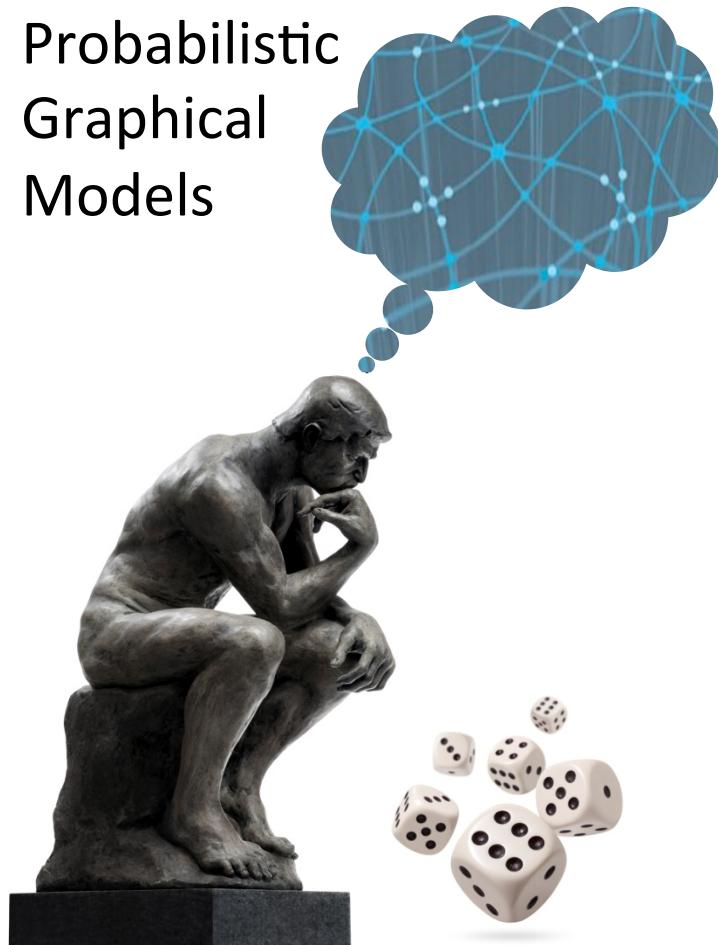


Probabilistic
Graphical
Models



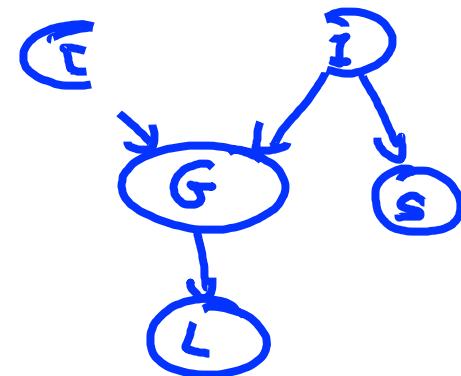
Representation

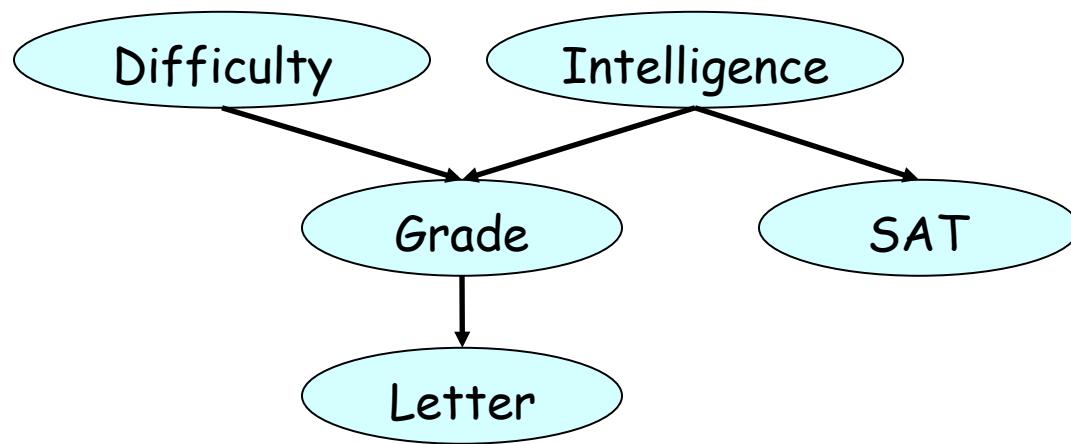
Bayesian Networks

Semantics &
Factorization

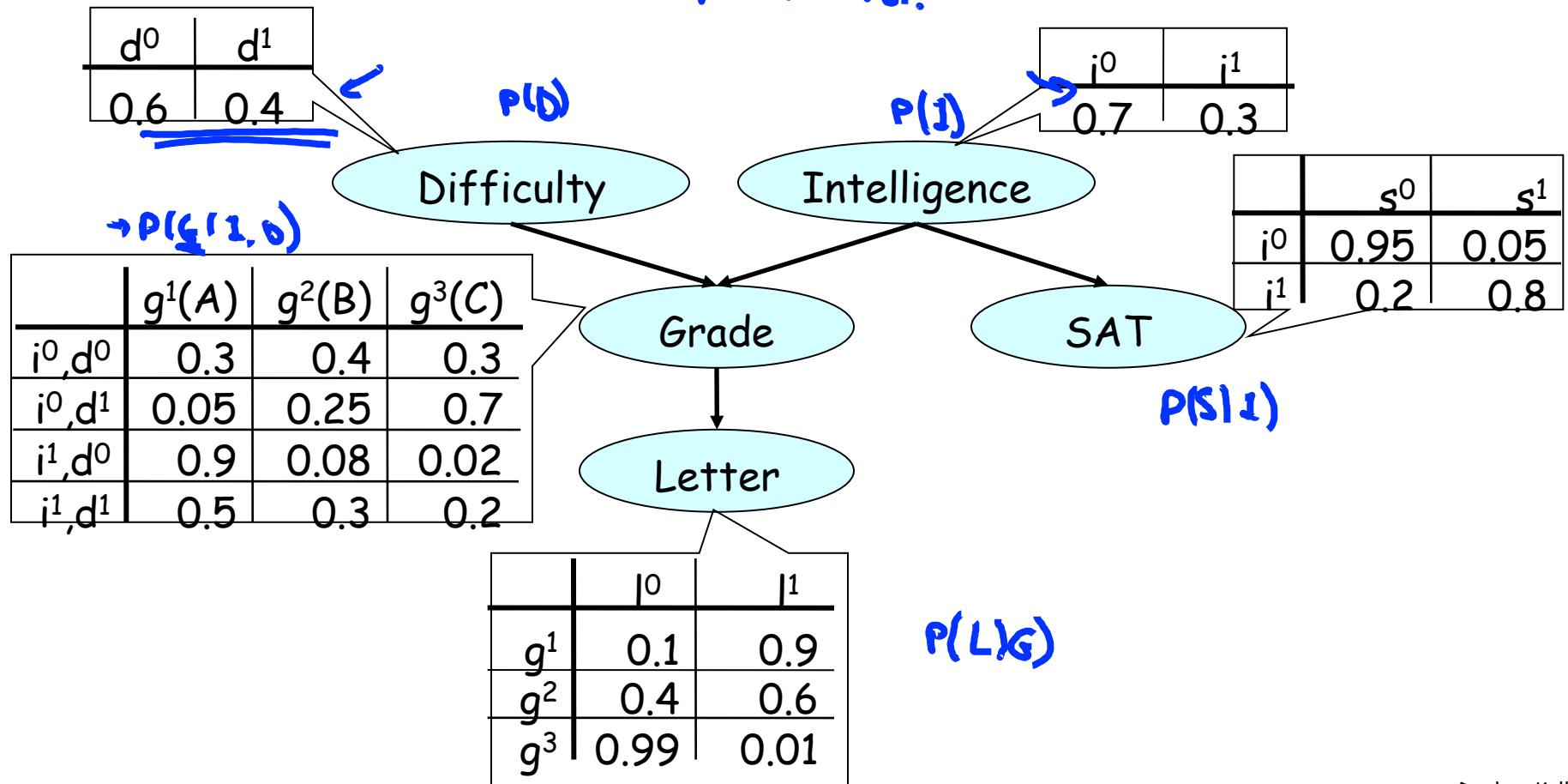
- Grade
- Course Difficulty
- Student Intelligence
- Student SAT
- Reference Letter

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

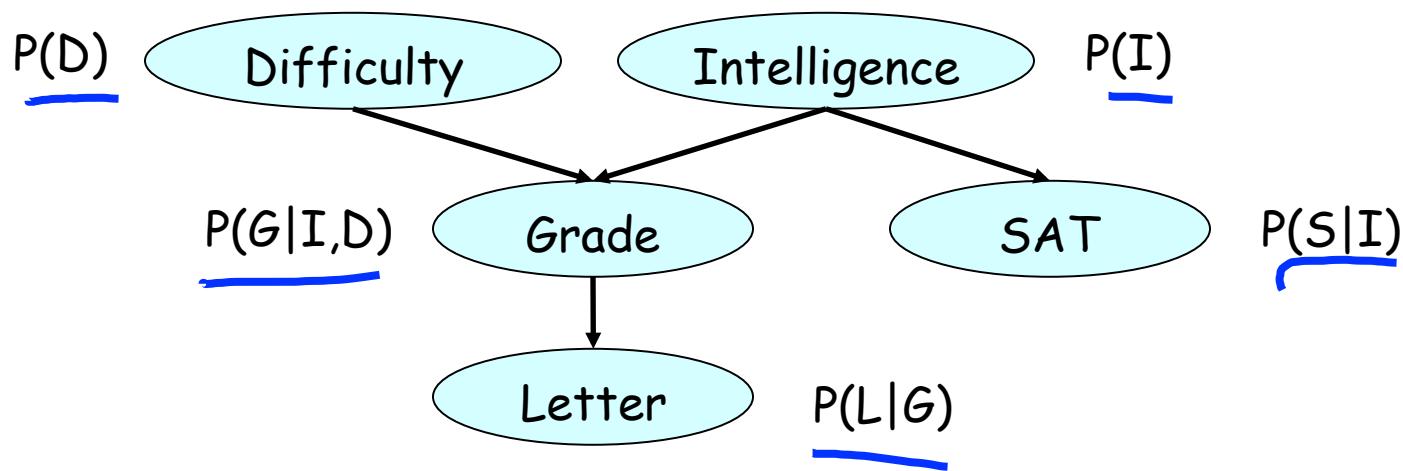




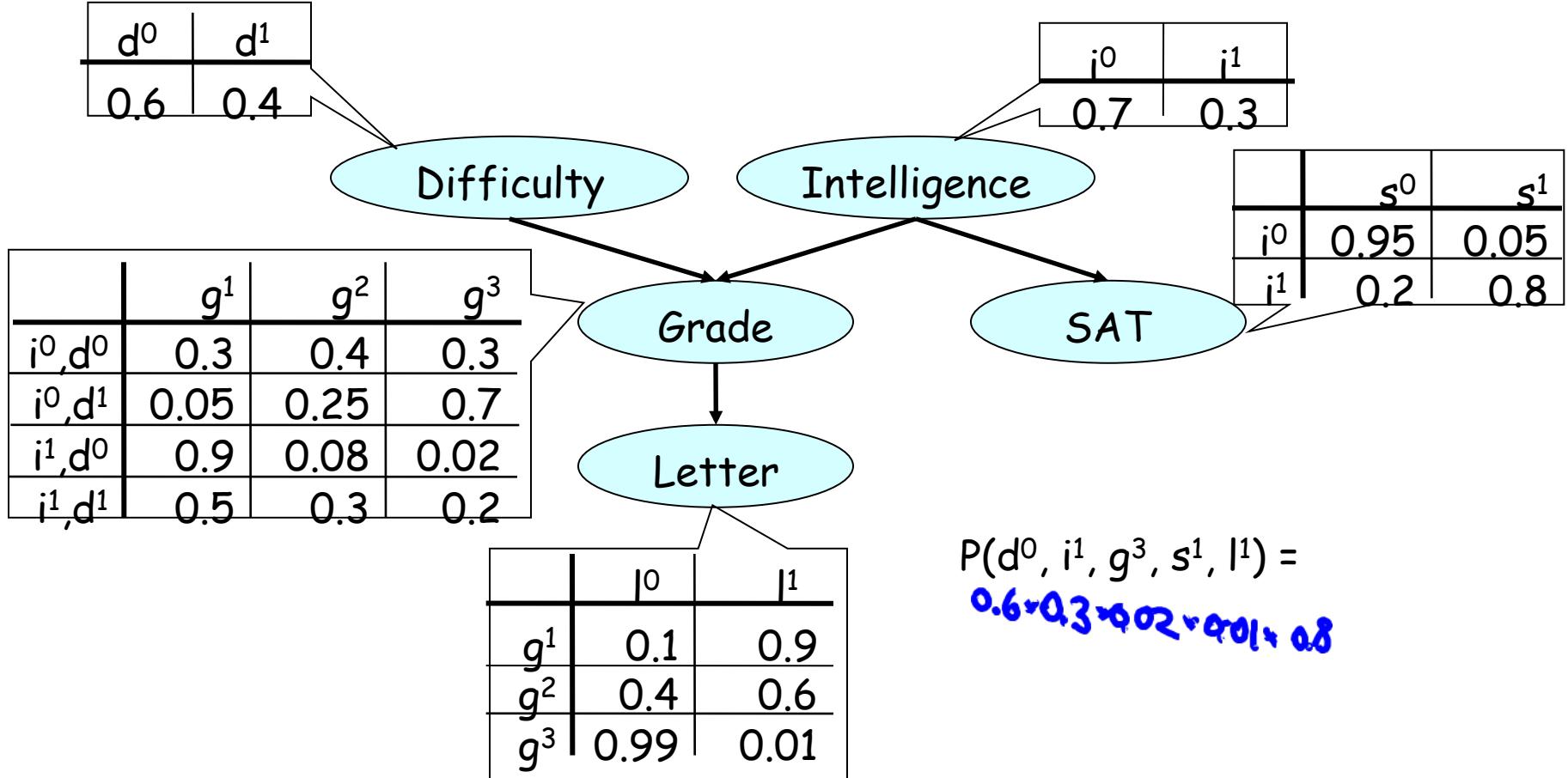
CPD = cond. prob. dist.



