

# Sequence Labeling (III)

## Conditional Random Fields

**CS 690N, Spring 2018**

Advanced Natural Language Processing

<http://people.cs.umass.edu/~brenocon/anlp2018/>

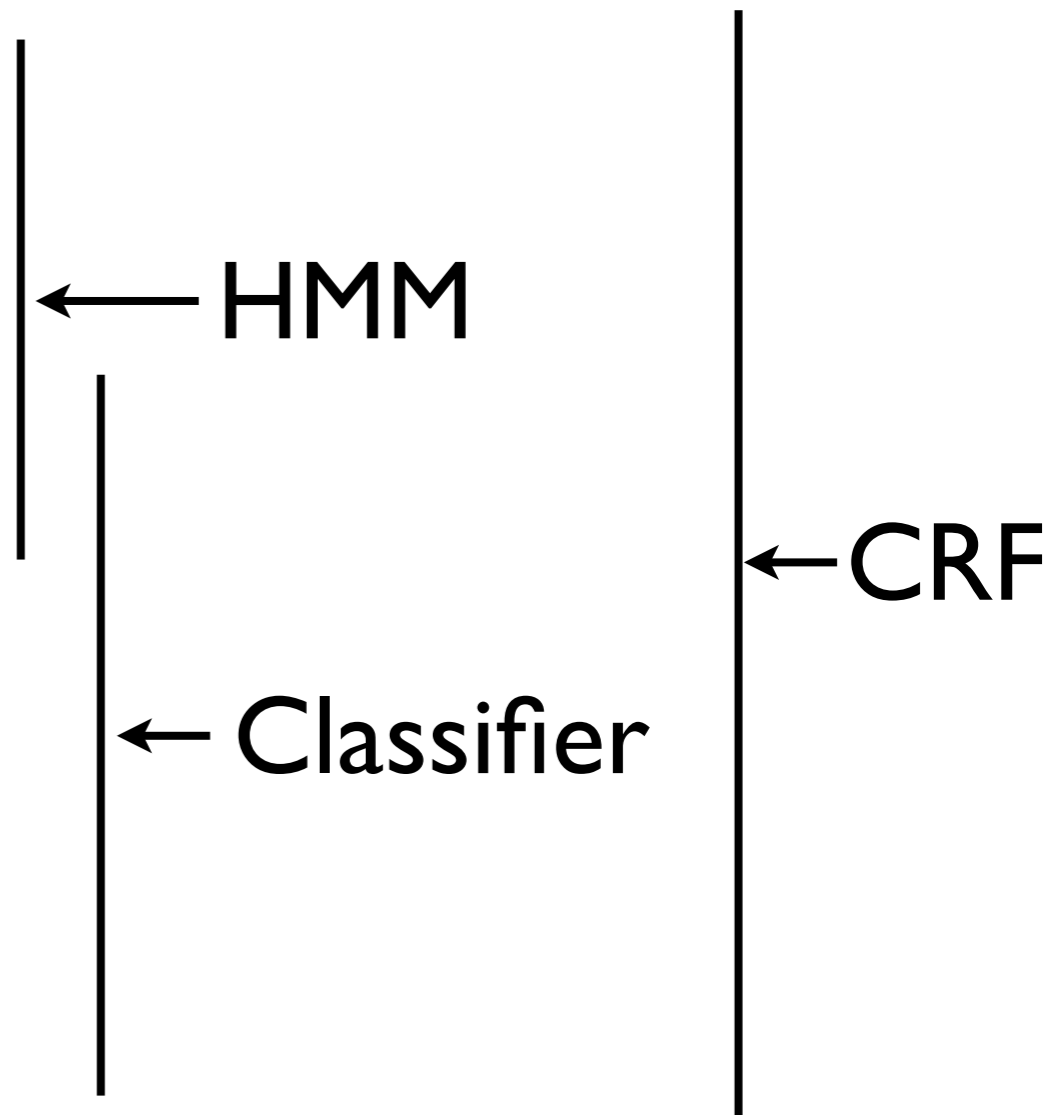
**Brendan O'Connor**

College of Information and Computer Sciences

University of Massachusetts Amherst

# How to build a POS tagger?

- Sources of information:
  - POS tags of surrounding words: syntactic context
  - The word itself
  - Features, etc.!
    - Word-internal information
    - Features from surrounding words
    - External lexicons
    - Embeddings, LSTM states



