



Royal Institute of
Technology

STATISTICAL METHODS IN CS, CH 10

Lecture 5

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*.

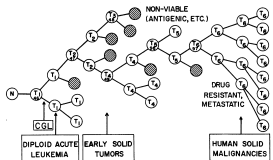
SOMATIC EVOLUTION

The Clonal Evolution of Tumor Cell Populations

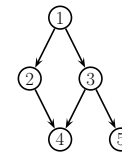
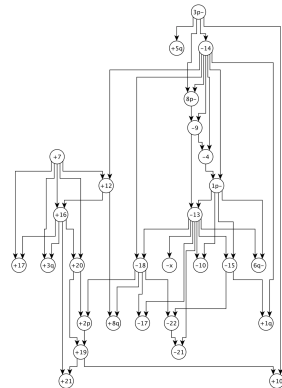
Acquired genetic lability permits stepwise selection of variant sublines and underlies tumor progression.

Peter C. Nowell

The author is professor of pathology, School of Medicine, University of Pennsylvania, Philadelphia 19174.
1 OCTOBER 1976 SCIENCE, VOL. 194

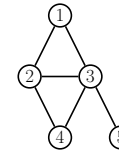


Oncogenetic network



Directed graphical model

- DAG
- vertices r.v.s
- equipped with local CPDs

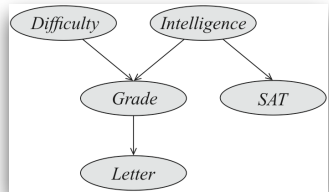


Undirected graphical model

- graph
- vertices r.v.s
- equipped with local "factors"

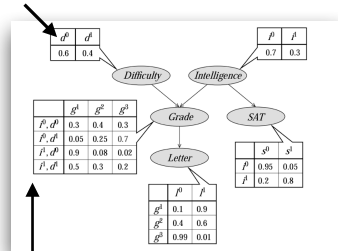
GRAPHICAL MODELS

DGM - GRAPH AND CPDS VS JOINT



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

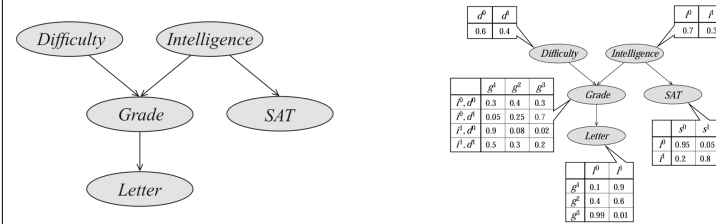
d has value 0



CPD

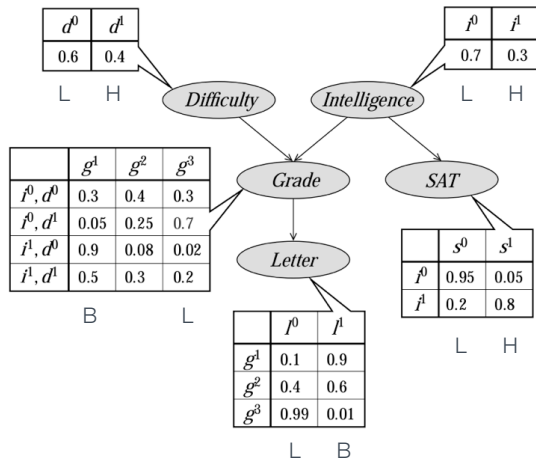
- * CPT - table, i.e., categorical
- * Gaussian

THREE LEVELS OF COMPUTATIONAL PROBLEMS



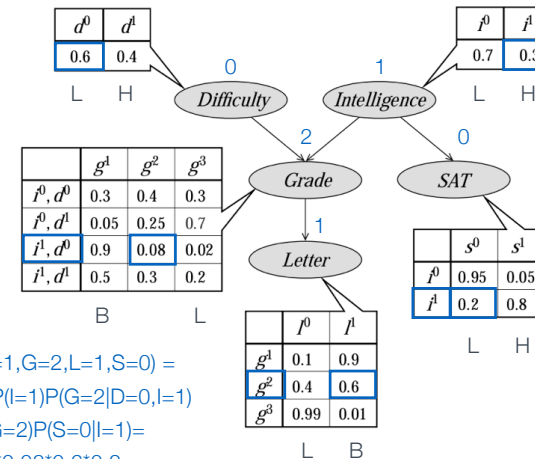
- Inference: given G and θ , compute probabilities or marginalize
- Parameter learning: given G and D, learn θ
- Structure learning: given D, learn G and θ

