

# **Fundamentals of Artificial Intelligence and Knowledge Representation (Module 3)**

Last update: 08 December 2023

Academic Year 2023 – 2024  
Alma Mater Studiorum · University of Bologna

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Uncertainty . . . . .	1
1.1.1	Handling uncertainty . . . . .	1
<b>2</b>	<b>Probability</b>	<b>2</b>
2.1	Inference with full joint distributions . . . . .	3
<b>3</b>	<b>Bayesian networks</b>	<b>6</b>
3.1	Bayes' rule . . . . .	6
3.2	Bayesian network reasoning . . . . .	6
3.3	Building Bayesian networks . . . . .	10
3.3.1	Algorithm . . . . .	10
3.3.2	Structure learning . . . . .	11
3.4	Causal networks . . . . .	11
3.5	Compact conditional distributions . . . . .	12
3.5.1	Noisy-OR . . . . .	12
3.5.2	Hybrid Bayesian networks . . . . .	13
3.5.3	Other methods . . . . .	14
<b>4</b>	<b>Exact inference</b>	<b>15</b>
4.1	Inference by enumeration . . . . .	15
4.2	Inference by variable elimination . . . . .	15
4.2.1	Irrelevant variables . . . . .	17
4.2.2	Complexity . . . . .	17
4.3	Clustering algorithm . . . . .	17
<b>5</b>	<b>Approximate inference</b>	<b>18</b>
5.1	Sampling from an empty network . . . . .	18
5.2	Rejection sampling . . . . .	19
5.3	Likelihood weighting . . . . .	19
5.4	Markov chain Monte Carlo . . . . .	20

# 1 Introduction

## 1.1 Uncertainty

**Uncertainty** A task is uncertain if it has:

Uncertainty

- Partial observations
- Noisy or wrong information
- Uncertain outcomes of the actions
- Complex models

A purely logic approach leads to:

- Risks falsehood: unreasonable conclusion when applied in practice.
- Weak decisions: too many conditions required to make a conclusion.

### 1.1.1 Handling uncertainty

**Default/non-monotonic logic** Works on assumptions. An assumption can be contradicted by the evidence.

Default/non-monotonic logic

**Rule-based systems with fudge factors** Formulated as premise  $\rightarrow_{\text{prob.}}$  effect. Have the following issues:

Rule-based systems with fudge factors

- Locality: how can the probability account all the evidence.
- Combination: chaining of unrelated concepts.

**Probability** Assign a probability given the available known evidence.

Probability

Note: fuzzy logic handles the degree of truth and not the uncertainty.

**Decision theory** Defined as:

Decision theory

Decision theory = Utility theory + Probability theory

where the utility theory depends on one's preferences.

## 2 Probability

**Sample space** Set  $\Omega$  of all possible worlds.

Sample space

**Event** Subset  $A \subseteq \Omega$ .

Event

**Sample point/Possible world/Atomic event** Element  $\omega \in \Omega$ .

Sample point

**Probability space** A probability space/model is a function  $\mathcal{P}(\cdot) : \Omega \rightarrow [0, 1]$  assigned to a sample space such that:

Probability space

- $0 \leq \mathcal{P}(\omega) \leq 1$
- $\sum_{\omega \in \Omega} \mathcal{P}(\omega) = 1$
- $\mathcal{P}(A) = \sum_{\omega \in A} \mathcal{P}(\omega)$

**Random variable** A function from an event to some range (e.g. reals, booleans, ...).

Random variable

**Probability distribution** For any random variable  $X$ :

Probability distribution

$$\mathcal{P}(X = x_i) = \sum_{\omega \text{ s.t. } X(\omega) = x_i} \mathcal{P}(\omega)$$

**Proposition** Event where a random variable has a certain value.

Proposition

$$a = \{\omega \mid A(\omega) = \text{true}\}$$

$$\neg a = \{\omega \mid A(\omega) = \text{false}\}$$

$$(\text{Weather} = \text{rain}) = \{\omega \mid B(\omega) = \text{rain}\}$$

**Prior probability** Prior/unconditional probability of a proposition based on known evidence.

Prior probability

**Probability distribution (all)** Gives all the probabilities of a random variable.

Probability distribution (all)

$$\mathbf{P}(A) = \langle \mathcal{P}(A = a_1), \dots, \mathcal{P}(A = a_n) \rangle$$

**Joint probability distribution** The joint probability distribution of a set of random variables gives the probability of all the different combinations of their atomic events.

Joint probability distribution

Note: Every question on a domain can, in theory, be answered using the joint distribution. In practice, it is hard to apply.

**Example.**  $\mathbf{P}(\text{Weather}, \text{Cavity}) =$

	Weather=sunny	Weather=rain	Weather=cloudy	Weather=snow
Cavity=true	0.144	0.02	0.016	0.02
Cavity=false	0.576	0.08	0.064	0.08

**Probability density function** The probability density function (PDF) of a random variable  $X$  is a function  $p : \mathbb{R} \rightarrow \mathbb{R}$  such that:

Probability density function

$$\int_{\mathcal{T}_X} p(x) dx = 1$$

## Uniform distribution

Uniform distribution

$$p(x) = \text{Unif}[a, b](x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

## Gaussian (normal) distribution

Gaussian (normal) distribution

$$\mathcal{N}(\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\mathcal{N}(0, 1)$  is the standard Gaussian.

**Conditional probability** Probability of a prior knowledge with new evidence:

Conditional probability

$$\mathcal{P}(a|b) = \frac{\mathcal{P}(a \wedge b)}{\mathcal{P}(b)}$$

The product rule gives an alternative formulation:

$$\mathcal{P}(a \wedge b) = \mathcal{P}(a|b)\mathcal{P}(b) = \mathcal{P}(b|a)\mathcal{P}(a)$$

**Chain rule** Successive application of the product rule:

Chain rule

$$\begin{aligned} \mathbf{P}(X_1, \dots, X_n) &= \mathbf{P}(X_1, \dots, X_{n-1})\mathbf{P}(X_n|X_1, \dots, X_{n-1}) \\ &= \mathbf{P}(X_1, \dots, X_{n-2})\mathbf{P}(X_{n-1}|X_1, \dots, X_{n-2})\mathbf{P}(X_n|X_1, \dots, X_{n-1}) \\ &= \prod_{i=1}^n \mathbf{P}(X_i|X_1, \dots, X_{i-1}) \end{aligned}$$

