words aren't black boxes:
343,127.23 11-year Minteria reintroducibly
 - Solution: unknown words classes (affixes or shapes)
D⁺, D⁺.D⁺ D⁺-x⁺ Xx⁺ x⁺"ly"
- [Brants 00] used a suffix trie as its emission model

Disambiguation

- Given these two multinomials, we can score any word / tag sequence pair

$\langle \diamond, \diamond \rangle$ $\langle \diamond, \text{NNP} \rangle$ $\langle \text{NNP}, \text{VBZ} \rangle$ $\langle \text{VBZ}, \text{NN} \rangle$ $\langle \text{NN}, \text{NNS} \rangle$ $\langle \text{NNS}, \text{CD} \rangle$ $\langle \text{CD}, \text{NN} \rangle$ $\langle \text{STOP} \rangle$
 NNP VBZ NN NNS CD NN .
 Fed raises interest rates 0.5 percent .

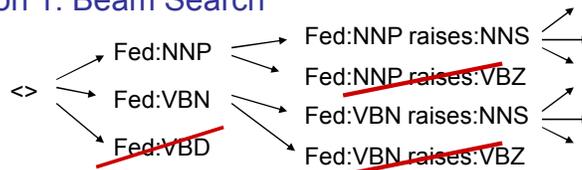
$P(\text{NNP}|\langle \diamond, \diamond \rangle)$ $P(\text{Fed}|\text{NNP})$ $P(\text{VBZ}|\langle \text{NNP}, \diamond \rangle)$ $P(\text{raises}|\text{VBZ})$ $P(\text{NN}|\text{VBZ}, \text{NNP})$

- 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 \Rightarrow $\log P = -23$
 NNP NNS NN NNS CD NN \Rightarrow $\log P = -29$
 NNP VBZ VB NNS CD NN \Rightarrow $\log P = -27$

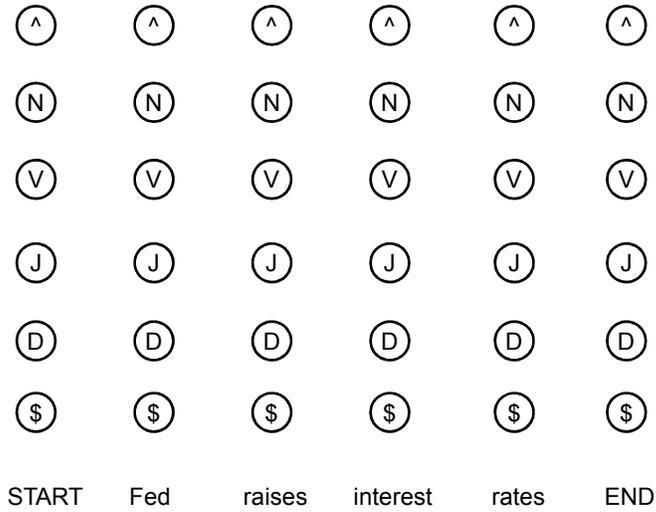
Finding the Best Trajectory

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

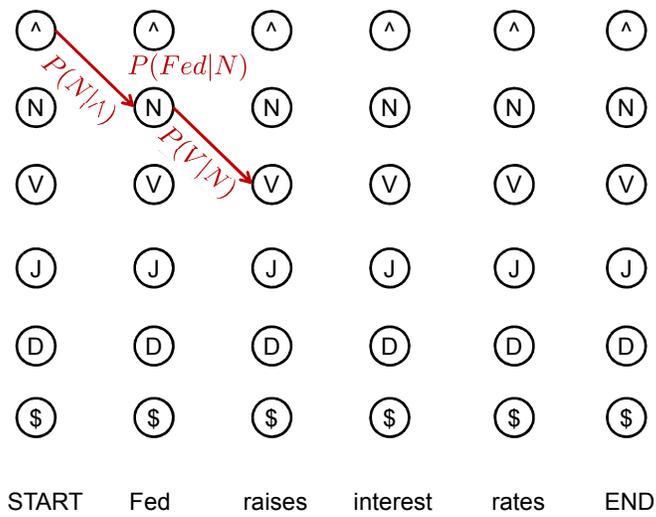


- A beam is a set of partial hypotheses
- Start with just the single empty trajectory
- At each derivation step:
 - Consider all continuations of previous hypotheses
 - Discard most, keep top k, or those within a factor of the best, (or some combination)
- Beam search works relatively well in practice
 - ... but sometimes you want the optimal answer
 - ... and you need optimal answers to validate your beam search

