

# Lecture 7: Probabilistic Graphical Models

André Martins



Deep Structured Learning Course, Fall 2020

# Announcements

- Homework 2 is due today!
- Project midterm report is due next week!
- Homework 3 is out, the deadline is December 9. Start early!

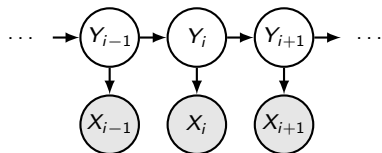
# Slide Credits

- Vlad Niculae (co-instructor of DSL last year)

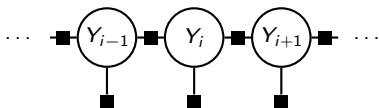
# Graphical Models

In this unit, we will formalize & extend these graphical representations encountered in previous lectures.

**Directed**



**Undirected**



# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

Markov random fields

Factor graphs

# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

Markov random fields

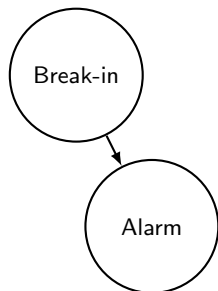
Factor graphs

# Bayes (belief) networks

- Common task: Characterize how some related events co-occur.  
Specifically, in terms of **probabilities!**
- A car alarm is going off. Was there a break-in?

# Bayes (belief) networks

- Common task: Characterize how some related events co-occur. Specifically, in terms of **probabilities!**
- A car alarm is going off. Was there a break-in?



| $P(B)$ | B=yes | B=no |
|--------|-------|------|
|        | .05   | .95  |

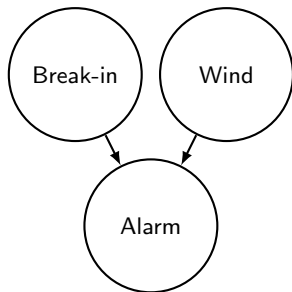
| $P(A   B)$ | A=on | A=off |
|------------|------|-------|
| B=yes      | .99  | .01   |
| B=no       | .10  | .90   |

- $\mathbb{P}(B | A) = ?$



# Bayes (belief) networks

- Common task: Characterize how some related events co-occur. Specifically, in terms of **probabilities!**
- A car alarm is going off. Was there a break-in?

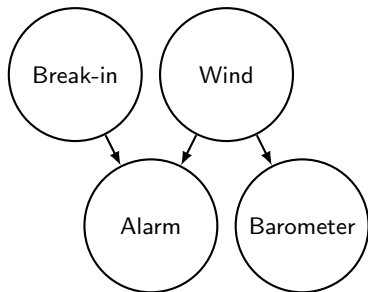


|               |       | $P(B)$ |       |
|---------------|-------|--------|-------|
|               |       | B=yes  | B=no  |
|               |       | .05    | .95   |
| $P(A   B, W)$ |       | A=on   | A=off |
| B=yes         | W=lo  | .99    | .01   |
| B=yes         | W=med | .99    | .01   |
| B=yes         | W=hi  | .999   | .001  |
| B=no          | W=lo  | .01    | .99   |
| B=no          | W=med | .05    | .95   |
| B=no          | W=hi  | .25    | .75   |

- $\mathbb{P}(B | A) = ?$  Can we observe wind?  $\mathbb{P}(B | A, W) = ?$

# Bayes (belief) networks

- Common task: Characterize how some related events co-occur. Specifically, in terms of **probabilities!**
- A car alarm is going off. Was there a break-in?

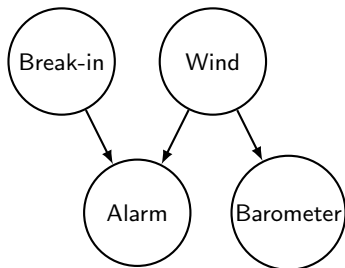


|               |       | $P(B)$ |       |
|---------------|-------|--------|-------|
|               |       | B=yes  | B=no  |
|               |       | .05    | .95   |
| $P(A   B, W)$ |       | A=on   | A=off |
| B=yes         | W=lo  | .99    | .01   |
| B=yes         | W=med | .99    | .01   |
| B=yes         | W=hi  | .999   | .001  |
| B=no          | W=lo  | .01    | .99   |
| B=no          | W=med | .05    | .95   |
| B=no          | W=hi  | .25    | .75   |

- $\mathbb{P}(B | A) = ?$  Can we observe wind?  $\mathbb{P}(B | A, W) = ?$   
Maybe we're in the basement, but have a barometer.

# Bayes networks

Toolkit for encoding **knowledge** about **interaction structures** between rv's.



Directed acyclic graph (DAG). Nodes = variables. Arrows = statistical dependencies.

$$\text{In general: } \mathbb{P}(X_1, \dots, X_n) = \prod_i \mathbb{P}(X_i \mid \text{parents}(X_i))$$

$$\begin{aligned} &\text{For example: } \mathbb{P}(\text{Break-in, Wind, Alarm, Barometer}) \\ &= \mathbb{P}(\text{Break-in})\mathbb{P}(\text{Wind})\mathbb{P}(\text{Alarm} \mid \text{Break-in, Wind})\mathbb{P}(\text{Barometer} \mid \text{Wind}) \end{aligned}$$

Without any structure,  $\mathbb{P}(\text{Break-in, Wind, Alarm, Barometer})$   
 would have to be stored & estimated like

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| yes  | lo   | on    | lo   | 0.0243   |
| yes  | lo   | on    | med  | 0.0002   |
| yes  | lo   | on    | hi   | 0.0002   |
| yes  | lo   | off   | lo   | 0.0002   |
| yes  | lo   | off   | med  | 2.50e-06 |
| yes  | lo   | off   | hi   | 2.50e-06 |
| yes  | med  | on    | lo   | 0.0001   |
| yes  | med  | on    | med  | 0.0146   |
| yes  | med  | on    | hi   | 0.0001   |
| yes  | med  | off   | lo   | 1.50e-06 |
| yes  | med  | off   | med  | 0.0001   |
| yes  | med  | off   | hi   | 1.50e-06 |
| yes  | hi   | on    | lo   | 9.99e-05 |
| yes  | hi   | on    | med  | 9.99e-05 |
| yes  | hi   | on    | hi   | 0.0098   |
| yes  | hi   | off   | lo   | 1.00e-07 |
| yes  | hi   | off   | med  | 1.00e-07 |
| yes  | hi   | off   | hi   | 9.80e-06 |

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| no   | lo   | on    | lo   | 0.0047   |
| no   | lo   | on    | med  | 4.75e-05 |
| no   | lo   | on    | hi   | 4.75e-05 |
| no   | lo   | off   | lo   | 0.4608   |
| no   | lo   | off   | med  | 0.0047   |
| no   | lo   | off   | hi   | 0.0047   |
| no   | med  | on    | lo   | 0.0001   |
| no   | med  | on    | med  | 0.0140   |
| no   | med  | on    | hi   | 0.0001   |
| no   | med  | off   | lo   | 0.0027   |
| no   | med  | off   | med  | 0.2653   |
| no   | med  | off   | hi   | 0.0027   |
| no   | hi   | on    | lo   | 0.0005   |
| no   | hi   | on    | med  | 0.0005   |
| no   | hi   | on    | hi   | 0.0466   |
| no   | hi   | off   | lo   | 0.0014   |
| no   | hi   | off   | med  | 0.0014   |
| no   | hi   | off   | hi   | 0.1397   |

Without any structure,  $\mathbb{P}(\text{Break-in, Wind, Alarm, Barometer})$   
would have to be stored & estimated like

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| yes  | lo   | on    | lo   | 0.0243   |
| yes  | lo   | on    | med  | 0.0002   |
| yes  | lo   | on    | hi   | 0.0002   |
| yes  | lo   | off   | lo   | 0.0002   |
| yes  | lo   | off   | med  | 2.50e-06 |
| yes  | lo   | off   | hi   | 2.50e-06 |
| yes  | med  | on    | lo   | 0.0001   |
| yes  | med  | on    | med  | 0.0146   |
| yes  | med  | on    | hi   | 0.0001   |
| yes  | med  | off   | lo   | 1.50e-06 |
| yes  | med  | off   | med  | 0.0001   |
| yes  | med  | off   | hi   | 1.50e-06 |
| yes  | hi   | on    | lo   | 9.99e-05 |
| yes  | hi   | on    | med  | 9.99e-05 |
| yes  | hi   | on    | hi   | 0.0098   |
| yes  | hi   | off   | lo   | 1.00e-07 |
| yes  | hi   | off   | med  | 1.00e-07 |
| yes  | hi   | off   | hi   | 9.80e-06 |

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| no   | lo   | on    | lo   | 0.0047   |
| no   | lo   | on    | med  | 4.75e-05 |
| no   | lo   | on    | hi   | 4.75e-05 |
| no   | lo   | off   | lo   | 0.4608   |
| no   | lo   | off   | med  | 0.0047   |
| no   | lo   | off   | hi   | 0.0047   |
| no   | med  | on    | lo   | 0.0001   |
| no   | med  | on    | med  | 0.0140   |
| no   | med  | on    | hi   | 0.0001   |
| no   | med  | off   | lo   | 0.0027   |
| no   | med  | off   | med  | 0.2653   |
| no   | med  | off   | hi   | 0.0027   |
| no   | hi   | on    | lo   | 0.0005   |
| no   | hi   | on    | med  | 0.0005   |
| no   | hi   | on    | hi   | 0.0466   |
| no   | hi   | off   | lo   | 0.0014   |
| no   | hi   | off   | med  | 0.0014   |
| no   | hi   | off   | hi   | 0.1397   |

$$\mathbb{P}(\text{Break-in}=\text{yes}, \text{Alarm}=\text{on}) = 0.0496$$

Without any structure,  $\mathbb{P}(\text{Break-in, Wind, Alarm, Barometer})$   
would have to be stored & estimated like

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| yes  | lo   | on    | lo   | 0.0243   |
| yes  | lo   | on    | med  | 0.0002   |
| yes  | lo   | on    | hi   | 0.0002   |
| yes  | lo   | off   | lo   | 0.0002   |
| yes  | lo   | off   | med  | 2.50e-06 |
| yes  | lo   | off   | hi   | 2.50e-06 |
| yes  | med  | on    | lo   | 0.0001   |
| yes  | med  | on    | med  | 0.0146   |
| yes  | med  | on    | hi   | 0.0001   |
| yes  | med  | off   | lo   | 1.50e-06 |
| yes  | med  | off   | med  | 0.0001   |
| yes  | med  | off   | hi   | 1.50e-06 |
| yes  | hi   | on    | lo   | 9.99e-05 |
| yes  | hi   | on    | med  | 9.99e-05 |
| yes  | hi   | on    | hi   | 0.0098   |
| yes  | hi   | off   | lo   | 1.00e-07 |
| yes  | hi   | off   | med  | 1.00e-07 |
| yes  | hi   | off   | hi   | 9.80e-06 |

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| no   | lo   | on    | lo   | 0.0047   |
| no   | lo   | on    | med  | 4.75e-05 |
| no   | lo   | on    | hi   | 4.75e-05 |
| no   | lo   | off   | lo   | 0.4608   |
| no   | lo   | off   | med  | 0.0047   |
| no   | lo   | off   | hi   | 0.0047   |
| no   | med  | on    | lo   | 0.0001   |
| no   | med  | on    | med  | 0.0140   |
| no   | med  | on    | hi   | 0.0001   |
| no   | med  | off   | lo   | 0.0027   |
| no   | med  | off   | med  | 0.2653   |
| no   | med  | off   | hi   | 0.0027   |
| no   | hi   | on    | lo   | 0.0005   |
| no   | hi   | on    | med  | 0.0005   |
| no   | hi   | on    | hi   | 0.0466   |
| no   | hi   | off   | lo   | 0.0014   |
| no   | hi   | off   | med  | 0.0014   |
| no   | hi   | off   | hi   | 0.1397   |

$$\mathbb{P}(\text{Break-in=yes, Alarm=on}) = 0.0496$$

$$\mathbb{P}(\text{Break-in=no, Alarm=on}) = 0.0665$$

Without any structure,  $\mathbb{P}(\text{Break-in, Wind, Alarm, Barometer})$   
would have to be stored & estimated like

| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| yes  | lo   | on    | lo   | 0.0243   |
| yes  | lo   | on    | med  | 0.0002   |
| yes  | lo   | on    | hi   | 0.0002   |
| yes  | lo   | off   | lo   | 0.0002   |
| yes  | lo   | off   | med  | 2.50e-06 |
| yes  | lo   | off   | hi   | 2.50e-06 |
| yes  | med  | on    | lo   | 0.0001   |
| yes  | med  | on    | med  | 0.0146   |
| yes  | med  | on    | hi   | 0.0001   |
| yes  | med  | off   | lo   | 1.50e-06 |
| yes  | med  | off   | med  | 0.0001   |
| yes  | med  | off   | hi   | 1.50e-06 |
| yes  | hi   | on    | lo   | 9.99e-05 |
| yes  | hi   | on    | med  | 9.99e-05 |
| yes  | hi   | on    | hi   | 0.0098   |
| yes  | hi   | off   | lo   | 1.00e-07 |
| yes  | hi   | off   | med  | 1.00e-07 |
| yes  | hi   | off   | hi   | 9.80e-06 |

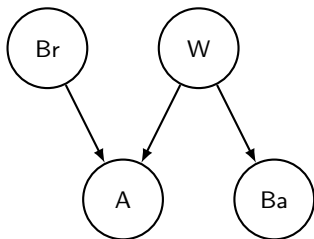
| Brk. | Wind | Alarm | Bar. | P        |
|------|------|-------|------|----------|
| no   | lo   | on    | lo   | 0.0047   |
| no   | lo   | on    | med  | 4.75e-05 |
| no   | lo   | on    | hi   | 4.75e-05 |
| no   | lo   | off   | lo   | 0.4608   |
| no   | lo   | off   | med  | 0.0047   |
| no   | lo   | off   | hi   | 0.0047   |
| no   | med  | on    | lo   | 0.0001   |
| no   | med  | on    | med  | 0.0140   |
| no   | med  | on    | hi   | 0.0001   |
| no   | med  | off   | lo   | 0.0027   |
| no   | med  | off   | med  | 0.2653   |
| no   | med  | off   | hi   | 0.0027   |
| no   | hi   | on    | lo   | 0.0005   |
| no   | hi   | on    | med  | 0.0005   |
| no   | hi   | on    | hi   | 0.0466   |
| no   | hi   | off   | lo   | 0.0014   |
| no   | hi   | off   | med  | 0.0014   |
| no   | hi   | off   | hi   | 0.1397   |

$$\mathbb{P}(\text{Break-in=yes, Alarm=on}) = 0.0496$$

$$\mathbb{P}(\text{Break-in=no, Alarm=on}) = 0.0665$$

$$\begin{aligned} \mathbb{P}(\text{Break-in=yes} \mid \text{Alarm=on}) &= \frac{\mathbb{P}(\text{Break-in=yes, Alarm=on})}{\sum_b \mathbb{P}(\text{Break-in}=b, \text{Alarm=on})} \\ &= .427 \end{aligned}$$

Knowing the model structure (statistical dependencies), complicated models become manageable.



$$\mathbb{P}(\text{Br}, W, A, \text{Ba})$$

$$= \mathbb{P}(\text{Br})\mathbb{P}(W)\mathbb{P}(A \mid \text{Br}, W)\mathbb{P}(\text{Ba} \mid W)$$

| $P(\text{Br})$ | yes | no  |
|----------------|-----|-----|
|                | .05 | .95 |

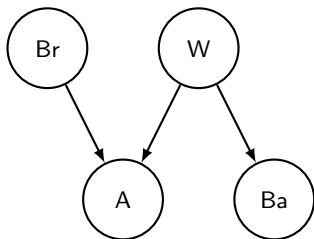
| $P(W)$ | lo | mid | hi |
|--------|----|-----|----|
|        | .5 | .3  | .2 |

| $P(A \mid \text{Br}, W)$ |       | on   | off  |
|--------------------------|-------|------|------|
| Br=yes                   | W=lo  | .99  | .01  |
| Br=yes                   | W=med | .99  | .01  |
| Br=yes                   | W=hi  | .999 | .001 |
| Br=no                    | W=lo  | .01  | .99  |
| Br=no                    | W=med | .05  | .95  |
| Br=no                    | W=hi  | .25  | .75  |

| $P(\text{Ba} \mid W)$ | lo  | mid | hi  |
|-----------------------|-----|-----|-----|
| W=lo                  | .98 | .01 | .01 |
| W=mid                 | .01 | .98 | .01 |
| W=hi                  | .01 | .01 | .98 |



Knowing the model structure (statistical dependencies), complicated models become manageable.



$$\mathbb{P}(\text{Br}, W, A, \text{Ba}) \\ = \mathbb{P}(\text{Br})\mathbb{P}(W)\mathbb{P}(A \mid \text{Br}, W)\mathbb{P}(\text{Ba} \mid W)$$

- Can estimate parts in isolation  
e.g.  $\mathbb{P}(\text{Wind})$  from weather history.