**Independence** Two random variables  $A$  and  $B$  are independent ( $A \perp B$ ) iff:

Independence

$$\mathbf{P}(A|B) = \mathbf{P}(A) \text{ or } \mathbf{P}(B|A) = \mathbf{P}(B) \text{ or } \mathbf{P}(A, B) = \mathbf{P}(A)\mathbf{P}(B)$$

**Conditional independence** Two random variables  $A$  and  $B$  are conditionally independent iff:

Conditional independence

$$\mathbf{P}(A|C, B) = \mathbf{P}(A|C)$$

## 2.1 Inference with full joint distributions

Given a joint distribution, the probability of any proposition  $\phi$  can be computed as the sum of the atomic events where  $\phi$  is true:

$$\mathcal{P}(\phi) = \sum_{\omega: \omega \models \phi} \mathcal{P}(\omega)$$

**Example.** Given the following joint distribution:

	toothache		$\neg$ toothache	
	catch	$\neg$ catch	catch	$\neg$ catch
cavity	0.108	0.012	0.072	0.008
$\neg$ cavity	0.016	0.064	0.144	0.576

We have that:

- $\mathcal{P}(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$
- $\mathcal{P}(\text{cavity} \vee \text{toothache}) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$
- $\mathcal{P}(\neg \text{cavity} | \text{toothache}) = \frac{\mathcal{P}(\neg \text{cavity} \wedge \text{toothache})}{\mathcal{P}(\text{toothache})} = \frac{0.016 + 0.064}{0.2} = 0.4$

**Marginalization** The probability that a random variable assumes a specific value is given by the sum off all the joint probabilities where that random variable assumes the given value.

Marginalization

**Example.** Given the joint distribution:

	Weather=sunny	Weather=rain	Weather=cloudy	Weather=snow
Cavity=true	0.144	0.02	0.016	0.02
Cavity=false	0.576	0.08	0.064	0.08

We have that  $\mathcal{P}(\text{Weather} = \text{sunny}) = 0.144 + 0.576$

**Conditioning** Adding a condition to a probability (reduction and renormalization).

Conditioning

**Normalization** Given a conditional probability distribution  $\mathbf{P}(A|B)$ , it can be formulated as:

Normalization

$$\mathbf{P}(A|B) = \alpha \mathbf{P}(A, B)$$

where  $\alpha$  is a normalization constant. In fact, fixed the evidence  $B$ , the denominator to compute the conditional probability is the same for each probability.

**Example.** Given the joint distribution:

	toothache		$\neg$ toothache	
	catch	$\neg$ catch	catch	$\neg$ catch
cavity	0.108	0.012	0.072	0.008
$\neg$ cavity	0.016	0.064	0.144	0.576

We have that:

$$\mathbf{P}(\text{Cavity} | \text{toothache}) = \left\langle \frac{\mathcal{P}(\text{cavity}, \text{toothache}, \text{catch})}{\mathcal{P}(\text{toothache})}, \frac{\mathcal{P}(\neg \text{cavity}, \text{toothache}, \neg \text{catch})}{\mathcal{P}(\text{toothache})} \right\rangle$$

**Probability query** Given a set of query variables  $\mathbf{Y}$ , the evidence variables  $\mathbf{e}$  and the other hidden variables  $\mathbf{H}$ , the probability of the query can be computed as:

Probability query

$$\mathbf{P}(\mathbf{Y} | \mathbf{E} = \mathbf{e}) = \alpha \mathbf{P}(\mathbf{Y}, \mathbf{E} = \mathbf{e}) = \alpha \sum_{\mathbf{h}} \mathbf{P}(\mathbf{Y}, \mathbf{E} = \mathbf{e}, \mathbf{H} = \mathbf{h})$$

The problem of this approach is that it has exponential time and space complexity that makes it not applicable in practice.

To reduce the size of the variables, conditional independence can be exploited.

**Example.** Knowing that  $\mathbf{P} \models (\text{Catch} \perp \text{Toothache} | \text{Cavity})$ , we can compute the distribution  $\mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity})$  as follows:

$$\begin{aligned} \mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity}) &= \\ &= \mathbf{P}(\text{Toothache} | \text{Catch}, \text{Cavity}) \mathbf{P}(\text{Catch} | \text{Cavity}) \mathbf{P}(\text{Cavity}) \\ &= \mathbf{P}(\text{Toothache} | \text{Cavity}) \mathbf{P}(\text{Catch} | \text{Cavity}) \mathbf{P}(\text{Cavity}) \end{aligned}$$

$\mathbf{P}(\text{Toothache}, \text{Catch}, \text{Cavity})$  has 7 independent values that grows exponentially ( $2 \cdot 2 \cdot 2 = 8$  values, but one of them can be omitted as a probability always sums up to 1).

$\mathbf{P}(\text{Toothache} | \text{Cavity})\mathbf{P}(\text{Catch} | \text{Cavity})\mathbf{P}(\text{Cavity})$  has 5 independent values that grows linearly ( $4 + 4 + 2 = 10$ , but a value of  $\mathbf{P}(\text{Cavity})$  can be omitted. The conditional probabilities require two tables (one for each prior) each with 2 values, but for each table a value can be omitted, therefore requiring 2 independent values per conditional probability instead of 4).

## 3 Bayesian networks

### 3.1 Bayes' rule

Bayes' rule

Bayes' rule

$$\mathcal{P}(a|b) = \frac{\mathcal{P}(b|a)\mathcal{P}(a)}{\mathcal{P}(b)}$$

**Bayes' rule and conditional independence** Given the random variables **Cause** and  $\text{Effect}_1, \dots, \text{Effect}_n$ , with  $\text{Effect}_i$  independent from each other, we can compute  $\mathbf{P}(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n)$  as follows:

$$\mathbf{P}(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = \left( \prod_i \mathbf{P}(\text{Effect}_i | \text{Cause}) \right) \mathbf{P}(\text{Cause})$$

The number of parameters is linear.

**Example.** Knowing that  $\mathbf{P} \models (\text{Catch} \perp \text{Toothache} | \text{Cavity})$ :

$$\begin{aligned} & \mathbf{P}(\text{Cavity} | \text{toothache} \wedge \text{catch}) \\ &= \alpha \mathbf{P}(\text{toothache} \wedge \text{catch} | \text{Cavity}) \mathbf{P}(\text{Cavity}) \\ &= \alpha \mathbf{P}(\text{toothache} | \text{Cavity}) \mathbf{P}(\text{catch} | \text{Cavity}) \mathbf{P}(\text{Cavity}) \end{aligned}$$

### 3.2 Bayesian network reasoning

**Bayesian network** Graph for conditional independence assertions and a compact specification of full joint distributions.

