

# CS5340

## Uncertainty Modeling in AI

### Lecture 3: Markov Random Fields (Undirected Graphical Models)

Assoc. Prof. Lee Gim Hee

AY 2022/23

Semester 1

# Course Schedule

Week	Date	Topic	Remarks
1	10 Aug	Introduction to probabilistic reasoning	<b>Assignment 0:</b> Python Numpy Tutorial (Ungraded)
2	17 Aug	Bayesian networks (Directed graphical models)	
3	24 Aug	Markov random Fields (Undirected graphical models)	
4	31 Aug	Variable elimination and belief propagation	<b>Assignment 1:</b> Belief propagation and maximal probability (15%)
5	07 Sep	Factor graph and the junction tree algorithm	
6	14 Sep	Parameter learning with complete data	<b>Assignment 1:</b> Due <b>Assignment 2:</b> Junction tree and parameter learning (15%)
-	21 Sep	Recess week	<b>No lecture</b>
7	28 Sep	Mixture models and the EM algorithm	<b>Assignment 2:</b> Due
8	05 Oct	Hidden Markov Models (HMM)	<b>Assignment 3:</b> Hidden Markov model (15%)
9	12 Oct	Monte Carlo inference (Sampling)	
*	15 Oct	Variational inference	Makeup Lecture (Venue TBD) Time: 9.30am – 12.30pm (Saturday)
10	19 Oct	Variational Auto-Encoder and Mixture Density Networks	<b>Assignment 3:</b> Due <b>Assignment 4:</b> MCMC Sampling (15%)
11	26 Oct	No Lecture	I will be traveling
12	02 Nov	Graph-cut and alpha expansion	<b>Assignment 4:</b> Due
13	09 Nov	-	

# Acknowledgements

- A lot of slides and content of this lecture are adopted from:
  1. "Machine learning - a probabilistic approach", Kevin Murphy (Chapter 19)
  2. "Probabilistic graphical models", Koller and Friedman (Chapter 4)
  3. "An introduction to probabilistic graphical models", Michael I. Jordan, 2002 (Section 2.2)  
<http://people.eecs.berkeley.edu/~jordan/prelims/chapter2.pdf>
  4. "Pattern recognition and machine learning", Christopher Bishop (Chapter 8, Section 8.3).
  5. <http://www.cs.cmu.edu/~epxing/Class/10708/lectures/lecture3-MRFrepresentation.pdf>, Eric Xing

# Learning Outcomes

- Students should be able to:
  1. Explain the concepts of **Markov properties (global, local and pairwise)** and use it to find all conditional independences in an UGM.
  2. Use **clique potential functions** to parameterize a Markov Random Field, i.e. to represent the joint distribution with clique potential functions.
  3. Describe the differences and similarities between a **Markov Random Field** and **Conditional Random Field**.

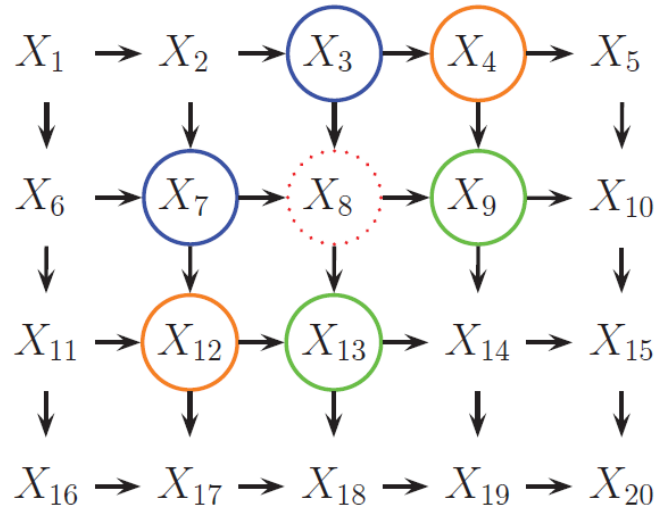
# Why Undirected Graphical Models?

- We discussed the **Directed Graphical Models** (DGMs) or Bayesian Networks in the last lecture.
- However, for some domains, the requirement for a directed edge is **rather awkward**.

# Why Undirected Graphical Models?

**Example:**

Causal MRF



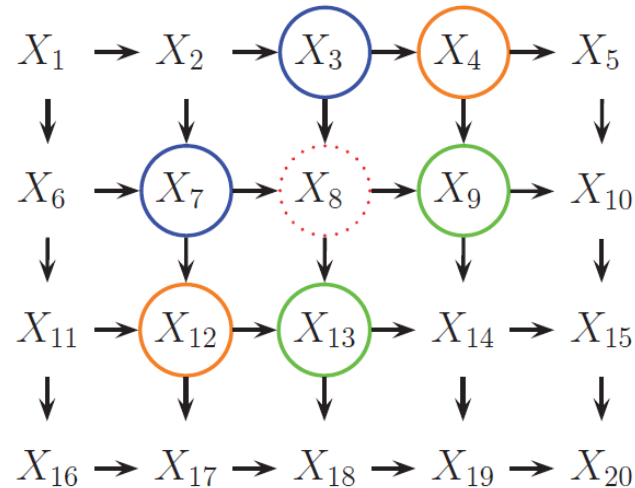
- Modeling a 2D image where the intensity of neighboring pixels are correlated.
- We can create a DAG model with a 2d lattice topology known as a **causal MRF** or a **Markov mesh**.

Image Source: “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

# Why Undirected Graphical Models?

## Example:

# Causal MRF



- However, its conditional independence properties are rather **unnatural**.
- The **Markov blanket** of the node  $X_8$  in the middle is the other colored nodes (3, 4, 7, 9, 12 and 13) rather than just its 4 nearest neighbors as one might expect.

Image Source: “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

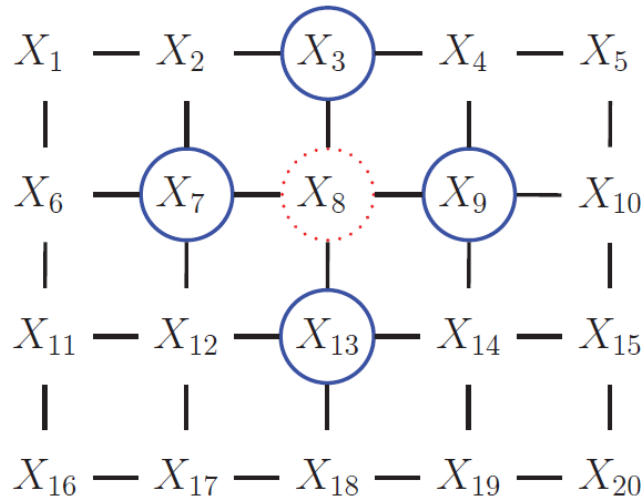
# Why Undirected Graphical Models?

- An alternative is to use an **Undirected Graphical model (UGM)**, also called a **Markov Random Field (MRF)** or **Markov network**.
- Formally, an UGM is a graph  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ , where:
  - $\mathcal{V}$  is a set of **nodes** that are in one-to-one correspondence with a set of random variables.
  - $\mathcal{E}$  is a set of **undirected** edges.
- No edge orientations, hence **more natural** for some problems such as **image analysis** and **spatial statistics**.



# Why Undirected Graphical Models?

## Example:



- We use an **undirected 2d lattice** to model a 2D image where the intensity of neighboring pixels are correlated.
- Now the **Markov blanket** of each node is just its nearest neighbors (more on Markov blanket for UGMs later).

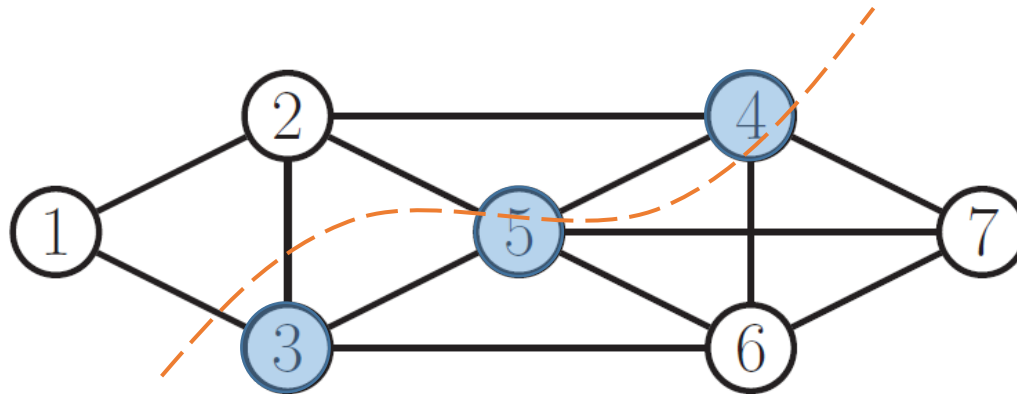
Image Source: "Machine Learning – A Probabilistic Perspective", Kevin Murphy

# Conditional Independence

## 1. Global Markov Property

- Given the sets of nodes  $A, B$  and  $C$ ,  $X_A \perp X_B \mid X_C$  if and only if  $C$  separates  $A$  from  $B$  in the graph  $\mathcal{G}$ .
- This means that there are **no paths** connecting any node in  $A$  to any node in  $B$  when we remove all nodes in  $C$ .

**Example:**



$$\{X_1, X_2\} \perp \{X_6, X_7\} \mid \{X_3, X_4, X_5\}$$

Image Source: Modified from “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

# Conditional Independence

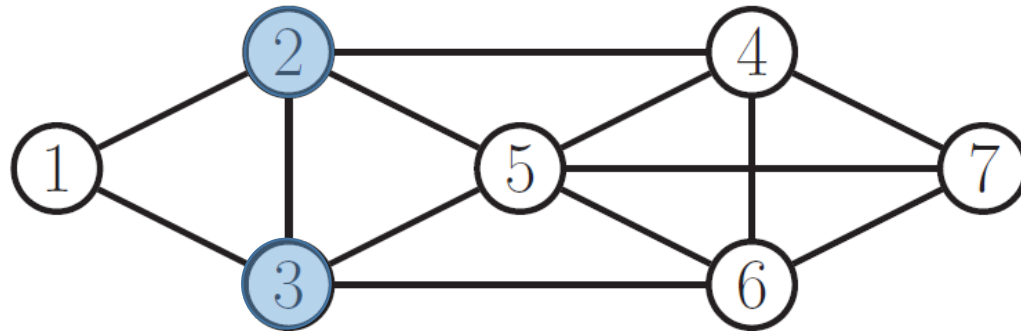
## 2. Local Markov Property

- The set of nodes that renders a node  $X_s$  conditionally independent of all the other nodes in  $\mathcal{G}$  :

$$X_s \perp \mathcal{V} \setminus \{\text{mb}(X_s), X_s\} \mid \text{mb}(X_s)$$

- where  $\text{mb}(X_s)$  is called the **Markov blanket** of  $X_s$ .

**Example:**



$$\text{mb}(X_1) = \{X_2, X_3\}, \text{ i.e. } X_1 \perp \{X_4, X_5, X_6, X_7\} \mid \{X_2, X_3\}$$

Image Source: Modified from “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

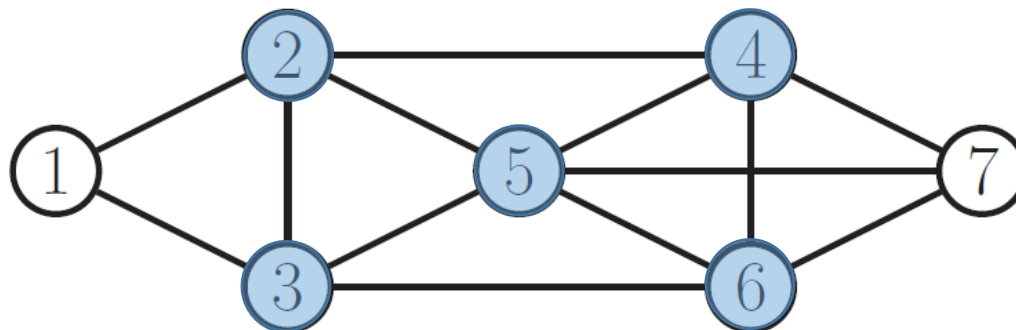
# Conditional Independence

## 3. Pairwise Markov Property

- Two nodes  $X_s$  and  $X_t$  are conditionally independent given the rest if there is **no direct edge** between them:

$$X_s \perp X_t \mid \mathcal{V} \setminus \{X_s, X_t\}, \text{ where } \mathcal{E}_{st} = \emptyset$$

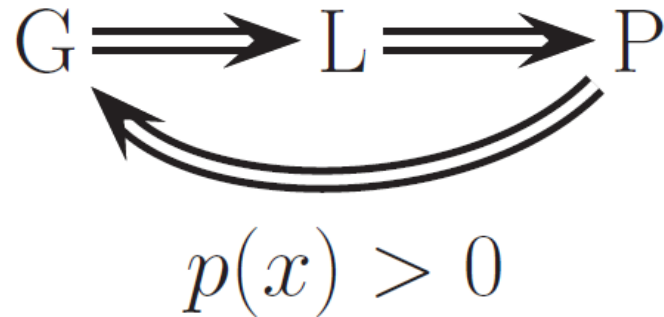
**Example:**



$$X_1 \perp X_7 \mid \{X_2, X_3, X_4, X_5, X_6\}$$

Image Source: Modified from “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

# Conditional Independence

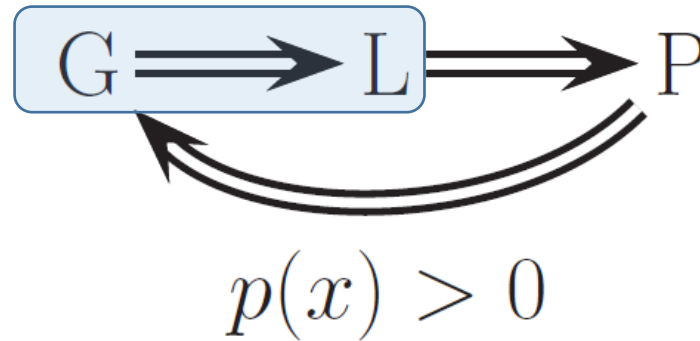


- Obvious that global Markov implies local Markov which implies pairwise Markov.
- What is less obvious, but true (assuming  $p(\mathbf{x}) > 0$  for all  $\mathbf{x}$ ), is that **pairwise Markov implies global Markov**.

