



Royal Institute of  
Technology

# STATISTICAL METHODS IN CS, CH 10

## Lecture 5

# LAST LECTURE

- ★ The Bayesian approach
- ★ How to analyse and use the posterior
- ★ Empirical Bayes
- ★ Cancer incidents - sharing parameters

# THIS LECTURE

- ★ DGMs
- ★ Basic example
- ★ Computing the likelihood
- ★ Factorizing (decomposing) the likelihood
- ★ A decomposable prior
- ★ Categorical CPDs - easy to get a posterior

# REPRESENTING AND WORKING WITH DISTRIBUTIONS

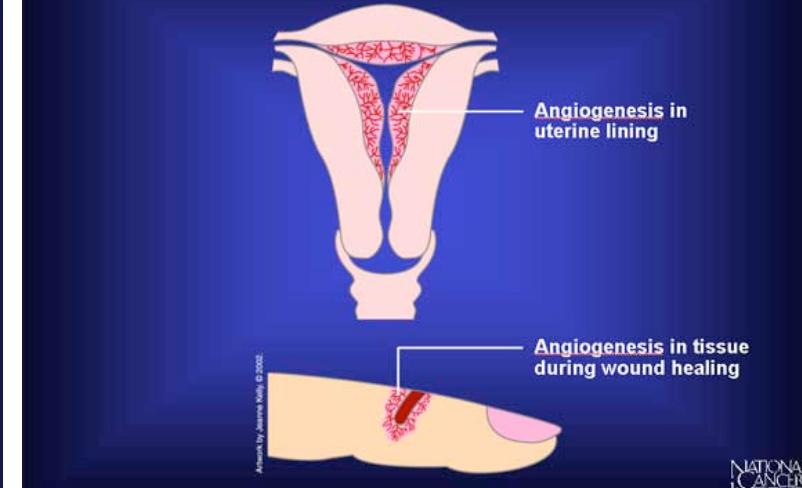
- ★ For all but the smallest  $n$ , the explicit representation of the joint distribution is *unmanageable from every perspective*.
  - Computationally, it is very *expensive to manipulate* and generally *too large to store in memory*.
  - Cognitively, it is *impossible to acquire so many numbers* from a human expert; moreover, the numbers are very small and *do not correspond to events that people can reasonably contemplate*.
  - Statistically, if we want to learn the distribution from data, we would *need ridiculously large amounts of data to estimate* this many parameters robustly.
- ★ These problems were the *main barrier* to the adoption of probabilistic methods for expert systems *until the development of the methodologies we now will consider*.

# Normal Angiogenesis in Children



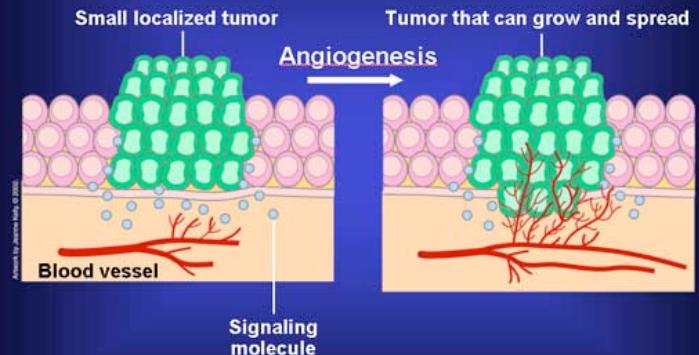
Artwork by Jeannine Kelly © 2002.

# Normal Angiogenesis in Adults



NATIONAL CANCER INSTITUTE

## What Is Tumor Angiogenesis?



NATIONAL CANCER INSTITUTE

ABERRATION DEPENDENCIES -  
EX. ANGIOGENESIS

# SOMATIC EVOLUTION

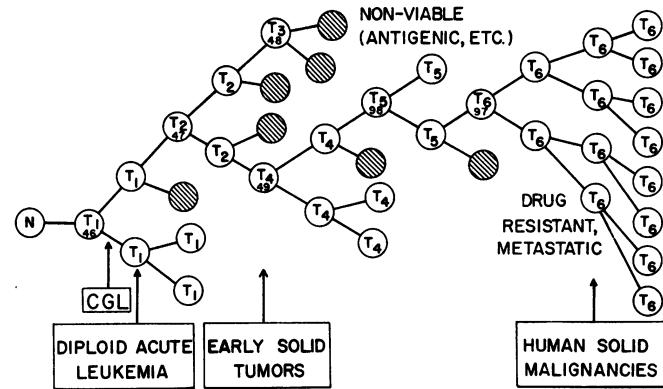
## The Clonal Evolution of Tumor Cell Populations

Acquired genetic lability permits stepwise selection of variant sublines and underlies tumor progression.

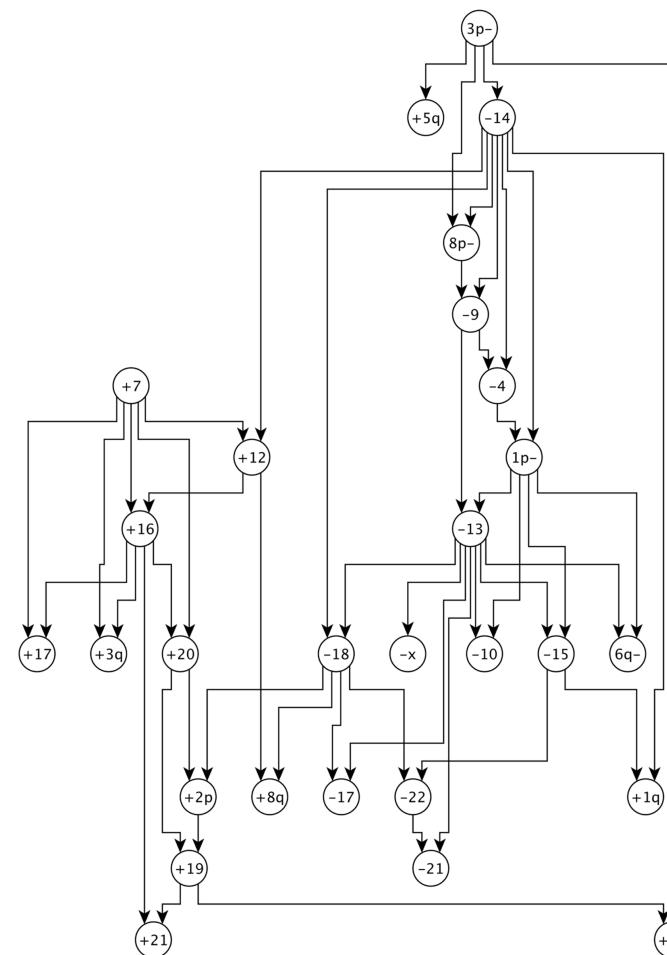
Peter C. Nowell

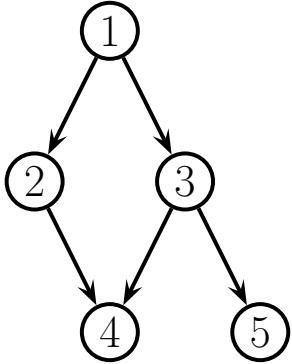
The author is professor of pathology, School of Medicine, University of Pennsylvania, Philadelphia 19174.

1 OCTOBER 1976 SCIENCE, VOL. 194



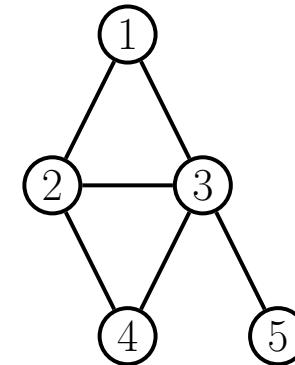
Oncogenetic network





Directed graphical model

- DAG
- vertices r.v.s
- equipped with local CPDs
- allows causal like dependencies



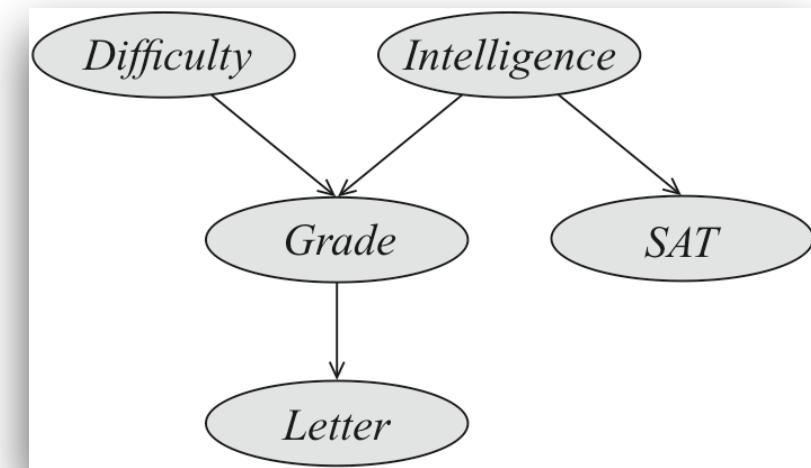
Undirected graphical model - Markov

Random Fields

- graph
- vertices r.v.s
- equipped with local “factors”

