# Statistical NLP
## Spring 2010

University of California
Berkeley

## Lecture 7: POS / NER Tagging

Dan Klein – UC Berkeley

---

# Feature-Rich Sequence Models

- Problem: HMMs make it hard to work with arbitrary features of a sentence

- Example: name entity recognition (NER)

PER PER O    O  O  O      O     O     ORG       O  O  O  O  O LOC LOC O

Tim Boon has signed a contract extension with Leicestershire which will keep him at Grace Road .

### Local Context

|       | Prev | Cur   | Next |
|-------|------|-------|------|
| State | Other | ??? | ??? |
| Word  | at   | Grace | Road |
| Tag   | IN   | NNP   | NNP  |
| Sig   | x    | Xx    | Xx   |

# MEMM Taggers

- Idea: left-to-right local decisions, condition on previous tags and also entire input

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

  - Train up $P(t_i|w, t_{i-1}, t_{i-2})$ as a normal maxent model, then use to score sequences
  - This is referred to as an MEMM tagger [Ratnaparkhi 96]
  - Beam search effective!  (Why?)
  - What about beam size 1?

# Decoding

- Decoding MEMM taggers:
  - Just like decoding HMMs, different local scores
  - 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 (MEMMs):

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

- General:

$$\delta_i(s) = \arg\max_{s'} \phi_i(s', s)\delta_{i-1}(s')$$

# Maximum Entropy II

- Remember: maximum entropy objective

$$L(\mathbf{w}) = \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

- Problem: lots of features allow perfect fit to training set
- Regularization (compare to smoothing)

$$\max_{\mathbf{w}} \ \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right) - k||\mathbf{w}||^2$$

# Derivative for Maximum Entropy

$$L(\mathbf{w}) = -k||\mathbf{w}||^2 + \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}_n} = -2k\mathbf{w}_n + \sum_i \left( \mathbf{f}_i(\mathbf{y}^i)_n - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i)\mathbf{f}_i(\mathbf{y})_n \right)$$

Big weights are bad

Total count of feature n
in correct candidates

Expected count of
feature n in predicted
candidates

# Example: NER Regularization

### Feature Weights

Because of regularization term, the more common prefixes have larger weights even though entire-word features are more specific.

| Feature Type | Feature | PERS | LOC |
|---|---|---|---|
| Previous word | *at* | -0.73 | 0.94 |
| Current word | *Grace* | 0.03 | 0.00 |
| Beginning bigram | *<G* | 0.45 | -0.04 |
| Current POS tag | NNP | 0.47 | 0.45 |
| Prev and cur tags | IN NNP | -0.10 | 0.14 |
| Previous state | Other | -0.70 | -0.92 |
| Current signature | Xx | 0.80 | 0.46 |
| Prev state, cur sig | O-Xx | 0.68 | 0.37 |
| Prev-cur-next sig | x-Xx-Xx | -0.69 | 0.37 |
| P. state - p-cur sig | O-x-Xx | -0.20 | 0.82 |
| … | | | |
| **Total:** | | **-0.58** | **2.68** |

### Local Context

| | Prev | Cur | Next |
|---|---|---|---|
| State | Other | ??? | ??? |
| Word | at | Grace | Road |
| Tag | IN | NNP | NNP |
| Sig | x | Xx | Xx |

---

# Perceptron Taggers

[Collins 01]

- Linear models:

$$\text{score}(\mathbf{t}|\mathbf{w}) = \lambda^\top f(\mathbf{t}, \mathbf{w})$$

- … that decompose along the sequence

$$= \lambda^\top \sum_i f(t_i, t_{i-1}, \mathbf{w}, i)$$

- … allow us to predict with the Viterbi algorithm

$$\mathbf{t}^* = \arg\max_\mathbf{t} \ \text{score}(\mathbf{t}|\mathbf{w})$$

- … which means we can train with the perceptron algorithm (or related updates, like MIRA)

# Conditional Random Fields

- Make a maxent model over entire taggings
    - MEMM

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

    - CRF

$$P(\mathbf{t}|\mathbf{w}) = \frac{1}{Z(\mathbf{w})} \exp\left(\lambda^\top f(\mathbf{t}, \mathbf{w})\right)$$

$$= \frac{1}{Z(\mathbf{w})} \exp\left(\lambda^\top \sum_i f(t_i, t_{i-1}, \mathbf{w}, i)\right)$$

$$= \frac{1}{Z(\mathbf{w})} \prod_i \phi_i(t_i, t_{i-1})$$

# 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 of each feature (for example the number of times the label pair *DT-NN* occurs, or the number of times *NN-interest* occurs)

- Critical quantity: counts of posterior marginals:

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

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

# Computing Posterior Marginals

- How many (expected) times is word w tagged with s?