Image Source: “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

# Conditional Independence

Global Markov implies local Markov:



## Proof Sketch:

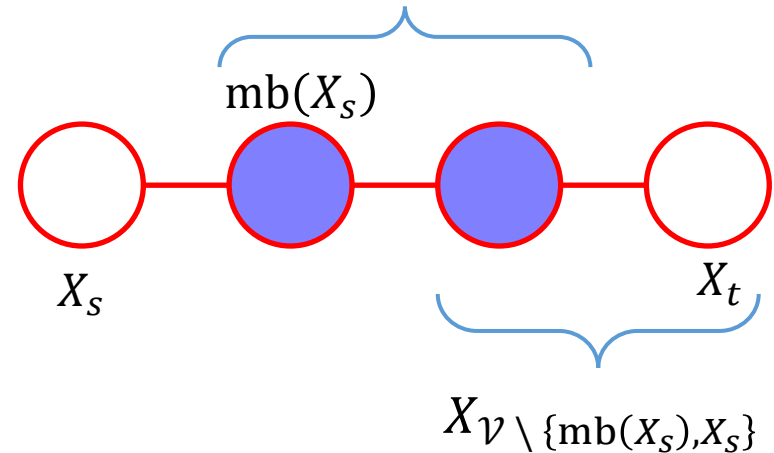
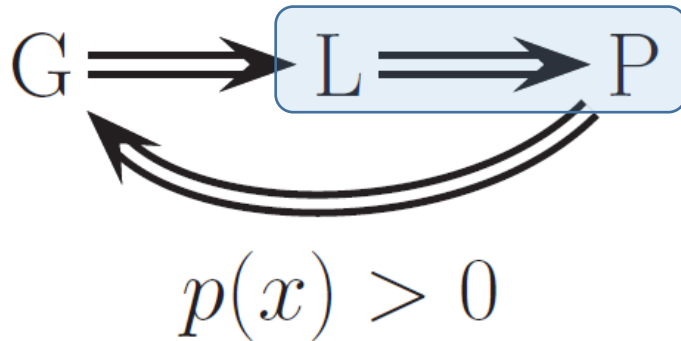
The **global Markov property** implies the **local Markov property**: this is the case when the sets  $X_A = X_S$ ,  $X_C = \text{mb}(X_S)$ , and  $X_B = \mathcal{V} \setminus \{\text{mb}(X_S), X_S\}$ .

$$X_A \perp X_B \mid X_C \Rightarrow X_S \perp \mathcal{V} \setminus \{\text{mb}(X_S), X_S\} \mid \text{mb}(X_S)$$

Image Source: "Machine Learning – A Probabilistic Perspective", Kevin Murphy

# Conditional Independence

Local Markov implies pairwise Markov:  $X_{\mathcal{V} \setminus \{X_s, X_t\}}$



## Proof Sketch :

Given any node  $X_t$  that is not adjacent to the node  $X_s$ , it follows from **local Markov property** that:

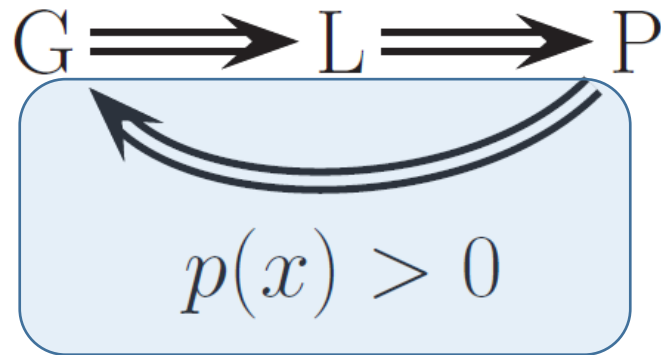
$$X_s \perp X_{\mathcal{V} \setminus \{mb(X_s), X_s\}} \mid mb(X_s)$$

This implies  $X_s \perp X_t \mid X_{\mathcal{V} \setminus \{X_s, X_t\}}$ , i.e. **pairwise Markov property**.

Image Source: "Machine Learning – A Probabilistic Perspective", Kevin Murphy

# Conditional independence

Pairwise Markov implies global Markov:



To prove this, we have to first understand the **Independence-Map** and the **Intersection Lemma**!

Image Source: “Machine Learning – A Probabilistic Perspective”, Kevin Murphy



# Independence-Map

- Also known as the **I-Map**.
- The I-Map of a **joint distribution**  $p(x_1, \dots, x_N)$ , often written as  $I(p)$  represents **all independencies** in  $p(x_1, \dots, x_N)$ .
- Similarly, the I-Map of a **directed/undirected graph**  $\mathcal{G}$ , i.e.  $I(\mathcal{G})$  represents **all independencies** encoded in  $\mathcal{G}$ .
- $\mathcal{G}$  is a valid representation of  $p$  if  $I(\mathcal{G}) \subseteq I(p)$ .

# Independence-Map

## Example:

- **Given:** 4 disjoint sets  $W, X, Y, Z$ , where the **non-zero** distribution  $p(X, Y, Z, W)$  contains **at least two conditional independences**.
- I-map implies that all the following are **valid factorizations** (encoded in  $\mathcal{G}$ ) of the joint distribution  $p(X, Y, Z, W)$ :

$$p(X, Y, Z, W) = p(X \mid Z, W)p(Y \mid Z, W)$$

for the conditional independence  $X \perp Y \mid \{Z, W\}$

$$p(X, Y, Z, W) = p(X \mid Z, Y)p(W \mid Z, Y)$$

for the conditional independence  $X \perp W \mid \{Z, Y\}$

# Independence-Map

## Example:

- Without a loss of generality, we replace the conditional probabilities with **functions**  $f(\cdot)$  and  $g(\cdot)$ :

$$\begin{aligned} p(X, Y, Z, W) &= p(X \mid Z, W) p(Y \mid Z, W) \\ &= f_{XWZ}(X, W, Z) f_{WYZ}(W, Y, Z) \end{aligned}$$

$$\begin{aligned} p(X, Y, Z, W) &= p(X \mid Z, Y) p(W \mid Z, Y) \\ &= g_{XYZ}(X, Y, Z) g_{W,Y,Z}(W, Y, Z) \end{aligned}$$

# Independence-Map

## Example:

- The factorizations are valid because  $p(X, Y, Z, W)$  contains **at least two conditional independences**, i.e.,  $X \perp Y \mid \{Z, W\}$  and  $X \perp W \mid \{Z, Y\}$ .
- It also illustrates the **non-uniqueness** of probability factorization.
- And the respective factorizations encodes only **one conditional independence** each, i.e.,  $I(f) \subseteq I(p)$  and  $I(g) \subseteq I(p)$ .

# Intersection Lemma

- **c.f. previous example:** We know that  $p(X, Y, Z, W)$  has **at least** two conditional independences, i.e.,  $X \perp Y \mid \{Z, W\}$  and  $X \perp W \mid \{Z, Y\}$ .
- This means that **another factorization that satisfies BOTH** the conditional independences is also valid and must exist!
- Since

$$\begin{aligned} p(X, Y, Z, W) &= f_{XWZ}(X, W, Z) f_{WYZ}(W, Y, Z), \\ p(X, Y, Z, W) &= g_{XYZ}(X, Y, Z) g_{W,Y,Z}(W, Y, Z) \end{aligned}$$

represent the same distribution  $p(X, Y, Z, W)$ , we can **equate them**, i.e.

$$f_{XWZ}(X, W, Z) f_{WYZ}(W, Y, Z) = g_{XYZ}(X, Y, Z) g_{W,Y,Z}(W, Y, Z)$$

# Intersection Lemma

$$f_{XWZ}(X, W, Z)f_{WYZ}(W, Y, Z) = g_{XYZ}(X, Y, Z)g_{W,Y,Z}(W, Y, Z)$$

- Through inspection, we see that  $\{X, Z\}$  and  $\{W, Y, Z\}$  must appear in two factors, respectively.

$$f_{XWZ}(X, W, Z)f_{WYZ}(W, Y, Z) = g_{XYZ}(X, Y, Z)g_{W,Y,Z}(W, Y, Z)$$



positive distributions

$$p(X, Y, Z, W) = \mu_{XZ}(X, Z)\mu_{W,Y,Z}(W, Y, Z)$$

# Intersection Lemma

- As a result, we get **both conditional independences**  $X \perp Y \mid \{Z, W\}$  and  $X \perp W \mid \{Z, Y\}$  encoded in the **same factorization**.
- In addition, we observe **an additional conditional independence**  $X \perp \{Y, W\} \mid Z$ , which is the **Intersection Lemma**!
- **Remark:** For the conditional independences to exist, the probability distributions **cannot be zero**, e.g. if  $p(Z) = 0$ , then

$$p(X, W, Y \mid Z) = \frac{p(X, W, Y, Z)}{p(Z)} \quad \text{Is undefined!!!}$$

# Intersection Lemma

For **positive distributions**, and for mutually disjoint sets  $X, Y, W, Z$  :

$$\text{If } X \perp Y \mid \{W, Z\} \text{ and } X \perp W \mid \{Y, Z\} \Rightarrow X \perp \{Y, W\} \mid Z$$

From  $X \perp Y \mid \{Z, W\}$  and  $X \perp W \mid \{Z, Y\}$ , we can write the joint distribution  $p(X, Y, Z, W)$  as:

$$\overbrace{f_{XWZ}(X, W, Z)f_{WYZ}(W, Y, Z)}^{X \perp Y \mid \{Z, W\}} = \overbrace{g_{XYZ}(X, Y, Z)g_{W,Y,Z}(W, Y, Z)}^{X \perp W \mid \{Z, Y\}}$$



**positive distributions**

$$p(X, Y, Z, W) = \mu_{XZ}(X, Z)\mu_{W,Y,Z}(W, Y, Z) \Rightarrow X \perp \{Y, W\} \mid Z$$

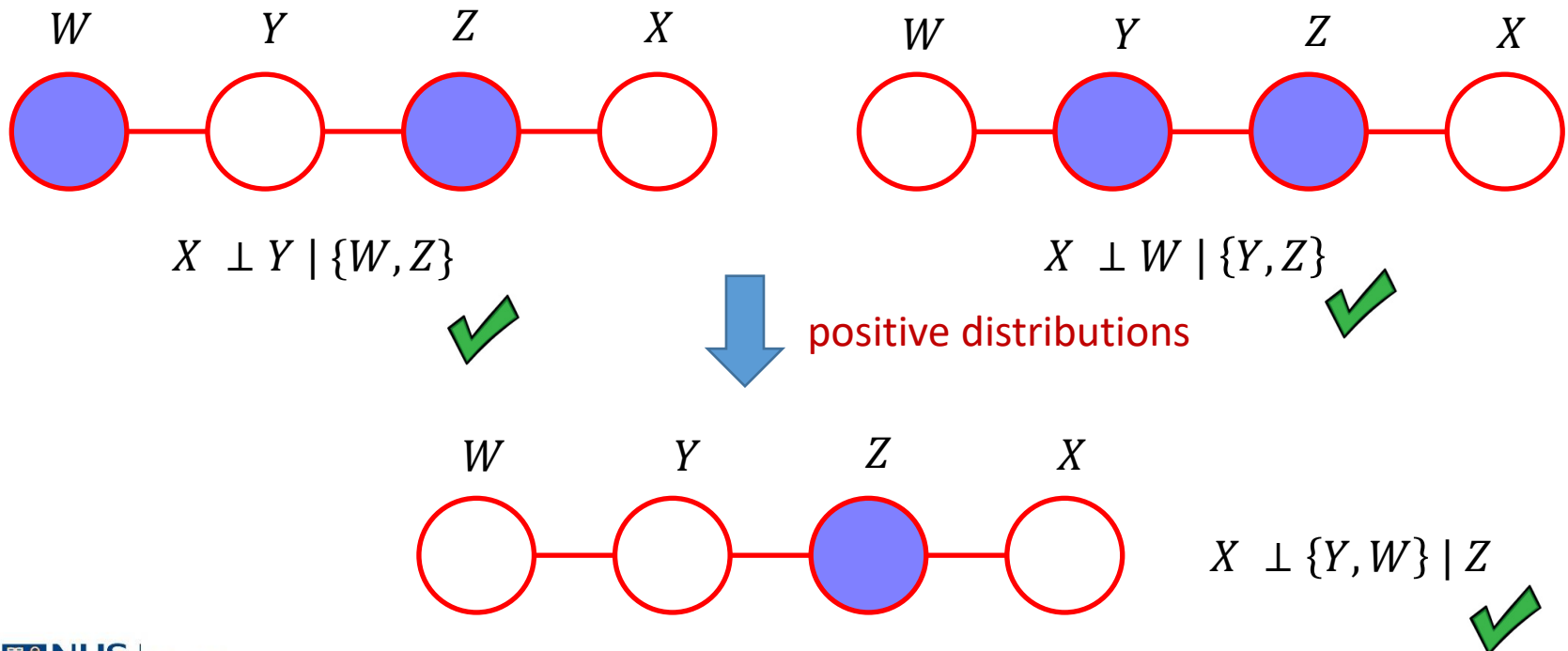


# Intersection Lemma

For **positive distributions**, and for mutually disjoint sets  $X, Y, W, Z$  :

$$\text{If } X \perp Y \mid \{W, Z\} \text{ and } X \perp W \mid \{Y, Z\} \Rightarrow X \perp \{Y, W\} \mid Z$$

