# Probabilistic Graphical Models : Hidden Markov Models and Conditional Random Fields

आई आई टी हैदराबाद
**IIT Hyderabad**

# Probabilistic Graphical Models

- Used in wide range of applications : natural language processing (POS tagging, named entity recognition, automatic speech recognition)
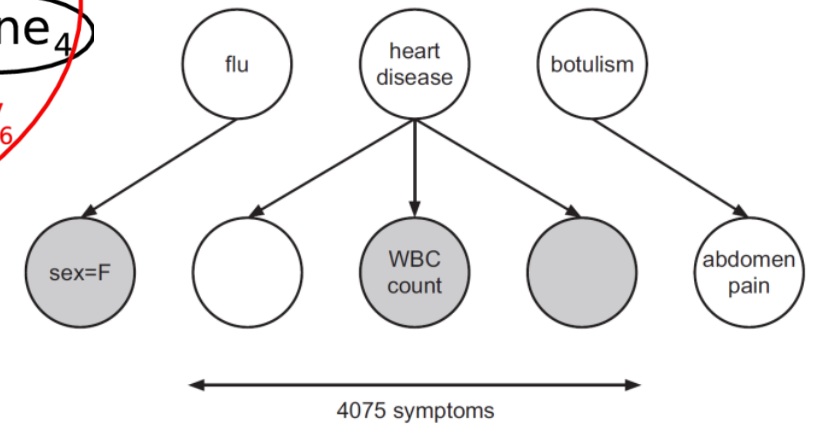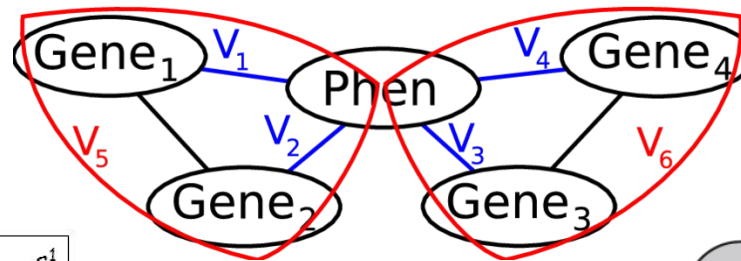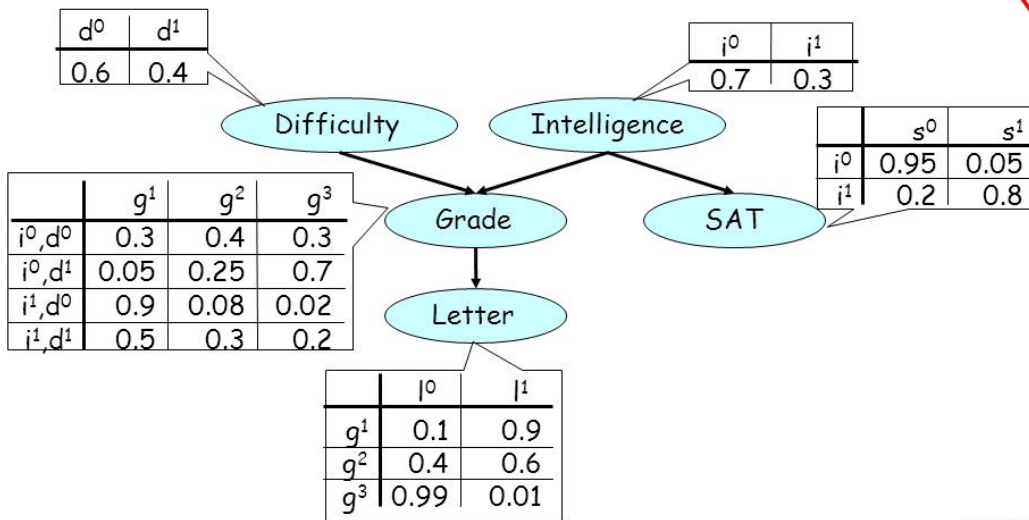  - They provide a simple way to visualize the structure of a probabilistic model and can be used to design and motivate new models.
  - Insights into the properties of the model, including conditional independence properties, can be obtained by inspection of the graph.
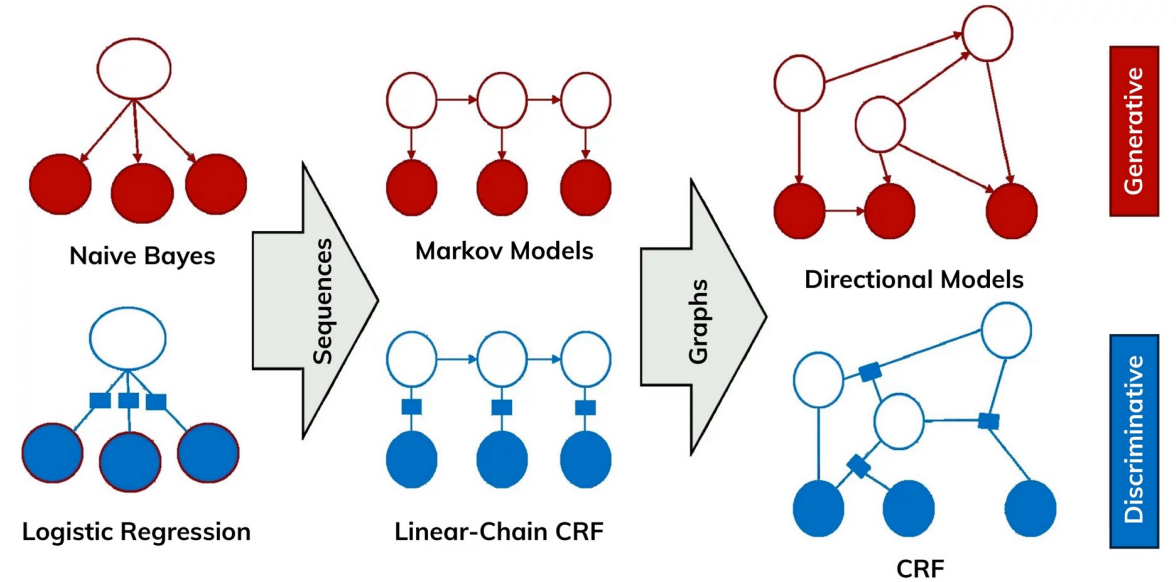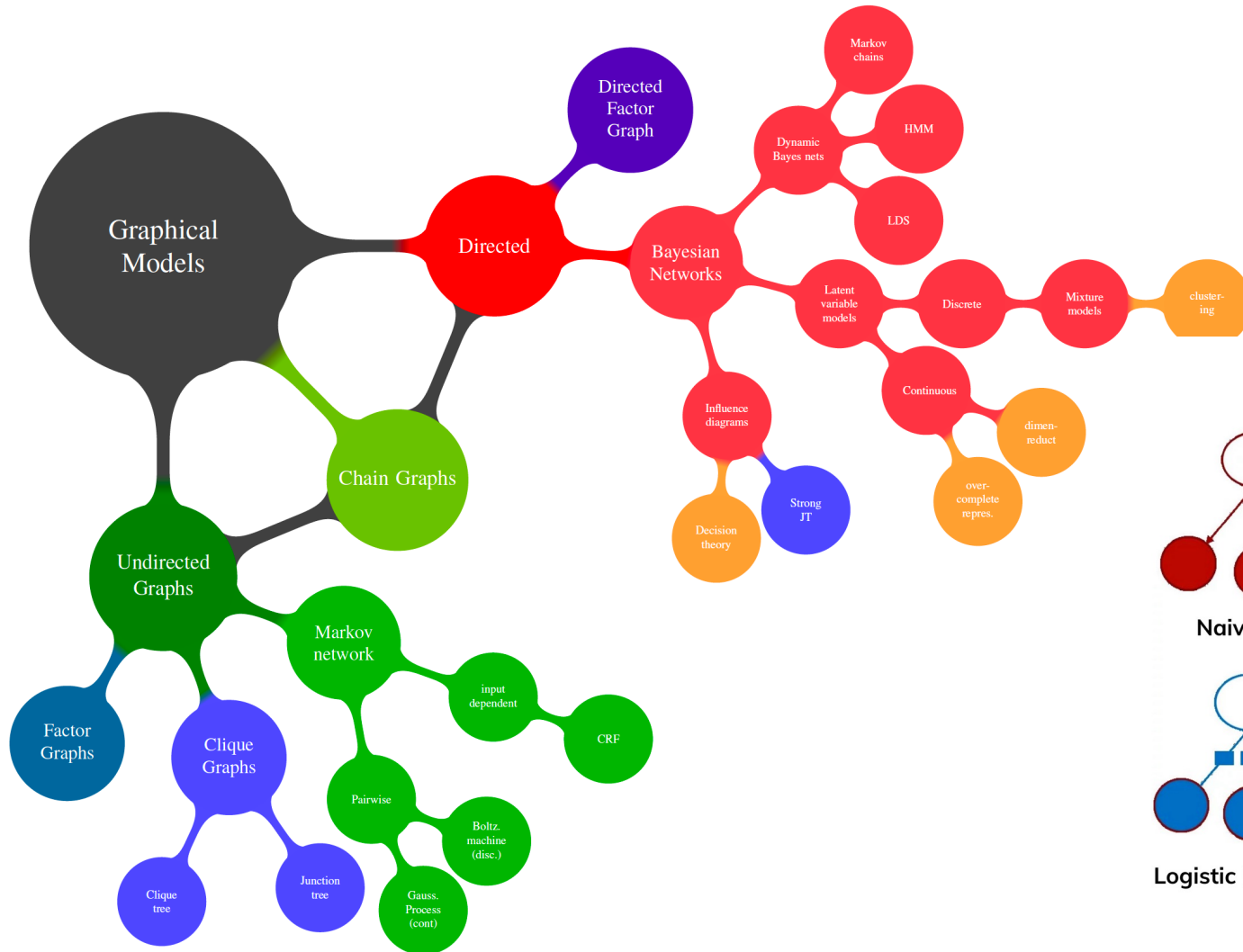
# Probabilistic Graphical Models

