

PROBABILISTIC INFERENCE AND LEARNING

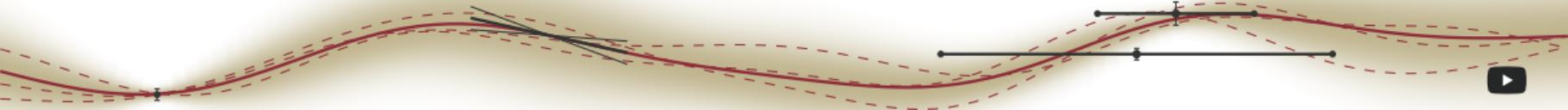
LECTURE 02

REASONING UNDER UNCERTAINTY

Philipp Hennig
21 April 2020



FACULTY OF SCIENCE
DEPARTMENT OF COMPUTER SCIENCE
CHAIR FOR THE METHODS OF MACHINE LEARNING





#	date	content	Ex	#	date	content	Ex
1	20.04.	Introduction		1	14	09.06.	Logistic Regression
2	21.04.	Reasoning under Uncertainty		15	15.06.	Exponential Families	8
3	27.04.	Continuous Variables	2	16	16.06.	Graphical Models	9
4	28.04.	Monte Carlo		17	22.06.	Factor Graphs	
5	04.05.	Markov Chain Monte Carlo	3	18	23.06.	The Sum-Product Algorithm	10
6	05.05.	Gaussian Distributions		19	29.06.	Example: Topic Models	
7	11.05.	Parametric Regression	4	20	30.06.	Mixture Models	11
8	12.05.	Understanding Deep Learning		21	06.07.	EM	
9	18.05.	Gaussian Processes	5	22	07.07.	Variational Inference	12
10	19.05.	An Example for GP Regression		23	13.07.	Example: Topic Models	
11	25.05.	Understanding Kernels	6	24	14.07.	Example: Inferring Topics	13
12	26.05.	Gauss-Markov Models		25	20.07.	Example: Kernel Topic Models	
13	08.06.	GP Classification	7	26	21.07.	Revision	



Life's most important problems are, for the most part, problems of probability.

Pierre-Simon, marquis de Laplace (1749-1827)

Catch-up from Last Time:

- ▶ Probabilities are the mathematical formalization of uncertainty
- ▶ Two basic rules

product rule: $P(A, B) = P(A | B) \cdot P(B) = P(B | A) \cdot P(A)$

sum rule: $P(A) = P(A, B) + P(A, \neg B)$

- ▶ Corollary: Bayes' Theorem

$$\underbrace{P(X | D)}_{\text{posterior for } X \text{ given } D} = \frac{\overbrace{P(X)}^{\text{prior for } X} \cdot \overbrace{P(D | X)}^{\text{likelihood for } X}}{\underbrace{P(D)}_{\text{evidence for the model}}} = \frac{P(X) \cdot P(D | X)}{\sum_{x \in \mathcal{X}} P(D | x) P(x)}$$

- ▶ This extends deductive reasoning to plausible reasoning

Today:

- ▶ Building an intuition for probability
- ▶ The computational complexity of probabilistic inference





Plausible reasoning extends Boolean Logic

Catch-Up from last time

Lemma (from Bayes' theorem)

$A \Rightarrow B: P(B | A) = 1$ implies

if A is true, then B is true

A is true implies B is true

- ▶ $P(B | A) = 1$ "modus ponens"

A is false implies B becomes less plausible

- ▶ $P(B | \neg A) \leq P(B)$

B is true implies A becomes more plausible

- ▶ $P(A | B) \geq P(A)$

B is false implies A is false

- ▶ $P(\neg A | \neg B) = 1$ "modus tollens"

Lemma (from Bayes' theorem)

$P(B | A) \geq P(B)$ implies

if A is true, then B becomes more plausible

- ▶ $P(B | A) \geq P(B)$

A is true implies B becomes more plausible

- ▶ $P(B | \neg A) \leq P(B)$

A is false implies B becomes less plausible

- ▶ $P(A | B) \geq P(A)$

B is true implies A becomes more plausible

- ▶ $P(\neg A | \neg B) \geq P(\neg A)$

B is false implies A becomes less plausible





Boole was a Bayesian

G. Boole. *An Investigation of the Laws of Thought*, 1854, §XVI, p. 249

PRINCIPLE 1st. If p be the probability of the occurrence of any event, $1 - p$ will be the probability of its non-occurrence.

