

Lecture 7: Probabilistic Graphical Models

Vlad Niculae & André Martins

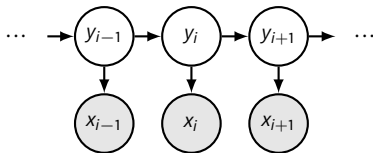


Deep Structured Learning Course, Fall 2019

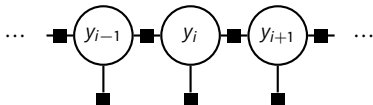
Graphical Models

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

Directed
(today)



Undirected
(next time)



1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

Factor graphs

1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

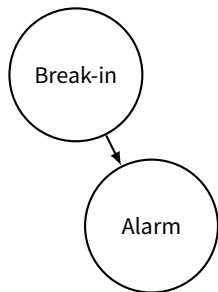
Factor graphs

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?

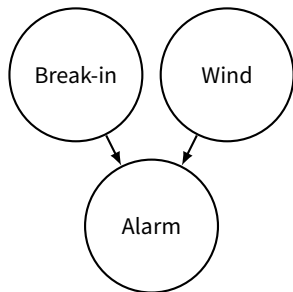


P(B)	B=yes	B=no
	.05	.95
P(A B)	A=on	A=off
B=yes	.99	.01
B=no	.10	.90

- $P(B | A) = ?$

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?

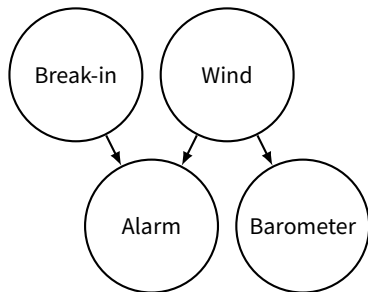


P(B)		B=yes	B=no
		.05	.95
P(A B, W)		A=on	A=off
B=yes	W=lo	.99	.01
B=yes	W=med	.99	.01
B=yes	W=hi	.999	.001
B=no	W=lo	.01	.99
B=no	W=med	.05	.95
B=no	W=hi	.25	.75

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

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?



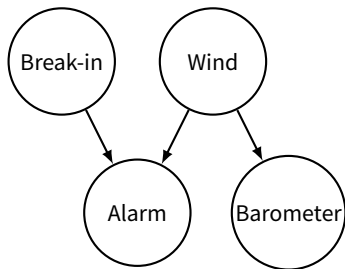
P(B)		B=yes	B=no
		.05	.95
P(A B, W)		A=on	A=off
B=yes	W=lo	.99	.01
B=yes	W=med	.99	.01
B=yes	W=hi	.999	.001
B=no	W=lo	.01	.99
B=no	W=med	.05	.95
B=no	W=hi	.25	.75

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

Maybe we're in the basement, but have a barometer.

Bayes networks

Toolkit for encoding **knowledge** about **interaction structures** between random variables.



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

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

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

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

Brk.	Wind	Alarm	Bar.	P	Brk.	Wind	Alarm	Bar.	P
yes	lo	on	lo	0.0243	no	lo	on	lo	0.0047
yes	lo	on	med	0.0002	no	lo	on	med	4.75e-05
yes	lo	on	hi	0.0002	no	lo	on	hi	4.75e-05
yes	lo	off	lo	0.0002	no	lo	off	lo	0.4608
yes	lo	off	med	2.50e-06	no	lo	off	med	0.0047
yes	lo	off	hi	2.50e-06	no	lo	off	hi	0.0047
yes	med	on	lo	0.0001	no	med	on	lo	0.0001
yes	med	on	med	0.0146	no	med	on	med	0.0140
yes	med	on	hi	0.0001	no	med	on	hi	0.0001
yes	med	off	lo	1.50e-06	no	med	off	lo	0.0027
yes	med	off	med	0.0001	no	med	off	med	0.2653
yes	med	off	hi	1.50e-06	no	med	off	hi	0.0027
yes	hi	on	lo	9.99e-05	no	hi	on	lo	0.0005
yes	hi	on	med	9.99e-05	no	hi	on	med	0.0005
yes	hi	on	hi	0.0098	no	hi	on	hi	0.0466
yes	hi	off	lo	1.00e-07	no	hi	off	lo	0.0014
yes	hi	off	med	1.00e-07	no	hi	off	med	0.0014
yes	hi	off	hi	9.80e-06	no	hi	off	hi	0.1397

Without any structure, $P(\text{Break-in, Wind, Alarm, Barometer})$ 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

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

Without any structure, $P(\text{Break-in, Wind, Alarm, Barometer})$ 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

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

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

Without any structure, P(Break-in, Wind, Alarm, Barometer) would have to be stored & estimated like

