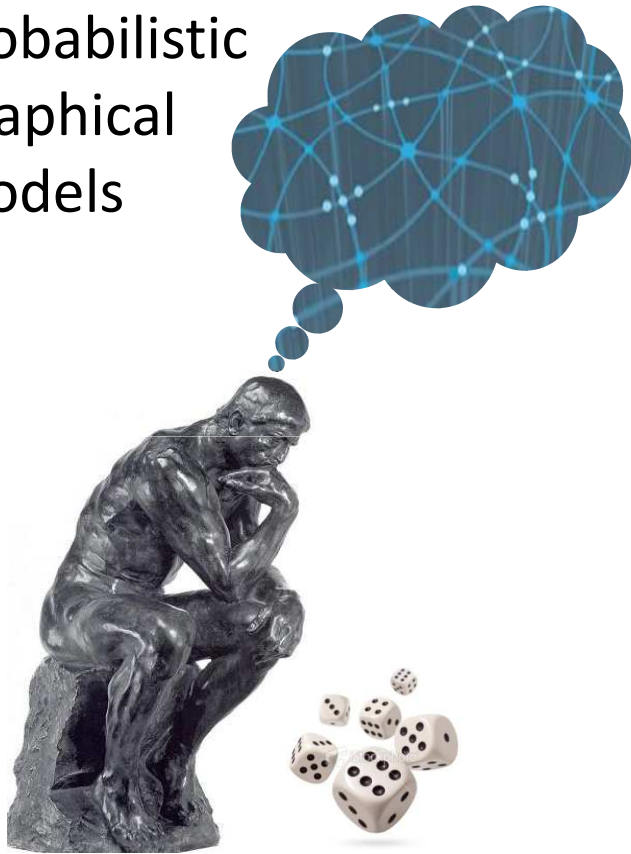


Probabilistic
Graphical
Models



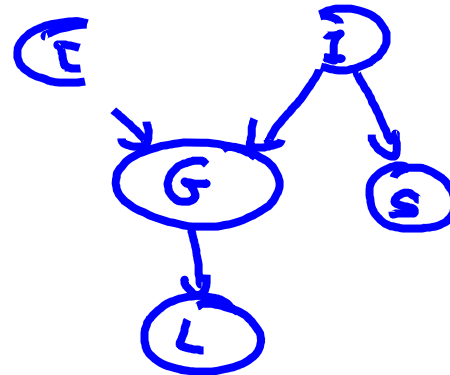
Representation

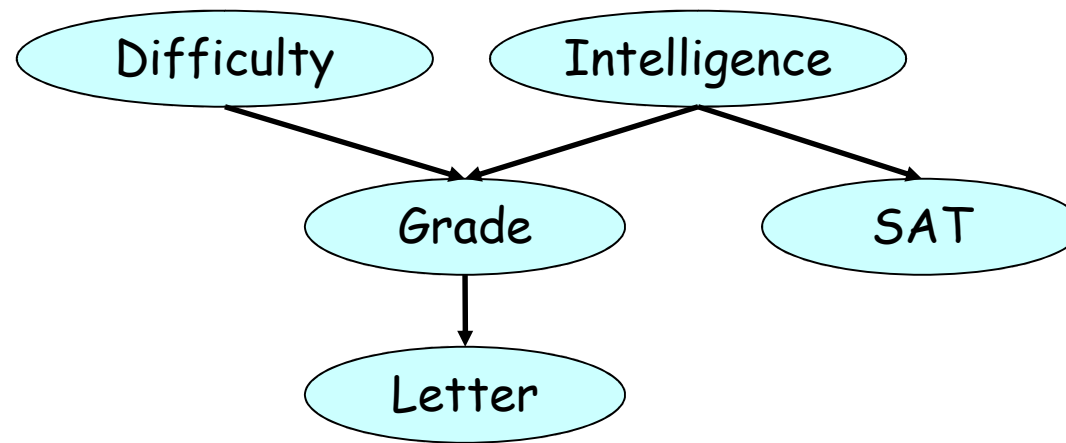
Bayesian Networks

Semantics &
Factorization

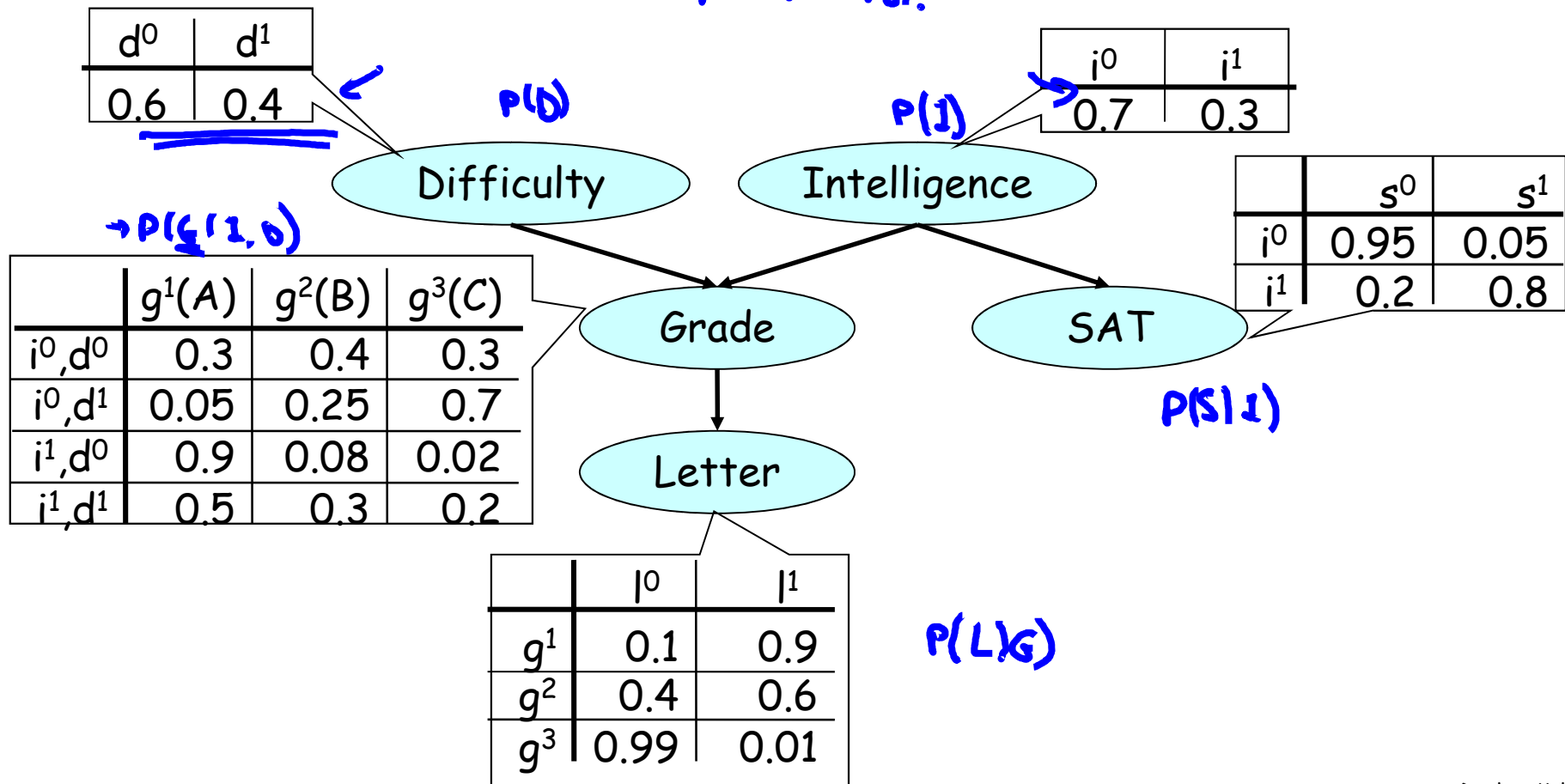
- Grade
- Course Difficulty
- Student Intelligence
- Student SAT
- Reference Letter

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

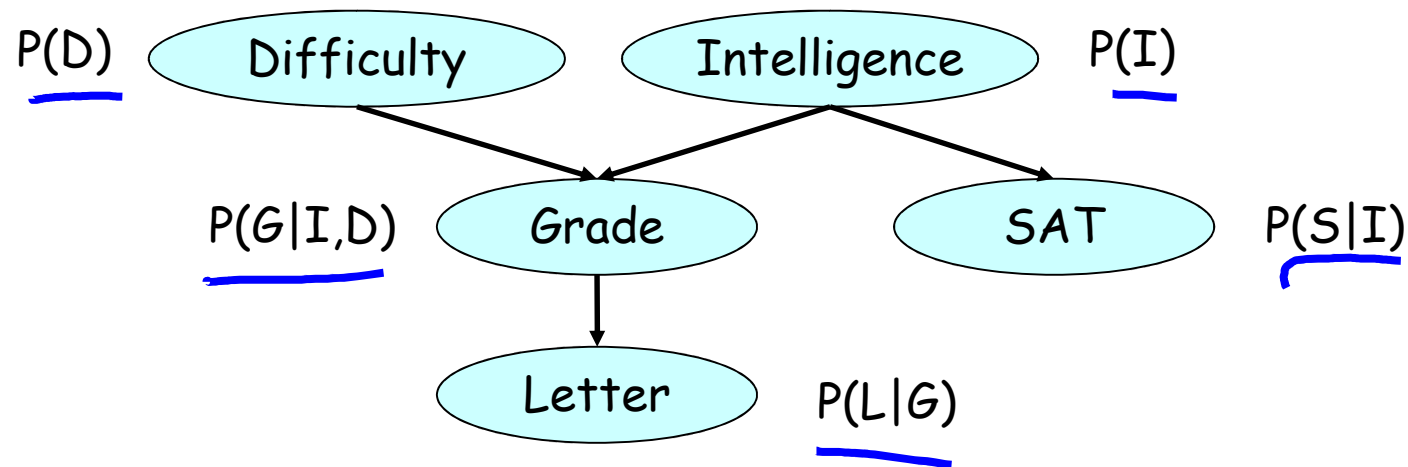




CPD = cond. prob. dist.

