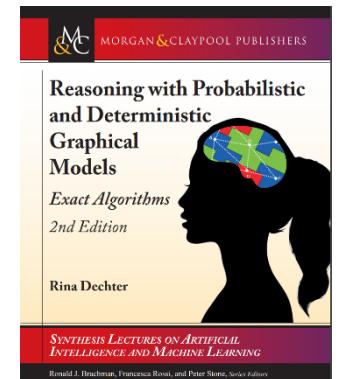


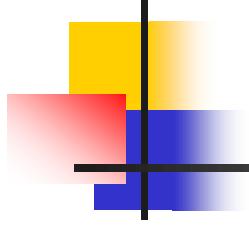
CS 295: Causal Reasoning

Rina Dechter

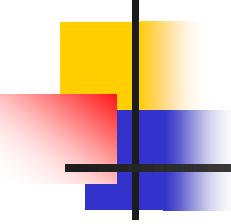
Exact Inference Algorithms Bucket-elimination

Chapter 4 Dechter's book





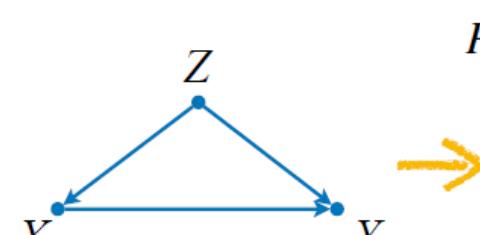
Factorizing Observational Distributions



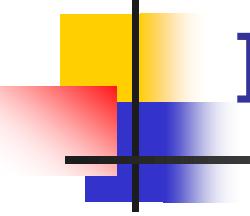
Markovian Case

- The distribution $P(\mathbf{v})$ decomposes as:

$$P(\mathbf{v}) = \sum_{\mathbf{u}} P(\mathbf{u}) \prod_{V_i \in \mathbf{V}} P(v_i | v_1, \dots, v_{i-1}, \mathbf{u}) = \sum_{\mathbf{u}} P(\mathbf{u}) \prod_{V_i \in \mathbf{V}} P(v_i | pa_i, u_i)$$


$$\begin{aligned} P(z, x, y) &= \sum_{\mathbf{u}} P(\mathbf{u}) P(z | u_z) P(x | z, u_x) P(y | x, z, u_y) \\ &= \left(\sum_{u_z} P(z | u_z) P(u_z) \right) \left(\sum_{u_x} P(x | z, u_x) P(u_x) \right) \left(\sum_{u_y} P(y | x, z, u_y) P(u_y) \right) \\ &= P(z) P(x | z) P(y | x, z) \end{aligned}$$

- In Markovian models, $P(v_i | pa_i)$ can be seen as “canonical factors”.

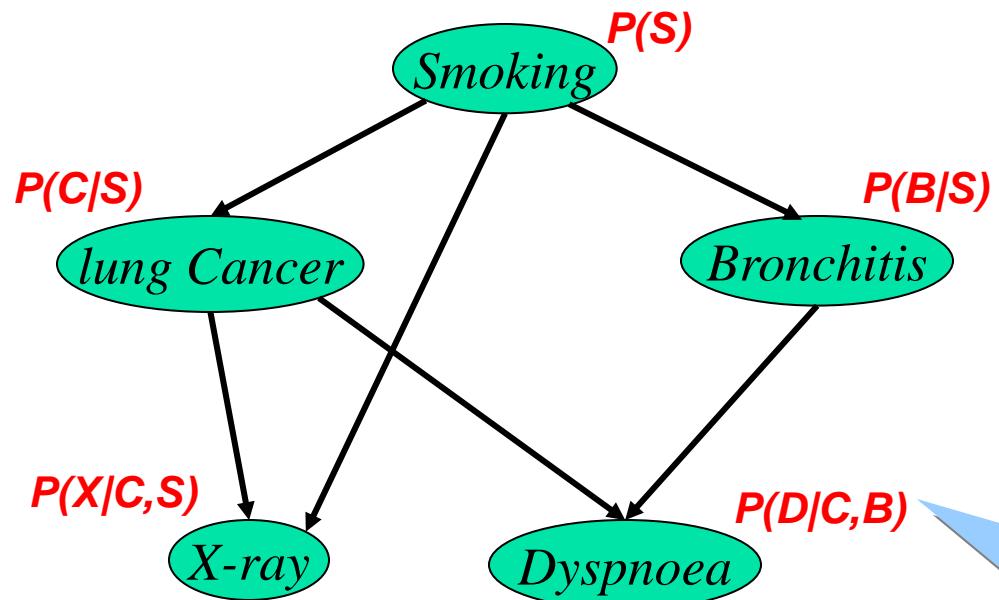


Inference for probabilistic networks

- Bucket elimination
 - Belief-updating, $P(e)$, partition function
 - Marginals, probability of evidence
 - The impact of evidence
 - for MPE (\rightarrow MAP)
 - for MAP (\rightarrow Marginal Map)
- Induced-Width

Bayesian Networks: Example

(Pearl, 1988)



$$\text{BN} = (\mathbf{G}, \Theta)$$

CPD:

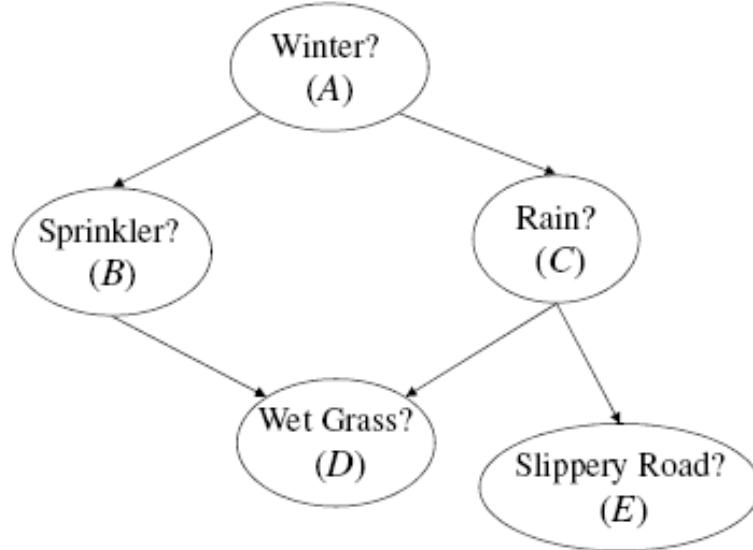
C	B	$P(D C,B)$	
0	0	0.1	0.9
0	1	0.7	0.3
1	0	0.8	0.2
1	1	0.9	0.1

$$P(S, C, B, X, D) = P(S) P(C|S) P(B|S) P(X|C,S) P(D|C,B)$$

Belief Updating:

$P(\text{lung cancer=yes} \mid \text{smoking=no, dyspnoea=yes}) = ?$

A Bayesian Network



A	C	$\Theta_{C A}$
true	true	.8
true	false	.2
false	true	.1
false	false	.9

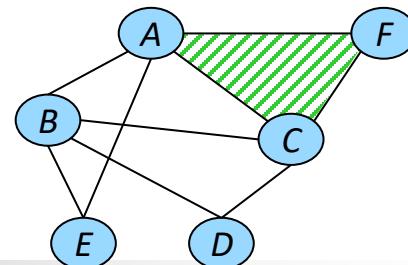
B	C	D	$\Theta_{D BC}$
true	true	true	.95
true	true	false	.05
true	false	true	.9
true	false	false	.1
false	true	true	.8
false	true	false	.2
false	false	true	0
false	false	false	1