Chain Rule for Bayesian Networks



$$\underbrace{P(D, I, G, S, L)}_{\text{Distribution defined as a product of factors!}} = P(D) P(I) P(G|I,D) P(S|I) P(L|G)$$



Bayesian Network

- A Bayesian network is:
 - A directed acyclic graph (DAG) G whose nodes represent the random variables X_1, \dots, X_n
 - For each node $\underline{X_i}$ a CPD $P(\underline{X_i} \mid \underline{\text{Par}_G(X_i)})$
- The BN represents a joint distribution via the chain rule for Bayesian networks

$$P(X_1, \dots, X_n) = \prod_i P(X_i \mid \text{Par}_G(X_i))$$

BN Is a Legal Distribution: $P \geq 0$

P is a product of CPDs

CPDs are non-negative

BN Is a Legal Distribution: $\sum P = 1$

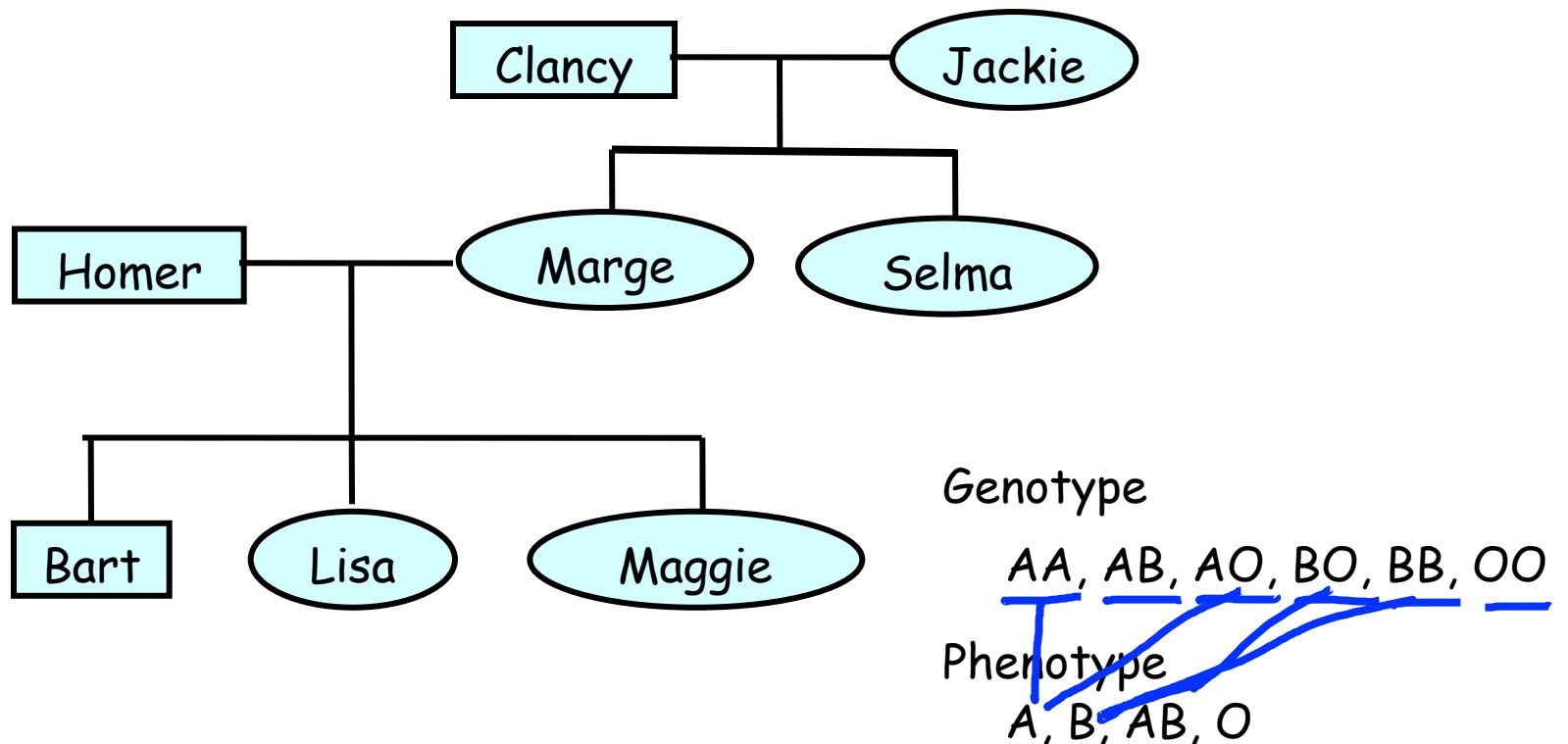
$$\begin{aligned}\sum_{D,I,G,S,L} P(D, I, G, S, L) &= \sum_{D,I,G,S,L} P(D) P(I) P(G|I,D) P(S|I) P(L|G) \\&= \sum_{D,I,G,S} P(D) P(I) P(G|I,D) P(S|I) \sum_L P(L|G) \\&= \sum_{D,I,G,S} P(D) P(I) P(G|I,D) P(S|I) \\&= \sum_{D,I,G} P(D) P(I) P(G|I,D) \sum_S P(S|I) \\&= \sum_{D,I} P(D) P(I) \sum_G P(G|I,D)\end{aligned}$$

P Factorizes over G

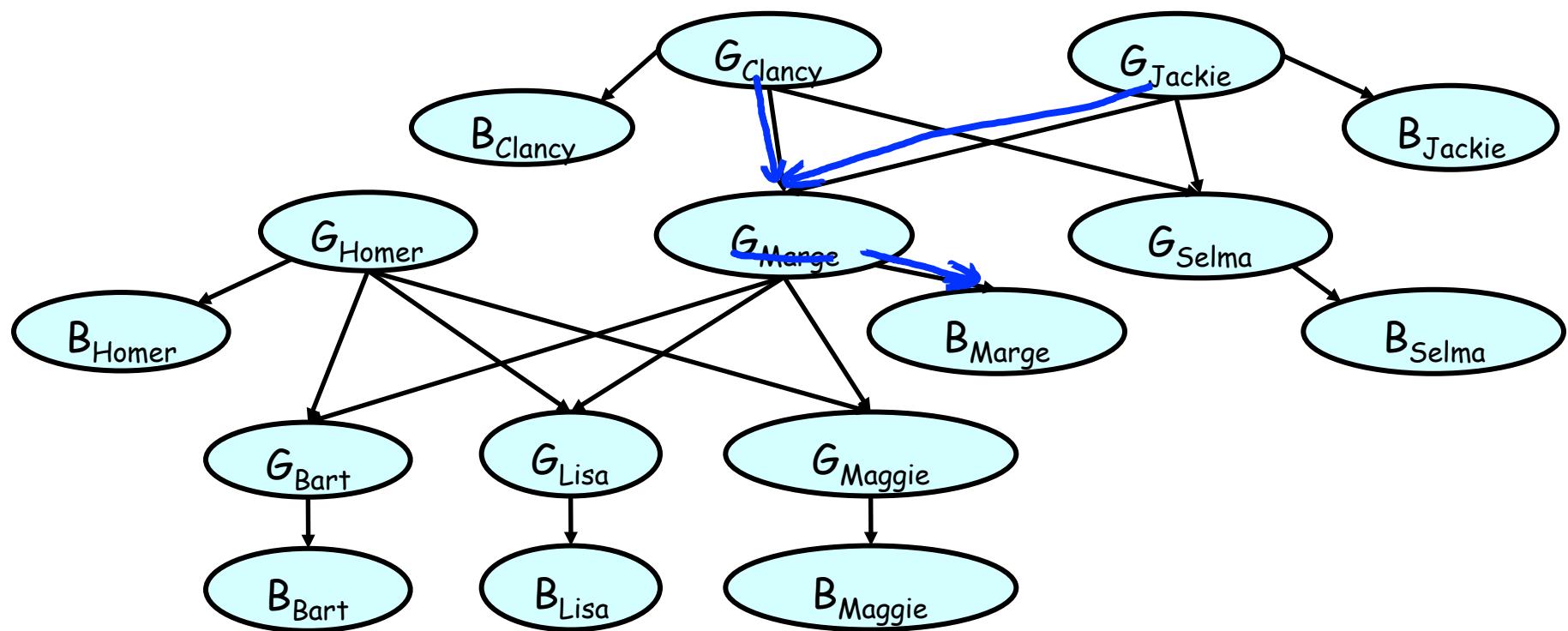
- Let G be a graph over X_1, \dots, X_n .
- P factorizes over G if

