# POS tagging

#### CMSC 723 / LING 723 / INST 725

Marine Carpuat

# POS tagging Sequence labeling with the perceptron

### Sequence labeling problem

- Input:
  - sequence of tokens  $x = [x_1 ... x_L]$
  - Variable length L
- Output (aka label):
  - sequence of tags  $y = [y_1 \dots y_L]$
  - # tags = K
  - Size of output space?

### **Structured Perceptron**

- Perceptron algorithm can be used for sequence labeling
- But there are challenges
  - How to compute argmax efficiently?
  - What are appropriate features?
- Approach: leverage structure of output space

Solving the argmax problem for sequences with dynamic programming

- x = " monsters eat tasty bunnies "
- y = noun verb adj noun

- Efficient algorithms possible if the feature function decomposes over the input
- This holds for unary and markov features used for POS tagging

### Feature functions for sequence labeling

- x = " monsters eat tasty bunnies "
- y = noun verb adj noun

- Standard features of POS tagging
  - Unary features: # times word w has been labeled with tag I for all words w and all tags I
  - Markov features: # times tag I is adjacent to tag I' in output for all tags I and I'

• Size of feature representation is constant wrt input length

## Solving the argmax problem for sequences



- Trellis sequence labeling
  - Any path represents a labeling of input sentence
  - Gold standard path in red
  - Each edge receives a weight such that adding weights along the path corresponds to score for input/ouput configuration
- Any max-weight max-weight path algorithm can find the argmax
  - e.g. Viterbi algorithm O(LK<sup>2</sup>)

## Defining weights of edge in treillis



 Weight of edge that goes from time l-1 to time l, and transitions from y to y'

$$w \cdot \phi_l(x, \cdots \circ y \circ y')$$

## Dynamic program

 Define: the score of best possible output prefix up to and including position I that labels the I-th word with label k

 $\alpha_{l,k} = \max_{\hat{y}_{1:l-1}} \boldsymbol{w} \cdot \phi_{1:l}(\boldsymbol{x}, \hat{\boldsymbol{y}} \circ \boldsymbol{k})$ 

• With decomposable features, alphas can be computed recursively

$$\alpha_{l+1,k} = \max_{k'} \left[ \alpha_{l,k'} + \boldsymbol{w} \cdot \phi_{l+1}(\boldsymbol{x}, \langle \dots, k', k \rangle) \right]$$

$$\alpha_{0,k} = 0 \quad \forall k \tag{17.41}$$

$$\zeta_{0,k} = \emptyset \quad \forall k \tag{17.42}$$

the score for any empty sequence is zero

$$\alpha_{l+1,k} = \max_{\hat{y}_{1:l}} w \cdot \phi_{1:l+1}(x, \hat{y} \circ k)$$
(17.43)

separate score of prefix from score of position I+1

$$= \max_{\hat{y}_{1:l}} w \cdot \left( \phi_{1:l}(x, \hat{y}) + \phi_{l+1}(x, \hat{y} \circ k) \right)$$
(17.44)

distributive law over dot products

$$= \max_{\hat{y}_{1:l}} \left[ \boldsymbol{w} \cdot \boldsymbol{\phi}_{1:l}(\boldsymbol{x}, \hat{\boldsymbol{y}}) + \boldsymbol{w} \cdot \boldsymbol{\phi}_{l+1}(\boldsymbol{x}, \hat{\boldsymbol{y}} \circ \boldsymbol{k}) \right]$$
(17.45)

separate out final label from prefix, call it k'

$$= \max_{\hat{y}_{1:l-1}} \max_{k'} \left[ w \cdot \phi_{1:l}(x, \hat{y} \circ k') + w \cdot \phi_{l+1}(x, \hat{y} \circ k' \circ k) \right]$$
(17.46)

swap order of maxes, and last term doesn't depend on prefix

$$= \max_{k'} \left[ \left[ \max_{\hat{y}_{1:l-1}} \boldsymbol{w} \cdot \phi_{1:l}(\boldsymbol{x}, \hat{\boldsymbol{y}} \circ k') \right] + \boldsymbol{w} \cdot \phi_{l+1}(\boldsymbol{x}, \langle \dots, k', k \rangle) \right]$$
(17.47)

apply recursive definition

$$= \max_{k'} \left[ \alpha_{l,k'} + \boldsymbol{w} \cdot \boldsymbol{\phi}_{l+1}(\boldsymbol{x}, \langle \dots, k', k \rangle) \right]$$
(17.48)