EXTENDED STUDENT EXAMPLE



B - better
H - higher
L - less

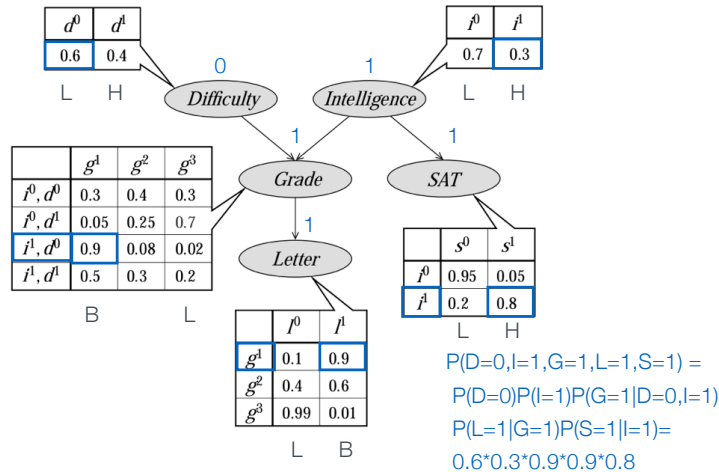
EXTENDED STUDENT EXAMPLE



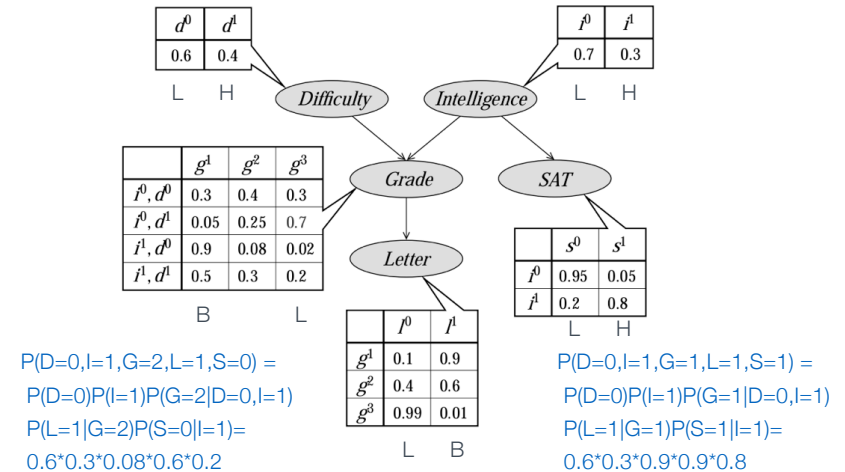
$$P(D=0, I=1, G=2, L=1, S=0) = P(D=0)P(I=1)P(G=2|D=0, I=1)P(L=1|G=2)P(S=0|I=1) = 0.6 * 0.3 * 0.08 * 0.6 * 0.2$$

B - better
H - higher
L - less

EXTENDED STUDENT EXAMPLE



EXTENDED STUDENT EXAMPLE



EXTENDED STUDENT EXAMPLE

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

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

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

INFERENCE – THE CHAIN RULE

$$p(\underbrace{x_{[V]}}_{x_1, \dots, x_V}) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \cdots p(x_V|x_{[V-1]})$$

- 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

CAT – NOTATION

- ★ For a $v \in [M]$,

values $k \in [K_v]$

combined values $c \in C_v = \prod_{s \in \text{pa}(v)} [K_s]$

↖ Cartesian product

- ★ Cat CPDs

where $P(x_v | \mathbf{x}_{\text{pa}(v)} = c) = \text{Cat}(\boldsymbol{\theta}_{vc})$

and $\boldsymbol{\theta}_{vck} = P(x_v = k | \mathbf{x}_{\text{pa}(v)} = c)$

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

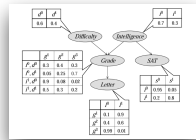
$$= \prod_{v=1}^V \prod_{n=1}^N p(\mathbf{x}_{nv}|\mathbf{x}_{n,\text{pa}(v)}, \boldsymbol{\theta}) = \prod_{v=1}^V p(\mathcal{D}_v|\boldsymbol{\theta}_v)$$

where \mathcal{D}_v is values of v together with its parents and $\boldsymbol{\theta}_v$ 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

