

Sequence labeling and hands on

Kan Min-Yen
Day 3 / Morning

(partial slide credits: Inna Weiner, Rongkun Shen, Jie Tang)

Recap

- **Information Retrieval**
 - Weighting words with respect to global and local importance
 - Represent both docs and queries as vectors
- **Introduced NLP as a machine learning problem**
- **Casts the problem as annotation and feature engineering**
 - Annotation requires clear policy and guidelines
 - Evaluation to assess performance and identify sources of error for improvement

Day Outline

Day 1

AM

- Applications' Input / Output
- Resources

PM

- Selected Toolkits
- Python Intro
- NLTK Hands-on

Day 2

AM

- Evaluation
- Annotation
- Information Retrieval
- ML Intro

PM

- Machine Learning
- SVM Hands-on

>>Day 3

AM

- Sequence Labeling
- CRF++ Hands-on

PM

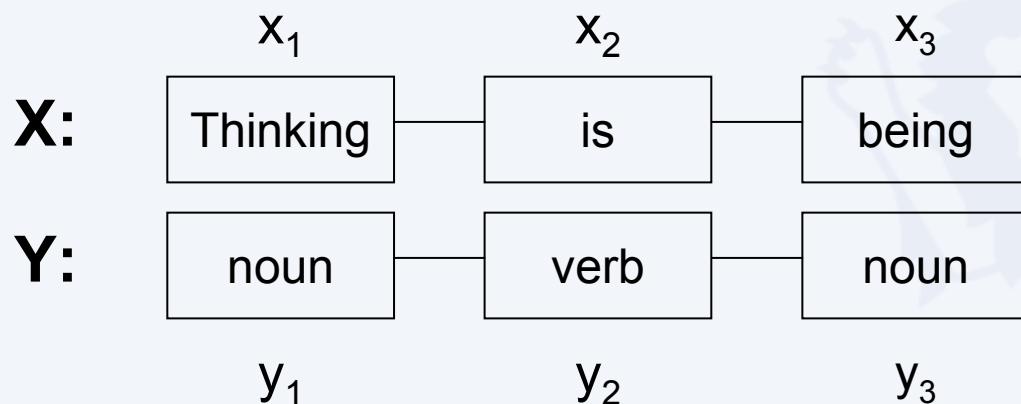
- Dimensionality Reduction
- Trends & Issues

Sequence Labeling Models

- **HMM**
 - Generative model
 - E.g. Ghahramani (1997), Manning and Schutze (1999)
- **MEMM**
 - Conditional model
 - E.g. Berger and Pietra (1996), McCallum and Freitag (2000)
- **CRFs**
 - Conditional model without label bias problem
 - Linear-Chain CRFs
 - E.g. Lafferty and McCallum (2001), Wallach (2004)
 - Non-Linear Chain CRFs
 - Modeling more complex interaction between labels: DCRFs, 2D-CRFs
 - E.g. Sutton and McCallum (2004), Zhu and Nie (2005)

Labeling Sequence Data Problem

- X is a random variable over data sequences
- Y is a random variable over label sequences
- Y_i is assumed to range over a finite label alphabet A
- **The problem:**
 - Learn how to give labels from a closed set Y to a data sequence X



Generative Probabilistic Models

- Learning problem:
Choose Θ to maximize *joint likelihood*:

$$L(\Theta) = \sum \log p_{\Theta}(y_i, x_i)$$

- The goal: maximization of the joint likelihood of training examples

$$y = \operatorname{argmax} p^*(y|x) = \operatorname{argmax} p^*(y,x)/p(x)$$

Joint likelihood

- Needs to enumerate all possible observation sequences

Markov Model

A **Markov process** or **model** assumes that we can predict the future based just on the present (or on a limited horizon into the past):

Let $\{X_1, \dots, X_T\}$ be a sequence of random variables taking values $\{1, \dots, N\}$ then the Markov properties are:

1. Limited Horizon:

$$P(X_{t+1}|X_1, \dots, X_t) = P(X_{t+1}|X_t) =$$

2. Time invariant (stationary):

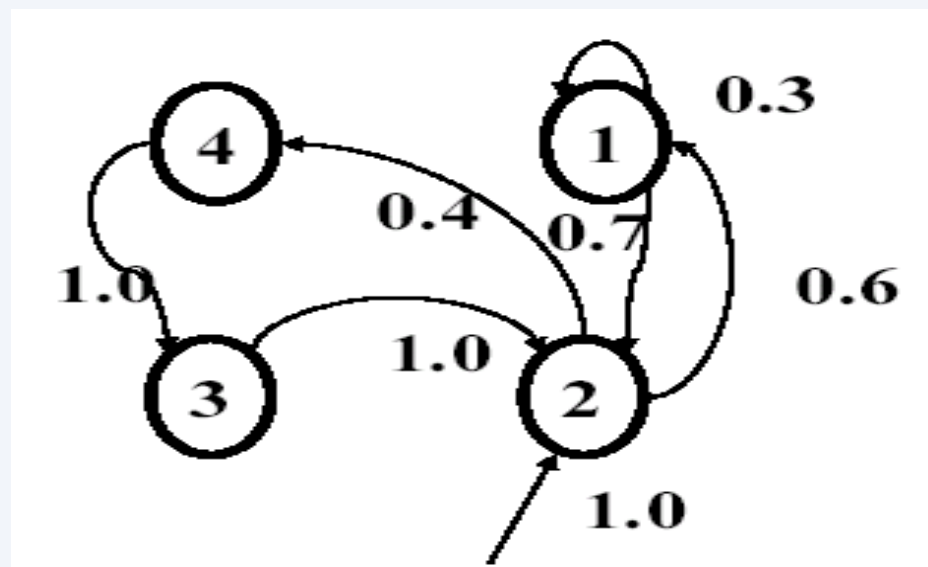
$$= P(X_2|X_1)$$

Describing a Markov Chain

Markov Chains can be described by the transition matrix **A** and the initial (start) probabilities **Q**:

$$A_{ij} = P(X_{t+1}=j|X_t=i)$$

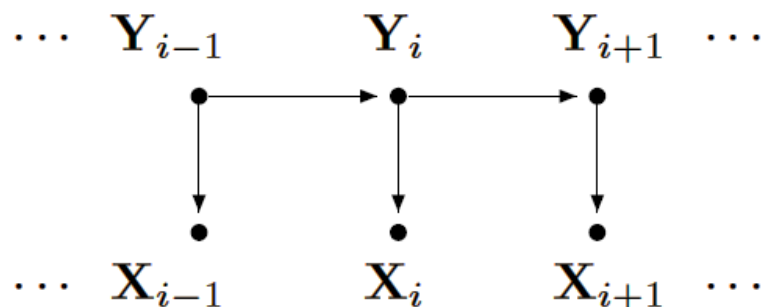
$$q_i = P(X_1=i)$$



Hidden Markov Model

- Do not observe the sequence that the model passes through (**X**) but only some probabilistic function of it (**Y**). Thus, it is a Markov model with the addition of *emission probabilities*:

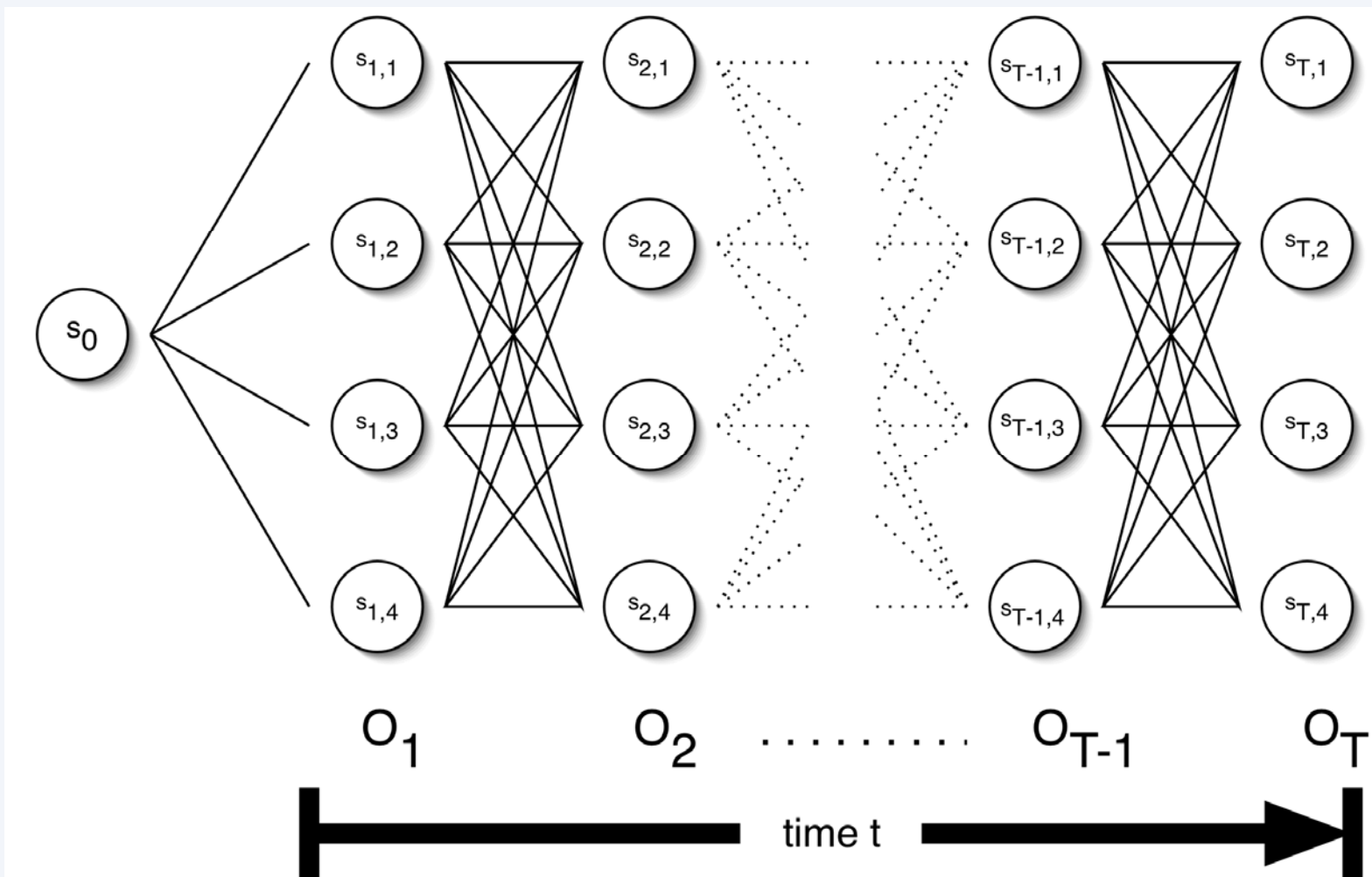
Standard tool is the hidden Markov Model (HMM).



Hidden = Latent

$$P(\mathbf{X}, \mathbf{Y}) = \prod_i P(\mathbf{X}_i | \mathbf{Y}_i) P(\mathbf{Y}_i | \mathbf{Y}_{i-1})$$

The Trellis



The Three Problems in HMMs

- **Likelihood/Evaluation:** Given a series of observations \mathbf{y} and a model $\lambda = \{\mathbf{A}, \mathbf{B}, \mathbf{q}\}$, compute the likelihood $p(\mathbf{y} | \lambda)$
 - >> Forward Algorithm
- **Inference/Decoding:** Given a series of observations \mathbf{y} and a model $\lambda = \{\mathbf{A}, \mathbf{B}, \mathbf{q}\}$, compute the most likely sequence of hidden states \mathbf{x}
 - >> Viterbi Algorithm (like forward algorithm but just do max instead of sum)
- **Learning:** Given a series of observations, learn the best model λ
 - >> Forward-Backward Algorithm (Baum Welch)
 - (Iterative algorithm to re-estimate parameters, like EM)

Likelihood in HMMs

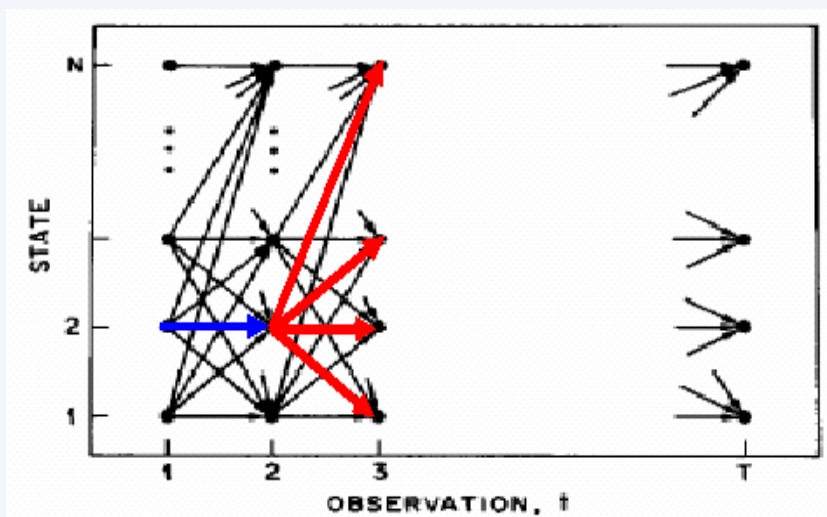
- Given a model $\lambda = \{\mathbf{A}, \mathbf{B}, \mathbf{q}\}$, we can compute the likelihood by

$$\begin{aligned} P(\mathbf{y}) &= p(\mathbf{y} | \lambda) \\ &= \sum p(\mathbf{x}) p(\mathbf{y} | \mathbf{x}) \\ &= q(x_1) \prod A(x_{t+1} | x_t) \prod B(y_t | x_t) \end{aligned}$$