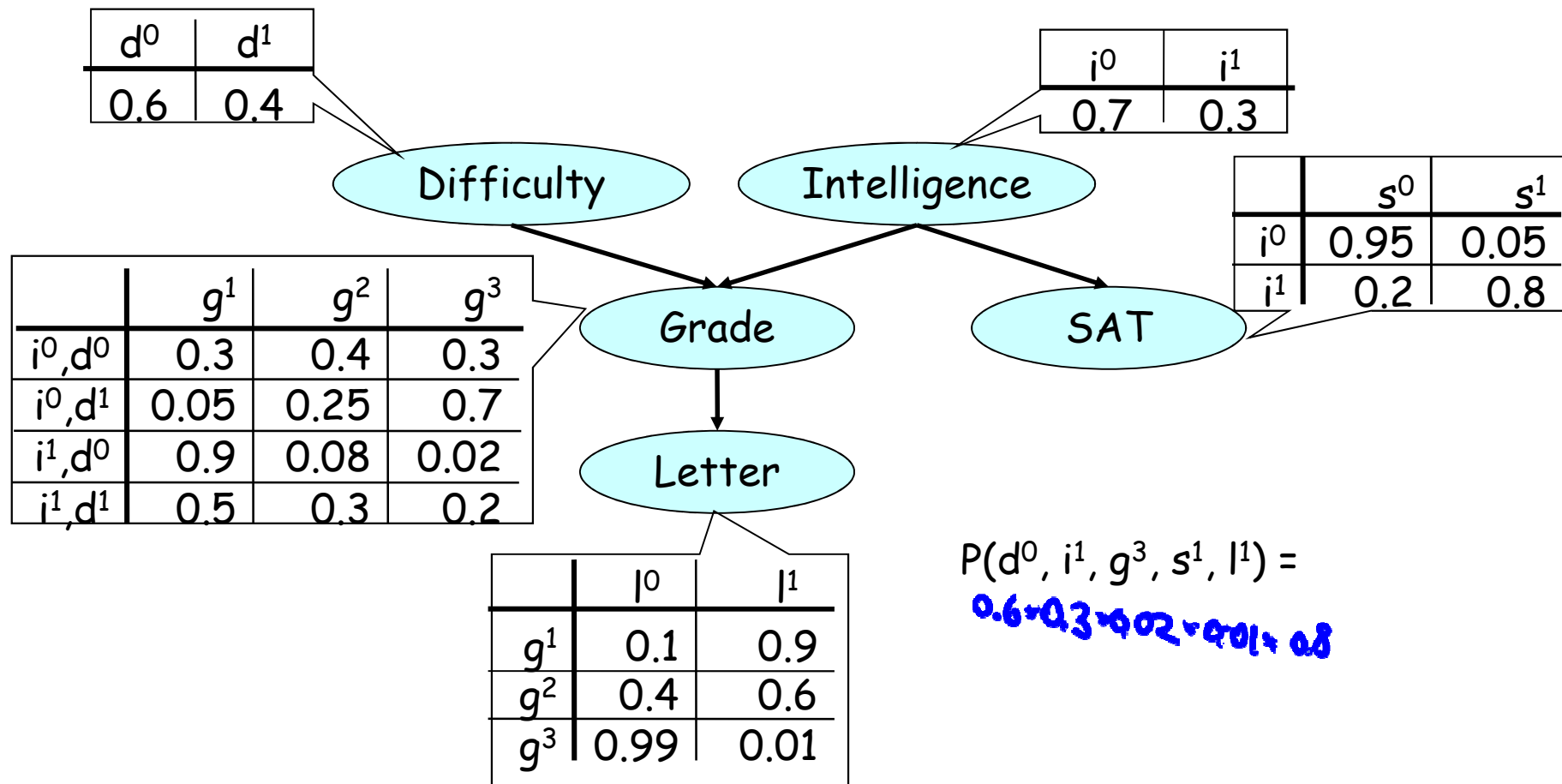


Chain Rule for Bayesian Networks



$$\underline{P(D,I,G,S,L)} = P(D) P(I) P(G|I,D) P(S|I) P(L|G)$$

Distribution defined as a product of factors!

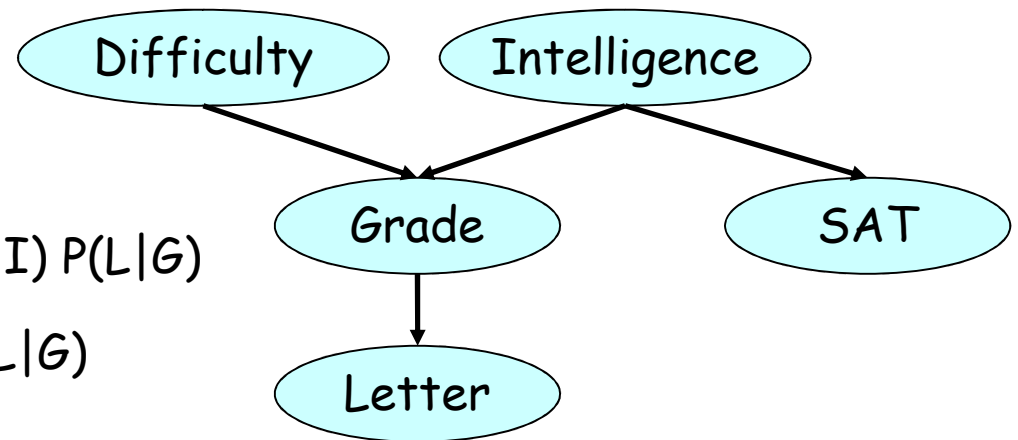


$$P(d^0, i^1, g^3, s^1, l^1) =$$

$$0.6 \times 0.3 \times 0.02 \times 0.01 \times 0.8$$

Defining a joint distribution

- What is the joint distribution $P(D, I, G, S, L)$?



- ☐ $P(D) P(I) P(G|I) P(G|D) P(S|I) P(L|G)$
- ☐ $P(D) P(I) P(G|I, D) P(S|I) P(L|G)$
- ☐ $P(D) P(I) P(G) P(S) P(L)$
- ☐ $P(D|G) P(I|D) P(S|I) P(G|L, I, D) P(L|G)$

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 X_i a CPD $P(X_i \mid \text{Par}_G(X_i))$
- 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: $\sum P = 1$

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

What is the value of $\sum_L P(L | G)$

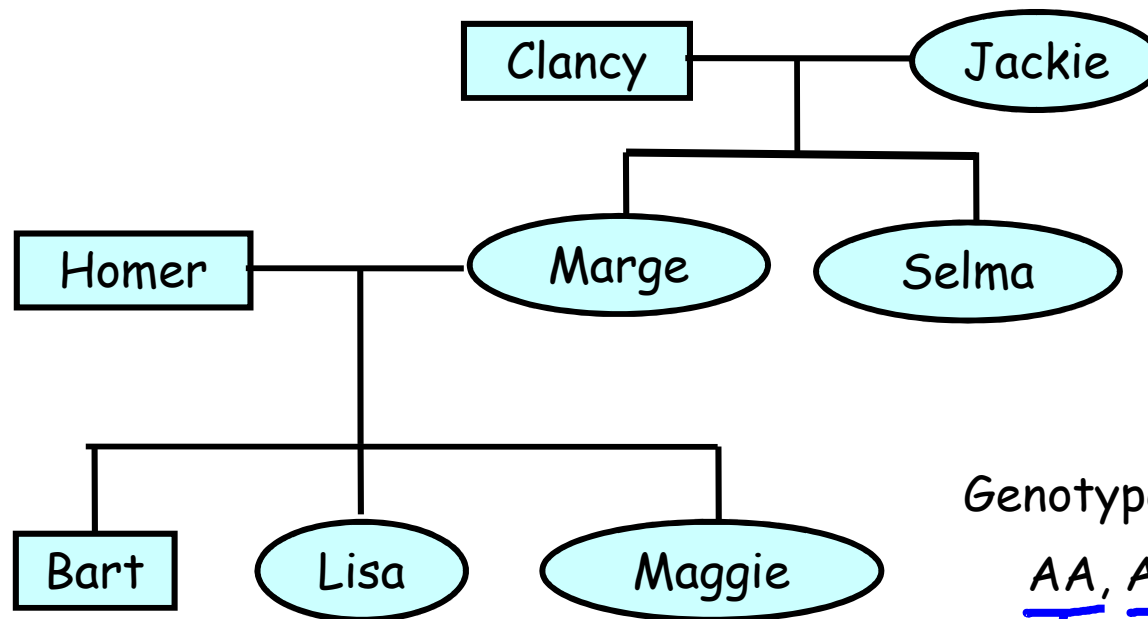
- ☐ 1
- ☐ $P(L)$
- ☐ $P(G)$
- ☐ None of the above