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

- How to compute that marginal?

$$\alpha_i(s) = \sum_{s'} \phi_i(s', s)\alpha_{i-1}(s')$$

$$\beta_i(s) = \sum_{s'} \phi_{i+1}(s, s')\beta_{i+1}(s')$$

$$P(t_i = s | \mathbf{w}) = \frac{\alpha_i(s)\beta_i(s)}{\alpha_N(\text{END})}$$

(^) (^) (^) (^) (^) (^)
(N) (N) (N) (N) (N) (N)
(V) (V) (V) (V) (V) (V)
(J) (J) (J) (J) (J) (J)
(D) (D) (D) (D) (D) (D)
($) ($) ($) ($) ($) ($)

START   Fed   raises   interest   rates   END

---

# 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 $\to$ NN : DT __
    - VBD $\to$ 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 definitely not the most accurate: 96.6% / 82.0 %

# TBL Tagger II

- **What gets learned? [from Brill 95]**

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

| | Change Tag | | |
|---|---|---|---|
| # | From | To | Condition |
| 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: 99% on *their* tag set
  - Linguistic representation matters…
  - … but it's easier to win when you make up the rules

```
walk
    walk <SV> <SVO> V SUBJUNCTIVE VFIN
    walk <SV> <SVO> V IMP VFIN
    walk <SV> <SVO> V INF
    walk <SV> <SVO> V PRES -SG3 VFIN
    walk N NOM SG
```

```
walk V-SUBJUNCTIVE V-IMP V-INF
    V-PRES-BASE N-NOM-SG
```

# 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:

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

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

- Same quantities we needed to train a CRF!

# EM for HMMs: Quantities

- Total path values (correspond to probabilities here):

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

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

# EM for HMMs: Process

- From these quantities, can compute expected transitions:

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

- 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* ♦

| | |
|---|---|
| *president* | *the __ of* |
| *president* | *the __ said* |
| *governor* | *the __ of* |
| *governor* | *the __ appointed* |
| *said* | *sources __* ♦ |
| *said* | *president __ that* |
| *reported* | *sources __* ♦ |

*president*
*governor*

*the*
*a*

*said*
*reported*

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

11

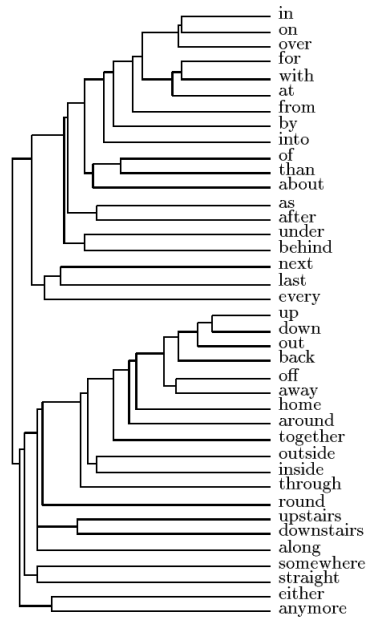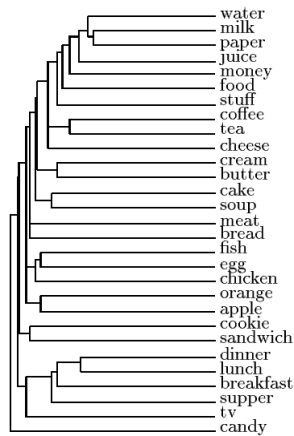# Distributional Clustering

- 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
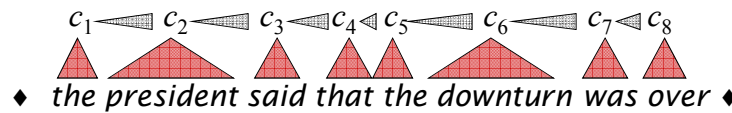    - Various formulations, e.g. [Lee and Pereira 99]

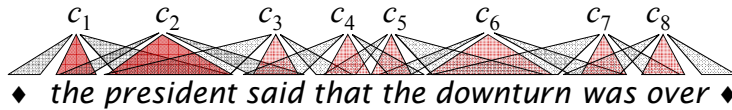# Nearest Neighbors

| 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



# A Probabilistic Version?

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



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



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

# What Else?

- Various newer ideas:
  - Context distributional clustering [Clark 00]
  - Morphology-driven models [Clark 03]
  - Contrastive estimation [Smith and Eisner 05]
  - Feature-rich induction [Haghighi and Klein 06]

- 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)