A	Θ_A
true	.6
false	.4

A	B	$\Theta_{B A}$
true	true	.2
true	false	.8
false	true	.75
false	false	.25

C	E	$\Theta_{E C}$
true	true	.7
true	false	.3
false	true	0
false	false	1

Types of queries

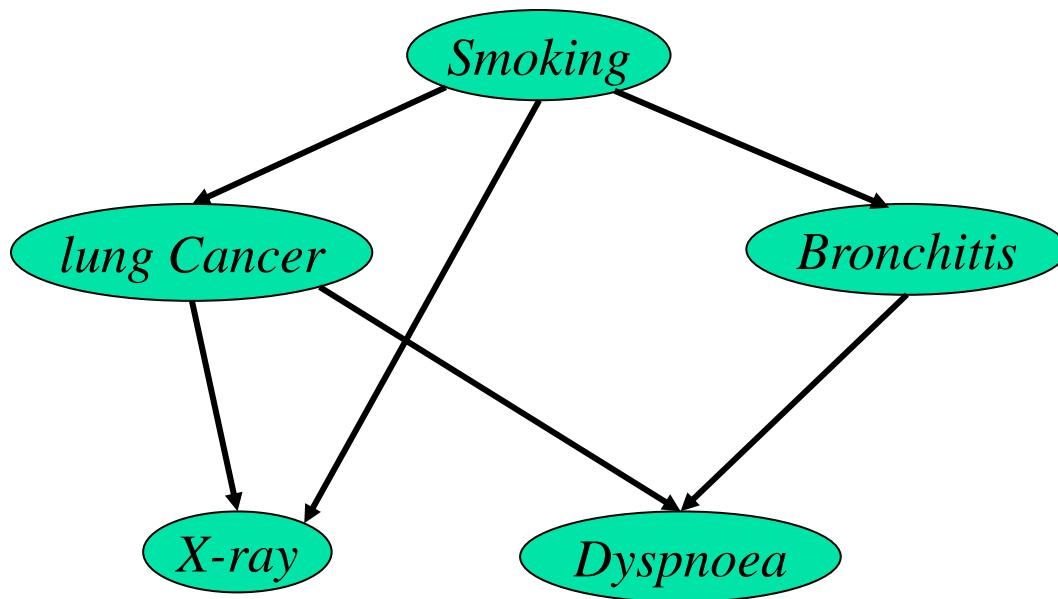


▶ Max-Inference	$f(\mathbf{x}^*) = \max_{\mathbf{x}} \prod_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha})$
▶ Sum-Inference	$Z = \sum_{\mathbf{x}} \prod_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha})$
▶ Mixed-Inference	$f(\mathbf{x}_M^*) = \max_{\mathbf{x}_M} \sum_{\mathbf{x}_S} \prod_{\alpha} f_{\alpha}(\mathbf{x}_{\alpha})$

Harder ↓

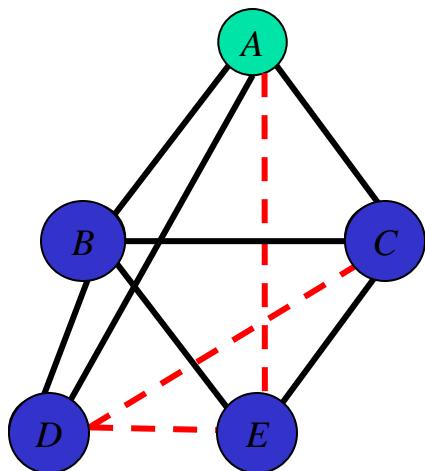
- **NP-hard**: exponentially many terms
- We will focus on exact and then on **approximation** algorithms
 - **Anytime**: very fast & very approximate ! Slower & more accurate

Belief Updating



$$P(\text{lung cancer=yes} \mid \text{smoking=no}, \text{dyspnoea=yes}) = ?$$

Belief updating: $P(X|\text{evidence})=?$



“Moral” graph

$$P(a|e=0) \propto P(a, e=0) =$$

$$\sum_{e=0,d,c,b} P(a) \underbrace{P(b|a)P(c|a)}_{b} \underbrace{P(d|b,a)P(e|b,c)}_{c} =$$

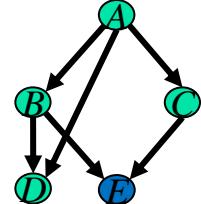
$$P(a) \sum_{e=0} \sum_d \sum_c P(c|a) \sum_b P(b|a) P(d|b,a) P(e|b,c)$$

Variable Elimination

$h^B(a, d, c, e)$

Bucket elimination

Algorithm *BE-bel* (Dechter 1996)



$$P(A | E = 0) = \alpha \sum_{E=0,D,C,B} P(A) \cdot P(B | A) \cdot P(C | A) \cdot P(D | A, B) \cdot P(E | B, C)$$

$\sum \prod_b$ ← *Elimination operator*

bucket B:

$$P(b|a) \quad P(d|b,a) \quad P(e|b,c)$$

bucket C:

$$P(c|a) \quad \lambda^B(a, d, c, e)$$

bucket D:

$$\lambda^C(a, d, e)$$

bucket E:

$$e=0 \quad \lambda^D(a, e)$$

bucket A:

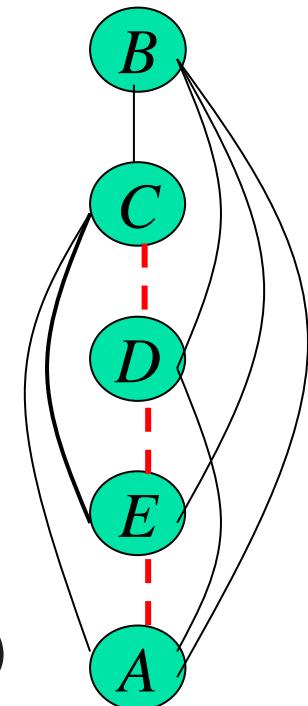
$$P(a) \quad \lambda^E(a)$$

$$P(a, e=0)$$

$$P(a | e=0) = \frac{P(a, e=0)}{P(e=0)}$$

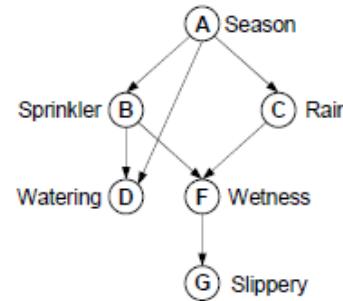
W = 4*

"induced width"
(max clique size)

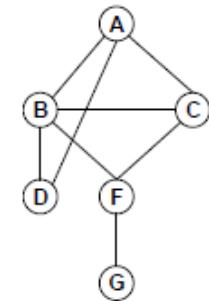


A Bayesian Network

Ordering: A,C,B,E,D,G



(a) Directed acyclic graph



(b) Moral graph

$$P(a, g=1) = \sum_{c,b,e,d,g=1} P(a, b, c, d, e, g) = \sum_{c,b,f,d,g=1} P(g|f)P(f|b,c)P(d|a,b)P(c|a)P(b|a)P(a).$$

$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \sum_f P(f|b,c) \sum_d P(d|b,a) \sum_{g=1} P(g|f). \quad (4.1)$$