$$P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Par}_G(X_i))$$

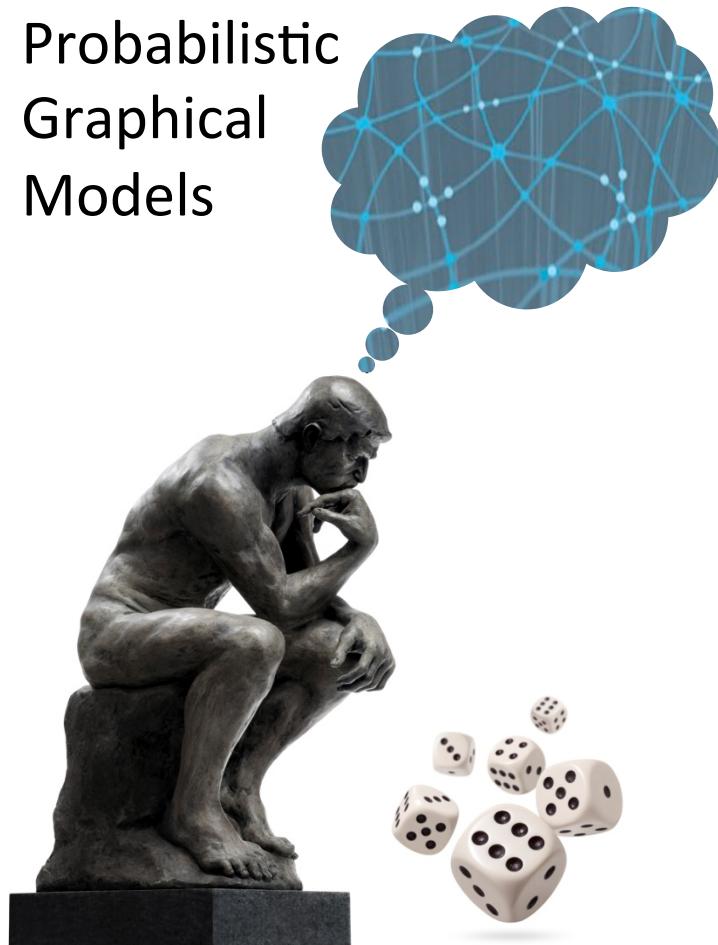
Genetic Inheritance



BNs for Genetic Inheritance



Probabilistic
Graphical
Models

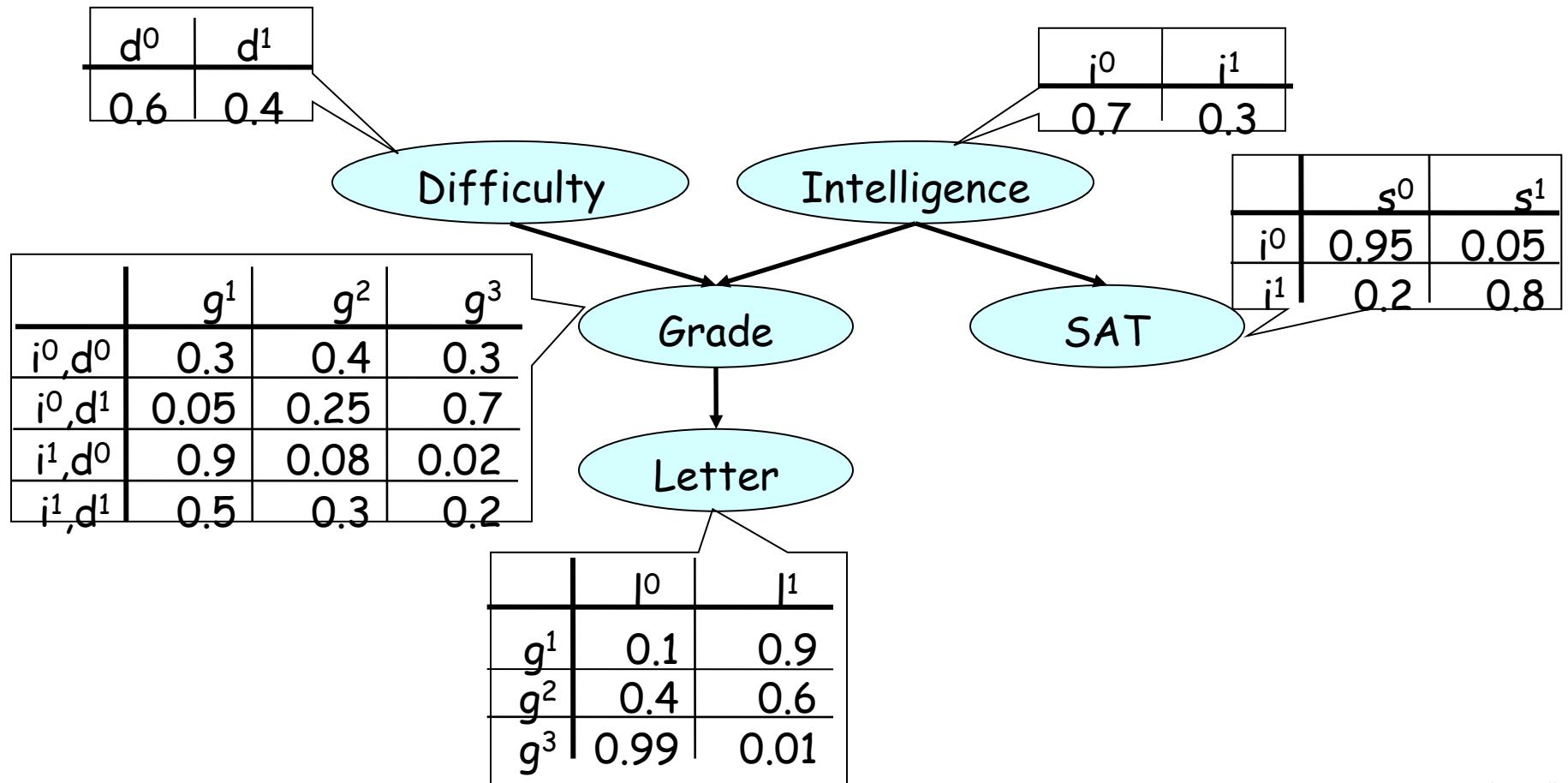


Representation

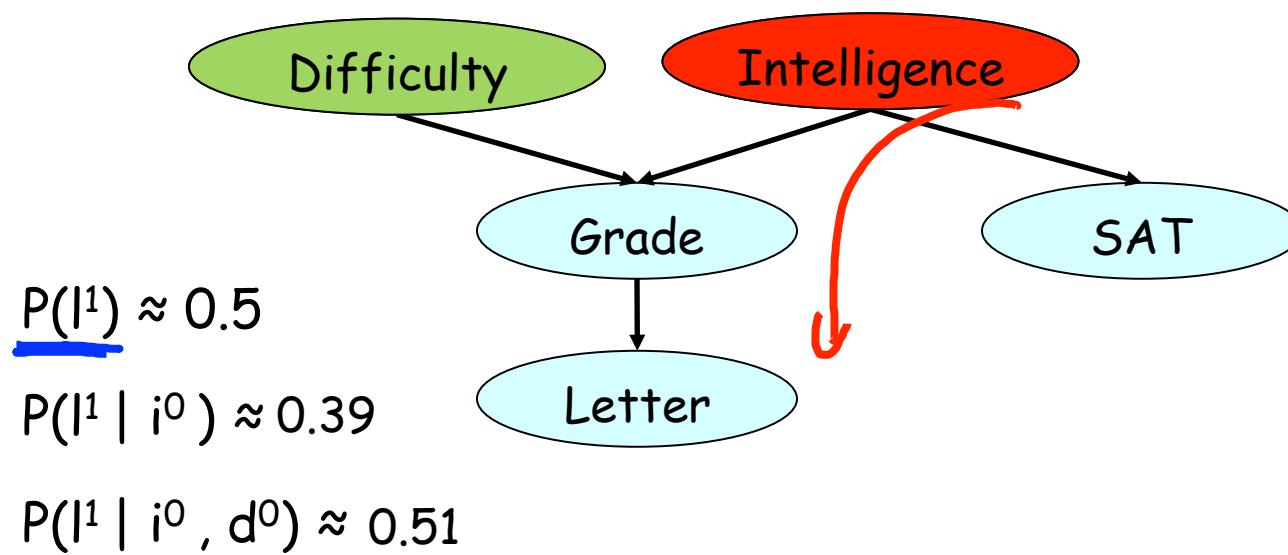
Bayesian Networks

Reasoning
Patterns

The Student Network



Causal Reasoning



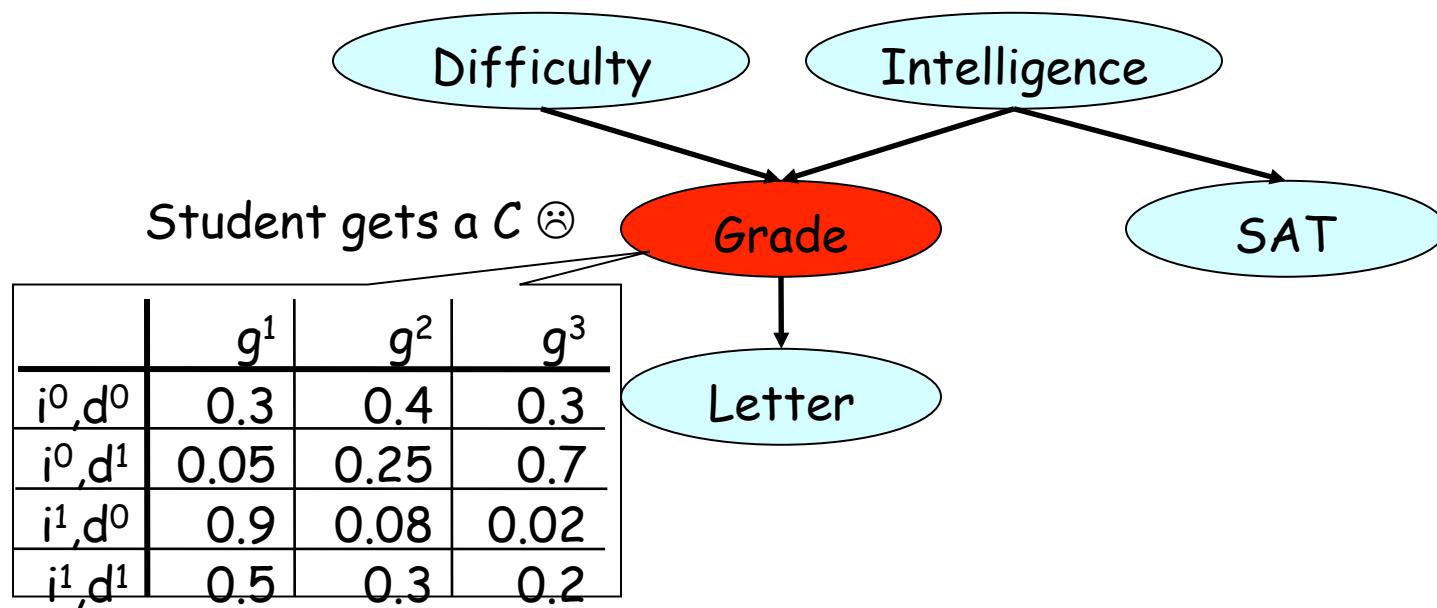
Evidential Reasoning

$$P(d^1) = 0.4$$

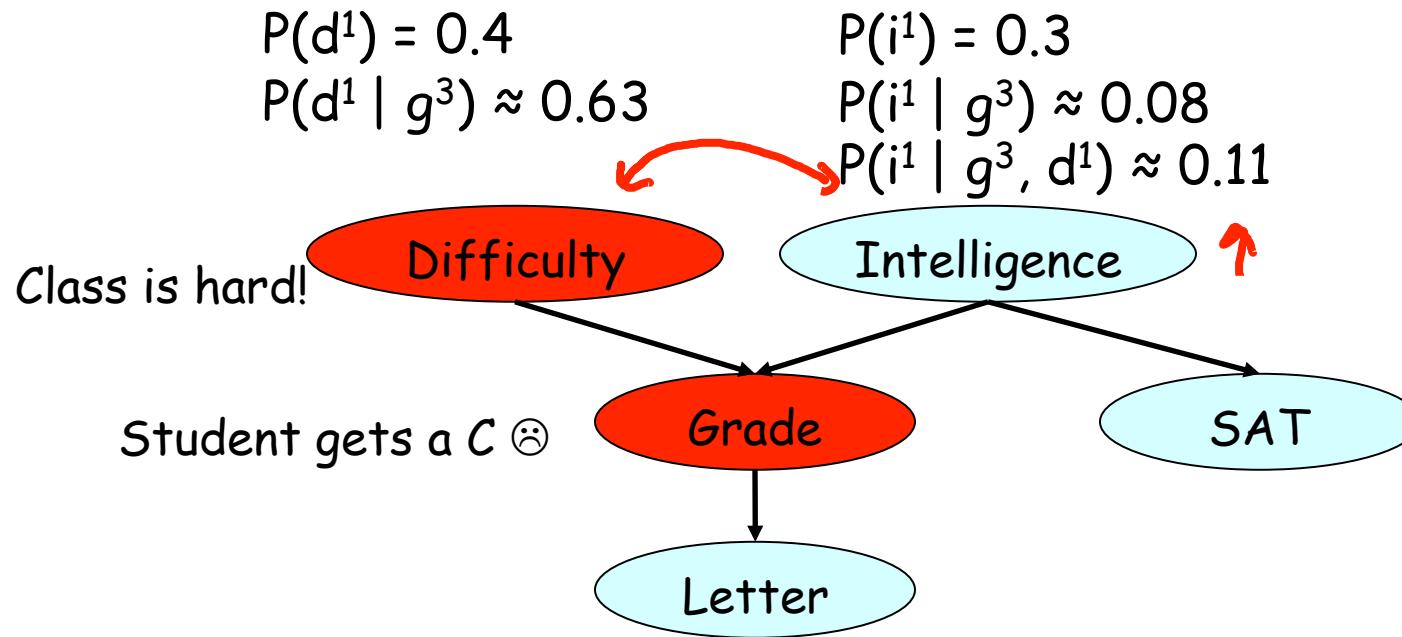
$$P(d^1 | g^3) \approx 0.63$$

$$P(i^1) = 0.3$$

$$P(i^1 | g^3) \approx 0.08$$

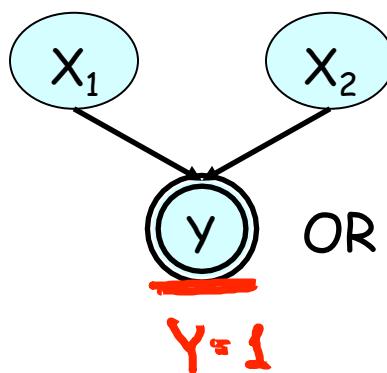


Intercausal Reasoning



Intercausal Reasoning Explained

explaining away



X_1	X_2	Y	Prob
0	0	0	0.25
0	1	1	0.25
1	0	1	0.25
1	1	1	0.25

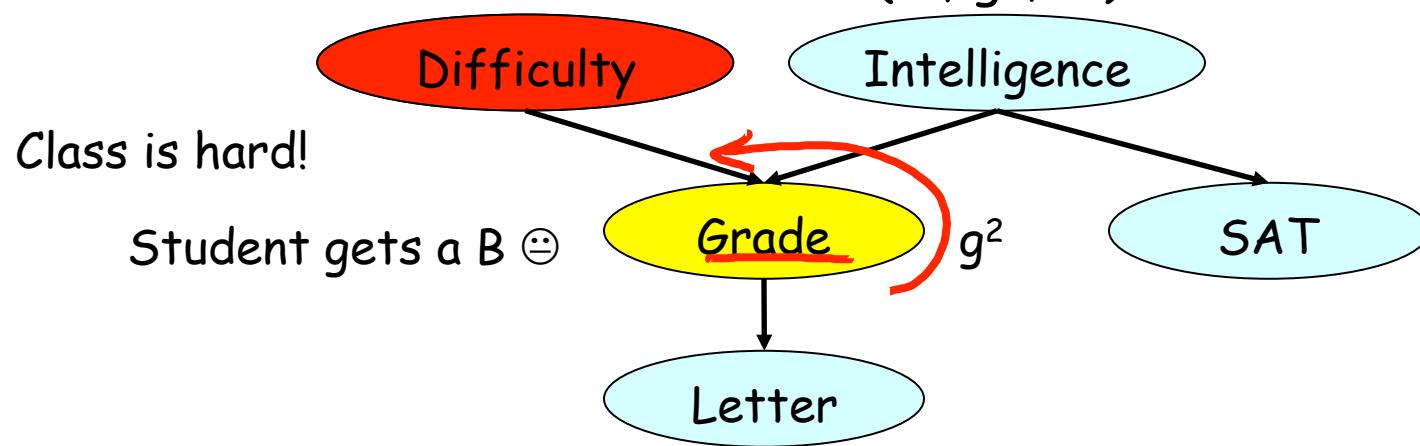
Annotations in red and blue highlight specific rows and columns of the table, corresponding to the causal paths from X_1 and X_2 to y .