P Factorizes over G

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

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

Genetic Inheritance



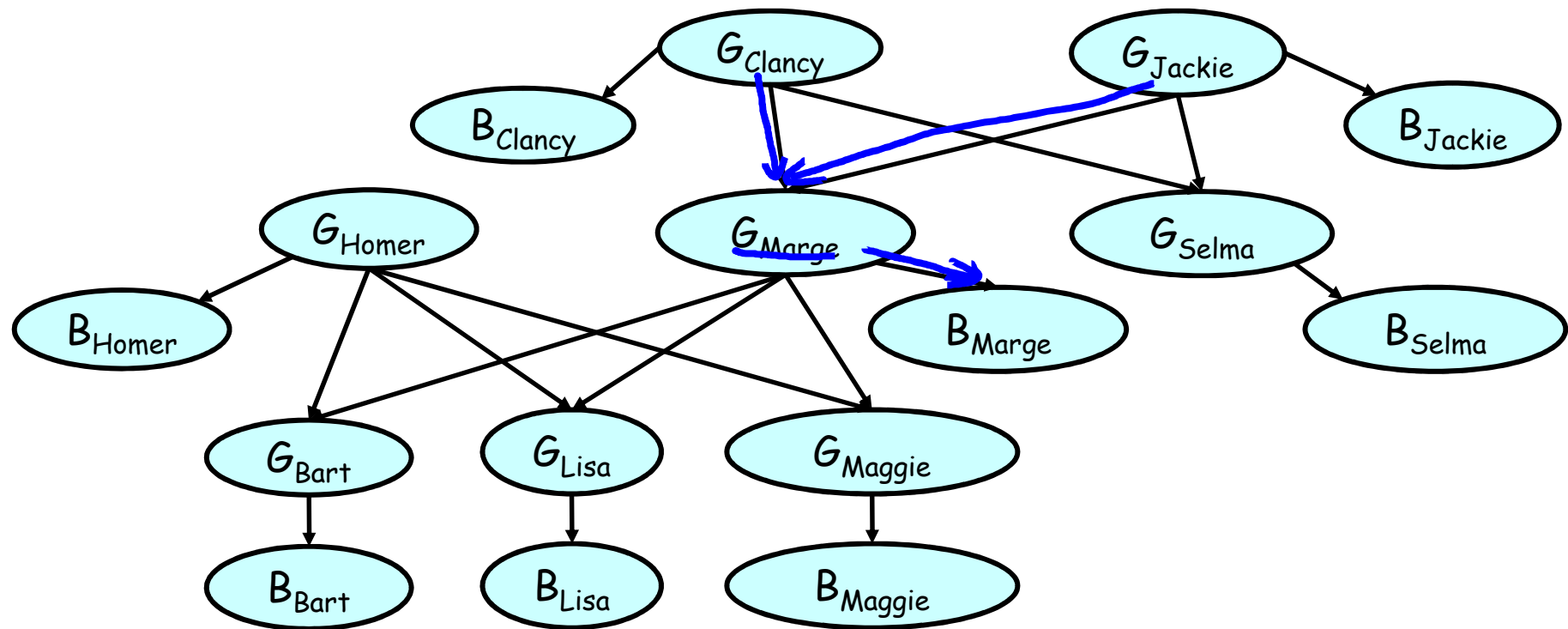
Genotype

AA, AB, AO, BO, BB, OO

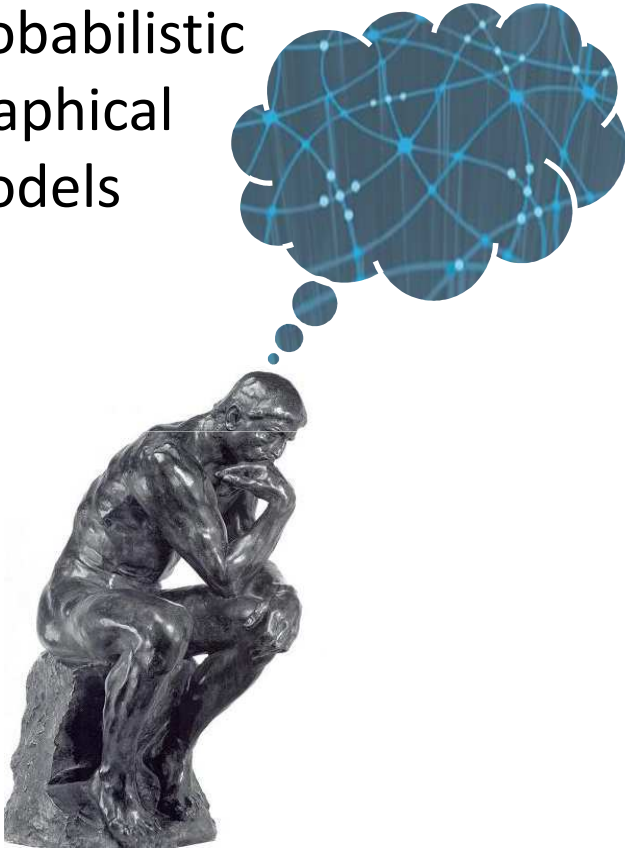
Phenotype

A, B, AB, O

BNs for Genetic Inheritance



Probabilistic
Graphical
Models



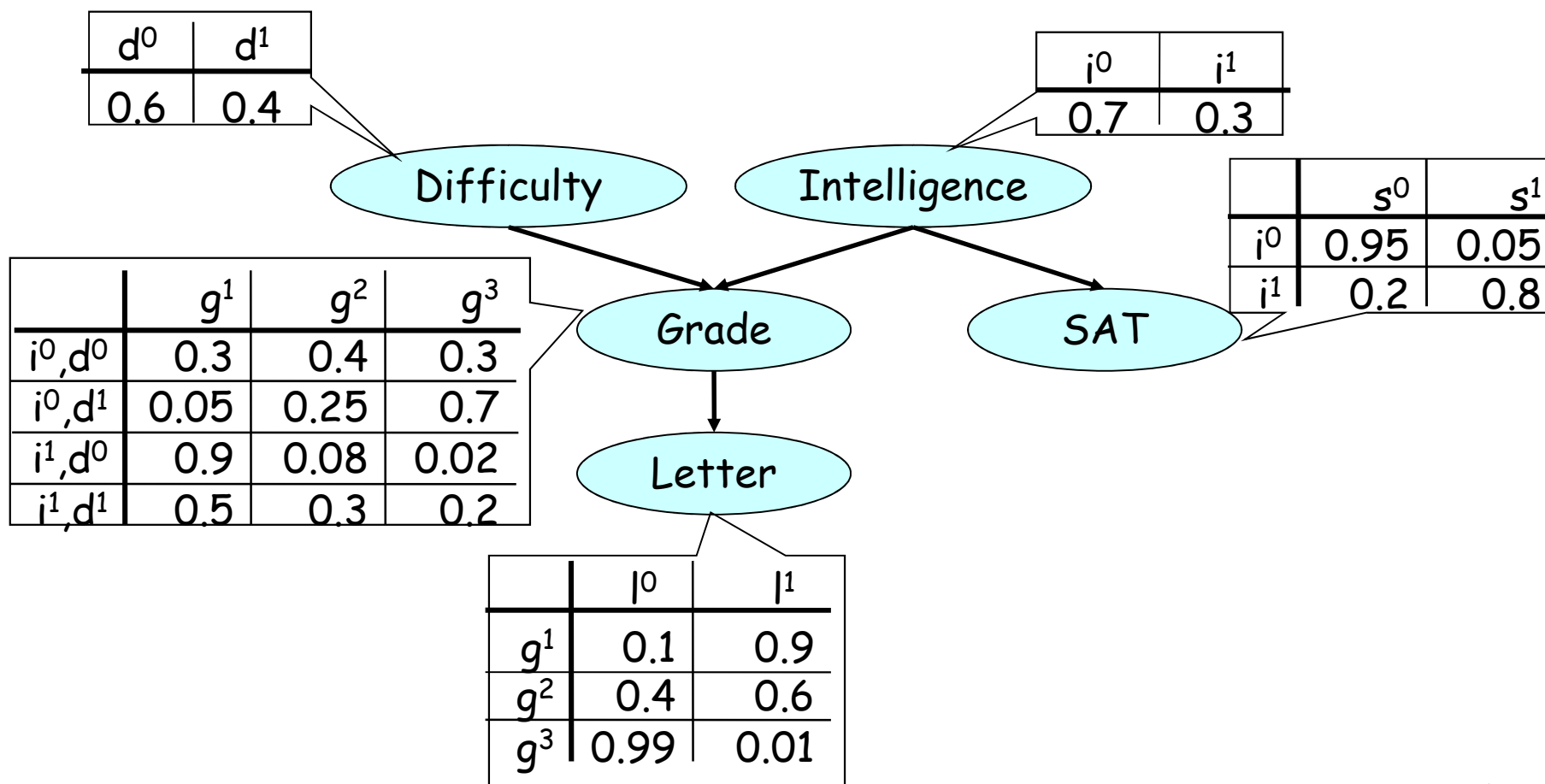
Representation

Bayesian Networks

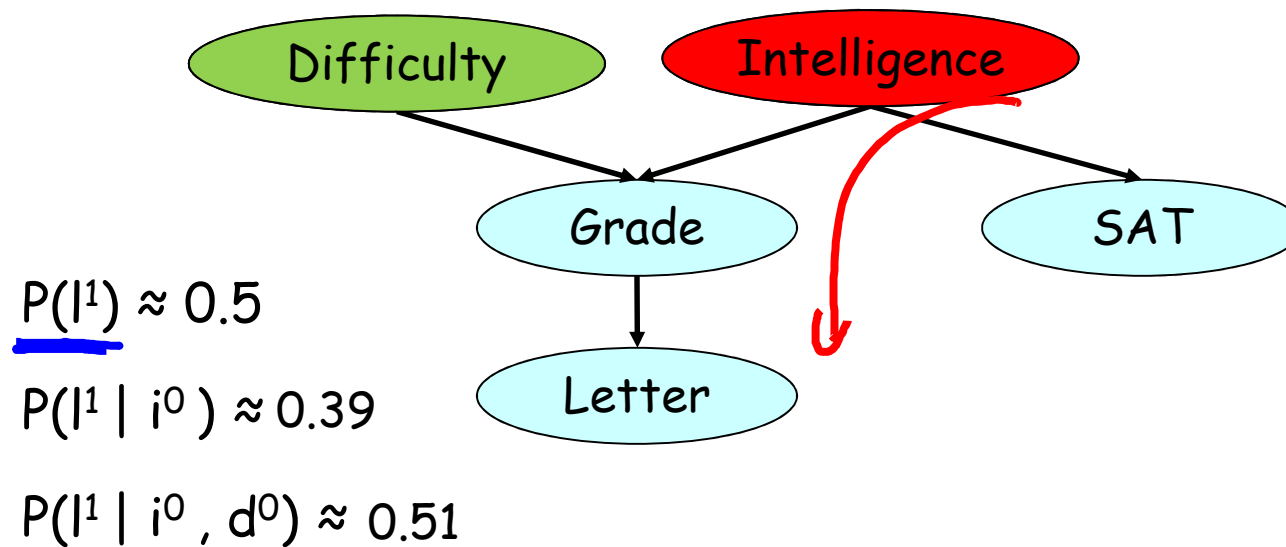
Reasoning

Patterns

The Student Network



Causal Reasoning



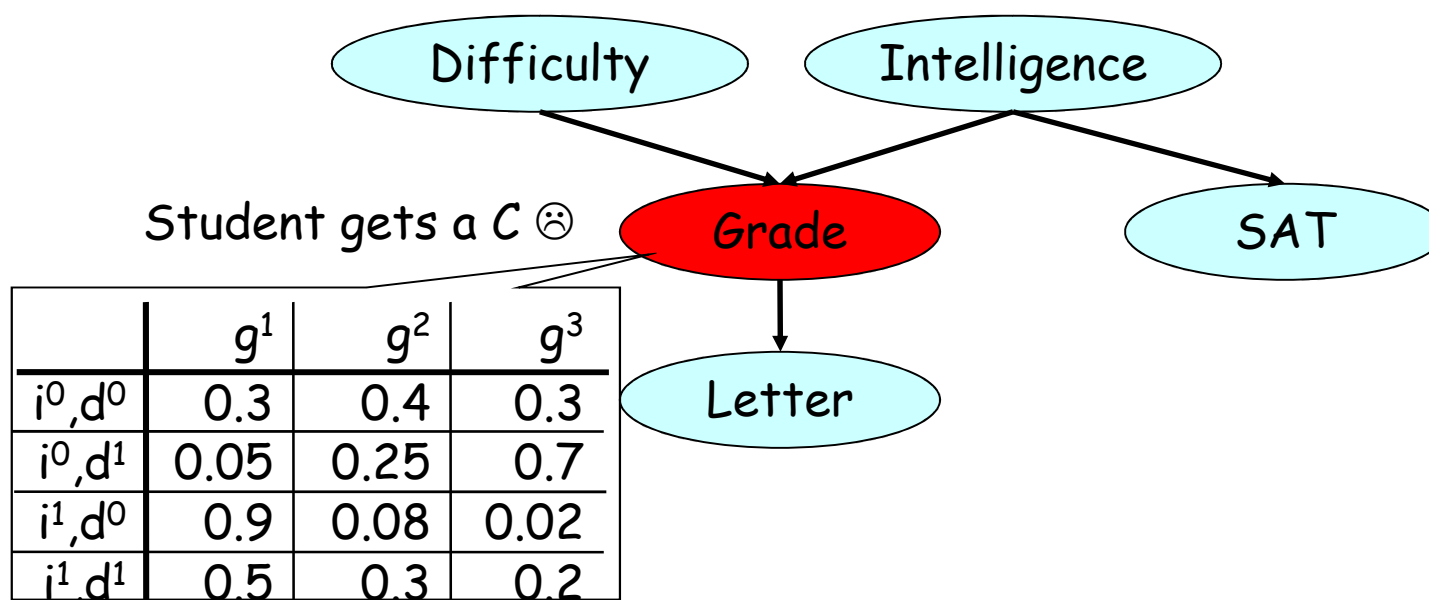
Evidential Reasoning

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

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