Brk.	Wind	Alarm	Bar.	P	Brk.	Wind	Alarm	Bar.	P
yes	lo	on	lo	0.0243	no	lo	on	lo	0.0047
yes	lo	on	med	0.0002	no	lo	on	med	4.75e-05
yes	lo	on	hi	0.0002	no	lo	on	hi	4.75e-05
yes	lo	off	lo	0.0002	no	lo	off	lo	0.4608
yes	lo	off	med	2.50e-06	no	lo	off	med	0.0047
yes	lo	off	hi	2.50e-06	no	lo	off	hi	0.0047
yes	med	on	lo	0.0001	no	med	on	lo	0.0001
yes	med	on	med	0.0146	no	med	on	med	0.0140
yes	med	on	hi	0.0001	no	med	on	hi	0.0001
yes	med	off	lo	1.50e-06	no	med	off	lo	0.0027
yes	med	off	med	0.0001	no	med	off	med	0.2653
yes	med	off	hi	1.50e-06	no	med	off	hi	0.0027
yes	hi	on	lo	9.99e-05	no	hi	on	lo	0.0005
yes	hi	on	med	9.99e-05	no	hi	on	med	0.0005
yes	hi	on	hi	0.0098	no	hi	on	hi	0.0466
yes	hi	off	lo	1.00e-07	no	hi	off	lo	0.0014
yes	hi	off	med	1.00e-07	no	hi	off	med	0.0014
yes	hi	off	hi	9.80e-06	no	hi	off	hi	0.1397

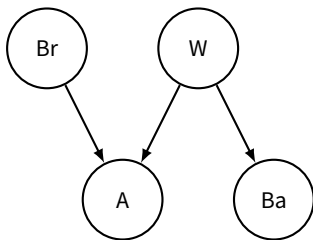
$$P(\text{Break-in=yes, Alarm=on}) = 0.0496$$

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

$$P(\text{Break-in=yes} \mid \text{Alarm=on}) = \frac{P(\text{Break-in=yes, Alarm=on})}{\sum_b P(\text{Break-in}=b, \text{Alarm=on})}$$

$$= .427$$

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



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

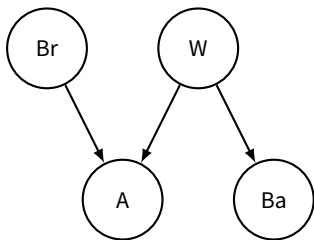
P(Br)	yes	no
	.05	.95

P(W)	lo	mid	hi
	.5	.3	.2

P(A 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(Ba 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.



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

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

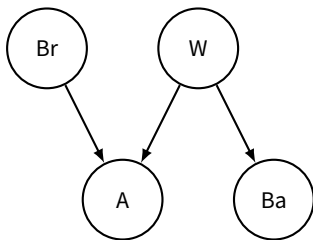
P(Br)	yes	no
	.05	.95

P(W)	lo	mid	hi
	.5	.3	.2

P(A 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(Ba 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.



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

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

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

P(Br)	yes	no
	.05	.95

P(W)	lo	mid	hi
	.5	.3	.2

P(A 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(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

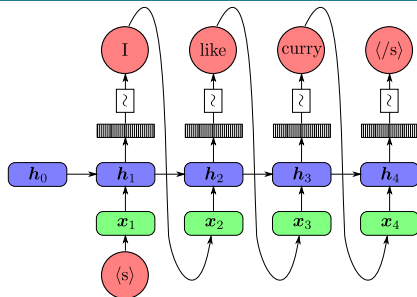
encode structure, not parametrization

are non-unique for a distribution

encode independence **requirements**, not necessarily all

BN are not neural net diagrams

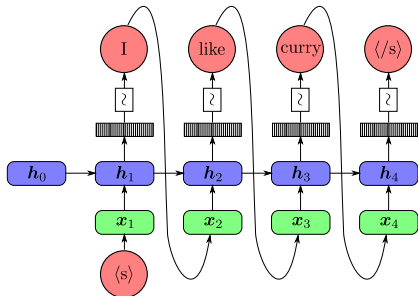
Recall the RNN language model:



- In statistical terms, what are we modeling?

BN are not neural net diagrams

Recall the RNN language model:

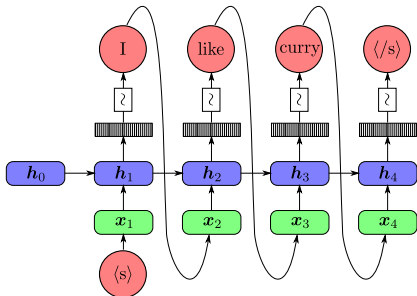


- In statistical terms, what are we modeling?

$$P(X_1, \dots, X_n) = P(X_1) P(X_2 \mid X_1) P(X_3 \mid X_1, X_2) \dots$$

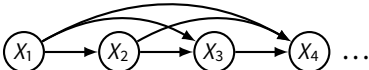
BN are not neural net diagrams

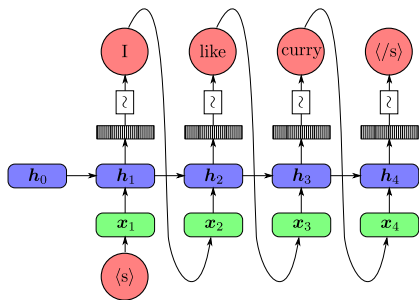
Recall the RNN language model:



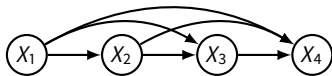
- In statistical terms, what are we modeling?

$$P(X_1, \dots, X_n) = P(X_1) P(X_2 \mid X_1) P(X_3 \mid X_1, X_2) \dots$$

- 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)

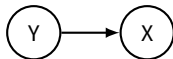
BN encode structure, not parametrization

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:

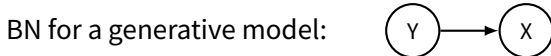


We must posit what are $P(Y)$ and $P(X | Y)$. **Many possible options!**

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, $Y \in \{\text{sports, music, ...}\}$.