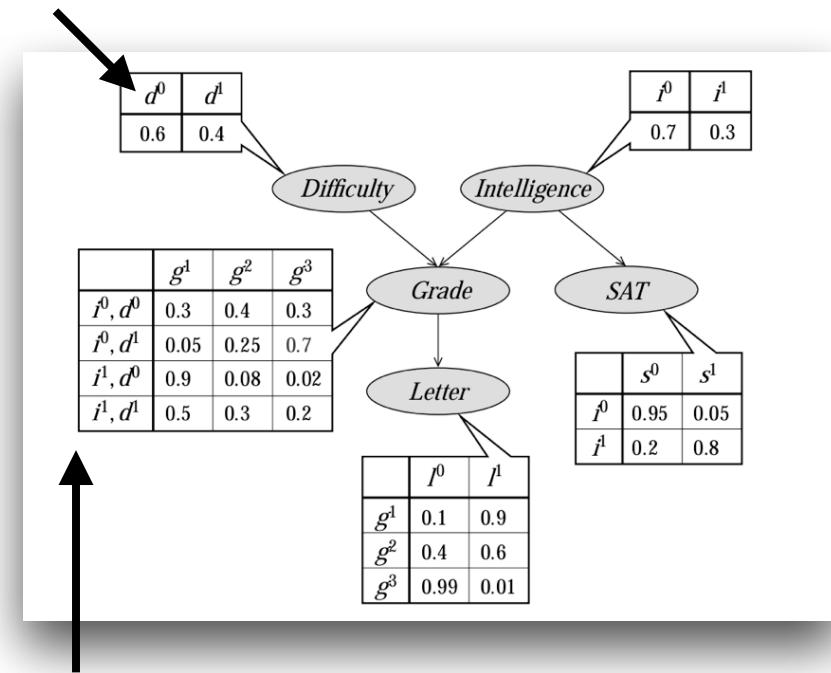
# GRAPHICAL MODELS

# DGM - GRAPH AND CPDS VS JOINT



$P(D, I, G, S, L)$

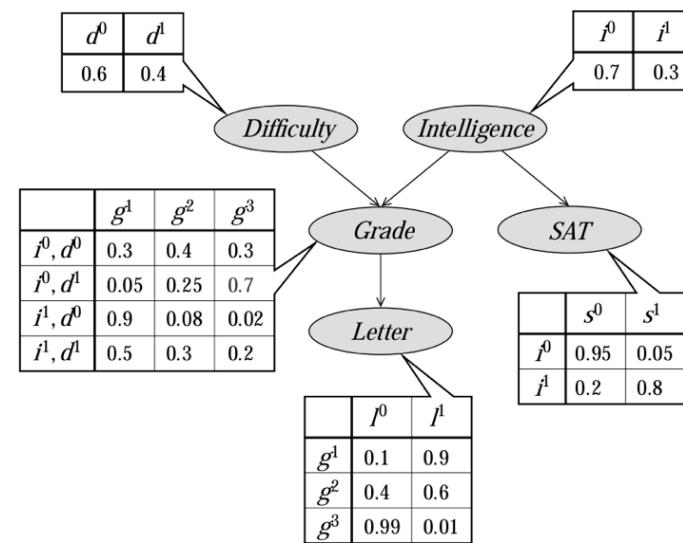
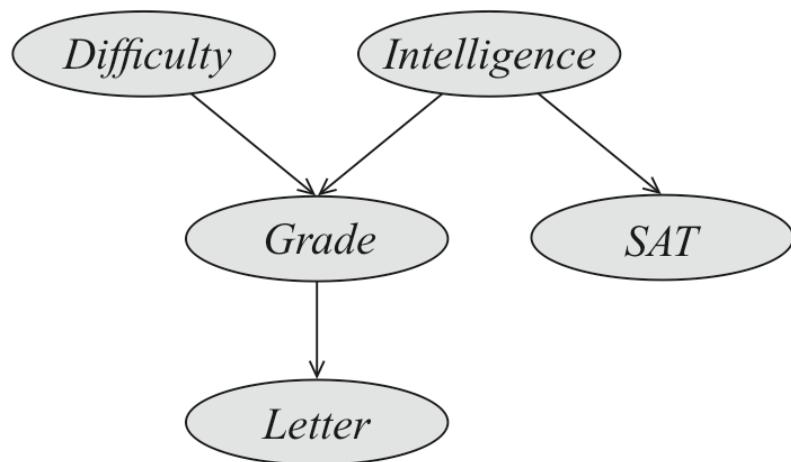
d has value 0



CPD

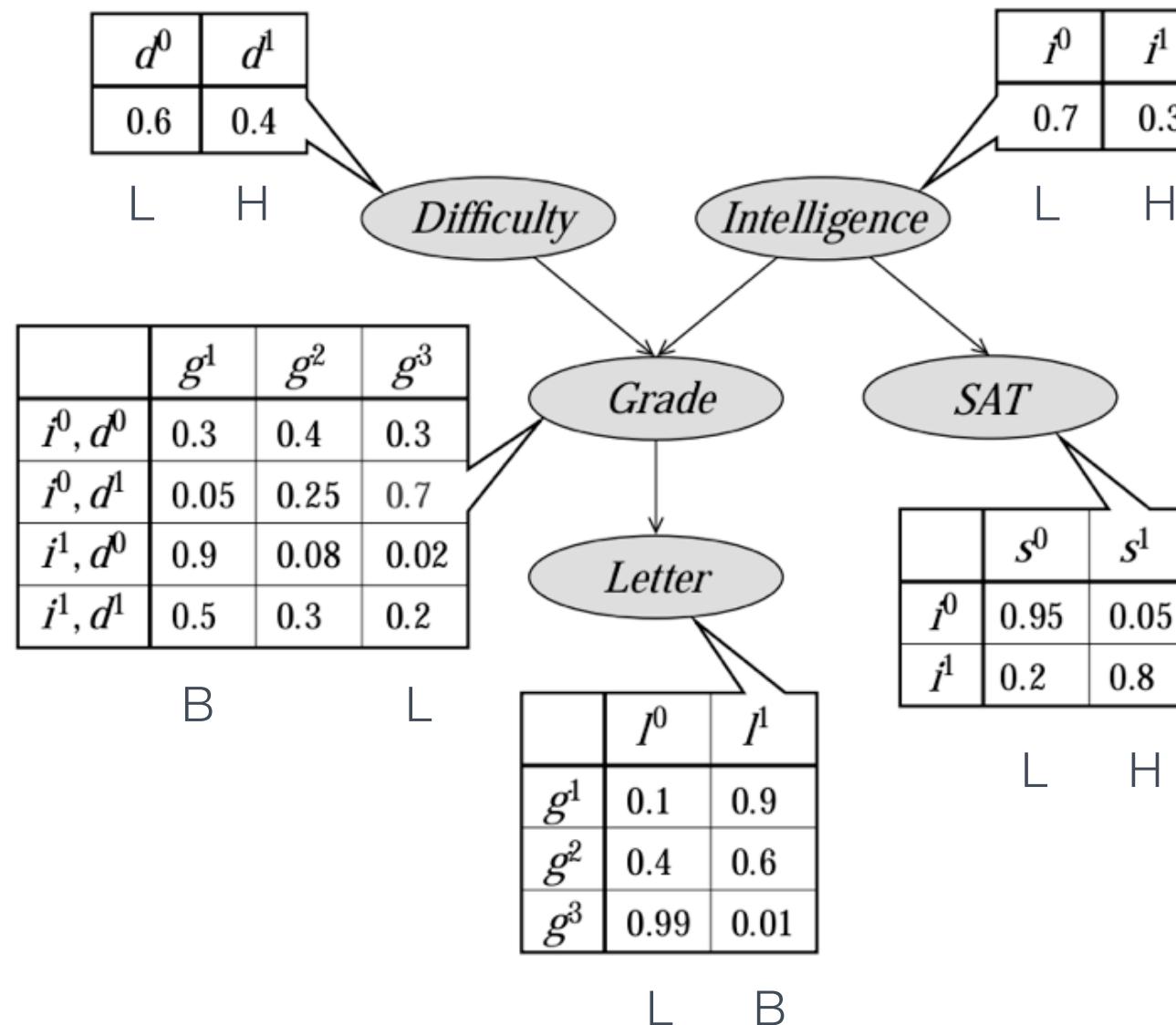
- ★ CPT - table, i.e., categorical
- ★ Gaussian

