



Integrating expert knowledge with data in Bayesian networks: Preserving data-driven expectations when the expert variables remain unobserved



Anthony Costa Constantinou^{a,*}, Norman Fenton^{a,b}, Martin Neil^{a,b}

^a Risk and Information Management (RIM) Research Group, School of Electronic Engineering and Computer Science, Queen Mary University of London, London E1 4NS, UK

^b Agena Ltd., Cambridge CB23 7NU, UK

ARTICLE INFO

Article history:

Received 8 November 2015

Revised 26 February 2016

Accepted 29 February 2016

Available online 18 March 2016

Keywords:

Bayesian networks

Belief networks

Causal inference

Expert knowledge

Knowledge elicitation

Probabilistic graphical models

ABSTRACT

When developing a causal probabilistic model, i.e. a Bayesian network (BN), it is common to incorporate expert knowledge of factors that are important for decision analysis but where historical data are unavailable or difficult to obtain. This paper focuses on the problem whereby the distribution of some continuous variable in a BN is known from data, but where we wish to explicitly model the impact of some additional expert variable (for which there is expert judgment but no data). Because the statistical outcomes are already influenced by the causes an expert might identify as variables missing from the dataset, the incentive here is to add the expert factor to the model in such a way that the distribution of the data variable is preserved when the expert factor remains unobserved. We provide a method for eliciting expert judgment that ensures the *expected values* of a data variable are preserved under all the known conditions. We show that it is generally neither possible, nor realistic, to preserve the variance of the data variable, but we provide a method towards determining the accuracy of expertise in terms of the extent to which the variability of the revised empirical distribution is minimised. We also describe how to incorporate the assessment of extremely rare or previously unobserved events.

© 2016 Elsevier Ltd. All rights reserved.

1. Introduction and motivation

Causal probabilistic networks, also known as *Bayesian networks* (BNs), are a well-established graphical formalism for encoding conditional probabilistic relationships among uncertain variables. The nodes of a BN represent variables and the arcs represent causal or influential relationships between them. BNs are based on sound foundations of causality and probability theory; namely Bayesian probability (Pearl, 2009).

It has been argued that developing an effective BN requires a combination of expert knowledge and data (Fenton & Neil, 2012). Yet, rather than combining both sources of information, in practice many BN models have been ‘learnt’ purely from data, while others have been built solely on expert knowledge. Apart from lack of data, one possible explanation for this phenomenon is that in order to be able to combine knowledge with data researchers typically require a strong background in both data mining and

expert systems, as well as to have access to, and time for, the actual domain expert elicitation.

Irrespective of the method used, building a BN involves the following two main steps:

1. *Determining the structure of the network*: many of the real-world application models that have been constructed solely based on expert elicitation are in areas where humans have a good understanding of the underlying causal factors. These include medicine, project management, sports, forensics, marketing and investment decision making (Heckerman, Horvitz, & Nathwani, 1992a; Heckerman & Nathwani, 1992b; Andreassen, Riekehr, Kristensen, Schönheyder, & Leibovici, 1999; Lucas et al., 2000; van der Gaag, Renooij, Witteman, Aleman, & Taal, 2002; Fenton & Neil, 2012; Constantinou, Fenton, & Neil, 2012; Constantinou, Freestone, Marsh, Fenton, & Coid, 2015; Yet et al., 2013; 2015; Kendrick, 2015).

In other applications such as bioinformatics, image processing and natural language processing, the task of determining the causal structure is generally too complex for humans. With the advent of big-data, much of the current

* Corresponding author. Tel.: +44 7903292836.

E-mail addresses: anthony@constantinou.info (A.C. Constantinou), n.fenton@qmul.ac.uk (N. Fenton), m.neil@qmul.ac.uk (M. Neil).

research on BN development assumes that sufficient data are available to learn the underlying BN structure (Spirites & Glymour, 1991; Verma & Pearl, 1991; Spirites, Glymour, & Scheines, 1993; Friedman, Geiger, & Goldszmidt, 1997; 2000; Jaakkola, Sontag, Globerson, & Meila, 2010; Nassif, Wu, Page, & Burnside, 2012; Nassif et al., 2013; Petitjean, Webb, & Nicholson, 2013), hence assuming the expert's input is minimal or even redundant. Recent relevant research does relax this impression and allows for some expert input to be incorporated in the form of constraints (de Campos & Ji, 2011; Zhou, Fenton, & Neil, 2014a). It is, however, increasingly widely understood that incorporating expert knowledge can result in significant model improvements (Spiegelhalter, Abrams, & Myles, 2004; Rebonato, 2010; Pearl, 2009; Fenton & Neil, 2012; Constantinou et al., 2012; Constantinou, Fenton, & Neil, 2013; Zhou, Fenton, & Neil, 2014b), and this becomes even more obvious when dealing with interventions and counterfactuals (Constantinou, Yet, Fenton, Neil, & Marsh, 2016).