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

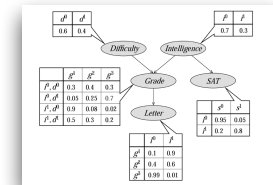
$$= \prod_{v=1}^V \prod_{n=1}^N p(\mathbf{x}_{nv}|\mathbf{x}_{n,\text{pa}(v)}, \theta) = \prod_{v=1}^V p(\mathcal{D}_v|\theta_v)$$

where \mathcal{D}_v is values of v together with its parents and θ_v is v 's CPD

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

MLE FOR CAT CPDS

- ★ Each $P(\mathcal{D}_v|\theta_v)$, i.e., here each $\theta_{vc} = (\theta_{vc1}, \dots, \theta_{vcK_v})$ can be maximized independently



- ★ So, MLE is

$$\theta_{vc\mathbf{c}} = N_{vc\mathbf{c}}/N_{vc}$$

- ★ where

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

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

c

BAYESIAN PARAMETER LEARNING

- ★ Decomposable prior

$$p(\theta) = \prod_{v=1}^V p(\theta_v) \quad \text{where } \theta_v = (\theta_{v1}, \dots, \theta_{vK_{\text{pa}(v)}})$$

- ★ Gives decomposable posterior

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

POSTERIOR

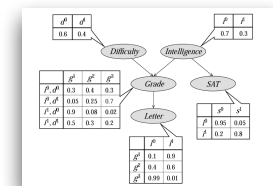
$$\theta_{vc} = (\theta_{vc1}, \dots, \theta_{vcK_v})$$

- ★ α_{vc} is a vector of hyperparameters, prior

$$\theta_{vc} \sim \text{Dir}(\alpha_{vc})$$

- ★ The posterior is

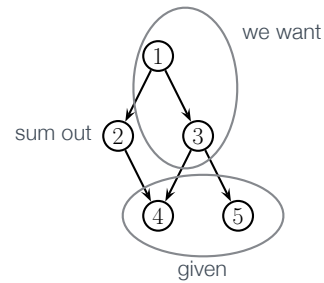
$$\theta_{vc}|\mathcal{D} \sim \text{Dir}(N_{vc} + \alpha_{vc})$$



X, X' two hidden variables
 X_h the other hidden variables
 X_v the visible variables

MARGINALIZE

$$\begin{aligned}
 p(X = k, X' = k' | \mathbf{x}_v, \boldsymbol{\theta}) &= \frac{p(X = k, X' = k', \mathbf{x}_v | \boldsymbol{\theta})}{P(\mathbf{x}_v | \boldsymbol{\theta})} \\
 &= \frac{\sum_{\mathbf{x}_h} p(X = k, X' = k', \mathbf{x}_h, \mathbf{x}_v | \boldsymbol{\theta})}{\sum_{\mathbf{x}, \mathbf{x}', \mathbf{x}_h} p(\mathbf{x}, \mathbf{x}', \mathbf{x}_h, \mathbf{x}_v | \boldsymbol{\theta})}
 \end{aligned}$$



MARGINALIZE

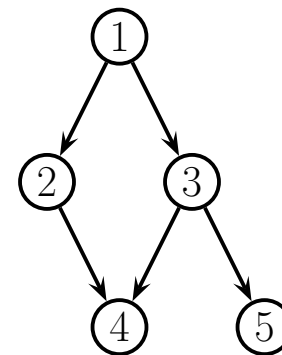
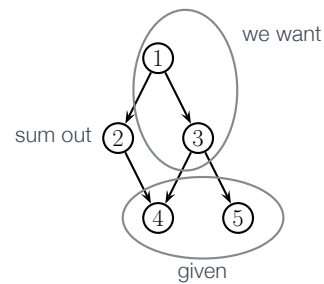
$$p(\mathbf{X}_m | \mathbf{x}_e, \boldsymbol{\theta}) = \frac{p(\mathbf{X}_m, \mathbf{x}_e | \boldsymbol{\theta})}{p(\mathbf{x}_e | \boldsymbol{\theta})} = \frac{\sum_{\mathbf{x}_{V \setminus (m \cup e)}} p(\mathbf{X}_m, \mathbf{x}_e | \boldsymbol{\theta})}{\sum_{\mathbf{x}_{V \setminus e}} p(\mathbf{x}_e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

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

X, X' two hidden variables
 X_h the other hidden variables
 X_v the visible variables