- But ... this computation complexity is $O(N^T)$, when $|\mathbf{x}_i| = N \rightarrow$ impossible in practice

Forward-Backward algorithm

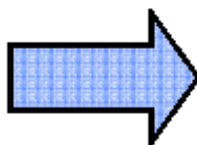
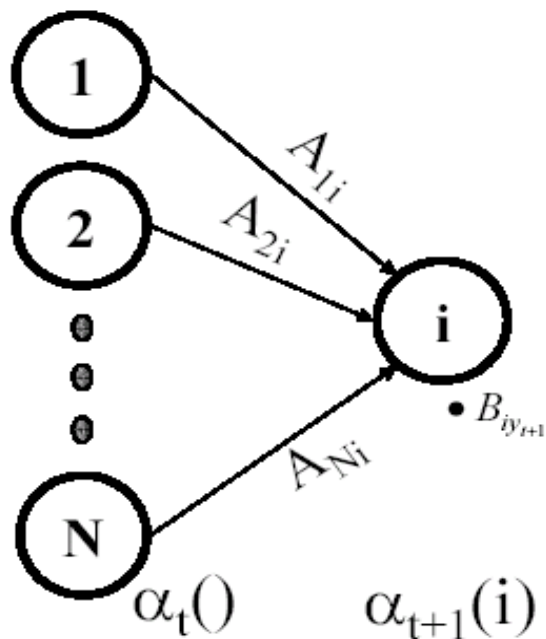
- To compute likelihood:
 - Need to enumerate over all paths in the lattice (all possible instantiations of $X_1 \dots X_T$). But ... some starting subpath (blue) is common to many continuing paths (blue+red)



The idea:
Use **dynamic programming**, calculate a path in terms of shorter sub-paths

Forward-Backward algorithm (cont'd)

- We build a matrix of the probability of being at time t at state i : $\alpha_t(i) = P(x_t=i, y_1 y_2 \dots y_t)$. This is a function of the previous column (forward procedure):



$$\alpha_1(i) = q_i B_{iy_1}$$

$$\alpha_{t+1}(i) = B_{iy_{t+1}} \sum_{j=1}^N \alpha_t(j) A_{ji}$$

$$P(Y) = \sum_{i=1}^N \alpha_T(i)$$

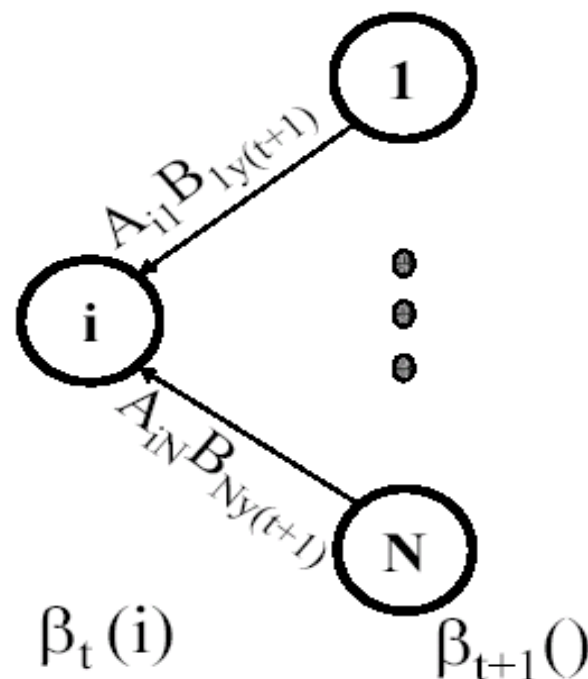
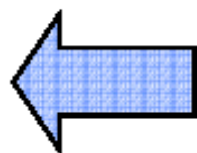
Forward-Backward algorithm (cont'd)