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

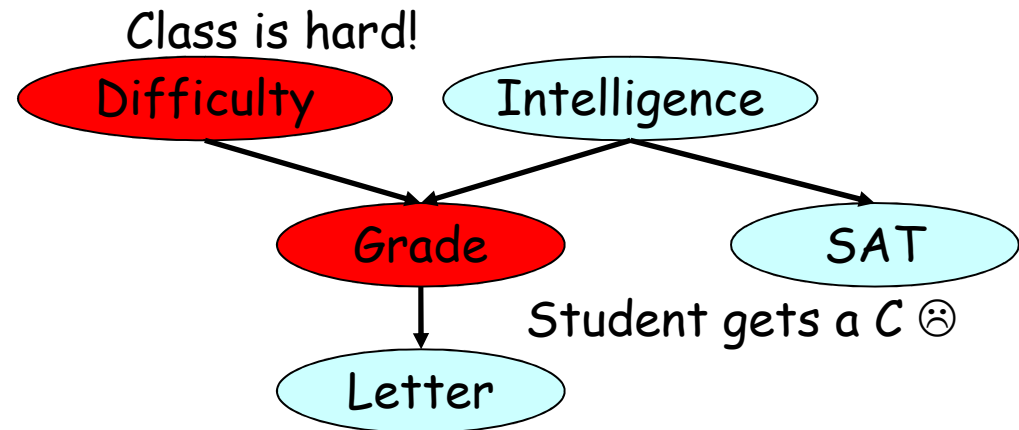
$$P(i^1 \mid g^3) \approx 0.08$$



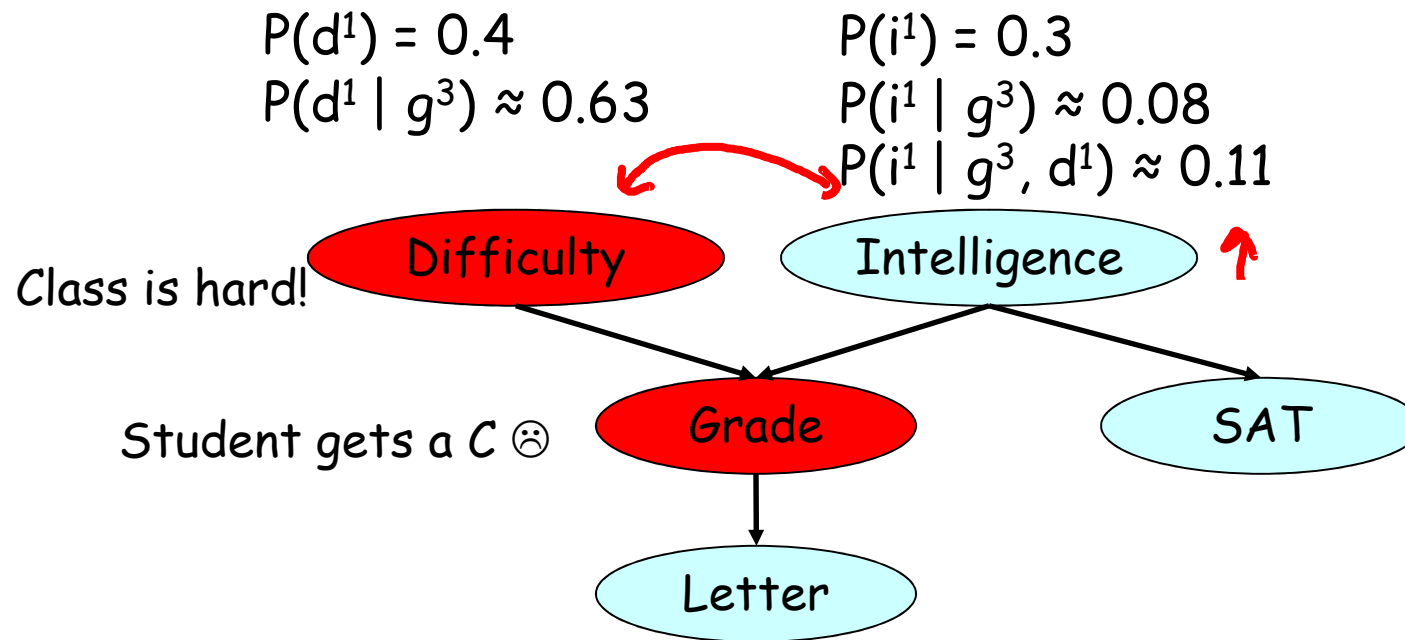
We find out that class is hard

- What happens to the posterior probability of high intelligence?

- ☐ Goes up
- ☐ Goes down
- ☐ Doesn't change
- ☐ We can't know

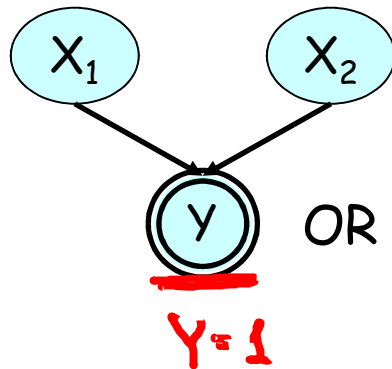


Intercausal Reasoning



Intercausal Reasoning Explained

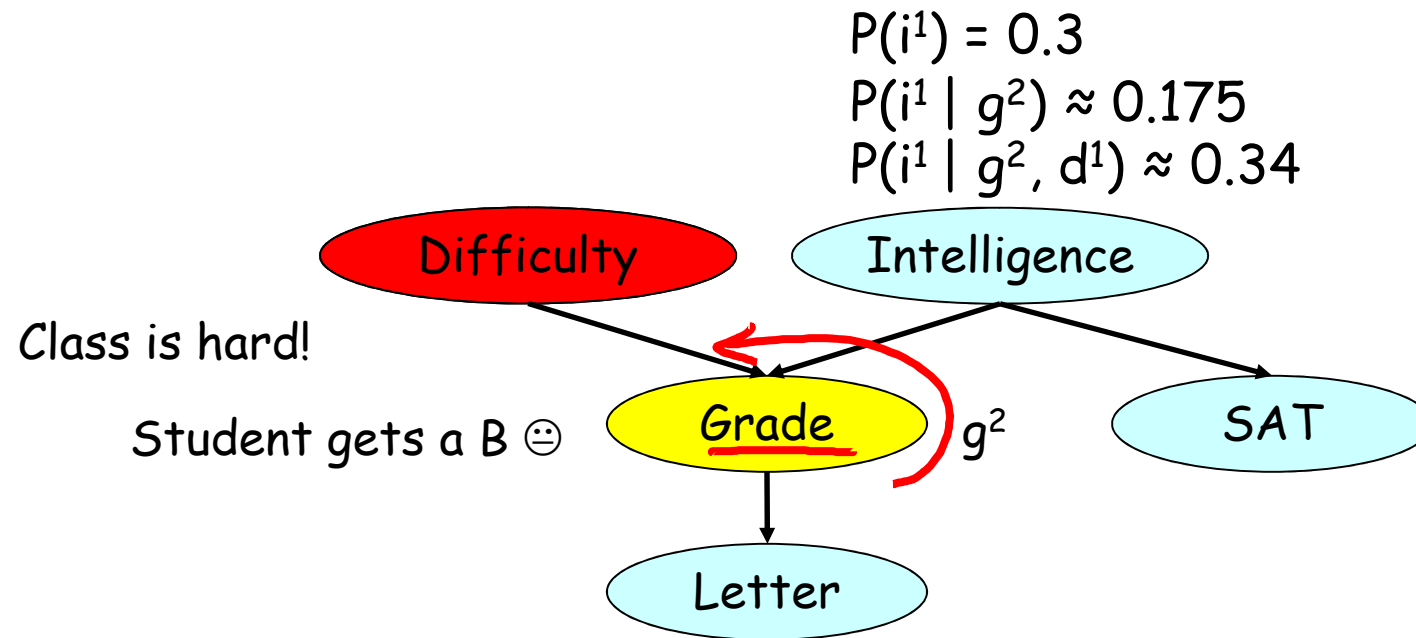
explaining away



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

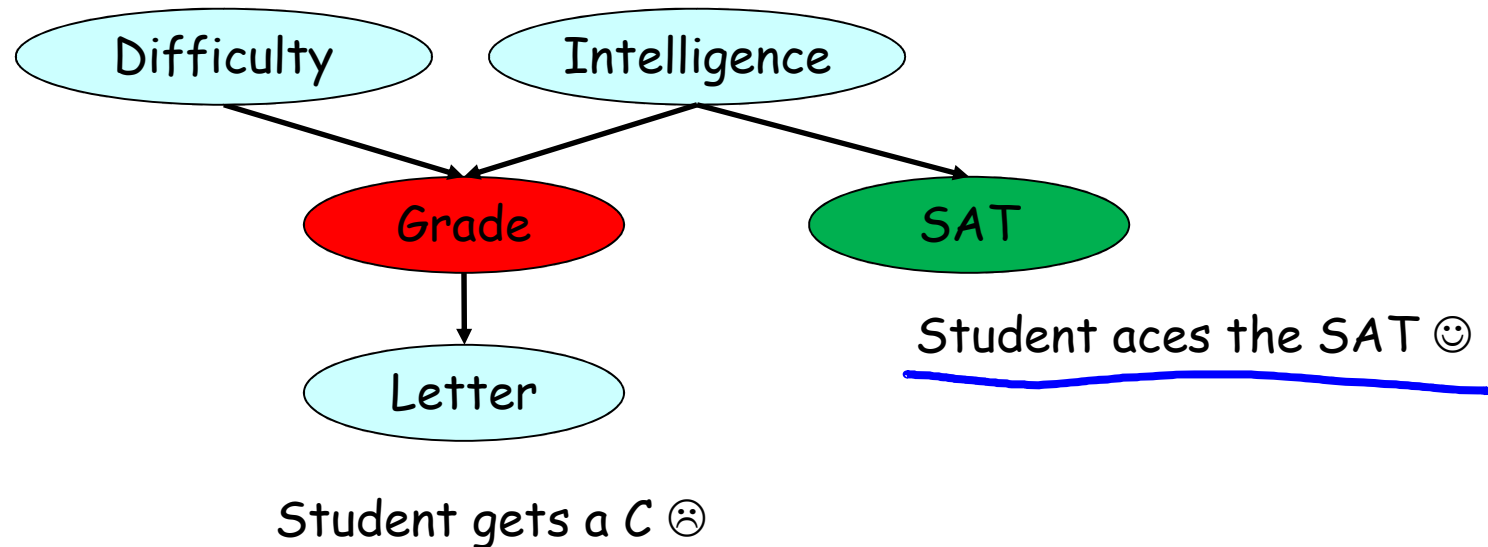
$P(X_1=1) = \frac{2}{3}$ $P(X_2=1) = \frac{2}{3}$
 condition X_1 $P(X_2=1|X_1=1) = 0.5$

