

Get started

Open in app

**towards**
data science

Follow

565K Followers

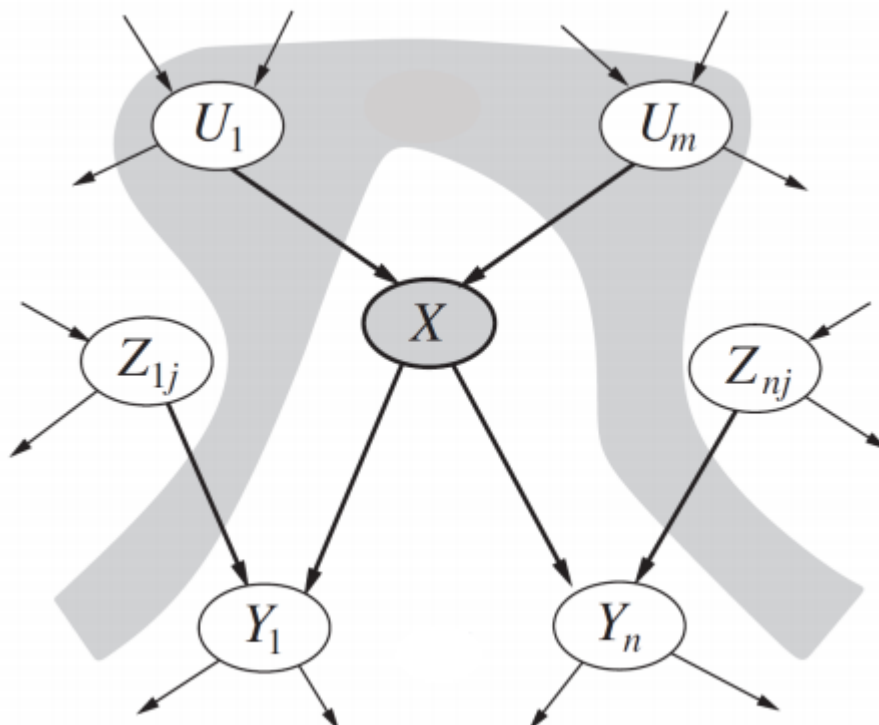


You have **2** free member-only stories left this month. [Sign up for Medium and get an extra one](#)

Introduction to Bayesian Networks



Devin Soni 📖 Jun 9, 2018 · 5 min read ★



Bayesian networks are a type of probabilistic graphical model that uses Bayesian inference for probability computations. Bayesian networks aim to model conditional dependence, and therefore causation, by representing conditional dependence by edges

in a directed graph. Through these relationships, one can efficiently conduct inference on the random variables in the graph through the use of factors.

• • •

Probability

Before going into exactly what a Bayesian network is, it is first useful to review probability theory.

First, remember that the joint probability distribution of random variables A_0, A_1, \dots, A_n , denoted as $P(A_0, A_1, \dots, A_n)$, is equal to $P(A_1 | A_2, \dots, A_n) * P(A_2 | A_3, \dots, A_n) * \dots * P(A_n)$ by the chain rule of probability. We can consider this a **factorized** representation of the distribution, since it is a product of N factors that are localized probabilities.