$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \sum_f P(f|b,c) \lambda_G(f) \sum_d P(d|b,a). \quad (4.2)$$

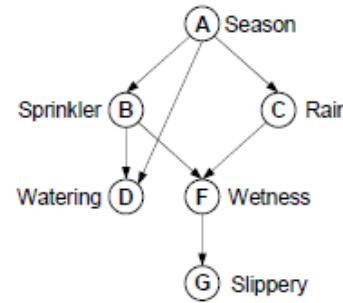
$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \lambda_D(a,b) \sum_f P(f|b,c) \lambda_G(f) \quad (4.3)$$

$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \lambda_D(a,b) \lambda_F(b,c) \quad (4.4)$$

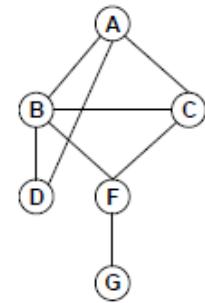
$$P(a, g=1) = P(a) \sum_c P(c|a) \lambda_B(a,c) \quad (4.5)$$

A Bayesian Network

Ordering: A,C,B,E,D,G



(a) Directed acyclic graph



(b) Moral graph

$$P(a, g=1) = \sum_{c,b,e,d,g=1} P(a, b, c, d, e, g) = \sum_{c,b,f,d,g=1} P(g|f)P(f|b,c)P(d|a,b)P(c|a)P(b|a)P(a).$$

$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \sum_f P(f|b,c) \sum_d P(d|b,a) \sum_{g=1} P(g|f). \quad (4.1)$$

$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \sum_f P(f|b,c) \lambda_G(f) \sum_d P(d|b,a). \quad (4.2)$$

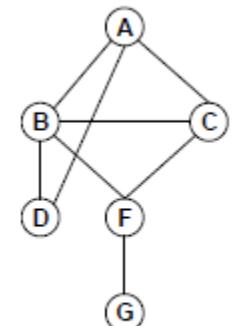
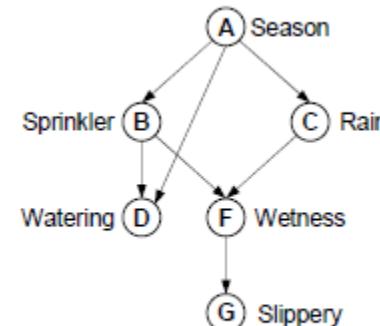
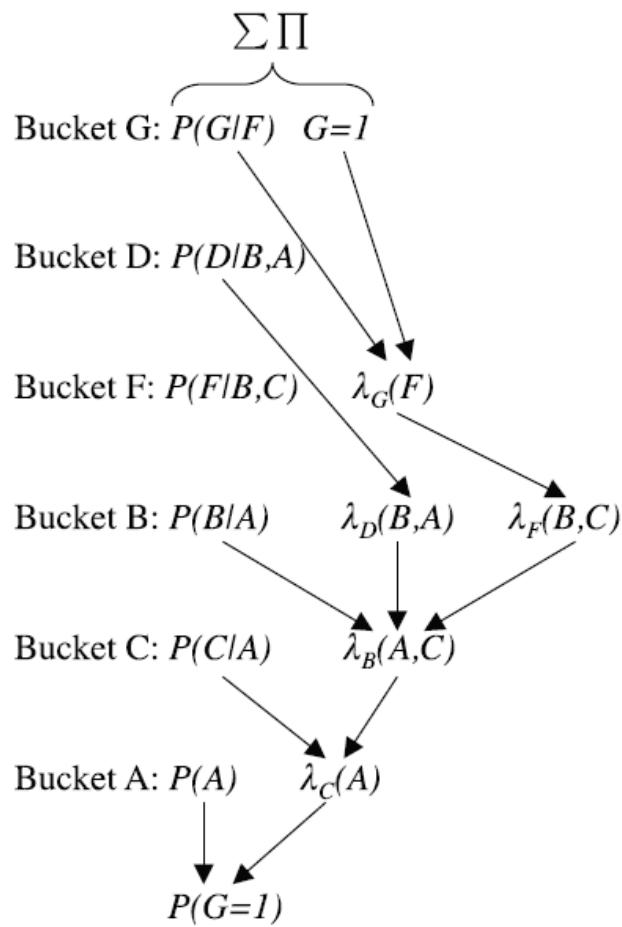
$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \lambda_D(a,b) \sum_f P(f|b,c) \lambda_G(f) \quad (4.3)$$

$$P(a, g=1) = P(a) \sum_c P(c|a) \sum_b P(b|a) \lambda_D(a,b) \lambda_F(b,c) \quad (4.4)$$

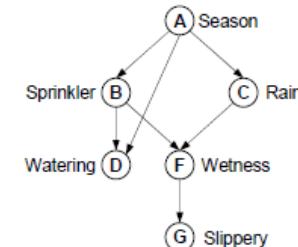
$$P(a, g=1) = P(a) \sum_c P(c|a) \lambda_B(a,c) \quad (4.5)$$

A Bayesian Network

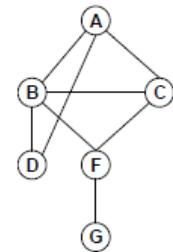
Ordering: A,C,B,F,D,G



A Different Ordering



(a) Directed acyclic graph



(b) Moral graph

Ordering: A,F,D,C,B,G

$$\begin{aligned}
 P(a, g=1) &= P(a) \sum_f \sum_d \sum_c P(c|a) \sum_b P(b|a) P(d|a,b) P(f|b,c) \sum_{g=1} P(g|f) \\
 &= P(a) \sum_f \lambda_G(f) \sum_d \sum_c P(c|a) \sum_b P(b|a) P(d|a,b) P(f|b,c) \\
 &= P(a) \sum_f \lambda_G(f) \sum_d \sum_c P(c|a) \lambda_B(a,d,c,f) \\
 &= P(a) \sum_f \lambda_g(f) \sum_d \lambda_C(a,d,f) \\
 &= P(a) \sum_f \lambda_G(f) \lambda_D(a,f) \\
 &= P(a) \lambda_F(a)
 \end{aligned}$$

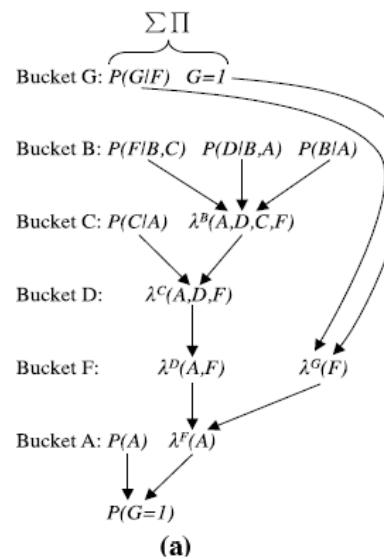
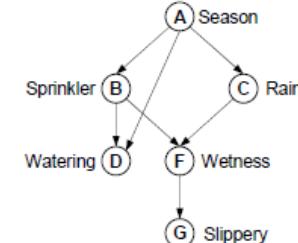
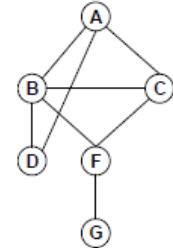


Figure 4.3: The bucket's output when processing along $d_2 = A, F, D, C, B, G$
CS295, Spring 2021