$$\text{Ans: } P(X_1=1) = \frac{2}{3} \quad P(X_2=1) = \frac{2}{3}$$

Condition $X_1=1 \quad P(Y=1 | X_2=1) = 0.5$

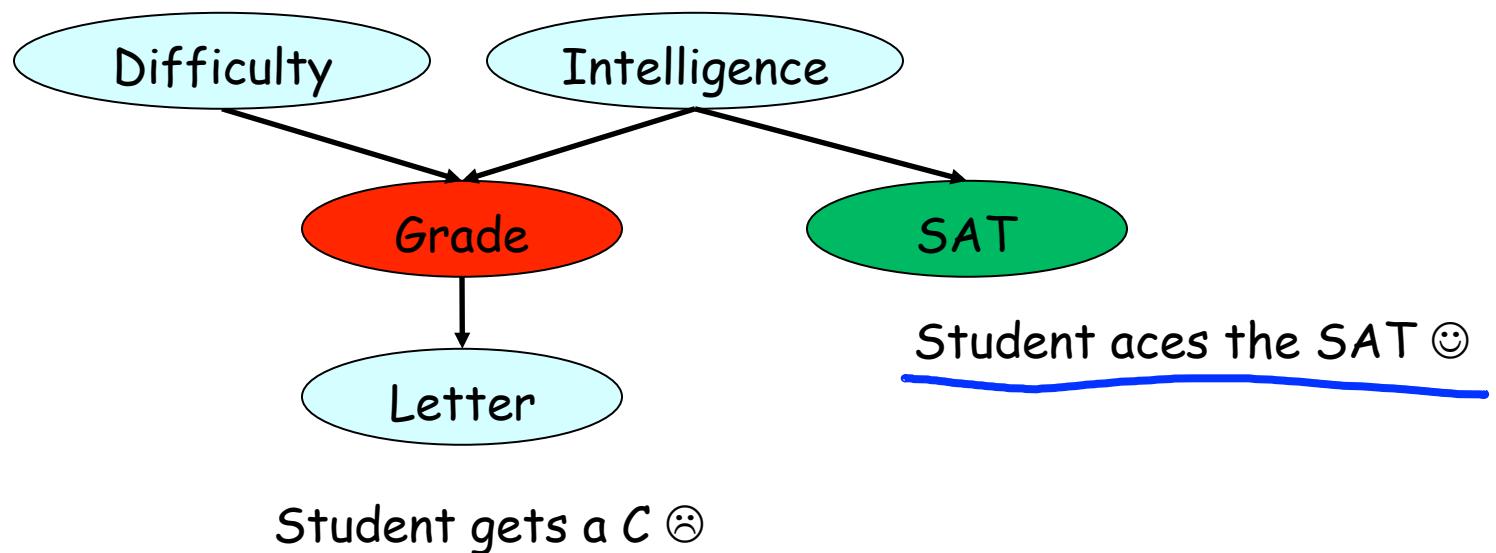
Intercausal Reasoning II

$$\begin{aligned}P(i^1) &= 0.3 \\P(i^1 | g^2) &\approx 0.175 \\P(i^1 | g^2, d^1) &\approx 0.34\end{aligned}$$



Student Aces the SAT

- What happens to the posterior probability that the class is hard?



Student Aces the SAT

$$P(d^1) = 0.4$$

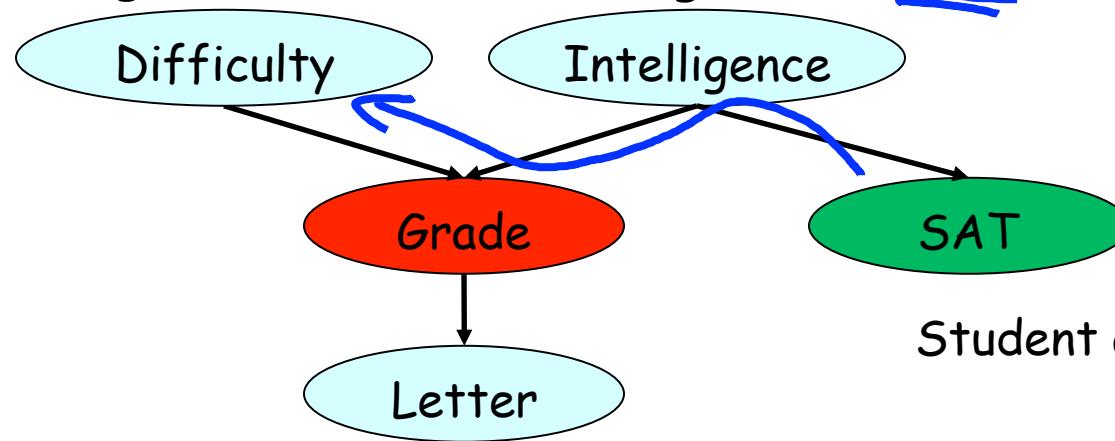
$$P(d^1 | g^3) \approx 0.63$$

$$P(d^1 | g^3, s^1) \approx \underline{0.76}$$

$$P(i^1) = 0.3$$

$$P(i^1 | g^3) \approx \underline{0.08}$$

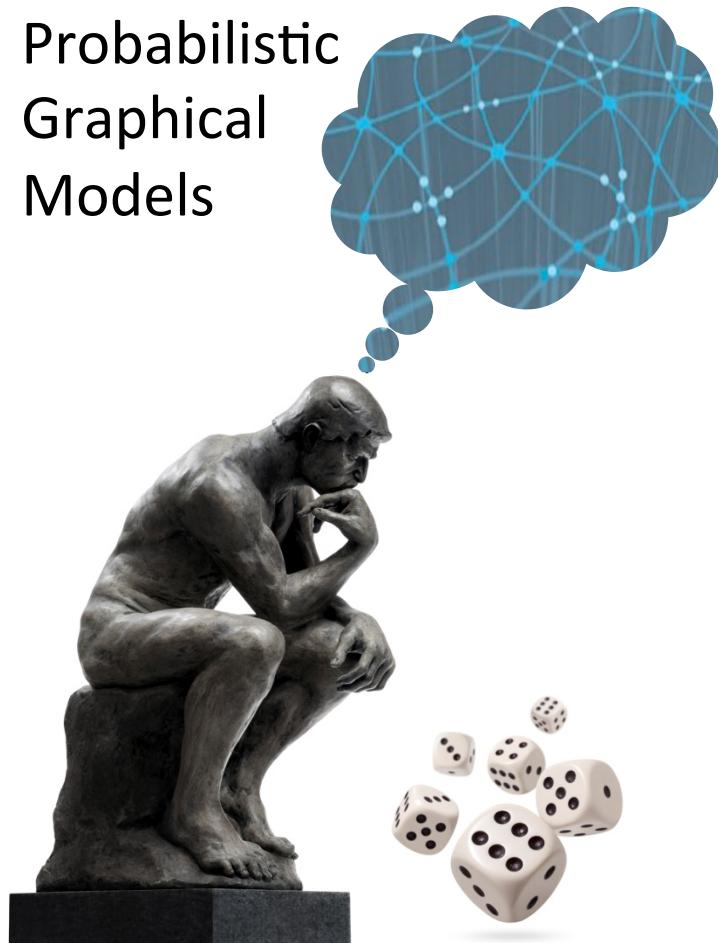
$$P(i^1 | g^3, s^1) \approx \underline{0.58}$$



Student aces the SAT ☺

Student gets a C ☹

Probabilistic
Graphical
Models



Representation

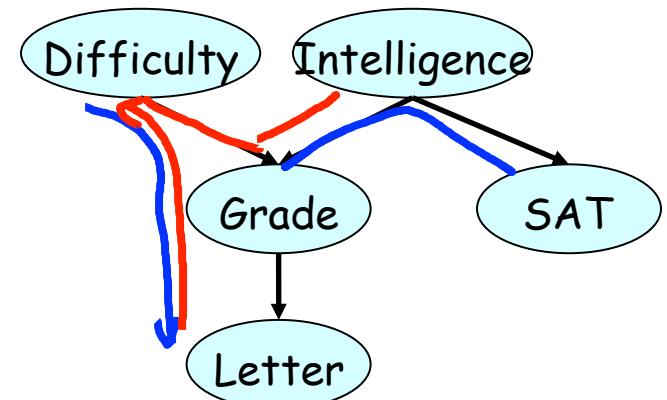
Bayesian Networks

Flow of
Probabilistic
Influence

When can X influence Y?