## Example 1:

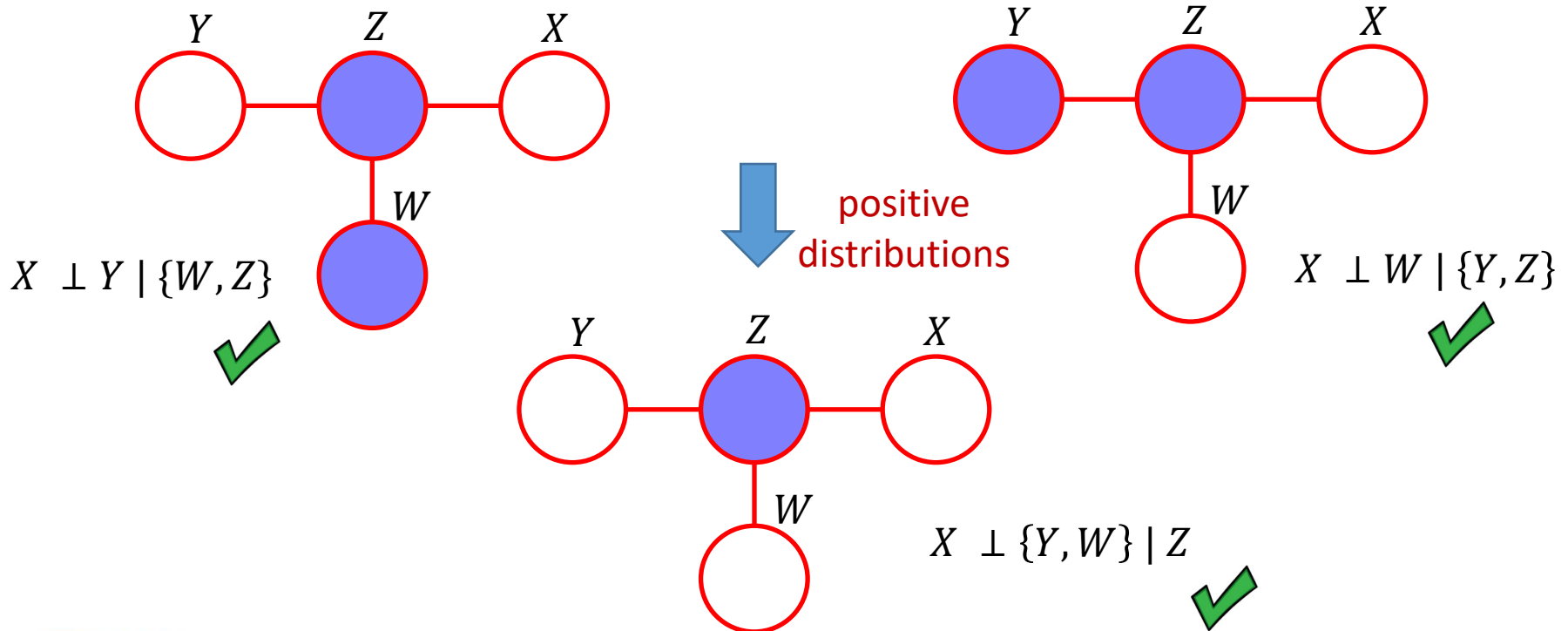


# Intersection Lemma

For **positive distributions**, and for mutually disjoint sets  $X, Y, W, Z$  :

$$\text{If } X \perp Y \mid \{W, Z\} \text{ and } X \perp W \mid \{Y, Z\} \Rightarrow X \perp \{Y, W\} \mid Z$$

## Example 2:

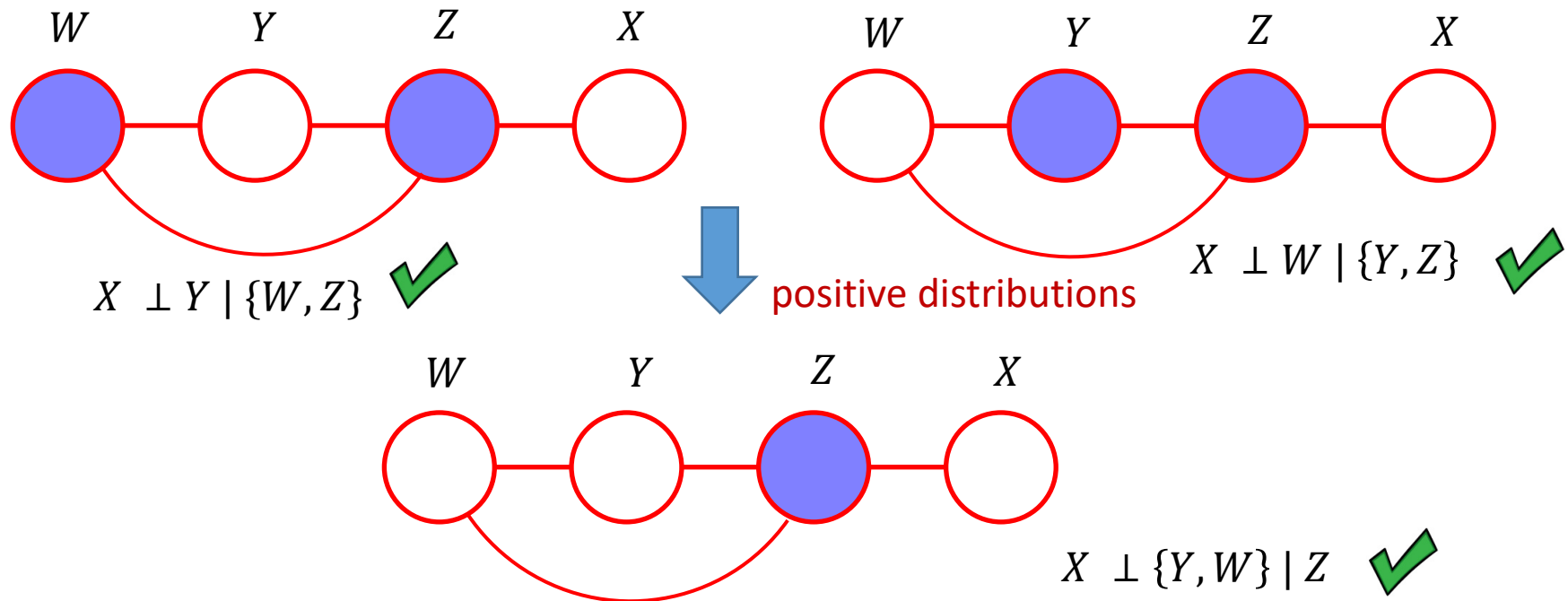


# Intersection Lemma

For **positive distributions**, and for mutually disjoint sets  $X, Y, W, Z$  :

$$\text{If } X \perp Y \mid \{W, Z\} \text{ and } X \perp W \mid \{Y, Z\} \Rightarrow X \perp \{Y, W\} \mid Z$$

## Example 3:

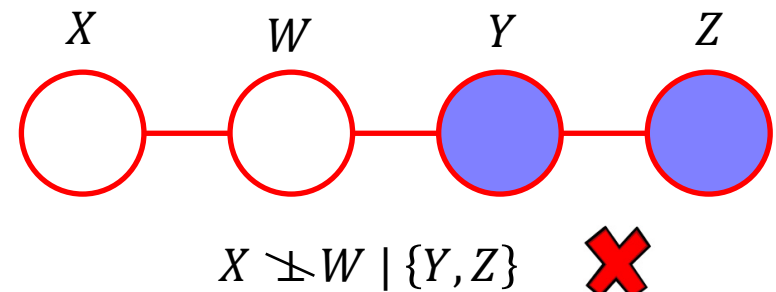
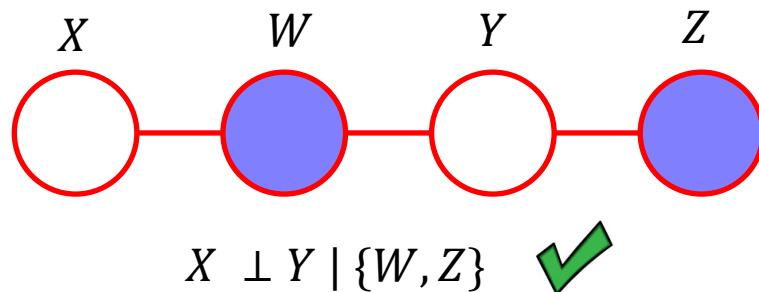


# Intersection Lemma

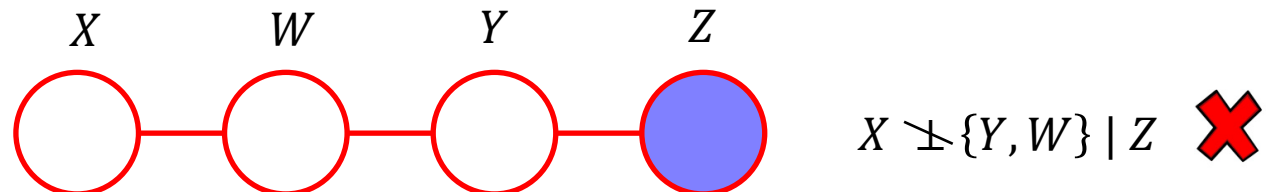
For **positive distributions**, and for mutually disjoint sets  $X, Y, W, Z$  :

$$\text{If } X \perp Y \mid \{W, Z\} \text{ and } X \perp W \mid \{Y, Z\} \Rightarrow X \perp \{Y, W\} \mid Z$$

## Counter-Example:



positive distributions



# Intersection Lemma

$$X \perp Y \mid \{W, Z\} \text{ and } X \perp W \mid \{Y, Z\} \Rightarrow X \perp \{Y, W\} \mid Z$$

**Proof:**

$$p(X, W, Y \mid Z) = p(X, Y \mid Z, W)p(W \mid Z) \quad (\text{chain rule})$$

$$= p(X \mid Z, W)p(Y \mid Z, W)p(W \mid Z) \quad (X \perp Y \mid \{W, Z\})$$

$$= p(X \mid Z) \frac{p(Y, W, Z) \cancel{p(Z, W)}}{\cancel{p(Z, W)}} \frac{1}{p(Z)}$$

$$(X \perp W \mid \{Y, Z\}, X \perp Y \mid \{W, Z\} \text{ and chain rule})$$

$$= p(X \mid Z)p(Y, W \mid Z) \quad (\text{chain rule})$$

$$\Rightarrow X \perp \{Y, W\} \mid Z$$

□

# Intersection Lemma

- We have  $p(X | Z, W) = p(X | Z)$  because:

$$\begin{aligned} p(X | Z, W) &= p(X | Z, W) \sum_Y p(Y | Z) \quad (\text{sum rule}) \\ &= \sum_Y p(X | Z, W) p(Y | Z) \end{aligned}$$