Bayesian network

- Directed acyclic graph.
- Nodes represent variables.
- The conditional distribution of a node is given by its parents

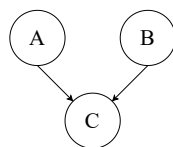
$$\mathbf{P}(X_i | \text{parents}(X_i))$$

In other words, if there is an edge from  $A$  to  $B$ , then  $A$  (cause) influences  $B$  (effect).

**Conditional probability table (CPT)** In the case of boolean variables, the conditional distribution of a node can be represented using a table by considering all the combinations of the parents.

Conditional probability table (CPT)

**Example.** Given the boolean variables  $A$ ,  $B$  and  $C$ , with  $C$  depending on  $A$  and  $B$ , we have that:



A	B	$\mathcal{P}(c A, B)$	$\mathcal{P}(\neg c A, B)$
a	b	$\alpha$	$1 - \alpha$
$\neg a$	b	$\beta$	$1 - \beta$
a	$\neg b$	$\gamma$	$1 - \gamma$
$\neg a$	$\neg b$	$\delta$	$1 - \delta$



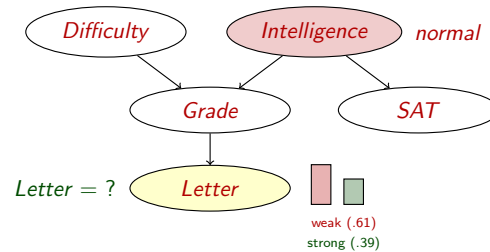
**Reasoning patterns** Given a Bayesian network, the following reasoning patterns can be used:

Reasoning patterns

**Causal** To make a prediction. From the cause, derive the effect.

Causal reasoning

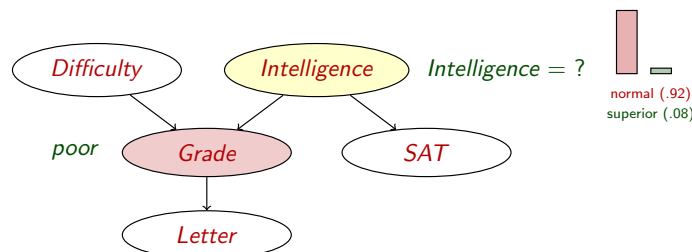
**Example.** Knowing **Intelligence**, it is possible to make a prediction of **Letter**.



**Evidential** To find an explanation. From the effect, derive the cause.

Evidential reasoning

**Example.** Knowing **Grade**, it is possible to explain it by estimating **Intelligence**.

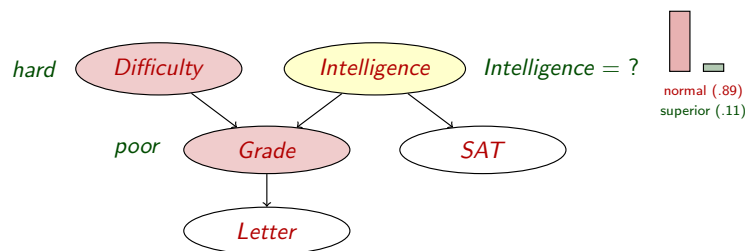


**Explain away** Observation obtained "passing through" other observations.

Explain away reasoning

**Example.** Knowing **Difficulty** and **Grade**, it is possible to estimate **Intelligence**.

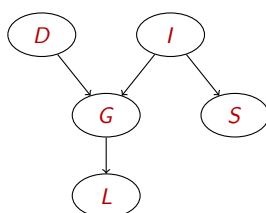
Note that if **Grade** was not known, **Difficulty** and **Intelligence** would have been independent.



**Independence** Intuitively, an effect is independent from a cause, if there is another cause in the middle whose value is already known.

Bayesian network independence

**Example.**



$$\mathbf{P} \models (L \perp D, I, S \mid G)$$

$$\mathbf{P} \models (S \perp L \mid G)$$

$$\mathbf{P} \models (S \perp D) \text{ but } \mathbf{P} \models (S \not\perp D \mid G) \text{ (explain away)}$$

**V-structure** Effect with two causes. If the effect is not in the evidence, the causes are independent. V-structure

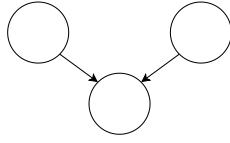


Figure 3.1: V-structure

**Active two-edge trail** The trail  $X \rightleftharpoons Z \rightleftharpoons Y$  is active either if:

Active two-edge trail

- $X, Z, Y$  is a v-structure  $X \rightarrow Z \leftarrow Y$  and  $Z$  or one of its children is in the evidence.
- $Z$  is not in the evidence.

In other words, influence can flow from  $X$  to  $Y$  passing by  $Z$ .

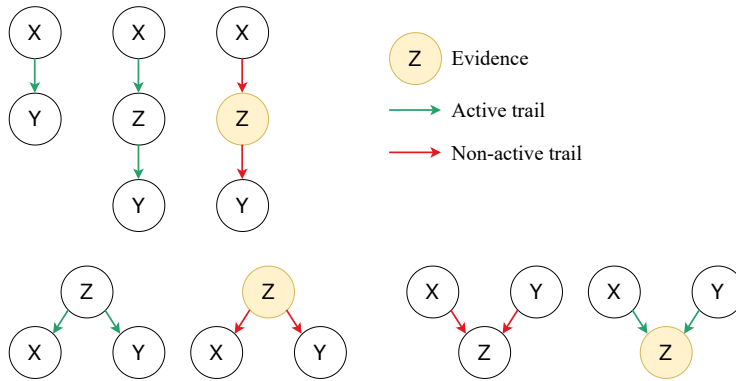


Figure 3.2: Example of active and non-active two-edge trails

**Active trail** A trail  $X_1 \rightleftharpoons \dots \rightleftharpoons X_n$  is active iff each two-edge trail  $X_{i-1} \rightleftharpoons X_i \rightleftharpoons X_{i+1}$  along the trail is active. Active trail

**D-separation** Two sets of nodes  $\mathbf{X}$  and  $\mathbf{Y}$  are d-separated given the evidence  $\mathbf{Z}$  if there is no active trail between any  $X \in \mathbf{X}$  and  $Y \in \mathbf{Y}$ . D-separation

**Theorem 3.2.1.** Two d-separated nodes are independent. In other words, two nodes are independent if there are no active trails between them.