Algorithm 42 ArgmaxForSequences(x, w)

1:  $L \leftarrow LEN(x)$ <sup>2:</sup>  $\alpha_{l,k} \leftarrow o, \quad \zeta_{k,l} \leftarrow o, \quad \forall k = 1 \dots K, \quad \forall l = 0 \dots L$ // initialize variables  $_{3:}$  for  $l = 0 \dots L^{-1}$  do for  $k = 1 \dots K$  do 4:  $\alpha_{l+1,k} \leftarrow \max_{k'} \left[ \alpha_{l,k'} + w \cdot \phi_{l+1}(x, \langle \dots, k', k \rangle) \right]$ // recursion: 5: // here,  $\phi_{l+1}(\ldots k', k \ldots)$  is the set of features associated with // output position l + 1 and two adjacent labels k' and k at that position  $\zeta_{l+1,k} \leftarrow$  the k' that achieves the maximum above // store backpointer 6: end for 7: 8: end for 9:  $\mathbf{y} \leftarrow \langle 0, 0, \dots, 0 \rangle$ // initialize predicted output to L-many zeros // extract highest scoring final label 10:  $y_L \leftarrow \operatorname{argmax}_k \alpha_{L,k}$ <sup>111</sup> for  $l = L - 1 \dots 1$  do // traceback  $\zeta$  based on  $y_{l+1}$  $y_l \leftarrow \zeta_{l,y_{l+1}}$ 13: end for 14: return y// return predicted output

monsters

eat

tastv

bunnies

# A more general approach for argmax Integer Linear Programming

• ILP: optimization problem of the form, for a fixed vector a

 $\max_{z} \quad a \cdot z \quad \text{subj. to} \quad \text{linear constraints on } z$ 

• With integer constraints

- Pro: can leverage well-engineered solvers (e.g., Gurobi)
- Con: not always most efficient

# POS tagging as ILP

• Markov features as binary indicator variables

 $z_{l,k',k} = \mathbf{1}[\text{label } l \text{ is } k \text{ and label } l - 1 \text{ is } k']$ 

- Output sequence: y(z) obtained by reading off variables z
- Define a such that a.z is equal to score

 $a_{l,k',k} = \boldsymbol{w} \cdot \phi_l(\boldsymbol{x}, \langle \dots, k', k \rangle)$ 

- Enforcing constraints for well formed solutions
- 1. That all the *z*s are binary. That's easy: just say  $z_{l,k',k} \in \{0,1\}$ , for all l,k',k.
- 2. That for a given position l, there is exactly one active z. We can do this with an equality constraint:  $\sum_{k} \sum_{k'} z_{l,k',k} = 1$  for all l.
- 3. That the *z*s are internally consistent: if the label at position 5 is supposed to be "noun" then both  $z_{5,...}$  and  $z_{6,...}$  need to agree on this. We can do this as:  $\sum_{k'} z_{l,k',k} = \sum_{k''} z_{l+1,k,k''}$  for all *l*, *k*. Effectively what this is saying is that  $z_{5,?,\text{verb}} = z_{6,\text{verb},?}$  where the "?" means "sum over all possibilities."

## Sequence labeling

- Structured perceptron
  - A general algorithm for structured prediction problems such as sequence labeling
- The Argmax problem
  - Efficient argmax for sequences with Viterbi algorithm, given some assumptions on feature structure
  - A more general solution: Integer Linear Programming
- Loss-augmented argmax
  - Hamming Loss

# POS tagging

#### CMSC 723 / LING 723 / INST 725

Marine Carpuat