- Since we have  $X \perp Y | \{W, Z\}$  and  $X \perp W | \{Y, Z\}$ , this implies  $p(X | Z, W) = p(X | Y, Z, W) = p(X | Z, Y)$ .  
Thus, we get:

$$p(X | Z, W) = \sum_Y p(X | Z, Y) p(Y | Z) \quad (\text{chain rule})$$

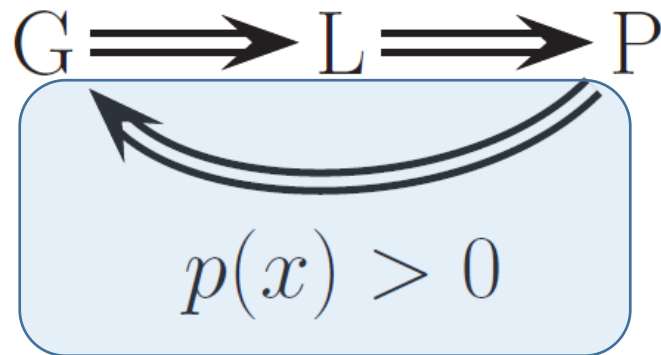
$$= \sum_Y \frac{p(X, Y, Z)}{p(Z, Y)} \frac{p(Z, Y)}{p(Z)} = \sum_Y \frac{p(X, Y, Z)}{p(Z)} \quad (\text{chain rule})$$

$$= \sum_Y p(X, Y | Z) = p(X | Z). \quad (\text{sum rule})$$

□

# Conditional independence

Pairwise Markov implies global Markov:



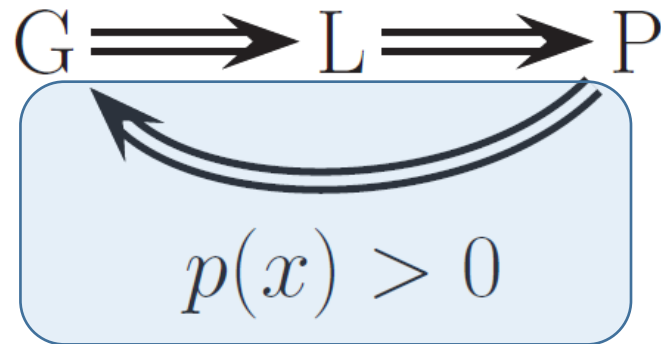
## Proof Sketch:

Let  $S, A, B, D \subset \mathcal{V}$  be disjoint sets of nodes with  $S$  separating  $A$  from  $B$  in the graph  $\mathcal{G}$ , where  $A \neq \emptyset$  and  $B \neq \emptyset$ . We will show a sketch of proof that **pairwise Markov** implies **global Markov** using backward induction.

Image Source: “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

# Conditional independence

Pairwise Markov implies global Markov:



**Proof Sketch :**

Let  $d = |\mathcal{V}|$ , when  $|S| = d - 2$ :

$A \perp B \mid S$ , where  $|A| = |B| = 1$

$\Rightarrow$  pairwise Markov

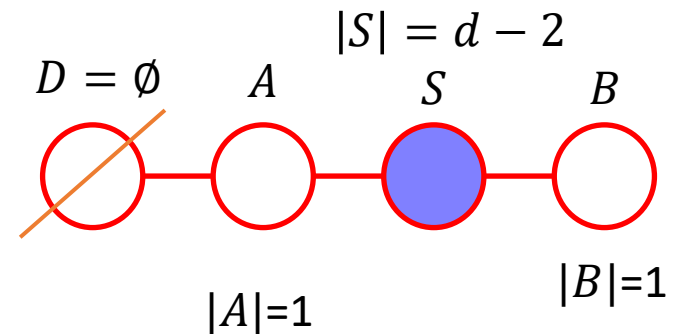
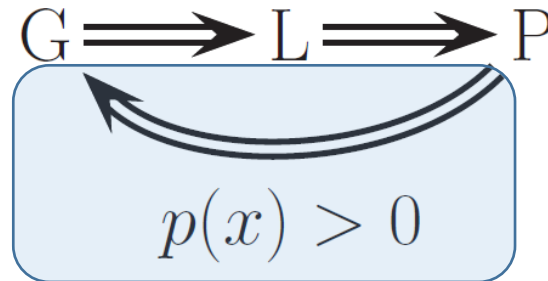


Image Source: "Machine Learning – A Probabilistic Perspective", Kevin Murphy



# Conditional independence

Pairwise Markov implies global Markov:



**Proof Sketch :**

For  $|S| < d - 2$ , WLOG, let us assume the set of nodes  $D$  is connected only to  $A$ , where  $|D| \geq 1$ ,  $|A| \geq 1$  and  $|B| \geq 1$ .

We have:

$$A \perp B \mid \{S, D\} \text{ and } B \perp D \mid \{A, S\}$$

Intersection  
Lemma



$$B \perp \{A, D\} \mid S \Rightarrow \text{global Markov}$$

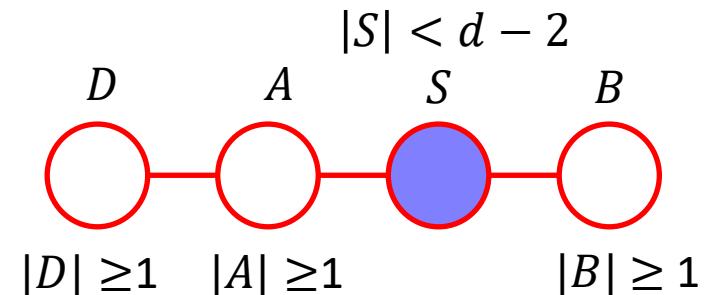
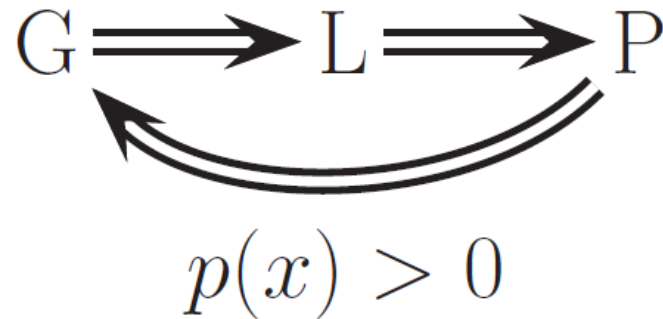


Image Source: "Machine Learning – A Probabilistic Perspective", Kevin Murphy

# Conditional independence



- The importance of this result is that it is usually **easier** to empirically **assess pairwise conditional independence**.
- Such pairwise CI statements **can be used to construct a graph** from which global CI statements can be extracted.

Image Source: “Machine Learning – A Probabilistic Perspective”, Kevin Murphy

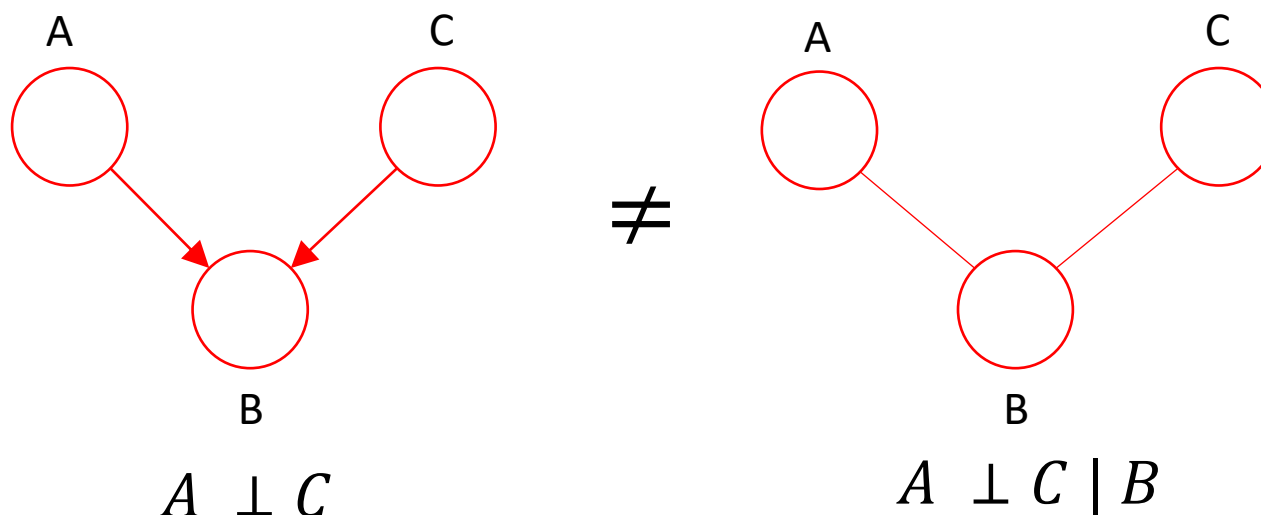
# Comparative Semantics

- We have seen that it is **easier** to determine conditional independence using UGMs than DGMs.
- **Question:** Can we determine conditional independence in a DGM using a UGM, or vice versa?
- This is **NOT possible in general!**

# Comparative Semantics

- It is tempting to simply convert the DGM to a UGM by **dropping the orientation of the edges**, but this is **not always** correct!

Example:

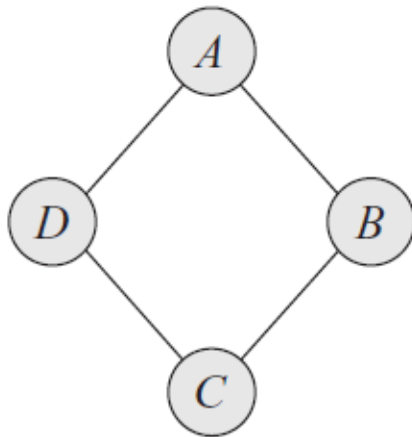


This conditional independence  
is **NOT** in the DGM!

# Comparative Semantics

- An example of some CI relationships that can be perfectly modeled by a UGM but not a DGM:

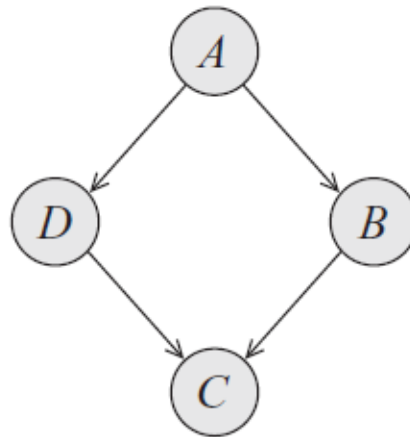
UGM



$$A \perp C \mid \{B, D\}$$

$$B \perp D \mid \{A, C\}$$

Attempt 1

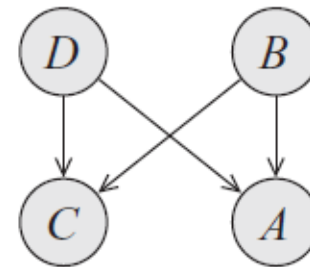


$$\checkmark A \perp C \mid \{B, D\}$$

$$\checkmark B \perp D \mid A$$

$$\times B \perp D \mid \{A, C\}$$

Attempt 2



$$\checkmark A \perp C \mid \{B, D\}$$

$$\checkmark B \perp D$$

$$\times B \perp D \mid \{A, C\}$$

Image source: Koller and Friedman 2009

# Parameterization of MRFs

- As in the case of DGMs, we would like to obtain a **local parameterization** for UGMs.
- We have seen earlier that for **DGMs**:
  - Parameterization was based on **local conditional probabilities** of a node and its parents, i.e.  $p(x_i|x_{\pi_i})$ .
  - Joint probability is a **product of local conditional probabilities** as a result of the chain rule, i.e.

$$p(x_1, \dots, x_N) = \prod_{i=1}^N p(x_i|x_{\pi_i})$$

# Parameterization of MRFs

- Difficult to do local parameterization based on conditional probabilities since **no topological ordering** associated with UGMs.
- It turns out that its better to **abandon conditional probabilities** altogether and **use some functions** instead.

# Parameterization of MRFs

- **Lose the ability** to give local probabilistic interpretation to the functions used to represent the joint probability.
- **Retain the ability** the all-important representation of the joint as a product of local functions.



# Parameterization of MRFs

How do we decide the domain of the **local functions**?

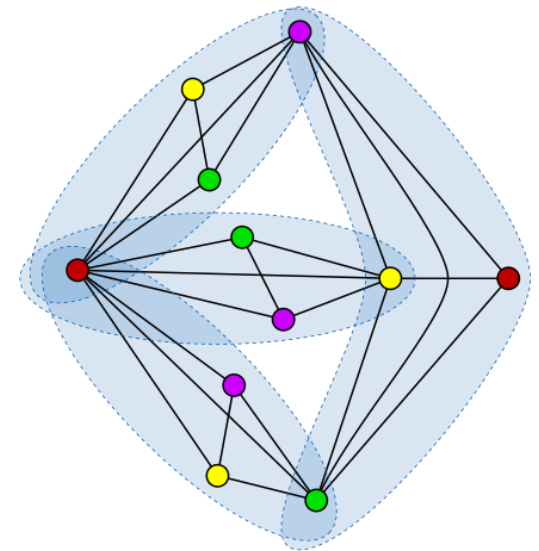
- Recall two nodes  $X_i$  and  $X_j$  that are not directly linked in an UGM are **conditionally independent** given all other nodes.
- Thus, it must be possible to obtain a factorization of the joint probability that **places  $X_i$  and  $X_j$  in different factors**.
- This implies that we **cannot have** a local function that depends on both  $X_i$  and  $X_j$ .

$$p(x_1, \dots, x_N) \neq \psi_1(x_i, x_j, \dots) \dots \psi_m(\dots)$$

# Parameterization of MRFs

How do we decide the domain of the **local functions**?