A Different Ordering



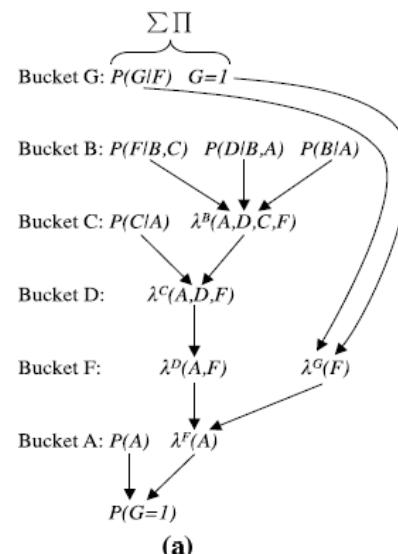
(a) Directed acyclic graph



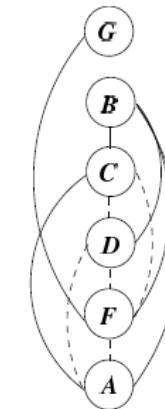
(b) Moral graph

Ordering: A,F,D,C,B,G

$$\begin{aligned}
 P(a, g=1) &= P(a) \sum_f \sum_d \sum_c P(c|a) \sum_b P(b|a) P(d|a,b) P(f|b,c) \sum_{g=1} P(g|f) \\
 &= P(a) \sum_f \lambda_G(f) \sum_d \sum_c P(c|a) \sum_b P(b|a) P(d|a,b) P(f|b,c) \\
 &= P(a) \sum_f \lambda_G(f) \sum_d \sum_c P(c|a) \lambda_B(a,d,c,f) \\
 &= P(a) \sum_f \lambda_g(f) \sum_d \lambda_C(a,d,f) \\
 &= P(a) \sum_f \lambda_G(f) \lambda_D(a,f) \\
 &= P(a) \lambda_F(a)
 \end{aligned}$$



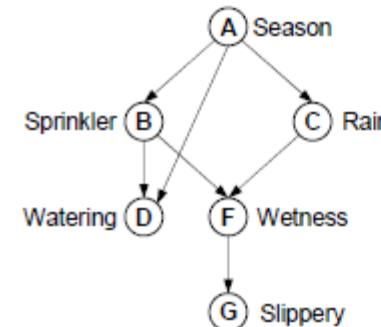
(a)



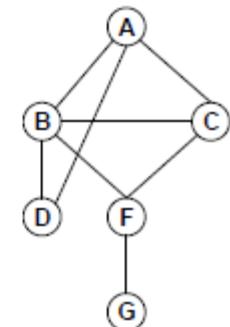
(b)

Figure 4.3: The bucket's output when processing along $d_2 = A, F, D, C, B, G$
CS295, Spring 2021

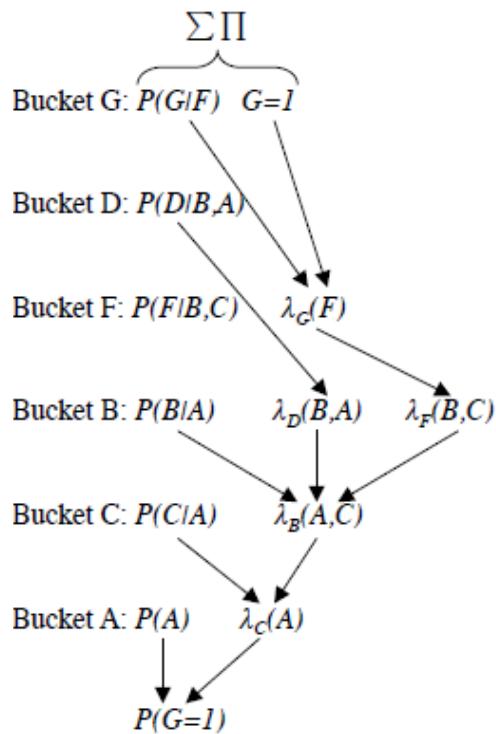
A Bayesian Network Processed Along 2 Orderings



(a) Directed acyclic graph



(b) Moral graph



$d_1 = A, C, B, F, D, G$

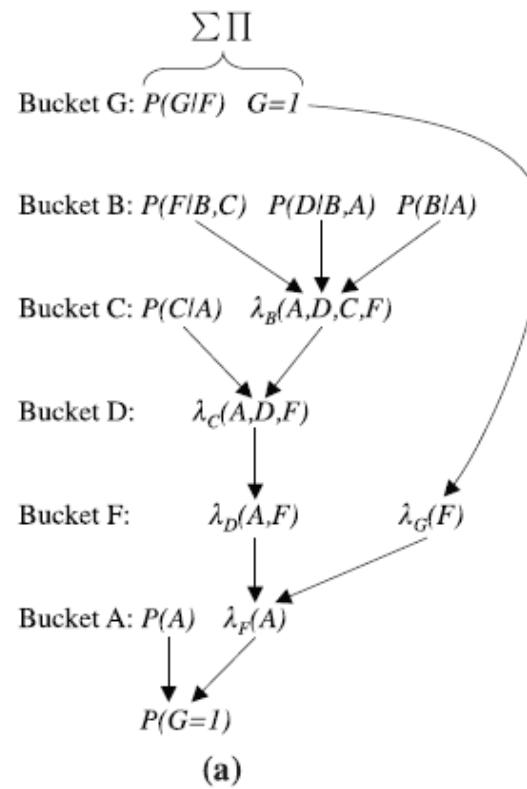
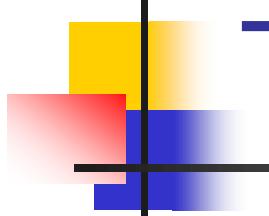


Figure 4.4: The bucket's output when processing along $d_2 = A, F, D, C, B, G$.
CS295, Spring 2021



The Operation In a Bucket

- Multiplying functions
- Marginalizing (summing-out) functions

Combination of Cost Functions

A	B	f(A,B)
b	b	0.4
b	g	0.1
g	b	0
g	g	0.5



A	B	C	f(A,B,C)
b	b	b	0.1
b	b	g	0
b	g	b	0
b	g	g	0.08
g	b	b	0
g	b	g	0
g	g	b	0
g	g	g	0.4

B	C	f(B,C)
b	b	0.2
b	g	0
g	b	0
g	g	0.8

$$= 0.1 \times 0.8$$

Factors: Sum-Out Operation

The result of **summing out** variable X from factor $f(\mathbf{X})$

is another factor over variables $\mathbf{Y} = \mathbf{X} \setminus \{X\}$:

$$\left(\sum_X f \right) (\mathbf{y}) \stackrel{\text{def}}{=} \sum_X f(x, \mathbf{y})$$

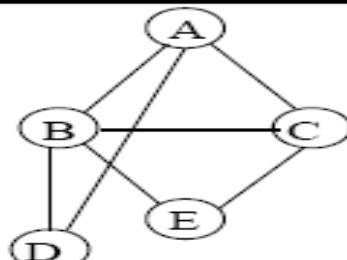
B	C	D	f_1
true	true	true	.95
true	true	false	.05
true	false	true	.9
true	false	false	.1
false	true	true	.8
false	true	false	.2
false	false	true	0
false	false	false	1