# THREE LEVELS OF COMPUTATIONAL PROBLEMS



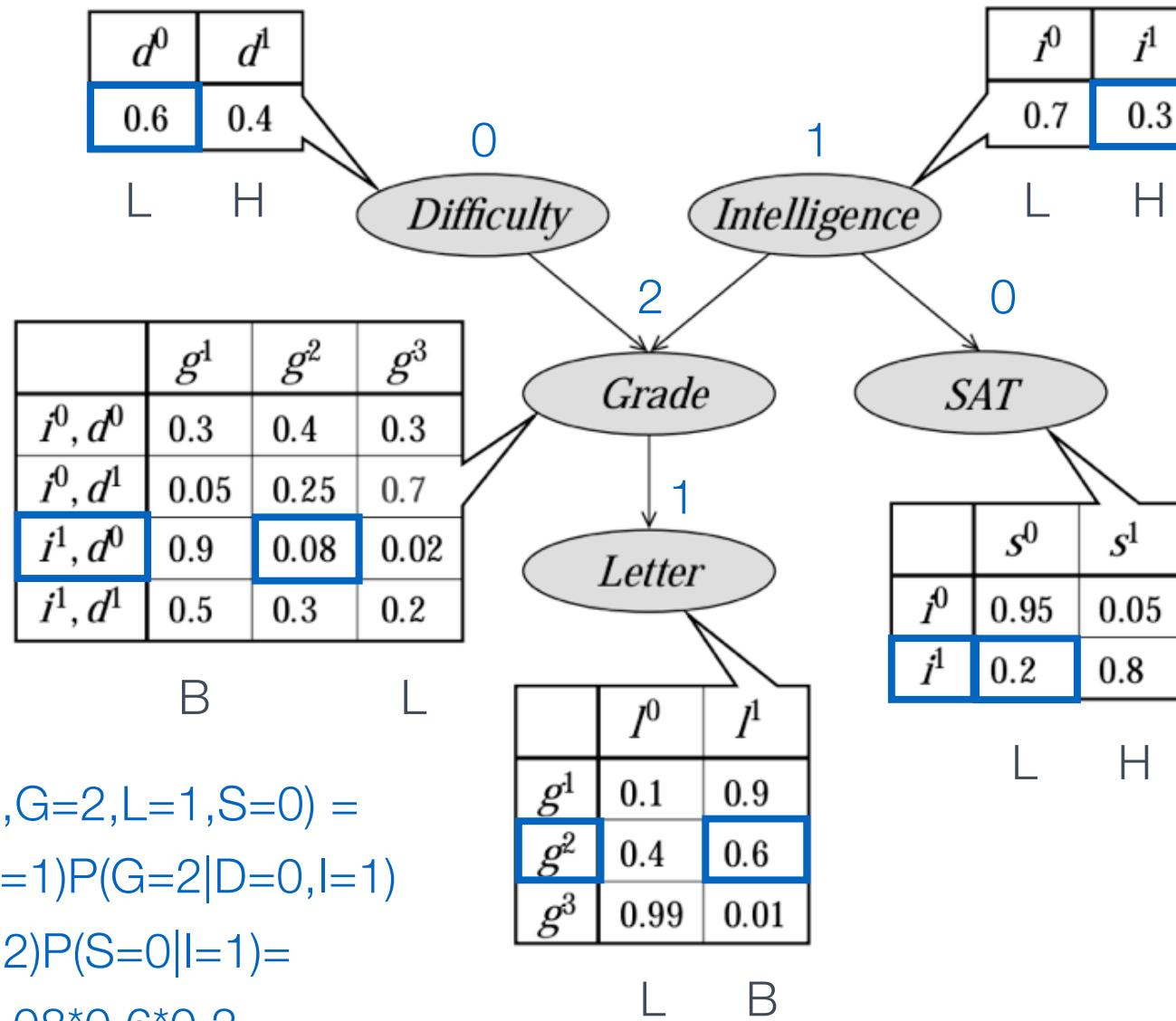
- Inference: given  $G$  and  $\theta$ , compute probabilities or marginalize
- Parameter learning: given  $G$  and  $D$ , learn  $\theta$
- Structure learning: given  $D$ , learn  $G$  and  $\theta$

# EXTENDED STUDENT-CATEGORICAL CPDS



B - better  
 H - higher  
 L - less

# COMPLETE DATA



$$P(D=0, I=1, G=2, L=1, S=0) =$$

$$P(D=0)P(I=1)P(G=2|D=0, I=1)$$

$$P(L=1|G=2)P(S=0|I=1) =$$

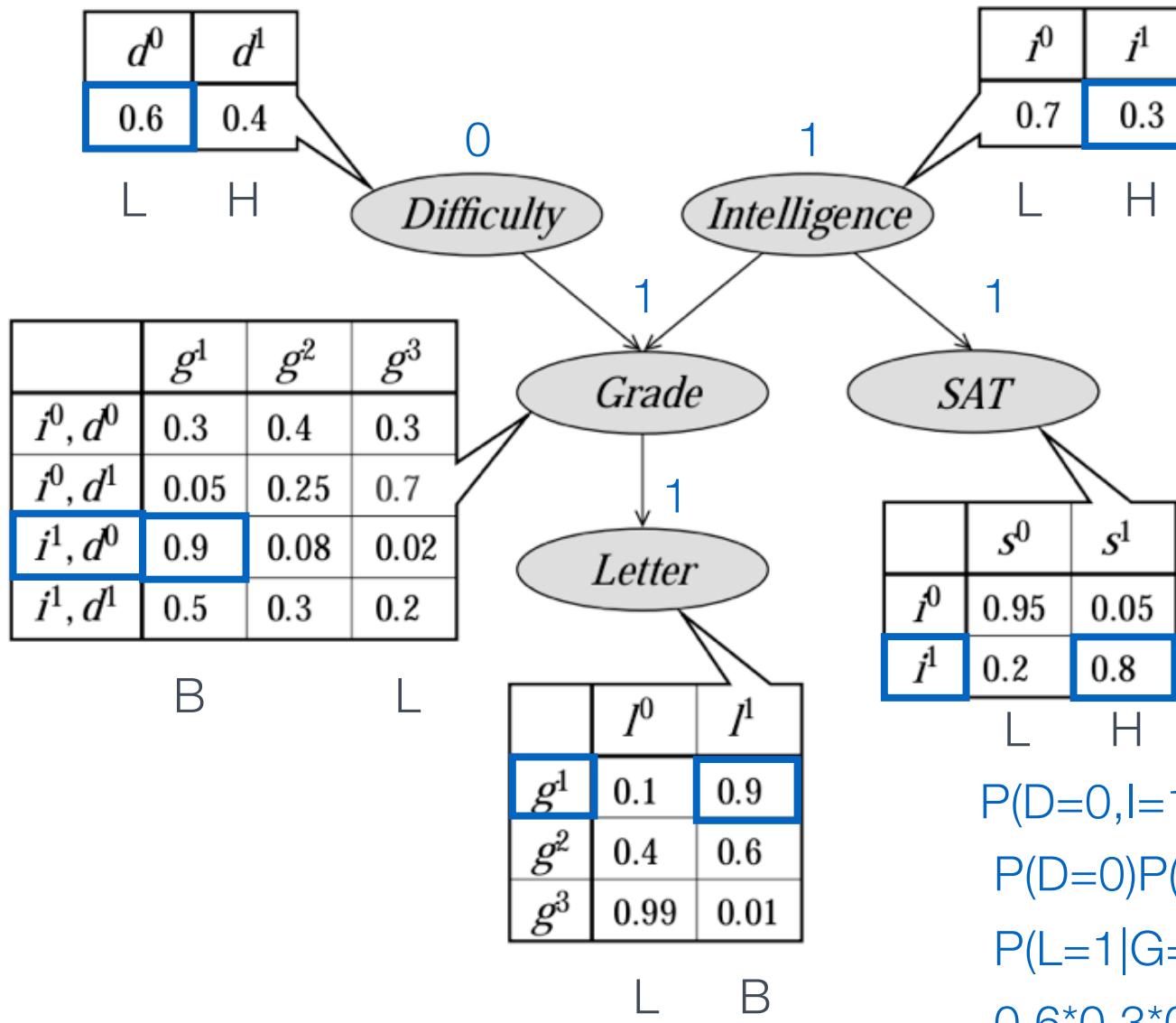
$$0.6 * 0.3 * 0.08 * 0.6 * 0.2$$

B - better

H - higher

L - less

# EXTENDED STUDENT EXAMPLE



$$\begin{aligned}
 P(D=0, I=1, G=1, L=1, S=1) &= \\
 P(D=0)P(I=1)P(G=1|D=0, I=1) & \\
 P(L=1|G=1)P(S=1|I=1)= & \\
 0.6 * 0.3 * 0.9 * 0.9 * 0.8
 \end{aligned}$$

# INFERENCE – THE CHAIN RULE

$$p(\underbrace{\boldsymbol{x}_{[V]}}_{\boldsymbol{x}_1, \dots, \boldsymbol{x}_V}) = p(\boldsymbol{x}_1)p(\boldsymbol{x}_2|\boldsymbol{x}_1)p(\boldsymbol{x}_3|\boldsymbol{x}_1, \boldsymbol{x}_2) \cdots p(\boldsymbol{x}_V|\boldsymbol{x}_{[V-1]})$$

- ★ Assuming binary r.v.,  $p(X_V | X_{[V-1]})$  has  $2^{V-1}$  parameters
- ★ Total # parameters  $\sum_{1 \leq i \leq V} 2^{i-1} = 2^V - 1$

# EX. WHERE IND. OBVIOUSLY FACILITATES

$$p(\underbrace{\mathbf{x}_{[V]}}_{\mathbf{x}_1, \dots, \mathbf{x}_V}) = p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1)p(\mathbf{x}_3|\mathbf{x}_1, \mathbf{x}_2) \cdots p(\mathbf{x}_V|\mathbf{x}_{[V-1]})$$

★ Assume first order Markov property  $\mathbf{x}_t \perp \mathbf{x}_{[t-2]} | \mathbf{x}_{t-1}$

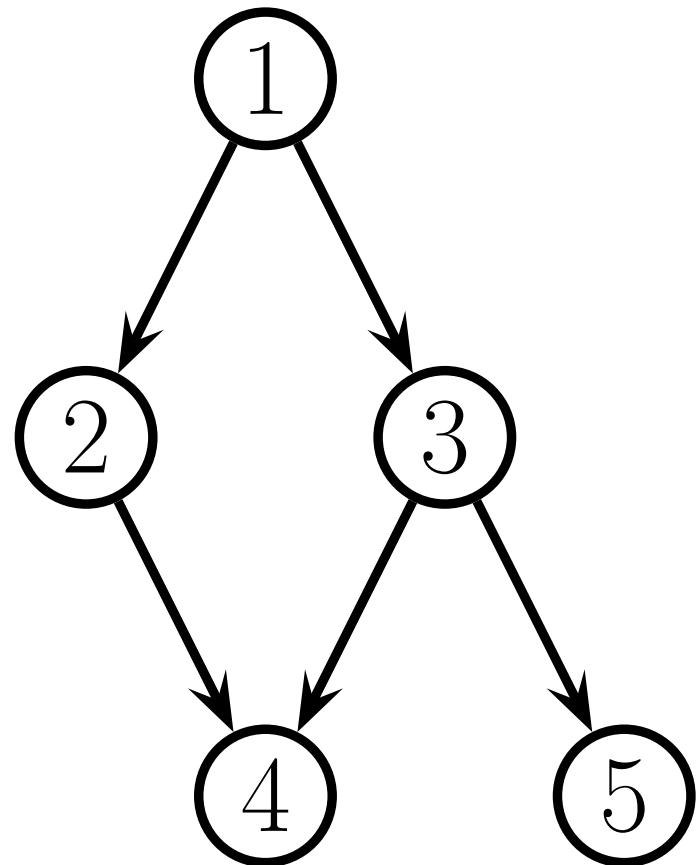
i.e., if time ordered, future independent of past given present

