

## Lecture 8 – Tagging

### Goals

1. Generative Tagging
  - a. Parameterization
  - b. Estimation
  - c. Inference
2. Discriminative tagging
  - a. Max Entropy
  - b. Simple RNN
  - c. Bidirectional RNN

### Generative Tagging

- Sentence: I love white dogs
- Tags: Pr, Vb Adj, N
- Baseline: most frequent tag – already gets 90%+ accuracy
- We are given
  - o  $S = w_1, \dots, w_n$
  - o  $T = y_1, \dots, y_n$
- Parameterization
  - o  $P(S, T) = p(w_1, \dots, w_n, y_1, \dots, y_n) = \prod p(w_i | w_1, \dots, w_n, y_1, \dots, y_i) p(y_i | w_1, \dots, w_n, y_1, \dots, y_{i-1})$
- Markov Assumption
  - o  $P(S, T) = \prod p(w_i | y_i) p(y_i | y_{i-1})$ 
    - $p(w_i | y_i)$  called the emission probability
    - $p(y_i | y_{i-1})$  called the transition probability
  - o Number of parameters:  $VT + T^2$
- Estimation
  - o MLE of transition probability:  $p(Vb | Pr) = \frac{\text{count}(Pr, Vb)}{\text{count}(Pr)}$
  - o MLE of emission probability:  $p(I | Pr) = \frac{\text{count}(I, Pr)}{\text{count}(Pr)}$
  - o Remember there are smoothing issues
  - o Want to use EM...but
- Inference Problem
  - o We are looking for  $T^* = \text{argmax}_T (P(S, T))$ 
    - Assume  $p(t_i | t_{i-1})$  and  $p(w_i | t_i)$  are given
  - o Solution: Viterbi Algorithm
    - Let  $n$  be the length of the sequence
    - $\pi[i, t] \rightarrow \max \log(\text{probability of a sequence that ends at position } i \text{ with tag } t)$
    - Goal:  $\max_{t \in T} \pi[n, t]$
    - Base Case:  $\pi[0, *] = \log(1), \pi[0, t] = \log(0) = -\infty$

- Recursive Case
  - $\pi[i, t] = \max_{t_{prev}} \pi[i_{prev}, t_{prev}] + \log(p(t|t_{prev})) + \log(p(i|t))$
- Complexity:  $O(nT^2)$
- Now we can do estimation with EM
  - Take a random guess of the parameters, and compute the MLE efficiently with the Viterbi algorithm
  - Bad results with random initialization
  - Can get good results with good initialization

## Discriminative Tagging

- Now we estimate
  - $p(T|S) = p(y_1, \dots, y_n | w_1, \dots, w_n) = \prod_i p(y_i | w_1, \dots, w_n, y_1, \dots, y_{i-1})$
- Max Entropy
  - $p(y|w) = \frac{1}{z(\theta, w)} e^{\theta f(y, w)}$ 
    - Indicator function  $f(y, w)$  of features 1 ...  $K$
    - i.e.  $f_1(y, w) = \begin{cases} 1, & \text{if } y = \text{noun}, w \text{ is capitalized} \\ 0, & \text{o.w.} \end{cases}$
- NN Structure for Max Entropy
  - Final layer: softmax, each unit represents  $p(y_t = k | w)$
  - Vector representing indicator function for  $y, w$