### Independence algorithm

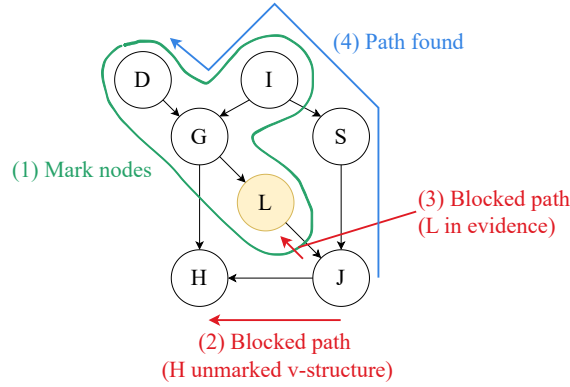
**Blocked node** A node is blocked if it blocks the flow. This happens if one and only one of the following conditions are met:

- The node is in the middle of an unmarked v-structure.
- The node is in the evidence.

To determine if  $X \perp Y$  given the evidence  $\mathbf{Z}$ :

1. Traverse the graph bottom-up marking all nodes in  $\mathbf{Z}$  or having a child in  $\mathbf{Z}$ .
2. Find a path from  $X$  to  $Y$  that does not pass through a blocked node.
3. If  $Y$  is not reachable from  $X$ , then  $X$  and  $Y$  are independent. Otherwise  $X$  and  $Y$  are dependent.

**Example.** To determine if  $J \perp D$ :

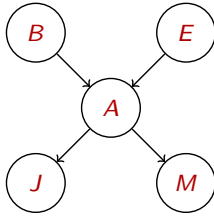


As a path has been found,  $J \not\perp D$ .

**Global semantics** Given a Bayesian network, the full joint distribution can be defined as the product of the local conditional distributions: Global semantics

$$\mathcal{P}(x_1, \dots, x_n) = \prod_{i=1}^n \mathcal{P}(x_i | \text{parents}(X_i))$$

**Example.** Given the following Bayesian network:



$$\begin{aligned} \mathcal{P}(j \wedge m \wedge a \wedge \neg b \wedge \neg e) \\ = \mathcal{P}(\neg b) \mathcal{P}(\neg e) \mathcal{P}(a | \neg b, \neg e) \mathcal{P}(j | a) \mathcal{P}(m | a) \end{aligned}$$

**Local semantics** Each node is conditionally independent of its non-descendants given its parents.

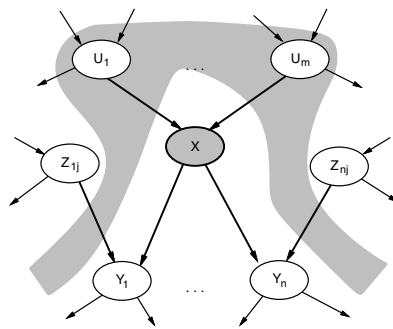


Figure 3.3: Local independence

**Theorem 3.2.2.** Local semantics  $\iff$  Global semantics

**Markov blanket** Each node is conditionally independent of all the other nodes if its Markov blanket (parents, children, children's parents) is in the evidence.

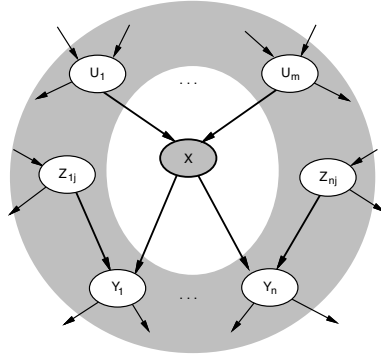


Figure 3.4: Markov blanket

## 3.3 Building Bayesian networks

### 3.3.1 Algorithm

The following algorithm can be used to construct a Bayesian network of  $n$  random variables:

1. Choose an ordering of the variables  $X_1, \dots, X_n$ .
2. For  $i = 1, \dots, n$ :
  - Add  $X_i$  to the network.
  - Select the parents of  $X_i$  from  $X_1, \dots, X_{i-1}$  such that:

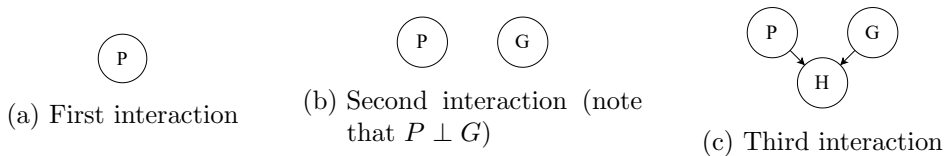
$$\mathbf{P}(X_i | \text{parents}(X_i)) = \mathbf{P}(X_i | X_1, \dots, X_{i-1})$$

By construction, this algorithm guarantees the global semantics.

**Example** (Monty Hall). The variables are:

- $G$ : the choice of the guest.
- $H$ : the choice of the host.
- $P$ : the position of the prize.

Note that  $P \perp G$ . Let the order be fixed as follows:  $P, G, H$ .



The nodes of the resulting network can be classified as:

**Initial evidence** The initial observation.

**Testable variables** Variables that can be verified.

**Operable variables** Variables that can be changed by intervening on them.

**Hidden variables** Variables that "compress" more variables to reduce the parameters.

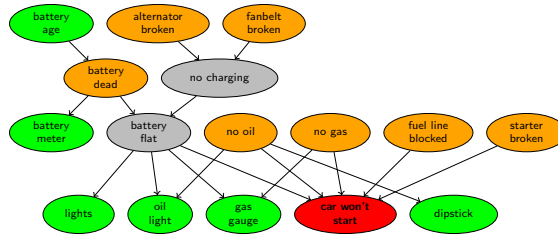
**Example.**

**Initial evidence** Red.

**Testable variables** Green.

**Operable variables** Orange.

**Hidden variables** Gray.



### 3.3.2 Structure learning

Learn the network from the available data.

Structure learning

**Constraint-based** Independence tests to identify the constraints of the edges.

**Score-based** Define a score to evaluate the network.

## 3.4 Causal networks

When building a Bayesian network, a correct ordering of the nodes that respects the causality allows to obtain more compact networks.

**Structural equation** Given a variable  $X_i$  with values  $x_i$ , its structural equation is a function  $f_i$  such that it represents all its possible values:

Structural equation

$$x_i = f_i(\text{other variables}, U_i)$$

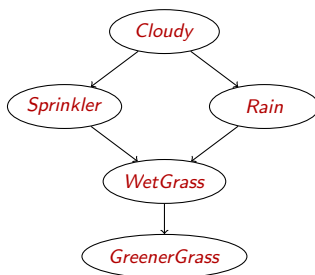
$U_i$  represents unmodeled variables or error terms.