★ Then  $p(\mathbf{x}_{[V]}) = p(\mathbf{x}_1) \prod_{t=1}^{V-1} p(\mathbf{x}_{t+1}|\mathbf{x}_t)$

# FACTORIZATION OVER G

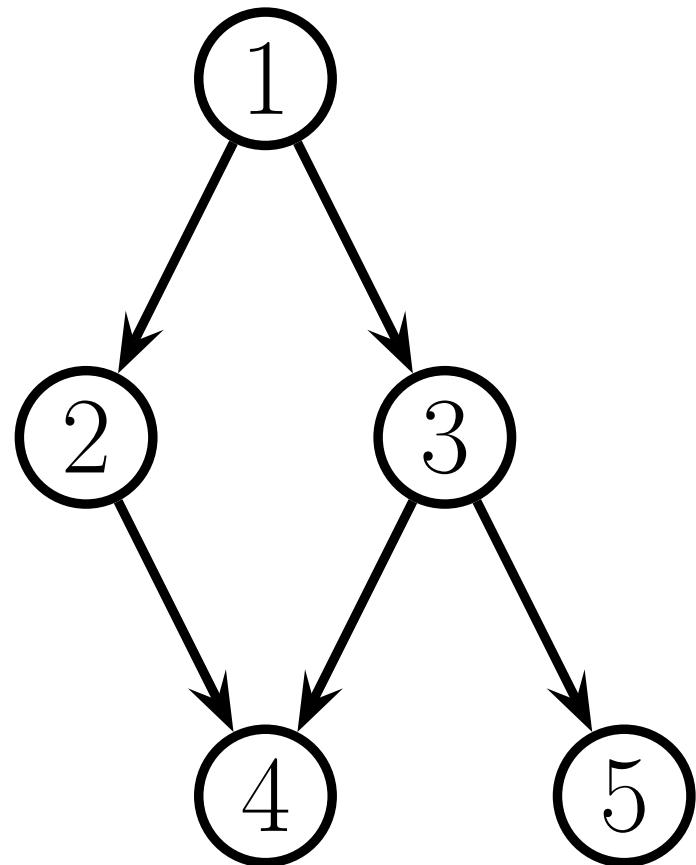
$$p(x_1, \dots, x_N) = \prod_{n=1}^N p(x_n | \mathbf{x}_{\text{pa}(x_n)})$$

p can be factorized over G if it can be expressed as above



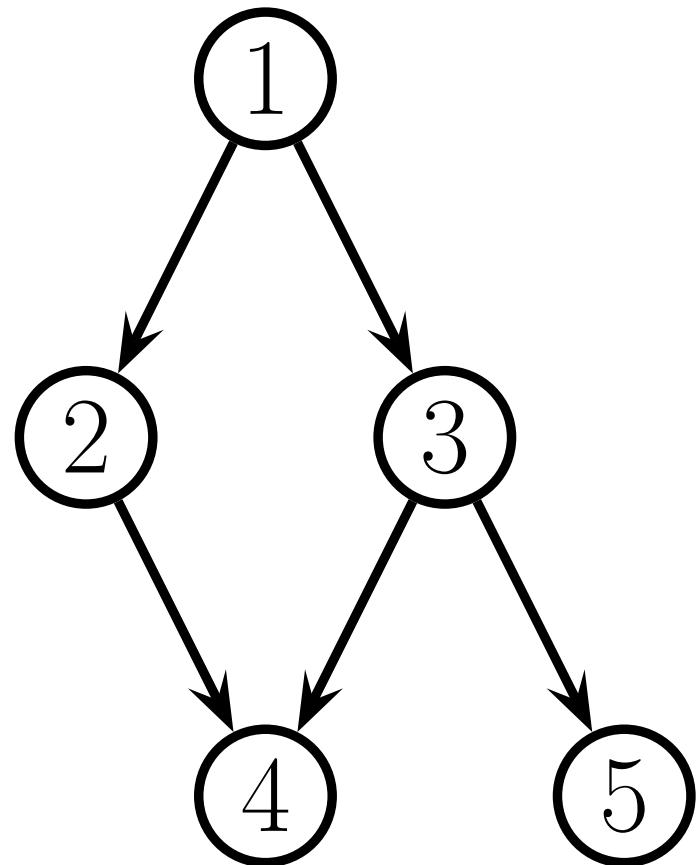
# TERMINOLOGY

- ★ Parent
- ★ Child
- ★ Family
- ★ Root
- ★ Leaf
- ★ Neighbor



# TERMINOLOGY

- ★ Degree (in and out)
- ★ Cycle (directed or not)
- ★ Directed Acyclic Graph (DAG)
- ★ Topological order (parents < child)
- ★ Path (directed or not)
- ★ Ancestors



# TERMINOLOGY

- ★ Tree
- ★ Polytree – directed tree with multiple parents for some vertices
- ★ Forest
- ★ Subgraph
- ★ Clique
- ★ Maximal clique

# CAT — NOTATION

- ★ For a  $v \in [M]$ ,

$$\text{values } k \in [K_v] \quad \xrightarrow{\text{Cartesian product}} \quad \text{combined values } c \in C_v = \prod_{s \in \text{pa}(v)} [K_s]$$

-  Cat CPDs

where  $P(x_v | x_{\text{pa}(v)} = c) = \text{Cat}(\theta_{vc})$

and  $\theta_{vck} = P(x_v = k | x_{\text{pa}(v)} = c)$

# NOTATION EXAMPLE

$$K_d = \{0, 1\}$$

$d^0$	$d^1$
0.6	0.4

d

*Difficulty*

$$K_i = \{0, 1\}$$

$i^0$	$i^1$
0.7	0.3

i

*Intelligence*

$$\begin{aligned} \text{Cat}(\theta_{g(0,0)}) \\ \text{Cat}(\theta_{g(0,1)}) \\ \text{Cat}(\theta_{g(1,0)}) \\ \text{Cat}(\theta_{g(1,1)}) \end{aligned}$$

	$g^1$	$g^2$	$g^3$
$i^0, d^0$	0.3	0.4	0.3
$i^0, d^1$	0.05	0.25	0.7
$i^1, d^0$	0.9	0.08	0.02
$i^1, d^1$	0.5	0.3	0.2