B	C	$\sum_D f_1$
true	true	1
true	false	1
false	true	1
false	false	1

$$\sum_B \sum_C \sum_D f_1$$

T 4

Bucket Elimination and Induced Width



Ordering: **a, e, d, c, b**

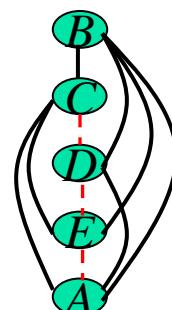
$$\text{bucket}(B) = P(e|b, c), P(d|a, b), P(b|a)$$

$$\text{bucket}(C) = P(c|a) \parallel \lambda_B(a, c, d, e)$$

$$\text{bucket}(D) = \parallel \lambda_C(a, d, e)$$

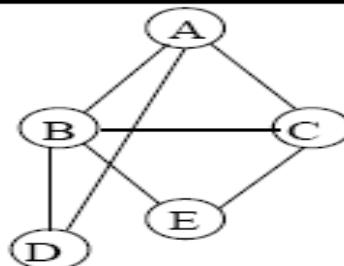
$$\text{bucket}(E) = e = 0 \parallel \lambda_D(a, c)$$

$$\text{bucket}(A) = P(a) \parallel \lambda_E(a)$$



24

Bucket Elimination and Induced Width



Ordering: a, b, c, d, e

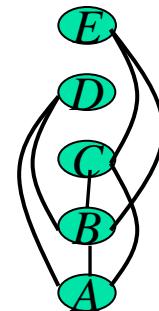
$$\text{bucket}(E) = P(e|b, c), \quad e = 0$$

$$\text{bucket}(D) = P(d|a, b)$$

$$\text{bucket}(C) = P(c|a) \parallel P(e = 0|b, c)$$

$$\text{bucket}(B) = P(b|a) \parallel \lambda_D(a, b), \lambda_C(b, c)$$

$$\text{bucket}(A) = P(a) \parallel \lambda_B(a)$$



Ordering: a, e, d, c, b

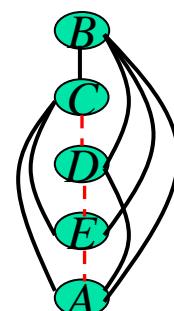
$$\text{bucket}(B) = P(e|b, c), P(d|a, b), P(b|a)$$

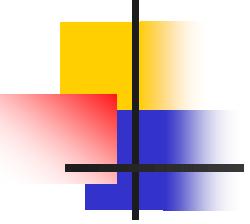
$$\text{bucket}(C) = P(c|a) \parallel \lambda_B(a, c, d, e)$$

$$\text{bucket}(D) = \parallel \lambda_C(a, d, e)$$

$$\text{bucket}(E) = e = 0 \parallel \lambda_D(a, c)$$

$$\text{bucket}(A) = P(a) \parallel \lambda_E(a)$$





ALGORITHM BE-BEL

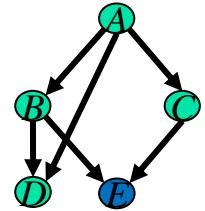
Input: A belief network $\mathcal{B} = \langle X, D, P_G, \prod \rangle$, an ordering $d = (X_1, \dots, X_n)$; evidence e
output: The belief $P(X_1|e)$ and probability of evidence $P(e)$

1. Partition the input functions (CPTs) into $bucket_1, \dots, bucket_n$ as follows:
 for $i \leftarrow n$ **downto** 1, put in $bucket_i$ all unplaced functions mentioning X_i .
 Put each observed variable in its bucket. Denote by ψ_i the product of input
 functions in $bucket_i$.
2. **backward:** **for** $p \leftarrow n$ **downto** 1 **do**
3. **for** all the functions $\psi_{S_0}, \lambda_{S_1}, \dots, \lambda_{S_j}$ in $bucket_p$ **do**
 If (observed variable) $X_p = x_p$ appears in $bucket_p$,
 assign $X_p = x_p$ to each function in $bucket_p$ and then
 put each resulting function in the bucket of the *closest* variable in its scope.
 else,
4. $\lambda_p \leftarrow \sum_{X_p} \psi_p \cdot \prod_{i=1}^j \lambda_{S_i}$
5. place λ_p in bucket of the latest variable in $scope(\lambda_p)$,
6. **return** (as a result of processing $bucket_1$):
 $P(e) = \alpha = \sum_{X_1} \psi_1 \cdot \prod_{\lambda \in bucket_1} \lambda$
 $P(X_1|e) = \frac{1}{\alpha} \psi_1 \cdot \prod_{\lambda \in bucket_1} \lambda$

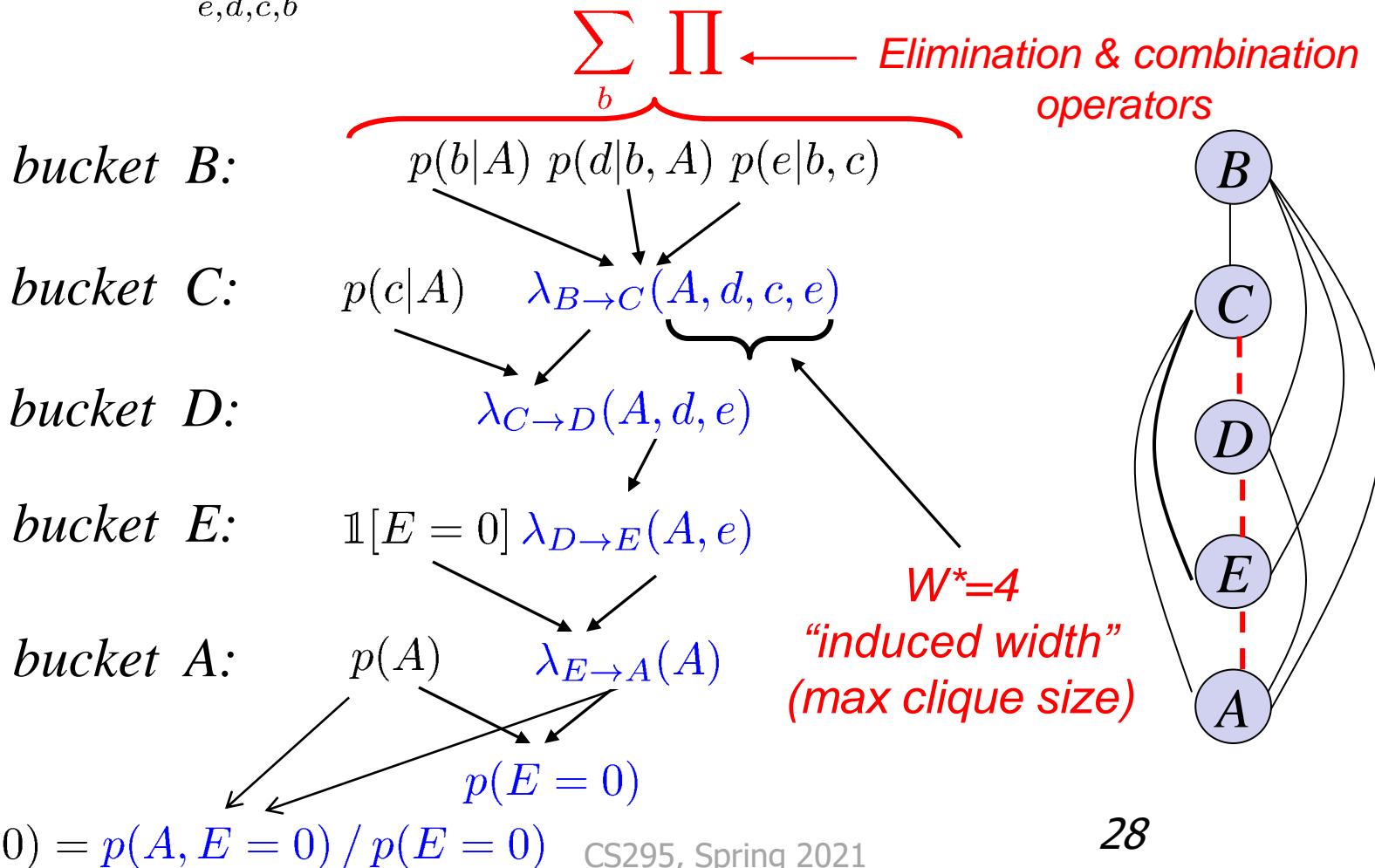
Figure 4.5: BE-bel: a sum-product bucket-elimination algorithm.