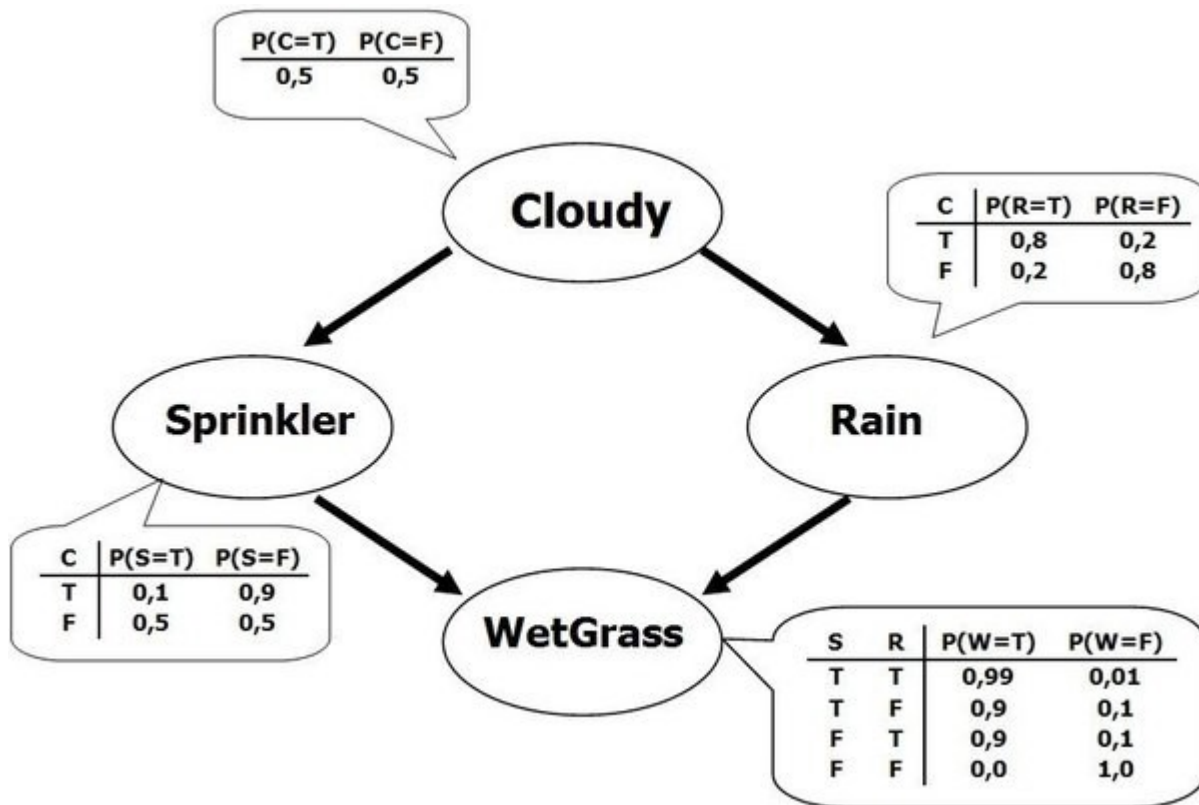
$$P\left(\bigcap_{k=1}^n A_k\right) = \prod_{k=1}^n P\left(A_k \mid \bigcap_{j=1}^{k-1} A_j\right)$$

Next, recall that **conditional independence** between two random variables, A and B , given another random variable, C , is equivalent to satisfying the following property: $P(A, B | C) = P(A | C) * P(B | C)$. In other words, as long as the value of C is known and fixed, A and B are independent. Another way of stating this, which we will use later on, is that $P(A | B, C) = P(A | C)$.

• • •

The Bayesian Network

Using the relationships specified by our Bayesian network, we can obtain a compact, factorized representation of the joint probability distribution by taking advantage of conditional independence.



A Bayesian network is a **directed acyclic graph** in which each edge corresponds to a conditional dependency, and each node corresponds to a unique random variable. Formally, if an edge (A, B) exists in the graph connecting random variables A and B , it means that $P(B|A)$ is a **factor** in the joint probability distribution, so we must know $P(B|A)$ for all values of B and A in order to conduct inference. In the above example, since Rain has an edge going into WetGrass, it means that $P(\text{WetGrass}|\text{Rain})$ will be a factor, whose probability values are specified next to the WetGrass node in a conditional probability table.