g

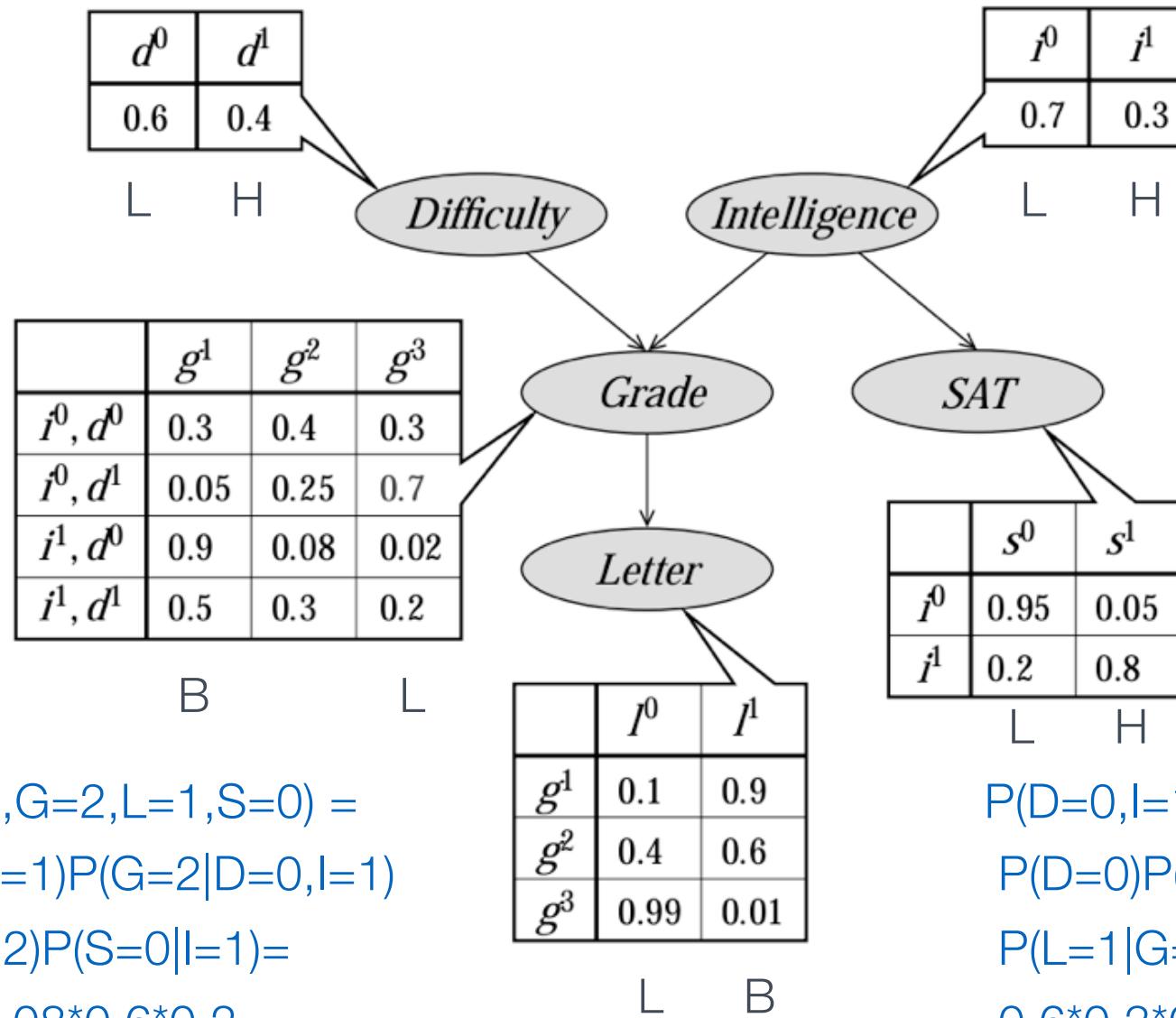
*Grade*

$$K_g = \{1, 2, 3\}$$

$$\theta_{g(1,1)2}$$

$$C_g = \prod_{s \in \text{pa}(g)} [K_s] = K_d \times K_i = \{\langle 0, 0 \rangle, \langle 0, 1 \rangle, \langle 1, 0 \rangle, \langle 1, 1 \rangle\}$$

# EXTENDED STUDENT EXAMPLE



$$P(D=0, I=1, G=2, L=1, S=0) =$$

$$P(D=0)P(I=1)P(G=2|D=0, I=1)$$

$$P(L=1|G=2)P(S=0|I=1)=$$

$$0.6 * 0.3 * 0.08 * 0.6 * 0.2$$

$$P(D=0, I=1, G=1, L=1, S=1) =$$

$$P(D=0)P(I=1)P(G=1|D=0, I=1)$$

$$P(L=1|G=1)P(S=1|I=1)=$$

$$0.6 * 0.3 * 0.9 * 0.9 * 0.8$$

# EXTENDED STUDENT EXAMPLE

$$\begin{aligned} & P(D=0, I=1, G=2, L=1, S=0) \quad P(D=0, I=1, G=1, L=1, S=0) \\ & = P(D=0)P(I=1)P(G=1|D=0, I=1) \quad P(L=1|G=1)P(S=1|I=1) \\ & \quad P(D=0)P(I=1)P(G=2|D=0, I=1)P(L=1|G=2)P(S=0|I=1) \\ & = P(D=0)^2P(I=1)^2P(G=1|D=0, I=1)P(G=2|D=0, I=1) \\ & \quad P(L=1|G=1)P(L=1|G=2)P(S=1|I=1)P(S=0|I=1) \end{aligned}$$

$$\begin{aligned} P(D=0, I=1, G=2, L=1, S=0) &= \\ P(D=0)P(I=1)P(G=2|D=0, I=1) & \\ P(L=1|G=2)P(S=0|I=1) &= \\ 0.6 * 0.3 * 0.08 * 0.6 * 0.2 & \end{aligned}$$

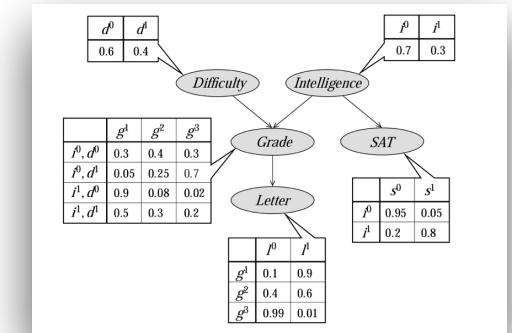
$$\begin{aligned} P(D=0, I=1, G=1, L=1, S=1) &= \\ P(D=0)P(I=1)P(G=1|D=0, I=1) & \\ P(L=1|G=1)P(S=1|I=1) &= \\ 0.6 * 0.3 * 0.9 * 0.9 * 0.8 & \end{aligned}$$

# THE LIKELIHOOD FACTORIZES

- ★ Complete data

$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

$$\mathbf{x}_n = \{\mathbf{x}_{n1}, \dots, \mathbf{x}_{nV}\}$$



- ★ Likelihood

$$\begin{aligned}
 p(\mathcal{D}|\boldsymbol{\theta}) &= \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{v=1}^V p(\mathbf{x}_{nv}|\mathbf{x}_{n,\text{pa}(v)}, \boldsymbol{\theta}) \\
 &= \prod_{v=1}^V \prod_{n=1}^N p(\mathbf{x}_{nv}|\mathbf{x}_{n,\text{pa}(v)}, \boldsymbol{\theta}) = \prod_{v=1}^V p(\mathcal{D}_v|\boldsymbol{\theta}_v)
 \end{aligned}$$

where D<sub>v</sub> is values of v together with its parents and θ<sub>v</sub> is v's CPD

- ★ Called: decomposable likelihood (factorizes into family-factors)

# THE LIKELIHOOD FACTORIZES

- ★ Complete data

$$\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

$$\mathbf{x}_n = \{\mathbf{x}_{n1}, \dots, \mathbf{x}_{nV}\}$$

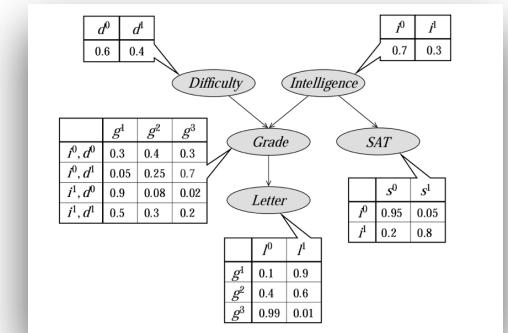
- ★ Likelihood

$$\begin{aligned}
 p(\mathcal{D}|\boldsymbol{\theta}) &= \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\theta}) = \prod_{n=1}^N \prod_{v=1}^V p(\mathbf{x}_{nv}|\mathbf{x}_{n,\text{pa}(v)}, \boldsymbol{\theta}) \\
 &= \prod_{v=1}^V \prod_{n=1}^N p(\mathbf{x}_{nv}|\mathbf{x}_{n,\text{pa}(v)}, \boldsymbol{\theta}) = \prod_{v=1}^V p(\mathcal{D}_v|\boldsymbol{\theta}_v)
 \end{aligned}$$

where  $D_v$  is values of  $v$  together with its parents and  $\theta_v$  is  $v$ 's CPD

- ★ Each  $p(\mathcal{D}_v|\boldsymbol{\theta}_v)$ , i.e., here each  $\boldsymbol{\theta}_{vc}$

can be maximized independently



# MLE FOR CAT CPDS

- ★ Since  $p(\mathcal{D}_v | \boldsymbol{\theta}_v)$ , i.e., here each  $\boldsymbol{\theta}_{vc}$

can be maximized independently

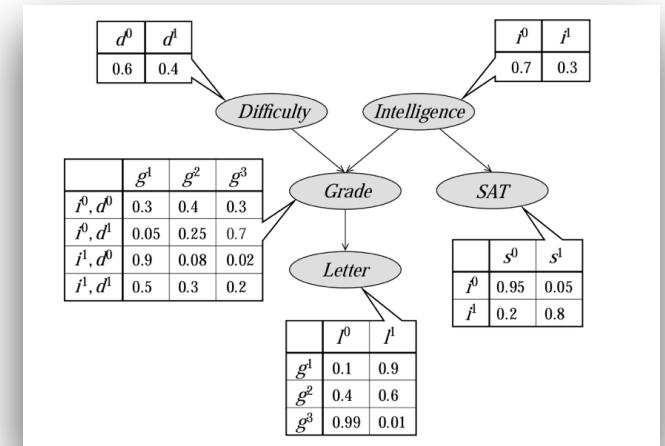
- ★ The MLE is

$$\boldsymbol{\theta}_{vc k} = N_{vc k} / N_{vc}$$

- ★ where

$$N_{vc k} = \sum_{n=1}^N I(x_{nv} = k, x_{n,\text{pa}(v)} = \mathbf{c})$$

$$N_{vc} = \sum_{n=1}^N I(x_{n,\text{pa}(v)} = \mathbf{c})$$



# BAYESIAN PARAMETER LEARNING

- ★ Decomposable prior

$$p(\boldsymbol{\theta}) = \prod_{v=1}^V p(\boldsymbol{\theta}_v) \quad \text{where} \quad \boldsymbol{\theta}_v = \{\boldsymbol{\theta}_{v\mathbf{c}}\}_{\mathbf{c} \in \mathbf{K}_{\text{pa}(\mathbf{v})}}$$

- ★ Gives decomposable posterior

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta})p(\boldsymbol{\theta}) = \prod_{v=1}^V p(\mathcal{D}_v|\boldsymbol{\theta}_v)p(\boldsymbol{\theta}_v)$$

