



COMPUTER SCIENCE
UNIVERSITY OF MARYLAND

Sequence Labeling with the Structured Perceptron

CMSC 470

Marine Carpuat

POS tagging

Sequence labeling with the perceptron

Sequence labeling problem

- Input:
 - sequence of tokens $x = [x_1 \dots x_L]$
 - 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

Perceptron algorithm remains the same as for multiclass classification

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$$

Note: CIML denotes

- the weight vector as w instead of θ
- The feature function as $\Phi(x, y)$ instead of $f(x, y)$

Algorithm 3 Perceptron learning algorithm

```
1: procedure PERCEPTRON( $\mathbf{x}^{(1:N)}, y^{(1:N)}$ )
2:    $t \leftarrow 0$ 
3:    $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ 
4:   repeat
5:      $t \leftarrow t + 1$ 
6:     Select an instance  $i$ 
7:      $\hat{y} \leftarrow \operatorname{argmax}_y \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$ 
8:     if  $\hat{y} \neq y^{(i)}$  then
9:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$ 
10:    else
11:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
12:    until tired
13:  return  $\boldsymbol{\theta}^{(t)}$ 
```

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 l for all words w and all tags l
 - **Markov features:** # times tag l is adjacent to tag l' in output for all tags l and l'
- Size of feature representation is constant wrt input length

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

Decomposition of structure

- Features decompose over the input if $\phi(x, y) = \sum_{l=1}^L \phi_l(x, y)$

Feature function that only includes features about position l

- If features decompose over the input, structures (x, y) can be scored incrementally

$$w \cdot \phi(x, y) = w \cdot \sum_{l=1}^L \phi_l(x, y)$$

decomposition of structure

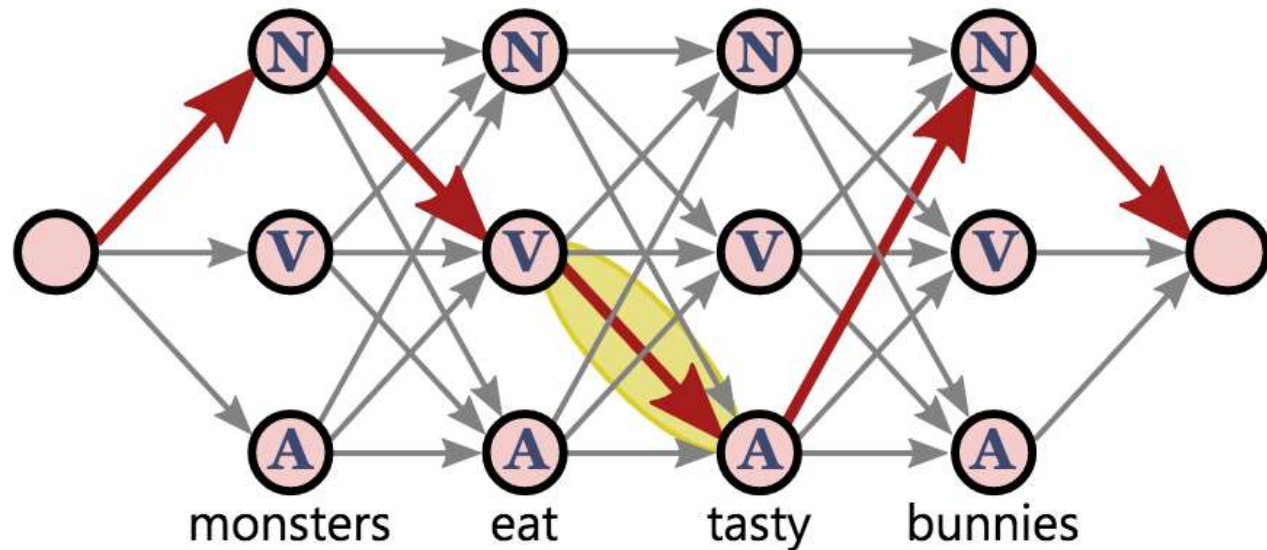
$$= \sum_{l=1}^L w \cdot \phi_l(x, y)$$

associative law

Decomposition of structure: Lattice/trellis representation

$x =$ " monsters eat tasty bunnies "

$y =$ noun verb adj noun



- 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/output configuration
- Any max-weight path algorithm can find the argmax
 - We'll describe the Viterbi algorithm

Dynamic programming solution relies on recursively computing prefix scores $\alpha_{l,k}$

Score of best possible output prefix, up to and including position l , that labels the l -th word as label k