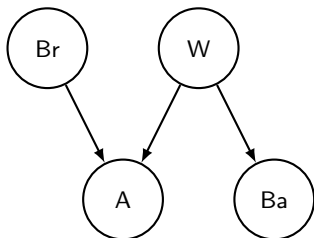
| $P(\text{Br})$ | yes | no  |
|----------------|-----|-----|
|                | .05 | .95 |

| $P(W)$ | lo | mid | hi |
|--------|----|-----|----|
|        | .5 | .3  | .2 |

| $P(A \mid \text{Br}, W)$ | on   | off  |
|--------------------------|------|------|
| Br=yes    W=lo           | .99  | .01  |
| Br=yes    W=med          | .99  | .01  |
| Br=yes    W=hi           | .999 | .001 |
| Br=no    W=lo            | .01  | .99  |
| Br=no    W=med           | .05  | .95  |
| Br=no    W=hi            | .25  | .75  |

| $P(\text{Ba} \mid W)$ | lo  | mid | hi  |
|-----------------------|-----|-----|-----|
| W=lo                  | .98 | .01 | .01 |
| W=mid                 | .01 | .98 | .01 |
| W=hi                  | .01 | .01 | .98 |

Knowing the model structure (statistical dependencies), complicated models become manageable.



$$\mathbb{P}(\text{Br}, W, A, \text{Ba}) \\ = \mathbb{P}(\text{Br})\mathbb{P}(W)\mathbb{P}(A | \text{Br}, W)\mathbb{P}(\text{Ba} | W)$$

- Can estimate parts in isolation  
e.g.  $\mathbb{P}(\text{Wind})$  from weather history.
- Can sample by following the graph  
from roots to leaves.

|                |     |     |    |
|----------------|-----|-----|----|
| $P(\text{Br})$ | yes | no  |    |
|                | .05 | .95 |    |
| $P(W)$         | lo  | mid | hi |
|                | .5  | .3  | .2 |

|                       |      |      |
|-----------------------|------|------|
| $P(A   \text{Br}, W)$ | on   | off  |
| Br=yes    W=lo        | .99  | .01  |
| Br=yes    W=med       | .99  | .01  |
| Br=yes    W=hi        | .999 | .001 |
| Br=no    W=lo         | .01  | .99  |
| Br=no    W=med        | .05  | .95  |
| Br=no    W=hi         | .25  | .75  |

|                    |     |     |     |
|--------------------|-----|-----|-----|
| $P(\text{Ba}   W)$ | lo  | mid | hi  |
| W=lo               | .98 | .01 | .01 |
| W=mid              | .01 | .98 | .01 |
| W=hi               | .01 | .01 | .98 |

# Bayes Nets:

reduce number of parameters & aid estimation

let us reason about **independencies** in a model

are a building-block for modeling **causality**

# Bayes Nets:

are not neural network diagrams

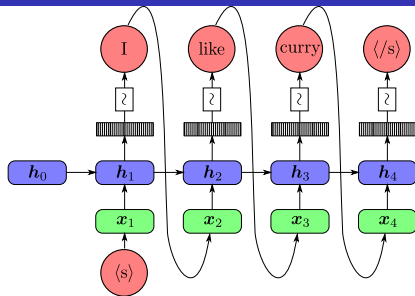
encode structure, not parametrization

are non-unique for a distribution

encode independence **requirements**, not necessarily all

# BN are not neural net diagrams

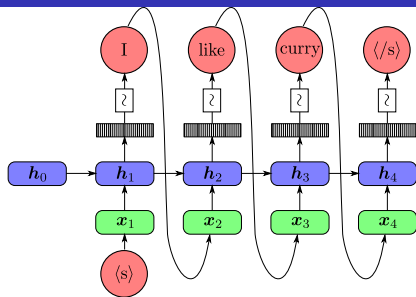
Recall the RNN language model:



- In statistical terms, what are we modeling?

# BN are not neural net diagrams

Recall the RNN language model:

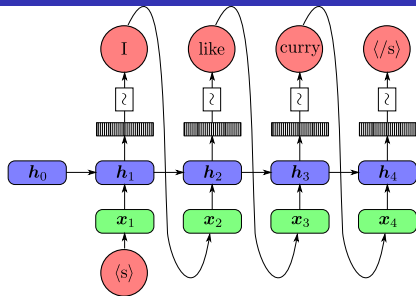


- In statistical terms, what are we modeling?

$$\mathbb{P}(X_1, \dots, X_n) = \mathbb{P}(X_1)\mathbb{P}(X_2 | X_1)\mathbb{P}(X_3 | X_1, X_2) \dots$$

# BN are not neural net diagrams

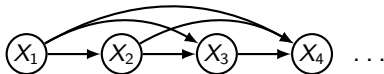
Recall the RNN language model:

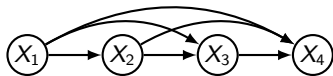
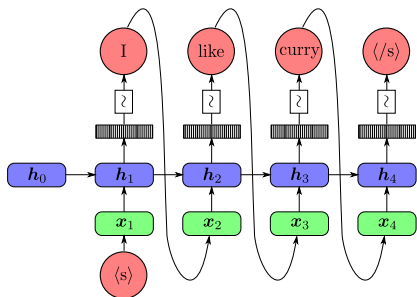


- In statistical terms, what are we modeling?

$$\mathbb{P}(X_1, \dots, X_n) = \mathbb{P}(X_1)\mathbb{P}(X_2 | X_1)\mathbb{P}(X_3 | X_1, X_2)\dots$$

- Bayes Net:
- Not useful! Everything conditionally depends on everything. (more later)





Neural net diagrams  
(and computation graphs)  
show **how to compute something**

Bayes networks  
show **how a distribution factorizes**  
(what is assumed independent)



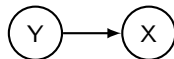
# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X} =$  all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .

BN for a generative model:



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

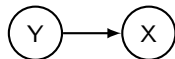
# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X} =$  all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .

BN for a generative model:



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

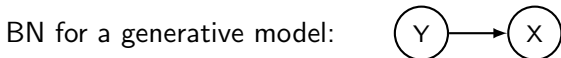
$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ ,

# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X} =$  all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

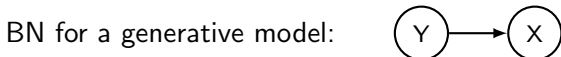
$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X} =$  all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

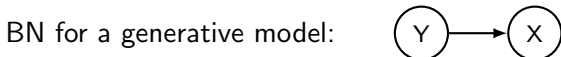
$\mathbb{P}(X | Y)$  (remember: values of  $X$  are sentences)

# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X} =$  all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

$\mathbb{P}(X | Y)$  (remember: values of  $X$  are sentences)

Naive Bayes

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | Y)$$

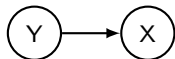
# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X}$  = all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .

BN for a generative model:



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

$\mathbb{P}(X | Y)$  (remember: values of  $X$  are sentences)

Naive Bayes

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | Y)$$

Per-class Markov language model

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | X_{j-1}, Y)$$

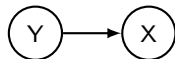
# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X}$  = all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .

BN for a generative model:



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

$\mathbb{P}(X | Y)$  (remember: values of  $X$  are sentences)

Naive Bayes

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | Y)$$

Per-class Markov language model

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | X_{j-1}, Y)$$

Per-class recurrent NN language model

$$\mathbb{P}(X | Y) = \text{LSTM}(x_1, \dots, x_L; \mathbf{w}_Y)$$

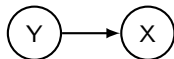
# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X}$  = all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .

BN for a generative model:



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

$\mathbb{P}(X | Y)$  (remember: values of  $X$  are sentences)

Naive Bayes

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | Y)$$

Per-class Markov language model

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | X_{j-1}, Y)$$

Per-class recurrent NN language model

$$\mathbb{P}(X | Y) = \text{LSTM}(x_1, \dots, x_L; \mathbf{w}_Y)$$

$\mathbb{P}(X | Y)$  need not be parametrized as a table.



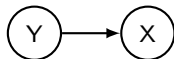
# BN encode structure, not parametrization

A BN tells us: **how the distribution decomposes**

A BN can't tell us: **what the probabilities are!**

Example:  $X \in \mathcal{X}$  = all English sentences,  $Y \in \{\text{sports, music, ...}\}$ .

BN for a generative model:



We must posit what are  $\mathbb{P}(Y)$  and  $\mathbb{P}(X | Y)$ . **Many possible options!**

$\mathbb{P}(Y)$ : uniform:  $\mathbb{P}(Y = \text{sports}) = \mathbb{P}(Y = \text{music}) = \frac{1}{|Y|}$ , or estimated from data.

$\mathbb{P}(X | Y)$  (remember: values of  $X$  are sentences)

Naive Bayes

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | Y)$$

Per-class Markov language model

$$\mathbb{P}(X | Y) = \prod_{j=1}^L \mathbb{P}(X_j | X_{j-1}, Y)$$

Per-class recurrent NN language model

$$\mathbb{P}(X | Y) = \text{LSTM}(x_1, \dots, x_L; \mathbf{w}_y)$$

$\mathbb{P}(X | Y)$  need not be parametrized as a table.

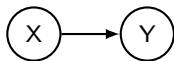
rv's need not be discrete! mixture of Gaussians:  $\mathbb{P}(X | Y = y) \sim \mathcal{N}(\mu_y, \Sigma_y)$ .

# Equivalent factorizations

There are many possible factorizations!  $\mathbb{P}(X, Y) =$

# Equivalent factorizations

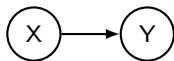
There are many possible factorizations!  $\mathbb{P}(X, Y) =$



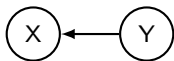
$$\mathbb{P}(X)\mathbb{P}(Y | X)$$

# Equivalent factorizations

There are many possible factorizations!  $\mathbb{P}(X, Y) =$



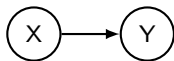
$$\mathbb{P}(X)\mathbb{P}(Y | X)$$



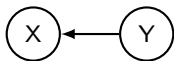
$$\mathbb{P}(Y)\mathbb{P}(X | Y)$$

# Equivalent factorizations

There are many possible factorizations!  $\mathbb{P}(X, Y) =$



$$\mathbb{P}(X)\mathbb{P}(Y | X)$$



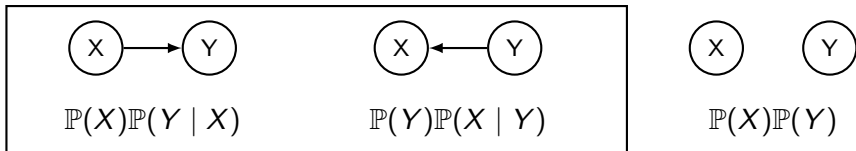
$$\mathbb{P}(Y)\mathbb{P}(X | Y)$$



$$\mathbb{P}(X)\mathbb{P}(Y)$$

# Equivalent factorizations

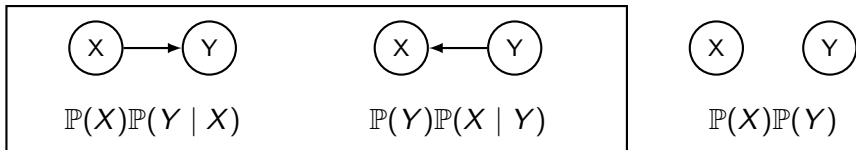
There are many possible factorizations!  $\mathbb{P}(X, Y) =$



The first two are valid Bayes nets for **any**  $\mathbb{P}(X, Y)$ !

# Equivalent factorizations

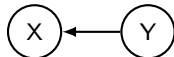
There are many possible factorizations!  $\mathbb{P}(X, Y) =$



The first two are valid Bayes nets for **any**  $\mathbb{P}(X, Y)$ !

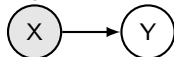
In fact, recall generative vs discriminative classifiers!

- Generative (e.g. naïve Bayes):



*To classify, we would compute  $\mathbb{P}(Y | X)$  via Bayes' rule.*

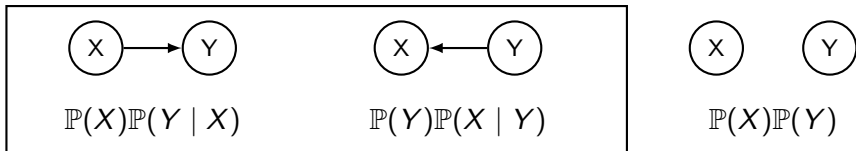
- Discriminative (e.g. logistic regression)



*in LR, we don't model  $\mathbb{P}(X)$ , we assume  $X$  is always observed (gray).*

# Equivalent factorizations

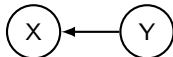
There are many possible factorizations!  $\mathbb{P}(X, Y) =$



The first two are valid Bayes nets for **any**  $\mathbb{P}(X, Y)$ !

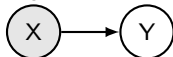
In fact, recall generative vs discriminative classifiers!

- Generative (e.g. naïve Bayes):



*To classify, we would compute  $\mathbb{P}(Y | X)$  via Bayes' rule.*

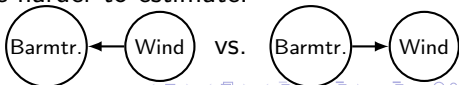
- Discriminative (e.g. logistic regression)



*in LR, we don't model  $\mathbb{P}(X)$ , we assume  $X$  is always observed (gray).*

Some arrow direction choices are harder to estimate.


Some make more sense (why?):






# Minimal independence assumptions

Recall, we say  $X \perp\!\!\!\perp Y$  iff.  $P(X, Y) = P(X)P(Y)$   
Let  $X =$  grade in DSL,  $Y =$  month you were born.

Bayes net (1): 

# Minimal independence assumptions

Recall, we say  $X \perp\!\!\!\perp Y$  iff.  $P(X, Y) = P(X)P(Y)$   
Let  $X =$  grade in DSL,  $Y =$  month you were born.

Bayes net (1):  A Bayesian network diagram consisting of two nodes, X and Y, each enclosed in a circle. There are no edges between the two nodes, representing independence.

Example parametrization:


| $P(X)$ | A+  | A   | B   | ... |
|--------|-----|-----|-----|-----|
|        | .01 | .02 | .04 |     |

---

| $P(Y)$ | Jan | Feb | Mar | ... |
|--------|-----|-----|-----|-----|
|        | .10 | .12 | .09 |     |

# Minimal independence assumptions

Recall, we say  $X \perp\!\!\!\perp Y$  iff.  $P(X, Y) = P(X)P(Y)$   
Let  $X =$  grade in DSL,  $Y =$  month you were born.

Bayes net (1): 

Example parametrization:

| $P(X)$ | A+  | A   | B   | ... |
|--------|-----|-----|-----|-----|
|        | .01 | .02 | .04 |     |


---

| $P(Y)$ | Jan | Feb | Mar | ... |
|--------|-----|-----|-----|-----|
|        | .10 | .12 | .09 |     |

BN (1) imposes  $X \perp\!\!\!\perp Y$   
in **any parametrization**.

# Minimal independence assumptions

Recall, we say  $X \perp\!\!\!\perp Y$  iff.  $P(X, Y) = P(X)P(Y)$   
Let  $X =$  grade in DSL,  $Y =$  month you were born.

Bayes net (1): 

Bayes net (2): 

Example parametrization:

| P(X) | A+  | A   | B   | ... |
|------|-----|-----|-----|-----|
|      | .01 | .02 | .04 |     |

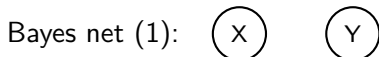
| P(Y) | Jan | Feb | Mar | ... |
|------|-----|-----|-----|-----|
|      | .10 | .12 | .09 |     |

Does it mean  $X \not\perp\!\!\!\perp Y$  necessarily?

BN (1) imposes  $X \perp\!\!\!\perp Y$   
in **any parametrization**.

# Minimal independence assumptions

Recall, we say  $X \perp\!\!\!\perp Y$  iff.  $P(X, Y) = P(X)P(Y)$   
Let  $X =$  grade in DSL,  $Y =$  month you were born.



Example parametrization:

|      |     |     |     |     |
|------|-----|-----|-----|-----|
| P(X) | A+  | A   | B   | ... |
|      | .01 | .02 | .04 |     |
| P(Y) | Jan | Feb | Mar | ... |
|      | .10 | .12 | .09 |     |

BN (1) imposes  $X \perp\!\!\!\perp Y$   
in **any parametrization.**

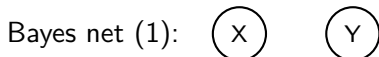
Does it mean  $X \not\perp\!\!\!\perp Y$  necessarily?