Condition on v-structure before inserting Y

- $X \rightarrow Y$ ✓
- $X \leftarrow Y$ ✓
- $X \rightarrow W \rightarrow Y$ ✓
- $X \leftarrow W \leftarrow Y$ ✓
- $X \leftarrow \underline{W} \rightarrow Y$ ✓
- $X \rightarrow \underline{W} \leftarrow Y$ ✗
v-structure



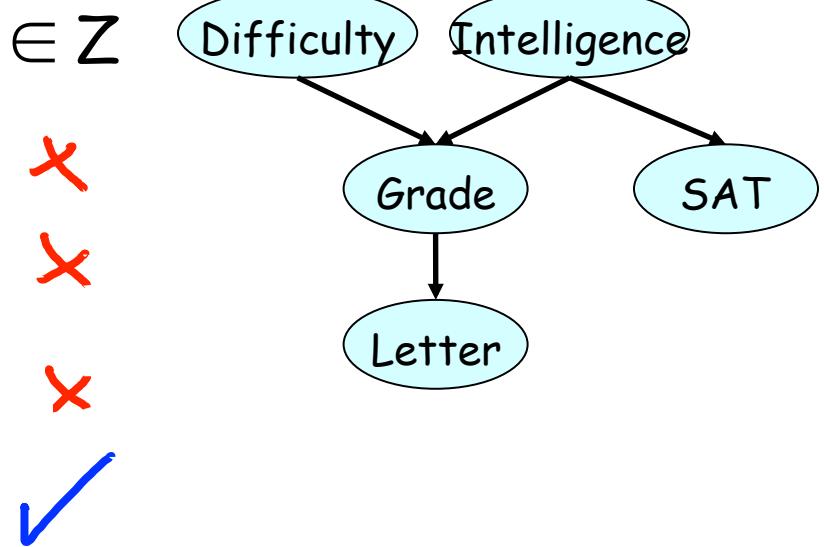
Active Trails

- A trail $X_1 - \dots - X_n$ is active if:
it has no v-structures $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$

When can X influence Y Given evidence about Z

- $X \rightarrow Y$
- $X \leftarrow Y$
- $X \rightarrow W \rightarrow Y$ $W \notin Z$ ✓
- $X \leftarrow W \leftarrow Y$ ✓
- $X \leftarrow W \rightarrow Y$ ✓
- $X \rightarrow W \leftarrow Y$ ✗

$W \in Z$



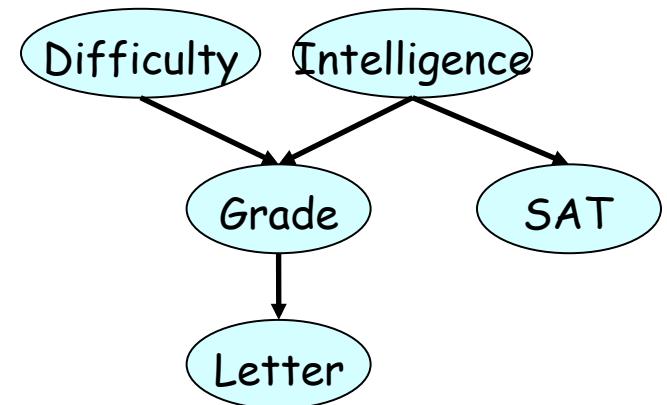
When can X influence Y given evidence about Z

- S – I – G – D allows influence to flow when:

I is observed X

I not observed,
nothing else X

I not observed
& G is observed

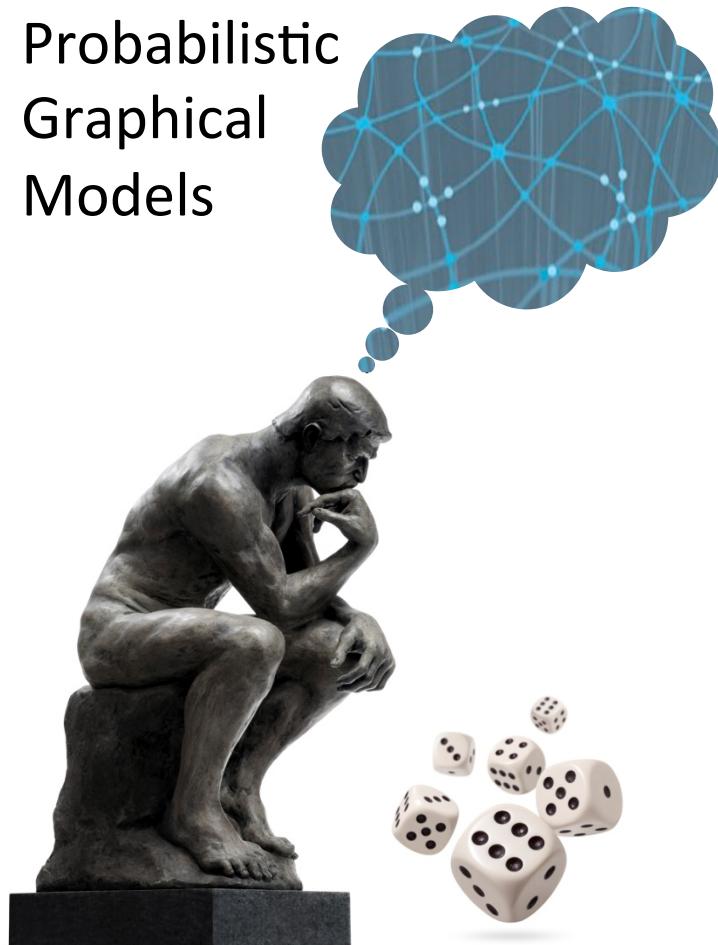


Active Trails

- A trail $X_1 - \dots - X_n$ is active given Z if:

- for any v-structure $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$ we have that X_i or one of its descendants $\in Z$
 - no other X_i is in Z
not in v-structure

Probabilistic
Graphical
Models



Representation

Independencies

Preliminaries

Independence

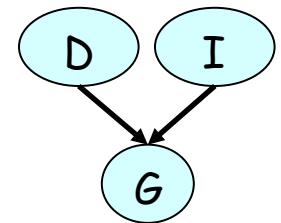
- For events α, β , $P \models \alpha \perp \beta$ if:
 - $P(\alpha \cap \beta) = P(\alpha) \cdot P(\beta)$ *satisfies independence*
 - $P(\alpha | \beta) = P(\alpha)$
 - $P(\beta | \alpha) = P(\beta)$
- For random variables X, Y , $P \models X \perp Y$ if:
 - $P(X, Y) = P(X) P(Y)$
 - $P(X | Y) = P(X)$ *universal*
 $\forall x, y \quad P(x, y) = P(x) \cdot P(y)$
 - $P(Y | X) = P(Y)$

Independence

I	D	G	Prob.
i ⁰	d ⁰	g ¹	0.126
i ⁰	d ⁰	g ²	0.168
i ⁰	d ⁰	g ³	0.126
i ⁰	d ¹	g ¹	0.009
i ⁰	d ¹	g ²	0.045
i ⁰	d ¹	g ³	0.126
i ¹	d ⁰	g ¹	0.252
i ¹	d ⁰	g ²	0.0224
i ¹	d ⁰	g ³	0.0056
i ¹	d ¹	g ¹	0.06
i ¹	d ¹	g ²	0.036
i ¹	d ¹	g ³	0.024

$$P(I, D) =$$

I	D	Prob
i ⁰	d ⁰	0.42
i ⁰	d ¹	0.18
i ¹	d ⁰	0.28
i ¹	d ¹	0.12



P(I)

I	Prob
i ⁰	0.6
i ¹	0.4

P(D)

D	Prob
d ⁰	0.7
d ¹	0.3

Conditional Independence

- For (sets of) random variables X, Y, Z

$P \models (X \perp Y \mid Z)$ if:

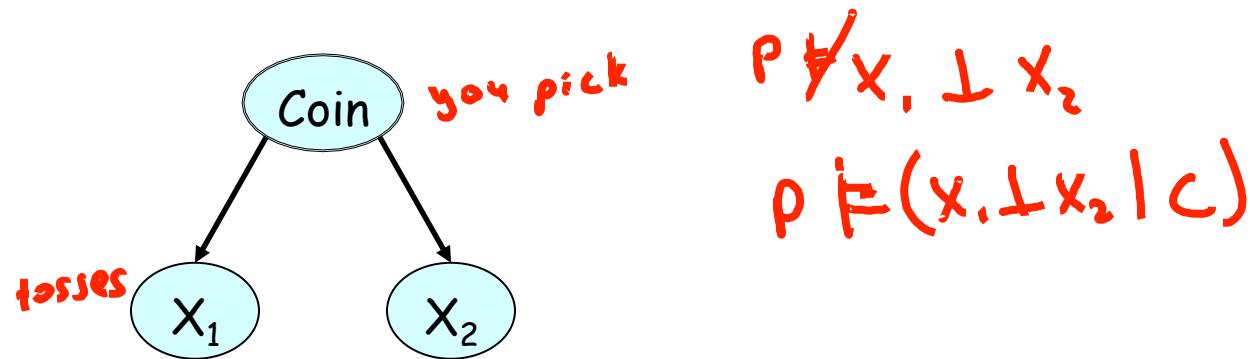
$$- P(X, Y \mid Z) = P(X \mid Z) P(Y \mid Z)$$

$$- P(X \mid Y, Z) = P(X \mid Z)$$

$$- P(Y \mid X, Z) = P(Y \mid Z)$$

$$- P(X, Y, Z) \propto \phi_1(X, Z) \phi_2(Y, Z)$$

Conditional Independence



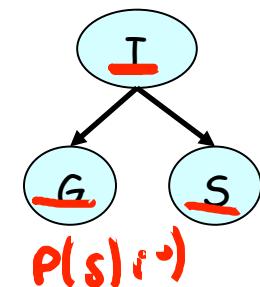
Conditional Independence

$P(I, S, G)$

I	S	G	Prob.
i^0	s^0	g^1	0.114
i^0	s^0	g^2	0.1938
i^0	s^0	g^3	0.2622
i^0	s^1	g^1	0.006
i^0	s^1	g^2	0.0102
i^0	s^1	g^3	0.0138
i^1	s^0	g^1	0.252
i^1	s^0	g^2	0.0224
i^1	s^0	g^3	0.0056
i^1	s^1	g^1	0.108
i^1	s^1	g^2	0.0096
i^1	s^1	g^3	0.0024

$P(S, G | \underline{i^0})$

S	G	Prob.
s^0	g^1	0.19
s^0	g^2	0.323
s^0	g^3	0.437
s^1	g^1	0.01
s^1	g^2	0.017
s^1	g^3	0.023



S	Prob.
s^0	0.95
s^1	0.05

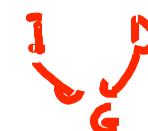
$P(g | i^0)$

G	Prob.
g^1	0.2
g^2	0.34
g^3	0.46

Daphne Koller

Conditioning can Lose Independences

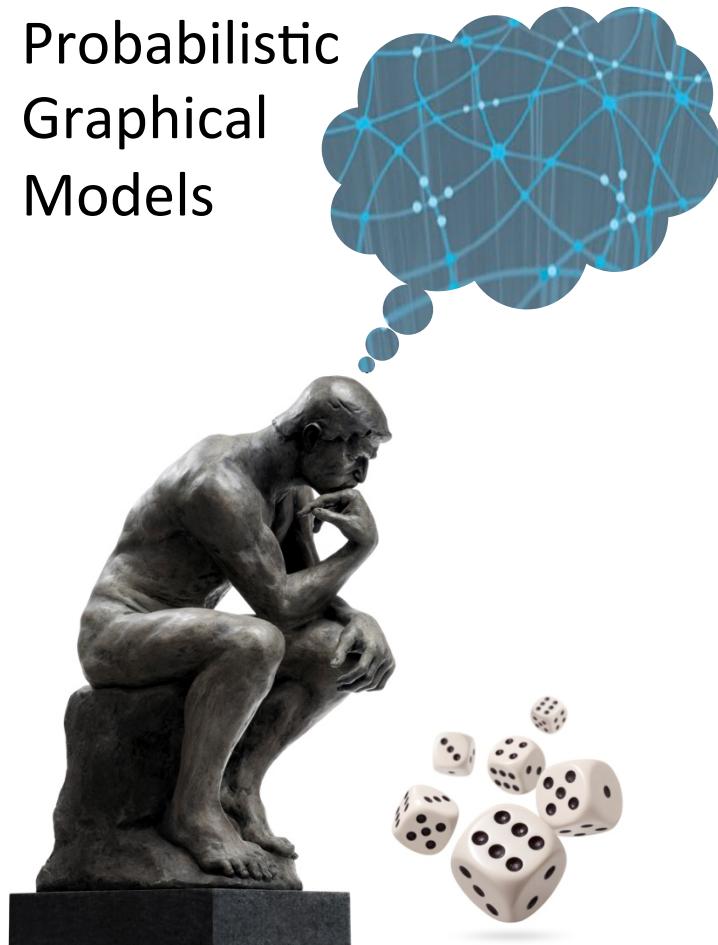
I	D	G	Prob.
i ⁰	d ⁰	g ¹	0.126
i ⁰	d ⁰	g ²	0.168
i ⁰	d ⁰	g ³	0.126
i ⁰	d ¹	g ¹	0.009
i ⁰	d ¹	g ²	0.045
i ⁰	d ¹	g ³	0.126
i ¹	d ⁰	g ¹	0.252
i ¹	d ⁰	g ²	0.0224
i ¹	d ⁰	g ³	0.0056
i ¹	d ¹	g ¹	0.06
i ¹	d ¹	g ²	0.036
i ¹	d ¹	g ³	0.024



$P(I, D | g^1)$

I	D	Prob.
i ⁰	d ⁰	0.282
i ⁰	d ¹	0.02
i ¹	d ⁰	0.564
i ¹	d ¹	0.134

Probabilistic
Graphical
Models



Representation

Independencies

Bayesian
Networks

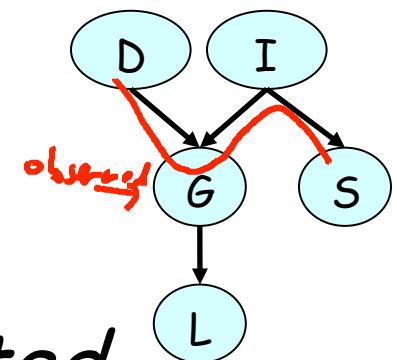
Independence & Factorization

$$P(X,Y) = P(X) P(Y) \quad X, Y \text{ independent}$$

$$P(X,Y,Z) \propto \phi_1(X,Z) \phi_2(Y,Z) \quad (X \perp Y \mid Z)$$

- Factorization of a distribution P implies independencies that hold in P
- If P factorizes over G , can we read these independencies from the structure of G ?