**Causal network** Restricted class of Bayesian networks that only allows causally compatible ordering.

Causal network

An edge exists between  $X_j \rightarrow X_i$  iff  $X_j$  is an argument of the structural equation  $f_i$  of  $X_i$ .

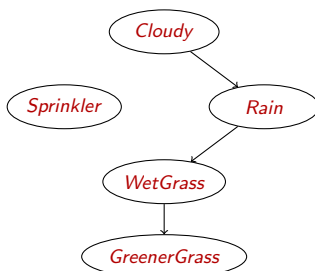
**Example.**



The structural equations are:

$$\begin{aligned} \text{cloudy} &= f_C(U_C) \\ \text{sprinkler} &= f_S(\text{Cloudy}, U_S) \\ \text{rain} &= f_R(\text{Cloudy}, U_R) \\ \text{wet\_grass} &= f_W(\text{Sprinkler}, \text{Rain}, U_W) \\ \text{greener\_grass} &= f_G(\text{WetGrass}, U_G) \end{aligned}$$

If the sprinkler is disabled, the network becomes:

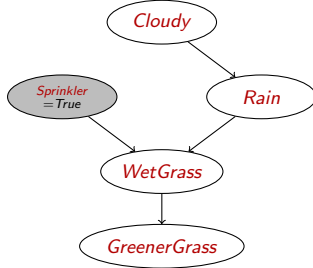


The structural equations become:

$$\begin{aligned} \text{cloudy} &= f_C(U_C) \\ \text{sprinkler} &= f_S(U_S) \\ \text{rain} &= f_R(\text{Cloudy}, U_R) \\ \text{wet\_grass} &= f_W(\text{Rain}, U_W) \\ \text{greener\_grass} &= f_G(\text{WetGrass}, U_G) \end{aligned}$$

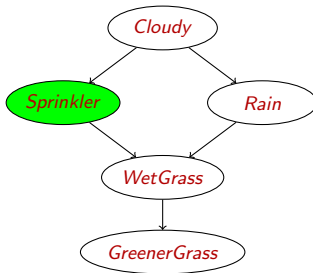
**do-operator** The do-operator allows to represent manual interventions on the network. The operation  $\text{do}(X_i = x_i)$  makes the structural equation of  $X_i$  constant (i.e.  $f_i = x_i$ , without arguments, so there won't be inward edges to  $X_i$ ). do-operator

**Example.**



By applying  $\text{do}(\text{Sprinkler} = \text{true})$ , the structural equations become:

$$\begin{aligned} \text{cloudy} &= f_C(U_C) \\ \text{sprinkler} &= \text{true} \\ \text{rain} &= f_R(\text{Cloudy}, U_R) \\ \text{wet\_grass} &= f_W(\text{Sprinkler}, \text{Rain}, U_W) \\ \text{greener\_grass} &= f_G(\text{WetGrass}, U_G) \end{aligned}$$



Note that Bayesian networks are not capable of modelling manual interventions. In fact, intervening and observing a variable are different concepts:

$$\begin{aligned} \mathcal{P}(\text{WetGrass} \mid \text{do}(\text{Sprinkler} = \text{true})) \\ \neq \\ \mathcal{P}(\text{WetGrass} \mid \text{Sprinkler} = \text{true}) \end{aligned}$$

## 3.5 Compact conditional distributions

Use canonical distributions (standard patterns) to reduce the number of variables in a conditional probability table.

### 3.5.1 Noisy-OR

Noisy-OR distributions model a network of non-interacting causes with a common effect. A node  $X$  has  $k$  parents  $U_1, \dots, U_k$  and possibly a leak node  $U_L$  to capture unmodeled concepts. Noisy-OR

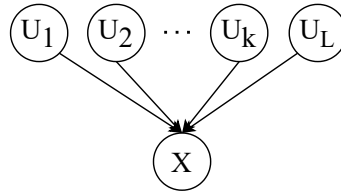


Figure 3.6: Example of noisy-OR network

Each node  $U_i$  has a failure (inhibition) probability  $q_i$ :

$$q_i = \mathcal{P}(\neg x \mid u_i, \neg u_j \text{ for } j \neq i)$$

The CPT can be built by computing the probabilities as:

$$\mathcal{P}(\neg x \mid \text{Parents}(X)) = \prod_{j: U_j = \text{true}} q_j$$

In other words:

$$\mathcal{P}(\neg x \mid u_1, \dots, u_n) = \mathcal{P}(\neg x \mid u_1) \cdot \mathcal{P}(\neg x \mid u_2) \cdot \dots \cdot \mathcal{P}(\neg x \mid u_n)$$

Because only the failure probabilities are required, the number of parameters is linear in the number of parents.

**Example.** We have as causes **Cold**, **Flu** and **Malaria** and as effect **Fever**. For simplicity there are no leak nodes. The failure probabilities are:

$$\begin{aligned} q_{\text{cold}} &= \mathcal{P}(\neg \text{fever} \mid \text{cold}, \neg \text{flu}, \neg \text{malaria}) = 0.6 \\ q_{\text{flu}} &= \mathcal{P}(\neg \text{fever} \mid \neg \text{cold}, \text{flu}, \neg \text{malaria}) = 0.2 \\ q_{\text{malaria}} &= \mathcal{P}(\neg \text{fever} \mid \neg \text{cold}, \neg \text{flu}, \text{malaria}) = 0.1 \end{aligned}$$

Known the failure probabilities, the entire CPT can be computed:

Cold	Flu	Malaria	$\mathcal{P}(\neg \text{fever})$	$1 - \mathcal{P}(\neg \text{fever})$
F	F	F	0.0	1.0
F	F	T	$q_{\text{malaria}} = 0.1$	0.9
F	T	F	$q_{\text{flu}} = 0.2$	0.8
F	T	T	$q_{\text{flu}} \cdot q_{\text{malaria}} = 0.02$	0.98
T	F	F	$q_{\text{cold}} = 0.6$	0.4
T	F	T	$q_{\text{cold}} \cdot q_{\text{malaria}} = 0.06$	0.94
T	T	F	$q_{\text{cold}} \cdot q_{\text{flu}} = 0.12$	0.88
T	T	T	$q_{\text{cold}} \cdot q_{\text{flu}} \cdot q_{\text{malaria}} = 0.012$	0.988

### 3.5.2 Hybrid Bayesian networks

Network with discrete and continuous random variables. Continuous variables must be converted into a finite representation. Possible approaches are:

Hybrid Bayesian networks

**Discretization** Values are divided into a fixed set of intervals. This approach may introduce large errors and large CPTs.

