

# Bayesian Networks in Data Science

Tairan Ye

March 6, 2018

## 1 Introduction

A Bayesian network, Bayesian network, belief network, Bayesian model or probabilistic directed acyclic graphical model is a probabilistic graphical model, a type of statistical model, that represents a set of variables and their conditional dependencies via a directed acyclic graph (DAG). For example, a Bayesian network could represent the probabilistic relationships between diseases and symptoms. Given symptoms, the network can be used to compute the probabilities of presence of various diseases.

### 1.1 Statistical Introduction

Given data  $x$  and parameter  $\theta$ , a simple Bayesian analysis starts with a prior probability or prior,  $p(x|\theta)$  to compute a posterior probability  $p(\theta|x) \propto p(x|\theta)p(\theta)$ . Most likely, the prior on  $\theta$  depends in turn on other parameters  $\varphi$  that are not mentioned in the likelihood. So, the prior  $p(\theta)$  must be replaced by a likelihood  $p(\theta | \varphi)$ , and a prior  $p(\varphi)$  on the newly introduced parameters  $\varphi$  is required, resulting in a posterior probability:

$$p(\theta, \varphi|x) \propto p(x|\theta)p(\theta|\varphi)p(\varphi) \quad (1)$$

which is the simplest example of a hierarchical Bayes model.

The process will be repeated – for example, the parameters  $\varphi$  may depend in turn on additional parameters  $\psi$ , which will require their own prior. Eventually the process must terminate, with priors that do not depend on any other unmentioned parameters.

### 1.2 Major Application

#### 1. *Inferring Unobserved Variables*

Since a Bayesian network is a complete model for the variables and their relationships, it can be used to answer probabilistic queries. The most common approximate inference algorithms are importance sampling, stochastic MCMC simulation, mini-bucket elimination, loopy belief propagation, generalized belief propagation, and variational methods.

### 2. *Parameter Learning*

In order to fully specify the Bayesian network and thus fully represent the joint probability distribution, it is necessary to specify for each node  $X$  the probability distribution for  $X$  conditional upon  $X$ 's parents.

### 3. *Structure Learning*

A particularly fast method for exact BN learning is to cast the problem as an optimization problem, and solve it using integer programming.

## 2 Data

The dataset is named "Breast Cancer Wisconsin (Diagnostic)", from the "UCI Machine Learning Repository". Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image. Separating plane described above was obtained using Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, pp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes. The actual linear program used to obtain the separating plane in the 3-dimensional space is that described in: [K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34]. The total instances is 569.

Table 1: Summary of variables

Variable	Description
ID	Identification of patient's cases
Diagnosis	M = malignant, B = benign
Radius	Mean of distances from center to points on the perimeter
Texture	Standard deviation of gray-scale values
Perimeter	The length of the boundary line
Area	The total area size
Smoothness	Local variation in radius lengths
Compactness	$\frac{Perimeter^2}{area} - 1$
Concavity	Severity of concave portions of the contour
Concave points	Number of concave portions of the contour
Symmetry	Whether the area is symmetric or not
Fractal dimension	Coastline approximation - 1

### **3 Method**

In the final project, I plan to apply the Bayesian Networks to predict the diagnosis in the data set of Breast Cancer in Wisconsin. As a comparison, I am going to conduct other classifier approach and then analyze the results.