# POSTERIOR

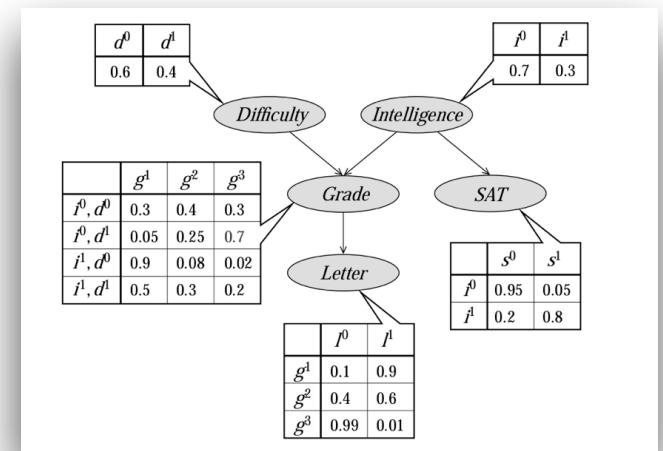
$$\boldsymbol{\theta}_v = \{\boldsymbol{\theta}_{v\mathbf{c}}\}_{\mathbf{c} \in \mathbf{K}_{\text{pa}(\mathbf{v})}}$$

- ★  $\alpha_{vc}$  is a vector of hyperparameters, prior

$$\boldsymbol{\theta}_{v\mathbf{c}} \sim \text{Dir}(\boldsymbol{\alpha}_{v\mathbf{c}})$$

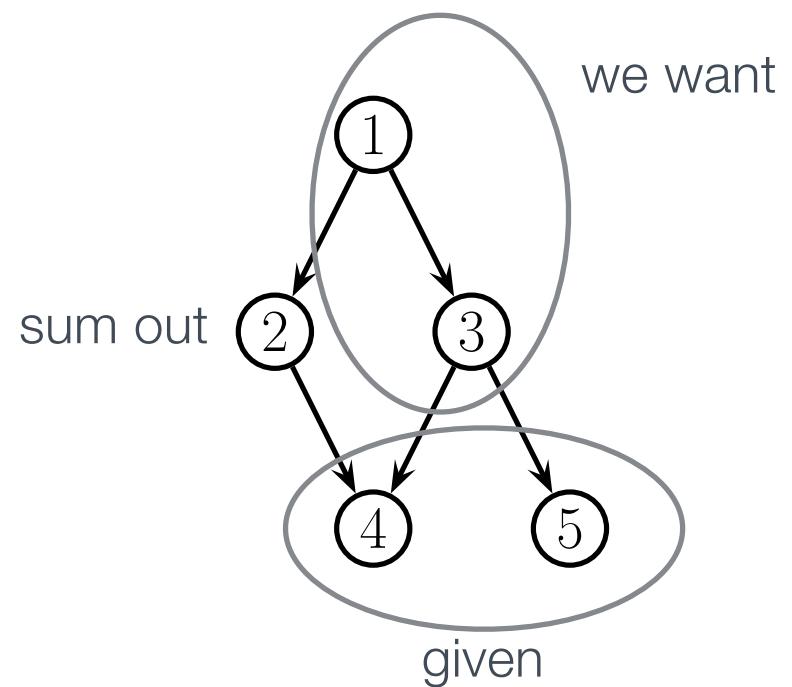
- ★ The posterior is

$$\boldsymbol{\theta}_{v\mathbf{c}} | \mathcal{D} \sim \text{Dir}(\mathbf{N}_{v\mathbf{c}} + \boldsymbol{\alpha}_{v\mathbf{c}})$$



# MARGINALIZE

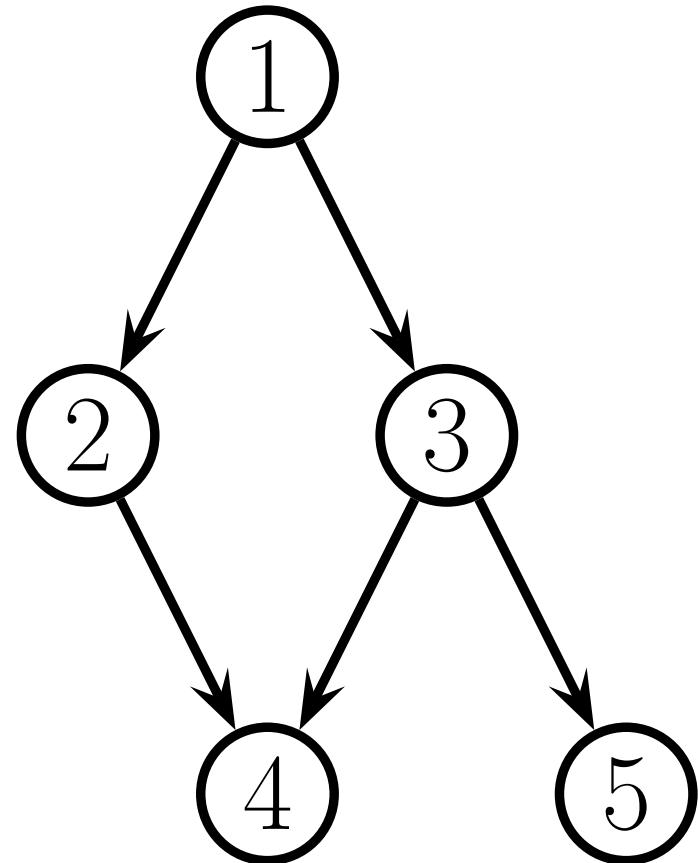
- $X_1, X_3$  two hidden variables that we want the marginal over
- $X_2$  the other hidden variables
- $x_e = \{x_4, x_5\}$  the visible variables



# MARGINALIZE

$$p(\mathbf{X}_m | \mathbf{x}_e, \boldsymbol{\theta}) = \frac{\sum_{\mathbf{x}_{V \setminus (m \cup e)}} p(\mathbf{X}_m, \mathbf{x}_{V \setminus (m \cup e)}, \mathbf{x}_e | \boldsymbol{\theta})}{\sum_{\mathbf{x}_{V \setminus e}} p(\mathbf{x}_{V \setminus e}, \mathbf{x}_e | \boldsymbol{\theta})}$$

- The denominator contains a marginal likelihood
- Summing out  $V$  binary hidden variables –  $O(2^V)$
- $K$  values –  $O(K^V)$



# DGM

- ★ What is the meaning of the underlying DAG? what is the semantics?
- ★ What does a DGM mean? what is the semantics?
- ★ Which DGMs represent a given distribution?