Discretization

**Finitely parametrized canonical families** There are two cases to handle using this approach:

Finitely parametrized canonical families

**Continuous child** Given the continuous variables  $X$  and  $C$  and a discrete (boolean, for simplicity) variable  $D$ , we want to compute the distribution  $\mathbf{P}(X \mid C, D)$ .

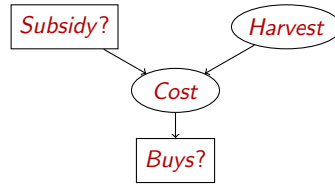
The discrete parent is handled by enumeration, by computing the probability over the domain of  $D$ .

For the continuous parent, an arbitrarily chosen distribution over the values of  $X$  is used. A common choice is the **linear Gaussian** whose mean is a linear combination of the values of the parents and the variance is fixed.

Linear Gaussian

A network with all continuous linear Gaussian distributions has the property of having a multivariate Gaussian distribution as joint distribution. Moreover, if a continuous variable has some discrete parents, it defines a conditional Gaussian distribution where, fixed the values of the discrete variables, the distribution over the continuous variable is a multivariate Gaussian.

**Example.** Let **Subsidy** and **Buys** be discrete variables and **Harvest** and **Cost** be continuous variables.



To compute  $\mathbf{P}(\text{Cost} \mid \text{Harvest}, \text{Subsidy})$ , we split the probabilities over the values of the discrete variable **Subsidy** and use a linear Gaussian for **Harvest**. We therefore have that:

$$\mathcal{P}(C = c \mid \text{Harvest} = h, \text{Subsidy} = \text{true}) = \mathcal{N}(a_t h + b_t, \sigma_t)(c)$$

$$\mathcal{P}(C = c \mid \text{Harvest} = h, \text{Subsidy} = \text{false}) = \mathcal{N}(a_f h + b_f, \sigma_f)(c)$$

where  $a_t, b_t, \sigma_t, a_f, b_f$  and  $\sigma_f$  are parameters.

**Discrete child with continuous parents** Given the continuous variable  $C$  and a discrete variable  $X$ , the probability of  $X$  given  $C$  is obtained by using a threshold function. For instance, probit or sigmoid distributions can be used.

### 3.5.3 Other methods

**Dynamic Bayesian network** Useful to model the evolution through time. A template variable  $X_i$  is instantiated as  $X_i^{(t)}$  at each time step.

Dynamic Bayesian network

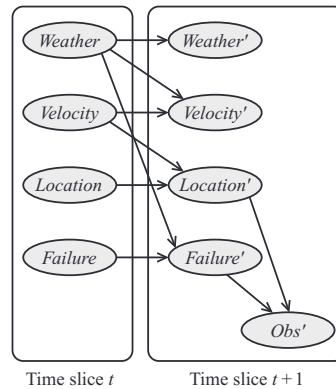


Figure 3.7: Example of dynamic Bayesian network

**Density estimation** Parameters of the conditional distribution can be learned using:

Density estimation

**Bayesian learning** calculate the probability of each hypothesis.

**Approximations** using the maximum-a-posteriori and maximum-likelihood hypothesis.

**Expectation-maximization algorithm.**

**Undirected graphical models** Markov networks are an alternative to probabilistic graphical models (as Bayesian networks). Markov networks are undirected graphs with factors (instead of probabilities) and are able to naturally capture independence relations.

Undirected graphical models



## 4 Exact inference

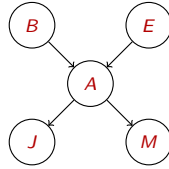
### 4.1 Inference by enumeration

Method to sum out a joint probability without explicitly representing it by using CPT entries.

Inference by enumeration

Enumeration follows a depth-first exploration and has a space complexity of  $O(n)$  and time complexity of  $O(d^n)$ . It must be noted that some probabilities appear multiple times but require to be recomputed because of the definition of the algorithm.

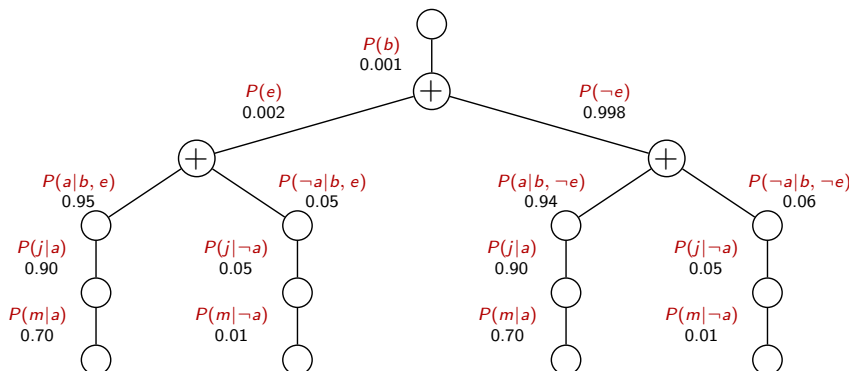
**Example** (Burglary). Given the Bayesian network:



We want to compute  $\mathbf{P}(B \mid j, m)$ :

$$\begin{aligned}
 \mathbf{P}(B \mid j, m) &= \frac{\mathbf{P}(B, j, m)}{\mathcal{P}(j, m)} \\
 &= \alpha \mathbf{P}(B, j, m) \\
 &= \alpha \sum_e \sum_a \mathbf{P}(B, j, m, e, a) \\
 &= \alpha \sum_e \sum_a \mathbf{P}(B) \mathcal{P}(e) \mathbf{P}(a \mid B, e) \mathcal{P}(j \mid a) \mathcal{P}(m \mid a) \\
 &= \alpha \mathbf{P}(B) \sum_e \mathcal{P}(e) \sum_a \mathbf{P}(a \mid B, e) \mathcal{P}(j \mid a) \mathcal{P}(m \mid a)
 \end{aligned}$$

This can be represented using a tree:



### 4.2 Inference by variable elimination

Method that carries out summations right-to-left and stores intermediate results (called factors).

Inference by variable elimination

**Pointwise product of factors**  $f(X, Y) \times g(Y, Z) = p(X, Y, Z)$

$X$	$Y$	$f(X, Y)$
0	0	1
0	1	3
1	0	2
1	1	1

$Y$	$Z$	$g(Y, Z)$
0	0	4
0	1	3
1	0	1
1	1	2

$X$	$Y$	$Z$	$f(X, Y) \times g(Y, Z)$
0	0	0	1 · 4
0	0	1	1 · 3
0	1	0	3 · 1
0	1	1	3 · 2
1	0	0	2 · 4
1	0	1	2 · 3
1	1	0	1 · 1
1	1	1	1 · 2

Figure 4.1: Example of pointwise product

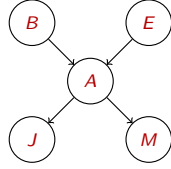
**Summing out** To sum out a variable  $X$  from a product of factors:

1. Move constant factors outside (i.e. factors that do not depend on  $X$ ).
2. Compute the pointwise product of the remaining terms.