- A graph comprises nodes (also called vertices) connected by links (also known as edges or arcs). In a probabilistic graphical model, each node represents a random variable (or group of random variables), and the links express probabilistic relationships between these variables.

- Bayesian networks, also known as directed graphical models

- Markov random fields, also known as undirected graphical models



The Student Network

| | g¹ | g² | g³ |
|---|---|---|---|
| i⁰,d⁰ | 0.3 | 0.4 | 0.3 |
| i⁰,d¹ | 0.05 | 0.25 | 0.7 |
| i¹,d⁰ | 0.9 | 0.08 | 0.02 |
| i¹,d¹ | 0.5 | 0.3 | 0.2 |

| d⁰ | d¹ |
|---|---|
| 0.6 | 0.4 |

| i⁰ | i¹ |
|---|---|
| 0.7 | 0.3 |

| | s⁰ | s¹ |
|---|---|---|
| i⁰ | 0.95 | 0.05 |
| i¹ | 0.2 | 0.8 |

| | l⁰ | l¹ |
|---|---|---|
| g¹ | 0.1 | 0.9 |
| g² | 0.4 | 0.6 |
| g³ | 0.99 | 0.01 |

Daphne Koller

# Probabilistic Graphical Models



Adapted from C.Sutton, A.McCallum, "An introduction to Conditional Random Fields"

# Probabilistic Graphical Models

- Independently specifying all the entries of a table p(x1; : : : ; xN) over binary variables xi takes O(2^N) space

- Structure is also important for computational tractability of inferring quantities of interest.

- Given a distribution on N binary variables, p(x1; : : : ; xN), computing a marginal such as p(x1) requires summing over the 2^(N-1) states of the other variables.

- Belief networks (also called Bayes' networks or Bayesian belief networks) are a way to depict the independence

- assumptions made in a distribution

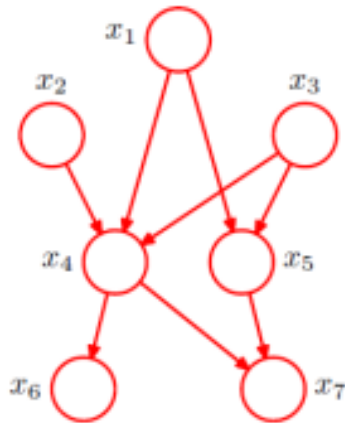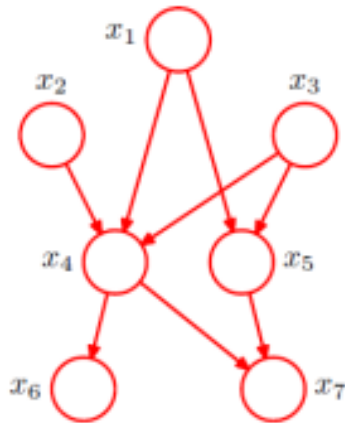$$p(a, b, c) = p(c|a, b)p(b|a)p(a).$$

# Probabilistic Graphical Models

- Independently specifying all the entries of a table p(x1; : : : ; xN) over binary variables xi takes O(2^N) space

- Structure is also important for computational tractability of inferring quantities of interest.

- Given a distribution on N binary variables, p(x1; : : : ; xN), computing a marginal such as p(x1) requires summing over the 2^(N-1) states of the other variables.

- Belief networks (also called Bayes' networks or Bayesian belief networks) are a way to depict the independence

- assumptions made in a distribution

# Probabilistic Graphical Models

- Independently specifying all the entries of a table p(x1; …; xN) over binary variables xi takes O(2^N) space

- Structure is also important for computational tractability of inferring quantities of interest.

- Given a distribution on N binary variables, p(x1; … ; xN), computing a marginal such as p(x1) requires summing over the 2^(N-1) states of the other variables.

- Belief networks (also called Bayes' networks or Bayesian belief networks) are a way to depict the independence

- assumptions made in a distribution



$$p(x_1)p(x_2)p(x_3)p(x_4|x_1,x_2,x_3)p(x_5|x_1,x_3)p(x_6|x_4)p(x_7|x_4,x_5).$$

$$p(\mathbf{x}) = \prod_{k=1}^{K} p(x_k|\mathrm{pa}_k)$$

# Probabilistic Graphical Models

- **Plate Notation**
- Bayesian Linear Regression

$$p(\mathcal{D}|\mu) = \prod_{n=1}^{N} p(x_n|\mu).$$

$$p(\mathbf{t}, \mathbf{w}|\mathbf{x}, \alpha, \sigma^2) = p(\mathbf{w}|\alpha) \prod_{n=1}^{N} p(t_n|\mathbf{w}, x_n, \sigma^2).$$

# Probabilistic Graphical Models

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph

$$p(a|b,c) = p(a|c).$$

$$p(a,b|c) = p(a|b,c)p(b|c)$$
$$= p(a|c)p(b|c).$$

$$p(a,b,c) = p(a|c)p(b|c)p(c).$$

- **Conditional Independence**

- *a* is conditionally independent of *b* given *c*.

  joint distribution of *a* and *b* factorizes into the product of the marginal distribution of *a* and the marginal distribution of *b* (again both conditioned on *c*).