Intercausal Reasoning II



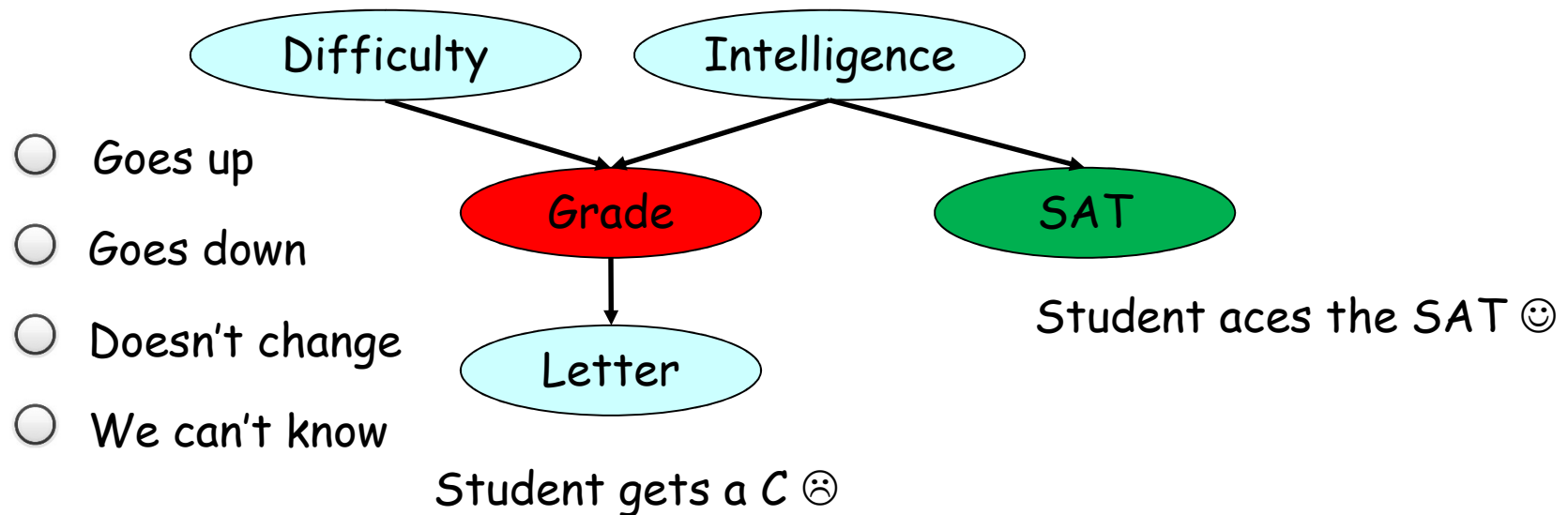
Student Aces the SAT

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



Student Aces the SAT

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



Student Aces the SAT

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

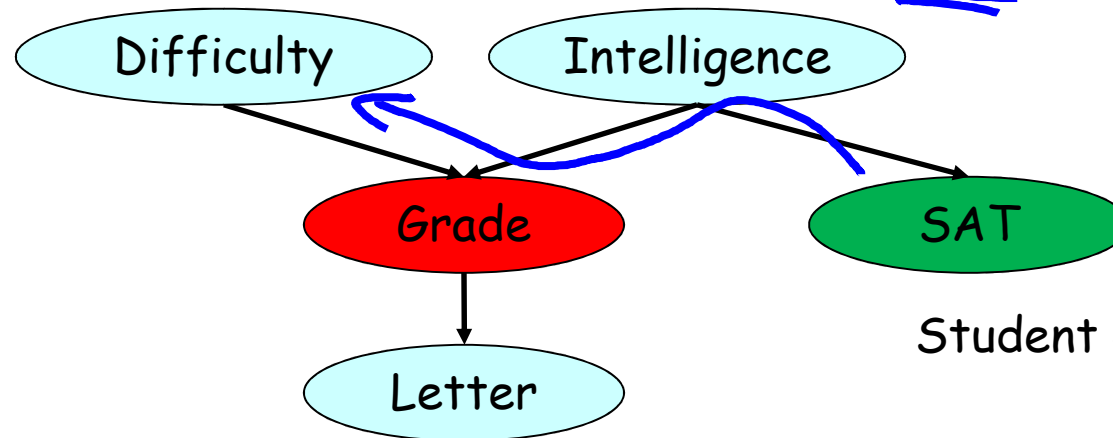
$$P(d^1 | g^3) \approx 0.63$$

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

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

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

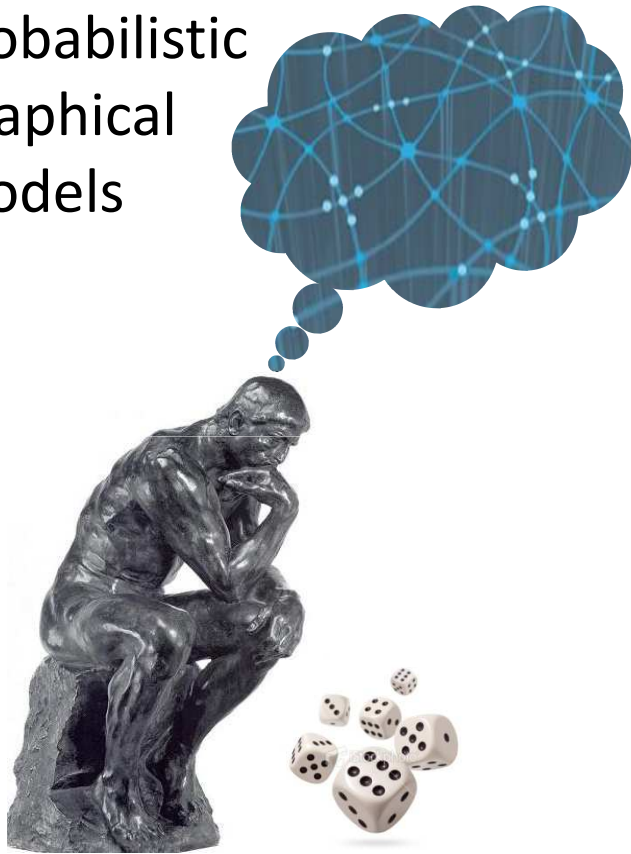
$$P(i^1 | g^3, s^1) \approx \underline{\underline{0.58}}$$



Student gets a C 😞

Student aces the SAT 😊

Probabilistic
Graphical
Models



Representation

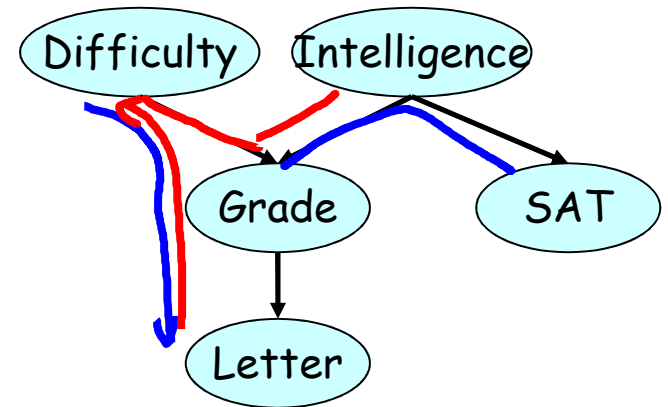
Bayesian Networks

Flow of
Probabilistic
Influence

When can X influence Y?

condition on v changes beliefs about 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

- $X \rightarrow Y$

- $X \leftarrow Y$

- $X \rightarrow W \rightarrow Y$

- $X \leftarrow W \leftarrow Y$

- $X \leftarrow W \rightarrow Y$

- $X \rightarrow W \leftarrow Y$

$W \notin Z$

✓

✓

✓

✗

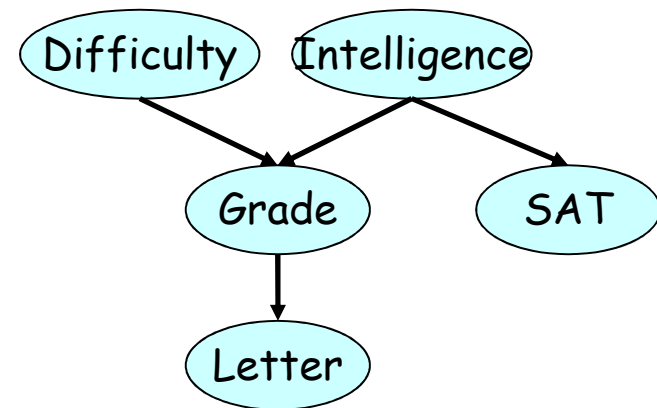
$W \in Z$

✗

✗

✗

✓



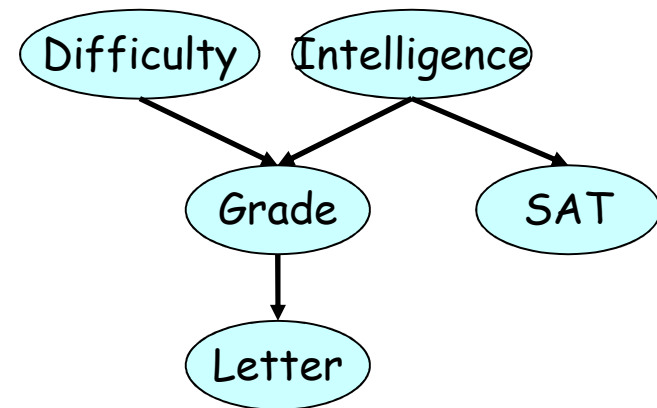
When can X influence Y given evidence about Z

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

I is observed X

I not observed,
nothing else X

I not observed
 G is observed



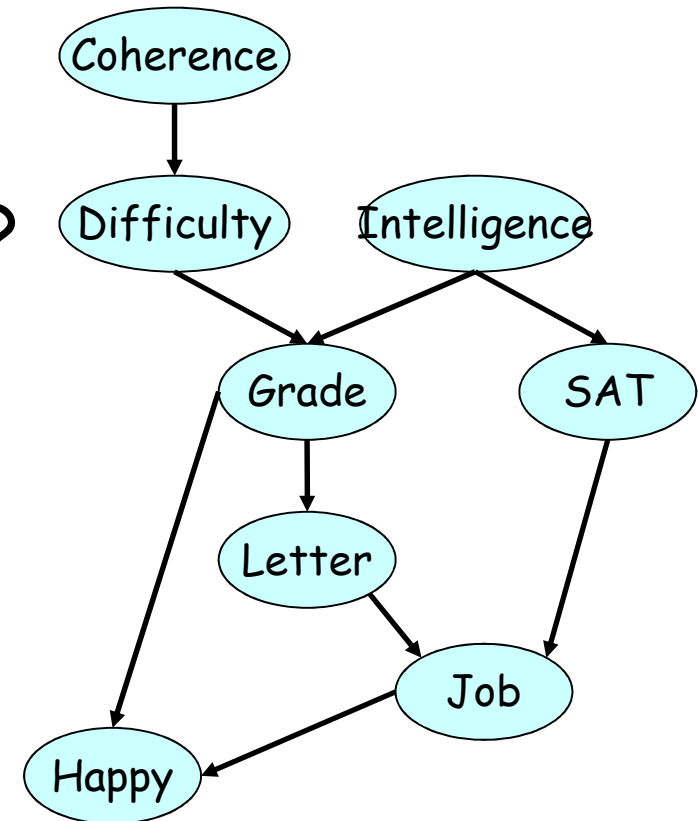
Active Trails

- 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
- activate v-structure*
- not in v-structure*

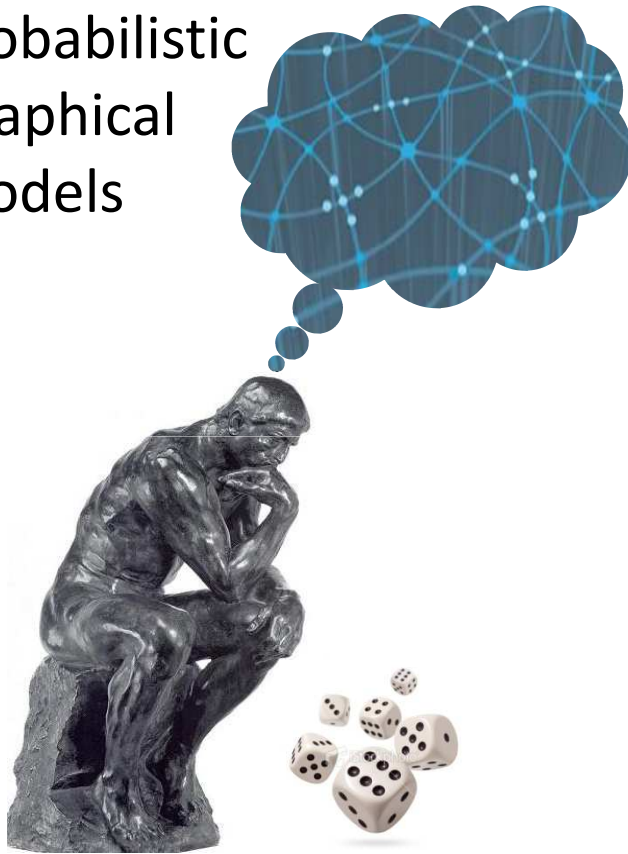
Active trails

- Which of the following are active trails if we observe G ? (Mark all that apply.)