The State Lattice / Trellis



The State Lattice / Trellis



The Viterbi Algorithm

- Dynamic program for computing

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

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

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

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

- Also store a backtrace

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

- 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 ~97.2%)

JJ JJ NN
chief executive officer

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

NN JJ NN
chief executive officer

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

JJ NN NN
chief executive officer

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

NN NN NN
chief executive officer

Overview: Accuracies

- Roadmap of (known / unknown) accuracies:

- Most freq tag: ~90% / ~50%
- Trigram HMM: ~95% / ~55%
- Maxent P(t|w): 93.7% / 82.6%
- TnT (HMM++): 96.2% / 86.0%
- MEMM tagger: 96.9% / 86.9%
- Cyclic tagger: 97.2% / 89.0%
- Upper bound: ~98%

Most errors on unknown words

Common Errors

- Common errors [from Toutanova & Manning 00]

	JJ	NN	NNP	NNPS	RB	RP	IN	VB	VBD	VCN	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	0	4	7	85	189
VBD	10	5	3	0	0	0	0	3	0	143	0	166
VCN	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

Better Features

- Can do surprisingly well just looking at a word by itself:
 - Word the: the → DT
 - Lowercased word Importantly: importantly → RB
 - Prefixes unfathomable: un- → JJ
 - Suffixes Surprisingly: -ly → RB
 - Capitalization Meridian: CAP → NNP
 - Word shapes 35-year: d-x → JJ
- Then build a maxent (or whatever) model to predict tag
- Maxent $P(t|w)$: 93.7% / 82.6%

Sequence-Free Tagging?

- What about looking at a word and its environment, but no sequence information?
 - Add in previous / next word the __
 - Previous / next word shapes X __ X
 - Occurrence pattern features [X: x X occurs]
 - Crude entity detection __ (Inc.|Co.)
 - Phrasal verb in sentence? put __
 - Conjunctions of these things
- All features except sequence: 96.6% / 86.8%
- Uses lots of features: > 200K
- Why isn't this the standard approach?

Why Linear Context is Useful

- Lots of local information!

RB
PRP VBD IN RB IN PRP VBD .
They left as soon as he arrived .

- We could fix this with a feature that looked at the next word

JJ
NNP NNS VBD VBN .
Intrinsic flaws remained undetected .

- We could fix this by linking capitalized words to their lowercase versions
- Solution: maximum entropy sequence models (MEMMs, CRFs)
- Reality check:
 - Taggers are already pretty good on WSJ journal text...
 - What the world needs is taggers that work on other text!

Maxent Taggers

- One step up: also condition on previous tags

$$P(t|\mathbf{w}) = \prod_i P_{ME}(t_i|\mathbf{w}, t_{i-1}, t_{i-2})$$

- Train up $P(t_i|\mathbf{w}, t_{i-1}, t_{i-2})$ as a normal maxent problem, then use to score sequences
- This is referred to as a *maxent tagger* [Ratnaparkhi 96]
- Beam search effective! (Why?)
- What's the advantage of beam size 1?

Decoding

- Decoding maxent taggers:
 - Just like decoding HMMs
 - Viterbi, beam search, posterior decoding
- Viterbi algorithm (HMMs):

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

- Viterbi algorithm (Maxent):

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

TBL Tagger

- [Brill 95] presents a *transformation-based tagger*
 - Label the training set with most frequent tags

DT MD VBD VBD .
The can was rusted .
 - Add transformation rules which reduce training mistakes
 - MD → NN : DT __
 - VBD → VBN : VBD __ .
 - Stop when no transformations do sufficient good
 - Does this remind anyone of anything?
- Probably the most widely used tagger (esp. outside NLP)
- ... but not the most accurate: 96.6% / 82.0 %

TBL Tagger II

- What gets learned? [from Brill 95]