$$a \perp\!\!\!\perp b \mid c$$

# Probabilistic Graphical Models

**D-separation** : An important and elegant
feature of graphical models is that conditional
independence properties of the joint
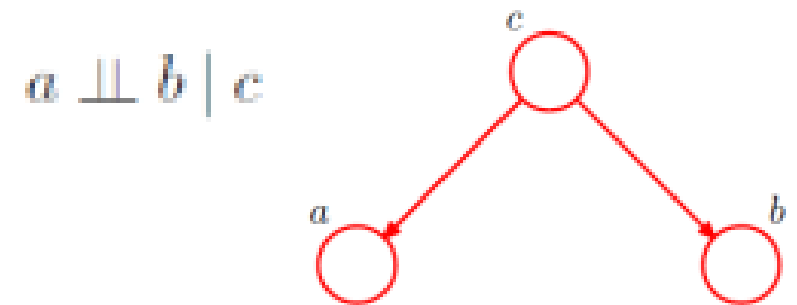distribution can be read directly from the graph

$$p(a, b, c) = p(a|c)p(b|c)p(c).$$

$$
\begin{aligned}
p(a, b|c) &= \frac{p(a, b, c)}{p(c)} \\
&= p(a|c)p(b|c)
\end{aligned}
$$

- **Conditional Independence**



$$a \perp\!\!\!\perp b \mid c.$$

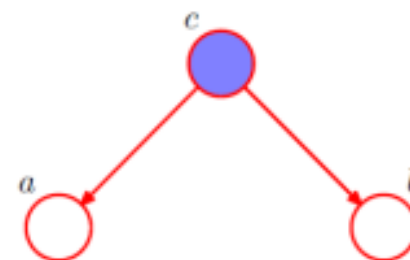Is a and b unconditionally independent ?

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph

- **Conditional Independence**

$$p(a, b, c) = p(a|c)p(b|c)p(c).$$

$$\begin{aligned} p(a, b|c) &= \frac{p(a, b, c)}{p(c)} \\ &= p(a|c)p(b|c) \end{aligned}$$

$$a \perp\!\!\!\perp b \mid c.$$

$$p(a, b) = \sum_c p(a|c)p(b|c)p(c).$$

$$a \not\!\perp\!\!\!\perp b \mid \emptyset$$

# Probabilistic Graphical Models

- **Conditional Independence**

**D-separation** : An important and elegant
feature of graphical models is that conditional
independence properties of the joint
distribution can be read directly from the graph

# Probabilistic Graphical Models

• **Conditional Independence**

**D-separation** : An important and elegant
feature of graphical models is that conditional
independence properties of the joint
distribution can be read directly from the graph
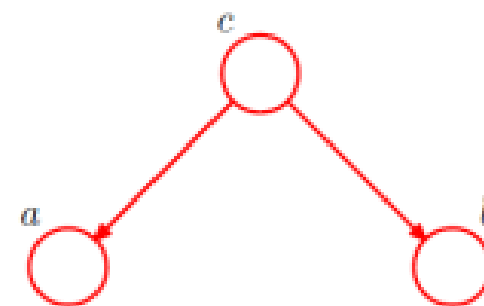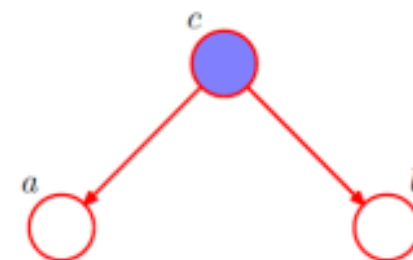


$$p(a, b, c) = p(a)p(c|a)p(b|c).$$

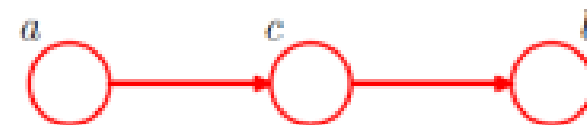Is a independent of b ?
Is a independent of b conditioned on c ?

# Probabilistic Graphical Models
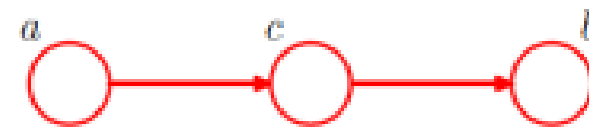
- **Conditional Independence**

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph



$$p(a, b, c) = p(a)p(c|a)p(b|c).$$

$$p(a, b) = p(a) \sum_c p(c|a)p(b|c) = p(a)p(b|a).$$

$$a \not\!\perp\!\!\!\perp b \mid \emptyset$$
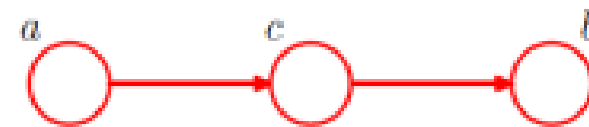
# Probabilistic Graphical Models