- Seq. labeling as log-linear **structured prediction**

$$\hat{\mathbf{y}}_{1:M} = \operatorname{argmax}_{\mathbf{y}_{1:M} \in \mathcal{Y}(\mathbf{w}_{1:M})} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}),$$

- Example: the **Hidden Markov model**

$$p(\mathbf{w}, \mathbf{y}) = \prod_t p(y_t | y_{t-1}) p(w_t | y_t)$$

- Efficiently supports operations via dynamic programming –  
because of **local (Markovian) assumptions**

- $P(\mathbf{w})$ : Likelihood (generative model)
  - Forward algorithm
- $P(\mathbf{y} | \mathbf{w})$ : Predicted sequence (“decoding”)
  - Viterbi algorithm
- $P(y_m | \mathbf{w})$ : Predicted tag marginals
  - Forward-Backward algorithm
  - Supports EM for unsupervised HMM learning

- Seq. labeling as log-linear **structured prediction**

$$\hat{\mathbf{y}}_{1:M} = \operatorname{argmax}_{\mathbf{y}_{1:M} \in \mathcal{Y}(\mathbf{w}_{1:M})} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}),$$

- Example: the **Hidden Markov model**

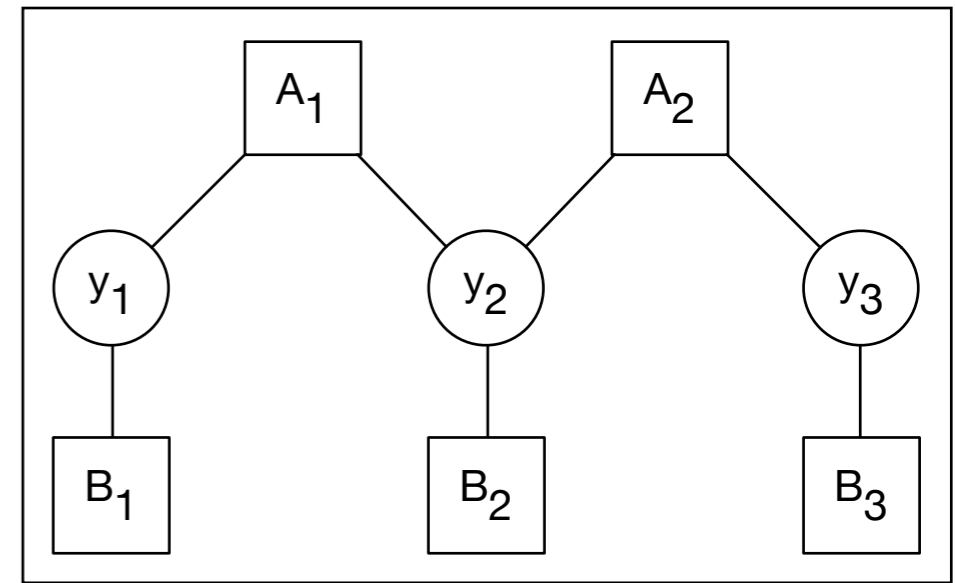
$$p(\mathbf{w}, \mathbf{y}) = \prod_t p(y_t | y_{t-1}) p(w_t | y_t)$$

- Today: **Conditional Random Fields**

$$p(\mathbf{y} | \mathbf{w}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}))}{\sum_{\mathbf{y}'_{1:M} \in \mathcal{Y}(\mathbf{w}_{1:M})} \exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, \mathbf{y}'_{1:M}))}$$