2nd. The probability of the concurrence of two independent events is the product of the probabilities of those events.

3rd. The probability of the concurrence of two dependent events is equal to the product of the probability of one of them by the probability that if that event occur, the other will happen also.

4th. The probability that if an event, E , take place, an event, F , will also take place, is equal to the probability of the concurrence of the events E and F , divided by the probability of the occurrence of E .



George Boole
(1815–1864)



Computational Difficulties of Probability Theory

Uncertainty is a global notion

- ▶ The joint distribution of $n = 26$ propositional variables A, B, \dots, Z has 2^n free parameters

$$\begin{array}{ll} [1] & P(A, B, \dots, Z) = \dots \\ [2] & P(\neg A, B, \dots, Z) = \dots \\ [3] & P(A, \neg B, \dots, Z) = \dots \\ \vdots & \vdots \\ [67\,108\,863] & P(\neg A, \neg B, \dots, Z) = \dots \\ [67\,108\,864] & P(\neg A, \neg B, \dots, \neg Z) = 1 - \sum P(\dots) \end{array}$$

- ▶ requires not just large memory, but computing marginals like $P(A)$ is also very expensive
- ▶ nb: just committing to a single guess is **much** (exponentially in n) cheaper
- ▶ can we specify the joint distribution with fewer numbers?





Computing with Probabilities

- ▶ Probabilistic reasoning extends propositional logic
- ▶ instead of tracking a single *true* value, we have to assign probabilities to *combinatorially many* hypotheses

Being uncertain is potentially *much* more expensive in terms of computation and memory than simply committing to a single hypothesis. This is *the challenge* of probabilistic reasoning in practice.





A note on notation

somewhat unfortunate, but very helpful in the remainder

So far, A was a propositional variable that forms formulae:

$P(A)$ = probability that formula A is true

$P(\neg A) = 1 - P(A)$ = probability that formula $\neg A$ is true

From now on A is a propositional variable with values in $\{0, 1\}$, i.e. $P(A)$ is a function of two possible input values $A = 1$ and $A = 0$, i.e. with slightly unusual notation:

$P(A = 1)$ = probability that formula A is true

$P(A = 0) = 1 - P(A = 1)$ = probability that formula A is false

Stating that $P(A, B) = P(A) \cdot P(B)$ means **all** of the following

$$P(A = 1, B = 1) = P(A = 1) \cdot P(B = 1) \quad P(A = 1, B = 0) = P(A = 1) \cdot P(B = 0)$$

$$P(A = 0, B = 1) = P(A = 0) \cdot P(B = 1) \quad P(A = 0, B = 0) = P(A = 0) \cdot P(B = 0)$$





Definition (independence)

Two variables A and B are **independent**, if and only if their joint distributions factorizes into so-called **marginal distributions**, i.e.

$$P(A, B) = P(A) P(B)$$

In that case $P(A|B) = P(A)$. Notation: $A \perp\!\!\!\perp B$. Information about B does not give information about A and vice versa.



Independence

Chiefly a computational concept

Definition (independence)

Two variables A and B are *independent*, if and only if their joint distributions factorizes into so-called marginal distributions, i.e.

$$P(A, B) = P(A) P(B)$$

In that case $P(A|B) = P(A)$. Notation: $A \perp\!\!\!\perp B$. Information about B does not give information about A and vice versa.

Example: Two coins.

A = coin 1 shows heads

B = coin 2 shows heads



Then $A \perp\!\!\!\perp B$.



Conditional Independence

Chiefly a computational concept



Definition (conditional independence)

Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$

In that case we have $P(A|B, C) = P(A|C)$, i.e. in light of information C , B provides no (further) information about A . Notation: $A \perp\!\!\!\perp B \mid C$



Conditional Independence

Chiefly a computational concept

[Example: Stefan Harmeling]

Definition (conditional independence)

Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$

In that case we have $P(A|B, C) = P(A|C)$, i.e. in light of information C , B provides no (further) information about A . Notation: $A \perp\!\!\!\perp B \mid C$

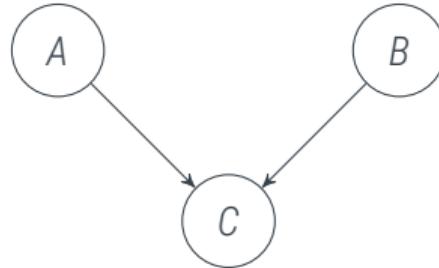
Example: Two coins and a bell.