**NO!**

|          |     |     |     |     |
|----------|-----|-----|-----|-----|
| P(Y)     | Jan | Feb | Mar | ... |
|          | .10 | .12 | .09 |     |
| P(X   Y) | 20  | 19  | 18  | ... |
| Y=Jan    | .01 | .02 | .04 |     |
| Y=Feb    | .01 | .02 | .04 |     |
| Y=Mar    | .01 | .02 | .04 |     |
| ...      |     |     |     |     |

# Minimal independence assumptions

Recall, we say  $X \perp\!\!\!\perp Y$  iff.  $P(X, Y) = P(X)P(Y)$   
Let  $X =$  grade in DSL,  $Y =$  month you were born.



Example parametrization:

|      |     |     |     |     |
|------|-----|-----|-----|-----|
| P(X) | A+  | A   | B   | ... |
|      | .01 | .02 | .04 |     |
| P(Y) | Jan | Feb | Mar | ... |
|      | .10 | .12 | .09 |     |

BN (1) imposes  $X \perp\!\!\!\perp Y$   
in **any parametrization**.

Does it mean  $X \not\perp\!\!\!\perp Y$  necessarily?

**NO!**

|          |     |     |     |     |
|----------|-----|-----|-----|-----|
| P(Y)     | Jan | Feb | Mar | ... |
|          | .10 | .12 | .09 |     |
| P(X   Y) | 20  | 19  | 18  | ... |
| Y=Jan    | .01 | .02 | .04 |     |
| Y=Feb    | .01 | .02 | .04 |     |
| Y=Mar    | .01 | .02 | .04 |     |
| ...      |     |     |     |     |

A BN expresses which independences **must exist**, but there can be additional ones.

# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

Markov random fields

Factor graphs

# Conditional independence in Bayes nets

Identifying independences in a distribution is generally hard.

Bayes nets let us reason about it via graph algorithms!

## Definition (conditional independence)

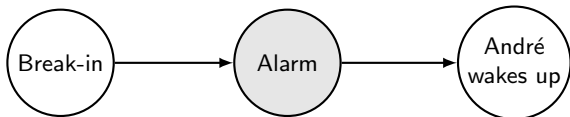
$A$  is independent of  $B$  given a set of variables  $C = \{C_1, \dots, C_n\}$ , denoted

$$A \perp\!\!\!\perp B \mid C,$$

iff  $\mathbb{P}(A, B \mid C_1, \dots, C_n) = \mathbb{P}(A \mid C_1, \dots, C_n)\mathbb{P}(B \mid C_1, \dots, C_n)$ .

**Note.** Equivalently,  $\mathbb{P}(A \mid B, C_1, \dots, C_n) = \mathbb{P}(A \mid C_1, \dots, C_n)$ .

Intuitively: if we observe  $C$ , does observing  $B$  too bring us more info about  $A$ ?

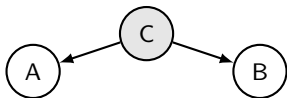


You want to assess if I'm awake. Does it matter if there really was a break-in?

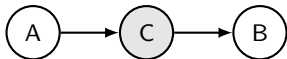


# Three fundamental relationships in BN

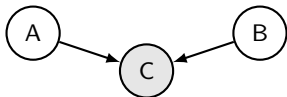
The Fork



The Chain

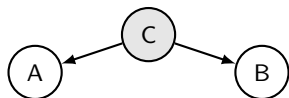


The Collider



# Three fundamental relationships in BN

The Fork

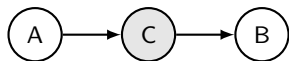


$$A \perp\!\!\!\perp B \mid C$$

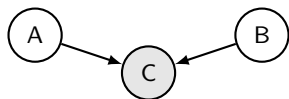
Given  $C$ ,  $A$  and  $B$  are independent.

Example: Alarm  $\leftarrow$  Wind  $\rightarrow$  Barometer

The Chain

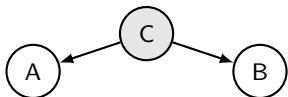


The Collider



# Three fundamental relationships in BN

The Fork

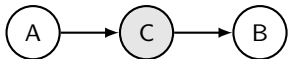


$$A \perp\!\!\!\perp B \mid C$$

Given  $C$ ,  $A$  and  $B$  are independent.

Example: Alarm  $\leftarrow$  Wind  $\rightarrow$  Barometer

The Chain



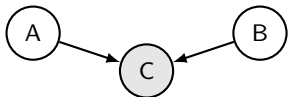
$$A \perp\!\!\!\perp B \mid C$$

After observing  $C$ ,

further observing  $A$  would not tell us about  $B$ .

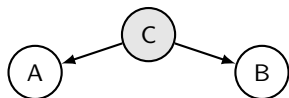
Example: Burglary  $\rightarrow$  Alarm  $\rightarrow$  André wakes up

The Collider



# Three fundamental relationships in BN

The Fork

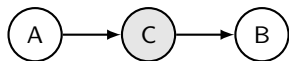


$$A \perp\!\!\!\perp B \mid C$$

Given  $C$ ,  $A$  and  $B$  are independent.

Example: Alarm  $\leftarrow$  Wind  $\rightarrow$  Barometer

The Chain



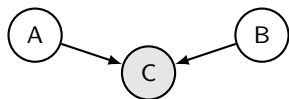
$$A \perp\!\!\!\perp B \mid C$$

After observing  $C$ ,

further observing  $A$  would not tell us about  $B$ .

Example: Burglary  $\rightarrow$  Alarm  $\rightarrow$  André wakes up

The Collider



**Surprisingly,**  $A \perp\!\!\!\perp B$   
but **not**  $A \perp\!\!\!\perp B \mid C$  !

Example: Burglary  $\rightarrow$  Alarm  $\leftarrow$  Wind

Burglaries occur regardless how windy it is.

If alarm rings, hearing wind makes burglary **less likely!**

Burglary is “explained away” by wind.

# Detecting independence: d-separation

**Definition:**  $A$  and  $B$  are **d-separated** given set  $C$  if **for any path**  $P$  from  $A$  to  $B$  **at least** one holds:

# Detecting independence: d-separation

**Definition:**  $A$  and  $B$  are **d-separated** given set  $C$  if **for any path**  $P$  from  $A$  to  $B$  **at least** one holds:

- 1  $P$  includes a fork with observed parent:

$$X \leftarrow Z \rightarrow Y \quad (\text{with } Z \in C)$$

# Detecting independence: d-separation

**Definition:**  $A$  and  $B$  are **d-separated** given set  $C$  if **for any path**  $P$  from  $A$  to  $B$  **at least** one holds:

- 1  $P$  includes a fork with observed parent:

$$X \leftarrow Z \rightarrow Y \quad (\text{with } Z \in C)$$

- 2  $P$  includes a chain with observed middle:

$$X \rightarrow Z \rightarrow Y \quad \text{or} \quad X \leftarrow Z \leftarrow Y \quad (\text{with } Z \in C)$$

# Detecting independence: d-separation

**Definition:**  $A$  and  $B$  are **d-separated** given set  $C$  if **for any path**  $P$  from  $A$  to  $B$  **at least** one holds:

- 1  $P$  includes a fork with observed parent:

$$X \leftarrow Z \rightarrow Y \quad (\text{with } Z \in C)$$

- 2  $P$  includes a chain with observed middle:

$$X \rightarrow Z \rightarrow Y \quad \text{or} \quad X \leftarrow Z \leftarrow Y \quad (\text{with } Z \in C)$$

- 3  $P$  includes a collider with unobserved descendants:

$$X \rightarrow Z \leftarrow Y \quad (\text{with neither } Z \text{ nor any of its descendants } \in C)$$



# Detecting independence: d-separation

**Definition:**  $A$  and  $B$  are **d-separated** given set  $C$  if **for any path**  $P$  from  $A$  to  $B$  **at least** one holds:

- 1  $P$  includes a fork with observed parent:

$$X \leftarrow Z \rightarrow Y \quad (\text{with } Z \in C)$$

- 2  $P$  includes a chain with observed middle:

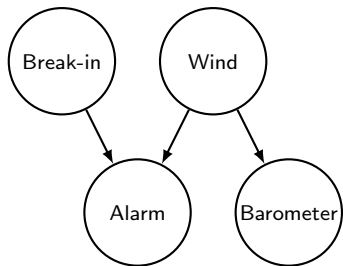
$$X \rightarrow Z \rightarrow Y \quad \text{or} \quad X \leftarrow Z \leftarrow Y \quad (\text{with } Z \in C)$$

- 3  $P$  includes a collider with unobserved descendants:

$$X \rightarrow Z \leftarrow Y \quad (\text{with neither } Z \text{ nor any of its descendants } \in C)$$

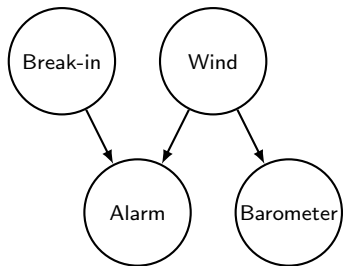
**Theorem:**  $A$  and  $B$  **d-separated** given  $C \implies A \perp\!\!\!\perp B \mid C$ .

# Examples



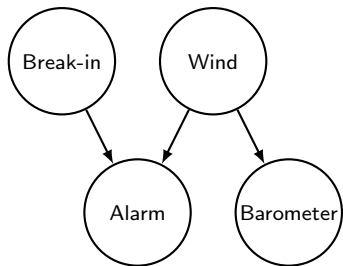
Wind  $\perp\!\!\!\perp$  Barometer?

# Examples



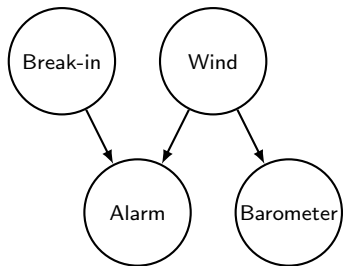
Wind  $\perp\!\!\!\perp$  Barometer? **No**

# Examples



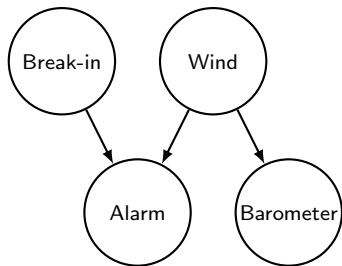
Wind  $\perp\!\!\!\perp$  Barometer? **No**  
Break-in  $\perp\!\!\!\perp$  Wind?

# Examples



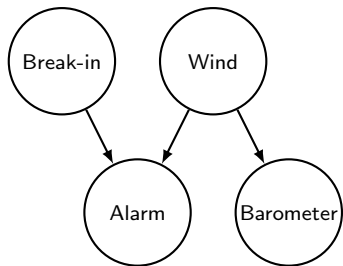
Wind  $\perp\!\!\!\perp$  Barometer? **No**  
Break-in  $\perp\!\!\!\perp$  Wind? **Yes**

# Examples



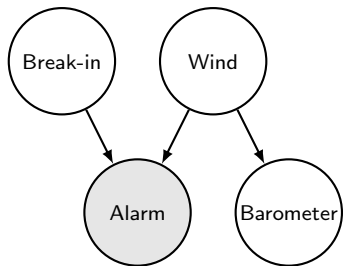
Wind  $\perp\!\!\!\perp$  Barometer? **No**  
Break-in  $\perp\!\!\!\perp$  Wind? **Yes**  
Break-in  $\perp\!\!\!\perp$  Barometer?

# Examples



Wind  $\perp\!\!\!\perp$  Barometer? **No**  
Break-in  $\perp\!\!\!\perp$  Wind? **Yes**  
Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**

# Examples



Wind  $\perp\!\!\!\perp$  Barometer? **No**

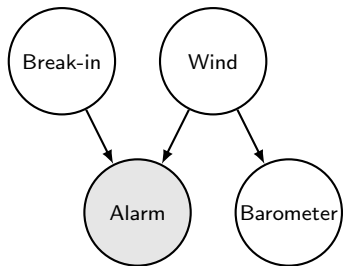
Break-in  $\perp\!\!\!\perp$  Wind? **Yes**

Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**

Break-in  $\perp\!\!\!\perp$  Barometer | Alarm?



# Examples



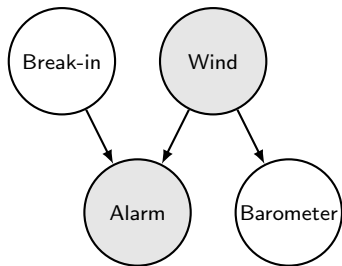
Wind  $\perp\!\!\!\perp$  Barometer? **No**

Break-in  $\perp\!\!\!\perp$  Wind? **Yes**

Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**

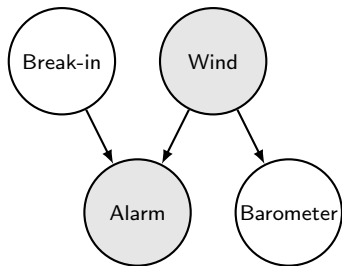
Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**

# Examples



- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind?

# Examples



Wind  $\perp\!\!\!\perp$  Barometer? **No**

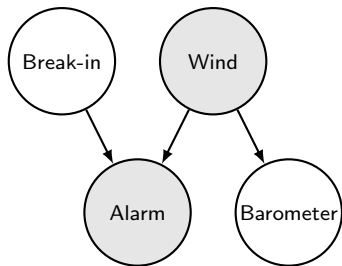
Break-in  $\perp\!\!\!\perp$  Wind? **Yes**

Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**

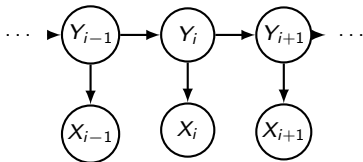
Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**

Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

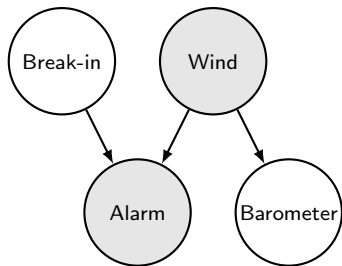
# Examples



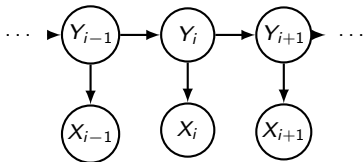
- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**



# Examples

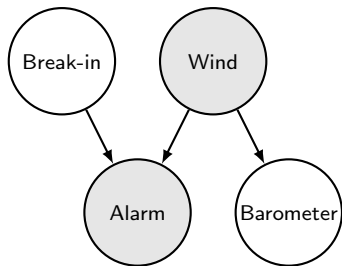


- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

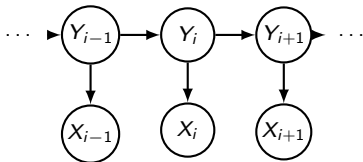


$$Y_{i+1} \perp\!\!\!\perp Y_{i-1}?$$

# Examples

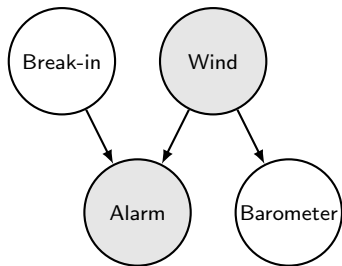


- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

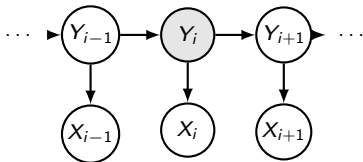


$Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**

# Examples

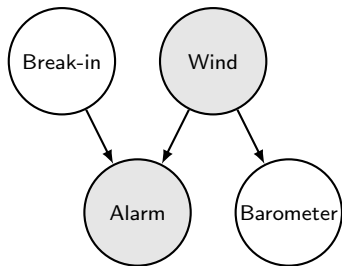


- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

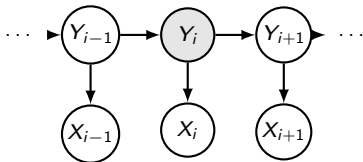


- $Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$ ?

# Examples



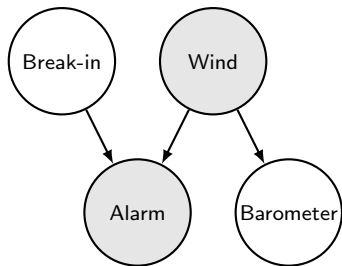
- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**



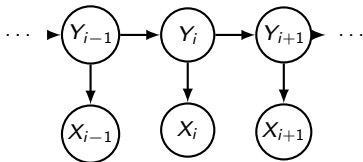
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$ ? **Yes**



# Examples

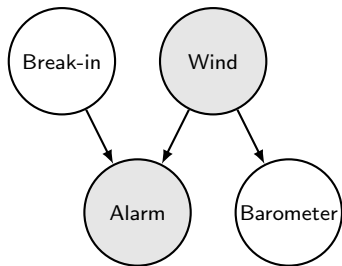


- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

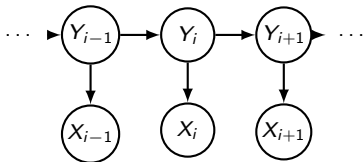


- $Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$ ? **Yes**
- $Y_{i+1} \perp\!\!\!\perp X_i$ ?

