

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

Last update: 11 November 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

# 1 Introduction

## 1.1 Uncertainty

**Uncertainty** A task is uncertain if we have:

Uncertainty

- Partial observations
- Noisy or wrong information
- Uncertain action outcomes
- 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/nonmonotonic logic** Works on assumptions. An assumption can be contradicted by an evidence.

Default/nonmonotonic 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

$$\mathcal{P}(X = x_i) = \sum_{\omega \text{ st } X(\omega) = x_i} \mathcal{P}(\omega)$$

**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**.

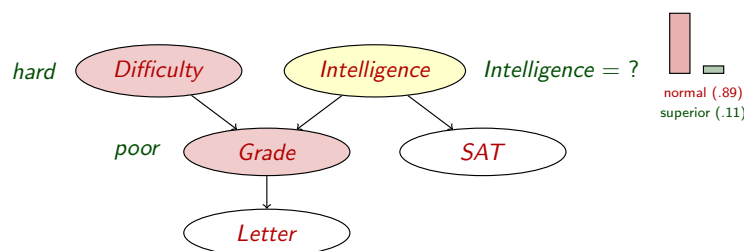


**Explain away** Observation obtained "passing through" other observations.

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 be 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  $Z$ .

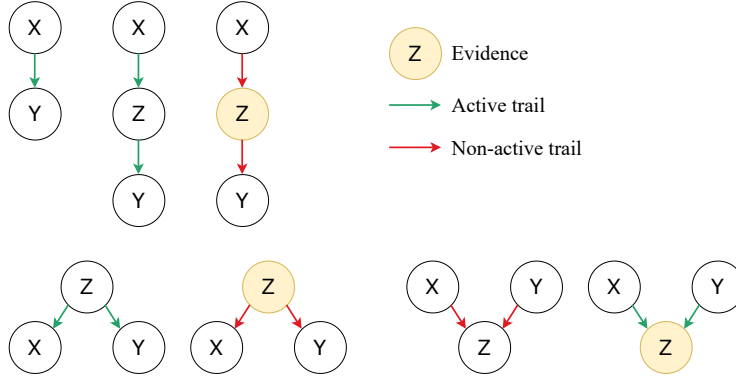


Figure 3.2: Example of active and non-active two-edge trails

**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 is no active trail between them.

### Independence algorithm

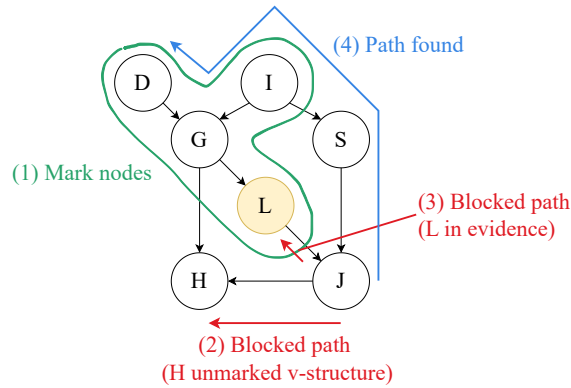
**Blocked node** A node is blocked if it blocks the flow. This happens if one and only one of the following conditions are met:

- 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  $J \perp D$ :

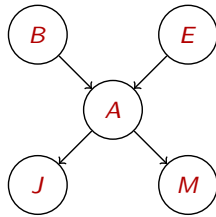


As a path has been found,  $J \not\perp D$ .

**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.

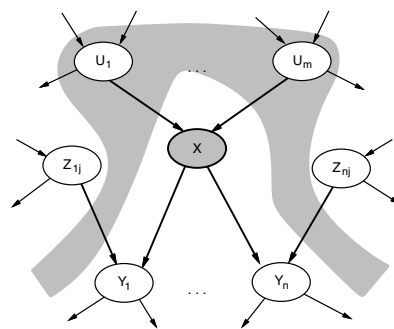


Figure 3.3: Local independence

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

**Markov blanket** Each node is conditionally independent of all 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:

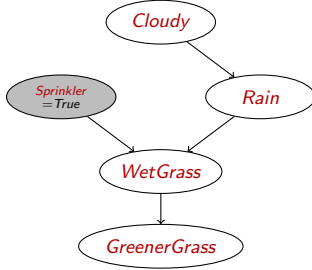


The structural equations become:

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

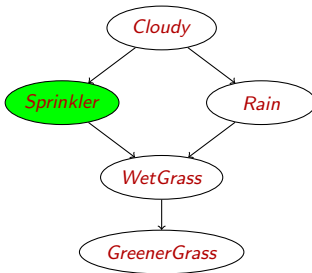
**do-operator** The do-operator allows to represent manual interventions on the network. The operation  $\text{do}(X_i = x_i)$  makes the structural equation of  $X_i$  constant (i.e.  $f_i = x_i$ , without arguments, so there won't be inward edges to  $X_i$ ). do-operator

**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:

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