A = coin 1 shows heads

B = coin 2 shows heads

C = bell rings if both coins show the same result

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





Conditional Independence

Chiefly a computational concept

[Example: Stefan Harmeling]

Definition (conditional independence)

Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$

In that case we have $P(A|B, C) = P(A|C)$, i.e. in light of information C , B provides no (further) information about A . Notation: $A \perp\!\!\!\perp B \mid C$

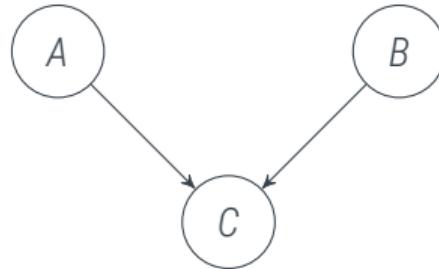
Example: Two coins and a bell.

A = coin 1 shows heads

B = coin 2 shows heads

C = bell rings if both coins show the same result

$A \perp\!\!\!\perp B$ and $A \perp\!\!\!\perp C$





Conditional Independence

Chiefly a computational concept

[Example: Stefan Harmeling]

Definition (conditional independence)

Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$

In that case we have $P(A|B, C) = P(A|C)$, i.e. in light of information C , B provides no (further) information about A . Notation: $A \perp\!\!\!\perp B \mid C$

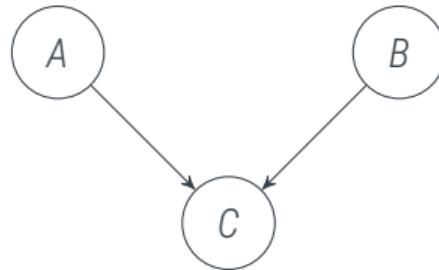
Example: Two coins and a bell.

A = coin 1 shows heads

B = coin 2 shows heads

C = bell rings if both coins show the same result

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





Conditional Independence

Chiefly a computational concept

[Example: Stefan Harmeling]

Definition (conditional independence)

Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$

In that case we have $P(A|B, C) = P(A|C)$, i.e. in light of information C , B provides no (further) information about A . Notation: $A \perp\!\!\!\perp B | C$

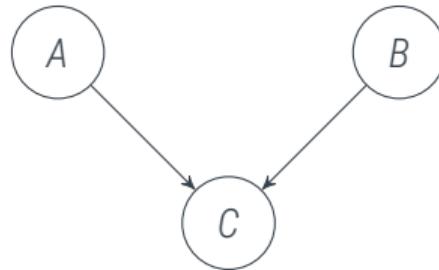
Example: Two coins and a bell.

A = coin 1 shows heads

B = coin 2 shows heads

C = bell rings if both coins show the same result

$A \perp\!\!\!\perp B$ and $A \perp\!\!\!\perp C$ and $B \perp\!\!\!\perp C$, but $A \not\perp\!\!\!\perp B | C$





Conditional Independence

Chiefly a computational concept

[Example: Stefan Harmeling]

Definition (conditional independence)

Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$

In that case we have $P(A|B, C) = P(A|C)$, i.e. in light of information C , B provides no (further) information about A . Notation: $A \perp\!\!\!\perp B | C$

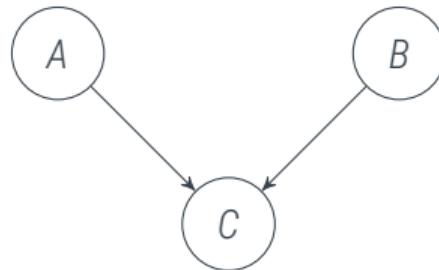
Example: Two coins and a bell.

A = coin 1 shows heads

B = coin 2 shows heads

C = bell rings if both coins show the same result

$A \perp\!\!\!\perp B$ and $A \perp\!\!\!\perp C$ and $B \perp\!\!\!\perp C$, but $A \not\perp\!\!\!\perp B | C$ and $A \not\perp\!\!\!\perp C | B$ and $B \not\perp\!\!\!\perp C | A$.