#	Change Tag		Condition
	From	To	
1	NN	VB	Previous tag is <i>TO</i>
2	VBP	VB	One of the previous three tags is <i>MD</i>
3	NN	VB	One of the previous two tags is <i>MD</i>
4	VB	NN	One of the previous two tags is <i>DT</i>
5	VBD	VBN	One of the previous three tags is <i>VBZ</i>
6	VBN	VBD	Previous tag is <i>PRP</i>
7	VBN	VBD	Previous tag is <i>NNP</i>
8	VBD	VBN	Previous tag is <i>VBD</i>
9	VBP	VB	Previous tag is <i>TO</i>
10	POS	VBZ	Previous tag is <i>PRP</i>
11	VB	VBP	Previous tag is <i>NNS</i>
12	VBD	VBN	One of previous three tags is <i>VBP</i>
13	IN	WDT	One of next two tags is <i>VB</i>
14	VBD	VBN	One of previous two tags is <i>VB</i>
15	VB	VBP	Previous tag is <i>PRP</i>
16	IN	WDT	Next tag is <i>VBZ</i>
17	IN	DT	Next tag is <i>NN</i>
18	JJ	NNP	Next tag is <i>NNP</i>
19	IN	WDT	Next tag is <i>VBD</i>
20	JJR	RBR	Next tag is <i>JJ</i>

#	Change Tag		Condition
	From	To	
1	NN	NNS	Has suffix -s
2	NN	CD	Has character .
3	NN	JJ	Has character -
4	NN	VBN	Has suffix -ed
5	NN	VBG	Has suffix -ing
6	??	RB	Has suffix -ly
7	??	JJ	Adding suffix -ly results in a word.
8	NN	CD	The word \$ can appear to the left.
9	NN	JJ	Has suffix -al
10	NN	VB	The word would can appear to the left.
11	NN	CD	Has character 0
12	NN	JJ	The word be can appear to the left.
13	NNS	JJ	Has suffix -us
14	NNS	VBZ	The word it can appear to the left.
15	NN	JJ	Has suffix -ble
16	NN	JJ	Has suffix -ic
17	NN	CD	Has character 1
18	NNS	NN	Has suffix -ss
19	??	JJ	Deleting the prefix un- results in a word
20	NN	JJ	Has suffix -ive

EngCG Tagger

- English constraint grammar tagger
 - [Tapanainen and Voutilainen 94]
 - Something else you should know about
 - Hand-written and knowledge driven
 - “Don’t guess if you know” (general point about modeling more structure!)
 - Tag set doesn’t make all of the hard distinctions as the standard tag set (e.g. JJ/NN)
 - They get stellar accuracies: 98.5% on *their* tag set
 - Linguistic representation matters...
 - ... but it’s easier to win when you make up the rules

CRF Taggers

- Newer, higher-powered discriminative sequence models
 - CRFs (also perceptrons, M3Ns)
 - Do not decompose training into independent local regions
 - Can be deathly slow to train – require repeated inference on training set
- Differences tend not to be too important for POS tagging
- Differences more substantial on other sequence tasks
- However: one issue worth knowing about in local models
 - “Label bias” and other explaining away effects
 - MEMM taggers’ local scores can be near one without having both good “transitions” and “emissions”
 - This means that often evidence doesn’t flow properly
 - Why isn’t this a big deal for POS tagging?
 - Also: in decoding, condition on predicted, not gold, histories

CRFs

- Make a maxent model over entire taggings

- MEMM

$$P(\mathbf{t}|\mathbf{w}) = \prod_i \frac{1}{Z(i)} \exp(\lambda^\top f(t_i, t_{i-1}, \mathbf{w}, i))$$

- CRF

$$\begin{aligned} P(\mathbf{t}|\mathbf{w}) &= \frac{1}{Z(\mathbf{w})} \exp(\lambda^\top f(\mathbf{t}, \mathbf{w})) \\ &= \frac{1}{Z(\mathbf{w})} \exp\left(\lambda^\top \sum_i f(t_i, t_{i-1}, \mathbf{w}, i)\right) \\ &= \frac{1}{Z(\mathbf{w})} \prod_i \phi_i(t_i, t_{i-1}) \end{aligned}$$

CRFs

- Like any maxent model, derivative is:

$$\frac{\partial L(\lambda)}{\partial \lambda} = \sum_k \left(\mathbf{f}_k(\mathbf{t}^k) - \sum_{\mathbf{t}} P(\mathbf{t}|\mathbf{w}_k) \mathbf{f}_k(\mathbf{t}) \right)$$

- So all we need is to be able to compute the expectation each feature, for example the number of times the label pair *DT-NN* occurs, or the number of times *NN-interest* occurs in a sentence
- How many times does, say, *DT-NN* occur at position 10? The ratio of the scores of trajectories with that configuration to the score of all
- This requires exactly the same forward-backward score ratios as for EM, but using the local potentials ϕ instead of the local probabilities

Domain Effects

- **Accuracies degrade outside of domain**
 - Up to triple error rate
 - Usually make the most errors on the things you care about in the domain (e.g. protein names)
- **Open questions**
 - How to effectively exploit unlabeled data from a new domain (what could we gain?)
 - How to best incorporate domain lexica in a principled way (e.g. UMLS specialist lexicon, ontologies)

Unsupervised Tagging?