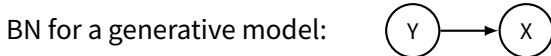
We must posit what are $P(Y)$ and $P(X | Y)$. **Many possible options!**

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

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, $Y \in \{\text{sports, music, ...}\}$.



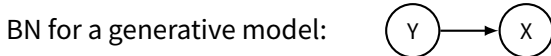
We must posit what are $P(Y)$ and $P(X | Y)$. **Many possible options!**

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



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

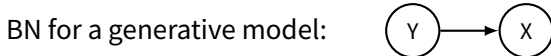
$P(Y)$: uniform: $P(Y = \text{sports}) = P(Y = \text{music}) = \frac{1}{|Y|}$, or estimated from data.

$P(X | Y)$ (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, $Y \in \{\text{sports}, \text{music}, \dots\}$.



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

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

$P(X | Y)$ (remember: values of X are sentences)

Naive Bayes

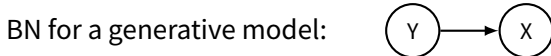
$$P(X | Y) = \prod_{j=1}^L P(X_j | Y)$$

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, $Y \in \{\text{sports, music, ...}\}$.



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

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

$P(X | Y)$ (remember: values of X are sentences)

Naive Bayes

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

Per-class Markov language model

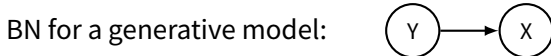
$$P(X | Y) = \prod_{j=1}^L P(X_j | X_{j-1}, Y)$$

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, $Y \in \{\text{sports, music, ...}\}$.



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

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

$P(X | Y)$ (remember: values of X are sentences)

Naive Bayes

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

Per-class Markov language model

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

Per-class recurrent NN language model

$$P(X | Y) = \text{LSTM}(x_1, \dots, x_L; w_y)$$

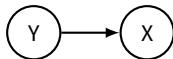
BN encode structure, not parametrization

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 $P(Y)$ and $P(X | Y)$. **Many possible options!**

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

$P(X | Y)$ (remember: values of X are sentences)

Naive Bayes

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

Per-class Markov language model

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

Per-class recurrent NN language model

$$P(X | Y) = \text{LSTM}(x_1, \dots, x_L; w_y)$$

$P(X | Y)$ need not be parametrized as a table.

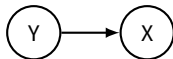
BN encode structure, not parametrization

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 $P(Y)$ and $P(X | Y)$. **Many possible options!**

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

$P(X | Y)$ (remember: values of X are sentences)

Naive Bayes

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

Per-class Markov language model

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

Per-class recurrent NN language model

$$P(X | Y) = \text{LSTM}(x_1, \dots, x_L; w_y)$$

$P(X | Y)$ need not be parametrized as a table.

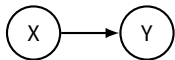
Variables need not be discrete! mixture of Gaussians: $P(X | Y = y) \sim \mathcal{N}(\mu_y, \Sigma_y)$.

Equivalent factorizations

There are many possible factorizations! $P(X, Y) =$

Equivalent factorizations

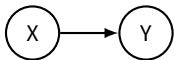
There are many possible factorizations! $P(X, Y) =$



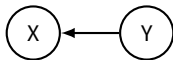
$$P(X) P(Y \mid X)$$

Equivalent factorizations

There are many possible factorizations! $P(X, Y) =$



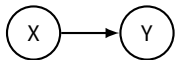
$$P(X) P(Y \mid X)$$



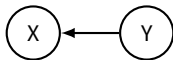
$$P(Y) P(X \mid Y)$$

Equivalent factorizations

There are many possible factorizations! $P(X, Y) =$



$$P(X) P(Y \mid X)$$



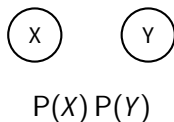
$$P(Y) P(X \mid Y)$$



$$P(X) P(Y)$$

Equivalent factorizations

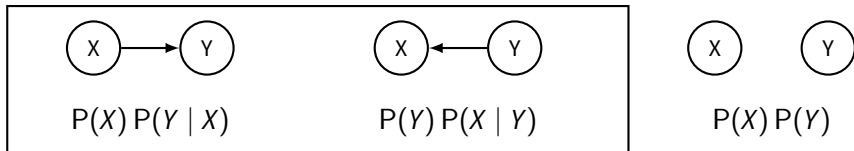
There are many possible factorizations! $P(X, Y) =$



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

Equivalent factorizations

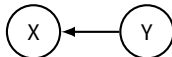
There are many possible factorizations! $P(X, Y) =$



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

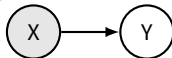
In fact, recall generative vs discriminative classifiers!

- Generative (e.g. naïve Bayes):



To classify, we would compute $P(Y | X)$ via Bayes' rule.

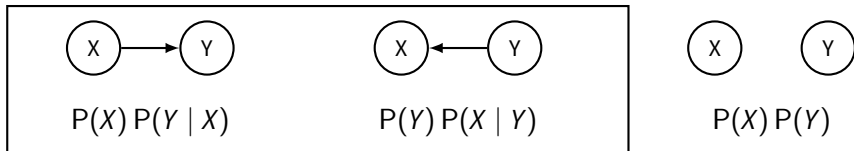
- Discriminative (e.g. logistic regression)



in LR, we don't model $P(X)$, we assume X is always observed (gray).

