

CS 294-5: Statistical Natural Language Processing



Unsupervised Tagging,
Word Clustering
Lecture 9: 9/28/05

Assignment 1 Honors

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: Quantities

- Remember from last time:

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

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

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

EM for HMMs: Process

- From these quantities, we can re-estimate 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:

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

- If you don't get these formulas immediately, just think about hard EM instead, where we re-estimate from the Viterbi sequences

Meritaldo: Setup

- Some (discouraging) experiments [Meritaldo 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

Meritaldo: 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* ◆

president	the __ of
president	the __ said
governor	the __ of
governor	the __ appointed
said	sources __ ◆
said	president __ that
reported	sources __ ◆



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

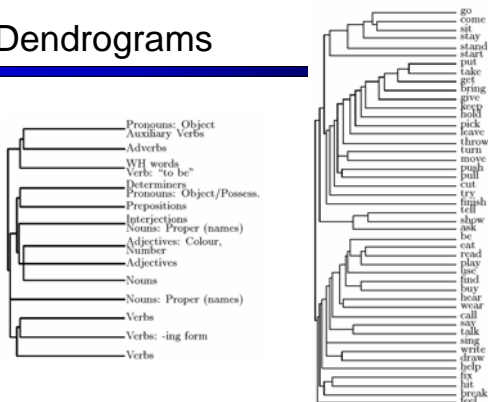
Distributional Clustering

- 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]

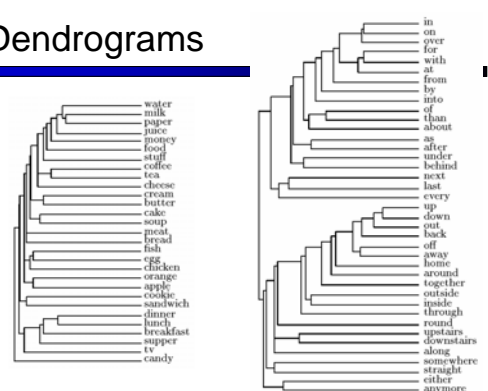
Nearest Neighbors

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

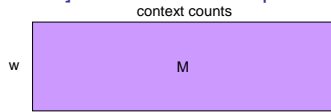


Dendrograms

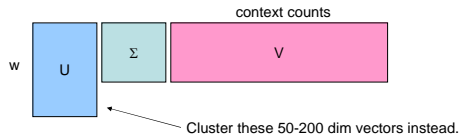


Vector Space Version

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

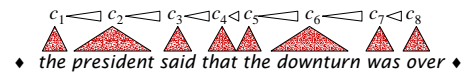
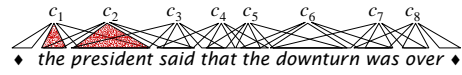


- Vectors too sparse, use SVD to reduce



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)