- ☐ $C - D - G - I - S$
- ☐ $I - G - L - J - H$
- ☐ $I - S - J - H$
- ☐ $C - D - G - I - S - J - L$



Probabilistic
Graphical
Models



Independencies

Preliminaries

Independence

- For events α, β , $P \models \alpha \perp \beta$ if:

– $P(\alpha, \beta) = P(\alpha) \cdot P(\beta)$ *satisfies independence*

→ ~~– $P(\alpha|\beta) = P(\alpha)$~~

– $P(\beta|\alpha) = P(\beta)$

- For random variables X, Y , $P \models X \perp Y$ if:

→ ~~– $P(X, Y) = P(X) P(Y)$~~

– $P(X|Y) = P(X)$

– $P(Y|X) = P(Y)$

*universal
e.g. $P(x, y) = P(x) \cdot P(y)$*

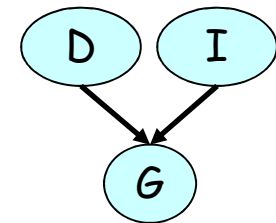
Independence

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

x



$P(I)$

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
satisfies indep conditioning

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

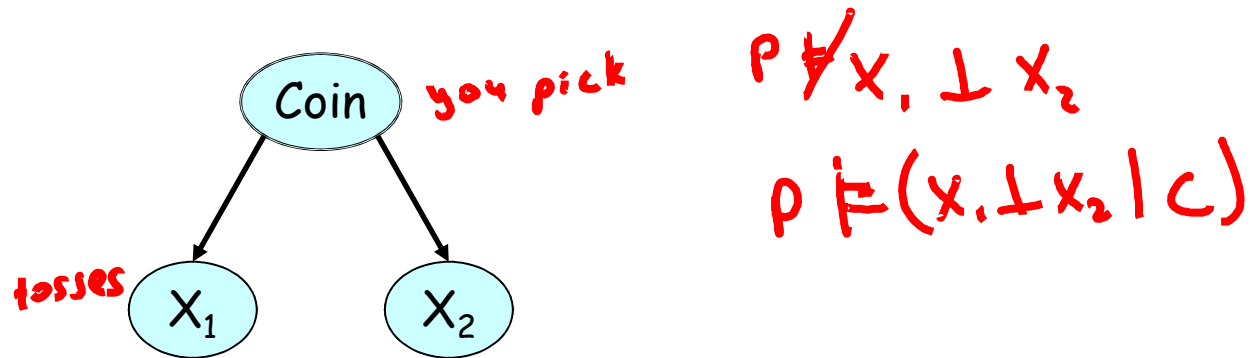
– $P(X, Y \mid Z) = P(X \mid Z) P(Y \mid Z)$
given

– $P(X \mid Y, Z) = P(X \mid Z)$

– $P(Y \mid X, Z) = P(Y \mid Z)$

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

Conditional Independence



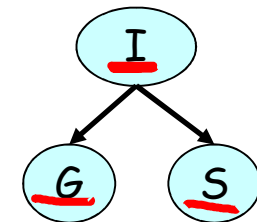
Conditional Independence

$P(I, S, G)$

I	S	G	Prob.
i^0	s^0	g^1	0.126
i^0	s^0	g^2	0.168
i^0	s^0	g^3	0.126
i^0	s^1	g^1	0.009
i^0	s^1	g^2	0.045
i^0	s^1	g^3	0.126
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	0.06
i^1	s^1	g^2	0.036
i^1	s^1	g^3	0.024

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

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



$P(S \mid \underline{i^0})$

S	Prob
s^0	0.95
s^1	0.05

$P(G \mid \underline{i^0})$

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

Conditioning can Lose Independences

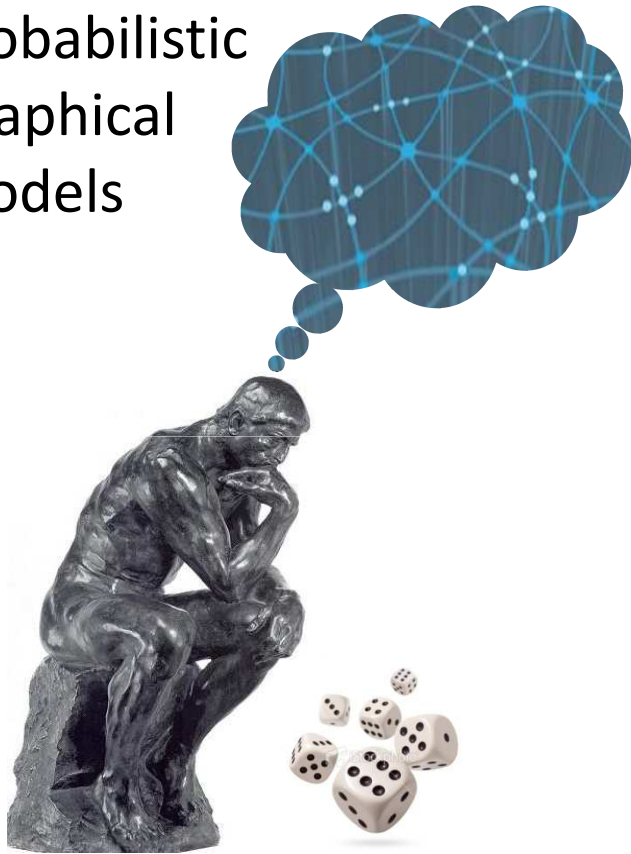


I	D	G	Prob.
i^0	d^0	g^1	0.126
i^0	d^0	g^2	0.168
i^0	d^0	g^3	0.126
i^0	d^1	g^1	0.009
i^0	d^1	g^2	0.045
i^0	d^1	g^3	0.126
i^1	d^0	g^1	0.252
i^1	d^0	g^2	0.0224
i^1	d^0	g^3	0.0056
i^1	d^1	g^1	0.06
i^1	d^1	g^2	0.036
i^1	d^1	g^3	0.024

$P(I, D \mid g^1)$

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

$$P(X,Y) = P(X) P(Y)$$

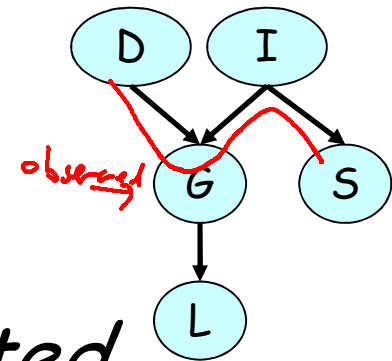
X,Y independent

$$P(\mathbf{X},\mathbf{Y},\mathbf{Z}) \propto \phi_1(\mathbf{X},\mathbf{Z}) \phi_2(\mathbf{Y},\mathbf{Z})$$

$(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z})$

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

d-separation

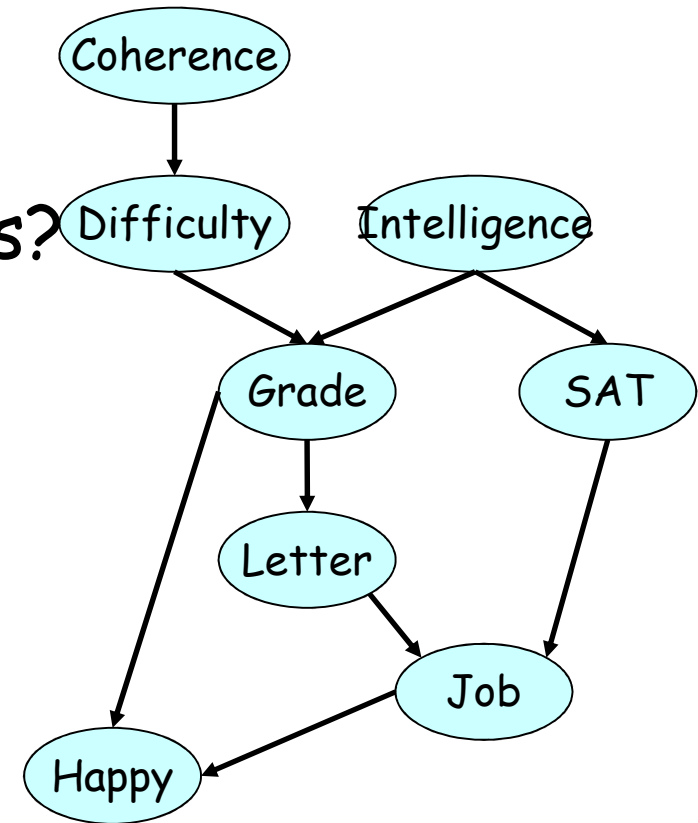
- Which of the following are true d-separation statements?
(Mark all that apply.)

d-sep(D,I | L)

d-sep(D,J | L)

d-sep(D,J | L,I)

d-sep(D,J | L,H,I)

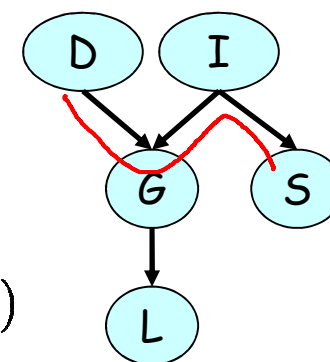


Factorization \Rightarrow Independence: BNs

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

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

$P \models D \perp S$



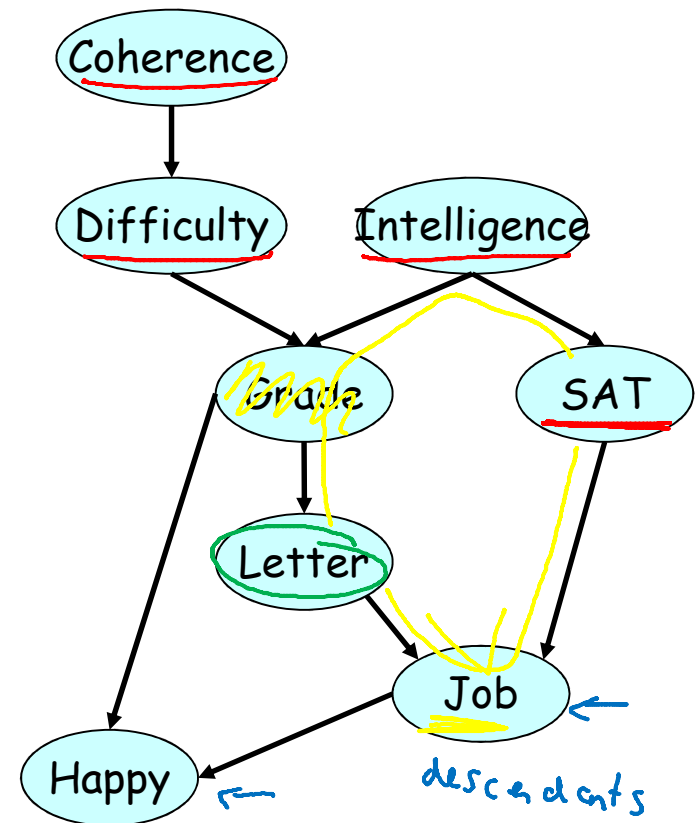
$$\begin{aligned}
 \underline{P(D, S)} &= \sum_{\underline{G, I, L}} P(D)P(I)P(G \mid D, I)P(S \mid I)P(L \mid G) \\
 &= \sum_I P(D)P(I)P(S \mid I) \sum_G \overbrace{P(G \mid D, I)}^1 \sum_L \overbrace{P(L \mid G)}^1 \\
 &= P(D) \underbrace{\left(\sum_I P(I)P(S \mid I) \right)}_{\phi_2(S)}
 \end{aligned}$$

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

Grade



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

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

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

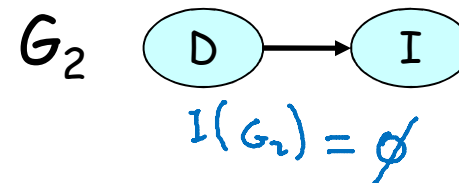
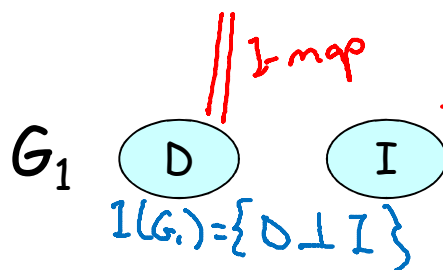
I-maps

$P_1 \neq D \perp I$

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

$P_2 \neq D \perp I$

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



Factorization \Rightarrow Independence: BNs

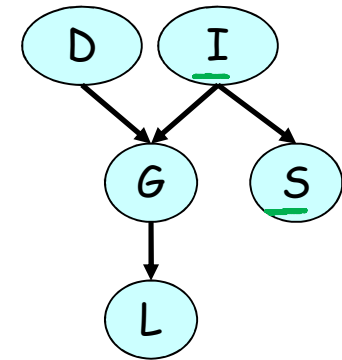
Theorem: If P factorizes over G , then G is
an I-map for P