Computing with Probabilities

- ▶ Probabilistic reasoning extends propositional logic
- ▶ instead of tracking a single *true* value, we have to assign probabilities to *combinatorially many* hypotheses
- ▶ Two variables A and B are **conditionally independent** given variable C , if and only if their conditional distribution factorizes,

$$P(A, B|C) = P(A|C) P(B|C)$$





Parameter Counting

a simple example

[adapted from Pearl, 1988 / MacKay, 2003 §21]



A = the alarm was triggered



E = there was an earthquake



B = there was a break-in



R = an announcement is made on the radio

Joint probability distribution has $2^4 - 1 = 15 = 8 + 4 + 2 + 1$ parameters

$$P(A, E, B, R) = P(A \mid R, E, B) \cdot P(R \mid E, B) \cdot P(E \mid B) \cdot P(B).$$

Removing irrelevant conditions (domain knowledge!) reduces to $8 = 4 + 2 + 1 + 1$ parameters:

$$P(A, E, B, R) = P(A \mid E, B) \cdot P(R \mid E) \cdot P(E) \cdot P(B)$$

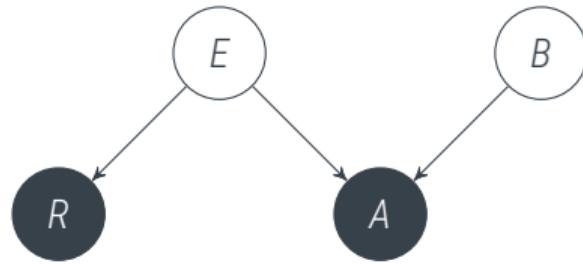


A Graphical Representation

Our first Bayesian network.

[adapted from Pearl, 1988 / MacKay, 2003 §21]

$$P(A, E, B, R) = P(A | E, B) \cdot P(R | E) \cdot P(E) \cdot P(B)$$



A = the alarm was triggered

E = there was an earthquake

B = there was a break-in

R = an announcement is made on the radio



Conditional Probability Tables

For the burglar/alarm problem

$$P(A, E, B, R) = P(A | E, B) \cdot P(R | E) \cdot P(E) \cdot P(B)$$

$$P(B = 1) = 10^{-3}$$

$$P(E = 1) = 10^{-3}$$

$$P(B = 0) = 1 - 10^{-3}$$

$$P(E = 0) = 1 - 10^{-3}$$

and

$$P(R = 1 | E = 1) = 1.0$$

$$P(R = 1 | E = 0) = 0.0$$

and, using $f = 10^{-3}$, $\alpha_b = 0.99$, $\alpha_e = 0.01$

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

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

$$P(A = 0 | B = 1, E = 0) = (1 - f)(1 - \alpha_b)$$

$$P(A = 1 | B = 1, E = 0) = 1 - (1 - f)(1 - \alpha_b)$$

$$P(A = 0 | B = 0, E = 1) = (1 - f)(1 - \alpha_e)$$

$$P(A = 1 | B = 0, E = 1) = 1 - (1 - f)(1 - \alpha_e)$$

$$P(A = 0 | B = 1, E = 1) = (1 - f)(1 - \alpha_b)(1 - \alpha_e)$$

$$P(A = 1 | B = 1, E = 1) = 1 - (1 - f)(1 - \alpha_b)(1 - \alpha_e)$$



Conditional Probability Tables

For the burglar/alarm problem

$$P(A, E, B, R) = P(A | E, B) \cdot P(R | E) \cdot P(E) \cdot P(B)$$

$$P(B = 1) = 10^{-3}$$

$$P(E = 1) = 10^{-3}$$

$$P(B = 0) = 1 - 10^{-3}$$

$$P(E = 0) = 1 - 10^{-3}$$

and

$$P(R = 1 | E = 1) = 1.0$$

$$P(R = 1 | E = 0) = 0.0$$

and, using $f = 10^{-3}$, $\alpha_b = 0.99$, $\alpha_e = 0.01$

$$P(A = 0 | B = 0, E = 0) = 0.999$$

$$P(A = 1 | B = 0, E = 0) = 0.001$$

$$P(A = 0 | B = 1, E = 0) = 0.00999$$

$$P(A = 1 | B = 1, E = 0) = 0.99001$$

$$P(A = 0 | B = 0, E = 1) = 0.98901$$

$$P(A = 1 | B = 0, E = 1) = 0.01099$$

$$P(A = 0 | B = 1, E = 1) = 0.0098901$$

$$P(A = 1 | B = 1, E = 1) = 0.9901099$$



What is the probability that there was a break-in and/or an earthquake, given that the alarm went off?