Flow of influence & d-separation



Definition: X and Y are d -separated in G given Z if there is no active trail in G between X and Y given Z

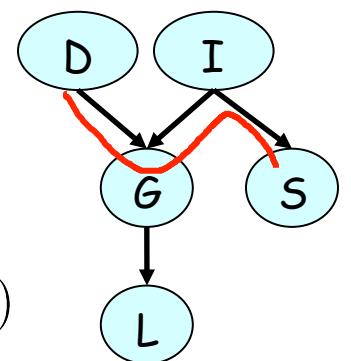
Notation: $d\text{-sep}_G(X, Y \mid Z)$

Factorization \Rightarrow Independence: BNs

Theorem: If P factorizes over G , and $d\text{-sep}_G(X, Y \mid Z)$
 then P satisfies $(X \perp Y \mid Z)$

$$P(D, I, G, S, L) = P(D)P(I)P(G \mid D, I)P(S \mid I)P(L \mid G) \quad \text{chain rule}$$

$P \models D \perp S$



$$\begin{aligned} P(D, S) &= \sum_{\substack{G, I}} P(D)P(I)P(G \mid D, I)P(S \mid I)P(L \mid G) \\ &= \sum_I P(D)P(I)P(S \mid I) \sum_G (P(G \mid D, I) \sum_L P(L \mid G)) \\ &= P(D) \left(\sum_I P(I)P(S \mid I) \right) \end{aligned}$$

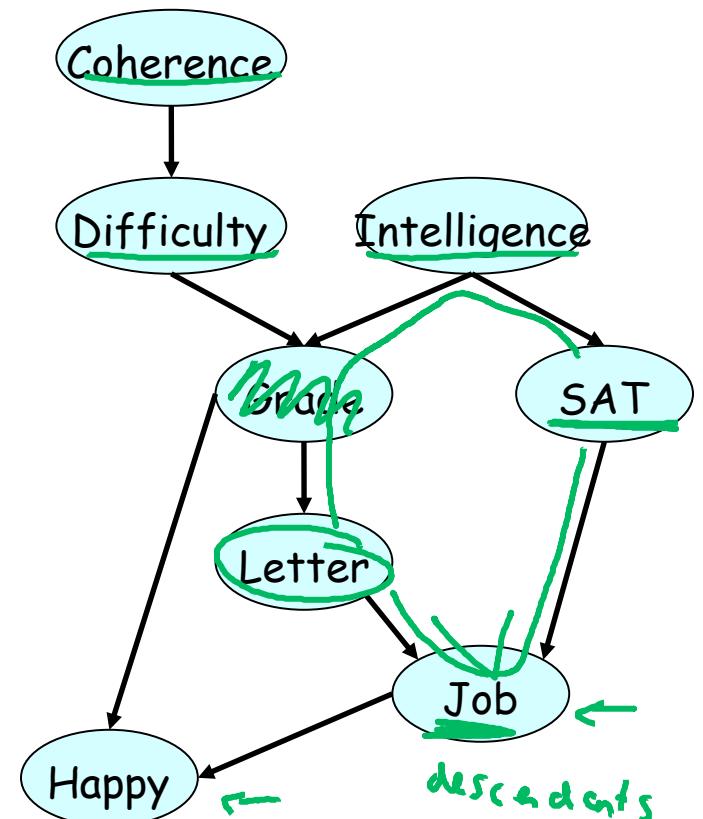
$\cancel{\phi_2(s)}$

J

Any node is d-separated from its non-descendants given its parents



If P factorizes over G, then in P, any variable is independent of its non-descendants given its parents



Daphne Koller

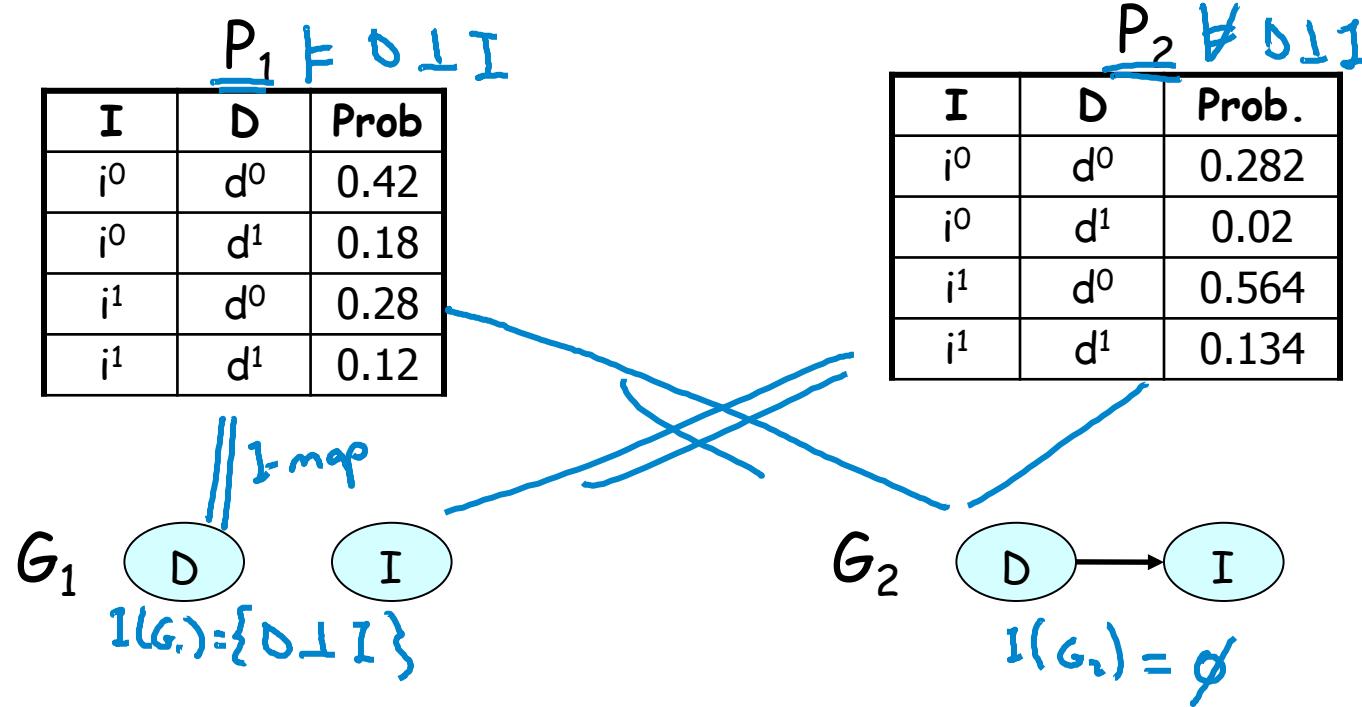
I-maps

- d-separation in $G \Rightarrow P$ satisfies corresponding independence statement

$$I(G) = \{(\underline{X} \perp \underline{Y} \mid \underline{Z}) : \text{d-sep}_G(\underline{X}, \underline{Y} \mid \underline{Z})\}$$

- Definition: If P satisfies $I(G)$, we say that G is an I-map (independency map) of P

I-maps



Factorization \Rightarrow Independence: BNs

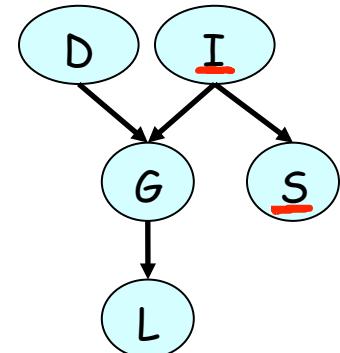
Theorem: If P factorizes over G, then G is
an I-map for P

Can read from G independencies in P
regardless of parameters

Independence \Rightarrow Factorization

Theorem: If G is an I-map for P, then P factorizes over G

IID



P(I,D) chain rule for probabilities

$$P(D, I, G, S, L) = \underbrace{P(D)}_{\text{P(I,D) chain rule for probabilities}} \underbrace{P(I | D)}_{\cancel{\text{P(I,D)}}} \underbrace{P(G | D, I)}_{\cancel{\text{P(G,I,D)}}} \underbrace{P(S | D, I, G)}_{\cancel{\text{P(S,I,G,D)}}} \underbrace{P(L | D, I, G, S)}_{\cancel{\text{P(L,I,G,D,S)}}}$$

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

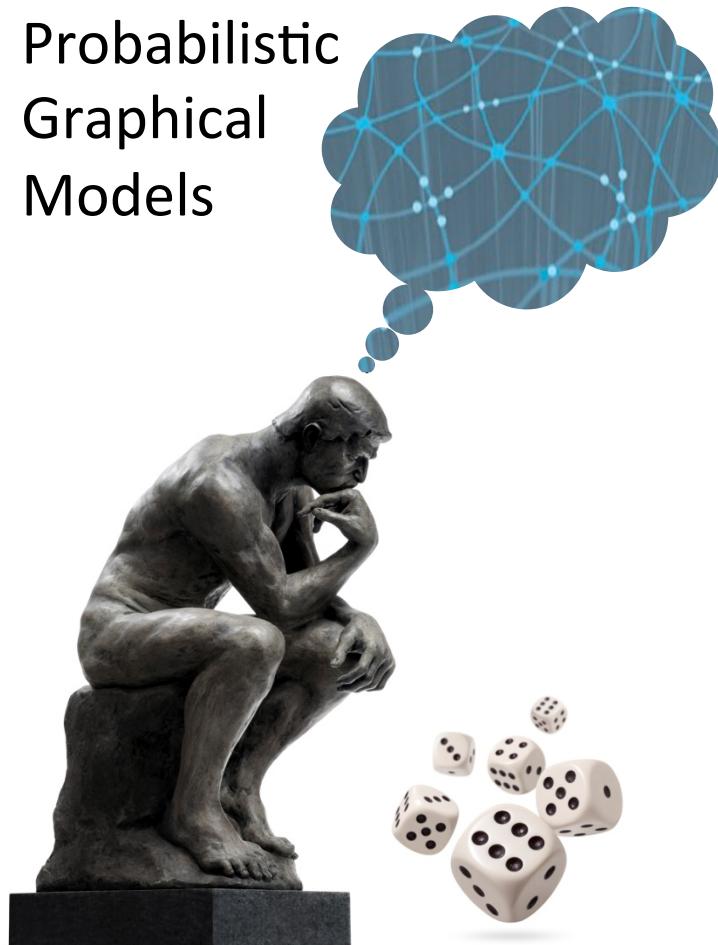
Summary

Two equivalent views of graph structure:

- Factorization: G allows P to be represented
- I-map: Independencies encoded by G hold in P

If P factorizes over a graph G , we can read from the graph independencies that must hold in P (an independency map)