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

110



$P(I, D)$ chain rule for probabilities

$$\underline{P(D, I, G, S, L)} = \underline{P(D)} \underline{P(I | D)} \underline{P(G | D, I)} \underline{P(S | D, I, G)} \underline{P(L | D, I, G, S)}$$

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

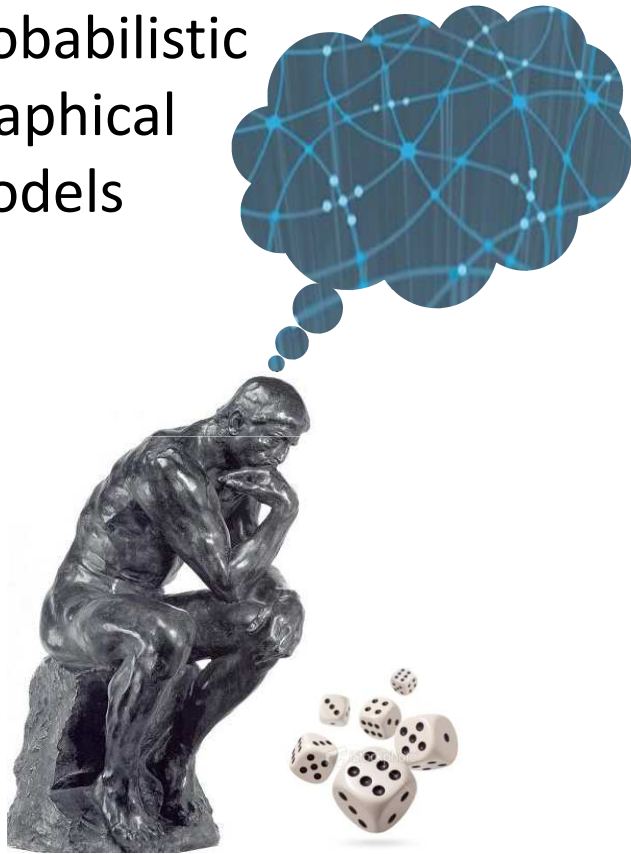
Summary

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

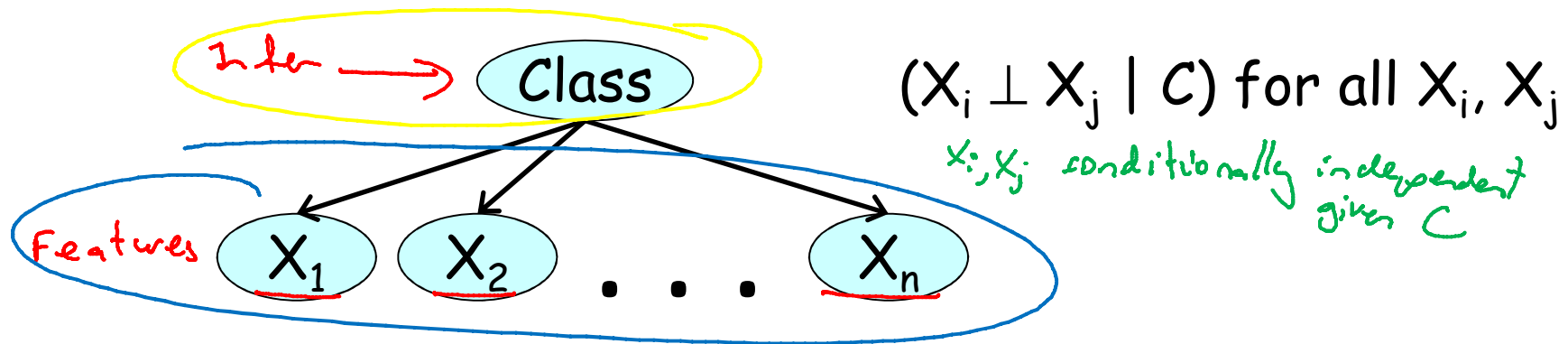


Representation

Bayesian Networks

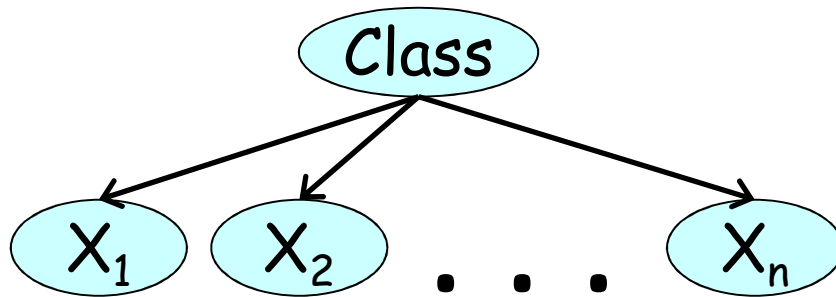
Naïve Bayes

Naïve Bayes Model



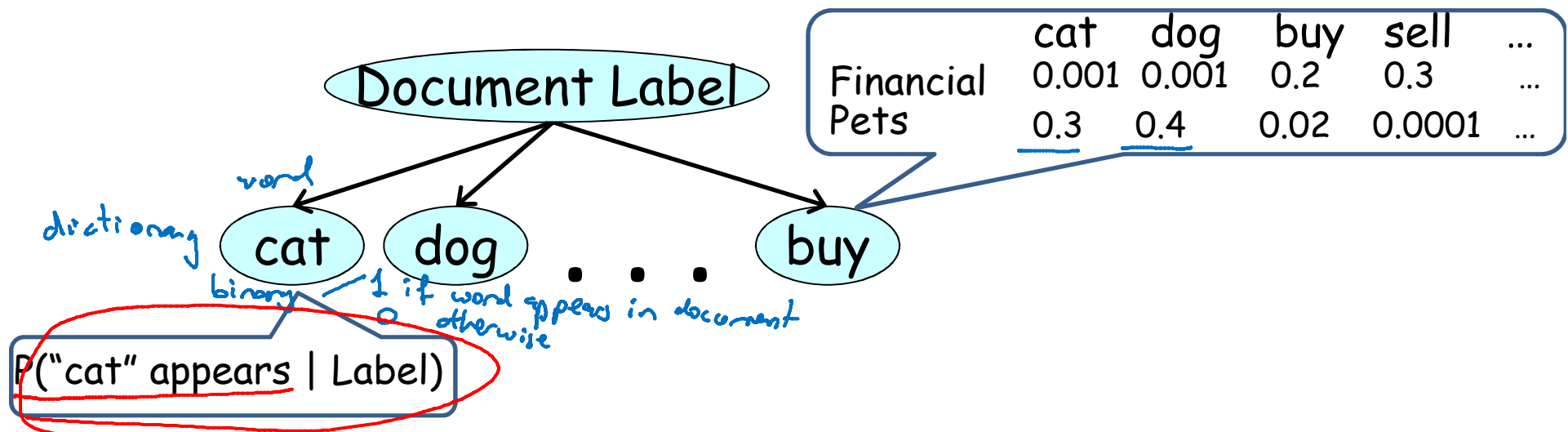
$$\underline{P(C, X_1, \dots, X_n)} = \underbrace{P(C) \prod_{i=1}^n P(X_i \mid C)}_{\text{Naïve Bayes Formula}}$$

Naïve Bayes Classifier



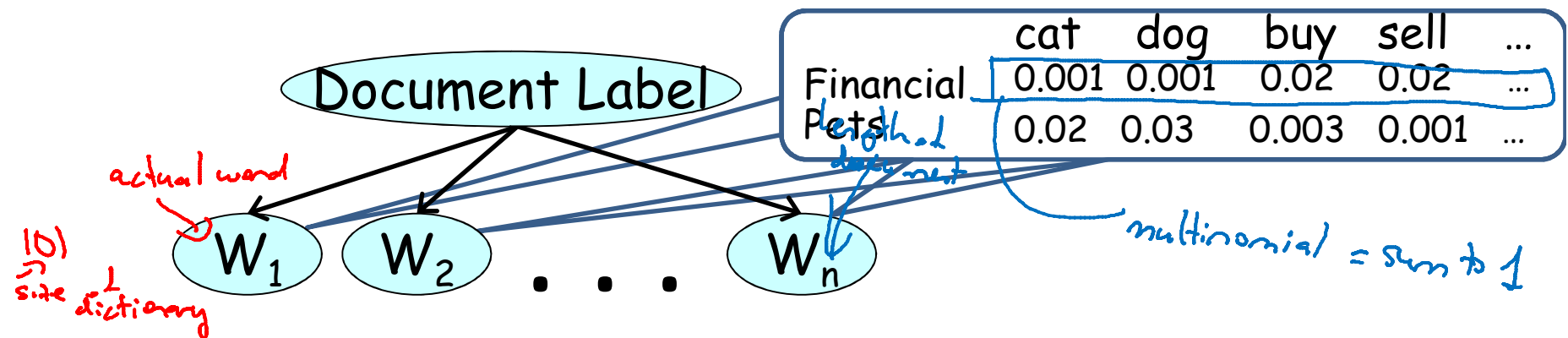
$$\frac{P(C = c^1 \mid x_1, \dots, x_n)}{P(C = c^2 \mid x_1, \dots, x_n)} = \underbrace{\frac{P(C = c^1)}{P(C = c^2)}}_{\text{odds ratios}} \prod_{i=1}^n \frac{P(\underline{x_i} \mid C = c^1)}{P(\underline{x_i} \mid C = c^2)}$$

Bernoulli Naïve Bayes for Text



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

Multinomial Naïve Bayes for Text

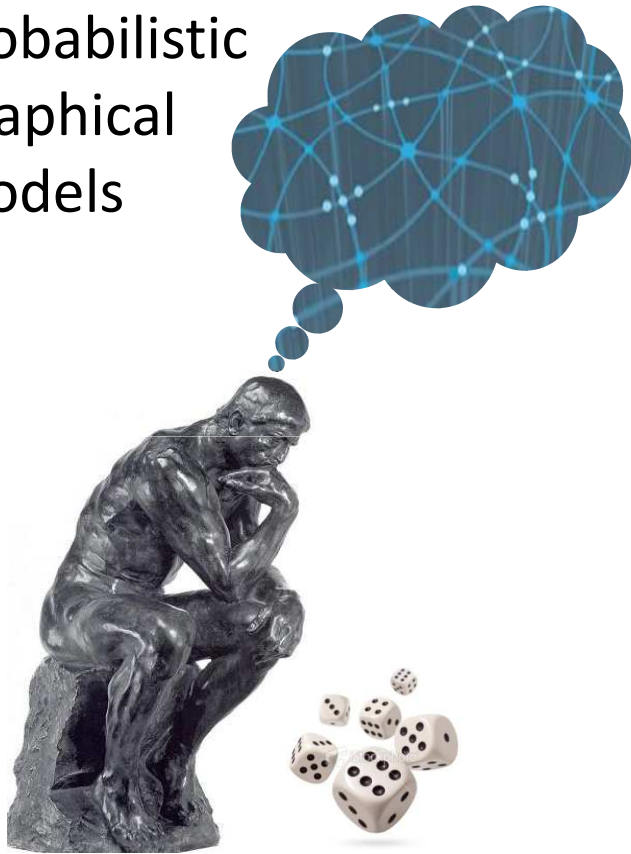


$$\frac{P(C = c^1 \mid x_1, \dots, x_n)}{P(C = c^2 \mid x_1, \dots, x_n)} = \frac{P(C = c^1)}{P(C = c^2)} \prod_{i=1}^n \frac{P(x_i \mid C = c^1)}{P(x_i \mid C = c^2)}$$

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.

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.