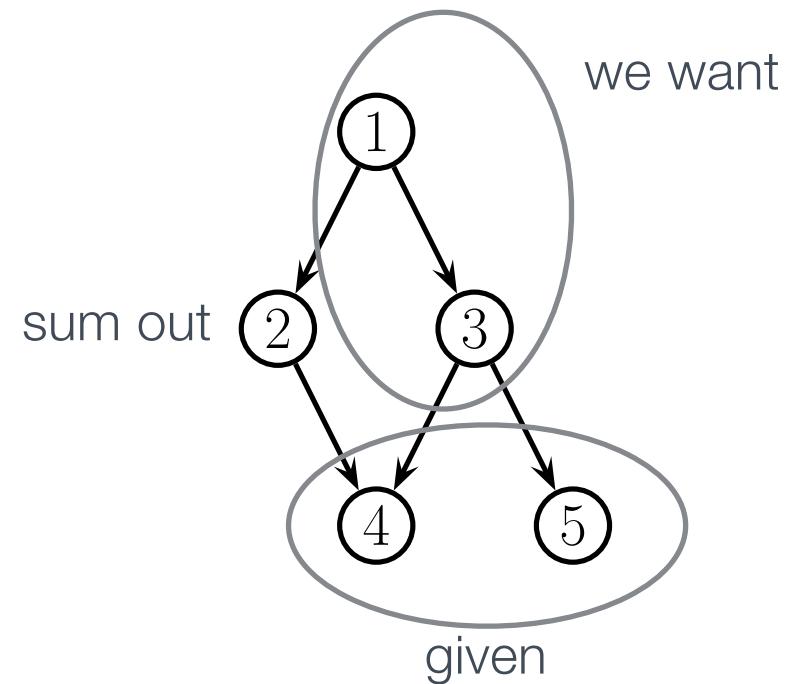
$X, X'$  two hidden variables

$X_h$  the other hidden variables

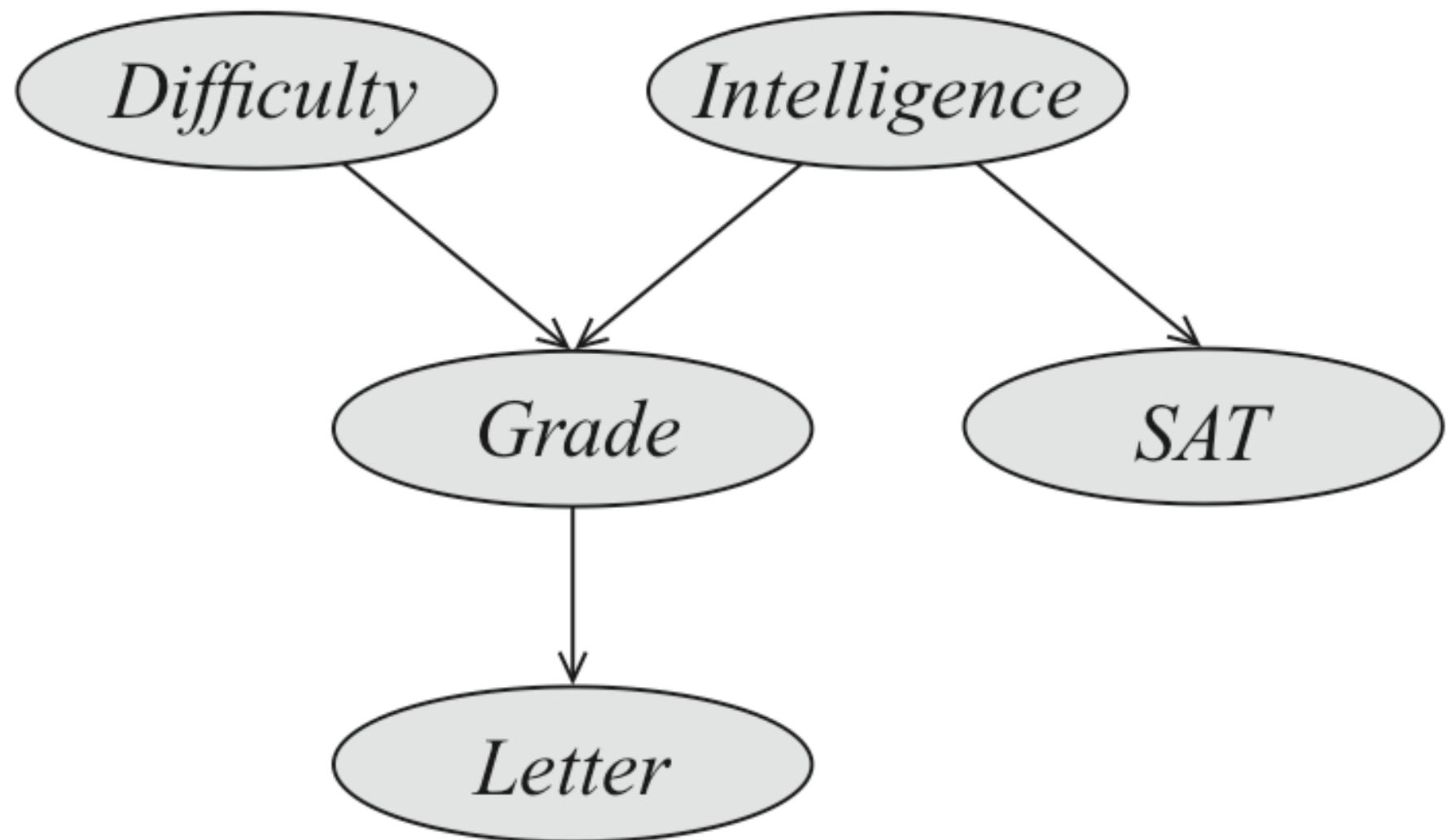
$X_v$  the visible variables

$$E[N_{k,k'}] = \sum_{\mathbf{x} \in \mathcal{D}} p(X = k, X' = k' | \mathbf{x}_v, \theta)$$

## EXPECTED SUFFICIENT STATISTICS - ESS



# EXTENDED STUDENT EXAMPLE



# INDEPENDENCE I-MAP

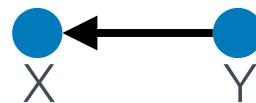
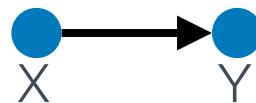
- ★  $I(G)$  (conditional) independences implied by  $G$  (not yet defined)
- ★  $I(P)$  (conditional) independences in the distribution  $P$
- ★  $G$  I-map for  $P$  in  $I(G) \subseteq I(P)$

p

$X$	$Y$	$P(X, Y)$
$x^0$	$y^0$	0.08
$x^0$	$y^1$	0.32
$x^1$	$y^0$	0.12
$x^1$	$y^1$	0.48

q

$X$	$Y$	$P(X, Y)$
$x^0$	$y^0$	0.4
$x^0$	$y^1$	0.3
$x^1$	$y^0$	0.2
$x^1$	$y^1$	0.1

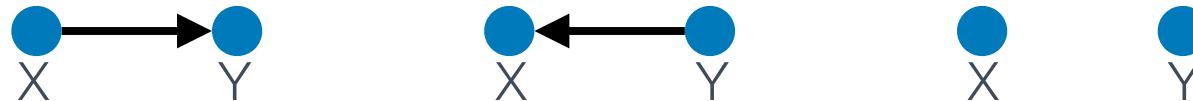


# INDEPENDENCE I-MAP

- ★  $I(G)$  independences implied by  $G$  (not yet defined)
- ★  $I(P)$  independences in the distribution  $P$
- ★  $G$  I-map for  $P$  in  $I(G) \subseteq I(P)$

p	$X$	$Y$	$P(X, Y)$
	$x^0$	$y^0$	0.08
	$x^0$	$y^1$	0.32
	$x^1$	$y^0$	0.12
	$x^1$	$y^1$	0.48

q	$X$	$Y$	$P(X, Y)$
	$x^0$	$y^0$	0.4
	$x^0$	$y^1$	0.3
	$x^1$	$y^0$	0.2
	$x^1$	$y^1$	0.1



- ★ p:  $X$  and  $Y$  ind. ex.  $p(X=1) = 0.48 + 0.12 = 0.6$ ,  $p(Y=1) = 0.8$ , and  $p(X=1, Y=1) = 0.48$
- ★ q:  $X$  and  $Y$  are dependent

# INDEPENDENCE I-MAP

- ★  $I(G)$  independences implied by  $G$  (not yet defined)
- ★  $I(P)$  independences in the distribution  $P$
- ★  $G$  I-map for  $P$  in  $I(G) \subseteq I(P)$

p	$X$	$Y$	$P(X, Y)$
	$x^0$	$y^0$	0.08
	$x^0$	$y^1$	0.32
	$x^1$	$y^0$	0.12
	$x^1$	$y^1$	0.48

q	$X$	$Y$	$P(X, Y)$
	$x^0$	$y^0$	0.4
	$x^0$	$y^1$	0.3
	$x^1$	$y^0$	0.2
	$x^1$	$y^1$	0.1

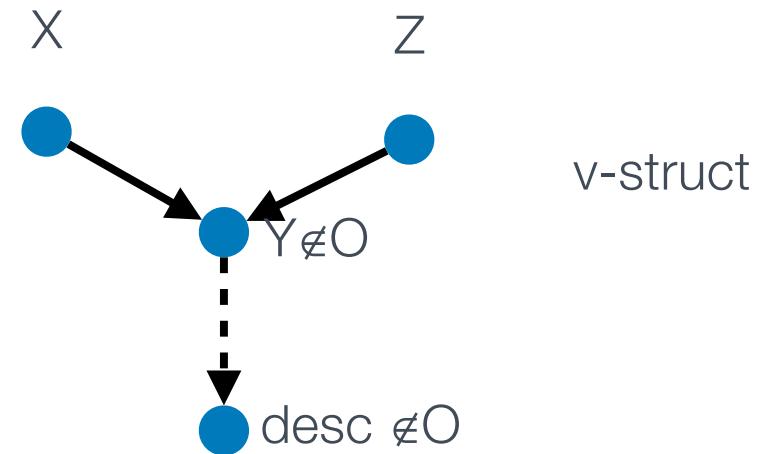
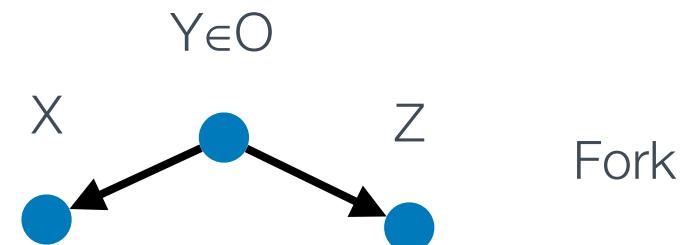
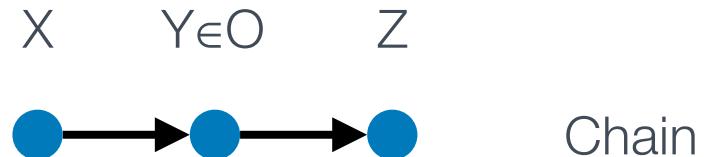


- ★ All three graphs are I-maps for  $p$
- ★  $G_1$  and  $G_2$  are I-maps for  $q$ , but  $G_3$  is not

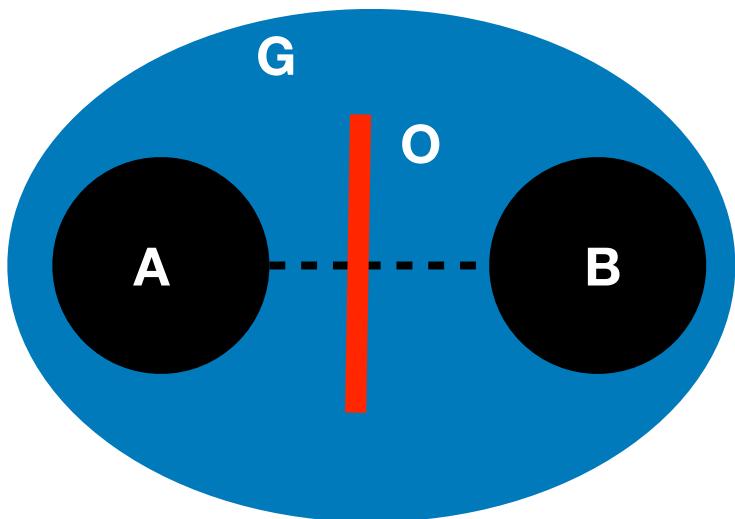
# D-SEPARATION

★ A path is d-separated by  $O$  if it has

- a chain  $X \rightarrow Y \rightarrow Z$  where  $Y \in O$
- a fork  $X \leftarrow Y \rightarrow Z$  where  $Y \in O$
- a v-structure  $X \rightarrow Y \leftarrow Z$  where  $(Y \cup \text{desc}(Y)) \cap O = \emptyset$