- Our argument thus far suggested that all nodes  $X_C$  that belong to a **maximal clique**  $C$  in the UGM appear together in a local function  $\psi(x_C)$ .
- A **clique** of a graph is a fully-connected subset of nodes.
- The **maximal cliques** of a graph are the cliques that cannot be extended to include additional nodes without losing the property of being fully connected.



**Conditional independence is impossible for any two nodes in maximal clique!**

# Parameterization of MRFs

## Hammersley-Clifford theorem:

A **positive distribution**  $p(y) > 0$  satisfies the CI properties of an undirected graph  $\mathcal{G}$  iff  $p$  can be represented as **a product of factors**, one per maximal clique:

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{y}_c|\boldsymbol{\theta}_c)$$

where

- $\mathcal{C}$  is the set of all the maximal cliques of  $\mathcal{G}$
- $\psi_c(\cdot)$  is the **factor** or **potential function** of clique  $c$
- $\theta$  is the parameter of the factors  $\psi_c(\cdot)$  for  $c \in \mathcal{C}$
- $Z(\theta)$  is the **partition function**

# Parameterization of MRFs

## Hammersley-Clifford theorem:

$Z(\theta)$  is the **partition function** given by:

$$Z(\theta) \triangleq \sum_{\mathbf{y}} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{y}_c | \theta_c)$$

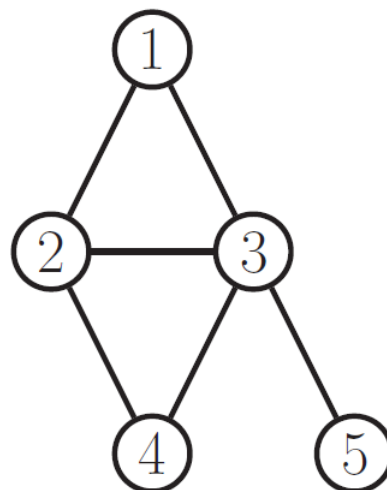
Since  $\psi(\cdot)$  can be any arbitrary positive function, the partition function  $Z(\theta)$  ensures the **overall distribution sums to 1**.

# Parameterization of MRFs

Example:

Assume:  $y_i \in \{0,1\}$

$y_2$	$y_3$	$y_4$	$\psi_{234}$
0	0	0	$a_1$
0	0	1	$a_2$
0	1	0	$a_3$
0	1	1	$a_4$
1	0	0	$a_5$
1	0	1	$a_6$
1	1	0	$a_7$
1	1	1	$a_8$



$y_1$	$y_2$	$y_3$	$\psi_{123}$
0	0	0	$b_1$
0	0	1	$b_2$
0	1	0	$b_3$
0	1	1	$b_4$
1	0	0	$b_5$
1	0	1	$b_6$
1	1	0	$b_7$
1	1	1	$b_8$

$y_3$	$y_5$	$\psi_{35}$
0	0	$c_1$
0	1	$c_2$
1	0	$c_3$
1	1	$c_4$

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \psi_{123}(y_1, y_2, y_3) \psi_{234}(y_2, y_3, y_4) \psi_{35}(y_3, y_5)$$

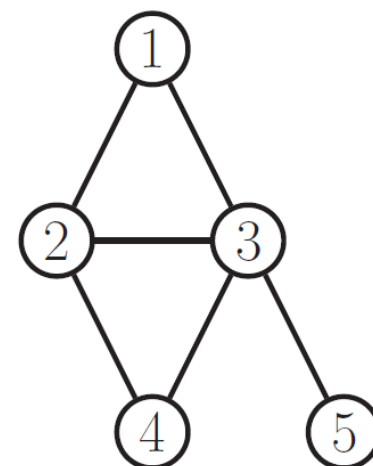
where

$$Z = \sum_{\mathbf{y}} \psi_{123}(y_1, y_2, y_3) \psi_{234}(y_2, y_3, y_4) \psi_{35}(y_3, y_5)$$

“An introduction to probabilistic graphical models”, Michael I. Jordan, 2002

# Pairwise Parameterization of MRFs

- Parameterization of MRFs is **not unique**!!!
- We are **free to relax** the parameterization to the edges of the graph, rather than the maximal cliques.
- This is **pairwise MRF**, and is widely used due to its simplicity, although it is not as general.



## Example:

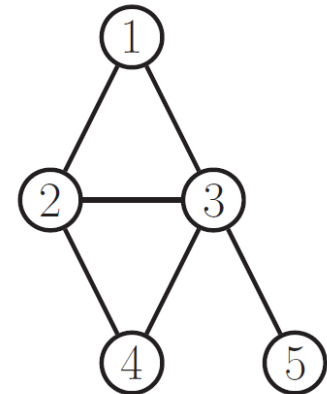
$$\begin{aligned} p(\mathbf{y}|\boldsymbol{\theta}) &\propto \psi_{12}(y_1, y_2)\psi_{13}(y_1, y_3)\psi_{23}(y_2, y_3)\psi_{24}(y_2, y_4)\psi_{34}(y_3, y_4)\psi_{35}(y_3, y_5) \\ &\propto \prod_{s \sim t} \psi_{st}(y_s, y_t) \end{aligned}$$

# Canonical Parameterization of MRFs

- Another way is to define the parameterization over **all cliques** in the graph, i.e., **canonical parameterization**.
- **Uniform prior** can be assumed on any potential function.

**Example:** Assume:  $y_i \in \{0,1\}$

Uniform Prior on pairwise potentials (**optional**)



$$p(\mathbf{y} | \theta) \propto \psi_1(y_1)\psi_2(y_2)\psi_3(y_3)\psi_4(y_4)\psi_5(y_5) \\
\cancel{\psi_{12}(y_1, y_2)}\cancel{\psi_{23}(y_2, y_3)}\cancel{\psi_{13}(y_1, y_3)} \\
\cancel{\psi_{24}(y_2, y_4)}\cancel{\psi_{34}(y_3, y_4)}\cancel{\psi_{35}(y_3, y_5)} \\
\psi_{123}(y_1, y_2, y_3)\psi_{234}(y_2, y_3, y_4)$$

$$= \psi_1(y_1)\psi_2(y_2)\psi_3(y_3)\psi_4(y_4)\psi_5(y_5) \\
\psi_{123}(y_1, y_2, y_3)\psi_{234}(y_2, y_3, y_4)$$

$y_i$	$y_j$	$\psi_{ij}$
0	0	1
0	1	1
1	0	1
1	1	1

# Gibbs Distribution

- There is a deep connection between UGMs and **statistical physics**.
- In particular, there is a model known as the **Gibbs distribution**, which can be written as follows:

$$p(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\left(-\sum_c E(\mathbf{y}_c|\boldsymbol{\theta}_c)\right)$$

- $E(\mathbf{y}_c) > 0$  is the **energy** associated with the variables in clique  $c$ .



# Gibbs Distribution

- We can convert the Gibbs distribution to a UGM by defining:

$$\psi_c(\mathbf{y}_c | \boldsymbol{\theta}_c) = \exp(-E(\mathbf{y}_c | \boldsymbol{\theta}_c))$$

$$\Rightarrow p(\mathbf{y} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp(-\sum_c E(\mathbf{y}_c | \boldsymbol{\theta}_c)) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{y}_c | \boldsymbol{\theta}_c)$$

- We see that **high probability states** correspond to **low energy configurations**.
- Also known as **energy based models**, hence the term **“potential function”** for  $\psi_c(\cdot)$ .

# Log-Linear Potential Functions

- Potentials represent the **relative “compatibility”** between the different assignments to the random variables.
- A general approach is to define the **log potentials** as a **linear function of the parameters**  $\theta_c \in \mathbb{R}^M$ :

$$\log \psi_c(\mathbf{y}_c) \triangleq \phi_c(\mathbf{y}_c)^T \boldsymbol{\theta}_c$$

- $\phi_c: y_c \mapsto \mathbb{R}^M$  is a **feature vector** derived from the values of the variables  $y_c$ .

# Log-Linear Potential Functions

- The resulting **log probability** has the form:

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \sum_c \phi_c(\mathbf{y}_c)^T \boldsymbol{\theta}_c - \log Z(\boldsymbol{\theta})$$

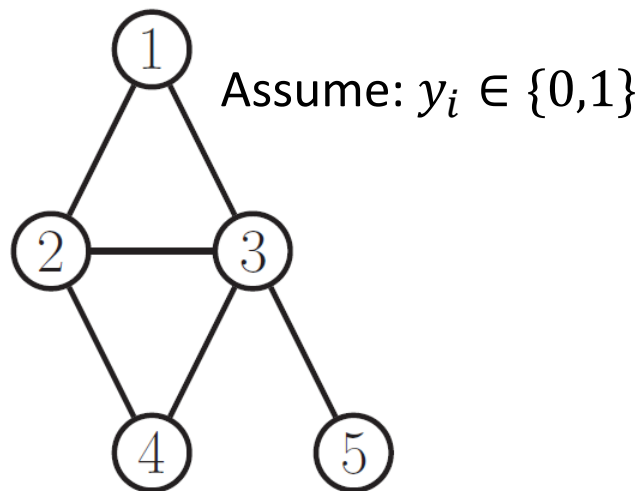
- This is also known as a **maximum entropy** or a **log-linear** model.

# Log-Linear Potential Functions

**Example:**  $\log \psi_c(\mathbf{y}_c) \triangleq \phi_c(\mathbf{y}_c)^T \boldsymbol{\theta}_c$

$$\Rightarrow \psi_c(\mathbf{y}_c) = \exp\{\phi_c(\mathbf{y}_c)^T \boldsymbol{\theta}_c\}$$

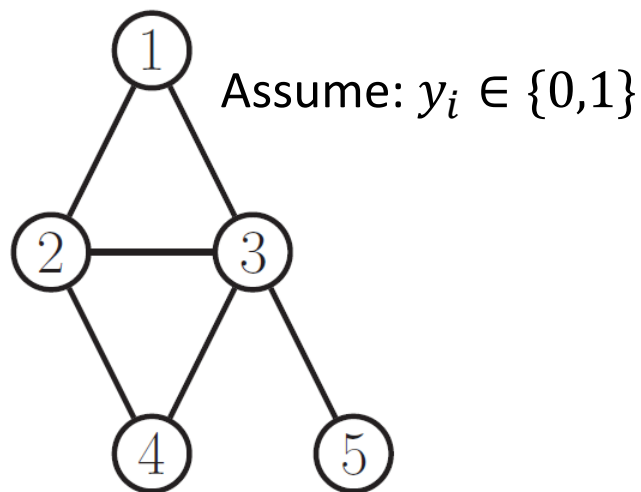
Let's look at  $\mathbf{y}_c = \{y_1, y_2, y_3\}$ , we set  $M = |c|$ , such that  $\boldsymbol{\theta}_c \in \mathbb{R}^M$  and  $\phi_c: \mathbf{y}_c \mapsto \mathbb{R}^M$ .



$y_1$	$y_2$	$y_3$	$\psi_{123}$
0	0	0	$\phi_c(y_1 = 0, y_2 = 0, y_3 = 0)^T \boldsymbol{\theta}_c$
0	0	1	$\phi_c(y_1 = 0, y_2 = 0, y_3 = 1)^T \boldsymbol{\theta}_c$
0	1	0	$\phi_c(y_1 = 0, y_2 = 1, y_3 = 0)^T \boldsymbol{\theta}_c$
0	1	1	$\phi_c(y_1 = 0, y_2 = 1, y_3 = 1)^T \boldsymbol{\theta}_c$
1	0	0	$\phi_c(y_1 = 1, y_2 = 0, y_3 = 0)^T \boldsymbol{\theta}_c$
1	0	1	$\phi_c(y_1 = 1, y_2 = 0, y_3 = 1)^T \boldsymbol{\theta}_c$
1	1	0	$\phi_c(y_1 = 1, y_2 = 1, y_3 = 0)^T \boldsymbol{\theta}_c$
1	1	1	$\phi_c(y_1 = 1, y_2 = 1, y_3 = 1)^T \boldsymbol{\theta}_c$

# Log-Linear Potential Functions

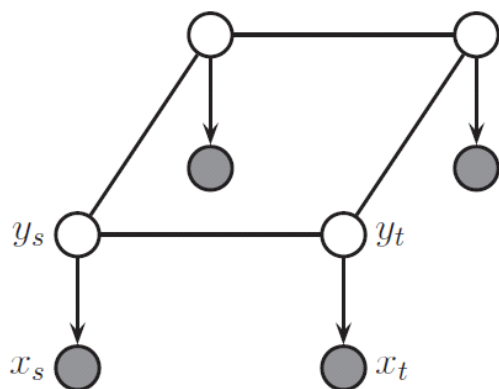
- In this case, we set  $M = 3$  and  $\phi_c: y_c \mapsto [y_1, y_2, y_3]^\top$ , we **only need 3 parameters**  $\theta_c = [\theta_{c1}, \theta_{c2}, \theta_{c3}]$  instead of  $2^3$  parameters.



