

Chapter 7.30

Retrieving Medical Records Using Bayesian Networks

Luis M. de Campos

Universidad de Granada, Spain

Juan M. Fernández Luna

Universidad de Granada, Spain

Juan F. Huete

Universidad de Granada, Spain

INTRODUCTION

Bayesian networks (Jensen, 2001) are powerful tools for dealing with uncertainty. They have been successfully applied in a wide range of domains where this property is an important feature, as in the case of information retrieval (IR) (Turtle & Croft, 1991). This field (Baeza-Yates & Ribeiro-Neto, 1999) is concerned with the representation, storage, organization, and accessing of information items (the textual representation of any kind of object). Uncertainty is also present in this field, and, consequently, several approaches based on these probabilistic graphical models have been designed in an attempt to represent documents and their contents (expressed by means of indexed terms), and the relationships between them, so as to retrieve as many relevant documents as possible, given a query submitted by a user.

Classic IR has evolved from flat documents (i.e., texts that do not have any kind of structure relating their contents) with all the indexing terms directly assigned to the document itself toward structured information retrieval (SIR) (Chiaromella, 2001), where the structure or the hierarchy of contents of a document is taken into account. For instance, a book can be divided into chapters, each chapter into sections, each section into paragraphs, and so on. Terms could be assigned to any of the parts where they occur. New standards, such as SGML or XML, have been developed to represent this type of document. Bayesian network models also have been extended to deal with this new kind of document.

In this article, a structured information retrieval application in the domain of a pathological anatomy service is presented. All the medical records that this service stores are represented

in XML, and our contribution involves retrieving records that are relevant for a given query that could be formulated by a Boolean expression on some fields, as well as using a text-free query on other different fields. The search engine that answers this second type of query is based on Bayesian networks.

BACKGROUND

Probabilistic retrieval models (Crestani et al., 1998) were designed in the early stages of this discipline to retrieve those documents relevant to a given query, computing the probability of relevance. The development of Bayesian networks and their successful application to real problems has caused several researchers in the field of IR to focus their attention on them as an evolution of probabilistic models. They realized that this kind of network model could be suitable for use in IR, specially designed to perform extremely well in environments where uncertainty is a very important feature, as is the case of IR, and also because they can properly represent the relationships between variables.

Bayesian networks are graphical models that are capable of representing and efficiently manipulating n -dimensional probability distributions. They use two components to codify qualitative and quantitative knowledge, respectively: first, a directed acyclic graph (DAG), $G=(V,E)$, where the nodes in V represent the random variables from the problem we want to solve, and set E contains the arcs that join the nodes. The topology of the graph (the arcs in E) encodes conditional (in)dependence relationships between the variables (by means of the presence or absence of direct connections between pairs of variables); and second, a set of conditional distributions drawn from the graph structure. For each variable $X_i \in V$, we therefore have a family of conditional probability distributions $P(X_i | pa(X_i))$, where $pa(X_i)$ represents any combination of the values of the variables

in $Pa(X_i)$, and $Pa(X_i)$ is the parent set of X_i in G . From these conditional distributions, we can recover the joint distribution over V .

This decomposition of the joint distribution gives rise to important savings in storage requirements. In many cases, it also enables probabilistic inference (propagation) to be performed efficiently (i.e., to compute the posterior probability for any variable, given some evidence about the values of other variables in the graph).

$$P(X_1, X_2, \dots, X_n) = \prod_{i=1}^n P(X_i | pa(X_i))$$

The first complete IR model based on Bayesian networks was the Inference Network Model (Turtle & Croft, 1991). Subsequently, two new models were developed: the Belief Network Model (Calado et al., 2001; Reis, 2000) and the Bayesian Network Retrieval Model (de Campos et al., 2003, 2003b, 2003c, 2003d). Of course, not only have complete models been developed in the IR context, but also solutions to specific problems (Dumais, et al., 1998; Tsirikika & Lalmas, 2002; Wong & Butz, 2000).

Structural document representation requires IR to design and implement new models and tools to index, retrieve, and present documents according to the given document structure. Models such as the previously mentioned Bayesian Network Retrieval Model have been adapted to cope with this new context (Crestani et al., 2003, 2003b), and others have been developed from scratch (Graves & Lalmas, 2002; Ludovic & Gallinari, 2003; Myaeng et al., 1998).

MAIN THRUST

The main purpose of this article is to present the guidelines for construction and use of a Bayesian-network-based information retrieval system. The source document collection is a set of medical records about patients and their medical tests stored

5 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/retrieving-medical-records-using-bayesian/26372

Related Content

Changing Healthcare Institutions with Large Information Technology Projects

Matthew W. Guah (2009). *Medical Informatics: Concepts, Methodologies, Tools, and Applications* (pp. 1689-1702).

www.irma-international.org/chapter/changing-healthcare-institutions-large-information/26330

A Framework for Privacy Assurance and Ubiquitous Knowledge Discovery in Health 2.0 Data Mashups

Jun Huan Liam Peyton (2010). *Ubiquitous Health and Medical Informatics: The Ubiquity 2.0 Trend and Beyond* (pp. 64-83).

www.irma-international.org/chapter/framework-privacy-assurance-ubiquitous-knowledge/42928

Innovative Smart Sensing Solutions for the Visually Impaired

Bruno Andò, Salvatore Baglio, Salvatore La Malfa and Vincenzo Marletta (2011). *Handbook of Research on Personal Autonomy Technologies and Disability Informatics* (pp. 60-74).

www.irma-international.org/chapter/innovative-smart-sensing-solutions-visually/48275

Web Portal for Genomic and Epidemiologic Medical Data

Mónica Miguélez Rico (2009). *Medical Informatics: Concepts, Methodologies, Tools, and Applications* (pp. 2351-2359).

www.irma-international.org/chapter/web-portal-genomic-epidemiologic-medical/26377

Motor Unit Synchronization as a Measure of Localized Muscle Fatigue

Sridhar P. Arjunanand Dinesh K. Kumar (2013). *International Journal of Biomedical and Clinical Engineering* (pp. 39-49).

www.irma-international.org/article/motor-unit-synchronization-as-a-measure-of-localized-muscle-fatigue/96827