Equivalent factorizations

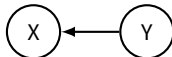
There are many possible factorizations! $P(X, Y) =$



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

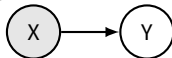
In fact, recall generative vs discriminative classifiers!

- Generative (e.g. naïve Bayes):



To classify, we would compute $P(Y | X)$ via Bayes' rule.

- Discriminative (e.g. logistic regression)



in LR, we don't model $P(X)$, we assume X is always observed (gray).


Some arrow direction choices are harder to estimate.

Some make more sense (why?): $\text{Barmtr.} \leftarrow \text{Wind}$ VS. $\text{Barmtr.} \rightarrow \text{Wind}$

Minimal independence assumptions

Recall, we say $X \perp\!\!\!\perp Y$ iff. $P(X, Y) = P(X)P(Y)$


Let $X = \text{grade in DSL}$, $Y = \text{month you were born}$.

Bayes net (1): 

Minimal independence assumptions

Recall, we say $X \perp\!\!\!\perp Y$ iff. $P(X, Y) = P(X)P(Y)$

Let $X = \text{grade in DSL}$, $Y = \text{month you were born}$.

Bayes net (1): 

Example parametrization:


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

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

Minimal independence assumptions

Recall, we say $X \perp\!\!\!\perp Y$ iff. $P(X, Y) = P(X)P(Y)$

Let $X = \text{grade in DSL}$, $Y = \text{month you were born}$.

Bayes net (1): 

Example parametrization:

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

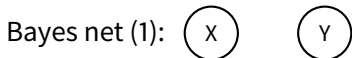
P(Y)	Jan	Feb	Mar	...
	.10	.12	.09	

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.



Example parametrization:

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

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

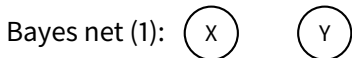
Does it mean we *must* have $X \not\perp\!\!\!\perp Y$?

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 we *must* have $X \not\perp\!\!\!\perp Y$? **NO!**

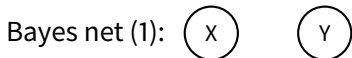
P(Y)	Jan	Feb	Mar	...
	.10	.12	.09	

P(X Y)	A+	A	B	...
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 we *must* have $X \not\perp\!\!\!\perp Y$? **NO!**

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

P(X Y)	A+	A	B	...
Y=Jan	.01	.02	.04	
Y=Feb	.01	.02	.04	
Y=Mar	.01	.02	.04	
...				

A BN constraints what independences **must be** in the model **as a minimum**.

1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

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 as

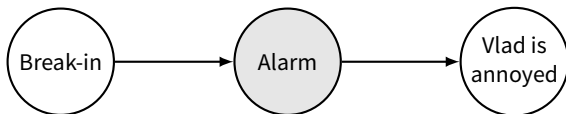
$$A \perp\!\!\!\perp B \mid C,$$

if and only if

$$P(A, B \mid C_1, \dots, C_n) = P(A \mid C_1, \dots, C_n) P(B \mid C_1, \dots, C_n).$$

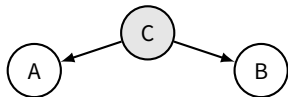
Note. Equivalently, $P(A \mid B, C_1, \dots, C_n) = P(A \mid C_1, \dots, C_n)$.

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

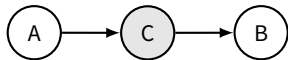


Three fundamental relationships in BN

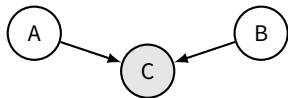
The Fork



The Chain

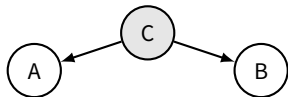


The Collider



Three fundamental relationships in BN

The Fork

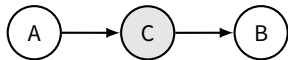


$$A \perp\!\!\!\perp B \mid C$$

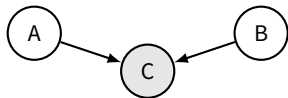
Given C , A and B are independent.

Example: Alarm \leftarrow Wind \rightarrow Barometer

The Chain

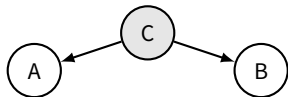


The Collider



Three fundamental relationships in BN

The Fork

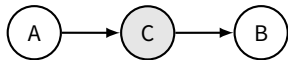


$$A \perp\!\!\!\perp B \mid C$$

Given C , A and B are independent.

Example: Alarm \leftarrow Wind \rightarrow Barometer

The Chain



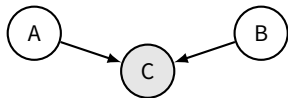
$$A \perp\!\!\!\perp B \mid C$$

After observing C ,

further observing A would not tell us about B .

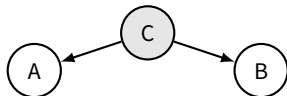
Example: Burglary \rightarrow Alarm \rightarrow Vlad distracted

The Collider



Three fundamental relationships in BN

The Fork

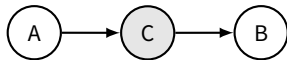


$$A \perp\!\!\!\perp B \mid C$$

Given C , A and B are independent.