- **Conditional Independence**

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph

$$p(a, b, c) = p(a)p(c|a)p(b|c).$$

$$p(a, b) = p(a) \sum_c p(c|a)p(b|c) = p(a)p(b|a).$$

$$a \not\perp\!\!\!\perp b \mid \emptyset$$

$$p(a, b|c) = \frac{p(a, b, c)}{p(c)}$$

$$= \frac{p(a)p(c|a)p(b|c)}{p(c)}$$

$$= p(a|c)p(b|c)$$

$$a \perp\!\!\!\perp b \mid c.$$
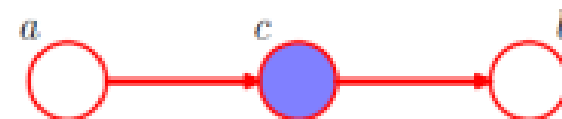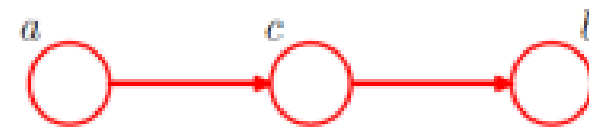
- **Conditional Independence**

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph

The node *c* is said to be *head-to-tail* with respect to the path from node *a* to node *b*. Such a path connects nodes *a* and *b* and c blocks them renders them dependent.

$$p(a, b) = p(a) \sum_c p(c|a)p(b|c) = p(a)p(b|a).$$

$$a \not\!\perp\!\!\!\perp b \mid \emptyset$$

$$p(a, b|c) = \frac{p(a, b, c)}{p(c)}$$

$$= \frac{p(a)p(c|a)p(b|c)}{p(c)}$$

$$= p(a|c)p(b|c)$$

$$a \perp\!\!\!\perp b \mid c.$$

# Probabilistic Graphical Models
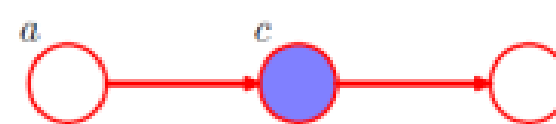
- **Conditional Independence**

**D-separation** : An important and elegant
feature of graphical models is that conditional
independence properties of the joint
distribution can be read directly from the graph

Is a independent of b ?

Is a independent of b conditioned on c ?
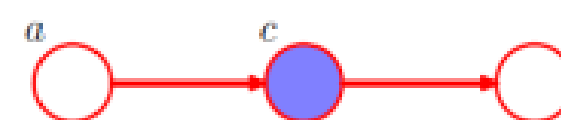
# Probabilistic Graphical Models

- **Conditional Independence**

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph

$$p(a, b, c) = p(a)p(b)p(c|a, b).$$

$$p(a, b) = p(a)p(b)$$

$$a \perp\!\!\!\perp b \mid \emptyset.$$

Is a independent of b conditioned on c ?

# Probabilistic Graphical Models

- **Conditional Independence**

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph

$$p(a, b, c) = p(a)p(b)p(c|a, b).$$

$$p(a, b) = p(a)p(b)$$

$$a \perp\!\!\!\perp b \mid \emptyset.$$

$$p(a, b|c) = \frac{p(a, b, c)}{p(c)}$$
$$= \frac{p(a)p(b)p(c|a, b)}{p(c)}$$

$$a \not\perp\!\!\!\perp b \mid c.$$

- **Conditional Independence**

**D-separation** : An important and elegant feature of graphical models is that conditional independence properties of the joint distribution can be read directly from the graph
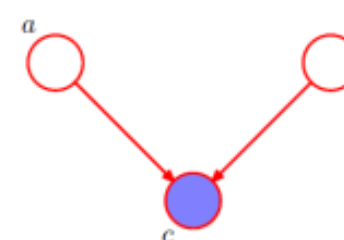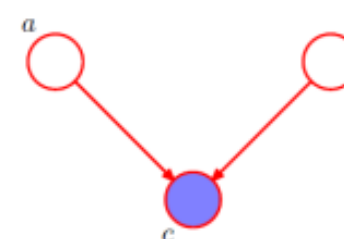
node $c$ is *head-to-head* with respect to the path from $a$ to $b$ because it connects to the heads of the two arrows. When node $c$ is unobserved, it 'blocks' the path, and the variables $a$ and $b$ are independent. However, conditioning on $c$ 'unblocks' the path and renders $a$ and $b$ dependent.

$$p(a, b) = p(a)p(b)$$

$$a \perp\!\!\!\perp b \mid \emptyset.$$

$$
\begin{aligned}
p(a, b|c) &= \frac{p(a, b, c)}{p(c)} \\
&= \frac{p(a)p(b)p(c|a, b)}{p(c)}
\end{aligned}
$$

$$a \not\perp\!\!\!\perp b \mid c.$$



explaining away'.

# Sequential Data

- Sequential data : rainfall measurements on successive days at a particular location, or the daily values of a currency exchange rate (time series data) , sequence of nucleotide base pairs along a strand of DNA or the sequence of characters in an English sentence



*In March 2005, the New York Times acquired About, Inc.*

| IN | NNP | CD | DT | NNP | NNP | NNP | VBD | NNP | NNP |

NP     NP     VP     NP

TEMP     ORG     ORG

# Markov Model

- binary variable denoting whether on a particular day it rained or not.

$$p(\mathbf{x}_1, \ldots, \mathbf{x}_N) = \prod_{n=1}^{N} p(\mathbf{x}_n | \mathbf{x}_1, \ldots, \mathbf{x}_{n-1}).$$