2. *Determining the conditional probabilities (CPTs) for each node (also referred to as the parameters of the model):* if the structure of the BN is learnt purely from data, then it is usual also for the parameter learning to be performed during that process. On the other hand, if expert knowledge is incorporated into a BN then parameter learning is, most typically, performed (or finalised) after the network structure has been determined.

The parameters can be learnt from data and/or expert judgments. If the data has missing values, then parameter learning is usually performed by the use of the Expectation Maximisation algorithm (Lauritzen, 1995), or other variations of this algorithm (Jamshidian & Jennrich, 1997; Jordan, 1999; Matsuyama, 2003; Hunter & Lange, 2004; Jiangtao, Yanfeng, & Lixin, 2012), which represent a likelihood-based iterative method for approximating the parameters of a BN. Other, much less popular methods, include restricting the parameter learning process only to cases with complete data, or using imputation-based approaches to fill the missing data points with the most probable values (Enders, 2006).

When developing BNs for practical applications, it is common to incorporate expert knowledge of factors that are important for decision analysis but where historical data is unavailable or difficult to obtain. That is the context for this paper. Previous related research in expert elicitation extensively covers:

1. *Accuracy in eliciting experts' beliefs:* it is often unrealistic to expect precise probability values to be provided by the expert. It is shown that participants with mathematical (or relevant) background tend to provide more accurate quantitative descriptions of their beliefs (Murphy & Winkler, 1977; Wallsten & Budescu, 1983). However, only few experts have sufficient mathematical experience and as a result, various probability elicitation methods have been proposed. These include probability scales with verbal and/or numerical anchors (Kuipers, Moskowitz, & Kassirer, 1988; van der Gaag, Renooij, Witteman, Aleman, & Taal, 1999; van der Gaag et al., 2002; Renooij, 2001), iterative processes which combine whatever the expert is willing to state (Druzdzel & van der Gaag, 1995), use of frequencies such as "1 in 10" in situations where events are believed to be based on extreme probabilities (Gigerenzer & Hoffrage, 1995), visual aids (Korb & Nicholson, 2011), as well as estimating the probabilities based on the lower and upper extremes of the experts' belief (Hughes, 1991).

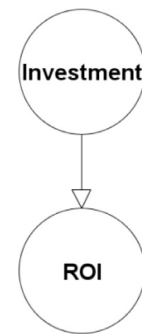


Fig. 1. Purely data-driven BN model M of the investment problem.

2. *Biases in experts' beliefs:* It has been demonstrated that limited knowledge of probability and statistics threatens the validity and reliability of expert judgments, leading to a number of biases (Johnson, Tomlinson, Hawker, Granton, & Feldman, 2010a). Various techniques for dealing with potential biases have been proposed. According to (Johnson et al., 2010b), these include provision of an example (Bergus, Chapman, Gjerde, & Elstein, 1995; Evans, Brooks, & Pollard, 1985; Evans, Handley, Over, & Perham, 2002; White, Pocock, & Wang, 2005), training exercises (Van der Fels-Klerx et al., 2002), use of clear instructions (Li & Krantz, 2005) or a standardised script (Chaloner, 1996), avoidance of scenarios or summaries of data, provision of feedback, verification, and opportunity for revision (O'Hagan, 1998; Normand, Frank, & McGuire, 2002), and a statement of the baseline rate or outcome in untreated patients (Evans et al., 2002). Further general guidelines in terms of how to reliably elicit expert judgments and minimise potential biases are provided in (Druzdzel & van der Gaag, 1995; O'Hagan et al., 2006; Johnson et al., 2010b).

While the above previous relevant research deals extensively with the process by which expert judgments are elicited, it does so under the assumption that any resulting CPTs will solely be based on expert knowledge as elicited. This paper tackles a problem which does not seem to have been addressed previously. Specifically, we are interested in preserving some aspects of a pure data-driven model when incorporating expert knowledge.

For example, we may have extensive historical data about *Return on Investment (ROI)* (we will call this the *dependent data node*) given different types of investment (such as properties, bonds, shares), as captured in the very simple BN model shown in Fig. 1.

If the data-driven *ROI* distribution given *Investment* is based on rich and accurate data that is fully representative of the context and is without bias, then we can be confident that the resulting marginal *ROI* distribution represents the *true* distribution. However, this distribution actually incorporates multiple dependent factors other than *Investment* type. If there is available expert knowledge about such factors such as, for example, *Economic growth*, then it is desirable to be able to incorporate such factors into an extended version of the BN as show in Fig. 2.

A logical and reasonable requirement is to preserve in M' as much as possible of the marginal distribution for the dependent data node (*ROI* in the example) when the expert variables (*Economic growth* in the example) remain unobserved. The paper describes a method to do this. In fact, for reasons explained in Section 2, it turns out that while it is possible to preserve the expected values of the marginal distribution under each of the known dependent scenarios, it is infeasible and unrealistic to preserve the variance. In Section 3, which describes the generic problem, we provide a method showing how to preserve the

Download English Version:

<https://daneshyari.com/en/article/382332>

Download Persian Version:

<https://daneshyari.com/article/382332>

[Daneshyari.com](https://daneshyari.com)