Medical Diagnosis: Pathfinder (1992)

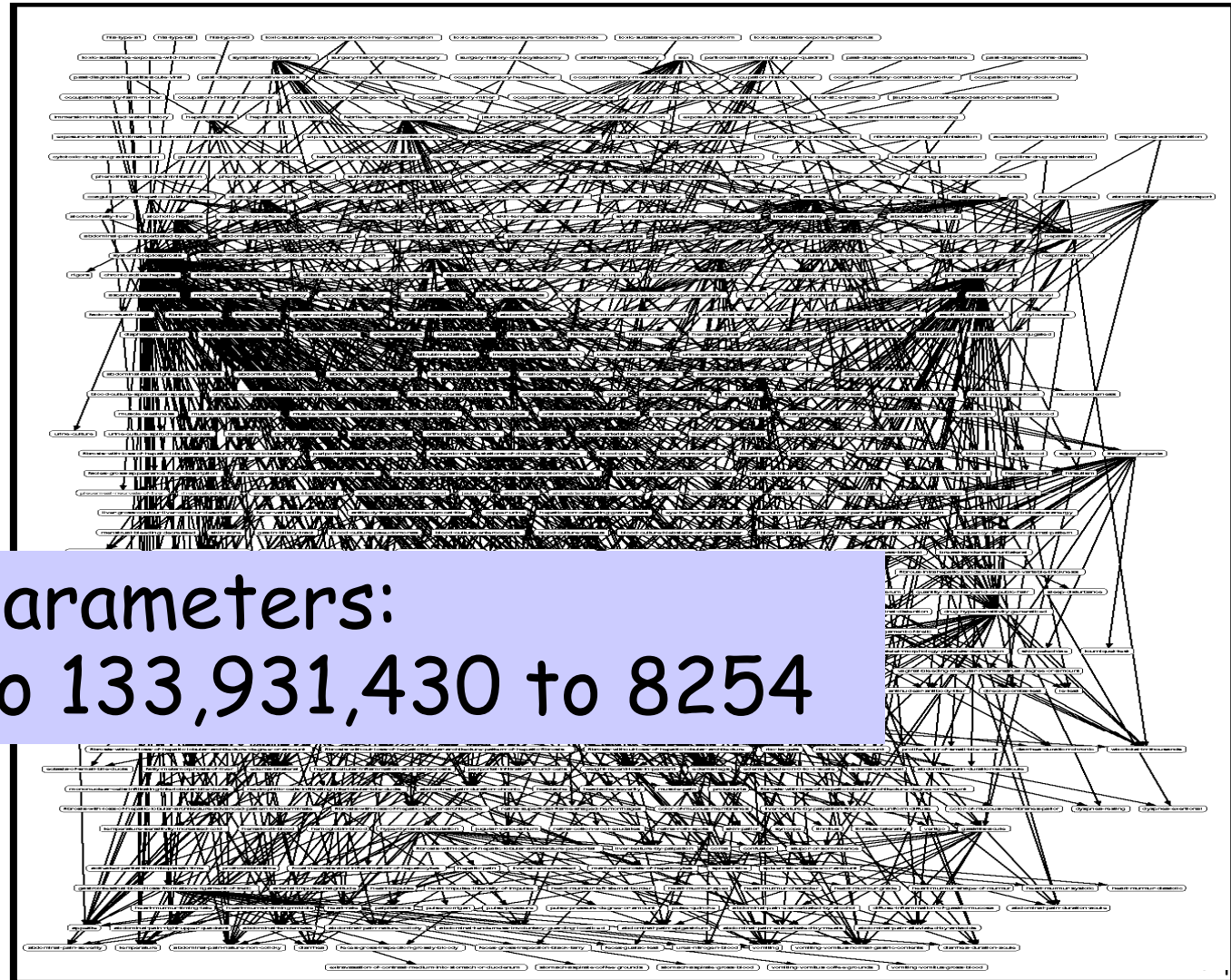
- Pathfinder IV: Full Bayesian network
 - Removed incorrect independencies
 - Additional parents led to more accurate estimation of probabilities
- BN model agreed with expert panel in 50/53 cases, vs 47/53 for naïve Bayes model
- Accuracy as high as expert that designed the model

Heckerman et al.

CPCS

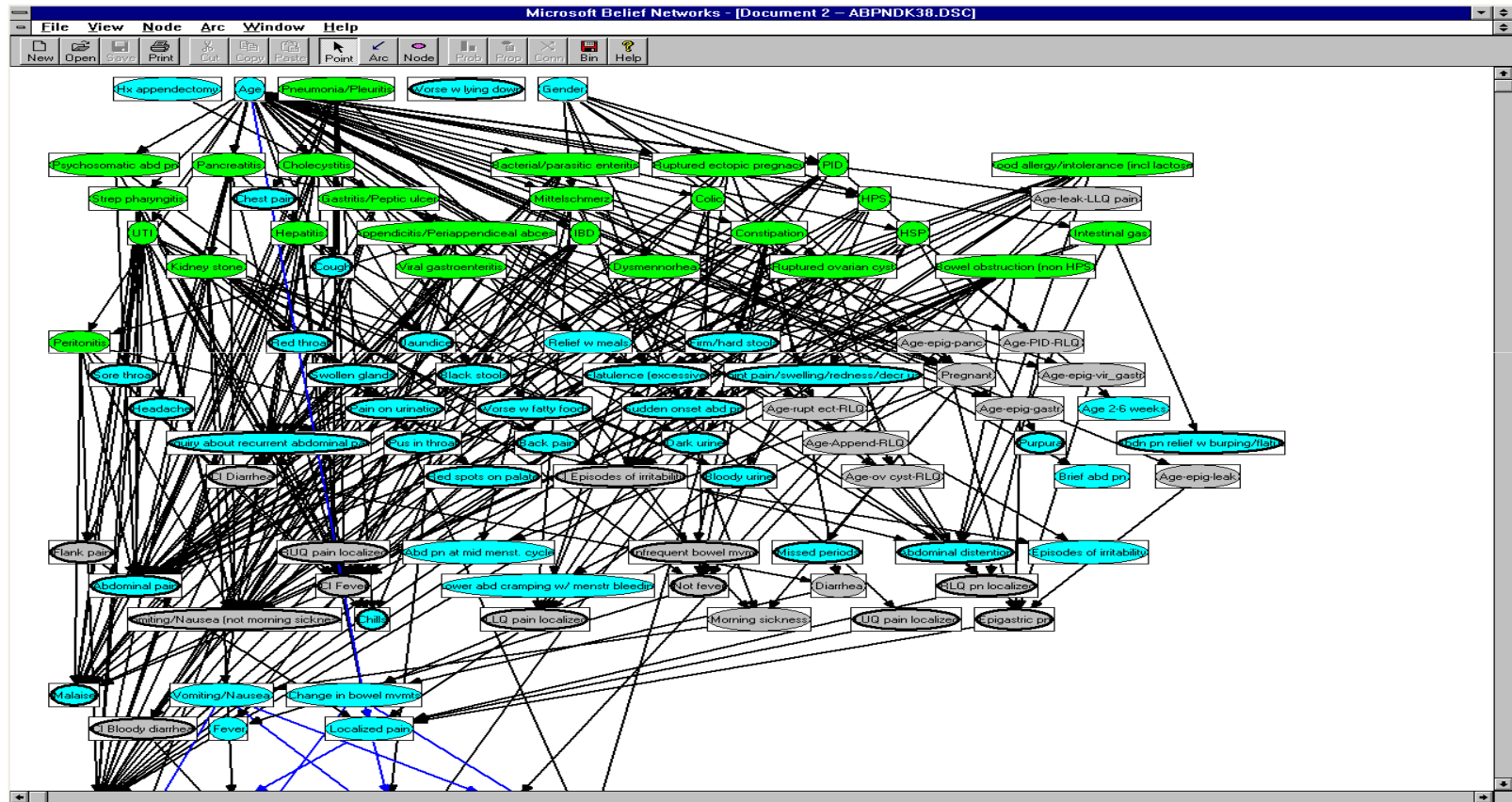
of parameters:
 2^{1000} to 133,931,430 to 8254

Pradhan et al.



John Koller

Medical Diagnosis (Microsoft)



Medical Diagnosis (Microsoft)

Applet started

ON STAGE ESSENTIALS COMMUNICATE FIND

OnParenting
May 14 - May 20, 1997

Fidelity Investments
Fidelity Distributors Corporation

Our home on the web [is where] click here

cover contents news experts fun handbook talk **find** help feedback

There are two ways to search for specific information in **OnParenting**. In **Find by Word**, type the word(s) you want to find and get a list of titles relevant to that word. **Find by Symptom** will help you get information about children's symptoms. [Help](#) has tips to target your search.

Find by Word

Find by Symptom ▶

Describe the child
in the drop-down boxes at the right. Relevant information will appear below.

Age: Sex:

Complaint:

Localized pain: Can the child localize, or point to, the site of the pain?

- ☐ No, unable to localize
- ☐ Below the navel to the child's left
- ☐ Above the child's navel
- ☐ Either of the child's sides
- ☐ Below the navel to the child's right
- ☐ Above the navel to the child's right
- ☐ Above the navel to the child's left
- ☐ Don't Know

Results so far

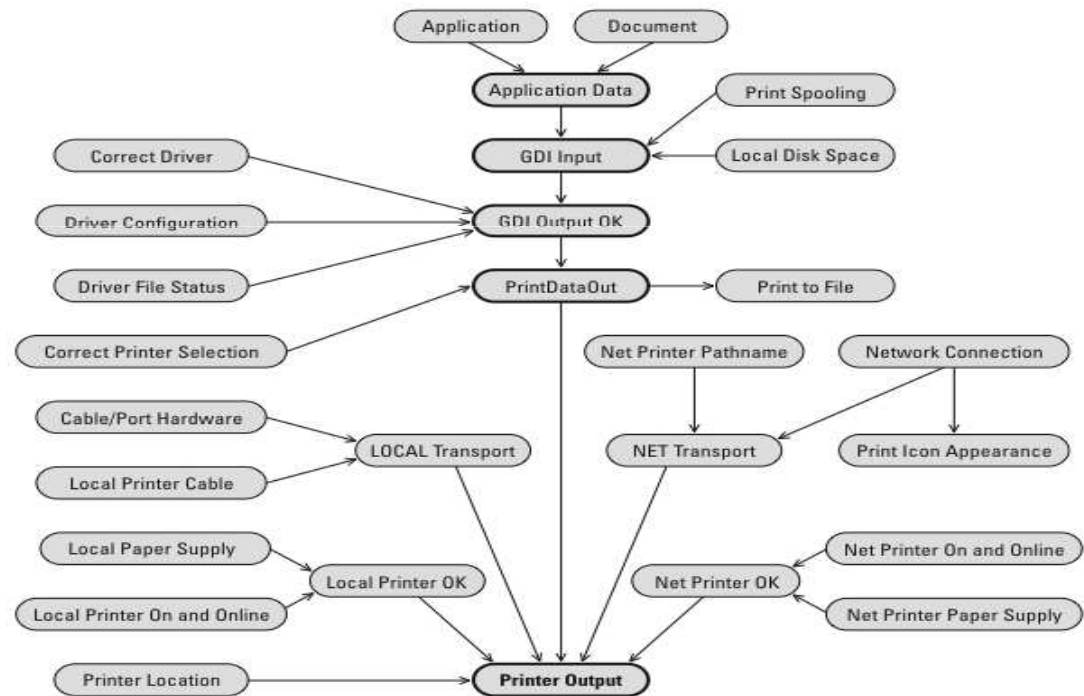
Disorder	Relevance
Viral gastroenteritis	<div><div></div></div>
Psychosomatic pain	<div><div></div></div>
Urinary tract infection	<div><div></div></div>
Other	<div><div></div></div>

Start Over **Review**

Next>> **Finish**

Fault Diagnosis

- Microsoft troubleshooters



Fault Diagnosis

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