We can similarly define a backwards procedure for filling the matrix $\beta_t(i) = P(y_{t+1} \dots y_T | x_t = i)$

$$\beta_T(i) = 1$$

$$\beta_t(i) = \sum_{j=1}^n A_{ij} B_{jy_{t+1}} \beta_{t+1}(j)$$

$$P(Y) = \sum_{i=1}^N q_i B_{iy_1} \beta_1(i)$$



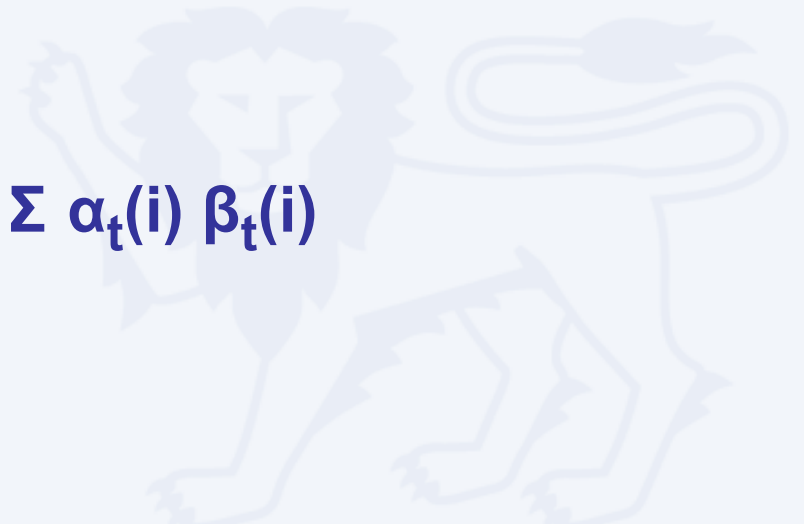
Combine both ...

- Combine both processes to arrive at likelihood:

$$\begin{aligned} P(y, x_t=i) &= P(x_t=i, y_1 y_2 \dots y_t) * P(y_{t+1} \dots y_T | x_t=i) \\ &= \alpha_t(i) \beta_t(i) \end{aligned}$$

- And then we get:

$$P(y) = \sum P(y, x_t=i) = \sum \alpha_t(i) \beta_t(i)$$

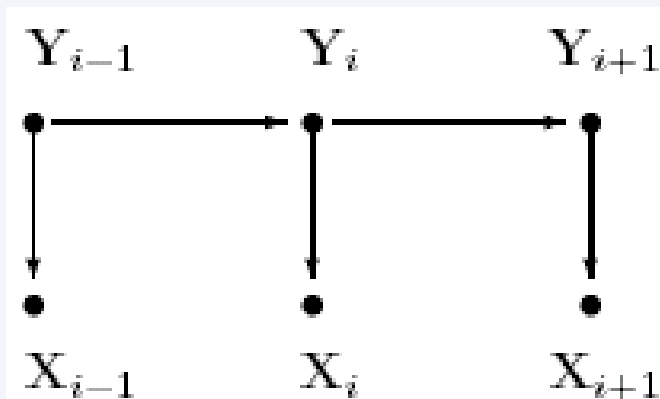


HMM Summary

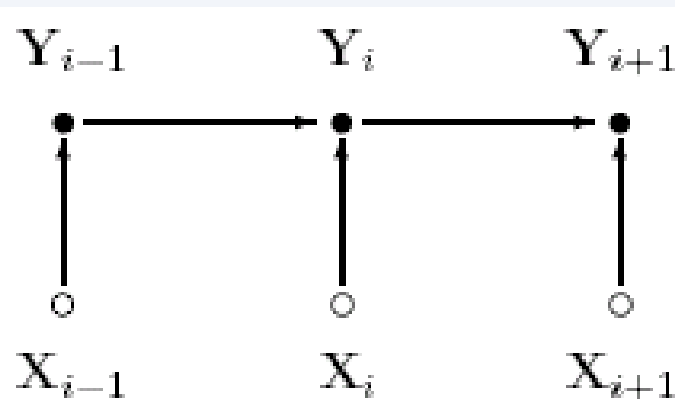
- **Advantages:**
 - Estimation very easy
 - Closed form solution
 - The parameters can be estimated with relatively high confidence from small samples
- **But:**
 - The model represents all possible (x,y) sequences and defines **joint probability** over all possible observation and label sequences
 - Need to enumerate all possible observation sequences
 - Impossible to represent multiple interacting features
 - Difficult to model long-range dependencies of the observations
 - Very strict independence assumptions on the observations

Discriminative Probabilistic Models

Generative



Discriminative



“Solve the problem you need to solve”:

The traditional approach inappropriately uses a generative *joint* model in order to solve a *conditional* problem in which the observations are given.