# Examples

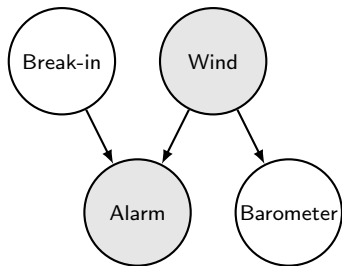


- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

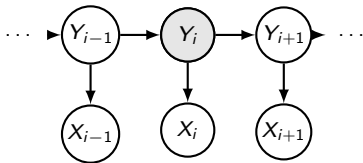


- $Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$ ? **Yes**
- $Y_{i+1} \perp\!\!\!\perp X_i$ ? **No**

# Examples

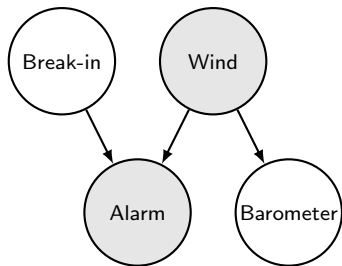


- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**

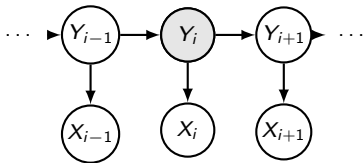


- $Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$ ? **Yes**
- $Y_{i+1} \perp\!\!\!\perp X_i$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp X_i \mid Y_i$ ?

# Examples



- Wind  $\perp\!\!\!\perp$  Barometer? **No**
- Break-in  $\perp\!\!\!\perp$  Wind? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer? **Yes**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm? **No**
- Break-in  $\perp\!\!\!\perp$  Barometer | Alarm, Wind? **Yes**



- $Y_{i+1} \perp\!\!\!\perp Y_{i-1}$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$ ? **Yes**
- $Y_{i+1} \perp\!\!\!\perp X_i$ ? **No**
- $Y_{i+1} \perp\!\!\!\perp X_i \mid Y_i$ ? **Yes**

# Generative stories and plate notation

In papers, you'll see statistical models defined through *generative stories*:

$$\mu \sim \text{Uniform}([-1, 1])$$

$$\sigma \sim \text{Uniform}([1, 2])$$

$$X \mid \mu, \sigma \sim \text{Normal}(\mu, \sigma)$$

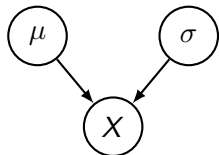
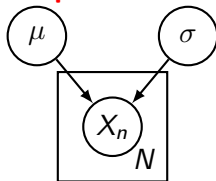


Plate notation is a way to denote **repetition of templates**:

$$\mu \sim \text{Uniform}([-1, 1])$$

$$\sigma \sim \text{Uniform}([1, 2])$$

$$X_n \mid \mu, \sigma \sim \text{Normal}(\mu, \sigma) \quad i = 1, \dots, N$$



# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

Markov random fields

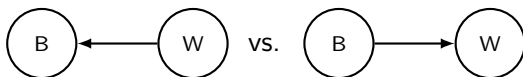
Factor graphs

Correlation does not imply causation;  
but then, *what does?*

# Seeing versus doing

Bayes nets only model independence assumptions.

The correlation between the a barometer reading  $B$  and wind strength  $W$  can be represented either way:

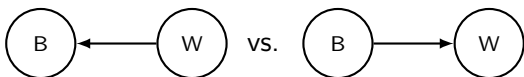




# Seeing versus doing

Bayes nets only model independence assumptions.

The correlation between the a barometer reading  $B$  and wind strength  $W$  can be represented either way:



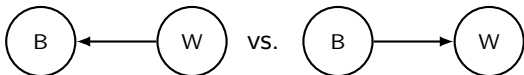
**Seeing** that the barometer reading is high, we can forecast wind.

| $\mathbb{P}(W   B)$ | lo  | mid | hi  |
|---------------------|-----|-----|-----|
| $B = \text{lo}$     | .98 | .01 | .01 |
| $B = \text{mid}$    | .01 | .98 | .01 |
| $B = \text{hi}$     | .01 | .01 | .98 |

# Seeing versus doing

Bayes nets only model independence assumptions.

The correlation between the a barometer reading  $B$  and wind strength  $W$  can be represented either way:



**Seeing** that the barometer reading is high, we can forecast wind.

| $\mathbb{P}(W   B)$ | lo  | mid | hi  |
|---------------------|-----|-----|-----|
| $B = \text{lo}$     | .98 | .01 | .01 |
| $B = \text{mid}$    | .01 | .98 | .01 |
| $B = \text{hi}$     | .01 | .01 | .98 |

But **setting** the barometer needle to high manually **won't cause wind!**

We write:  $\mathbb{P}(W | \text{do}(B = \text{hi})) = ?$

# Seeing versus doing

**Setting** the barometer needle to high manually **won't cause wind!**

# Seeing versus doing

**Setting** the barometer needle to high manually **won't cause wind!**

Two reasons why doing  $\neq$  seeing:

- the direction does not express a causal relationship
- we missed some confounding factor

If we created wind with a ceiling fan, does it alter the barometer?

# Seeing versus doing

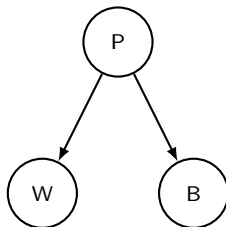
**Setting** the barometer needle to high manually **won't cause wind!**

Two reasons why doing  $\neq$  seeing:

- the direction does not express a causal relationship
- we missed some confounding factor

If we created wind with a ceiling fan, does it alter the barometer?

No! **Pressure** is a confounding factor.



# Causal models

## Definition (Pearl 2000)

A causal model is a DAG  $\mathcal{G}$  with vertices  $X_1, \dots, X_N$  representing events. Almost like a BN. However, paths are **causal**.

- $A$  causes  $B$  only if  $A$  is an ancestor of  $B$  in  $\mathcal{G}$ .
- $A \rightarrow B$  means  $A$  is a direct cause of  $B$ .

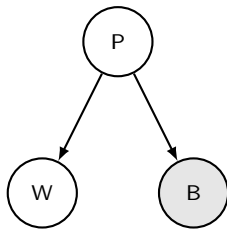
A good model is essential.

Wrong causal assumptions  $\Rightarrow$  wrong conclusions.

(We won't cover how to assess if the model is right. This is a bit *chicken-and-egg*, but domain knowledge helps, and we are allowed to reason about *unobserved* causes.)

# Seeing versus doing, more rigorously

**Seeing** (*observational*):  $\mathbb{P}(W \mid B = \text{hi})$



# Seeing versus doing, more rigorously

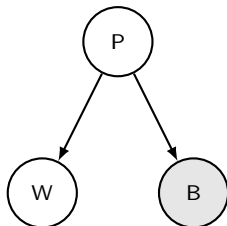
**Seeing** (*observational*):  $\mathbb{P}(W \mid B = \text{hi})$

Measure the world for a while (or call IPMA)

| Date       | Pressure | Wind | Barometer |
|------------|----------|------|-----------|
| 1977-01-01 | hi       | hi   | hi        |
| 1977-01-02 | hi       | mid  | hi        |
| 1977-01-02 | mid      | mid  | mid       |
| ...        |          |      |           |
| 2019-11-03 | hi       | hi   | hi        |

gives:

| $\mathbb{P}(W \mid B)$ | lo  | mid | hi  |
|------------------------|-----|-----|-----|
| $B = \text{hi}$        | .01 | .01 | .98 |





# Seeing versus doing, more rigorously

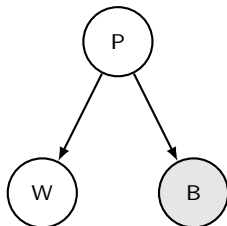
**Seeing** (*observational*):  $\mathbb{P}(W \mid B = \text{hi})$

Measure the world for a while (or call IPMA)

| Date       | Pressure | Wind | Barometer |
|------------|----------|------|-----------|
| 1977-01-01 | hi       | hi   | hi        |
| 1977-01-02 | hi       | mid  | hi        |
| 1977-01-02 | mid      | mid  | mid       |
| ...        |          |      |           |
| 2019-11-03 | hi       | hi   | hi        |

gives:

| $\mathbb{P}(W \mid B)$ | lo  | mid | hi  |
|------------------------|-----|-----|-----|
| $B = \text{hi}$        | .01 | .01 | .98 |



**Doing** (*interventional*):  $\mathbb{P}(W \mid \text{do}(B = \text{hi}))$

**Set** the needle to high, **breaking inbound arrows**;  
re-generate **new** data in this **new** DAG  
(or estimate what that would give.)

# Seeing versus doing, more rigorously

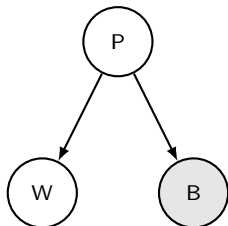
**Seeing** (*observational*):  $\mathbb{P}(W \mid B = \text{hi})$

Measure the world for a while (or call IPMA)

| Date       | Pressure | Wind | Barometer |
|------------|----------|------|-----------|
| 1977-01-01 | hi       | hi   | hi        |
| 1977-01-02 | hi       | mid  | hi        |
| 1977-01-02 | mid      | mid  | mid       |
| ...        |          |      |           |
| 2019-11-03 | hi       | hi   | hi        |

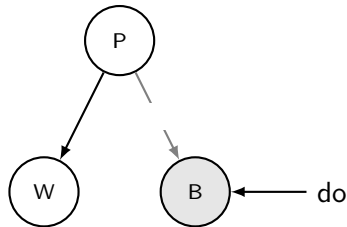
gives:

| $\mathbb{P}(W \mid B)$ | lo  | mid | hi  |
|------------------------|-----|-----|-----|
| $B = \text{hi}$        | .01 | .01 | .98 |



**Doing** (*interventional*):  $\mathbb{P}(W \mid \text{do}(B = \text{hi}))$

**Set** the needle to high, **breaking inbound arrows**;  
re-generate **new** data in this **new** DAG  
(or estimate what that would give.)



# Seeing versus doing, more rigorously

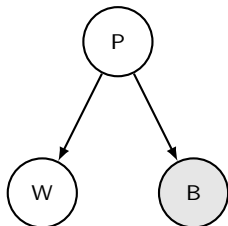
**Seeing** (*observational*):  $\mathbb{P}(W \mid B = \text{hi})$

Measure the world for a while (or call IPMA)

| Date       | Pressure | Wind | Barometer |
|------------|----------|------|-----------|
| 1977-01-01 | hi       | hi   | hi        |
| 1977-01-02 | hi       | mid  | hi        |
| 1977-01-02 | mid      | mid  | mid       |
| ...        |          |      |           |
| 2019-11-03 | hi       | hi   | hi        |

gives:

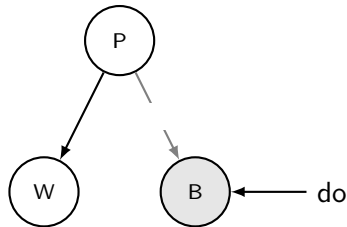
| $\mathbb{P}(W \mid B)$ | lo  | mid | hi  |
|------------------------|-----|-----|-----|
| $B = \text{hi}$        | .01 | .01 | .98 |



**Doing** (*interventional*):  $\mathbb{P}(W \mid \text{do}(B = \text{hi}))$

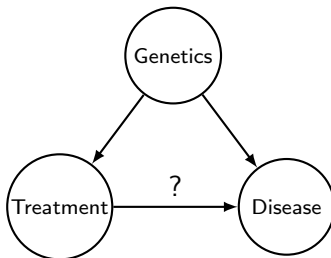
**Set** the needle to high, **breaking inbound arrows**;  
re-generate **new** data in this **new** DAG  
(or estimate what that would give.)

$$\mathbb{P}(W \mid \text{do}(B = \text{hi})) = \mathbb{P}(W)$$



# Randomized controlled trials

Try to actually implement the *do* operator in real life.

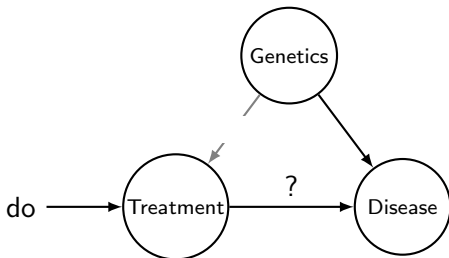


| Patient | Treatment | Genetics | Disease   |
|---------|-----------|----------|-----------|
| #42     | real      | ?        | cured     |
| #68     | placebo   | ?        | not cured |
| ...     |           |          |           |

No need to be able to measure genetics  
as long as we can sample A LOT OF test subjects with no/little bias.

# Randomized controlled trials

Try to actually implement the *do* operator in real life.



| Patient | Treatment | Genetics | Disease   |
|---------|-----------|----------|-----------|
| #42     | real      | ?        | cured     |
| #68     | placebo   | ?        | not cured |
| ...     |           |          |           |

No need to be able to measure genetics  
as long as we can sample A LOT OF test subjects with no/little bias.

RCTs are powerful, but often unethical, always expensive.

**Do-calculus**: use the **causal DAG assumptions** to draw causal conclusions from observational data.

- Apply transformations to  $\mathbb{P}(X \mid \text{do}(Y))$  until the “do” goes away.  
(Not always possible!)
- Quantities without “do” can be estimated observationally.
- Transformation: 3 rules.

# Pearl's 3 rules

**Notation:**  $X, Y, Z, W$  disjoint sets of events (sets of nodes); may be empty  
 $\mathcal{G}_{\bar{X}}$  the graph with all edges **into**  $X$  removed.  
 $\mathcal{G}_X$  the graph with all edges **out of**  $X$  removed.  
 $Z(X)$  subset of nodes in  $Z$  which are not ancestors of  $X$ .  
 $y; \text{do}(x)$  shorthand for  $Y = y$ ; respectively  $\text{do}(X = x)$ .

## 1 Ignoring observations:

$$\mathbb{P}(y \mid \text{do}(x), z, w) = \mathbb{P}(y \mid \text{do}(x), w) \quad \text{if} \quad (Y \perp\!\!\!\perp Z \mid X, W)_{\mathcal{G}_{\bar{X}}}$$

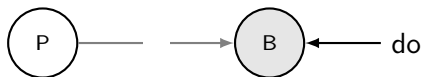
## 2 Action/observation exchange: the back-door criterion

$$\mathbb{P}(y \mid \text{do}(x), \text{do}(z), w) = \mathbb{P}(y \mid \text{do}(x), z, w) \quad \text{if} \quad (Y \perp\!\!\!\perp Z \mid X, W)_{\mathcal{G}_{\bar{X}, Z}}$$

## 3 Ignoring actions

$$\mathbb{P}(y \mid \text{do}(x), \text{do}(z), w) = \mathbb{P}(y \mid \text{do}(x), w) \quad \text{if} \quad (Y \perp\!\!\!\perp Z \mid X, W)_{\mathcal{G}_{\bar{X}, Z(\bar{w})}}$$

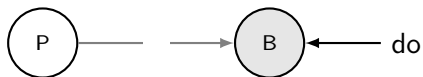
# Examples 1,2: Pressure and barometer



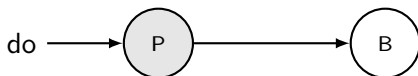
Rule 3:  $\mathbb{P}(P = hi \mid do(B = hi)) = \mathbb{P}(P = hi)$  since  $(P \perp\!\!\!\perp B)_{\mathcal{G}_B}$



# Examples 1,2: Pressure and barometer

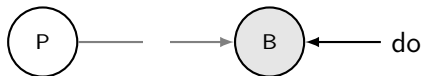


Rule 3:  $\mathbb{P}(P = \text{hi} \mid \text{do}(B = \text{hi})) = \mathbb{P}(P = \text{hi})$  since  $(P \perp\!\!\!\perp B)_{\mathcal{G}_B}$

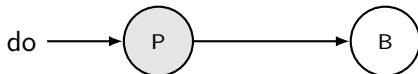


Rule 2:  $\mathbb{P}(B = \text{hi} \mid \text{do}(P = \text{lo})) = \mathbb{P}(B = \text{hi} \mid P = \text{lo})$  since  $(B \perp\!\!\!\perp P)_{\mathcal{G}_P}$

# Examples 1,2: Pressure and barometer



Rule 3:  $\mathbb{P}(P = \text{hi} \mid \text{do}(B = \text{hi})) = \mathbb{P}(P = \text{hi})$  since  $(P \perp\!\!\!\perp B)_{\mathcal{G}_B}$

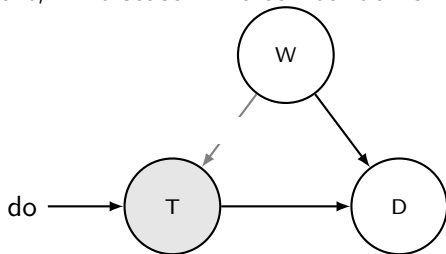


Rule 2:  $\mathbb{P}(B = \text{hi} \mid \text{do}(P = \text{lo})) = \mathbb{P}(B = \text{hi} \mid P = \text{lo})$  since  $(B \perp\!\!\!\perp P)_{\mathcal{G}_P}$

Good check: we get the intuitively correct results.