Bayesian networks satisfy the **local Markov property**, which states that a node is conditionally independent of its non-descendants given its parents. In the above example, this means that $P(\text{Sprinkler}|\text{Cloudy}, \text{Rain}) = P(\text{Sprinkler}|\text{Cloudy})$ since Sprinkler is conditionally independent of its non-descendant, Rain, given Cloudy. This property allows us to simplify the joint distribution, obtained in the previous section using the chain rule, to a smaller form. After simplification, the joint distribution for a Bayesian network is equal to the product of $P(\text{node}|\text{parents}(\text{node}))$ for all nodes, stated below:

In larger networks, this property allows us to greatly reduce the amount of required computation, since generally, most nodes will have few parents relative to the overall size of the network.

• • •

Inference

Inference over a Bayesian network can come in two forms.

The first is simply evaluating the joint probability of a particular assignment of values for each variable (or a subset) in the network. For this, we already have a factorized form of the joint distribution, so we simply evaluate that product using the provided conditional probabilities. If we only care about a subset of variables, we will need to marginalize out the ones we are not interested in. In many cases, this may result in underflow, so it is common to take the logarithm of that product, which is equivalent to adding up the individual logarithms of each term in the product.

The second, more interesting inference task, is to find $P(x|e)$, or, to find the probability of some assignment of a subset of the variables (x) given assignments of other variables (our evidence, e). In the above example, an example of this could be to find $P(\text{Sprinkler}, \text{WetGrass} | \text{Cloudy})$, where $\{\text{Sprinkler}, \text{WetGrass}\}$ is our x , and $\{\text{Cloudy}\}$ is our e . In order to calculate this, we use the fact that $P(x|e) = P(x, e) / P(e) = \alpha P(x, e)$, where α is a normalization constant that we will calculate at the end such that $P(x|e) + P(\neg x | e) = 1$. In order to calculate $P(x, e)$, we must marginalize the joint probability distribution over the variables that do not appear in x or e , which we will denote as Y .

$$P(x|e) = \alpha \sum_{\forall y \in Y} P(x, e, Y)$$

For the given example, we can calculate $P(\text{Sprinkler}, \text{WetGrass} | \text{Cloudy})$ as follows:

$$\begin{aligned} P(\text{Sprinkler}, \text{WetGrass} | \text{Cloudy}) &= \alpha \sum_{\text{Rain}} P(\text{WetGrass} | \text{Sprinkler}, \text{Rain}) P(\text{Sprinkler} | \text{Cloudy}) P(\text{Rain} | \text{Cloudy}) P(\text{Cloudy}) \\ &= \alpha P(\text{WetGrass} | \text{Sprinkler}, \text{Rain}) P(\text{Sprinkler} | \text{Cloudy}) P(\text{Rain} | \text{Cloudy}) P(\text{Cloudy}) + \alpha P(\text{WetGrass} | \text{Sprinkler}, \neg \text{Rain}) P(\text{Sprinkler} | \text{Cloudy}) P(\neg \text{Rain} | \text{Cloudy}) P(\text{Cloudy}) \end{aligned}$$

We would calculate $P(\neg x \mid e)$ in the same fashion, just setting the value of the variables in x to false instead of true. Once both $P(x \mid e)$ and $P(\neg x \mid e)$ are calculated, we can solve for α , which equals $1 / (P(x \mid e) + P(\neg x \mid e))$.

Note that in larger networks, Y will most likely be quite large, since most inference tasks will only directly use a small subset of the variables. In cases like these, exact inference as shown above is very computationally intensive, so methods must be used to reduce the amount of computation. One more efficient method of exact inference is through variable elimination, which takes advantage of the fact that each factor only involves a small number of variables. This means that the summations can be rearranged such that only factors involving a given variable are used in the marginalization of that variable. Alternatively, many networks are too large even for this method, so approximate inference methods such as MCMC are instead used; these provide probability estimations that require significantly less computation than exact inference methods.

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

 Get this newsletter

Data Science

Machine Learning

Artificial Intelligence

Big Data

Analytics



About Write Help Legal

Get the Medium app