To classify we actually need $p(y|x)$ – there’s no need to implicitly approximate $p(x,y)$.

Discriminative / Conditional Models

- Conditional probability **$P(\text{label seq } y \mid \text{observed seq } x)$** rather than joint probability $P(y, x)$
 - Specify the probability of possible label sequences given an observation sequence
- Allow arbitrary, non-independent features on the observation sequence X
- The probability of a transition between labels may depend on **past** and **future** observations
 - Relax strong independence assumptions in HMM

Maximum Entropy Markov Models (MEMMs)

a.k.a. Conditional Markov Models CMMs

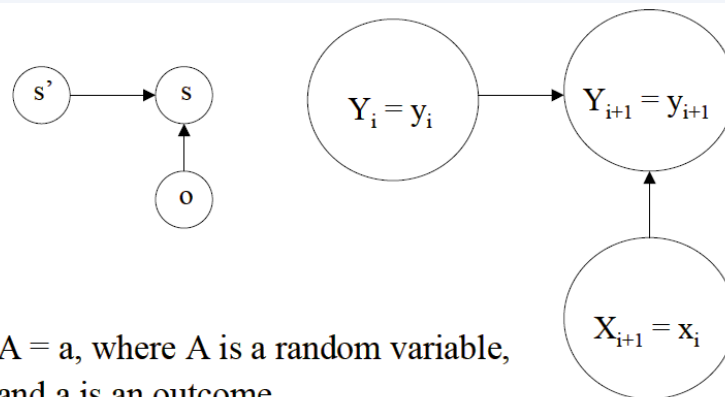
- Models probability of a state **given** an observation and just the previous state

- Conditional probs are represented as exponential models based on arbitrary observation features

- Given training set **X** with label sequence **Y**:

- Train a model θ that maximizes $P(Y|X, \theta)$

- For a new data sequence x , predict label y that maximizes $P(y|x, \theta)$



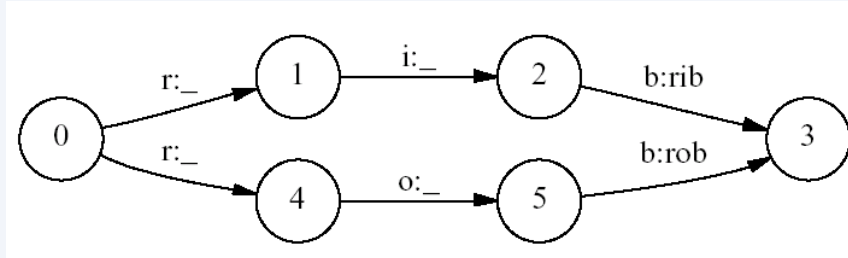
$A = a$, where A is a random variable, and a is an outcome

$$P(y' | y, x) = \frac{1}{Z(y, x)} \exp \left(\sum_k \underbrace{\lambda_k}_{\text{weight}} \underbrace{f_k(x, y, y')}_{\text{feature}} \right)$$

Per state normalization: all prob mass that arrives is distributed among its successor states

The Label Bias Problem

- In MaxEnt's formulation, the prob mass that arrives at the state must be distributed among the possible successor states



- If one of transitions leaving state 0 occurs more frequently in training, its transition prob is greater, irrespective of the observation sequence
 - Especially in cases where there are few outgoing transitions (as in states 1, 2, 4 and 5).
- In the example, say that 'rib' is slightly more common than 'rob' in the training data. Then in the test data, if 'rob' occurs it will be classified as 'rib' as the the transition to 1 is more likely than to 4; the observation of the 'o' is effectively ignored as that it is only observed later at state 1.

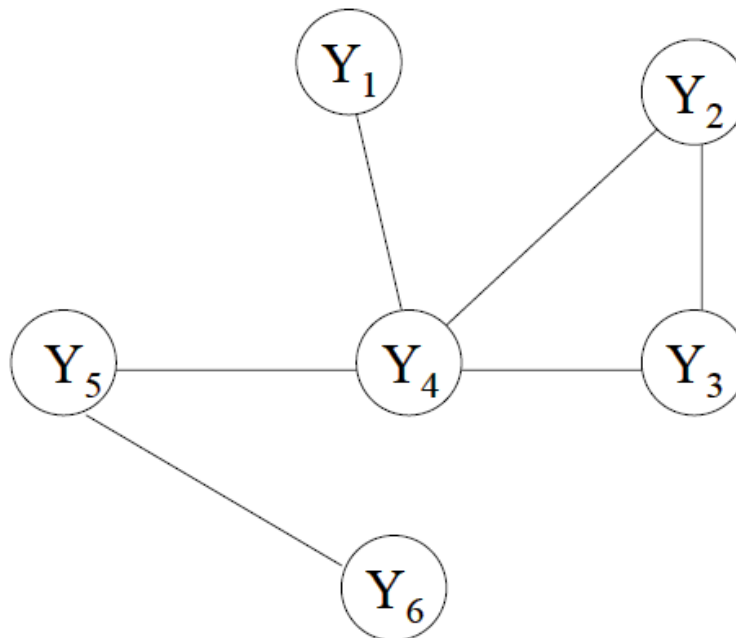
Conditional Random Fields (CRFs)