- Markov Model (first Order) : conditional distributions on the right-hand side is independent of all previous observations except the most recent

$$p(\mathbf{x}_n | \mathbf{x}_1, \ldots, \mathbf{x}_{n-1}) = p(\mathbf{x}_n | \mathbf{x}_{n-1})$$

$$p(\mathbf{x}_1, \ldots, \mathbf{x}_N) = p(\mathbf{x}_1) \prod_{n=2}^{N} p(\mathbf{x}_n | \mathbf{x}_{n-1}).$$

$M^{\text{th}}$ order Markov chain,
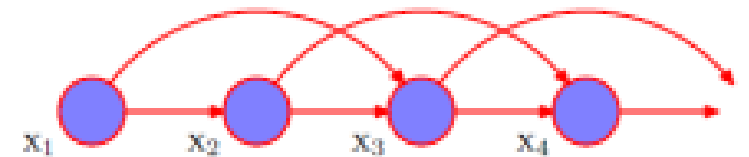$$p(\mathbf{x}_n | \mathbf{x}_{n-M}, \ldots, \mathbf{x}_{n-1}).$$

$$p(\mathbf{x}_1, \ldots, \mathbf{x}_N) = p(\mathbf{x}_1) p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{n=3}^{N} p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2}).$$



$x_1 \quad x_2 \quad x_3 \quad x_4$

$x_1 \quad x_2 \quad x_3 \quad x_4$

$K(K-1)$ parameters.

$x_1 \quad x_2 \quad x_3 \quad x_4$

$K^{M-1}(K-1)$ parameters.

# Probabilistic Graphical Models

- Inference in Graphical Models
  - Message passing algorithms
  - Max-sum algorithm                                  Sum product algorithm

$$\mathbf{x}^{\max} = \arg\max_{\mathbf{x}} p(\mathbf{x})$$

$$\max_{\mathbf{x}} p(\mathbf{x}) = \max_{x_1} \dots \max_{x_M} p(\mathbf{x})$$

$$\max_{\mathbf{x}} p(\mathbf{x}) = \frac{1}{Z} \max_{x_1} \cdots \max_{x_N} \left[ \psi_{1,2}(x_1, x_2) \cdots \psi_{N-1,N}(x_{N-1}, x_N) \right]$$

$$= \frac{1}{Z} \max_{x_1} \left[ \psi_{1,2}(x_1, x_2) \left[ \cdots \max_{x_N} \psi_{N-1,N}(x_{N-1}, x_N) \right] \right].$$

$$p(x_n) = \sum_{x_1} \cdots \sum_{x_{n-1}} \sum_{x_{n+1}} \cdots \sum_{x_N} p(\mathbf{x}).$$

$$p(\mathbf{x}) = \prod_{s \in \mathrm{ne}(x)} F_s(x, X_s)$$

$$p(x) = \prod_{s \in \mathrm{ne}(x)} \left[ \sum_{X_s} F_s(x, X_s) \right]$$

# Hidden Markov Model

- HMM is widely used in speech. recognition (Jelinek, 1997; Rabiner and Juang, 1993), natural language modelling (Manning and Sch¨utze, 1999), on-line handwriting recognition (Nag et al., 1986), and for the analysis of biological sequences such as proteins and DNA

- Standard classification problem assumes individual cases are disconnected and independent (i.i.d.: independently and identically distributed).

- Each token in a sequence is assigned a label. Labels of tokens are dependent on the labels of other tokens in the sequence, particularly their neighbors (not i.i.d).

- A given sentence, *"Time flies like an arrow"*
- Represent the input sentence with a token vector $\boldsymbol{x}$

| $t$ | 1 | 2 | 3 | 4 | 5 |
|-----|------|-------|------|-----|-------|
| $\boldsymbol{x}$ | *Time* | *flies* | *like* | *an* | *arrow* | (T = 5) |
| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |

(Bold italic)
(*NOTE: This does not present a feature vector*)

- Predict part-of-speech (a vector $\boldsymbol{y}$) tags for the tokens $\boldsymbol{x}$

| $t$ | 1 | 2 | 3 | 4 | 5 |
|-----|------|-------|------|-----|-------|
| $\boldsymbol{x}$ | *Time* | *flies* | *like* | *an* | *arrow* |
| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
| $\boldsymbol{y}$ | *NN* | *VBZ* | *IN* | *DT* | *NN* |
| | $y_1$ | $y_2$ | $y_3$ | $y_4$ | $y_5$ |

Predict

- *Modeling*: how to build (assume) $P(y|x)$
  - Hidden Markov Model (HMM), Structured Perceptron, Conditional Random Fields (CRFs), etc

- *Training*: how to determine unknown parameters in the model so that they fit to a training data
  - Maximum Likelihood (ML), Maximum a Posteriori (MAP), etc
  - Gradient-based method, Stochastic Gradient Descent (SGD), etc

- *Inference*: how to compute $\mathrm{argmax}\, P(y|x)$ efficiently
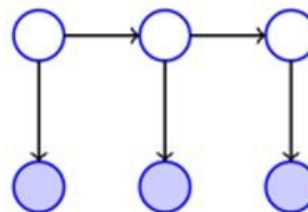  - Viterbi algorithm

# Probabilistic Sequence Models

- Probabilistic sequence models allow integrating uncertainty over multiple, interdependent classifications and collectively determine the most likely global assignment.