$y_1$	$y_2$	$y_3$	$\psi_{123}$
0	0	0	$\phi_c(y_1 = 0, y_2 = 0, y_3 = 0)^\top \theta_c$
0	0	1	$\phi_c(y_1 = 0, y_2 = 0, y_3 = 1)^\top \theta_c$
0	1	0	$\phi_c(y_1 = 0, y_2 = 1, y_3 = 0)^\top \theta_c$
0	1	1	$\phi_c(y_1 = 0, y_2 = 1, y_3 = 1)^\top \theta_c$
1	0	0	$\phi_c(y_1 = 1, y_2 = 0, y_3 = 0)^\top \theta_c$
1	0	1	$\phi_c(y_1 = 1, y_2 = 0, y_3 = 1)^\top \theta_c$
1	1	0	$\phi_c(y_1 = 1, y_2 = 1, y_3 = 0)^\top \theta_c$
1	1	1	$\phi_c(y_1 = 1, y_2 = 1, y_3 = 1)^\top \theta_c$

# Combining UGM and DGM

## Example: Ising and Potts Models



Depth Map from Stereo Images



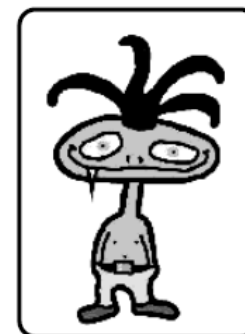
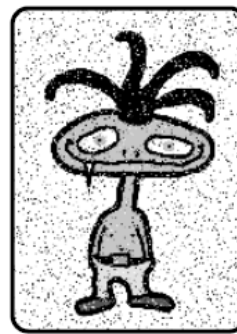
Observed Variables  
 $x \in \{1, \dots, L\}$

Latent Variables  
 $y \in \{1 \dots L\}$

$$p(y, x | J, \theta) = p(\mathbf{y} | J) \prod_t p(x_t | y_t, \theta)$$

$$= \underbrace{\left[ \frac{1}{Z(J)} \prod_{s \sim t} \psi(y_s, y_t; J) \right]}_{\text{Pairwise potential}} \underbrace{\prod_t p(x_t | y_t, \theta)}_{\text{Unary potential}}$$

Binary Image Denoising



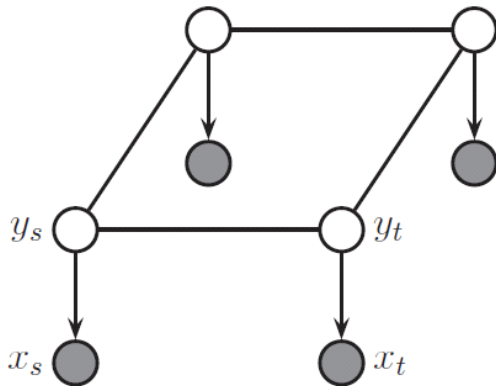
Observed Variables  
 $x \in \{0, 1\}$

Latent Variables  
 $y \in \{0, 1\}$

Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

# Combining UGM and DGM

## Example: Ising and Potts Models



$$\begin{aligned} p(y, x | J, \theta) &= p(y | J) \prod_t p(x_t | y_t, \theta) \\ &= \left[ \frac{1}{Z(J)} \prod_{s \sim t} \psi(y_s, y_t; J) \right] \prod_t p(x_t | y_t, \theta) \end{aligned}$$

Pairwise potentialUnary potential

### Ising Model:

$$y_i \in \{0, 1\}, \quad x_i \in \{0, 1\}$$

$$E(y_s, y_t; J) = J |y_s - y_t|,$$
$$J > 0$$

### Potts Model:

$$y_i \in \{1, \dots, L\}, \quad x_i \in \{1, \dots, L\}$$

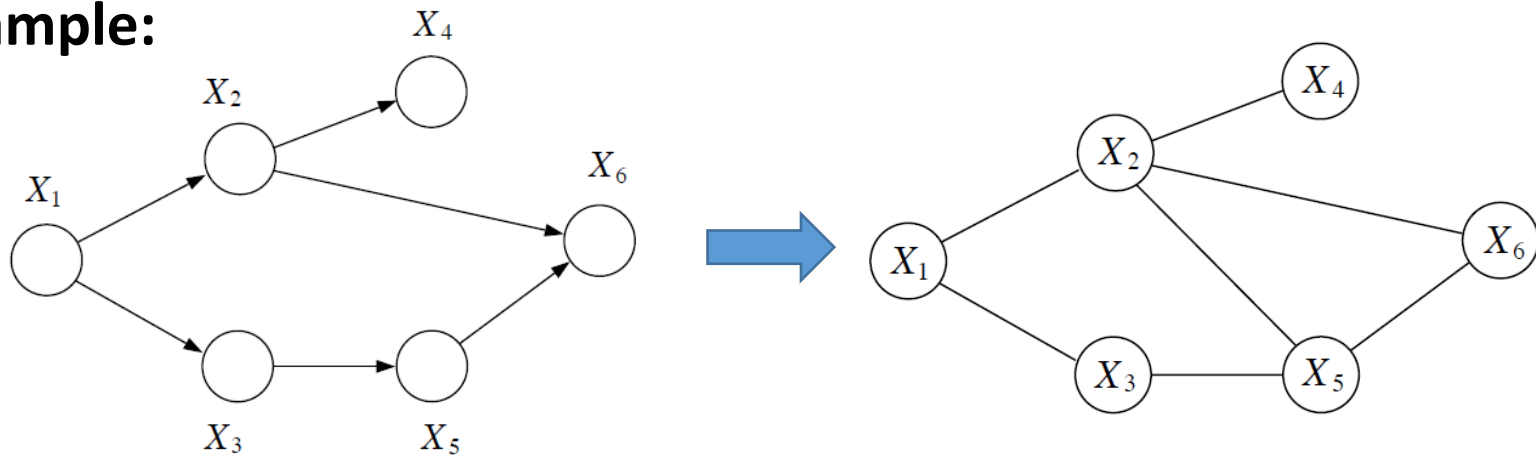
$$E(y_s, y_t; J) = J \min(|y_s - y_t|, 1),$$
$$J > 0$$

Image Source: "Computer Vision: Models, Learning, and Inference", Simon Prince

# Representing Potential Functions

- We can also use **local conditional probabilities** from a DGM to represent the potential functions in a UGM.

**Example:**



$$p(x) = \frac{1}{Z} \underbrace{\varphi_{12}(x_1, x_2)}_{p(x_2|x_1)} \underbrace{\varphi_{13}(x_1, x_3)}_{p(x_3|x_1)} \underbrace{\varphi_{24}(x_2, x_4)}_{p(x_4|x_2)} \underbrace{\varphi_{35}(x_3, x_5)}_{p(x_5|x_3)} \underbrace{\varphi_{256}(x_2, x_5, x_6)}_{p(x_6|x_2, x_5)}$$

$$Z = \sum_x \varphi_{12}(x_1, x_2) \varphi_{13}(x_1, x_3) \varphi_{24}(x_2, x_4) \varphi_{35}(x_3, x_5) \varphi_{256}(x_2, x_5, x_6) = \frac{1}{p(x_1)}$$

Image source: "An introduction to probabilistic graphical models", Michael I. Jordan, 2002.



# Moralization

- A DGM can be converted into a UGM by “marrying” the unmarried parents of a node, i.e. **moralization**.
- This process **preserves the joint distribution**, but **conditional independence is lost!**
- Moralization is important for **exact inference** (next lectures).

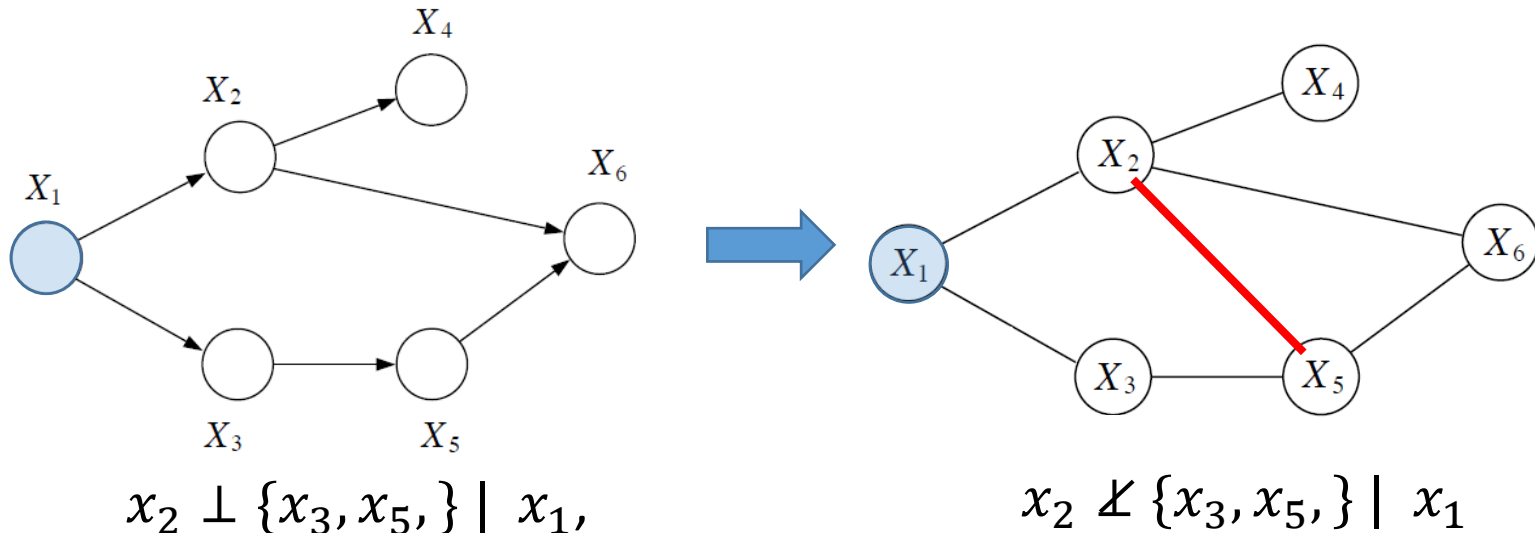


Image source: “An introduction to probabilistic graphical models”, Michael I. Jordan, 2002.

# Discriminative Vs Generative Models

- **Generative models:** Approaches that explicitly or implicitly model the distribution of inputs and outputs.
- Sampling from the distribution it is possible to generate synthetic data points in the input space.

**Likelihood:**  $p(\mathbf{x}|\mathcal{C}_k)$

- **Discriminative models:** Approaches that model the posterior probabilities directly.

**Posterior:**  $p(\mathcal{C}_k|\mathbf{x})$

# Conditional Random Fields

- A **CRF** or **discriminative random field**, is just a version of an MRF where all the clique potentials are **conditioned on input  $X$** :

$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{1}{Z(\mathbf{x}, \mathbf{w})} \prod_c \psi_c(\mathbf{y}_c|\mathbf{x}, \mathbf{w})$$

- We will usually assume a **log-linear representation** of the potentials:

$$\psi_c(\mathbf{y}_c|\mathbf{x}, \mathbf{w}) = \exp(\mathbf{w}_c^T \phi(\mathbf{x}, \mathbf{y}_c))$$

- where  $\phi(\mathbf{y}_c, \mathbf{x})$  is a **feature vector** derived from the **local set of labels  $Y_c$**  and **the global inputs  $X$** , and  $\mathbf{w}_c$  is the parameters.

# CRF vs MRF

## Advantages:

1. **No need to “waste resources”** modeling things that we always observe.

Focus our attention on **modeling what we care about**, i.e. the distribution of labels given the data.

2. We can make the potentials (or factors) of the model be **data-dependent**.

e.g. in natural language processing problems, we can make the **latent labels depend on global properties** of the sentence, such as which language it is written in.

# CRF vs MRF

Disadvantage:

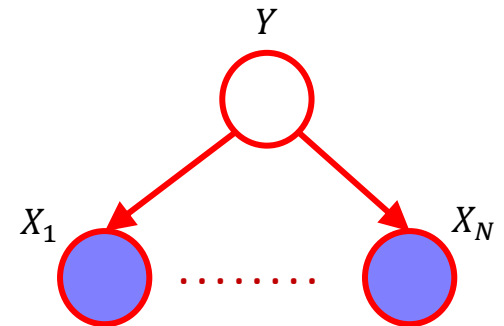
1. Not generative, this means we **cannot generate new samples**.
2. Learning is **slower** (more detail in the coming lectures).

# CRF: Logistic Regression

- Naives Bayes classifier we have seen previously **ignores dependencies** among the observations  $X_1, \dots, X_N$ .