$$P(B, E | A = 1) = \frac{P(A = 1 | B, E)P(B)P(E)}{P(A = 1)}$$

$$\begin{aligned} P(A = 1) &= P(A = 1 | B = 0, E = 0)P(B = 0)P(E = 0) \\ &\quad + P(A = 1 | B = 0, E = 1)P(B = 0)P(E = 1) \\ &\quad + P(A = 1 | B = 1, E = 0)P(B = 1)P(E = 0) \\ &\quad + P(A = 1 | B = 1, E = 1)P(B = 1)P(E = 1) \\ &= 0.000\,998 + 0.000\,989 + 0.000\,010\,979 + 0.000\,000\,99 = 0.002 \end{aligned}$$

thus – note conditional dependence!

$$P(B = 0, E = 0 | A = 1) = 0.4993 \quad P(B = 1, E = 0 | A = 1) = 0.4947$$

$$P(B = 0, E = 1 | A = 1) = 0.0055 \quad P(B = 1, E = 1 | A = 1) = 0.0005$$

$$P(B = 0 | A = 1) = P(B = 0, E = 0 | A = 1) + P(B = 0, E = 1 | A = 1) = 0.505$$

$$P(B = 1 | A = 1) = P(B = 1, E = 0 | A = 1) + P(B = 1, E = 1 | A = 1) = 0.495$$



What is the probability for a break-in, given alarm *and* radio announcement?

$$\begin{aligned} P(B = 0 \mid E = 1, A = 1) &= \frac{P(B = 0, E = 1 \mid A = 1)}{P(E = 1 \mid A = 1)} \\ &= \frac{P(B = 0, E = 1 \mid A = 1)}{P(B = 0, E = 1 \mid A = 1) + P(B = 1, E = 1 \mid A = 1)} = 0.92 \\ P(B = 1 \mid E = 1, A = 1) &= \frac{P(B = 1, E = 1 \mid A = 1)}{P(E = 1 \mid A = 1)} \\ &= \frac{P(B = 1, E = 1 \mid A = 1)}{P(B = 0, E = 1 \mid A = 1) + P(B = 1, E = 1 \mid A = 1)} = 0.08 \end{aligned}$$

The radio announcement is **explaining away** the break-in as the explanation for the alarm.

What is Probabilistic Reasoning?

One recipe for all your inference needs!



Always write down the probability of everything.

David JC MacKay (1967–2016)

- ▶ identify all relevant variables: A, R, E, B
- ▶ define **joint probability** $P(A, R, E, B)$ aka. the generative model
- ▶ **observations** fix certain variables: $A = 1$
- ▶ **inference** takes place exclusively by Bayes' Theorem
n.b.: this requires integrating out (**marginalizing**) latent variables not being inferred.



Directed Graphical Models

aka. Bayesian networks, Bayes nets, belief networks, ...



[Judea Pearl, *Probabilistic Reasoning in Intelligent Systems*, 1988]

Definition (Bayesian Network, preliminary definition – more in later lectures)

A **Directed Graphical Model (DGM)**, aka. **Bayesian Network** is a probability distribution over variables $\{X_1, \dots, X_D\}$ that can be written as

$$P(X_1, X_2, \dots, X_D) = \prod_{i=1}^D p(X_i | \text{pa}(X_i))$$

where $\text{pa}(X_i)$ are the parental variables of X_i , that is, $X_i \not\in \text{pa}(X_j) \forall X_j \in \text{pa}(X_i)$. A DGM can be represented by a **Directed Acyclic Graph (DAG)** with the propositional variables as nodes, and arrows from parents to children.

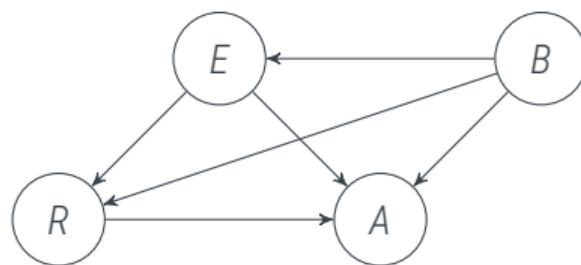
Every Probability Distribution is a DAG

It's just not always a helpful concept



By the Product Rule, every joint can be factorized into a (dense) DAG.

$$P(A, E, B, R) = P(A | E, B, R) \cdot P(R | E, B) \cdot P(E | B) \cdot P(B)$$



A = the alarm was triggered

E = there was an earthquake

B = there was a break-in

R = an announcement is made on the radio

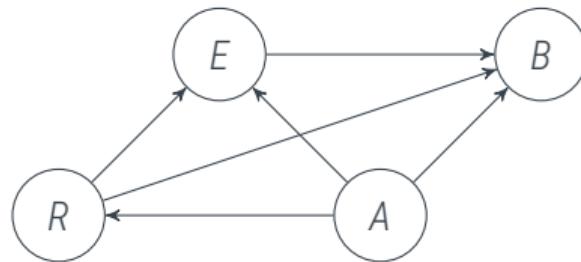
Every Probability Distribution is a DAG

It's just not always a helpful concept



The direction of the arrows is **not** a causal statement.