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

Sequence of labels of length $l-1$

Features for sequence starting at position 1 up to and including position l

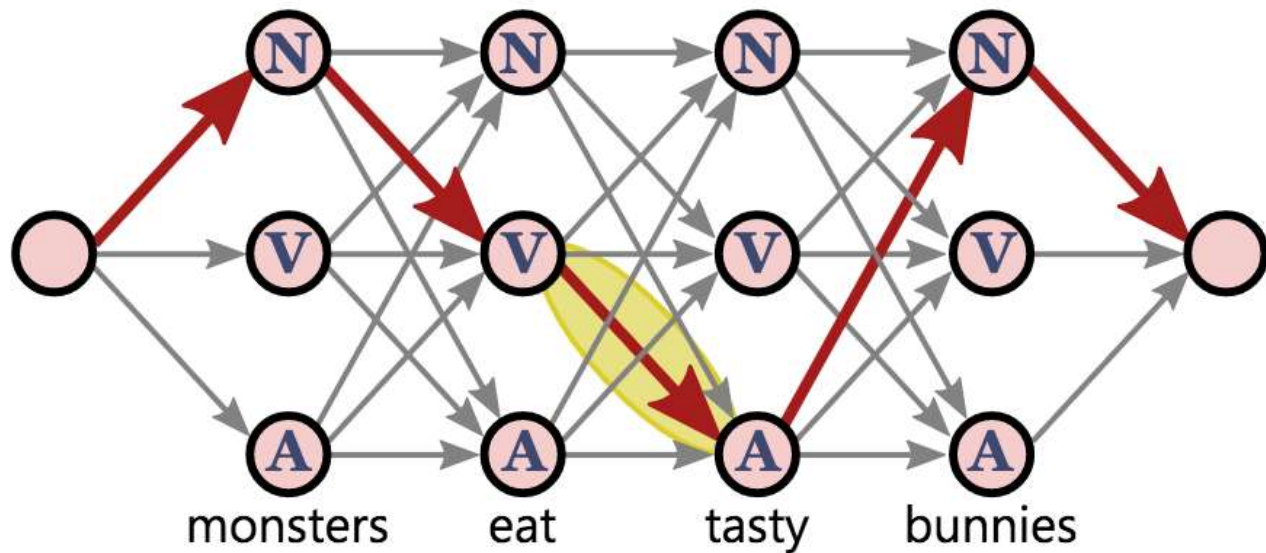
Sequence of length l obtained by adding k at the end.

Computing prefix scores $\alpha_{l,k}$

Example

$x =$ " monsters eat tasty bunnies "

$y =$ noun verb adj noun



Let's compute $\alpha_{3,A}$ given

- Prefix scores for length 2

$$\alpha_{2,N} = 2, \alpha_{2,V} = 9, \alpha_{2,A} = -1$$

- Unary feature weights

$$w_{tasty/A} = 1.2$$

- Markov feature weights

$$w_{N,A} = -5, w_{V,A} = 2.5, w_{A,A} = 2.2$$

Dynamic programming solution relies on recursively computing prefix scores $\alpha_{l,k}$

Score of best possible output prefix, up to and including position $l+1$, that labels the $(l+1)$ -th word as label k

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

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

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

Backpointer to the label that achieves the above maximum

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

$$\zeta_{l+1,k} = \operatorname{argmax}_{k'} \left[\alpha_{l,k'} + w \cdot \phi_{l+1}(\mathbf{x}, \langle \dots, k', k \rangle) \right]$$

Viterbi algorithm

Assumptions:

- Unary features
- Markov features based on 2 adjacent labels

Runtime: $O(LK^2)$

Algorithm 42 ARGMAXFORSEQUENCES(x, w)

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

Exercise: Impact of feature definitions

- Consider a structured perceptron with the following features
 - # times word w has been labeled with tag l for all words w and all tags l
 - # times word w has been labeled with tag l when it follows word w' for all words w, w' and all tags l
 - # times tag l occurs in the sequence (l', l'', l) in the output for all tags l, l', l''
- What is the dimension of the perceptron weight vector?
- Can we use dynamic programming to compute the argmax?

Recap: POS tagging

- An example of sequence labeling tasks
- Requires a predefined set of POS tags
 - Penn Treebank commonly used for English
 - Encodes some distinctions and not others
- Given annotated examples, we can address sequence labeling with multiclass perceptron
 - but computing the argmax naively is expensive
 - constraints on the feature definition make efficient algorithms possible
 - Viterbi algorithm for unary and markov features