Probabilistic
Graphical
Models

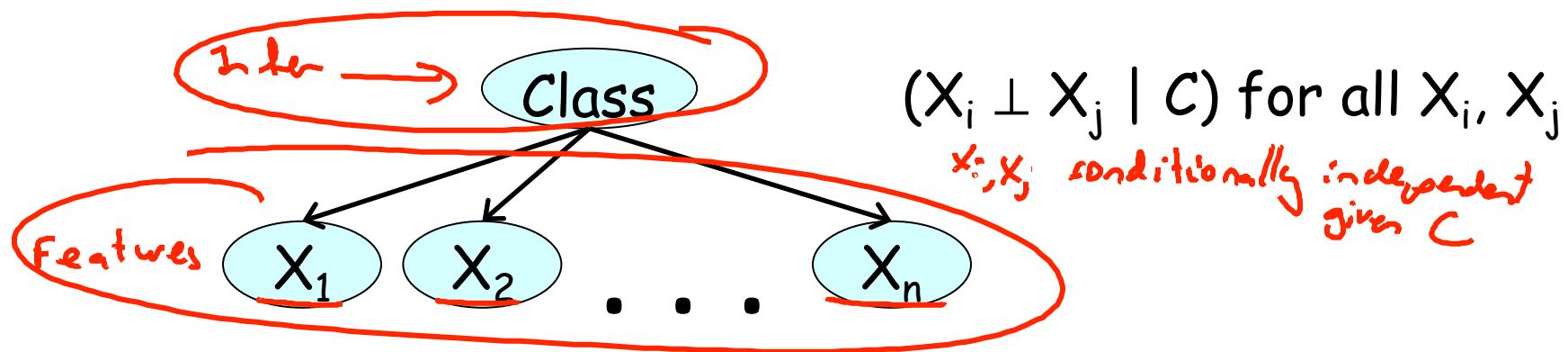


Representation

Bayesian Networks

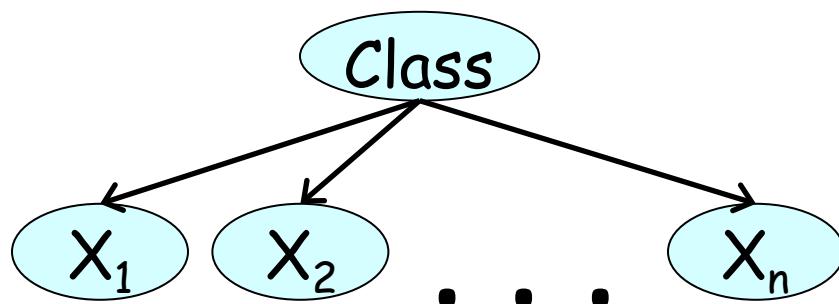
Naïve Bayes

Naïve Bayes Model



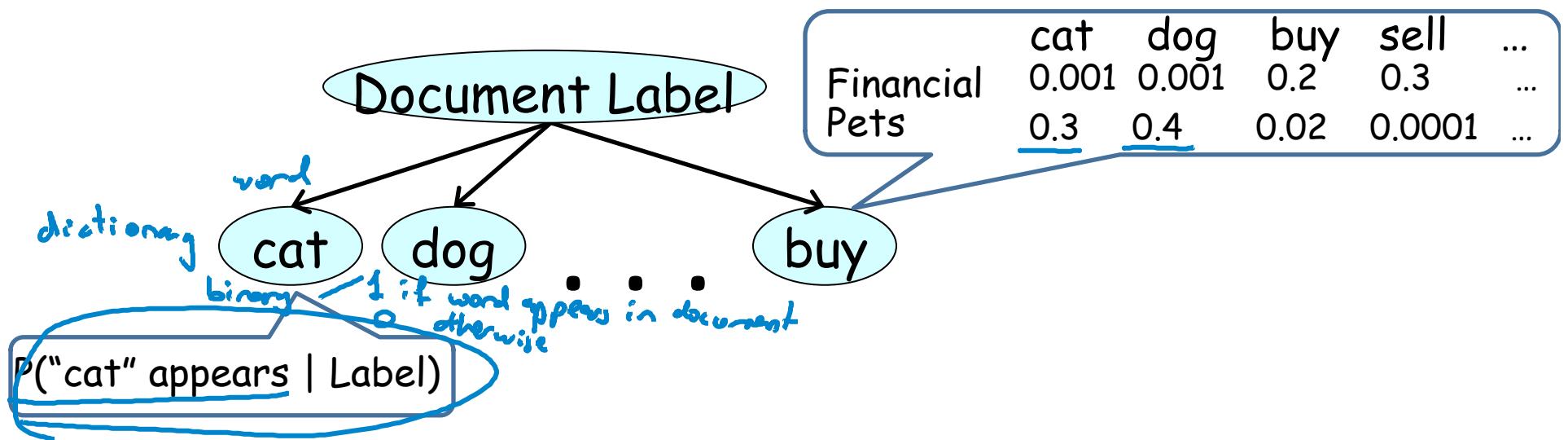
$$\underline{P(C, X_1, \dots, X_n)} = \underbrace{P(C)}_{i=1} \prod_{i=1}^n P(X_i \mid C)$$

Naïve Bayes Classifier



$$\frac{P(C = c^1 \mid x_1, \dots, x_n)}{P(C = c^2 \mid x_1, \dots, x_n)} = \underbrace{\frac{P(C = c^1)}{P(C = c^2)}}_{\text{odds ratios}} \prod_{i=1}^n \frac{P(x_i \mid C = c^1)}{P(x_i \mid C = c^2)}$$

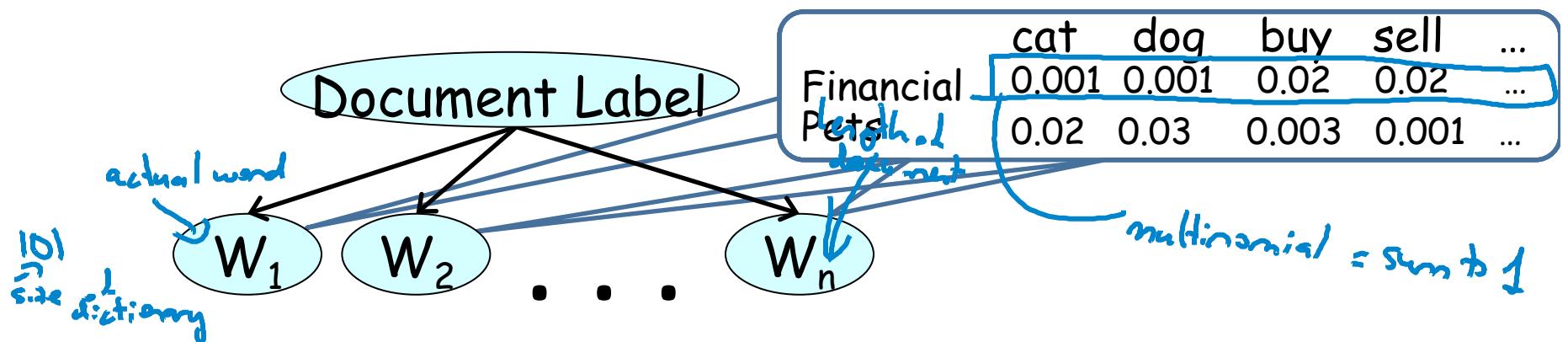
Bernoulli Naïve Bayes for Text



$$\frac{P(C = c^1 \mid x_1, \dots, x_n)}{P(C = c^2 \mid x_1, \dots, x_n)} = \frac{P(C = c^1)}{P(C = c^2)} \prod_{i=1}^n \frac{P(x_i \mid C = c^1)}{P(x_i \mid C = c^2)}$$

Daphne Koller

Multinomial Naïve Bayes for Text



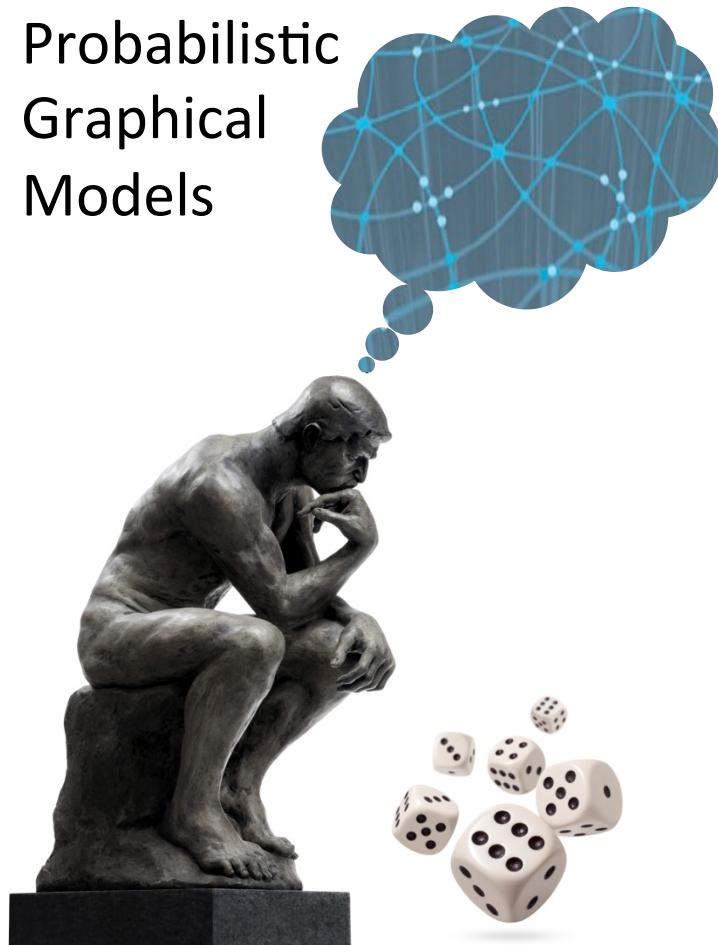
$$\frac{P(C = c^1 \mid x_1, \dots, x_n)}{P(C = c^2 \mid x_1, \dots, x_n)} = \frac{P(C = c^1)}{P(C = c^2)} \prod_{i=1}^n \frac{P(x_i \mid C = c^1)}{P(x_i \mid C = c^2)}$$

Daphne Koller

Summary

- Simple approach for classification
 - Computationally efficient
 - Easy to construct
- Surprisingly effective in domains with many weakly relevant features
- Strong independence assumptions reduce performance when many features are strongly correlated

Probabilistic
Graphical
Models



Representation

Bayesian Networks

Application:
Diagnosis

Medical Diagnosis: Pathfinder (1992)

- Help pathologist diagnose lymph node pathologies (60 different diseases)
- Pathfinder I: Rule-based system
- Pathfinder II used naïve Bayes and got superior performance

Heckerman et al.

Daphne Koller

Medical Diagnosis: Pathfinder (1992)

- Pathfinder III: Naïve Bayes with better knowledge engineering
- No incorrect zero probabilities
- Better calibration of conditional probabilities
 - $P(\text{finding} \mid \text{disease}_1)$ to $P(\text{finding} \mid \text{disease}_2)$
 - Not $P(\text{finding}_1 \mid \text{disease})$ to $P(\text{finding}_2 \mid \text{disease})$

Heckerman et al.

Daphne Koller

Medical Diagnosis: Pathfinder (1992)

- Pathfinder IV: Full Bayesian network
 - Removed incorrect independencies
 - Additional parents led to more accurate estimation of probabilities
- BN model agreed with expert panel in 50/53 cases, vs 47/53 for naïve Bayes model
- Accuracy as high as expert that designed the model