$$P(A, E, B, R) = P(B | A, E, R) \cdot P(E | A, R) \cdot P(R | A) \cdot P(A)$$



A = the alarm was triggered

E = there was an earthquake

B = there was a break-in

R = an announcement is made on the radio

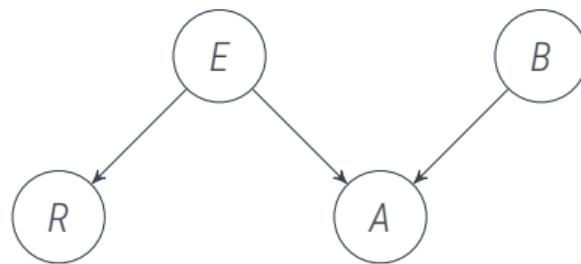
Every Probability Distribution is a DAG

It's just not always a helpful concept



But the representation is particularly interesting when it reveals **independence**.

$$P(A, E, B, R) = P(A | E, B) \cdot P(R | E) \cdot P(E) \cdot P(B)$$



A = the alarm was triggered

E = there was an earthquake

B = there was a break-in

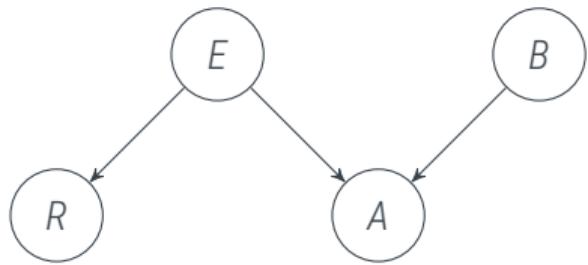
R = an announcement is made on the radio



Deducing Conditional Independencies

back to our example

[adapted from Pearl, 1988 / MacKay, 2003 §21]



$$P(A, E, B, R) = P(A | E, B) \cdot P(R | E) \cdot P(E) \cdot P(B)$$

A = the alarm was triggered

E = there was an earthquake

B = there was a break-in

R = an announcement is made on the radio

Which independencies can we infer only from the graph?



Atomic Independence Structures

DAGs imply conditional independence, but not dependence!

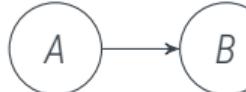
For uni- and bi-variate graphs, conditional independence is trivial.

For tri-variate sub-graphs, there are three possible structures:

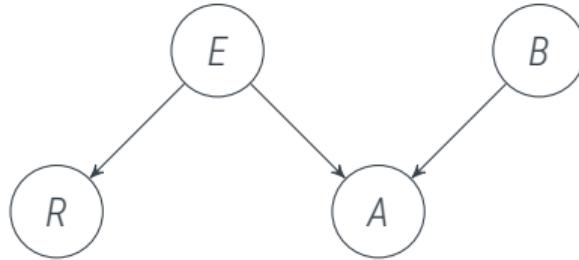
graph	factorization	implications
(i)	$P(A, B, C) = P(C B) \cdot P(B A) \cdot P(A)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(ii)	$P(A, B, C) = P(A B) \cdot P(C B) \cdot P(B)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(iii)	$P(A, B, C) = P(B A, C) \cdot P(C) \cdot P(A)$	$A \perp\!\!\!\perp C$ but not, i.g., $A \not\perp\!\!\!\perp C B$

Deducing Conditional Independencies

back to our example

	graph	factorization	implications
(i)		$P(A, B, C) = P(C B) \cdot P(B A) \cdot P(A)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(ii)		$P(A, B, C) = P(A B) \cdot P(C B) \cdot P(B)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(iii)		$P(A, B, C) = P(B A, C) \cdot P(C) \cdot P(A)$	$A \perp\!\!\!\perp C$ but not, i.g., $A \not\perp\!\!\!\perp C B$

Which independencies can we infer only from the graph?