- AKA part-of-speech induction
- Task:
 - Raw sentences in
 - Tagged sentences out
- Obvious thing to do:
 - Start with a (mostly) uniform HMM
 - Run EM
 - Inspect results

EM for HMMs: Process

- Alternate between recomputing distributions over hidden variables (the tags) and reestimating parameters
- Crucial step: we want to tally up how many (fractional) counts of each kind of transition and emission we have under current params:

$$\text{count}(s \rightarrow s') = \sum_i P(t_{i-1} = s, t_i = s' | \mathbf{w})$$

$$\text{count}(w, s) = \sum_{i:w_i=w} P(t_i = s | \mathbf{w})$$

- But we need a dynamic program to help, because there are too many sequences to sum over

EM for HMMs: Quantities

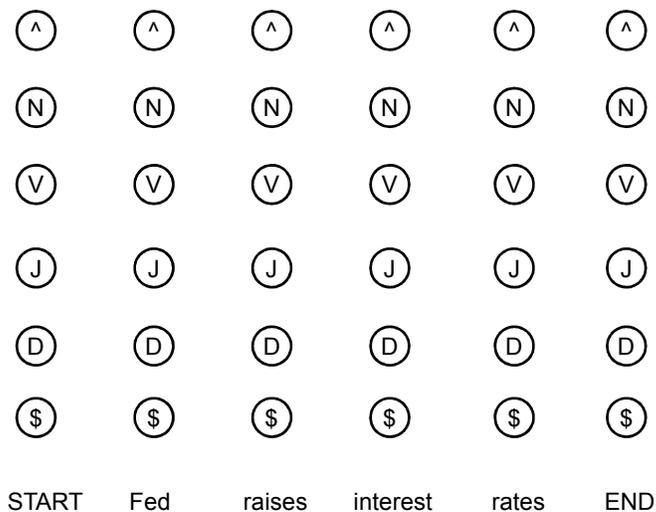
- Cache total path values:

$$\begin{aligned}\alpha_i(s) &= P(w_0 \dots w_i, s_i) \\ &= \sum_{s_{i-1}} P(s_i | s_{i-1}) P(w_i | s_i) \alpha_{i-1}(s_{i-1})\end{aligned}$$

$$\begin{aligned}\beta_i(s) &= P(w_{i+1} \dots w_n | s_i) \\ &= \sum_{s_{i+1}} P(s_{i+1} | s_i) P(w_{i+1} | s_{i+1}) \beta_{i+1}(s_{i+1})\end{aligned}$$

- Can calculate in $O(s^2n)$ time (why?)

The State Lattice / Trellis



EM for HMMs: Process

- From these quantities, can compute expected transitions:

$$\text{count}(s \rightarrow s') = \frac{\sum_i \alpha_i(s) P(s'|s) P(w_i|s) \beta_{i+1}(s')}{P(\mathbf{w})}$$

- And emissions:

$$\text{count}(w, s) = \frac{\sum_{i:w_i=w} \alpha_i(s) \beta_{i+1}(s)}{P(\mathbf{w})}$$

Merialdo: Setup

- Some (discouraging) experiments [Merialdo 94]
- Setup:
 - You know the set of allowable tags for each word
 - Fix k training examples to their true labels
 - Learn $P(w|t)$ on these examples
 - Learn $P(t|t_1, t_2)$ on these examples
 - On n examples, re-estimate with EM
- Note: we know allowed tags but not frequencies

Merialdo: Results

Number of tagged sentences used for the initial model							
	0	100	2000	5000	10000	20000	all
Iter	Correct tags (% words) after ML on 1M words						
0	77.0	90.0	95.4	96.2	96.6	96.9	97.0
1	80.5	92.6	95.8	96.3	96.6	96.7	96.8
2	81.8	93.0	95.7	96.1	96.3	96.4	96.4
3	83.0	93.1	95.4	95.8	96.1	96.2	96.2
4	84.0	93.0	95.2	95.5	95.8	96.0	96.0
5	84.8	92.9	95.1	95.4	95.6	95.8	95.8
6	85.3	92.8	94.9	95.2	95.5	95.6	95.7
7	85.8	92.8	94.7	95.1	95.3	95.5	95.5
8	86.1	92.7	94.6	95.0	95.2	95.4	95.4
9	86.3	92.6	94.5	94.9	95.1	95.3	95.3
10	86.6	92.6	94.4	94.8	95.0	95.2	95.2