- **CRFs have all the advantages of MEMMs without label bias problem**
 - MEMM uses **per-state exponential** model for the conditional probabilities of next states given the current state
 - CRF has a **single exponential** model for the joint probability of the entire sequence of labels given the observation sequence
 - This difference means that some transitions have more influence than others depending on the corresponding observations in our previous example
- **Undirected acyclic graph**

Random Fields – Undirected Graphical Models

Let $G = (Y, E)$ be a graph where each vertex Y_v is a random variable
Suppose $P(Y_v | \text{all other } Y) = P(Y_v | \text{neighbors}(Y_v))$ then Y is a
random field

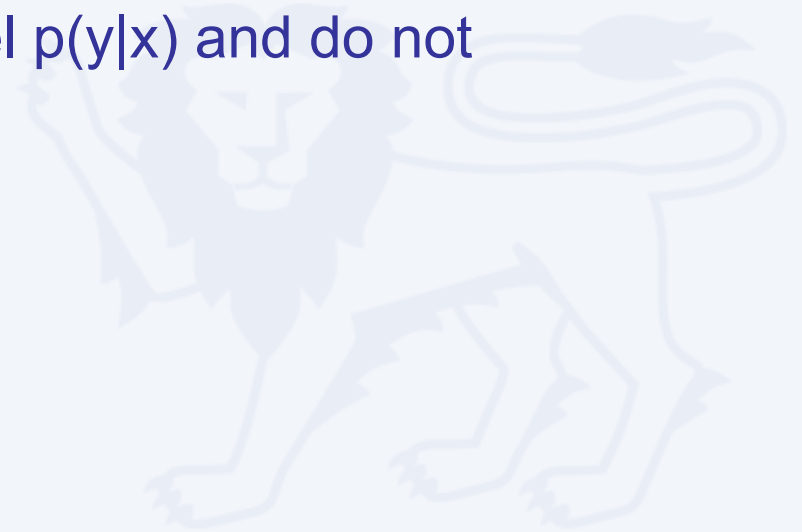
Example:



- $P(Y_5 | \text{all other } Y) = P(Y_5 | Y_4, Y_6)$

Conditional Random Field: Definition

- X – random variable over data sequences
- Y – random variable over label sequences
- Y_i is assumed to range over a finite label alphabet A
- **Discriminative approach:**
 - We construct a conditional model $p(y|x)$ and do not explicitly model marginal $p(x)$



CRF Distribution Function

$$p_{\theta}(y | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left(\sum_{e \in E, k} \lambda_k f_k(e, y |_e, \mathbf{x}) + \sum_{v \in V, k} \mu_k g_k(v, y |_v, \mathbf{x}) \right)$$

Where :

V = Set of random variables (observed and hidden)

f_k and **g_k** = features

g_k = State feature

f_k = Edge feature

$\theta = (\lambda_1, \lambda_2, \dots, \lambda_n; \mu_1, \mu_2, \dots, \mu_n); \lambda_k$ and μ_k