\* for carefully chosen **f**

# HMM as log-linear



$$p(y, w) = \prod p(w_y | y_t) p(y_t | y_{t-1})$$

$$\log p(y, w) = \sum_t \log p(w_t | y_t) + \log p(y_t | y_{t-1})$$

↑  
 $G(y)$   
goodness

↑  
 $B_t(y_t)$   
emission factor score

↑  
 $A(y_{t-1}, y_t)$   
transition factor score

$$\log p(y, w) = \sum_t \phi_t(y_{t-1}, y_t)$$

↑  
pair factor score

Decoding problem (Viterbi algorithm)  $\arg \max_{y^* \in outputs(x)} G(y^*)$

# HMM as log-linear

- HMM as a joint log-linear model

$$P(y, w) = \prod_t P(y_t | y_{t-1}) P(w_t | y_t)$$

$$P(y, w) = \exp(\theta^\top f(y, w))$$

$$f(y, w) = \sum_t f(y_{t-1}, y_t, w_t) \quad \begin{array}{l} \text{Local features only!} \\ \text{(Allows efficient inference)} \end{array}$$

↓  
e.g.  $\{(N, V):1, (V, \text{dog}):1\}$   
What are the weights?

- This implies the conditional is also log-linear

$$P(y | w) \propto \exp(\theta^\top f(y, w))$$

# From HMMs to CRFs

- **1. Discriminative learning:** take HMM features, but set weights to maximize *conditional LL* of labels
- **2. More features:** affix, positional, feature templates, embeddings, etc.
- For efficient inference: make sure to **preserve Markovian structure** within the feature function (e.g. first-order CRF)

# Learning a CRF

- Gradient descent on negative **conditional LL**
  - Log-linear gradient:  
sum over all possible predicted structures  
(Forward-Backward for marginalization)
- Non-probabilistic losses: compare gold structure to only one predicted structure
  - Structured perceptron algorithm:  
Collins, 2002 (just got Test of Time award)
  - Structured SVM (hinge loss)
  - (Viterbi for best-structure)



# Learning a CRF: max CLL

$$\log p_{\theta}(y | w) = \theta^{\top} f(y, w) - \log \sum_{y'} \exp(\theta^{\top} f(y', w))$$

$$\frac{\partial \log p_{\theta}(\dots)}{\partial \theta_j} = f_j(y, w) - \sum_{y'} p_{\theta}(y' | w) f_j(y', w)$$

- Apply local decomposition

# Learning a CRF: max CLL

$$\log p_{\theta}(y | w) = \theta^{\top} f(y, w) - \log \sum_{y'} \exp(\theta^{\top} f(y, w))$$

$$\frac{\partial \log p_{\theta}(\dots)}{\partial \theta_j} = f_j(y, w) - \sum_{y'} p_{\theta}(y' | w) f_j(y', w)$$

- Apply local decomposition

$$= \left( \sum_t f_j(y_{t-1}, y_t, w_t) \right) - \sum_{y'} p_{\theta}(y' | w) \sum_t f_j(y'_{t-1}, y'_t, w_t)$$

# Learning a CRF: max CLL

$$\log p_{\theta}(y | w) = \theta^{\top} f(y, w) - \log \sum_{y'} \exp(\theta^{\top} f(y, w))$$

$$\frac{\partial \log p_{\theta}(\dots)}{\partial \theta_j} = f_j(y, w) - \sum_{y'} p_{\theta}(y' | w) f_j(y', w)$$

- Apply local decomposition

$$= \left( \sum_t f_j(y_{t-1}, y_t, w_t) \right) - \sum_{y'} p_{\theta}(y' | w) \sum_t f_j(y'_{t-1}, y'_t, w_t)$$

$$= \sum_t \left( f_j(y_{t-1}, y_t, w_t) - \sum_{y'_t, y'_{t-1}} p_{\theta}(y'_{t-1}, y'_t | w) f_j(y'_{t-1}, y'_t, w_t) \right)$$