**Example.**

$$\begin{aligned} \sum_X f_1 \times \cdots \times f_k &= f_1 \times \cdots \times f_i \sum_X f_{i+1} \times \cdots \times f_k \\ &= f_1 \times \cdots \times f_i \times f_X \end{aligned}$$

**Example (Burglary).** Given the Bayesian network:



We want to compute  $\mathbf{P}(B \mid j, m) = \alpha \mathbf{P}(B) \sum_e \mathcal{P}(e) \sum_a \mathbf{P}(a \mid B, e) \mathcal{P}(j \mid a) \mathcal{P}(m \mid a)$ .

We first work on the summation on  $A$ . We introduce as factors the entries of the CPT:

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{P}(B) \sum_e \mathcal{P}(e) \sum_a f_A(a, b, e) f_J(a) f_M(a)$$

Note that  $j$  and  $m$  are not parameters of the factors  $f_J$  and  $f_M$  because they are already given. We then sum out on  $A$ :

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{P}(B) \sum_e \mathcal{P}(e) f_{AJM}(b, e)$$

Now, we repeat the same process and sum out  $E$ :

$$\mathbf{P}(B \mid j, m) = \alpha \mathbf{P}(B) f_{EAJM}(b)$$

At last, we factor  $\mathbf{P}(B)$ :

$$\mathbf{P}(B \mid j, m) = \alpha f_B(b) f_{EAJM}(b)$$

### 4.2.1 Irrelevant variables

A variable  $X$  is irrelevant if summing over it results in a probability of 1.

Irrelevant variables

**Theorem 4.2.1.** Given a query  $X$ , the evidence  $\mathbf{E}$  and a variable  $Y$ :

$$Y \notin (\text{Ancestors}(\{X\}) \cup \text{Ancestors}(\mathbf{E})) \rightarrow Y \text{ is irrelevant}$$

**Theorem 4.2.2.** Given a query  $X$ , the evidence  $\mathbf{E}$  and a variable  $Y$ :

$$Y \text{ d-separated from } X \text{ by } \mathbf{E} \rightarrow Y \text{ is irrelevant}$$

### 4.2.2 Complexity

**Singly connected networks** Network where any two nodes are connected with at most one undirected path. Time and space complexity is  $O(d^k n)$ .

**Multiply connected networks** The problem is NP-hard.

## 4.3 Clustering algorithm

Method that joins individual nodes to form clusters. Allows to estimate the posterior probabilities for  $n$  variables with complexity  $O(n)$ .

Clustering algorithm

## 5 Approximate inference

**Stochastic simulation** Class of methods that draw  $N$  samples from the distribution and estimate an approximate posterior  $\hat{\mathcal{P}}$ .

Stochastic simulation

**$\delta$ -stochastic absolute approximation** Given  $\delta \in ]0, 0.5[$  and  $\varepsilon \in ]0, 0.5[$ , a  $\delta$ -stochastic absolute approximation has error:

$$\left| \mathcal{P}(X|\mathbf{E}) - \hat{\mathcal{P}}(X|\mathbf{E}) \right| \leq \varepsilon$$

Moreover, the method might fail (with greater error) with probability  $\delta$ .

**$\delta$ -stochastic relative approximation** Given  $\delta \in ]0, 0.5[$  and  $\varepsilon \in ]0, 0.5[$ , a  $\delta$ -stochastic relative approximation has error:

$$\frac{\left| \mathcal{P}(X|\mathbf{E}) - \hat{\mathcal{P}}(X|\mathbf{E}) \right|}{\mathcal{P}(X|\mathbf{E})} \leq \varepsilon$$

Moreover, the method might fail (with greater error) with probability  $\delta$ .

**Theorem 5.0.1.** Approximate inference is NP-hard for any  $\delta, \epsilon < 0.5$ .

**Consistency** A sampling method is consistent if:

Consistency

$$\lim_{N \rightarrow \infty} \hat{\mathcal{P}}(x) = \mathcal{P}(x)$$

### 5.1 Sampling from an empty network

Sample each variable in topological order (i.e. from parents to children).

The probability  $\mathcal{S}$  of sampling a specific event  $x_1, \dots, x_n$  is given by the probability of the single events knowing their parents:

Sampling from an empty network

$$\mathcal{S}(x_1, \dots, x_n) = \prod_{i=1}^n \mathcal{P}(x_i | \text{parents}(X_i)) = \mathcal{P}(x_1, \dots, x_n)$$

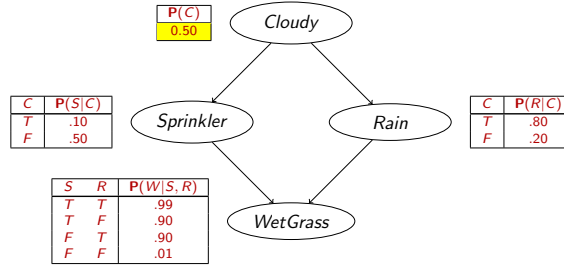
**Theorem 5.1.1.** Sampling from an empty network is consistent.

*Proof.* Let  $N$  be the number of samples and  $\mathcal{N}(x_1, \dots, x_n)$  the number of times the event  $x_1, \dots, x_n$  has been sampled.

$$\begin{aligned} \lim_{N \rightarrow \infty} \hat{\mathcal{P}}(x_1, \dots, x_n) &= \lim_{N \rightarrow \infty} \frac{\mathcal{N}(x_1, \dots, x_n)}{N} \\ &= \mathcal{S}(x_1, \dots, x_n) = \mathcal{P}(x_1, \dots, x_n) \end{aligned}$$

□

**Example.** Given the following Bayesian network:



A possible sampling order is **Cloudy**, **Sprinkler**, **Rain**, **WetGrass**.