Example: Alarm \leftarrow Wind \rightarrow Barometer

The Chain



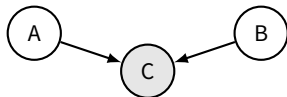
$$A \perp\!\!\!\perp B \mid C$$

After observing C ,

further observing A would not tell us about B .

Example: Burglary \rightarrow Alarm \rightarrow Vlad distracted

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

Algorithm for deciding if A and B are **d-separated** given set C , implying:

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

For all paths P from A to B in the **skeleton**¹ of the BN, at least one holds:

¹skeleton = the graph with undirected edges replacing the directed arcs

Detecting independence: d-separation

Algorithm for deciding if A and B are **d-separated** given set C , implying:

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

For all paths P from A to B in the **skeleton**¹ of the BN, at least one holds:

- 1 P includes a fork with observed parent:

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

¹skeleton = the graph with undirected edges replacing the directed arcs

Detecting independence: d-separation

Algorithm for deciding if A and B are **d-separated** given set C , implying:

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

For all paths P from A to B in the **skeleton**¹ of the BN, at least one holds:

- 1 P includes a fork with observed parent:

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

- 2 P includes a chain with observed middle:

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

¹skeleton = the graph with undirected edges replacing the directed arcs

Detecting independence: d-separation

Algorithm for deciding if A and B are **d-separated** given set C , implying:

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

For all paths P from A to B in the **skeleton**¹ of the BN, at least one holds:

- 1 P includes a fork with observed parent:

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

- 2 P includes a chain with observed middle:

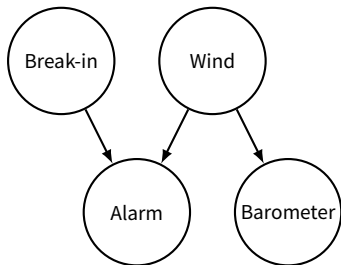
$$X \rightarrow C \rightarrow Y \quad \text{or} \quad X \leftarrow C \leftarrow Y \quad (\text{with } C \in C)$$

- 3 P includes a collider

$$X \rightarrow U \leftarrow Y \quad (\text{with } U \notin C)$$

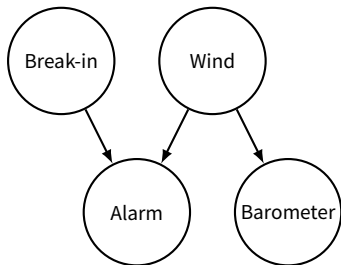
¹skeleton = the graph with undirected edges replacing the directed arcs

Examples



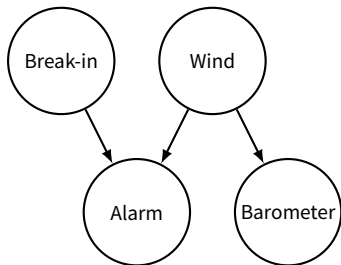
Wind $\perp\!\!\!\perp$ Barometer?

Examples



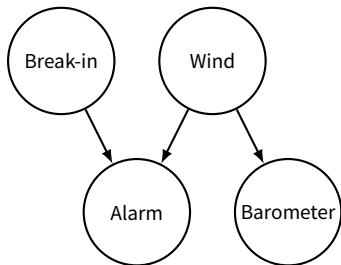
Wind $\perp\!\!\!\perp$ Barometer? **No**

Examples



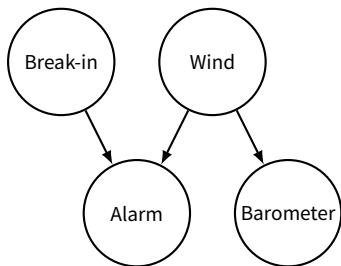
Wind $\perp\!\!\!\perp$ Barometer? **No**
Break-in $\perp\!\!\!\perp$ Wind?

Examples



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

Examples

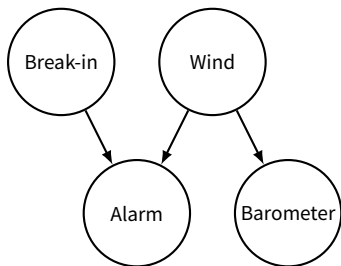


Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

Examples

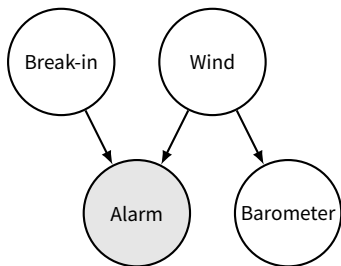


Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

Examples



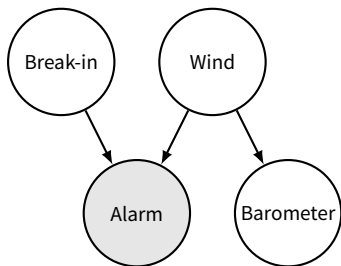
Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

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

Examples



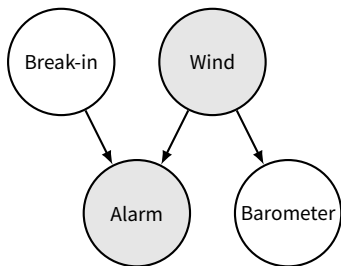
Wind $\perp\!\!\!\perp$ Barometer? **No**