# Learning a CRF: max CLL

$$\log p_{\theta}(y | w) = \theta^{\top} f(y, w) - \log \sum_{y'} \exp(\theta^{\top} f(y, w))$$

$$\frac{\partial \log p_{\theta}(\dots)}{\partial \theta_j} = f_j(y, w) - \sum_{y'} p_{\theta}(y' | w) f_j(y', w)$$

- Apply local decomposition

$$= \left( \sum_t f_j(y_{t-1}, y_t, w_t) \right) - \sum_{y'} p_{\theta}(y' | w) \sum_t f_j(y'_{t-1}, y'_t, w_t)$$

Real feature value



Expected feature value



$$= \sum_t \left( f_j(y_{t-1}, y_t, w_t) - \sum_{y'_t, y'_{t-1}} p_{\theta}(y'_{t-1}, y'_t | w) f_j(y'_{t-1}, y'_t, w_t) \right)$$



Tag marginals (to compute: forward-backward)

- stopped here 3/27

# (Log?-)linear Viterbi

$$\hat{y} = \operatorname{argmax}_y \theta^\top f(w, y) \quad f(w, y) = \sum_{m=1}^M f(w, y_m, y_{m-1}, m).$$

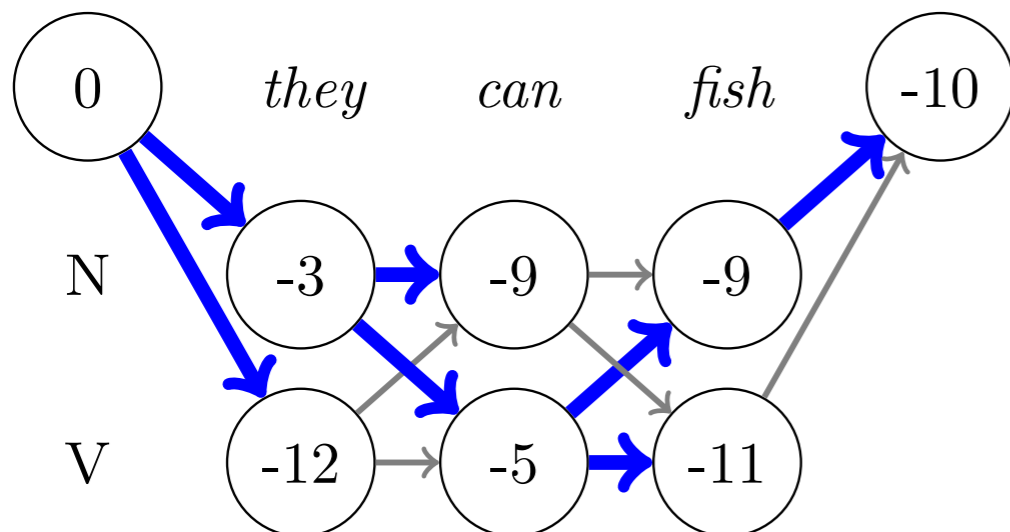
$$\max_y \theta^\top f(w, y) = \max_k v_M(k)$$

Score of best sequence ending in  $k$

$$v_m(k) \triangleq \max_{y_{1:m-1}} \theta^\top f(w, k, y_{m-1}, m) + \sum_{n=1}^{m-1} \theta^\top f(w, y_n, y_{n-1}, n)$$

$$= \max_{y_{m-1}} \theta^\top f(w, k, y_{m-1}, m)$$

$$+ \underbrace{\max_{y_{1:m-2}} \theta^\top f(w, y_{m-1}, y_{m-2}) + \sum_{n=1}^{m-2} \theta^\top f(w, y_n, y_{n-1}, n)}_{v_{m-1}(y_{m-1})}$$



	<i>they</i>	<i>can</i>	<i>fish</i>		N	V	◆
N	-2	-3	-3	◆	-1	-2	$-\infty$
V	-10	-1	-3	N	-3	-1	-12
				V	-1	-3	-1

(a) Weights for emission features.

(b) Weights for transition features.