EXPECTED SUFFICIENT STATISTICS - ESS

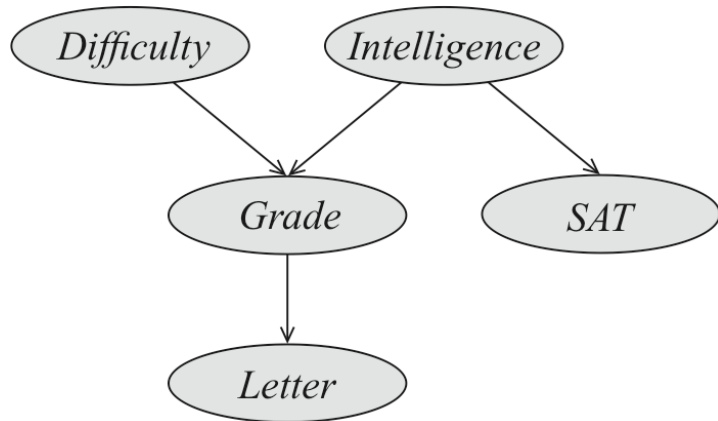
$$E[N_{k,k'}] = \sum_{\mathbf{x} \in \mathcal{D}} p(X = k, X' = k' | \mathbf{x}_v, \boldsymbol{\theta})$$



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?

EXTENDED STUDENT EXAMPLE



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 $I(G) \subseteq I(P)$

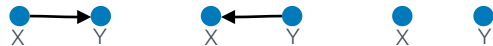
p	X	Y	$P(X,Y)$	q	X	Y	$P(X,Y)$
	x^0	y^0	0.08		x^0	y^0	0.4
	x^0	y^1	0.32		x^0	y^1	0.3
	x^1	y^0	0.12		x^1	y^0	0.2
	x^1	y^1	0.48		x^1	y^1	0.1



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	$P(X,Y)$	q	X	Y	$P(X,Y)$
	x^0	y^0	0.08		x^0	y^0	0.4
	x^0	y^1	0.32		x^0	y^1	0.3
	x^1	y^0	0.12		x^1	y^0	0.2
	x^1	y^1	0.48		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

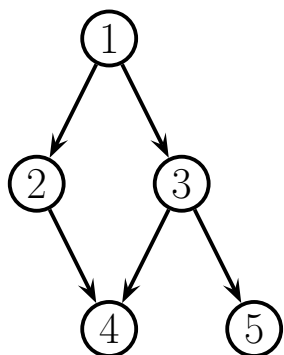
INDEPENDENCE I-MAP

- * $I(G)$ 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	$P(X,Y)$	q	X	Y	$P(X,Y)$
	x^0	y^0	0.08		x^0	y^0	0.4
	x^0	y^1	0.32		x^0	y^1	0.3
	x^1	y^0	0.12		x^1	y^0	0.2
	x^1	y^1	0.48		x^1	y^1	0.1

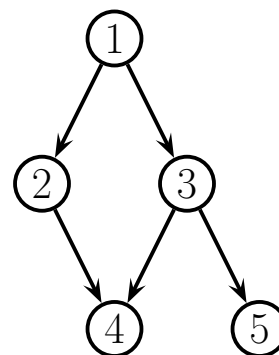


- * All three graphs are I-maps for p
- * G_1 and G_2 are I-maps for q , but G_3 is not



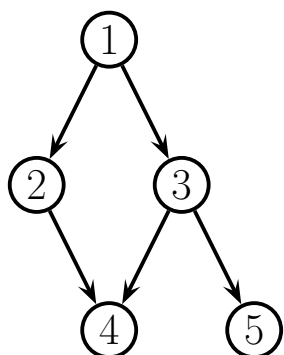
TERMINOLOGY

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



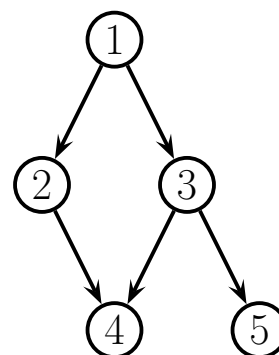
TERMINOLOGY

- ★ 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



ORDERED MARKOV PROPERTY

- ★ The directed local Markov property.

$$\mathbf{x}_t \perp \mathbf{x}_{V \setminus \text{desc}(t)} \mid \mathbf{x}_{\text{pa}(t)}$$

- ★ In this case

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

★ Global (G): d-separation

★ Local (L): $\mathbf{X}_t \perp \mathbf{X}_{V \setminus \text{desc}(t)} \mid \mathbf{X}_{\text{pa}(t)}$

★ Ordered (O): $\mathbf{X}_t \perp \mathbf{X}_{\text{pred}(t)} \mid \mathbf{X}_{\text{pa}(t)}$

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