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**

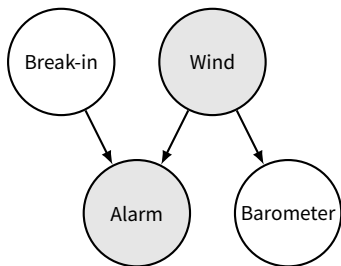
Break-in $\perp\!\!\!\perp$ Wind? **Yes**

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**

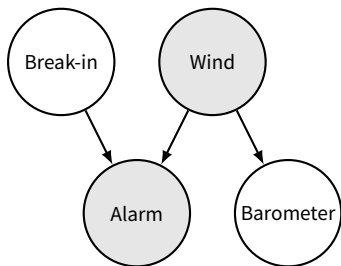
Break-in $\perp\!\!\!\perp$ Wind? **Yes**

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

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

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

Examples



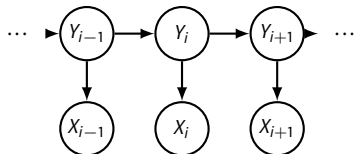
Wind $\perp\!\!\!\perp$ Barometer? **No**

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

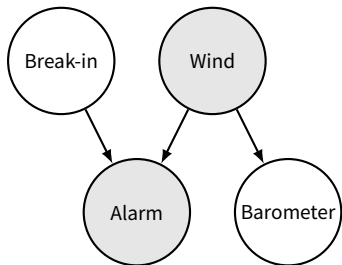
Break-in $\perp\!\!\!\perp$ Barometer? **Yes**

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

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



Examples



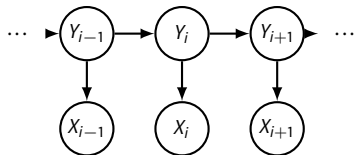
Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

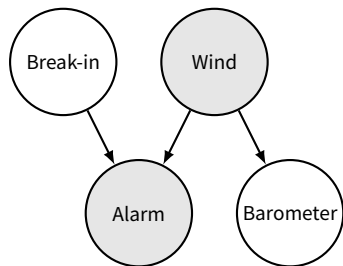
Break-in $\perp\!\!\!\perp$ Barometer | Alarm? **No**

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



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

Examples



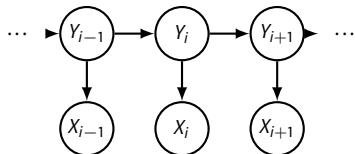
Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

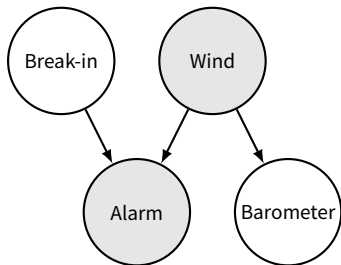
Break-in $\perp\!\!\!\perp$ Barometer | Alarm? **No**

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



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

Examples



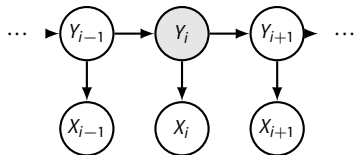
Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

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

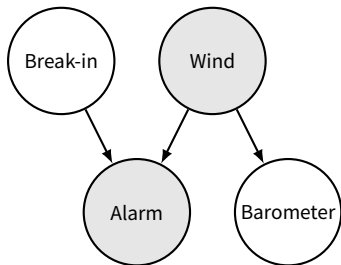
Break-in $\perp\!\!\!\perp$ Barometer | Alarm, Wind? **Yes**



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

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

Examples



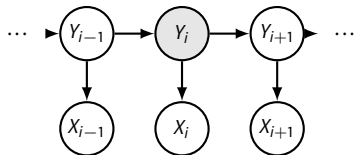
Wind $\perp\!\!\!\perp$ Barometer? **No**

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

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

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

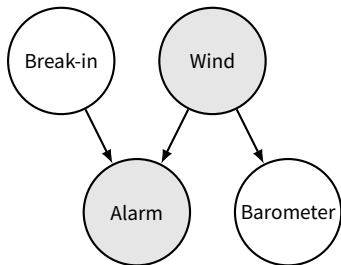
Break-in $\perp\!\!\!\perp$ Barometer | Alarm, Wind? **Yes**



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

$Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$? **Yes**

Examples



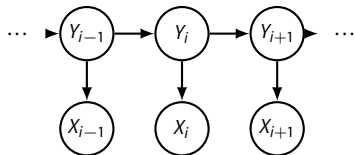
Wind $\perp\!\!\!\perp$ Barometer? **No**

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**

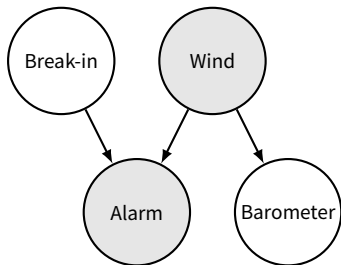


$Y_{i+1} \perp\!\!\!\perp Y_{i-1}$? **No**

$Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$? **Yes**

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

Examples



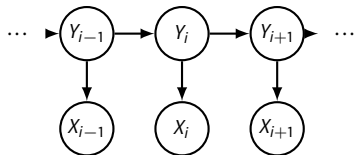
Wind $\perp\!\!\!\perp$ Barometer? **No**