are parameters to be estimated

y|_e = Set of components of y defined by edge e

y|_v = Set of components of y defined by vertex v

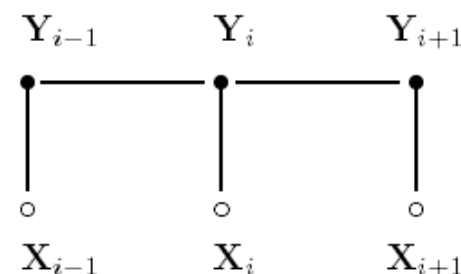
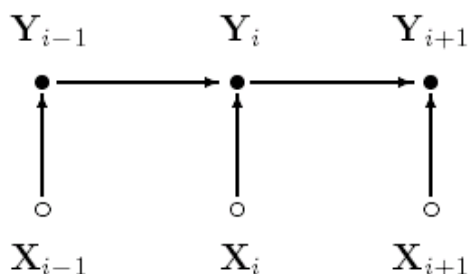
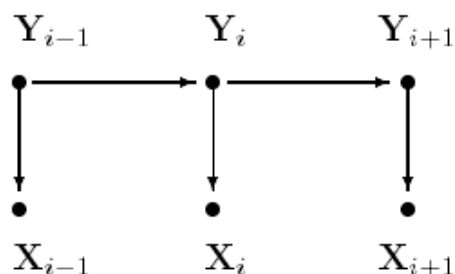
CRF on the linear chain graph

- We will handle the case when G is a simple chain: $G = (V = \{1, \dots, m\}, E = \{ (l, i+1) \})$

HMM (Generative)

MEMM (Discriminative)

CRF

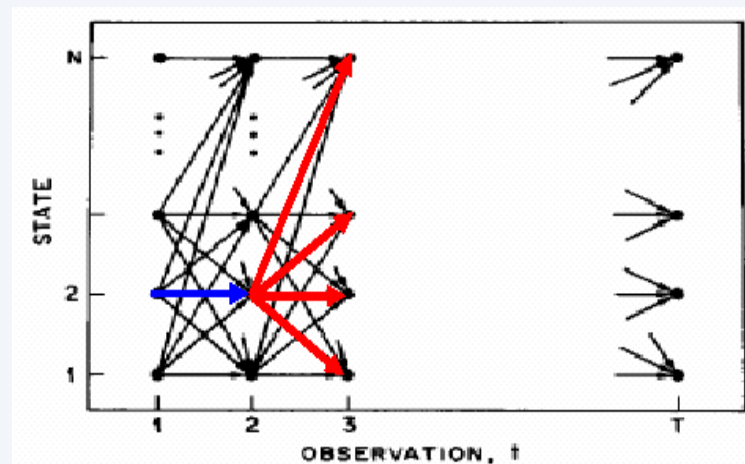


CRF – the Learning Problem

- **Assumption: the *features* f_k and g_k are given and fixed.**
 - For example, a boolean feature g_k is TRUE if the word X_i is upper case and the label Y_i is a “noun”.
- **The learning problem**
 - We need to determine the parameters $\Theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$ from training data $D = \{(x(i), y(i))\}$ with empirical distribution $p\tilde{(x, y)}$.

Parameter Estimation for CRFs

- The parameter vector Θ that maximizes the log-likelihood is found using an iterative scaling algorithm.
- We define standard HMM-like **forward** and **backward vectors** α and β , which allow polynomial time calculations.
- However as the normalization is conditioned over the entire CRF (not over single vertices), it is expensive to compute \rightarrow **slow training time**



Experiment Validation of the models

- Part-of-speech (POS) tagging experiments

<i>model</i>	<i>error</i>	<i>oov error</i>
HMM	5.69%	45.99%
MEMM	6.37%	54.61%
CRF	5.55%	48.05%

MEMM ⁺	4.81%	26.99%
CRF ⁺	4.27%	23.76%

⁺Using spelling features

Summary

- **HMM:**
 - Basic sequence labeling framework where labels are considered hidden and generate the observed words, and are just dependent on the previous state (Markovian assumption)
 - Fast algorithms for three classic problems
- **CRF:**
 - Discriminatively trained models $P(y|x)$ as compared to modeling joint $P(x,y)$ probability
 - Allows combination of arbitrary and overlapping observation features from both the past and future
 - main current limitation is the slow convergence of the training algorithm relative to MEMMs or HMMs, for which training is efficient.

Hands on with CRF++

Reference string labeling over the Cora
dataset

Cora Dataset

- 200 reference strings taken from articles in Computer Science
- Labeled with fields in XML style

Isaac G. Councill, C. Lee Giles, Min-Yen Kan. (2008) *ParsCit: An open-source CRF reference string parsing package*. To appear in the proceedings of the Language Resources and Evaluation Conference (LREC 08), Marrakesh, Morocco, May.

CRF++

- Implementation of CRFs in C++
- Built with multithreading
- Generates individual binary feature functions from feature templates (each which describe a class of features)
- Uses `conlleval.pl` script to assess performance

Six Steps

1. Convert data to CRF++ format
2. Create basic data and template file
3. Create basic word features
4. Integrate lexicon features
5. Error analysis – Inspect results more closely
6. Create punctuation, numeric features
7. Create global features