## Example 3: Measurable confounder

$T$ : treatment,  $D$ : disease. The confounder is  $W$ : wealth.



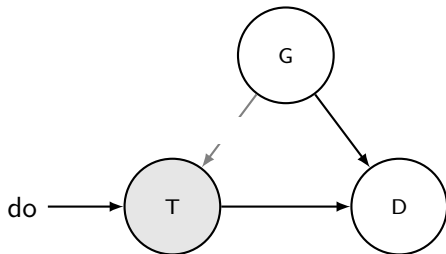
Condition on wealth (which thus needs to be measurable)

$$\begin{aligned}\mathbb{P}(D = \text{cured} \mid \text{do}(T = y)) &= \mathbb{P}(D = \text{cured} \mid \text{do}(T = y), W = y)\mathbb{P}(W = y \mid \text{do}(T = y)) \\ &\quad + \mathbb{P}(D = \text{cured} \mid \text{do}(T = y), W = n)\mathbb{P}(W = n \mid \text{do}(T = y)) \\ &= \mathbb{P}(D = \text{cured} \mid \text{do}(T = y), W = y)\mathbb{P}(W = y) \\ &\quad + \mathbb{P}(D = \text{cured} \mid \text{do}(T = y), W = n)\mathbb{P}(W = n) \quad (\text{R3}) \\ &= \mathbb{P}(D = \text{cured} \mid T = y, W = y)\mathbb{P}(W = y) \\ &\quad + \mathbb{P}(D = \text{cured} \mid T = y, W = n)\mathbb{P}(W = n) \quad (\text{R2})\end{aligned}$$

## Example 3: an impossible one

$T$ : treatment,  $D$ : disease.

The confounder is  $G$ : genetics (impractical to measure and estimate)

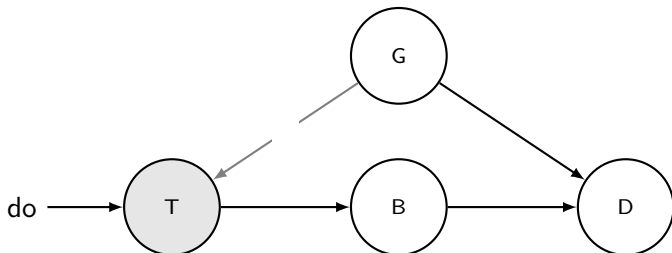


Without more info or more assumptions, we're stuck!

## Example 4: a surprisingly possible one

$T$ : treatment,  $D$ : disease,  $B$ : blood cell count.

The confounder is  $G$ : genetics (still hidden)



“The front-door criterion:” conditioning on  $B$  lets us remove dos!

(I won't show you how, derivation is a bit longer. Try it at home.)

$$\mathbb{P}(D = \text{cured} \mid \text{do}(T = y)) =$$

$$\sum_b \mathbb{P}(B = b \mid T = y) \sum_t \mathbb{P}(D = \text{cured} \mid T = t, B = b) \mathbb{P}(T = t)$$

# Directed models: summary

- Bayes nets: specify & estimate **fine-grained distributions** over **interdependent events**.
- Under a specified model, algorithm to decide conditional independence: **d-separation**
- Bestowing a DAG with **causal assumptions** lets us reason about **interventions**.

Further reading: (Pearl, 1988; Koller and Friedman, 2009; Pearl, 2000, 2012; Dawid, 2010)

Slides on causal inference and learning causal structure (links):

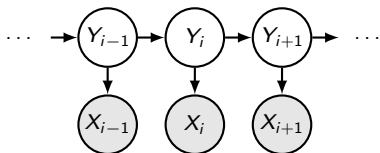
- Sanna Tyrväinen, Introduction to Causal Calculus
- Ricardo Silva, Causality
- Dominik Janzing & Bernhard Schölkopf, Causality

**Highly recommended online course:** <https://www.bradyneal.com/causal-inference-course>

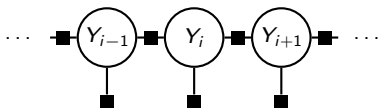
# Graphical Models

In this unit, we will formalize & extend these graphical representations encountered in previous lectures.

**Directed**



**Undirected**



# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

Markov random fields

Factor graphs



# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

Markov random fields

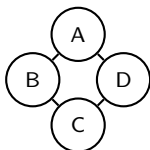
Factor graphs

# Modeling friendships

- Four students: An, Bo, Chris, Dee are voting on a Yes/No ballot.
- Friendship pairs: An–Bo, Bo–Chris, Chris–Dee, Dee–An.
- Friends are 100x more likely to vote the same way.

# Modeling friendships

- Four students: An, Bo, Chris, Dee are voting on a Yes/No ballot.
- Friendship pairs: An–Bo, Bo–Chris, Chris–Dee, Dee–An.
- Friends are 100x more likely to vote the same way.



- An's vote is a random variable  $A$  with values  $a \in \{Y, N\}$ , and so on.

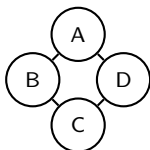
$$\mathbb{P}(a, b, c, d) \propto f(a, b) \cdot f(b, c) \cdot f(c, d) \cdot f(d, a)$$

For any  $X, Y \in \{A, B, C, D\}$ ,  $f$  is the **compatibility function**:

$$f(x, y) = \begin{cases} 100 & \text{if } x = y = \text{Yes or } x = y = \text{No} \\ 1 & \text{otherwise.} \end{cases}$$

# Modeling friendships

- Four students: An, Bo, Chris, Dee are voting on a Yes/No ballot.
- Friendship pairs: An–Bo, Bo–Chris, Chris–Dee, Dee–An.
- Friends are 100x more likely to vote the same way.



- An's vote is a random variable  $A$  with values  $a \in \{Y, N\}$ , and so on.

$$\mathbb{P}(a, b, c, d) \propto f(a, b) \cdot f(b, c) \cdot f(c, d) \cdot f(d, a)$$

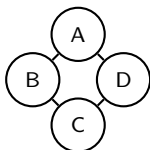
For any  $X, Y \in \{A, B, C, D\}$ ,  $f$  is the **compatibility function**:

$$f(x, y) = \begin{cases} 100 & \text{if } x = y = \text{Yes or } x = y = \text{No} \\ 1 & \text{otherwise.} \end{cases}$$

- Can we represent this exact factorization in a Bayes net?

# Modeling friendships

- Four students: An, Bo, Chris, Dee are voting on a Yes/No ballot.
- Friendship pairs: An–Bo, Bo–Chris, Chris–Dee, Dee–An.
- Friends are 100x more likely to vote the same way.



- An's vote is a random variable  $A$  with values  $a \in \{Y, N\}$ , and so on.

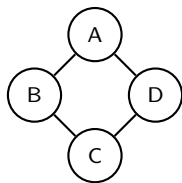
$$\mathbb{P}(a, b, c, d) \propto f(a, b) \cdot f(b, c) \cdot f(c, d) \cdot f(d, a)$$

For any  $X, Y \in \{A, B, C, D\}$ ,  $f$  is the **compatibility function**:

$$f(x, y) = \begin{cases} 100 & \text{if } x = y = \text{Yes or } x = y = \text{No} \\ 1 & \text{otherwise.} \end{cases}$$

- Can we represent this exact factorization in a Bayes net? **No!**

# Markov random fields



## Definition

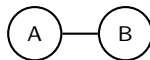
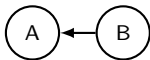
Let  $\mathcal{G}$  be an *undirected* graph with nodes corresponding to random variables  $X_1, \dots, X_N$ . Let  $C(\mathcal{G})$  denote the set of *cliques* (fully connected subgraphs) of  $\mathcal{G}$ . A MRF is a distribution of the form

$$\mathbb{P}(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} f_c(\mathbf{x}_c)$$

where for each clique  $c$ ,  $f_c$  is a non-negative compatibility function.

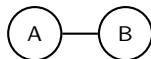
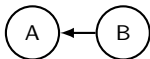
# Any BN can be encoded in a MRF

- Convert all arcs  $A \rightarrow B$  or  $A \leftarrow B$  into undirected edges  $A - B$ .



# Any BN can be encoded in a MRF

- 2 Convert all arcs  $A \rightarrow B$  or  $A \leftarrow B$  into undirected edges  $A - B$ .



| A | B | $\mathbb{P}(a   b)$ |
|---|---|---------------------|
| Y | Y | .9                  |
| N | Y | .1                  |
| Y | N | .1                  |
| N | N | .9                  |

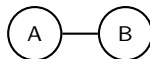
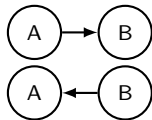
| B | $\mathbb{P}(b)$ |
|---|-----------------|
| Y | .75             |
| N | .25             |

| A | B | $f(a, b)$ |
|---|---|-----------|
| Y | Y | .9 · .75  |
| N | Y | .1 · .75  |
| Y | N | .1 · .25  |
| N | N | .9 · .25  |



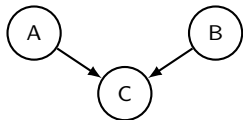
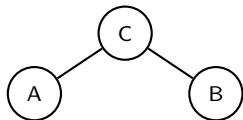
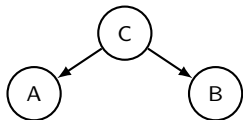
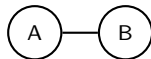
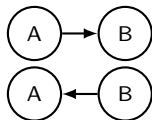
# Any BN can be encoded in a MRF

- 2 Convert all arcs  $A \rightarrow B$  or  $A \leftarrow B$  into undirected edges  $A - B$ .



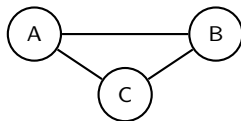
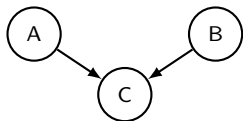
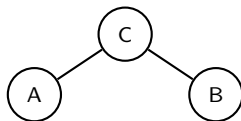
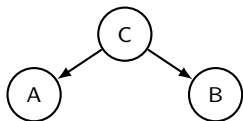
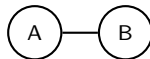
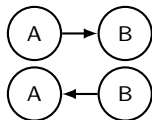
# Any BN can be encoded in a MRF

- 2 Convert all arcs  $A \rightarrow B$  or  $A \leftarrow B$  into undirected edges  $A - B$ .



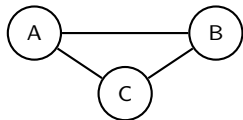
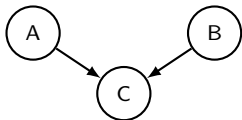
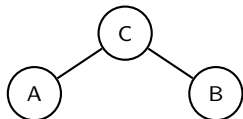
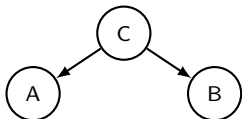
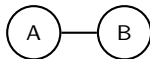
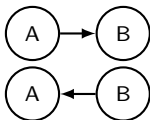
# Any BN can be encoded in a MRF

- 2 Convert all arcs  $A \rightarrow B$  or  $A \leftarrow B$  into undirected edges  $A - B$ .



# Any BN can be encoded in a MRF

- 1 First, add edge  $A - C$  for any collider structure  $A \rightarrow B \leftarrow C$ ;
- 2 Convert all arcs  $A \rightarrow B$  or  $A \leftarrow B$  into undirected edges  $A - B$ .

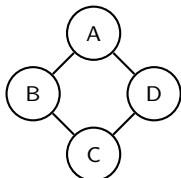


# Loose conversion

Similarly, we can convert a MRF to a BN (we won't cover it.)

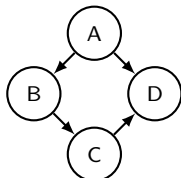
However, **independences may be lost** in either direction.

**From**

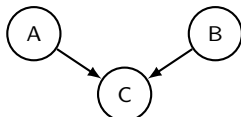


$$A \perp\!\!\!\perp C \mid B, D$$
$$B \perp\!\!\!\perp D \mid A, C$$

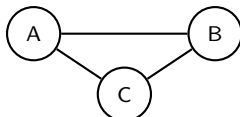
**To**



~~$$A \perp\!\!\!\perp C \mid B, D$$~~
$$B \perp\!\!\!\perp D \mid A, C$$



$$A \perp\!\!\!\perp B$$

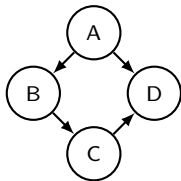


~~$$A \perp\!\!\!\perp B$$~~

# Bayes vs Markov

## Bayes network

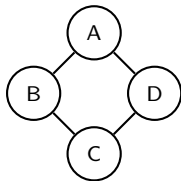
- Factors are conditionals (normalized)
- Easy to sample
- Can be made causal
- Can easily find  $\mathbb{P}(x_1, \dots, x_n)$ .



$$\mathbb{P}(a, b, c, d) = \mathbb{P}(a)\mathbb{P}(b | a)\mathbb{P}(c | b)\mathbb{P}(d | a, c)$$

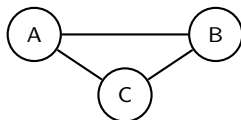
## Markov networks

- Factors are cliques (unnormalized)
- No directional ambiguity
- Often more compact
- More symmetric notation

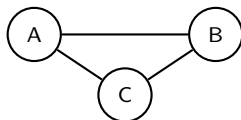


$$\mathbb{P}(a, b, c, d) = 1/Z f_1(a, b)f_2(b, c)f_3(c, d)f_4(d, a)$$

# What are the factors in a MRF?



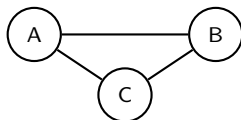
# What are the factors in a MRF?



Single clique:  $\{A, B, C\}$ , so  $\mathbb{P}(a, b, c) = \frac{1}{Z} f(a, b, c)$ .



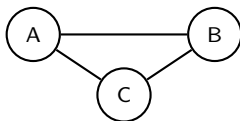
# What are the factors in a MRF?



Single clique:  $\{A, B, C\}$ , so  $\mathbb{P}(a, b, c) = \frac{1}{Z} f(a, b, c)$ .

No way to represent  $\mathbb{P}(a, b, c) = 1/Z f_1(a, b) f_2(b, c) f_3(c, a)$ .

# What are the factors in a MRF?

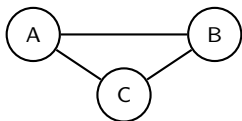


Single clique:  $\{A, B, C\}$ , so  $\mathbb{P}(a, b, c) = \frac{1}{Z} f(a, b, c)$ .

No way to represent  $\mathbb{P}(a, b, c) = 1/Z f_1(a, b) f_2(b, c) f_3(c, a)$ .

**Pairwise MRF:** Like a MRF, but factors are edges rather than cliques.

# What are the factors in a MRF?

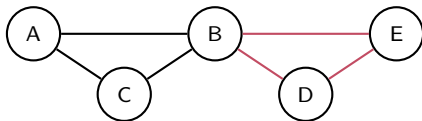


Single clique:  $\{A, B, C\}$ , so  $\mathbb{P}(a, b, c) = \frac{1}{Z} f(a, b, c)$ .

No way to represent  $\mathbb{P}(a, b, c) = 1/Z f_1(a, b) f_2(b, c) f_3(c, a)$ .

**Pairwise MRF:** Like a MRF, but factors are edges rather than cliques.

But what if we want to mix them?



$$\mathbb{P}(a, b, c, d, e) = 1/Z f_1(a, b) f_2(b, c) f_3(c, a) f_4(b, d, e)$$

# Outline

## ① Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

## ② Undirected Models

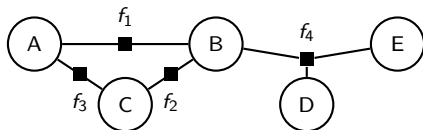
Markov random fields

Factor graphs

# Factor graphs

Explicitly represent factors in the graph to remove ambiguity.

$$\mathbb{P}(a, b, c, d, e) = 1/Z f_1(a, b)f_2(b, c)f_3(c, a)f_4(b, d, e)$$



## Definition (Factor graph)

A FG is a bipartite graph  $\mathcal{G}$  with vertices in  $\mathcal{V} \cup \mathcal{F}$ , where  $X_1, \dots, X_n \in \mathcal{V}$  are random variables and  $\alpha \in \mathcal{F}$  are factors, inducing a distribution

$$\mathbb{P}(x_1, \dots, x_n) = \frac{1}{Z} \prod_{\alpha \in \mathcal{F}} f_{\alpha}(\mathbf{x}_{\alpha})$$

where  $f_{\alpha} \geq 0$ , and  $\mathbf{X}_{\alpha}$  is the set of variables with an edge to factor  $\alpha$ .

# Factor graphs

- Any MRF can be mapped exactly to a FG (clique  $\rightarrow$  factor).
- Any Pairwise MRF can be mapped exactly to a FG (edge  $\rightarrow$  factor).
- FGs are more general / more *fine-grained*.