Joint distribution:

$$p(x_1, \dots, x_N, Y) = \prod_{n=1}^N p(x_n | Y)$$



**Predict** with Bayes rule:

$$p(Y = 1 | x_1, \dots, x_N) = \frac{p(y = 1) \prod_{n=1}^N p(x_n | Y = 1)}{\sum_{y=\{0,1\}} p(Y = y) \prod_{n=1}^N p(x_n | Y = y)}$$

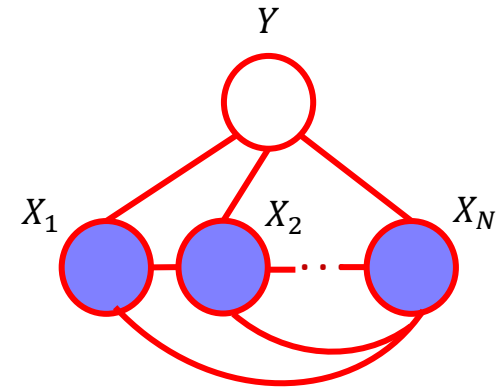
# CRF: Logistic Regression

- In CRF, we **model the full dependencies** among the observations  $X_1, \dots, X_N$ .

$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{1}{Z(\mathbf{x}, \mathbf{w})} \psi(\mathbf{y} | \mathbf{x}, \mathbf{w})$$

$$= \frac{\exp(\phi(\mathbf{y}, \mathbf{x})^\top \mathbf{w})}{\sum_{\mathbf{y}} \exp(\phi(\mathbf{y}, \mathbf{x})^\top \mathbf{w})}$$

$$= \frac{\exp(\sum_n \phi_n(\mathbf{y}, \mathbf{x}) w_n)}{\sum_{\mathbf{y}} \exp(\sum_n \phi_n(\mathbf{y}, \mathbf{x}) w_n)} \quad (\text{Log-linear model})$$

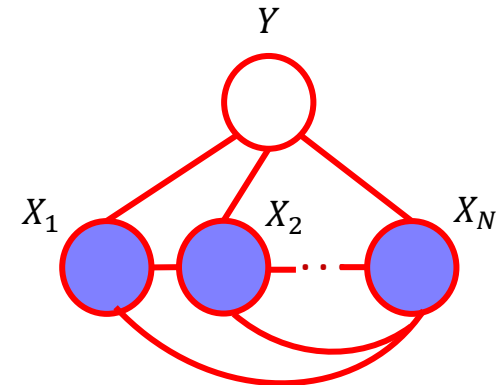


# CRF: Logistic Regression

- In CRF, we **model the full dependencies** among the observations  $X_1, \dots, X_N$ .

For  $y \in \{0,1\}$ , e.g., email spam classifier, we can define:

$$\begin{aligned}\phi_n(y = 1, \mathbf{x})w_n &= \mathbb{I}(y^+)x_nw_n^+, \text{ and} \\ \phi_n(y = 0, \mathbf{x})w_n &= \mathbb{I}(y^-)x_nw_n^-\end{aligned}$$



We get:

$$p(y^+ | \mathbf{x}) = \frac{\exp(\mathbf{w}^{+\top} \mathbf{x})}{\exp \mathbf{w}^{+\top} \mathbf{x} + \exp \mathbf{w}^{-\top} \mathbf{x}} = \frac{1}{1 + \exp(-\mathbf{w}'^{\top} \mathbf{x})}$$

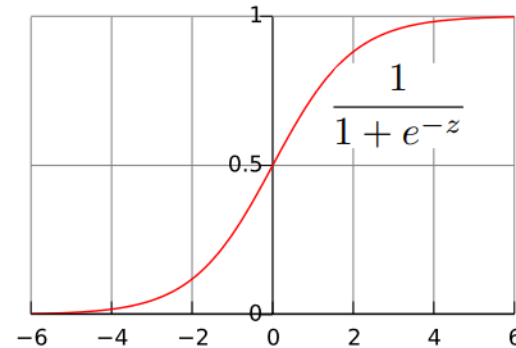
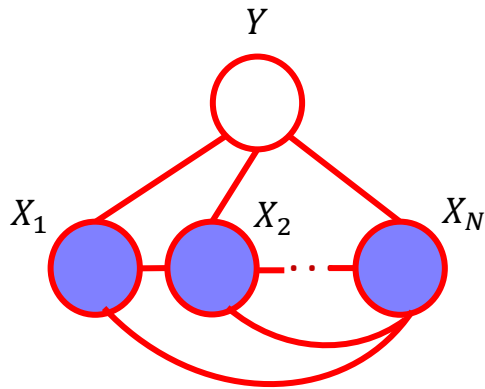
Logistic Regression!

where  $\mathbf{w}' = \mathbf{w}^+ - \mathbf{w}^-$ .

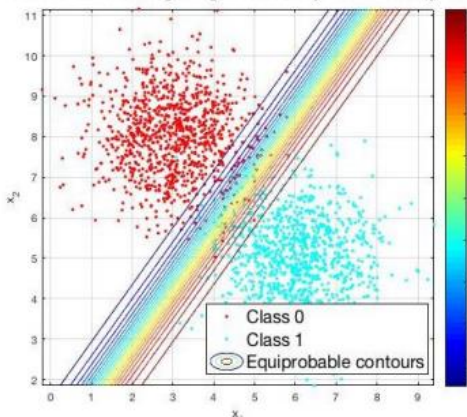
**Note:** there's just 1 set of parameters, i.e.,  $\mathbf{w}'$ , we don't care about modeling the respective  $\mathbf{w}^+$  and  $\mathbf{w}^-$ .



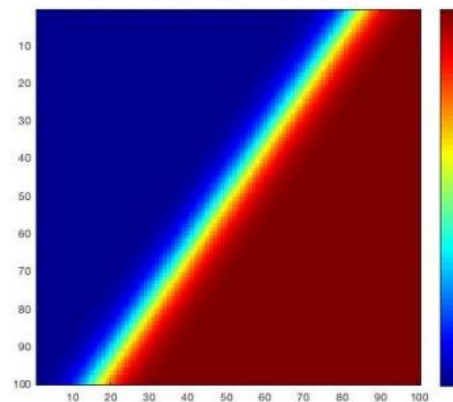
# CRF: Logistic Regression



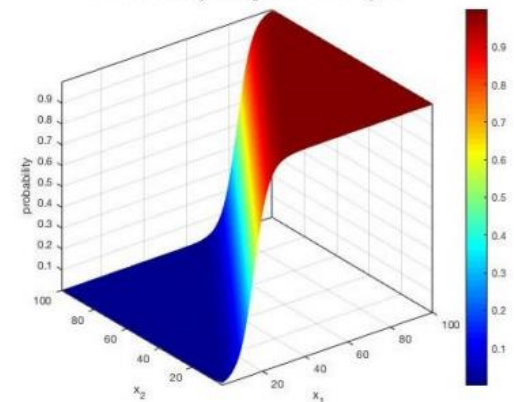
Contours of equal probability defined by  $\theta$



Probability map defined by  $\theta$



Probability map defined by  $\theta$

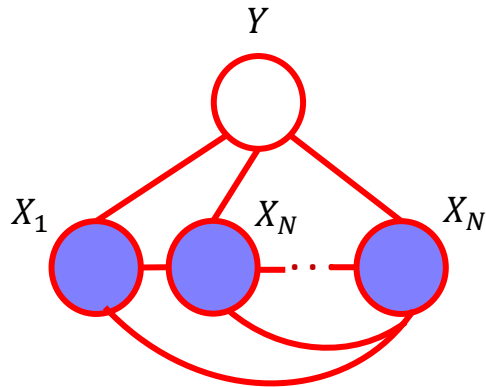


- **Decision boundary:**  $p(y = 1 \mid \mathbf{x}, \mathbf{w}) > 0.5$  is linear in  $\mathbf{x}$ .

Image source:

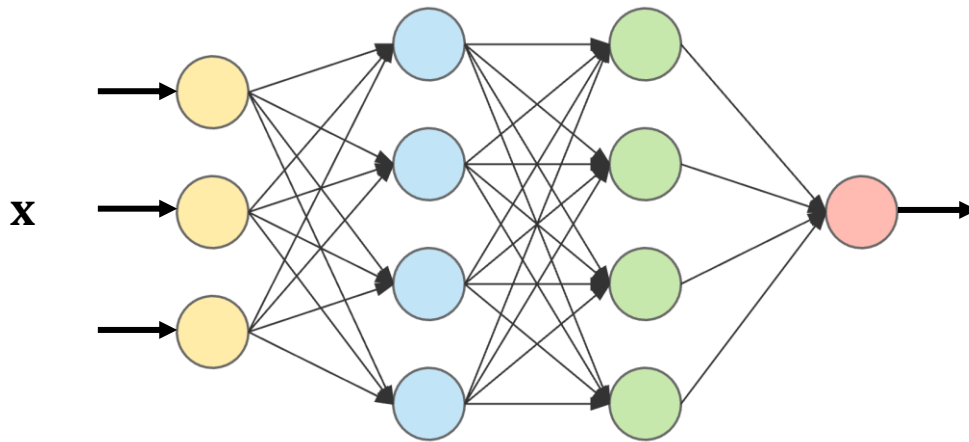
[https://deepgenerativemodels.github.io/assets/slides/cs236\\_lecture2.pdf](https://deepgenerativemodels.github.io/assets/slides/cs236_lecture2.pdf)

# Classification Deep Network



$$p(y^+ | \mathbf{x}) = \frac{1}{1 + \exp \mathbf{w}'^T \mathbf{x}}$$

- The CRF becomes a deep network when  $\mathbf{w}'^T \mathbf{x}$  is a **nonlinear function**  $f(\mathbf{x}; \mathbf{w})$ !



$\mathbf{w}$ : weight parameters

$$p(y^+ | \mathbf{x}) = \frac{1}{1 + \exp f(\mathbf{x}; \mathbf{w})}$$

$$= \sigma(f(\mathbf{x}; \mathbf{w})),$$

where  $\sigma(\cdot)$  is the sigmoid function.

**Note:**  $\sigma(\cdot)$  is the **softmax function** for multiclass classification.

# CRF: Linear Chain CRF

## Models for sequential data

Hidden Markov model:

$$p(\mathbf{x}, \mathbf{y} | \mathbf{w}) = \prod_{t=1}^T p(y_t | y_{t-1}, \mathbf{w}) \underbrace{p(\mathbf{x}_t | y_t, \mathbf{w})}_{\text{Likelihood, i.e. generative}}$$

Likelihood, i.e. generative

Chain structure MRF:

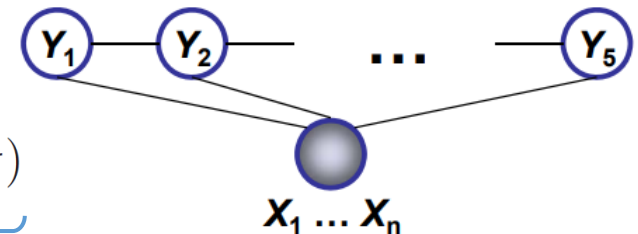
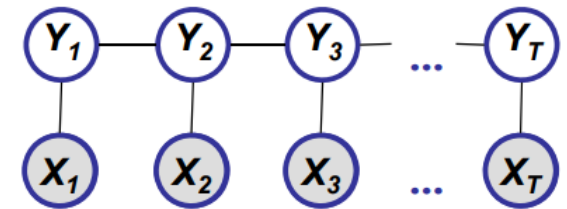
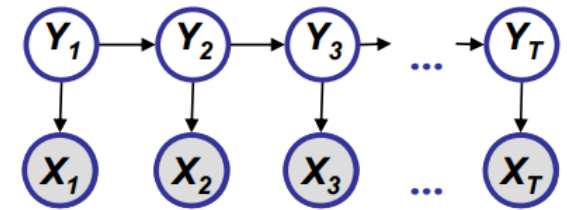
$$p(\mathbf{x}, \mathbf{y} | \mathbf{w}) = \prod_{t=1}^T \underbrace{\psi(y_t; \mathbf{x}_t, \mathbf{w})}_{\text{Likelihood, i.e. generative}} \prod_{t=1}^{T-1} \psi(y_t, y_{t+1}; \mathbf{w})$$

Likelihood, i.e. generative

Chain structure CRF:

$$p(\mathbf{y} | \mathbf{x}, \mathbf{w}) = \frac{1}{Z(\mathbf{x}, \mathbf{w})} \prod_{t=1}^T \underbrace{\psi(y_t; \mathbf{x}, \mathbf{w})}_{\text{Posterior, i.e. discriminative}} \prod_{t=1}^{T-1} \underbrace{\psi(y_t, y_{t+1}; \mathbf{x}, \mathbf{w})}_{\text{Posterior, i.e. discriminative}}$$