Assuming that a random generator gives the sequence of probabilities (0.4, 0.8, 0.1, 0.5), the sample will be:

$$\begin{aligned}
 &\langle \mathcal{P}(C), \mathcal{P}(S|C), \mathcal{P}(R|C), \mathcal{P}(W|S, R) \rangle \\
 &\langle C = \text{false}, \mathcal{P}(S|C = \text{false}), \mathcal{P}(R|C = \text{false}), \mathcal{P}(W|S, R) \rangle \\
 &\langle C = \text{false}, S = \text{false}, R = \text{true}, \mathcal{P}(W|S = \text{false}, R = \text{true}) \rangle \\
 &\langle C = \text{false}, S = \text{false}, R = \text{true}, W = \text{true} \rangle
 \end{aligned}$$

Note that the adopted convention is the following: if  $r$  is the probability given by a random generator and  $\mathcal{P}(X) = p$ ,  $X = \text{true}$  if  $r \leq p$ .

## 5.2 Rejection sampling

Given a known evidence  $\mathbf{E}$ , rejection sampling works as sampling from an empty network but removes any sample that does not agree with the evidence.

Rejection sampling

Obviously if  $\mathcal{P}(\mathbf{E})$  is low, the majority of the samples will be discarded and more iterations are required to reach the desired number of samples.

**Theorem 5.2.1.** Rejection sampling is consistent.

*Proof.* Let  $\mathcal{N}(\mathbf{X})$  be the number of times the event  $\mathbf{X}$  has been sampled.

$$\begin{aligned}
 \hat{\mathcal{P}}(\mathbf{X}|\mathbf{E}) &= \frac{\mathcal{N}(\mathbf{X}, \mathbf{E})}{\mathcal{N}(\mathbf{E})} \\
 &\approx \frac{\mathcal{P}(\mathbf{X}, \mathbf{E})}{\mathcal{P}(\mathbf{E})} = \mathcal{P}(\mathbf{X}|\mathbf{E})
 \end{aligned}$$

The approximation derives from the consistency of sampling from an empty network.  $\square$

## 5.3 Likelihood weighting

Given a known evidence  $\mathbf{E}$ , likelihood weighting samples non-evidence variables and weights each sample by the likelihood of the evidence.

Likelihood weighting

The probability  $\mathcal{S}$  of sampling a specific event  $\mathbf{Z}$  and evidence  $\mathbf{E}$  is given by the probability of the single events in  $\mathbf{Z}$  knowing their parents:

$$\mathcal{S}(\mathbf{Z}, \mathbf{E}) = \prod_{z_i \in \mathbf{Z}} \mathcal{P}(z_i | \text{parents}(z_i))$$

The weight of a sample  $(\mathbf{Z}, \mathbf{E})$  is given by the probability of the single events in  $\mathbf{E}$  knowing their parents:

$$w(\mathbf{Z}, \mathbf{E}) = \prod_{e_i \in \mathbf{E}} \mathcal{P}(e_i | \text{parents}(e_i))$$

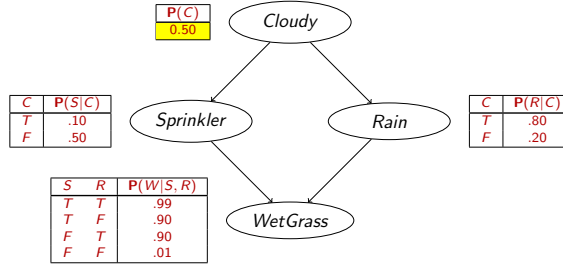
**Theorem 5.3.1.** Likelihood weighting is consistent.

*Proof.* The weighted sampling probability is given by:

$$\begin{aligned} \mathcal{S}(\mathbf{Z}, \mathbf{E}) \cdot w(\mathbf{Z}, \mathbf{E}) &= \prod_{z_i \in \mathbf{Z}} \mathcal{P}(z_i | \text{parents}(z_i)) \cdot \prod_{e_i \in \mathbf{E}} \mathcal{P}(e_i | \text{parents}(e_i)) \\ &= \mathcal{P}(\mathbf{Z}, \mathbf{E}) \end{aligned}$$

This is a consequence of the global semantics of Bayesian networks.  $\square$

**Example.** Given the following Bayesian network:



Knowing that  $S = \text{true}$  and  $W = \text{false}$ , we sample in the order: Cloudy, Rain.

Assuming that a random generator gives the sequence of probabilities (0.4, 0.1), the sample will be:

$$\begin{aligned} &\langle \mathcal{P}(C), S = \text{true}, \mathcal{P}(R|C), W = \text{false} \rangle \\ &\langle C = \text{true}, S = \text{true}, \mathcal{P}(R|C = \text{true}), W = \text{false} \rangle \\ &\langle C = \text{true}, S = \text{true}, R = \text{true}, W = \text{false} \rangle \end{aligned}$$

The weight associated to the sample is given by the probability of the evidence:

$$\begin{aligned} w &= \mathcal{P}(S = \text{true} | C = \text{true}) \cdot \mathcal{P}(W = \text{false} | S = \text{true}, R = \text{true}) \\ &= 0.1 \cdot (1 - 0.99) = 0.001 \end{aligned}$$

## 5.4 Markov chain Monte Carlo

Sampling on a Markov chain where states contain an assignment to all variables.

Adjacent states of the Markov chain differ by only one variable. Therefore, the probability of an edge connecting two states is given by the probability of the updated variable known its Markov blanket:

Markov chain Monte Carlo

$$\mathcal{P}(x_i | \text{markov\_blanket}(X_i)) = \mathcal{P}(x_i | \text{parents}(X_i)) \cdot \prod_{Z_j \in \text{children}(x_i)} \mathcal{P}(z_j | \text{parents}(Z_j))$$

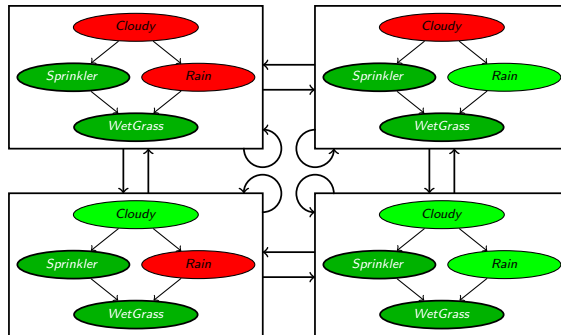
**Theorem 5.4.1.** Markov chain Monte Carlo is consistent.

Note: nevertheless, it is difficult to tell if convergence has been achieved.

*Proof.* Consequence of the fact that a long-run on a Markov chain converges to the posterior probability of the states.  $\square$

**Compiled network** A naive implementation of Markov chain Monte Carlo requires to repeatedly compute the probabilities with the Markov blanket. A solution is to compile the network into a model-specific inference code.

**Example.** Given the evidence  $S = \text{true}$  and  $W = \text{true}$ , the structure of the Markov chain can be defined as follows:



<end of course>