# Algorithms

- **Inference:** Given a FG with compatibility functions, answer **queries**
  - Maximization: Find most likely assignment  $x_1, \dots, x_N$  (possibly given evidence  $x_i : i \in \mathcal{E}$ ).

$$\operatorname{argmax}_{x_1, \dots, x_N} \mathbb{P}(x_1, \dots, x_N \mid \mathbf{x}_{\mathcal{E}})$$

- Marginalization: Find the marginal probability of some partial assignment over  $x_j : j \in \mathcal{M}$  (possibly given evidence  $x_i : i \in \mathcal{E}$ )

$$\mathbb{P}(\mathbf{x}_{\mathcal{M}} \mid \mathbf{x}_{\mathcal{E}})$$

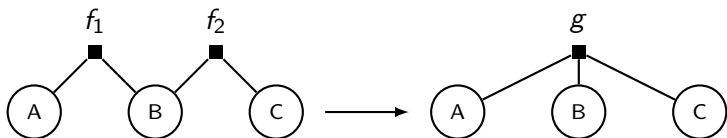
- **NP-hard** / **#P-hard** in general, but doable for tree-shaped graphs with dynamic programming.
- **Learning:** Given a dataset, estimate the compatibility tables (or, in general a model that produces them.)

Since  $\text{BN} \rightarrow \text{MRF} \rightarrow \text{FG}$ , it suffices to study inference algorithms for  $\text{FG}$ .<sup>1</sup>

<sup>1</sup>But not learning, since we cannot map back to BN losslessly!

# Multiplying factors

A core operation: combining factors by multiplying them.



| A | B | $f_1(a, b)$ |
|---|---|-------------|
| 0 | 0 | 3           |
| 0 | 1 | 1           |
| 1 | 0 | 2           |
| 1 | 1 | 8           |

| B | C | $f_2(a, b)$ |
|---|---|-------------|
| 0 | 0 | 5           |
| 0 | 1 | 4           |
| 1 | 0 | 1           |
| 1 | 1 | 1           |

→

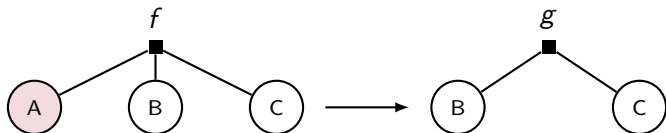
| A | B | C | $g(a, b, c)$     |
|---|---|---|------------------|
| 0 | 0 | 0 | $3 \cdot 5 = 15$ |
| 0 | 0 | 1 | $3 \cdot 4 = 12$ |
| 0 | 1 | 0 | $1 \cdot 1 = 1$  |
| 0 | 1 | 1 | $1 \cdot 1 = 1$  |
| 1 | 0 | 0 | $2 \cdot 5 = 10$ |
| 1 | 0 | 1 | $2 \cdot 4 = 8$  |
| 1 | 1 | 0 | $8 \cdot 1 = 8$  |
| 1 | 1 | 1 | $8 \cdot 1 = 8$  |

Distribution is preserved:

$$f_1(a, b) \cdot f_2(b, c) \cdot f_3(\dots) \cdot \dots = g(a, b, c) \cdot f_3(\dots) \cdot \dots$$



# Maximizing over a variable



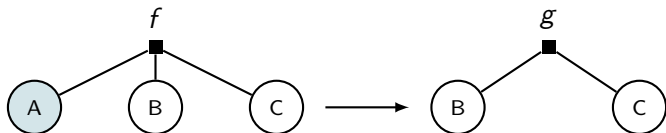
| A | B | C | $f(a, b, c)$ |
|---|---|---|--------------|
| 0 | 0 | 0 | 15           |
| 0 | 0 | 1 | 12           |
| 0 | 1 | 0 | 1            |
| 0 | 1 | 1 | 1            |
| 1 | 0 | 0 | 10           |
| 1 | 0 | 1 | 8            |
| 1 | 1 | 0 | 8            |
| 1 | 1 | 1 | 8            |

— maximizing over  $A \rightarrow$

| B | C | $g(b, c)$ |
|---|---|-----------|
| 0 | 0 | 15        |
| 0 | 1 | 12        |
| 1 | 0 | 8         |
| 1 | 1 | 8         |

$$\max_a f(a, b, c) \cdot \underbrace{f_4(\dots) \cdot \dots}_{A\text{-free}} = g(b, c) \cdot f_4(\dots) \cdot \dots$$

# Marginalizing over a variable



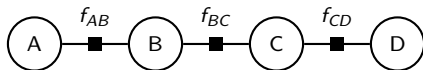
| A | B | C | $f(a, b, c)$ |
|---|---|---|--------------|
| 0 | 0 | 0 | 15           |
| 0 | 0 | 1 | 12           |
| 0 | 1 | 0 | 1            |
| 0 | 1 | 1 | 1            |
| 1 | 0 | 0 | 10           |
| 1 | 0 | 1 | 8            |
| 1 | 1 | 0 | 8            |
| 1 | 1 | 1 | 8            |

— summing over  $A \rightarrow$

| B | C | $g(b, c)$ |
|---|---|-----------|
| 0 | 0 | 25        |
| 0 | 1 | 20        |
| 1 | 0 | 9         |
| 1 | 1 | 9         |

$$\sum_a f(a, b, c) \cdot \underbrace{f_4(\dots)}_{A\text{-free}} \cdot \dots = g(b, c) \cdot f_4(\dots) \cdot \dots$$

# Variable elimination



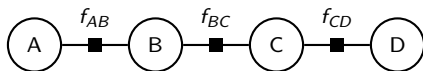
Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

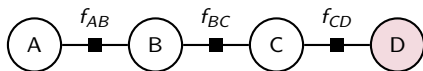
① Pick order: D, C, B, A

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

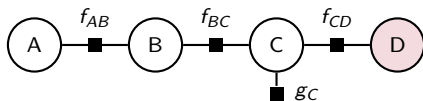
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

- 1 Pick order: D, C, B, A
- 2 Maximize over  $D$  ( $f_{CD} \rightarrow g_C$ )

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 Maximize over  $D$  ( $f_{CD} \rightarrow g_C$ )

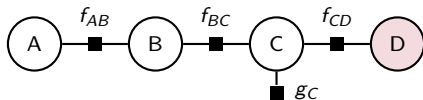
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 Maximize over  $D$  ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$

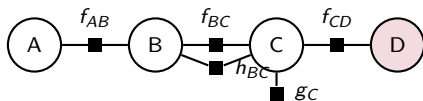
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 Maximize over  $D$  ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

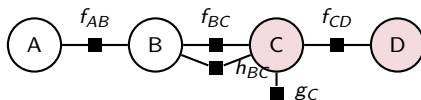
| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |



# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

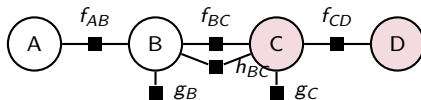
| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

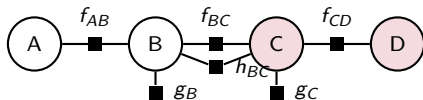
| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| B | $g_B(b)$  |
|---|-----------|
| 0 | $9^{C=1}$ |
| 1 | $6^{C=1}$ |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 Maximize over  $D$  ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over  $C$  ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

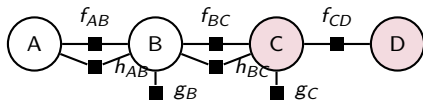
| B | $g_B(b)$  |
|---|-----------|
| 0 | $9^{C=1}$ |
| 1 | $6^{C=1}$ |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| B | $g_B(b)$  |
|---|-----------|
| 0 | $9^{C=1}$ |
| 1 | $6^{C=1}$ |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

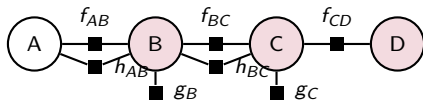
| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

| A | B | $h_{AB}(a, b)$          |
|---|---|-------------------------|
| 0 | 0 | $10 \cdot 9 = 90^{C=1}$ |
| 0 | 1 | $2 \cdot 6 = 12^{C=1}$  |
| 1 | 0 | $3 \cdot 9 = 27^{C=1}$  |
| 1 | 1 | $9 \cdot 6 = 54^{C=1}$  |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

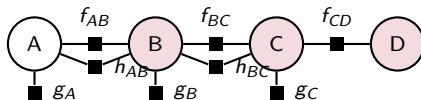
| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| A | B | $h_{AB}(a, b)$          |
|---|---|-------------------------|
| 0 | 0 | $10 \cdot 9 = 90^{C=1}$ |
| 0 | 1 | $2 \cdot 6 = 12^{C=1}$  |
| 1 | 0 | $3 \cdot 9 = 27^{C=1}$  |
| 1 | 1 | $9 \cdot 6 = 54^{C=1}$  |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Maximize over B ( $h_{AB} \rightarrow g_A$ )

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| A | $g_A(a)$   |
|---|------------|
| 0 | $90^{B=0}$ |
| 1 | $54^{B=1}$ |

| B | $g_B(b)$  |
|---|-----------|
| 0 | $9^{C=1}$ |
| 1 | $6^{C=1}$ |

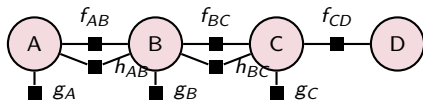
| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

| A | B | $h_{AB}(a, b)$          |
|---|---|-------------------------|
| 0 | 0 | $10 \cdot 9 = 90^{C=1}$ |
| 0 | 1 | $2 \cdot 6 = 12^{C=1}$  |
| 1 | 0 | $3 \cdot 9 = 27^{C=1}$  |
| 1 | 1 | $9 \cdot 6 = 54^{C=1}$  |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Maximize over B ( $h_{AB} \rightarrow g_A$ )

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ | A | $g_A(a)$   |
|---|---|----------------|---|------------|
| 0 | 0 | 10             | 0 | $90^{B=0}$ |
| 0 | 1 | 2              | 1 | $54^{B=1}$ |
| 1 | 0 | 3              |   |            |
| 1 | 1 | 9              |   |            |

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$  |
|---|---|----------------|---|-----------|
| 0 | 0 | 1              | 0 | $9^{C=1}$ |
| 0 | 1 | 3              | 1 | $6^{C=1}$ |
| 1 | 0 | 1              |   |           |
| 1 | 1 | 2              |   |           |

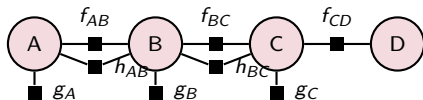
| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$  |
|---|---|----------------|---|-----------|
| 0 | 0 | 4              | 0 | $4^{D=0}$ |
| 0 | 1 | 2              | 1 | $3^{D=1}$ |
| 1 | 0 | 1              |   |           |
| 1 | 1 | 3              |   |           |

| A | B | $h_{AB}(a, b)$          |
|---|---|-------------------------|
| 0 | 0 | $10 \cdot 9 = 90^{C=1}$ |
| 0 | 1 | $2 \cdot 6 = 12^{C=1}$  |
| 1 | 0 | $3 \cdot 9 = 27^{C=1}$  |
| 1 | 1 | $9 \cdot 6 = 54^{C=1}$  |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Maximize over B ( $h_{AB} \rightarrow g_A$ )
- 7 Maximize over A ( $g_A \rightarrow \emptyset$ )

# Variable elimination



Query:  $\max_{a,b,c,d} \mathbb{P}(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| A | $g_A(a)$   |
|---|------------|
| 0 | $90^{B=0}$ |
| 1 | $54^{B=1}$ |

| B | $g_B(b)$  |
|---|-----------|
| 0 | $9^{C=1}$ |
| 1 | $6^{C=1}$ |

| C | $g_C(c)$  |
|---|-----------|
| 0 | $4^{D=0}$ |
| 1 | $3^{D=1}$ |

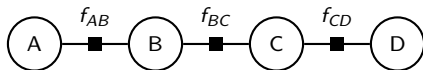
| A | B | $h_{AB}(a, b)$          |
|---|---|-------------------------|
| 0 | 0 | $10 \cdot 9 = 90^{C=1}$ |
| 0 | 1 | $2 \cdot 6 = 12^{C=1}$  |
| 1 | 0 | $3 \cdot 9 = 27^{C=1}$  |
| 1 | 1 | $9 \cdot 6 = 54^{C=1}$  |

| B | C | $h_{BC}(b, c)$        |
|---|---|-----------------------|
| 0 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 0 | 1 | $3 \cdot 3 = 9^{D=1}$ |
| 1 | 0 | $1 \cdot 4 = 4^{D=0}$ |
| 1 | 1 | $2 \cdot 3 = 6^{D=1}$ |

- 1 Pick order: D, C, B, A
- 2 Maximize over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Maximize over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Maximize over B ( $h_{AB} \rightarrow g_A$ )
- 7 Maximize over A ( $g_A \rightarrow \emptyset$ )
- 8 Just like Viterbi!  
The max is  $90/z$ .



# Variable elimination: sum



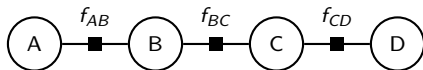
Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

# Variable elimination: sum



$$\text{Query: } Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$$

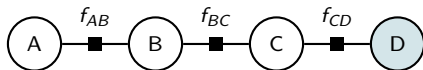
- 1 Pick order: D, C, B, A

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

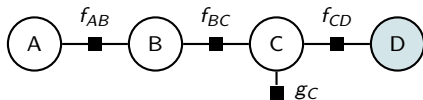
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

- 1 Pick order: D, C, B, A
- 2 **Sum** over D ( $f_{CD} \rightarrow g_C$ )

# Variable elimination: sum



$$\text{Query: } Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$$

- 1 Pick order: D, C, B, A
- 2 **Sum** over D ( $f_{CD} \rightarrow g_C$ )

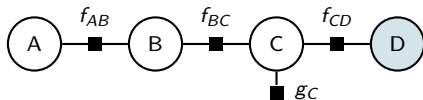
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 6        |
| 1 | 4        |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

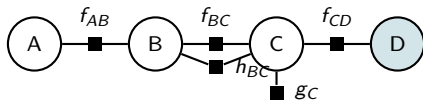
- 1 Pick order: D, C, B, A
- 2 **Sum** over  $D$  ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              |   |          |
| 0 | 1 | 2              | 0 | 6        |
| 1 | 0 | 1              | 1 | 4        |
| 1 | 1 | 3              |   |          |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 **Sum** over  $D$  ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

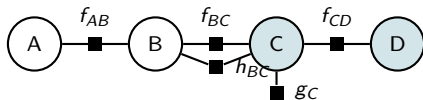
| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 6        |
| 1 | 4        |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 **Sum** over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 **Sum** over C ( $h_{BC} \rightarrow g_B$ )

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

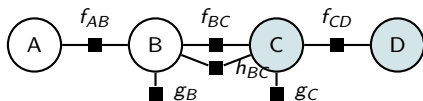
| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 6        |
| 1 | 4        |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 **Sum** over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 **Sum** over C ( $h_{BC} \rightarrow g_B$ )

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

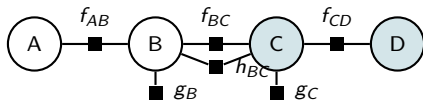
| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 6        |
| 0 | 1 | 2              | 1 | 4        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |



# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

- 1 Pick order: D, C, B, A
- 2 **Sum** over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 **Sum** over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$

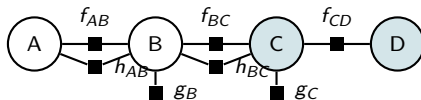
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              |   |          |
| 0 | 1 | 2              | 0 | 6        |
| 1 | 0 | 1              | 1 | 4        |
| 1 | 1 | 3              |   |          |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

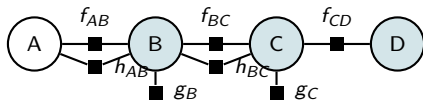
| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 6        |
| 0 | 1 | 2              | 1 | 4        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |

| A | B | $h_{AB}(a, b)$      |
|---|---|---------------------|
| 0 | 0 | $10 \cdot 18 = 180$ |
| 0 | 1 | $2 \cdot 14 = 28$   |
| 1 | 0 | $3 \cdot 18 = 54$   |
| 1 | 1 | $9 \cdot 14 = 126$  |

| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

- 1 Pick order: D, C, B, A
- 2 Sum over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Sum over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

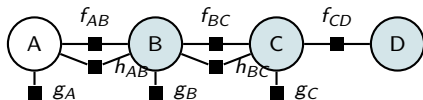
| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 6        |
| 0 | 1 | 2              | 1 | 4        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |

| A | B | $h_{AB}(a, b)$      |
|---|---|---------------------|
| 0 | 0 | $10 \cdot 18 = 180$ |
| 0 | 1 | $2 \cdot 14 = 28$   |
| 1 | 0 | $3 \cdot 18 = 54$   |
| 1 | 1 | $9 \cdot 14 = 126$  |

| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

- 1 Pick order: D, C, B, A
- 2 **Sum** over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 **Sum** over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 **Sum** over B ( $h_{AB} \rightarrow g_A$ )