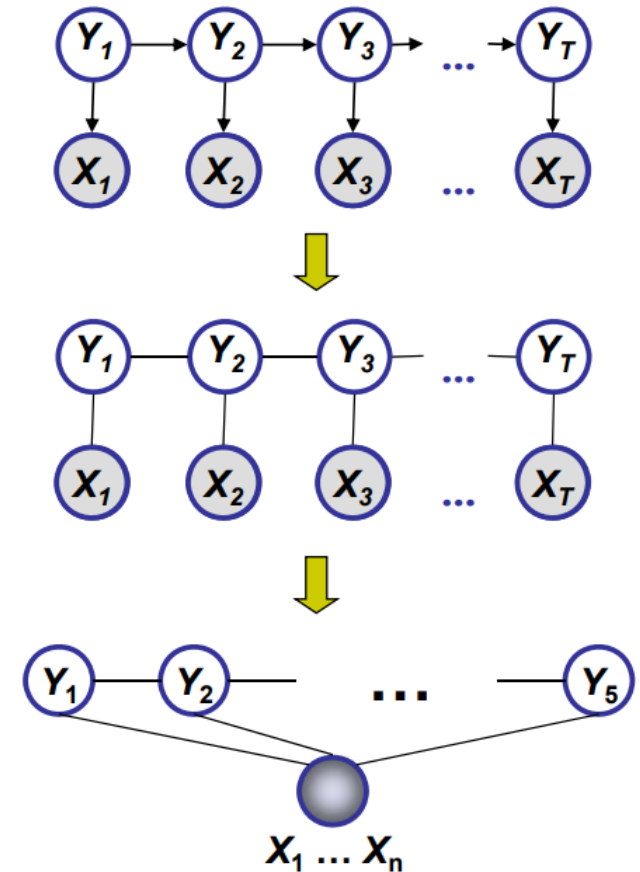
Posterior, i.e. discriminative



# CRF: Linear Chain CRF

## Models for sequential data

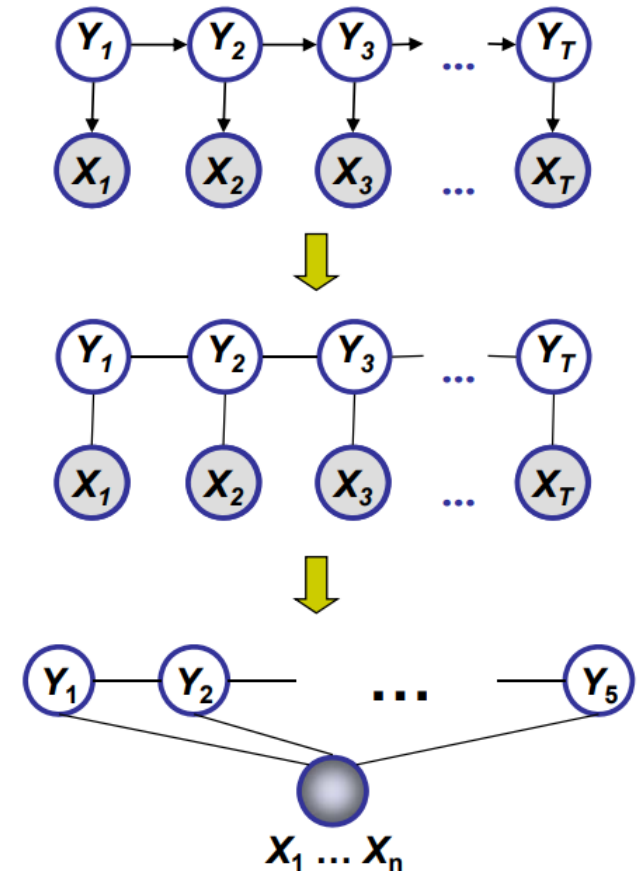
- HMM and MRF suffer from the **label bias problem**.
- Local features at time  $t$  do not influence states prior to time  $t$ .
- $X_t$  is **d-separated** from all other nodes at  $Y_t$  thus blocking the information flow.



# CRF: Linear Chain CRF

## Models for sequential data

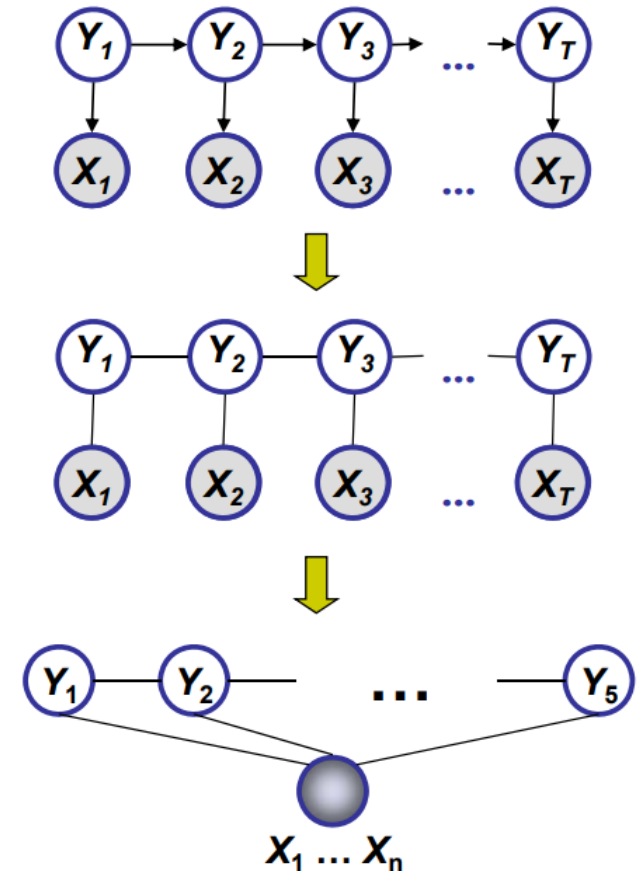
- Consider the **part of speech (POS) tagging** task.
- Suppose we see the word “banks”.
- This could be a **verb** (as in “he banks at DBS”), or a **noun** (as in “the river banks were overflowing”).
- Locally** the POS tag for the word is **ambiguous**.



# CRF: Linear Chain CRF

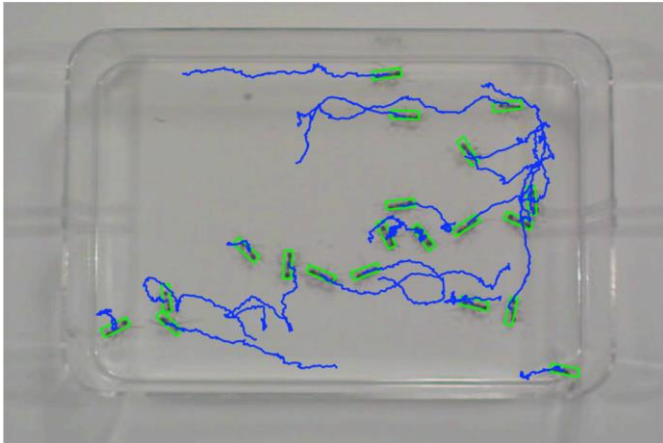
## Models for sequential data

- Suppose that later in the sentence, we see the word “fishing”.
- This gives us enough context to infer that the sense of “banks” is “river banks”.
- However, in HMM and MRF the “fishing” evidence is **d-separated**.
- Problem is alleviated in CRF.



# Example of MRF

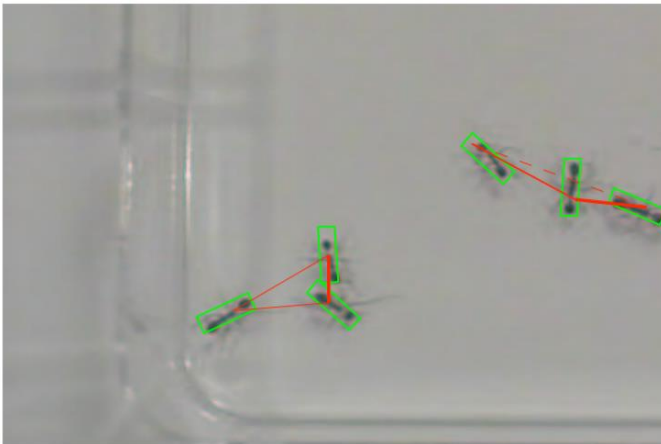
- Objective: Track every ant in the video



$$P(X_t|Z^t) = kP(Z_t|X_t) \int_{X_{t-1}} P(X_t|X_{t-1})P(X_{t-1}|Z^{t-1})$$

$x_t$ : current locations of the ants

$z_t$ : tracklets output (green boxes) from an ant detector



$$P(X_t|X_{t-1}) \propto \prod_i P(X_{it}|X_{i(t-1)}) \prod_{ij \in E} \psi(X_{it}, X_{jt})$$

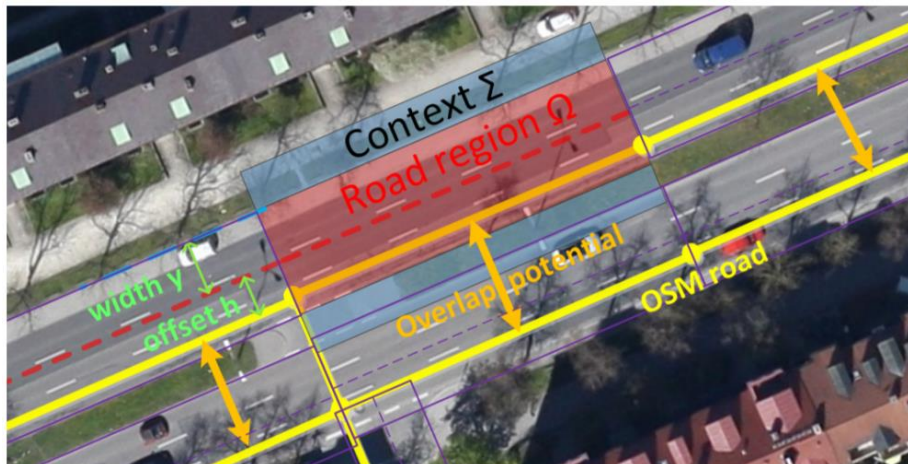
MRF used to model interactions between the ants.

Zia Khan, Tucker Balch, and Frank Dellaert, **An MCMC-based Particle Filter for Tracking Multiple Interacting Targets**, ECCV 2014



# Example of CRF

- **Objective:** To infer the location of the road-segment centerlines as well as their width.



Road classifier:  $\phi_{road}$

Appearance:  $\phi_{ap} = [\phi_{edge}, \phi_{hom}, \phi_{context}]$

Distance to edges:  $\phi_{edge}$

Homogeneity of the region:  $\phi_{hom}$

Appearance context:  $\phi_{context}$

$$\{\mathbf{h}^*, \mathbf{y}^*\} = \operatorname{argmin}_{\mathbf{h}, \mathbf{y}} E(\mathbf{h}, \mathbf{y})$$

$$\begin{aligned} E(\mathbf{h}, \mathbf{y}) = & \sum_{j=1}^L \sum_{i=1}^{l_j} \mathbf{w}_{road}^T \phi_{road}(h_i^j, y_i^j, \mathbf{x}) \\ & + \sum_{j=1}^L \sum_{i=1}^{l_j} \mathbf{w}_{ap}^T \phi_{ap}(h_i^j, y_i^j, \mathbf{x}) + \sum_{j=1}^L \sum_{i=1}^{l_j} \mathbf{w}_{car}^T \phi_{car}(h_i^j, y_i^j, \mathbf{x}) \\ & + \sum_{j=1}^L \sum_{i=1}^{l_j-1} \mathbf{w}_{sm}^T \phi_{sm}(h_i^j, y_i^j, h_{i+1}^j, y_{i+1}^j) \\ & + \sum_{i,j,k,m \in P} \phi_{ol}(h_i^j, y_i^j, h_k^m, y_k^m) \end{aligned}$$

Smoothness:  $\phi_{sm}$

Overlap:  $\phi_{ol}$

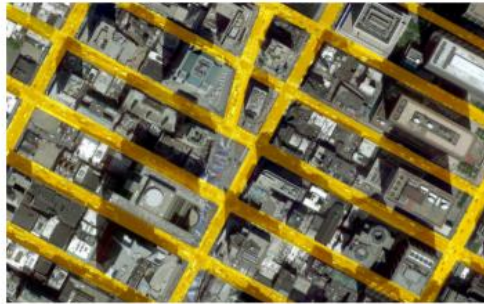
G. Mattyus, S. Wang, S. Fidler and R. Urtasun, **Enhancing World Maps by Parsing Aerial Images**, ICCV 2015



# Example of CRF



(Toronto: Pearson Airport)



(NYC: Times square)



(Nairobi, Kenya)\*



(Manila, Philippines)



(Mexico City)



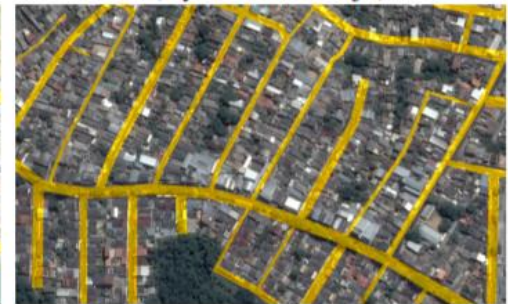
(Kyoto: Kinkakuji )



(Sydney: At Harbour bridge)



(St. Moritz, Switzerland)



(Manaus, Brazil)\*

G. Mattyus, S. Wang, S. Fidler and R. Urtasun, **Enhancing World Maps by Parsing Aerial Images**, ICCV 2015

# Summary

- You have learned how to:
  1. Explain the concepts of **Markov properties (global, local and pairwise)** and use it to find all conditional independences in an UGM.
  2. Use **clique potential functions** to parameterize a Markov Random Field, i.e. to represent the joint distribution with clique potential functions.
  3. Describe the differences and similarities between a **Markov Random Field** and **Conditional Random Field**.