Distributional Clustering

◆ *the president said that the downturn was over* ◆

<i>president</i>	<i>the __ of</i>
<i>president</i>	<i>the __ said</i>
<i>governor</i>	<i>the __ of</i>
<i>governor</i>	<i>the __ appointed</i>
<i>said</i>	<i>sources __ ◆</i>
<i>said</i>	<i>president __ that</i>
<i>reported</i>	<i>sources __ ◆</i>

*president
governor*

*the
a*

*said
reported*

[Finch and Chater 92, Shuetze 93, many others]

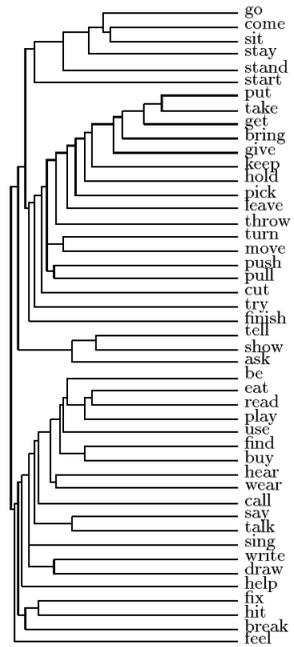
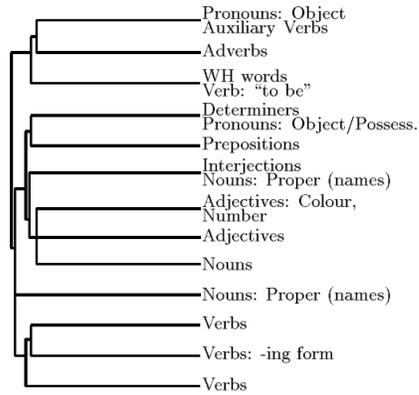
Distributional Clustering

- Three main variants on the same idea:
 - Pairwise similarities and heuristic clustering
 - E.g. [Finch and Chater 92]
 - Produces dendrograms
 - Vector space methods
 - E.g. [Shuetze 93]
 - Models of ambiguity
 - Probabilistic methods
 - Various formulations, e.g. [Lee and Pereira 99]

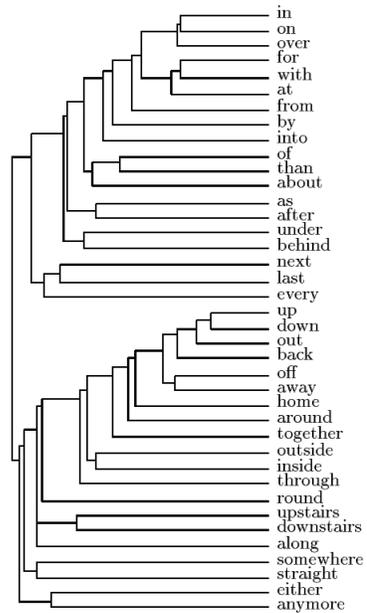
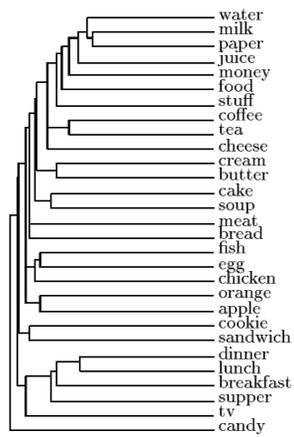
Nearest Neighbors

word	nearest neighbors
accompanied	submitted banned financed developed authorized headed canceled awarded barred
almost	virtually merely formally fully quite officially just nearly only less
causing	reflecting forcing providing creating producing becoming carrying particularly
classes	elections courses payments losses computers performances violations levels pictures
directors	professionals investigations materials competitors agreements papers transactions
goal	mood roof eye image tool song pool scene gap voice
japanese	chinese iraqi american western arab foreign european federal soviet indian
represent	reveal attend deliver reflect choose contain impose manage establish retain
think	believe wish know realize wonder assume feel say mean bet
york	angeles francisco sox rouge kong diego zone vegas inning layer
on	through in at over into with from for by across
must	might would could cannot will should can may does helps
they	we you i he she nobody who it everybody there