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**

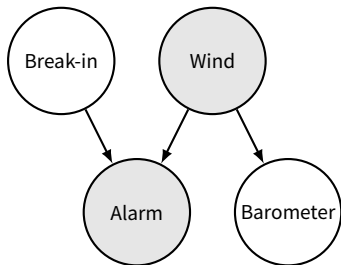


$Y_{i+1} \perp\!\!\!\perp Y_{i-1}$? **No**

$Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$? **Yes**

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

Examples



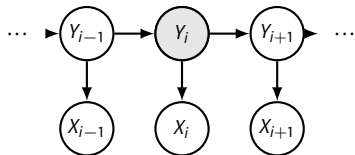
Wind $\perp\!\!\!\perp$ Barometer? **No**

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**



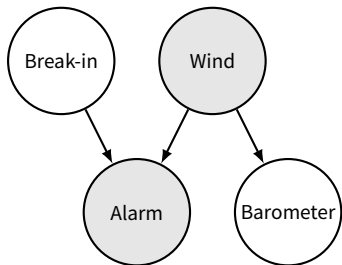
$Y_{i+1} \perp\!\!\!\perp Y_{i-1}$? **No**

$Y_{i+1} \perp\!\!\!\perp Y_{i-1} \mid Y_i$? **Yes**

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

$Y_{i+1} \perp\!\!\!\perp X_i \mid Y_i$?

Examples



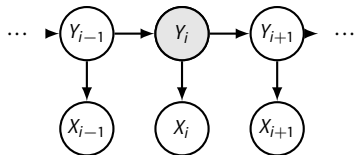
Wind $\perp\!\!\!\perp$ Barometer? **No**

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*:

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

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

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

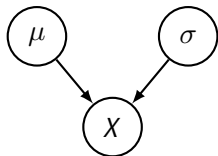
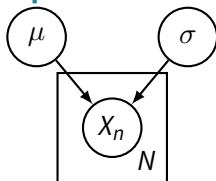


Plate notation is a way to denote **repetition of templates**:

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

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

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



1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

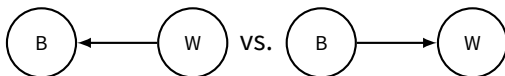
Factor graphs

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

Seeing versus doing

Bayes nets only model independence assumptions.

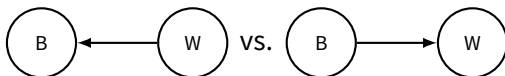
The correlation between the a barometer reading B and wind strength W can be represented either way:



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:



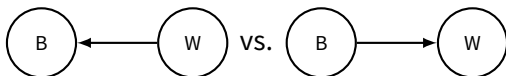
Seeing that the barometer reading is high, we can forecast wind.

$P(W \mid B)$	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.

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

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:

- we got the direction wrong
- we missed some confounding factor

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

Seeing versus doing

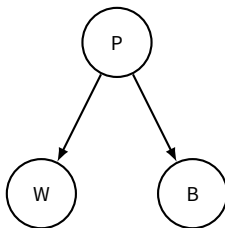
Setting the barometer needle to high manually **won't cause wind!**

Two reasons why doing \neq seeing:

- we got the direction wrong
- we missed some confounding factor

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

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



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**.

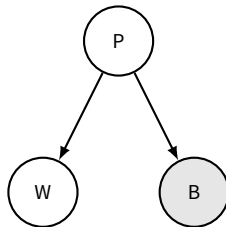
- A causes B only if A is an ancestor of B in \mathcal{G} .
- $A \rightarrow B$ means A is a direct cause of B .

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



Seeing versus doing, more rigorously

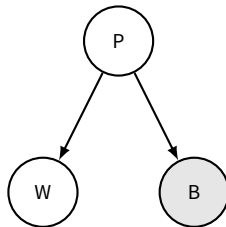
Seeing (*observational*): $P(W \mid B = \text{hi})$

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:

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



Seeing versus doing, more rigorously

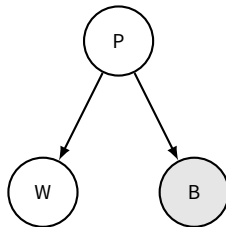
Seeing (*observational*): $P(W \mid B = \text{hi})$

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:

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



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

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

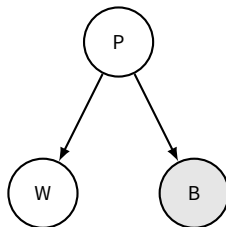
Seeing (*observational*): $P(W \mid B = \text{hi})$

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

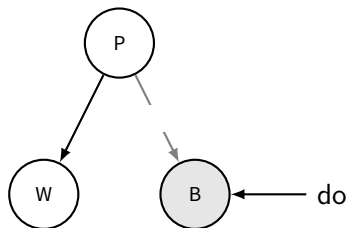
gives:

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



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

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

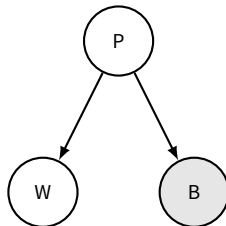
Seeing (*observational*): $P(W \mid B = \text{hi})$

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:

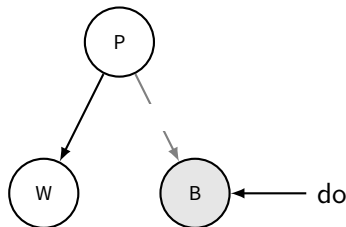
$P(W \mid B)$	lo	mid	hi
$B = \text{hi}$.01	.01	.98



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

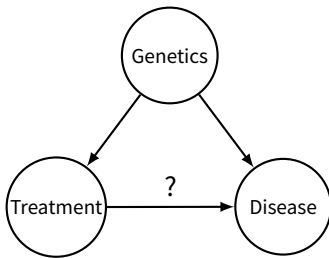
Set the needle to high **breaking inbound arrows**;
re-generate **new** data in this **new** DAG
(or estimate what that would give.)