Belief Updating

Algorithm BE-bel [Dechter 1996]



$$p(A|E=0) = \alpha \sum_{e,d,c,b} p(A) p(b|A) p(c|A) p(d|A,b) p(e|b,c) \mathbb{1}[e=0]$$

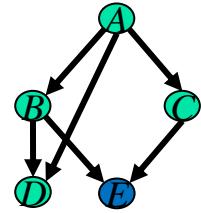




Bucket Elimination

Algorithm BE-bel

[Dechter 1996]

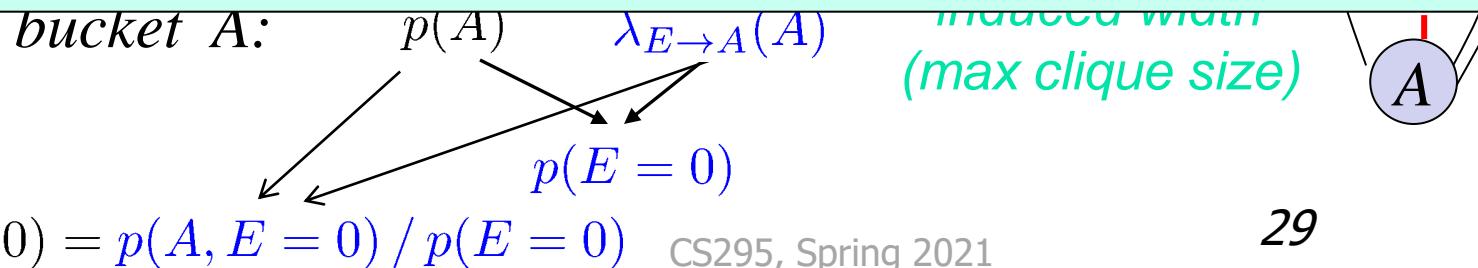


$$p(A|E = 0) = \alpha \sum_{e,d,c,b} p(A) p(b|A) p(c|A) p(d|A, b) p(e|b, c) \mathbb{1}[e = 0]$$

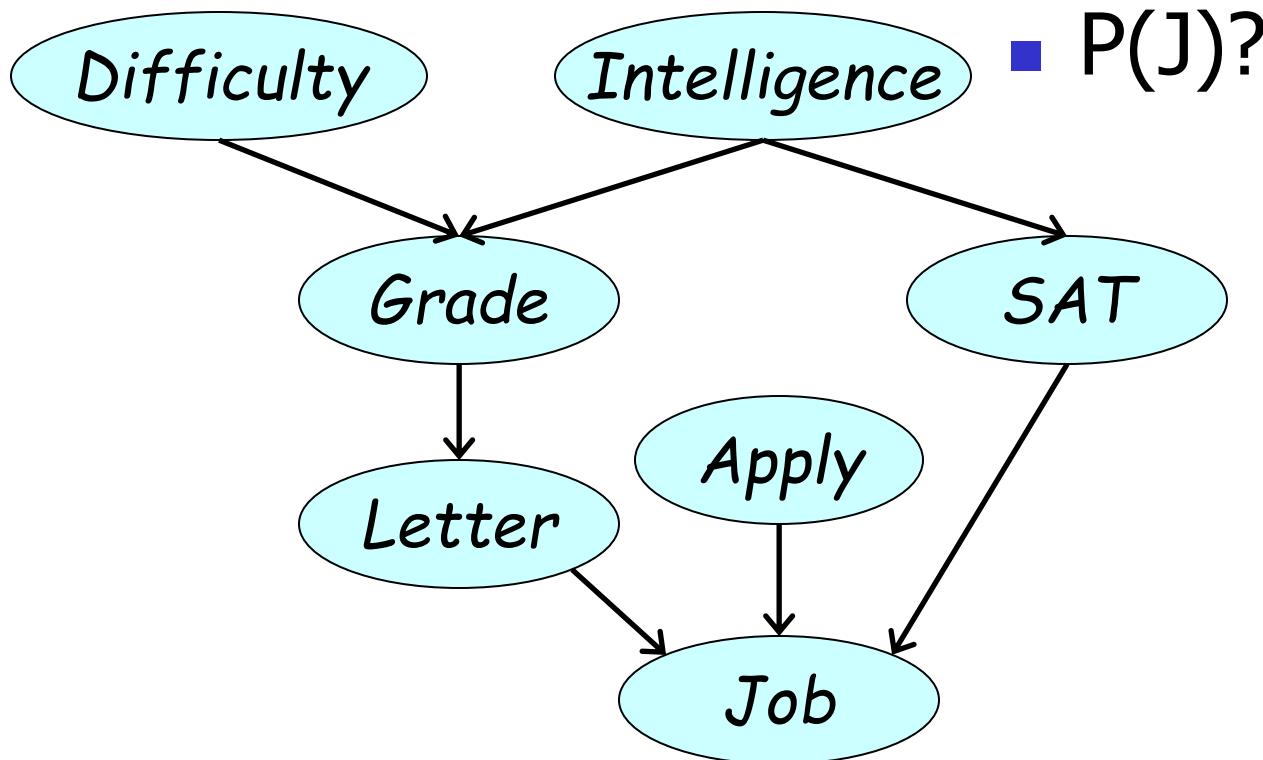
$$\sum_b \prod$$

*Elimination & combination
operators*

Time and space exponential in the induced-width / treewidth



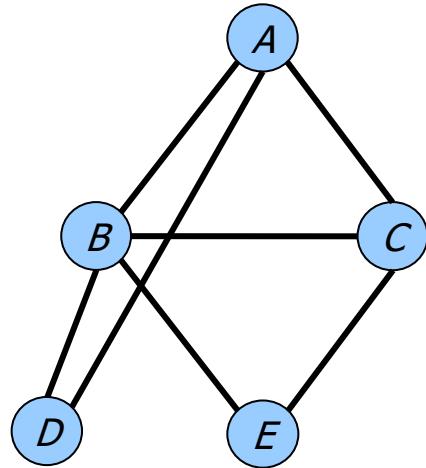
Student Network Example



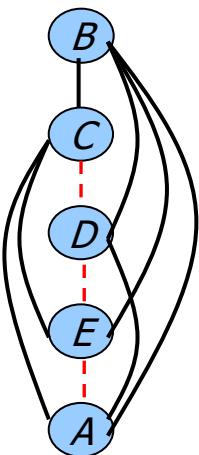
Induced Width (continued)