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ | A | $g_A(a)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 10             | 0 | 208      |
| 0 | 1 | 2              | 1 | 180      |
| 1 | 0 | 3              |   |          |
| 1 | 1 | 9              |   |          |

| A | B | $h_{AB}(a, b)$      |
|---|---|---------------------|
| 0 | 0 | $10 \cdot 18 = 180$ |
| 0 | 1 | $2 \cdot 14 = 28$   |
| 1 | 0 | $3 \cdot 18 = 54$   |
| 1 | 1 | $9 \cdot 14 = 126$  |

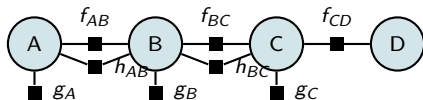
- 1 Pick order: D, C, B, A
- 2 Sum over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Sum over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Sum over B ( $h_{AB} \rightarrow g_A$ )

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 6        |
| 0 | 1 | 2              | 1 | 4        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ | A | $g_A(a)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 10             | 0 | 208      |
| 0 | 1 | 2              | 1 | 180      |
| 1 | 0 | 3              |   |          |
| 1 | 1 | 9              |   |          |

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

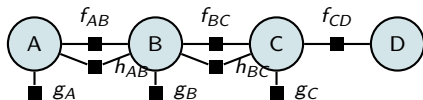
| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 6        |
| 0 | 1 | 2              | 1 | 4        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |

| A | B | $h_{AB}(a, b)$      |
|---|---|---------------------|
| 0 | 0 | $10 \cdot 18 = 180$ |
| 0 | 1 | $2 \cdot 14 = 28$   |
| 1 | 0 | $3 \cdot 18 = 54$   |
| 1 | 1 | $9 \cdot 14 = 126$  |

| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

- 1 Pick order: D, C, B, A
- 2 Sum over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Sum over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Sum over B ( $h_{AB} \rightarrow g_A$ )
- 7 Sum over A ( $g_A \rightarrow \emptyset$ )

# Variable elimination: sum



Query:  $Z = \sum_{a,b,c,d} f(a, b, c, d) = ?$

| A | B | $f_{AB}(a, b)$ | A | $g_A(a)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 10             | 0 | 208      |
| 0 | 1 | 2              | 1 | 180      |
| 1 | 0 | 3              |   |          |
| 1 | 1 | 9              |   |          |

| B | C | $f_{BC}(b, c)$ | B | $g_B(b)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 18       |
| 0 | 1 | 3              | 1 | 14       |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

| C | D | $f_{CD}(c, d)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 6        |
| 0 | 1 | 2              | 1 | 4        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |

| A | B | $h_{AB}(a, b)$      |
|---|---|---------------------|
| 0 | 0 | $10 \cdot 18 = 180$ |
| 0 | 1 | $2 \cdot 14 = 28$   |
| 1 | 0 | $3 \cdot 18 = 54$   |
| 1 | 1 | $9 \cdot 14 = 126$  |

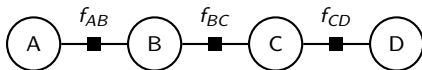
| B | C | $h_{BC}(b, c)$   |
|---|---|------------------|
| 0 | 0 | $1 \cdot 6 = 6$  |
| 0 | 1 | $3 \cdot 4 = 12$ |
| 1 | 0 | $1 \cdot 6 = 6$  |
| 1 | 1 | $2 \cdot 4 = 8$  |

- 1 Pick order: D, C, B, A
- 2 Sum over D ( $f_{CD} \rightarrow g_C$ )
- 3 Multiply  $f_{BC}$  with  $g_C$  giving  $h_{BC}$
- 4 Sum over C ( $h_{BC} \rightarrow g_B$ )
- 5 Multiply  $f_{AB}$  with  $g_B$  giving  $h_{AB}$
- 6 Sum over B ( $h_{AB} \rightarrow g_A$ )
- 7 Sum over A ( $g_A \rightarrow \emptyset$ )
- 8 Just like the Forward algorithm!  $Z = 388$ .

so  $\mathbb{P}(0, 0, 1, 1) = 90/Z \approx .23$

For free:  $\mathbb{P}(A = 0) = 208/388 \approx .54$ .

# Variable elimination: more complicated example



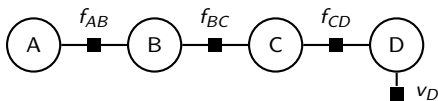
Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

1 Introduce evidence!

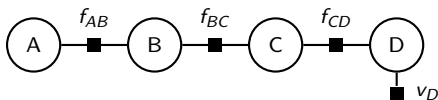
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              |   |          |
| 0 | 1 | 2              |   |          |
| 1 | 0 | 1              | 0 | 0        |
| 1 | 1 | 3              | 1 | 1        |



# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

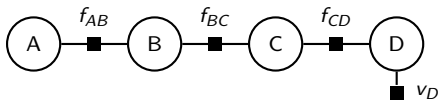
- 1 Introduce evidence!
- 2 Pick order: D, C, B, A

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 0        |
| 0 | 1 | 2              | 0 | 0        |
| 1 | 0 | 1              | 1 | 1        |
| 1 | 1 | 3              | 1 | 1        |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

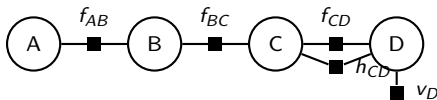
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              |   |          |
| 0 | 1 | 2              |   |          |
| 1 | 0 | 1              | 0 | 0        |
| 1 | 1 | 3              | 1 | 1        |

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors

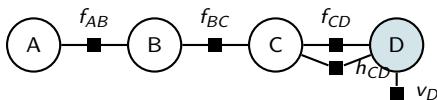
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              |   |          |
| 0 | 1 | 2              |   |          |
| 1 | 0 | 1              | 0 | 0        |
| 1 | 1 | 3              | 1 | 1        |

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 3              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors
- 4 Sum over  $D$  ( $h_{CD} \rightarrow g_C$ )

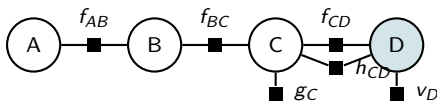
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              |   |          |
| 0 | 1 | 2              |   |          |
| 1 | 0 | 1              | 0 | 0        |
| 1 | 1 | 3              | 1 | 1        |

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 3              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors
- 4 Sum over  $D$  ( $h_{CD} \rightarrow g_C$ )

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

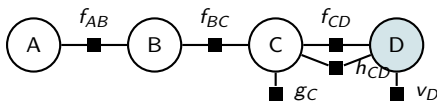
| C | $g_C(c)$ |
|---|----------|
| 0 | 2        |
| 1 | 3        |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| D | $v_D(d)$ |
|---|----------|
| 0 | 0        |
| 1 | 1        |

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 3              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors
- 4 Sum over  $D$  ( $h_{CD} \rightarrow g_C$ )
- 5 Multiply all  $C$  factors

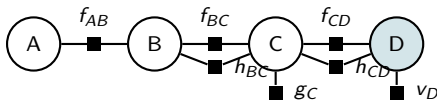
| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 2        |
| 0 | 1 | 3              | 1 | 3        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 2              |   |          |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 0        |
| 0 | 1 | 2              | 1 | 1        |
| 1 | 0 | 1              |   |          |
| 1 | 1 | 3              |   |          |

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 3              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors
- 4 Sum over  $D$  ( $h_{CD} \rightarrow g_C$ )
- 5 Multiply all  $C$  factors

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 2        |
| 1 | 3        |

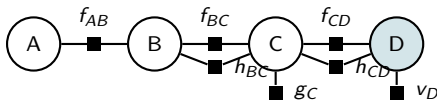
| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| D | $v_D(d)$ |
|---|----------|
| 0 | 0        |
| 1 | 1        |

| B | C | $h_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 2              |
| 0 | 1 | 9              |
| 1 | 0 | 2              |
| 1 | 1 | 6              |

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 3              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all  $D$  factors
- 4 Sum over  $D$  ( $h_{CD} \rightarrow g_C$ )
- 5 Multiply all  $C$  factors
- 6 Multiply all  $B$  factors

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 2        |
| 1 | 3        |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

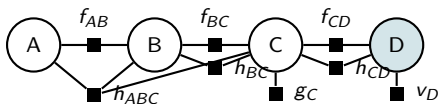
| D | $v_D(d)$ |
|---|----------|
| 0 | 0        |
| 1 | 1        |

| B | C | $h_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 2              |
| 0 | 1 | 9              |
| 1 | 0 | 2              |
| 1 | 1 | 6              |

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 0              |



# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 2        |
| 1 | 3        |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| D | $v_D(d)$ |
|---|----------|
| 0 | 0        |
| 1 | 1        |

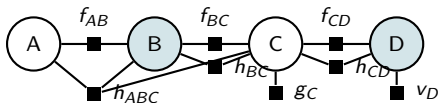
| A | B | C | $h_{ABC}(a, b, c)$ |
|---|---|---|--------------------|
| 0 | 0 | 0 | 20                 |
| 0 | 0 | 1 | 90                 |
| 0 | 1 | 0 | 4                  |
| 0 | 1 | 1 | 12                 |
| 1 | 0 | 0 | 6                  |
| 1 | 0 | 1 | 18                 |
| 1 | 1 | 0 | 18                 |
| 1 | 1 | 1 | 54                 |

| B | C | $h_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 2              |
| 0 | 1 | 9              |
| 1 | 0 | 2              |
| 1 | 1 | 6              |

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all D factors
- 4 Sum over D ( $h_{CD} \rightarrow g_C$ )
- 5 Multiply all C factors
- 6 Multiply all B factors

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 0              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

| A | B | $f_{AB}(a, b)$ |
|---|---|----------------|
| 0 | 0 | 10             |
| 0 | 1 | 2              |
| 1 | 0 | 3              |
| 1 | 1 | 9              |

| B | C | $f_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 1              |
| 0 | 1 | 3              |
| 1 | 0 | 1              |
| 1 | 1 | 2              |

| C | $g_C(c)$ |
|---|----------|
| 0 | 2        |
| 1 | 3        |

| C | D | $f_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 4              |
| 0 | 1 | 2              |
| 1 | 0 | 1              |
| 1 | 1 | 3              |

| D | $v_D(d)$ |
|---|----------|
| 0 | 0        |
| 1 | 1        |

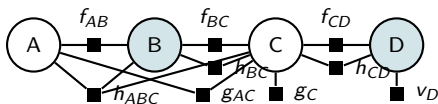
| A | B | C | $h_{ABC}(a, b, c)$ |
|---|---|---|--------------------|
| 0 | 0 | 0 | 20                 |
| 0 | 0 | 1 | 90                 |
| 0 | 1 | 0 | 4                  |
| 0 | 1 | 1 | 12                 |
| 1 | 0 | 0 | 6                  |
| 1 | 0 | 1 | 18                 |
| 1 | 1 | 0 | 18                 |
| 1 | 1 | 1 | 54                 |

| B | C | $h_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 2              |
| 0 | 1 | 9              |
| 1 | 0 | 2              |
| 1 | 1 | 6              |

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all D factors
- 4 Sum over D ( $h_{CD} \rightarrow g_C$ )
- 5 Multiply all C factors
- 6 Multiply all B factors
- 7 Sum over B.

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 0              |

# Variable elimination: more complicated example



Query:  $\mathbb{P}(a, c \mid D = 1) = ?$

| A | B | $f_{AB}(a, b)$ | A | C | $g_{AC}(a, c)$ |
|---|---|----------------|---|---|----------------|
| 0 | 0 | 10             | 0 | 0 | 24             |
| 0 | 1 | 2              | 0 | 1 | 102            |
| 1 | 0 | 3              | 1 | 0 | 24             |
| 1 | 1 | 9              | 1 | 1 | 72             |

| B | C | $f_{BC}(b, c)$ | C | $g_C(c)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 1              | 0 | 2        |
| 0 | 1 | 3              | 1 | 3        |
| 1 | 0 | 1              | 0 | 2        |
| 1 | 1 | 2              | 1 | 3        |

| C | D | $f_{CD}(c, d)$ | D | $v_D(d)$ |
|---|---|----------------|---|----------|
| 0 | 0 | 4              | 0 | 0        |
| 0 | 1 | 2              | 1 | 1        |
| 1 | 0 | 1              | 0 | 0        |
| 1 | 1 | 3              | 1 | 1        |

| A | B | C | $h_{ABC}(a, b, c)$ |
|---|---|---|--------------------|
| 0 | 0 | 0 | 20                 |
| 0 | 0 | 1 | 90                 |
| 0 | 1 | 0 | 4                  |
| 0 | 1 | 1 | 12                 |
| 1 | 0 | 0 | 6                  |
| 1 | 0 | 1 | 18                 |
| 1 | 1 | 0 | 18                 |
| 1 | 1 | 1 | 54                 |

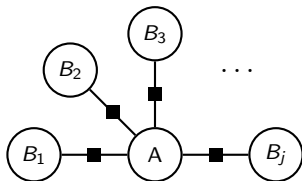
| B | C | $h_{BC}(b, c)$ |
|---|---|----------------|
| 0 | 0 | 2              |
| 0 | 1 | 9              |
| 1 | 0 | 2              |
| 1 | 1 | 6              |

- 1 Introduce evidence!
- 2 Pick order: D, C, B, A
- 3 Multiply all D factors
- 4 Sum over D ( $h_{CD} \rightarrow g_C$ )
- 5 Multiply all C factors
- 6 Multiply all B factors
- 7 Sum over B.

| C | D | $h_{CD}(c, d)$ |
|---|---|----------------|
| 0 | 0 | 0              |
| 0 | 1 | 2              |
| 1 | 0 | 0              |
| 1 | 1 | 0              |

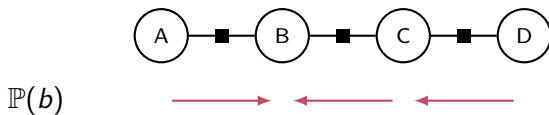
# Variable elimination

- Answer any query involving max, marginalization, evidence!
- Complexity depends on **elimination order**:  $\mathcal{O}(nk^M)$ 
  - where  $n$ =n. variables,  $k$ =dimension,  $M$ =size of largest intermediate factor.
  - Example: In chain, intuitive order has  $M = 2$ .  
eliminating from middle of chain gives  $M = 3$ .
  - Extreme example is a star graph. Best case  $M = 2$ , worst  $M = N!$



- In **chains** and **trees**: optimal order is easy. Not in general.
- When given a new query, need to restart algorithm from scratch!

# Variable elimination as message passing

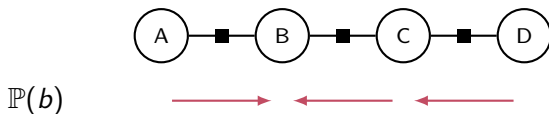


- Optimal order: A, D, C (or D, C, A)

---

<sup>2</sup>because it's a tree

# Variable elimination as message passing

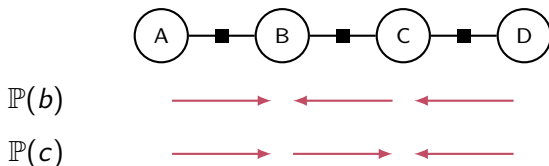


- Optimal order: A, D, C (or D, C, A)
- At each step, we eliminate a variable  $Y$  by multiplying (at most<sup>2</sup>) two factors and summing over  $Y$ :

$$g_{Y \rightarrow X}(x) = \sum_y f_{XY}(x, y) g_Y(y)$$

<sup>2</sup>because it's a tree

# Variable elimination as message passing



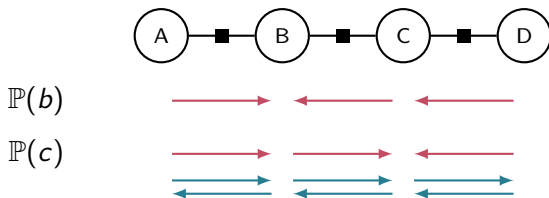
- Optimal order: A, D, C (or D, C, A)
- At each step, we eliminate a variable  $Y$  by multiplying (at most<sup>2</sup>) two factors and summing over  $Y$ :

$$g_{Y \rightarrow X}(x) = \sum_y f_{XY}(x, y) g_Y(y)$$

- These intermediate operations (“messages”) are shared for all queries,

<sup>2</sup>because it's a tree

# Variable elimination as message passing



- Optimal order: A, D, C (or D, C, A)
- At each step, we eliminate a variable  $Y$  by multiplying (at most<sup>2</sup>) two factors and summing over  $Y$ :

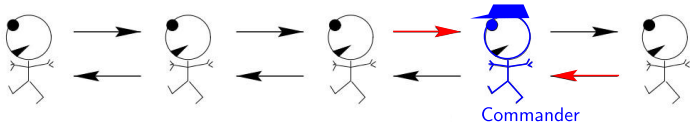
$$g_{Y \rightarrow X}(x) = \sum_y f_{XY}(x, y) g_Y(y)$$

- These intermediate operations (“messages”) are shared for all queries, so let’s compute **all messages** up front!

