

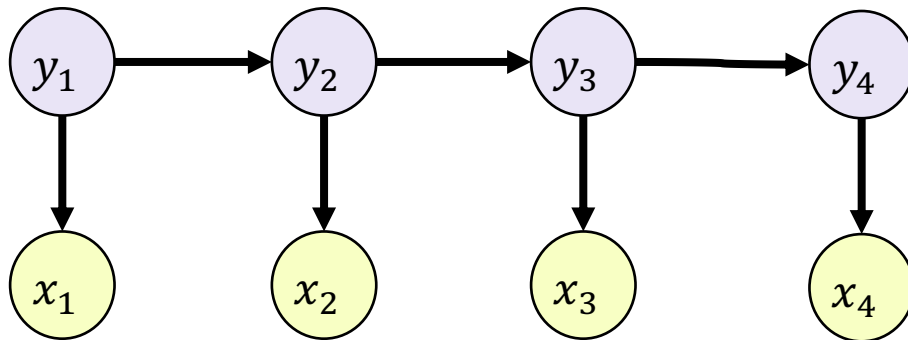
4.4 Conditional Random Fields

Edwin Simpson

Department of Computer Science,
University of Bristol, UK.

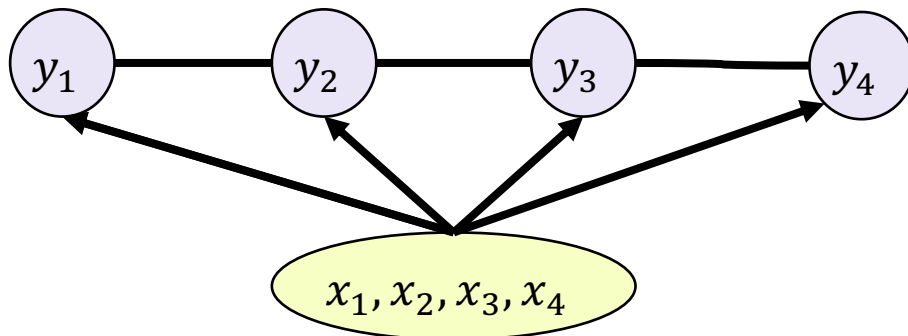
HMM as a Generative Model

- HMM is a **generative** model:
 - Learns the likelihoods of observations then apply Bayes' rule to predict tags
 - Benefits: closed form maximum likelihood estimates, interpretable, modular
 - Related generative approach for classification: Naïve Bayes



Discriminative Models: CRF

- Conditional Random Field (CRF) is **discriminative**:
 - Optimises predictive distribution $p(y|x)$
 - Related discriminative approach for classification: Logistic regression



CRF Prediction Function

Directly computes
probability of the sequence

Weights like in logistic regression

$$P(\mathbf{y}|\mathbf{x}) \propto \exp(\sum_{k=1}^K \theta_k F_k(\mathbf{x}, \mathbf{y}))$$

Global feature function to compute the feature from the sequence \mathbf{x}

Global Feature Function

Function for feature k

Compute a local function for each token in the sequence

$$F_k(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^N f_k(y_{i-1}, y_i, \mathbf{x}, i)$$

Local function can use previous tag, current tag, whole token sequence, and current position.

Local Feature Functions

- We can use any function that can extract a feature from $y_{i-1}, y_i, \mathbf{x}, i$
- Examples for POS tagging: $f_k(y_{i-1}, y_i, \mathbf{x}, i) =$

$[y_i = DET \text{ and } x_i = \text{"the"}]$

$[y_i = VB \text{ and } y_{i-1} = MD]$

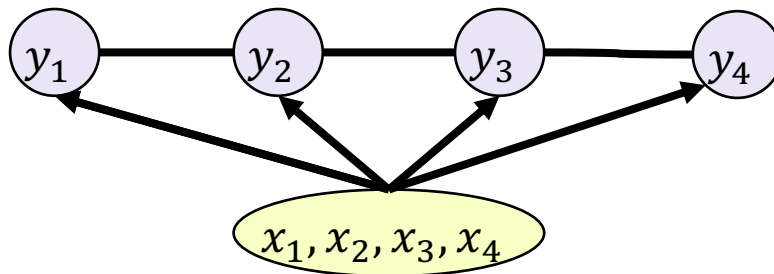
$[x_{i-1} = \text{"will"} \text{ and } x_{i+2} = \text{"bill"}]$

$[suffix(x_i) = \text{"ed"}]$

‘[...]’ notation means the value is 1 if $y_i = DET$ and $x_i = \text{"the"}$ and 0 otherwise

Training

- Stochastic gradient descent
- Forward-backward algorithm needed to compute gradients
- Training is expensive with computational complexity $\mathcal{O}(C^2N)$.
- Often more accurate than HMM but less suitable for online learning.



Summary

- Like NB, HMM is a generative model while CRF is discriminative, like logistic regression
- CRF is undirected, so can find globally optimal sequences
- It has higher training cost than HMMs