

Sequence Labeling (II)

CS 690N, Spring 2017

Advanced Natural Language Processing

<http://people.cs.umass.edu/~brenocon/anlp2017/>

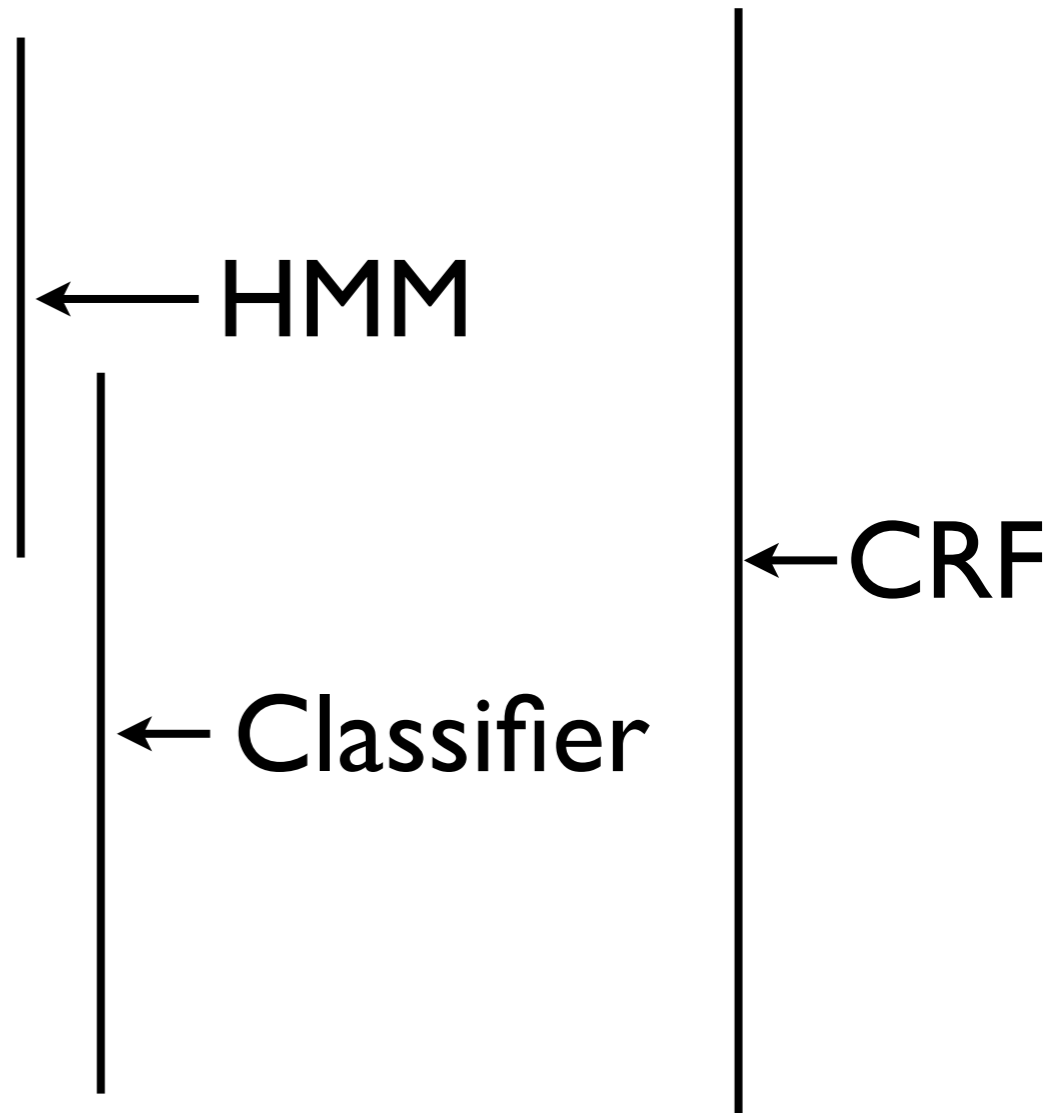
Brendan O'Connor

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!
 - Word-internal information
 - External lexicons
 - Features from surrounding words



- Sequence labeling as ***structured prediction***

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

- **Hidden Markov model**

- Fully generative, simple sequence model
- Supports many operations
 - $P(\mathbf{w})$: Likelihood (generative model)
 - Forward algorithm
 - $\operatorname{argmax}_y P(y \mid \mathbf{w})$: Predicted sequence (“decoding”)
 - Viterbi algorithm
 - $P(y_m \mid \mathbf{w})$: Predicted tag marginals
 - Forward-Backward algorithm
- The HMM is a type of log-linear model

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. $\{(N, V), (V, \text{dog})\}$

- The conditional is also log-linear: like we saw before, scoring just “outputs”:

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

(Log?-)linear Viterbi

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

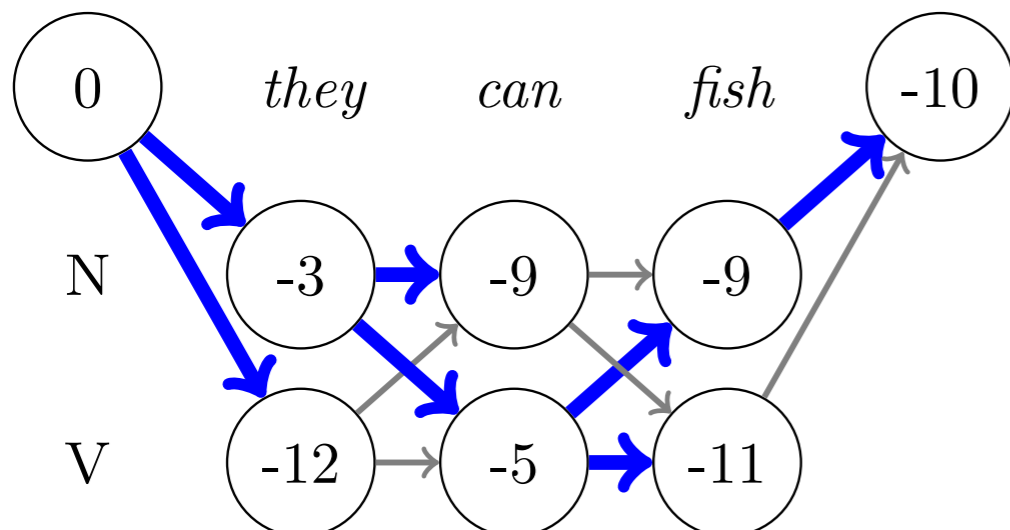
$$\max_y \theta^\top f(w, y) = \max_k v_M(k)$$

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.

Forward-Backward

- Purpose: compute
 - Tag marginals $p(y_t | w)$
 - Pair marginals $p(y_{t-1}, y_t | w)$
- Why?
 - Min Bayes Risk decoding
 - For each t , choose: $\operatorname{argmax}_k p(y_t=k | w)$
 - E-step for EM learning of unsupervised HMM
 - Feature expectations for supervised CRF

Learning a CRF

- Gradient descent on the neg. log-likelihood
 - Log-linear gradient: sum over all possible predicted structures
- Non-probabilistic losses: compare gold structure to only one predicted structure
 - Structured perceptron algorithm
 - Structured SVM (hinge loss)

Learning a CRF

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

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

- Apply local decomposition

Learning a CRF

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

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

- 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

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

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

- 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

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

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

- 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)