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



Randomized controlled trials

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

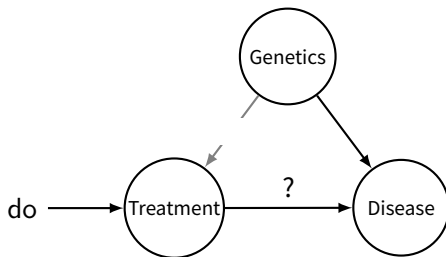


Patient	Treatment	Genetics	Disease
#42	real	?	cured
#68	placebo	?	not cured
...			

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 $P(X \mid \text{do}(Y))$ until 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:

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

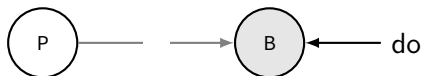
2 Action/observation exchange: the back-door criterion

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

3 Ignoring actions

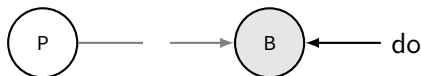
$$P(y \mid \text{do}(x), \text{do}(z), w) = P(y \mid \text{do}(x), w) \quad \text{if} \quad (Y \perp\!\!\!\perp Z \mid X, W)_{\mathcal{G}_{\bar{X}, Z(\bar{w})}}$$

Examples 1,2: Pressure and barometer

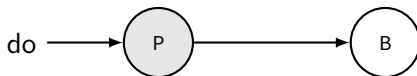


Rule 3: $P(P = \text{hi} \mid \text{do}(B = \text{hi})) = P(P = \text{hi})$ since $(P \perp\!\!\!\perp B)_{\mathcal{G}_{\bar{B}}}$

Examples 1,2: Pressure and barometer

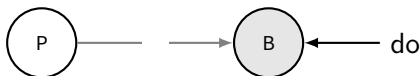


Rule 3: $P(P = \text{hi} \mid \text{do}(B = \text{hi})) = P(P = \text{hi})$ since $(P \perp\!\!\!\perp B)_{\mathcal{G}_{\bar{B}}}$

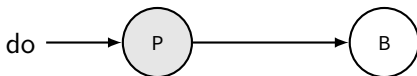


Rule 2: $P(B = \text{hi} \mid \text{do}(P = \text{lo})) = P(B = \text{hi} \mid P = \text{lo})$ since $(B \perp\!\!\!\perp P)_{\mathcal{G}_P}$

Examples 1,2: Pressure and barometer



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

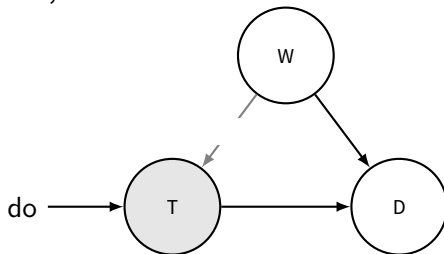


Rule 2: $P(B = \text{hi} \mid \text{do}(P = \text{lo})) = P(B = \text{hi} \mid P = \text{lo})$ since $(B \perp\!\!\!\perp P)_{\mathcal{G}_P}$

Good check: we get the intuitively correct results.

Example 3: Measurable confounder

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



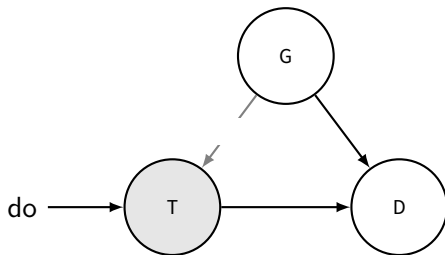
Condition on wealth (which thus needs to be measurable)

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

Example 3: an impossible one

T : treatment, D : disease.

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

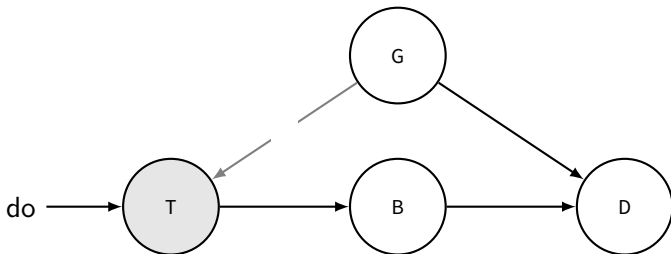


Without more info or more assumptions, we're stuck!

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.)

$$P(D = \text{cured} \mid \text{do}(T = y)) = \sum_{t,b} P(D = \text{cured} \mid T = t, B = b) P(B = b \mid T = t) P(T = t)$$

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

1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

Factor graphs

1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

Factor graphs

1 Directed Models

Bayes networks

Conditional independence and D-separation

Causal graphs & the *do* operator

2 Undirected Models

Markov networks

Factor graphs

References I

Dawid, A. P. (2010). Beware of the DAG! In *Causality: objectives and assessment*, pages 59–86.

Koller, D. and Friedman, N. (2009). *Probabilistic Graphical Models: Principles and Techniques*. The MIT Press.

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*.