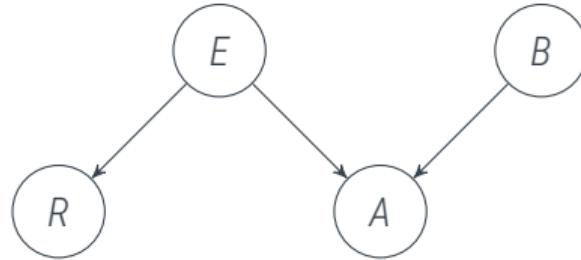


Deducing Conditional Independencies

back to our example

	graph	factorization	implications
(i)		$P(A, B, C) = P(C B) \cdot P(B A) \cdot P(A)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(ii)		$P(A, B, C) = P(A B) \cdot P(C B) \cdot P(B)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(iii)		$P(A, B, C) = P(B A, C) \cdot P(C) \cdot P(A)$	$A \perp\!\!\!\perp C$ but not, i.g., $A \not\perp\!\!\!\perp C B$

Which independencies can we infer only from the graph?



► $R \perp\!\!\!\perp A | E$ and $E \perp\!\!\!\perp B$

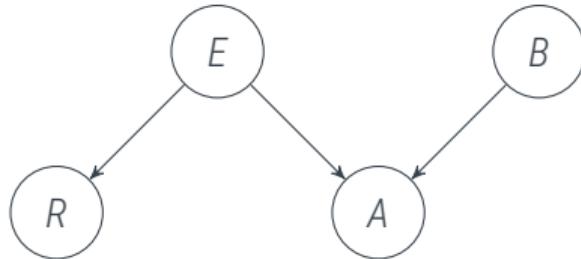


Deducing Conditional Independencies

back to our example

graph	factorization	implications
(i)	$P(A, B, C) = P(C B) \cdot P(B A) \cdot P(A)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(ii)	$P(A, B, C) = P(A B) \cdot P(C B) \cdot P(B)$	$A \perp\!\!\!\perp C B$ but not, i.g., $A \not\perp\!\!\!\perp C$
(iii)	$P(A, B, C) = P(B A, C) \cdot P(C) \cdot P(A)$	$A \perp\!\!\!\perp C$ but not, i.g., $A \not\perp\!\!\!\perp C B$

Which independencies can we infer only from the graph?



- ▶ $R \perp\!\!\!\perp A | E$ and $E \perp\!\!\!\perp B$
- ▶ but also $(R \perp\!\!\!\perp B | E)$, $(R \perp\!\!\!\perp B)$, $(R \perp\!\!\!\perp B | E, A)$, with more work



The Graph for Two Coins and a Bell

DAGs are not a perfect tool

$$P(A = 1) = 0.5$$

$$P(B = 1) = 0.5$$

$$P(C = 1 \mid A = 1, B = 1) = 1$$

$$P(C = 1 \mid A = 0, B = 1) = 0$$

$$P(C = 1 \mid A = 1, B = 0) = 0$$

$$P(C = 1 \mid A = 0, B = 0) = 1$$

These CPTs imply $P(A|B) = P(A)$, $P(B|C) = P(B)$ and $P(C|A) = P(C)$ and $P(C \mid B) = P(C)$.



The Graph for Two Coins and a Bell

DAGs are not a perfect tool

$$P(A = 1) = 0.5$$

$$P(C = 1 \mid A = 1, B = 1) = 1$$

$$P(C = 1 \mid A = 1, B = 0) = 0$$

$$P(B = 1) = 0.5$$

$$P(C = 1 \mid A = 0, B = 1) = 0$$

$$P(C = 1 \mid A = 0, B = 0) = 1$$

These CPTs imply $P(A|B) = P(A)$, $P(B|C) = P(B)$ and $P(C|A) = P(C)$ and $P(C \mid B) = P(C)$.

We thus have three factorizations:

1. $P(A, B, C) = P(C|A, B) \cdot P(A|B) \cdot P(B) = P(C|A, B) \cdot P(A) \cdot P(B)$
2. $P(A, B, C) = P(A|B, C) \cdot P(B|C) \cdot P(C) = P(A|B, C) \cdot P(B) \cdot P(C)$
3. $P(A, B, C) = P(B|C, A) \cdot P(C|A) \cdot P(A) = P(B|C, A) \cdot P(C) \cdot P(A)$