- Two standard models
  - Generative Model : Hidden Markov Model (HMM)
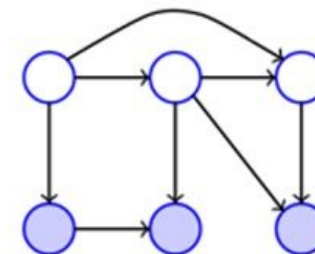  - Discriminative Model : Conditional Random Field (CRF)
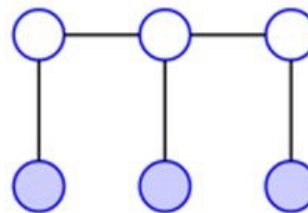


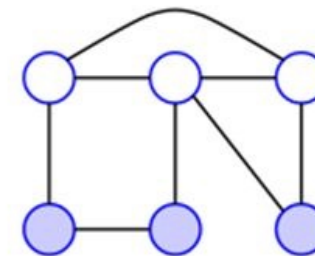Generative-Discriminative Pairs

Naïve Bayes    Hidden Markov Model    Generative Directed Model

Logistic Regression    Linear Chain CRF    Conditional Random Field

CRF

43

# Hidden Markov Model

- $x$: the sequence of tokens (words)
- $y$: the sequence of POS tags
- Bayes' theorem:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

- Bayesian inference: decompose $P(y|x)$ into two factors, $P(x|y)$ and $P(y)$, which might be easier to model

$$\hat{y} = \underset{y}{\text{argmax}}\, P(y|x) = \underset{y}{\text{argmax}} \frac{P(x|y)P(y)}{P(x)} = \underset{y}{\text{argmax}}\, P(x|y)P(y)$$

Bayes'
theorem

$P(x)$ is the same
for all $y$

# HMM

- Two Markov assumptions to simplify $P(x|y)$ and $P(y)$
  - A word appears depending only on its POS tag
    - Independently of other words around the word
    - Generated by emission probability distribution

$$P(x|y) \approx \prod_{t=1}^{T} P(x_t|y_t)$$

  - A POS tag is dependent only on the previous one
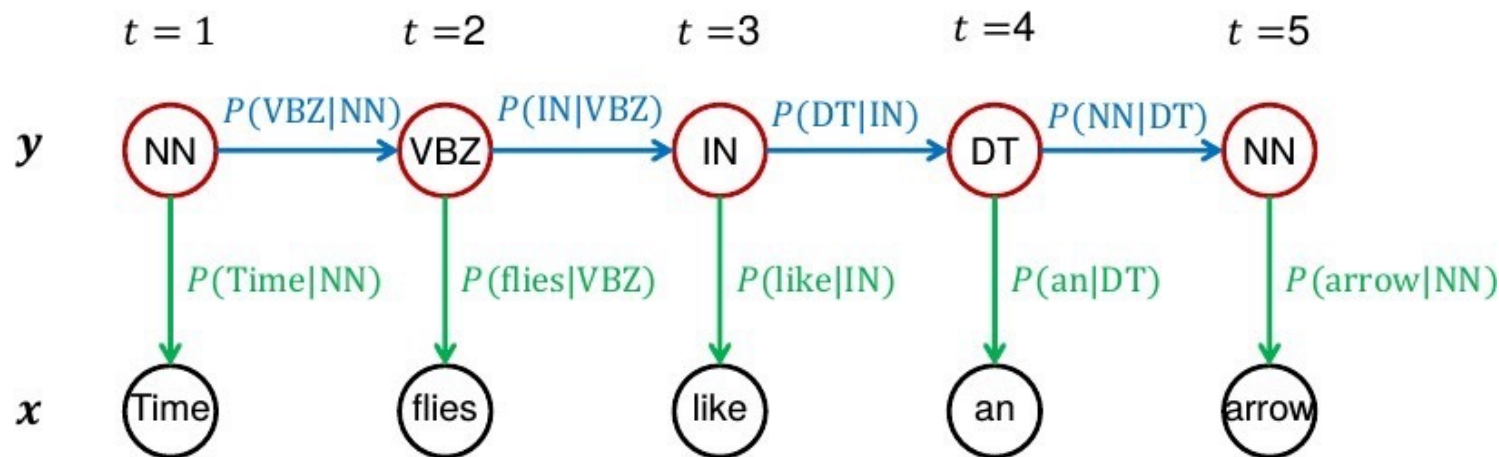    - Rather than the entire tag sequence
    - Generated by transition probability distribution

$$P(y) \approx \prod_{t=1}^{T} P(y_t|y_{t-1})$$

- Then, the most probable tag sequence $\hat{y}$ is computed by,

$$\hat{y} = \underset{y}{\text{argmax}}\, P(y|x) = \underset{y}{\text{argmax}}\, P(x|y)P(y) \approx \underset{y}{\text{argmax}} \prod_{t=1}^{T} P(x_t|y_t)P(y_t|y_{t-1})$$