# D-SEPARATION SETS AND CI OF DAGS



★  $A$  is d-separated from  $B$  given  $O$  if every undirected path between  $A$  and  $B$  is d-separated by  $O$

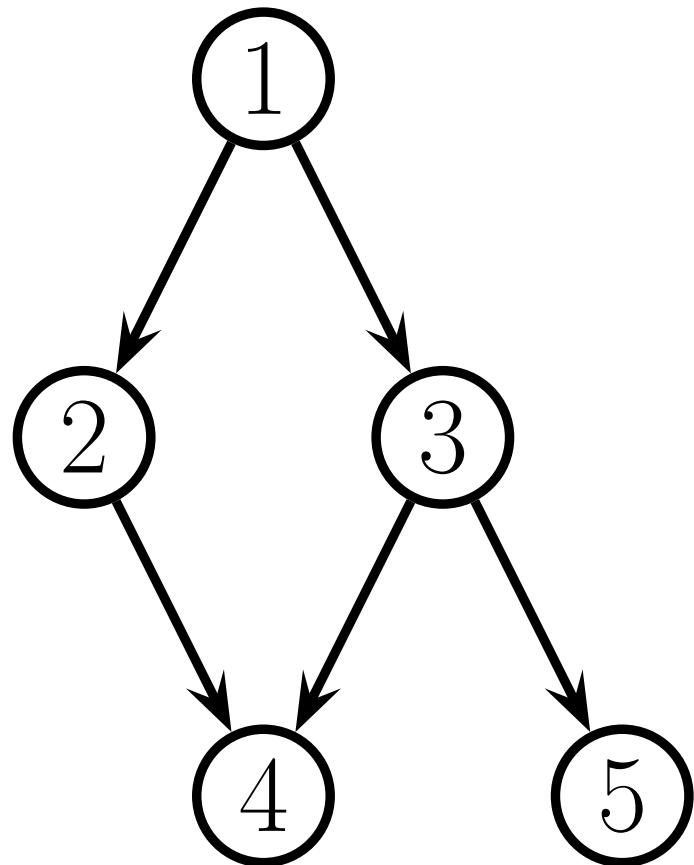
★ In a DAG  $G$ , global ind. def. by

$$x_A \perp_G x_B | x_O$$



$A$  is d-separated from  $B$  given  $O$

# ORDERED MARKOV PROPERTY



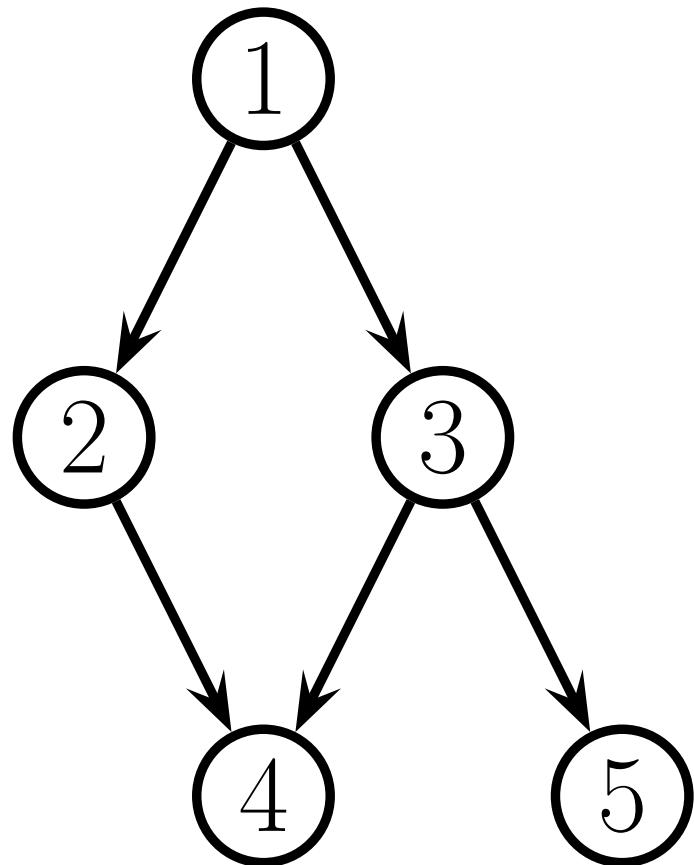
★ The directed local Markov property.

$$\mathbf{x}_t \perp_G \mathbf{x}_{\text{nd}(t) \setminus \text{pa}(t)} | \mathbf{x}_{\text{pa}(t)}$$

★ In this case

$$\begin{aligned} p(\mathbf{x}_{[5]}) &= p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1)p(\mathbf{x}_3|\mathbf{x}_1, \mathbf{x}_2) \\ &\quad p(\mathbf{x}_4|\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)p(\mathbf{x}_5|\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) \\ &= p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1)p(\mathbf{x}_3|\mathbf{x}_1) \\ &\quad p(\mathbf{x}_4|\mathbf{x}_2, \mathbf{x}_3)p(\mathbf{x}_5|\mathbf{x}_3) \end{aligned}$$

# ORDERED MARKOV PROPERTY



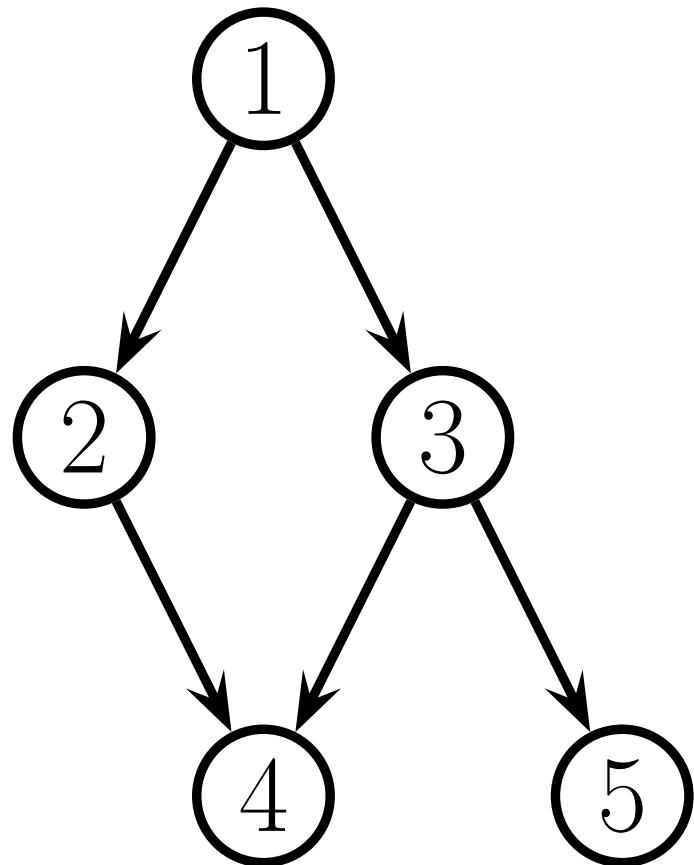
★ The directed local Markov property.

$$\mathbf{x}_t \perp_G \mathbf{x}_{\text{nd}(t) \setminus \text{pa}(t)} | \mathbf{x}_{\text{pa}(t)}$$

★ In general, if  $1, \dots, V$  topological order, the likelihood is decomposable (factorizes)

$$p(\mathbf{x}_{[V]} | G) = \prod_{t=1}^V p(\mathbf{x}_t | \mathbf{x}_{\text{pa}(t)})$$

# ORDERED MARKOV PROPERTY



★ The ordered Markov property.

$$\mathbf{x}_t \perp_G \mathbf{x}_{\text{pred}(t) \setminus \text{pa}(t)} | \mathbf{x}_{\text{pa}(t)}$$

topological order

★ In general, if  $1, \dots, V$  topological order,  
the likelihood is decomposable  
(factorizes)

$$p(\mathbf{x}_{[V]} | G) = \prod_{t=1}^V p(\mathbf{x}_t | \mathbf{x}_{\text{pa}(t)})$$

★ Global (G): d-separation

★ Local (L):

$$\mathbf{x}_t \perp_G \mathbf{x}_{\text{nd}(t) \setminus \text{pa}(t)} | \mathbf{x}_{\text{pa}(t)}$$

★ Ordered (O):  $\mathbf{x}_t \perp_G \mathbf{x}_{\text{pred}(t) \setminus \text{pa}(t)} | \mathbf{x}_{\text{pa}(t)}$

where pred is according to a topological order

★ Factorized (F): can be family-factorized

★ Theorem:  $G \Leftrightarrow L \Leftrightarrow O \Leftrightarrow F$

## EQUIVALENCE OF INDEPENDENCE DEFINITIONS

# SOUNDNESS AND COMPLETENESS

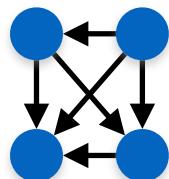
- ★ Theorem

If a distribution  $P$  factorizes according to  $G$ , then  $I(G) \subseteq I(P)$

- ★ Theorem

If  $X$  and  $Y$  are not d-separated given  $Z$  in  $G$ , then  $X$  and  $Y$  are dependent given  $Z$  in some distribution  $P$  that factorize over  $G$ .

We cannot have all. Ex. clique and independent distribution



# SKELETON AND EQUIVALENCE

- The skeleton is the underlying undirected graph
- Immorality is a pair of unmarried parents
- Theorem

Let  $G_1$  and  $G_2$  be two graphs over  $X$ . Then  $G_1$  and  $G_2$  have the same skeleton and the same set of immoralities if and only if

$$I(G_1) = I(G_2)$$

THE END