$w^*(d)$ – the induced width of the primal graph along ordering d

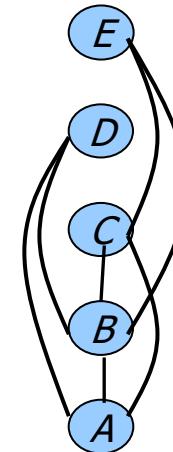
The effect of the ordering:



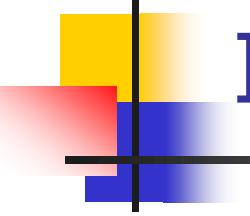
Primal (moral)
graph



$$w^*(d_1) = 4$$



$$w^*(d_2) = 2$$



Inference for probabilistic networks

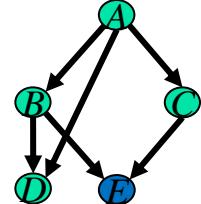
- Bucket elimination

- Belief-updating, $P(e)$, partition function
- Marginals, probability of evidence
- The impact of evidence
 - for MPE (\rightarrow MAP)
 - for MAP (\rightarrow Marginal Map)

- Induced-Width

The impact of evidence?

Algorithm *BE-bel*



$$P(A | E = 0) = \alpha \sum_{E=0,D,C,B} P(A) \cdot P(B | A) \cdot P(C | A) \cdot P(D | A, B) \cdot P(E | B, C)$$

$\sum \prod_b$ ← *Elimination operator*

bucket B:

$$P(b|a) \quad P(d|b,a) \quad P(e|b,c)$$

$B=1$

bucket C:

$$P(c|a) \quad \lambda^B(a, d, c, e)$$

bucket D:

$$\lambda^C(a, d, e)$$

bucket E:

$$e=0 \quad \lambda^D(a, e)$$

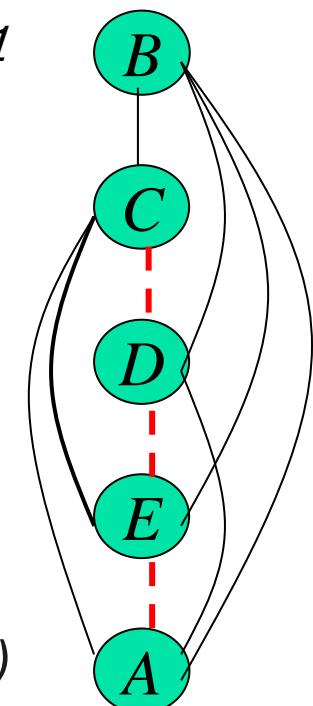
$W^*=4$

bucket A:

$$P(a) \quad \lambda^E(a)$$

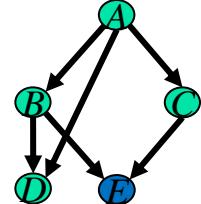
"induced width"
(max clique size)

$$P(a|e=0)$$



The impact of evidence?

Algorithm *BE-bel*



$$P(A | E = 0) = \alpha \sum_{E=0,D,C,B} P(A) \cdot P(B | A) \cdot P(C | A) \cdot P(D | A, B) \cdot P(E | B, C)$$

$P(A | E=0, B=1)?$

bucket B :

$$\underbrace{P(b/a) \quad P(d/b,a) \quad P(e/b,c)}_b$$

$B=1$

bucket C :

$$P(c/a) \quad P(e/b=1,c)$$

bucket D :

$$P(d/b=1,a)$$

bucket E :

$$e=0$$

bucket A :

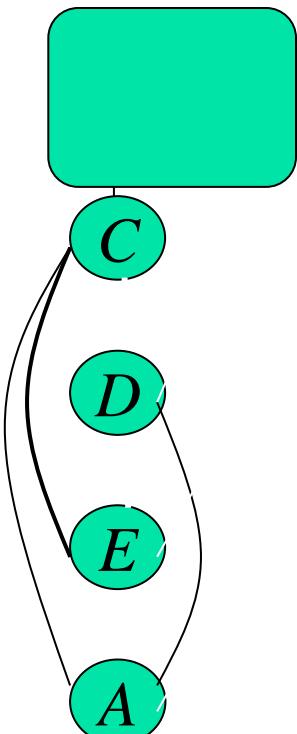
$$P(a)$$

$$P(b=1/a)$$

$$P(e=0)$$

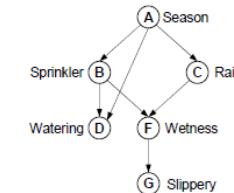
$$P(a|e=0)$$

$$P(a|e=0) = \frac{P(a,e=0)}{P(e=0)}$$

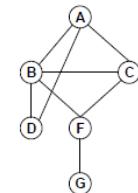


Elimination operator

The impact of observations



(a) Directed acyclic graph



(b) Moral graph

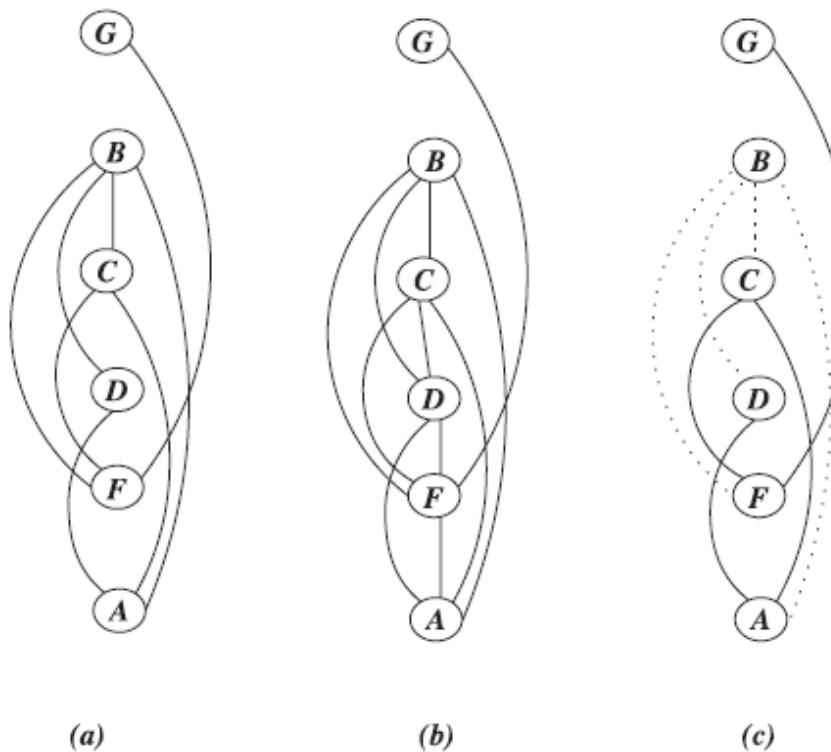


Figure 4.9: Adjusted induced graph relative to observing B .