Dendrograms

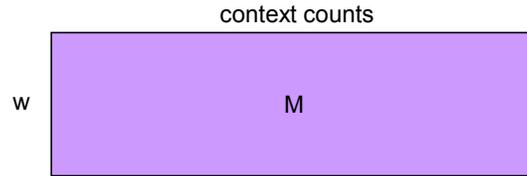


Dendrograms

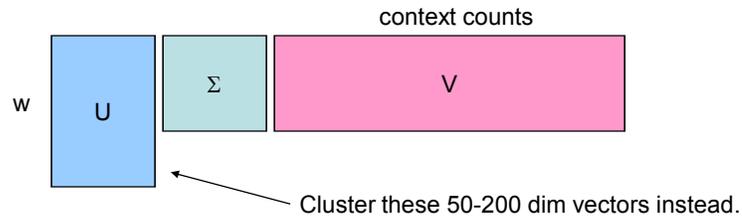


Vector Space Version

- [Shuetze 93] clusters words as points in R^n

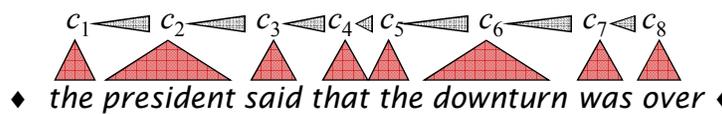
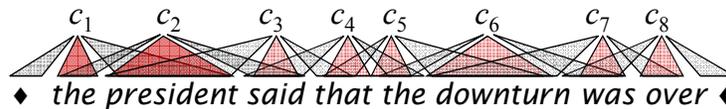


- Vectors too sparse, use SVD to reduce



A Probabilistic Version?

$$P(S, C) = \prod_i P(c_i)P(w_i | c_i)P(w_{i-1}, w_{i+1} | c_i)$$

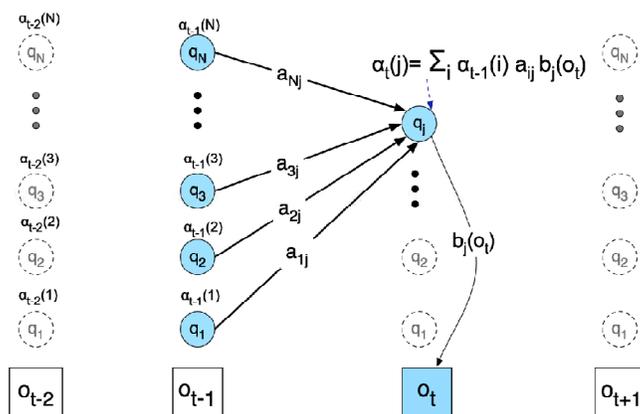


What Else?

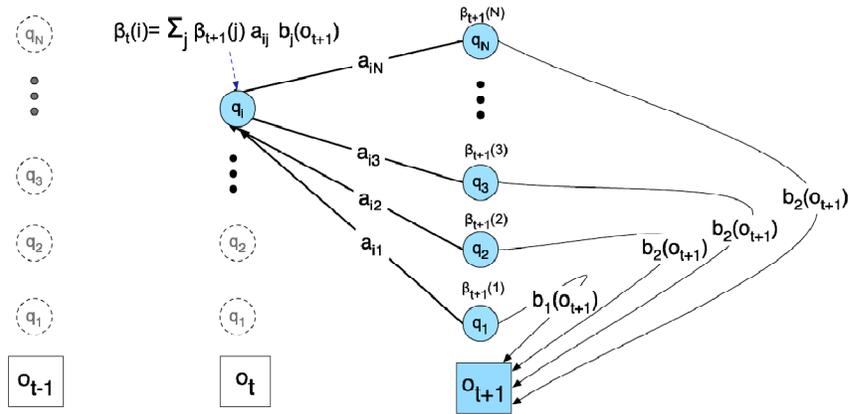
- Various newer ideas:
 - Context distributional clustering [Clark 00]
 - Morphology-driven models [Clark 03]
 - Contrastive estimation [Smith and Eisner 05]

- Also:
 - What about ambiguous words?
 - Using wider context signatures has been used for learning synonyms (what's wrong with this approach?)
 - Can extend these ideas for grammar induction (later)

Forward Recurrence



Backward Recurrence



Fractional Transitions