# POS Tagging



- We can compute $\phi(x, y)$ if we decide an assignment of $y$ for a given input $x$: $\prod_{t=1}^{T} P(x_t|y_t)P(y_t|y_{t-1})$

$$P(x_t|y_t) = \frac{C(x_t, y_t)}{C(y_t)} = \frac{\text{(the number of times where } x_t \text{ is annotated as } y_t)}{\text{(the number of occurrences of tag } y_t)}$$

$$P(y_t|y_{t-1}) = \frac{C(y_t, y_{t-1})}{C(y_{t-1})} = \frac{\text{(the number of occurrences of tag } y_t \text{ followd by } y_{t-1})}{\text{(the number of occurrences of tag } y_{t-1})}$$

# Viterbi Algorithm

- Given the observations, find the best possible tag sequence which generated it.

$$\text{Inference: } \hat{y} = \underset{y}{\text{argmax}} \prod_{t=1}^{T} P(x_t|y_t)P(y_t|y_{t-1})$$

- We cannot enumerate all possible $y$ for an input $x$
  - The number of candidate sequences is $|Y|^T$, where:
    - $|Y|$: the number of POS tags ($|Y| = 36$ for Penn Treebank)
    - $T$: the number of tokens in an input sentence
  - The number of candidates is too huge, $36^6 = 2176782336$, even for the short example sentence!
- Viterbi algorithm
  - An efficient algorithm for finding $\hat{y}$
  - Computational cost: $O(|Y|^2 T)$
  - Dynamic programing

# Conditional Random Field

$$Y = \bar{y}_1^n = y_1...yn.$$

$$X = x_1^n = x_1...x_n$$

HMM

CRF

$$\hat{Y} = \underset{Y}{\operatorname{argmax}}\, p(Y|X)$$

$$= \underset{Y}{\operatorname{argmax}}\, p(X|Y)p(Y)$$

$$= \underset{Y}{\operatorname{argmax}} \prod_i p(x_i|y_i) \prod_i p(y_i|y_{i-1})$$

$$\hat{Y} = \underset{Y \in \mathcal{Y}}{\operatorname{argmax}}\, P(Y|X)$$



Naïve Bayes     Sequences     Markov models     Graphs     Directional Models     Generative

Logistic Regression     Linear-chain CRF     CRF     Discriminative

# Conditional Random Field

- Conditional probability is defined,

$$P(y|x) = \frac{\exp((w \cdot F(x,y))}{\boxed{\sum_y \exp((w \cdot F(x,y))}} \leftarrow$$

Normalized by the sum of exp'd scores of all possible paths in the lattice

- The same inference algorithm (Viterbi)

- Input: sequence of tokens $x = (x_1 \; x_2 \; ... \; x_T)$
- Output: sequence of POS tags $\hat{y} = (\widehat{y_1} \; \widehat{y_2} \; ... \; \widehat{y_T})$
- Mapping to global feature vector: $F(x,y): (x,y) \to \mathcal{R}^m$

$$F(x,y) = \sum_{t=1}^{T} \boxed{\{u(x_t, y_t) \;,\; b(y_{t-1}, y_t)\}}$$

Local feature vector (at $t$):
- Unigram feature vector
- Bigram feature vector

- Each element of feature vector consists of a feature function, e.g.,
  - $u_{109}(x_t, y_t) = \{1 \; (\text{if } x_t = \text{Brown } and \; y_t = \text{Noun}); 0 \; (\text{otherwise})\}$
  - $b_2(y_{t-1}, y_t) = \{1 \; (\text{if } y_{t-1} = \text{Noun } and \; y_t = \text{Verb}); 0 \; (\text{otherwise})\}$

# CRF : Training and Inference

$$\hat{Y} = \underset{Y \in \mathcal{Y}}{\operatorname{argmax}} P(Y|X)$$

$$= \underset{Y \in \mathcal{Y}}{\operatorname{argmax}} \frac{1}{Z(X)} \exp\left(\sum_{k=1}^{K} w_k F_k(X,Y)\right)$$

$$= \underset{Y \in \mathcal{Y}}{\operatorname{argmax}} \exp\left(\sum_{k=1}^{K} w_k \sum_{i=1}^{n} f_k(y_{i-1}, y_i, X, i)\right)$$

$$= \underset{Y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{k=1}^{K} w_k \sum_{i=1}^{n} f_k(y_{i-1}, y_i, X, i)$$

$$= \underset{Y \in \mathcal{Y}}{\operatorname{argmax}} \sum_{i=1}^{n} \sum_{k=1}^{K} w_k f_k(y_{i-1}, y_i, X, i)$$

$$p(Y|X) = \frac{\exp\left(\sum_{k=1}^{K} w_k F_k(X,Y)\right)}{\sum_{Y' \in \mathcal{Y}} \exp\left(\sum_{k=1}^{K} w_k F_k(X,Y')\right)}$$

$$p(Y|X) = \frac{1}{Z(X)} \exp\left(\sum_{k=1}^{K} w_k F_k(X,Y)\right)$$

- Viterbi algorithm : like the HMM, the linearchain CRF depends at each timestep on only one previous output token y[i-1].

# Probabilistic Graphical Models

- References
  - *Pattern Recognition and Machine Learning* by Bishop
  - *Probabilistic Machine Learning* by Kevin Murphy
  - *Speech and Language processing* Jurafsky and Martin
  - McCallum, A.: Efficiently inducing features of conditional random fields. In: *Proc. 19th Conference on Uncertainty in Artificial Intelligence*. (2003)

IIT Hyderabad