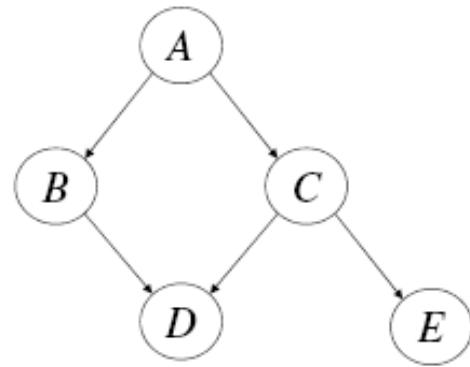
Ordered graph

Induced graph

Ordered conditioned graph

Pruning Nodes: Example

Example of pruning irrelevant subnetworks

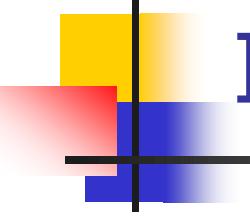


network structure



joint on B, E

joint on B

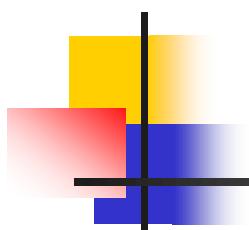


Inference for probabilistic networks

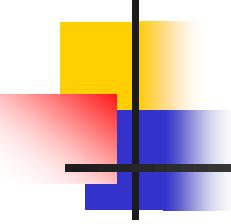
- Bucket elimination

- Belief-updating, $P(e)$, partition function
- Marginals, probability of evidence
- The impact of evidence
- for MPE (\rightarrow MAP)
- for MAP (\rightarrow Marginal Map)

- Induced-Width



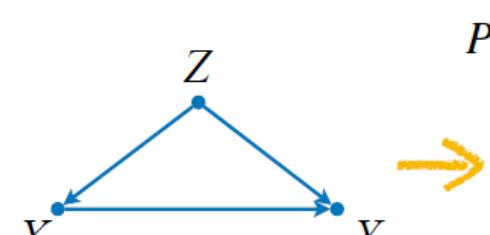
Back to SCM



Markovian Case

- The distribution $P(\mathbf{v})$ decomposes as:

$$P(\mathbf{v}) = \sum_{\mathbf{u}} P(\mathbf{u}) \prod_{V_i \in \mathbf{V}} P(v_i | v_1, \dots, v_{i-1}, \mathbf{u}) = \sum_{\mathbf{u}} P(\mathbf{u}) \prod_{V_i \in \mathbf{V}} P(v_i | pa_i, u_i)$$


$$\begin{aligned} P(z, x, y) &= \sum_{\mathbf{u}} P(\mathbf{u}) P(z | u_z) P(x | z, u_x) P(y | x, z, u_y) \\ &= \left(\sum_{u_z} P(z | u_z) P(u_z) \right) \left(\sum_{u_x} P(x | z, u_x) P(u_x) \right) \left(\sum_{u_y} P(y | x, z, u_y) P(u_y) \right) \\ &= P(z) P(x | z) P(y | x, z) \end{aligned}$$

- In Markovian models, $P(v_i | pa_i)$ can be seen as “canonical factors”.

Markovian Case



$v_z : P(v_z), P(z|v_z) \rightarrow$

$v_x : P(v_x), P(x|v_x, z)$

$v_y : P(v_y), P(y|z, x, v_y)$

$$x, y, z : \begin{cases} \lambda_{v_z}(z) = \sum_{v_z} P(z|v_z) \cdot P(v_z) = ? \\ \lambda_{v_x}(x, z) = \sum_{v_x} P(x|v_x, z) \cdot P(v_x) = ? \\ \lambda_{v_y}(y, x, z) = \sum_{v_y} P(y|z, x, v_y) \cdot P(v_y) = ? \end{cases}$$

$$P(x, y, z) = \lambda_{v_z}(z) \cdot \lambda_{v_x}(x, z) \cdot \lambda_{v_y}(y, x, z)$$

Markovian Case



$v_z : P(v_z), P(z|v_z) \rightarrow$

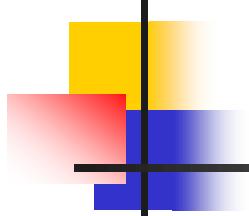
$v_x : P(v_x), P(x|v_x, z)$

$v_y : P(v_y), P(y|z, x, v_y)$

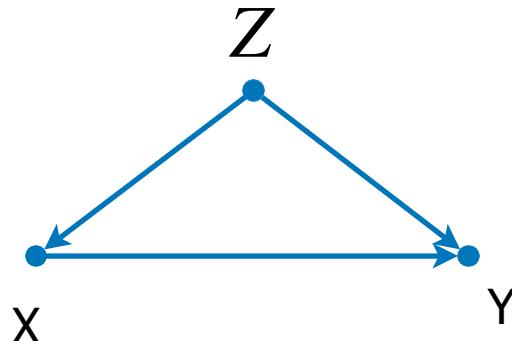
$$x, y, z : \begin{cases} \lambda_{v_z}(z) = \sum_{v_z} P(z|v_z) \cdot P(v_z) = ? \\ \lambda_{v_x}(x, z) = \sum_{v_x} P(x|v_x, z) \cdot P(v_x) = ? \\ \lambda_{v_y}(y, x, z) = \sum_{v_y} P(y|z, x, v_y) \cdot P(v_y) = ? \end{cases}$$

$$P(x, y, z) = \lambda_{v_z}(z) \cdot \lambda_{v_x}(x, z) \cdot \lambda_{v_y}(y, x, z)$$

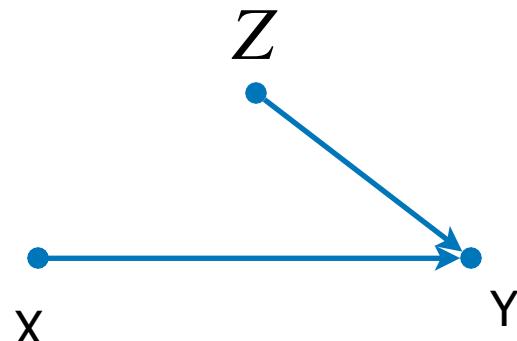
$$P(X, Y, Z) = P(Z)P(X|Z)P(Y|X, Y, Z)$$



Markovian Case



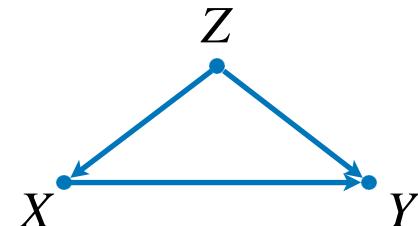
$\text{Do}(X)$ →



Markovian Case

- Every $P(v_i/pa_i)$ is computable from $P(v)$, i.e.,

$$P(v_i/pa_i) = \frac{\sum_{v|v_i, pa_i} P(v)}{\sum_{v|pa_i} P(v)}$$

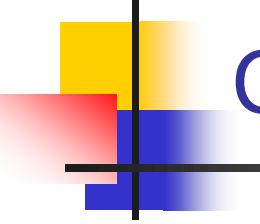


$$P(z, x, y) = \underbrace{P(z)}_{\text{blue}} \underbrace{P(x|z)}_{\text{green}} \underbrace{P(y|x,z)}_{\text{yellow}}$$

$$P(z) = \sum_{x,y} P(v)$$

$$P(y|x,z) = \frac{P(v)}{\sum_y P(v)}$$

$$P(x|z) = \frac{\sum_y P(v)}{\sum_{x,y} P(v)}$$



Complexity of Bucket-Elimination

- **Theorem:**

BE is $O(n \exp(w^*+1))$ time and $O(n \exp(w^*))$ space, when w^* is the induced-width of the moral graph along d when evidence nodes are processed (edges from evidence nodes to earlier variables are removed.)

More accurately: $O(r \exp(w^(d)))$ where r is the number of CPTs.
For Bayesian networks $r=n$. For Markov networks?*