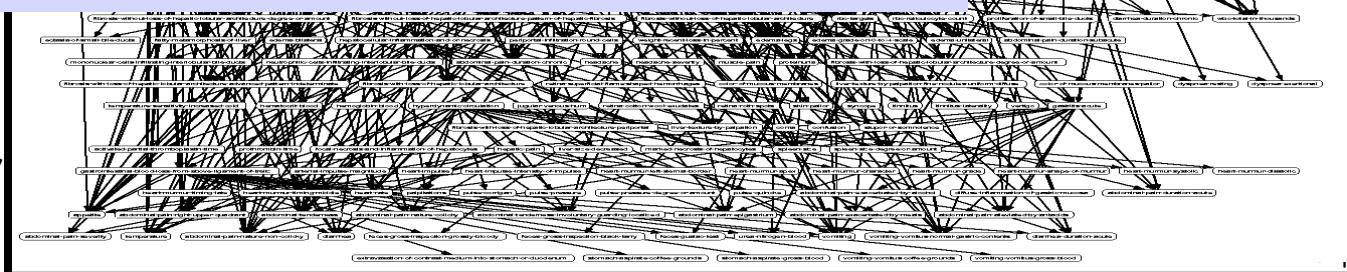
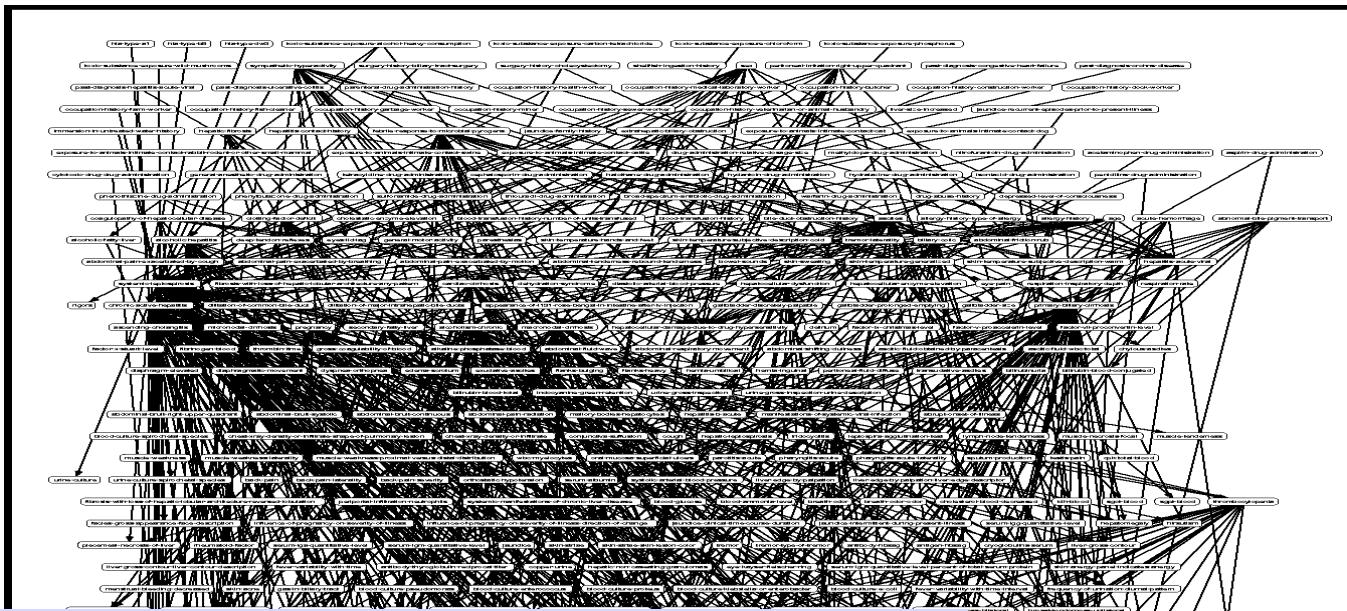
Heckerman et al.

Daphne Koller

CPCS

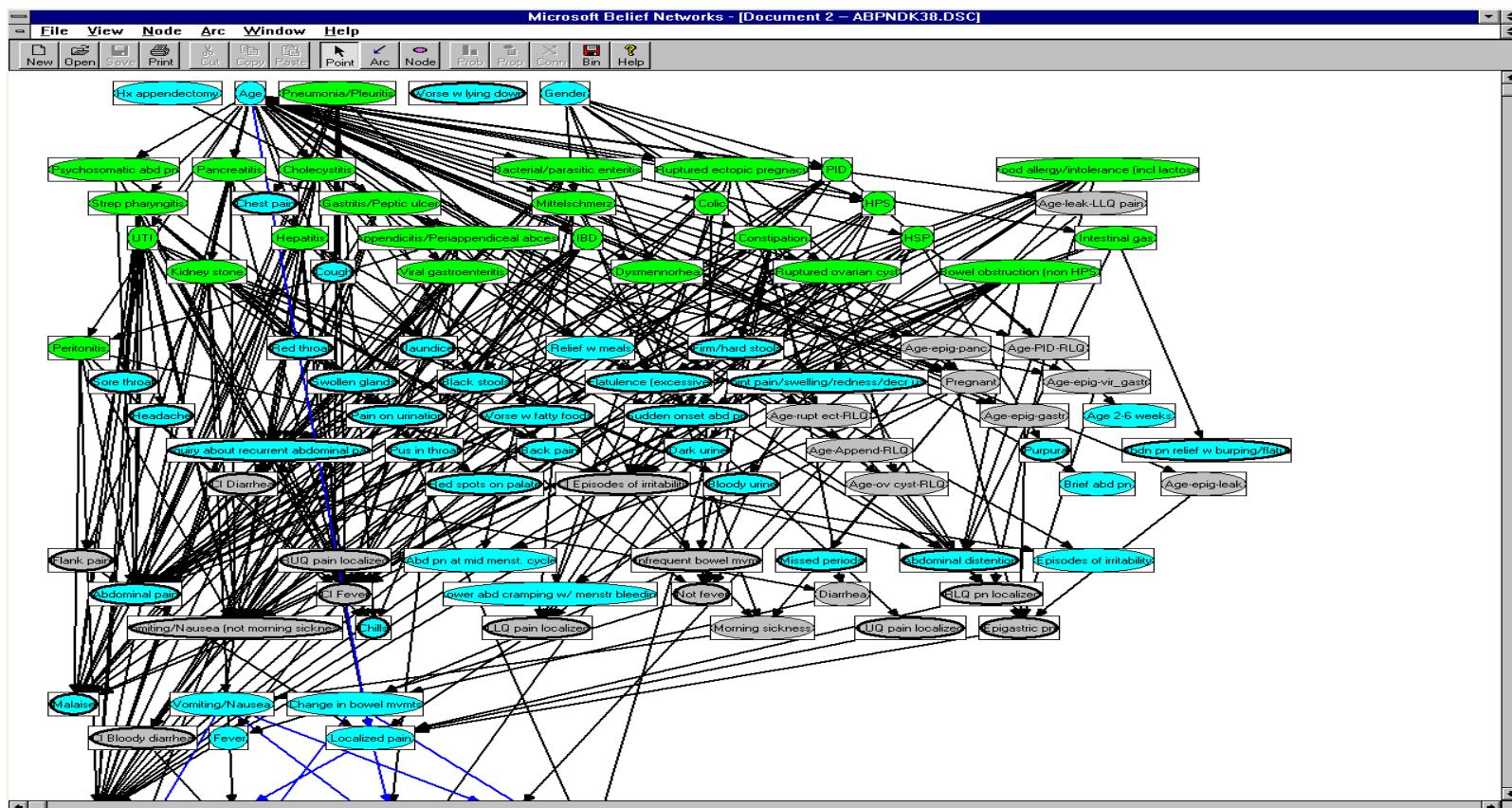
of parameters:
 2^{1000} to 133,931,430 to 8254

M. Pradhan , G. Provan ,
B. Middleton , M. Henrion,
UAI 94



hne Koller

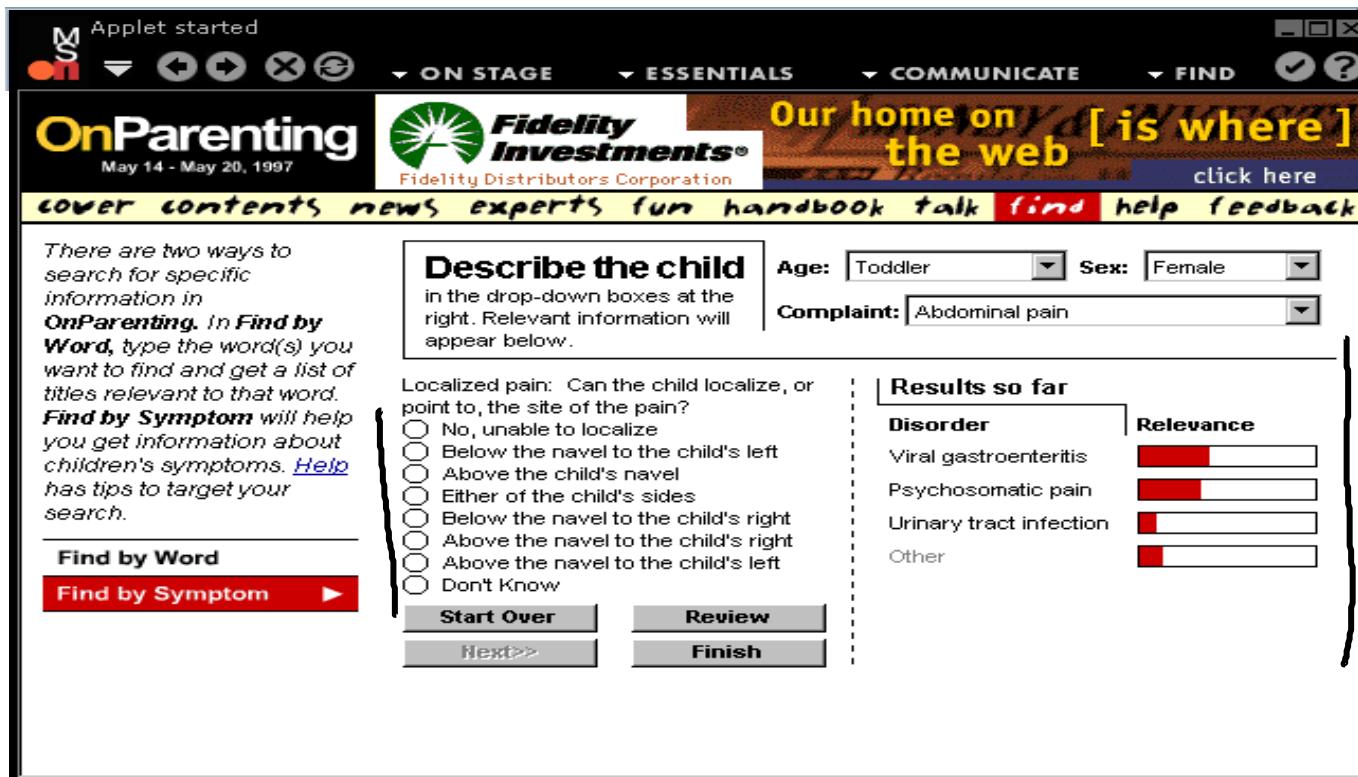
Medical Diagnosis (Microsoft)



Thanks to: Eric Horvitz, Microsoft Research

Daphne Koller

Medical Diagnosis (Microsoft)

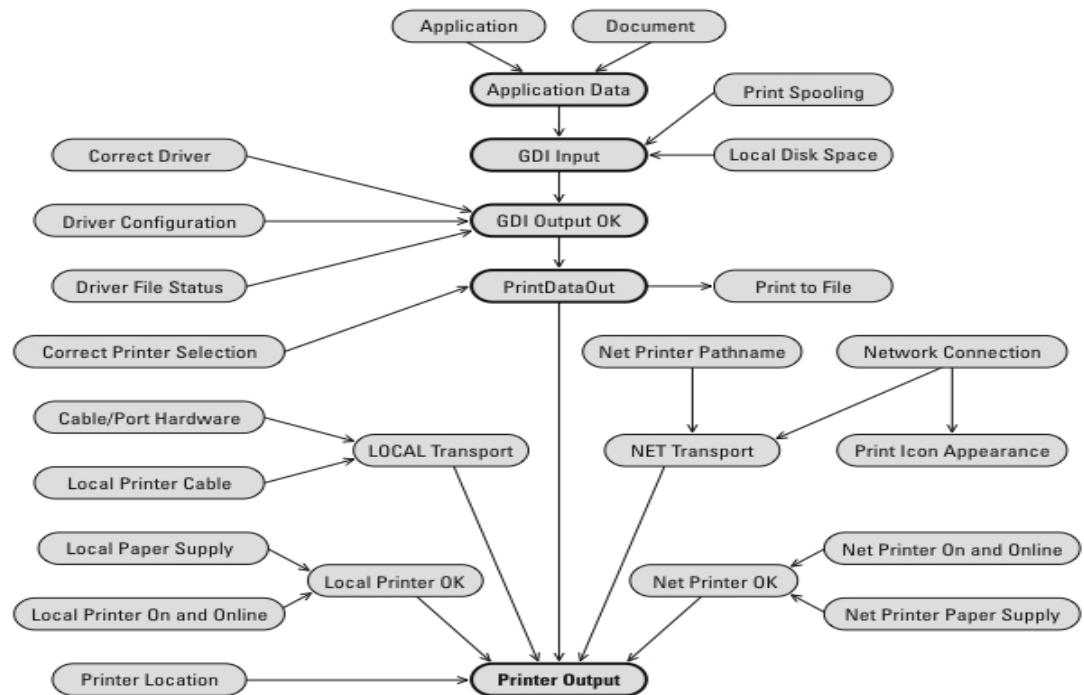


Thanks to: Eric Horvitz, Microsoft Research

Daphne Koller

Fault Diagnosis

- Microsoft troubleshooters



Daphne Koller

Fault Diagnosis

- Many examples:
 - Microsoft troubleshooters
 - Car repair
- Benefits:
 - Flexible user interface
 - Easy to design and maintain ←