The Graph for Two Coins and a Bell

DAGs are not a perfect tool

$$P(A = 1) = 0.5$$

$$P(C = 1 | A = 1, B = 1) = 1$$

$$P(C = 1 | A = 1, B = 0) = 0$$

$$P(B = 1) = 0.5$$

$$P(C = 1 | A = 0, B = 1) = 0$$

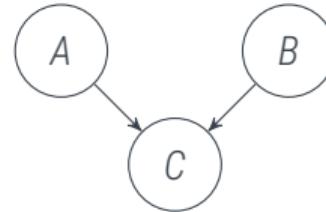
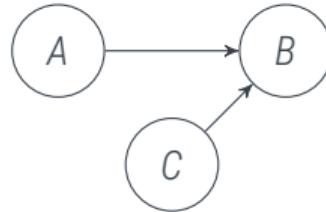
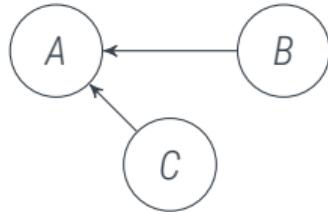
$$P(C = 1 | A = 0, B = 0) = 1$$

These CPTs imply $P(A|B) = P(A)$, $P(B|C) = P(B)$ and $P(C|A) = P(C)$ and $P(C | B) = P(C)$.

We thus have three factorizations:

1. $P(A, B, C) = P(C|A, B) \cdot P(A|B) \cdot P(B) = P(C|A, B) \cdot P(A) \cdot P(B)$
2. $P(A, B, C) = P(A|B, C) \cdot P(B|C) \cdot P(C) = P(A|B, C) \cdot P(B) \cdot P(C)$
3. $P(A, B, C) = P(B|C, A) \cdot P(C|A) \cdot P(A) = P(B|C, A) \cdot P(C) \cdot P(A)$

Each corresponds to a graph. Note that each can only express some of the independencies:





Graphical Models and Conditional Independence

- ▶ Multivariate distributions can have **exponentially** many degrees of freedom.
- ▶ **(conditional) independence** helps reduce this complexity to make things tractable
- ▶ **(directed) graphical models** provide a notation from which conditional independence can be read off using simple rules.
- ▶ Every probability distribution is a DAG, but not every independence structure of a distribution is captured by a DAG of it.
- ▶ We will return to graphs later in the course.

Conditional independence is a tool (and may be required)
to keep inference tractable in multi-variate problems.





§ 5. Unabhängigkeit.

Der Begriff der gegenseitigen *Unabhängigkeit* zweier oder mehrerer Versuche nimmt eine in gewissem Sinne zentrale Stellung in der Wahrscheinlichkeitsrechnung ein. In der Tat haben wir schon gesehen, daß die Wahrscheinlichkeitsrechnung vom mathematischen Standpunkte aus als eine spezielle Anwendung der allgemeinen Theorie der additiven Mengenfunktionen betrachtet werden kann. Man kann sich natürlich fragen, wie ist es dann möglich, daß die Wahrscheinlichkeitsrechnung sich in eine große, ihre eigenen Methoden besitzende selbständige Wissenschaft entwickelt hat?

Geschichtlich ist die Unabhängigkeit von Versuchen und zufälligen Größen derjenige mathematische Begriff, welcher der Wahrscheinlichkeitsrechnung ihr eigenartiges Gepräge gibt. Die klassischen Arbeiten von LAPLACE, POISSON, TCHEBYCHEFF, MARKOFF, LIAPOUNOFF, v. MISES und BERNSTEIN sind in der Tat im wesentlichen der Untersuchung von Reihen unabhängiger Größen gewidmet. Wenn man in den neueren Untersuchungen (MARKOFF, BERNSTEIN usw.) öfters die Forderung der vollständigen Unabhängigkeit ablehnt, so sieht man sich immer gezwungen, um hinreichend inhaltreiche Resultate zu erhalten, abgeschwächte analoge Forderungen einzuführen. (Vgl. in diesem Kap. § 6 — MARKOFFsche Ketten.)



Man kommt also dazu, im Begriffe der Unabhängigkeit wenigstens den ersten Keim der eigenartigen Problematik der Wahrscheinlichkeitsrechnung zu erblicken. [...] Es ist dementsprechend eine der wichtigsten Aufgaben der Philosophie der Naturwissenschaften, nachdem sie die vielumstrittene Frage über das Wesen des Wahrscheinlichkeitsbegriffes selbst erklärt hat, die Voraussetzungen zu präzisieren. bei denen man irgendwelche gegebene reelle Erscheinungen für gegenseitig unabhängig halten kann.

A. N. Kolmogorov. Grundbegriffe der Wahrscheinlichkeitsrechnung. §I.5