<sup>2</sup>because it’s a tree

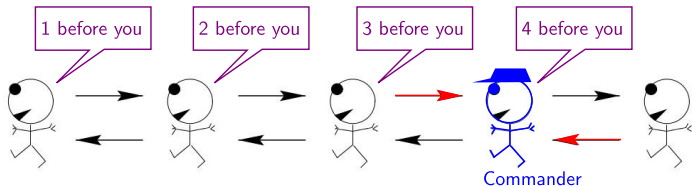


# Motivating Example: Counting Soldiers



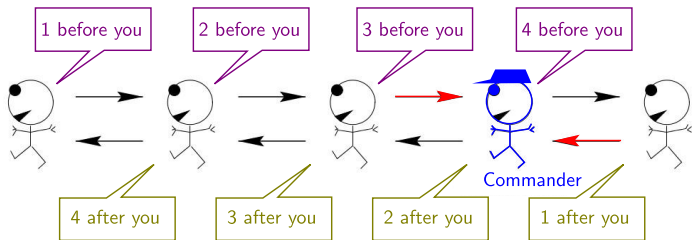
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

# Motivating Example: Counting Soldiers



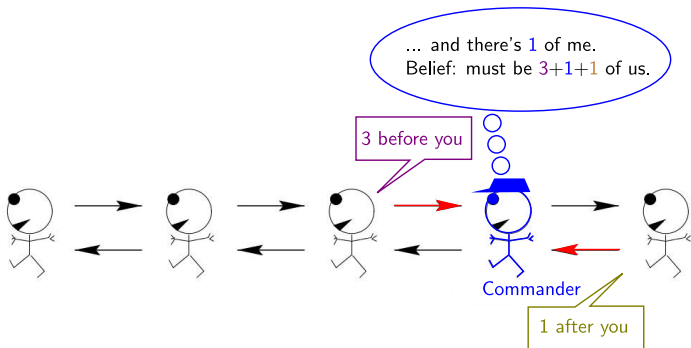
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

# Motivating Example: Counting Soldiers



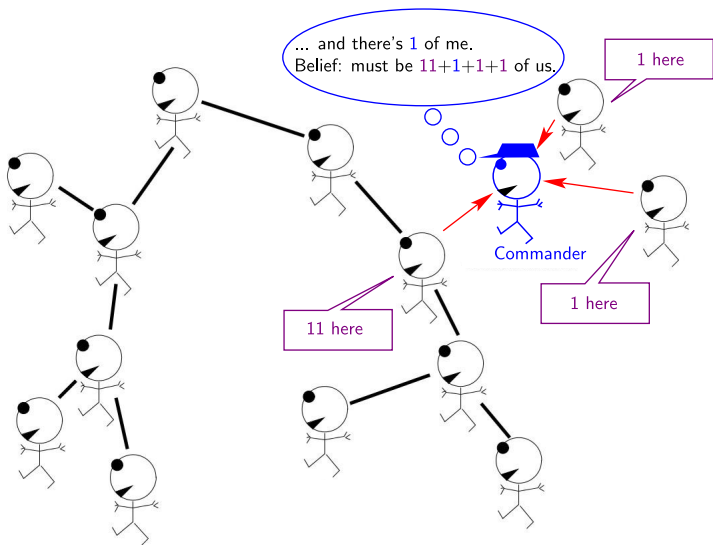
(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

# Motivating Example: Counting Soldiers



(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

# Motivating Example: Counting Soldiers



(Adapted from MacKay 2003 and Gormley & Eisner ACL'14 tutorial.)

# Message passing in a tree FG

- Messages from variable  $X$  to factor  $\alpha$ : aggregate variable beliefs from any other factors. (For leaves, this message is  $\mathbf{1}$ ).

$$\nu_{X \rightarrow \alpha}(x) = \prod_{\beta \in \mathcal{N}(X) - \alpha} \mu_{\beta \rightarrow X}(x)$$

- Messages from factor  $\alpha$  to variable  $X$ : marginalizes over all assignments  $y_1, \dots, y_k$  for  $Y_1, \dots, Y_k$  neighboring  $\alpha$

$$\mu_{\alpha \rightarrow X}(x) = \sum_{\substack{y_1, \dots, y_k \\ \{Y_1, \dots, Y_k\} = \mathcal{N}(\alpha) - X}} f_{\alpha}(x, y_1, \dots, y_k) \prod_{Y_i \in \mathcal{N}(\alpha) - X} \nu_{Y_i \rightarrow \alpha}(y_i)$$

- A message is sent once all messages it depends on have been received.
- For chain: **forward-backward!** For tree: leaves-to-root and back.
- If new evidence is added, many messages don't change.
- Replace sum with max for maximization.

# From messages to beliefs

- Once we collected all the messages, we can compute local beliefs.
- Variable beliefs:

$$p_X(x) \propto \prod_{\alpha \in \mathcal{N}(X)} \mu_{\alpha \rightarrow X}(x)$$

- Factor beliefs:

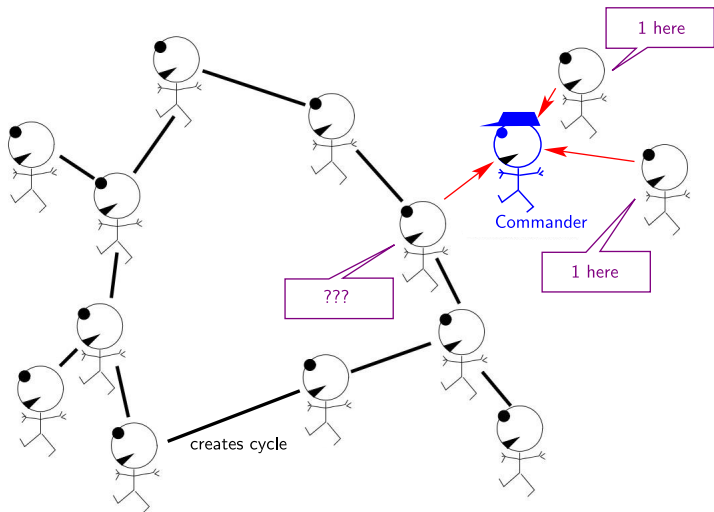
$$p_\alpha(x_1, \dots, x_k) \propto f_\alpha(x_1, \dots, x_k) \prod_{X_i \in \mathcal{N}(\alpha)} \nu_{X_i \rightarrow \alpha}(x_i)$$

- If no cycles, once all messages are passed, beliefs are true marginals:

$$p_X(x) = \mathbb{P}(x), \quad p_\alpha(x_1, \dots, x_k) = \mathbb{P}(x_1, \dots, x_k).$$

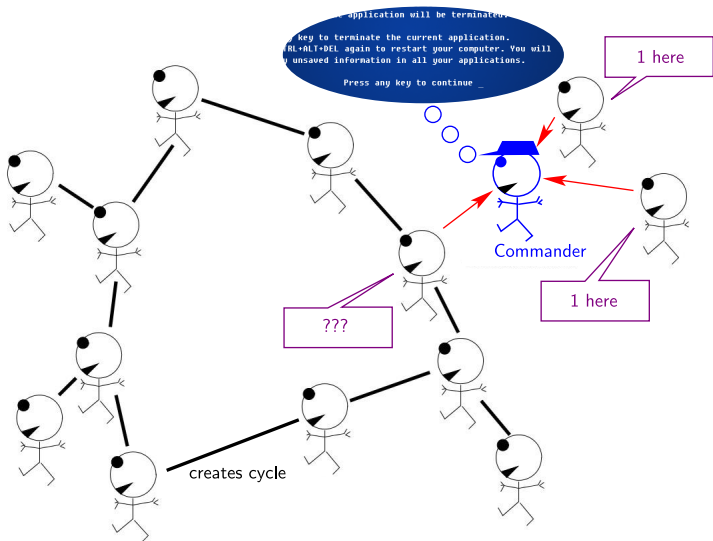
- What to do if there are cycles?

# Counting Soldiers with Loops



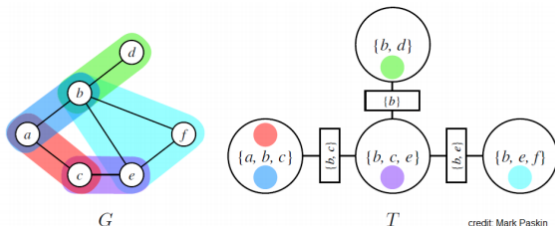


# Counting Soldiers with Loops



# Inference in loopy graphs

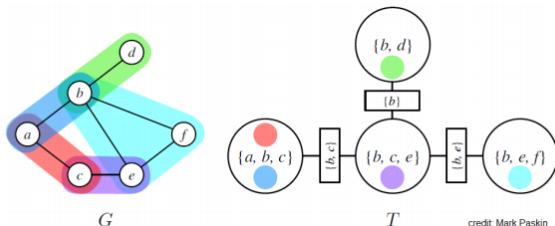
- Exact solution: **Junction Tree** algorithm:
  - convert the graph into a tree, by merging cliques!



- Complexity: like variable elimination. Finding the best tree is NP-hard. (corresponds to finding an ordering for variable elimination.)
- Better than VE because we get all marginals at once.

# Inference in loopy graphs

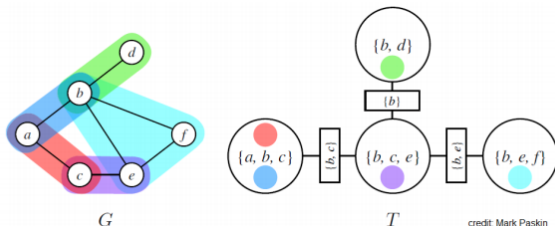
- Exact solution: **Junction Tree** algorithm:
  - convert the graph into a tree, by merging cliques!



- Complexity: like variable elimination. Finding the best tree is NP-hard. (corresponds to finding an ordering for variable elimination.)
- Better than VE because we get all marginals at once.
- Approximate solution: **Loopy Belief Propagation**:
  - initialize all messages;
  - pass messages in some order until convergence.
  - (may not terminate, result not guaranteed correct, but works ok.)

# Inference in loopy graphs

- Exact solution: **Junction Tree** algorithm:
  - convert the graph into a tree, by merging cliques!



- Complexity: like variable elimination. Finding the best tree is NP-hard. (corresponds to finding an ordering for variable elimination.)
- Better than VE because we get all marginals at once.
- Approximate solution: **Loopy Belief Propagation**:
  - initialize all messages;
  - pass messages in some order until convergence.
  - (may not terminate, result not guaranteed correct, but works ok.)
  - Many recent algorithms (early 2010s).

# CRFs for any factor graph

Above, we took the factor scores for granted. We can learn to model them:

# CRFs for any factor graph

Above, we took the factor scores for granted. We can learn to model them:

Use some model (neural or feature-based) to produce **unary scores**:

$$f_A(y) = \exp s_{A,y} = \text{(for example)} \exp w_{A,y} \cdot \phi_A(x)$$

and **pairwise scores**:

$$f_{AB}(y, y') = \exp s_{AB,y,y'} = \text{(for example)} \exp w_{A,B,y,y'} \cdot \phi_{A,B}(x)$$

# CRFs for any factor graph

Above, we took the factor scores for granted. We can learn to model them:

Use some model (neural or feature-based) to produce **unary scores**:

$$f_A(y) = \exp s_{A,y} = \text{(for example)} \exp w_{A,y} \cdot \phi_A(x)$$

and **pairwise scores**:

$$f_{AB}(y, y') = \exp s_{AB,y,y'} = \text{(for example)} \exp w_{A,B,y,y'} \cdot \phi_{A,B}(x)$$

(In general, **factor scores**  $f_\alpha(y_\alpha) = \exp s_{\alpha,y_\alpha}$ )

# CRFs for any factor graph

Above, we took the factor scores for granted. We can learn to model them:

Use some model (neural or feature-based) to produce **unary scores**:

$$f_A(y) = \exp s_{A,y} = \text{(for example)} \exp w_{A,y} \cdot \phi_A(x)$$

and **pairwise scores**:

$$f_{AB}(y, y') = \exp s_{AB,y,y'} = \text{(for example)} \exp w_{A,B,y,y'} \cdot \phi_{A,B}(x)$$

(In general, **factor scores**  $f_\alpha(y_\alpha) = \exp s_{\alpha,y_\alpha}$ )

The probability of an entire labeling  $y$  is then

$$\mathbb{P}(y | x) = \frac{\prod_\alpha f_\alpha(y_\alpha)}{Z} \quad \text{meaning} \quad \log \mathbb{P}(y | x) = \sum_\alpha s_{\alpha,y_\alpha} - \log Z$$



# CRFs for any factor graph

Above, we took the factor scores for granted. We can learn to model them:

Use some model (neural or feature-based) to produce **unary scores**:

$$f_A(y) = \exp s_{A,y} = \text{(for example)} \exp w_{A,y} \cdot \phi_A(x)$$

and **pairwise scores**:

$$f_{AB}(y, y') = \exp s_{AB,y,y'} = \text{(for example)} \exp w_{A,B,y,y'} \cdot \phi_{A,B}(x)$$

(In general, **factor scores**  $f_\alpha(y_\alpha) = \exp s_{\alpha,y_\alpha}$ )

The probability of an entire labeling  $y$  is then

$$\mathbb{P}(y | x) = \frac{\prod_\alpha f_\alpha(y_\alpha)}{Z} \quad \text{meaning} \quad \log \mathbb{P}(y | x) = \sum_\alpha s_{\alpha,y_\alpha} - \log Z$$

Gradient updates wrt a factor's scores:

$$\frac{\partial \log \mathbb{P}(y | x)}{\partial s_{\alpha,y_\alpha}} = \mathbb{I}[y_\alpha = y_\alpha^{\text{true}}] - \mathbb{P}(y_\alpha | x)$$

# CRFs for any factor graph

Above, we took the factor scores for granted. We can learn to model them:

Use some model (neural or feature-based) to produce **unary scores**:

$$f_A(y) = \exp s_{A,y} = \text{(for example)} \exp w_{A,y} \cdot \phi_A(x)$$

and **pairwise scores**:

$$f_{AB}(y, y') = \exp s_{AB,y,y'} = \text{(for example)} \exp w_{A,B,y,y'} \cdot \phi_{A,B}(x)$$

(In general, **factor scores**  $f_\alpha(y_\alpha) = \exp s_{\alpha,y_\alpha}$ )

The probability of an entire labeling  $y$  is then

$$\mathbb{P}(y \mid x) = \frac{\prod_\alpha f_\alpha(y_\alpha)}{Z} \quad \text{meaning} \quad \log \mathbb{P}(y \mid x) = \sum_\alpha s_{\alpha,y_\alpha} - \log Z$$

Gradient updates wrt a factor's scores:

$$\frac{\partial \log \mathbb{P}(y \mid x)}{\partial s_{\alpha,y_\alpha}} = \mathbb{1}[y_\alpha = y_\alpha^{\text{true}}] - \mathbb{P}(y_\alpha \mid x)$$

The updates use the factor beliefs  $\mathbb{P}(y_\alpha \mid x) = p_\alpha(y_\alpha)$  for each factor!

# Undirected models: summary

- MRFs and pairwise MRFs, both special cases of FGs.
- Powerful, expressive, widely used for discriminative modelling.
- Exact inference when not loopy.
  - We've seen some ideas of what to do when loopy
  - We did not cover more advanced approaches, relating message passing and dual decomposition: (Martins et al., 2015; Kolmogorov, 2006; Komodakis et al., 2007; Globerson and Jaakkola, 2007)
- For learning: a generalization of linear-chain CRFs

# References I

- Dawid, A. P. (2010). Beware of the DAG! In *Causality: objectives and assessment*, pages 59–86.
- Globerson, A. and Jaakkola, T. (2007). Fixing Max-Product: Convergent message passing algorithms for MAP LP-relaxations. In *Proc. of NeurIPS*.
- Koller, D. and Friedman, N. (2009). *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press.
- Kolmogorov, V. (2006). Convergent Tree-Reweighted Message Passing for energy minimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(10):1568–1583.
- Komodakis, N., Paragios, N., and Tziritas, G. (2007). MRF optimization via dual decomposition: Message-Passing revisited. In *Proc. of ICCV*.
- MacKay, D. (2003). *Information Theory, Inference, and Learning Algorithms*, volume 7. Cambridge University Press.
- Martins, A. F., Figueiredo, M. A., Aguiar, P. M., Smith, N. A., and Xing, E. P. (2015). AD3: Alternating directions dual decomposition for MAP inference in graphical models. *JMLR*, 16(1):495–545.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann.
- Pearl, J. (2000). *Causality: models, reasoning and inference*, volume 29. Springer.
- Pearl, J. (2012). The do-calculus revisited. *arXiv preprint arXiv:1210.4852*.