

知识获取 知识工程 知识获取研究

45-47, 63

基于模型的知识获取及其一个范式*)

杜劲松 张尧庭

TP18

(武汉大学管理学院 武汉 430072)

摘要 Model-based knowledge acquisition is discussed and a paradigm is presented in this paper. Automatic acquisition of production rules of diagnostic expert systems is taken as examples. According to this paradigm, deep domain knowledge is represented by multivariate belief functions, while causes are obtained via abduction with uncertainties given by belief functions. This paradigm deserves robustness and nonmonotonicity.

关键词 Knowledge acquisition, Belief function, Uncertainty.

一、引言

知识获取被认为是知识工程的瓶颈,正如 Nilsson^[1]所指出的:……系统所需的成百条规则和上千个事实往往都是靠访问有关应用领域的专家来获取的,要把专家们的知识表示成事实和规则(或者任何其它形式的表达式)往往是一个枯燥而费时的过程,使这种知识获取过程自动化将是人工智能技术的一大进步。回顾历史,第一代知识获取技术是七十年代所广泛采用的“会晤”方式,即由计算机专家与领域专家长期交流;第二代知识获取技术则以八十年代计算机的辅助应用为特色;进入九十年代后,人们提出了高度集成的知识获取环境,这种环境提供从编辑、有效性证明到分析、综合等多种功能。

广义上看,学习过程也属于知识获取,而狭义地讲,知识获取仅指从已存在的来源中转移知识而不是学习,即便如此,这种转移也不是简单地从一种表示形式到另一种表示形式。专家们也许能够完成复杂的任务,但是他们到底是如何完成的,专家自己有时也难以理解,更不用说与知识工程师交流了,过去的作法是将专家作为黑箱,知识获取过程也就是对黑箱的行为进行整理与归纳。

基于模型的知识获取试图打破这一黑箱,它是由系统当前的领域知识引导的。这种领域知识用领域模型来表示,当具有完备的领域模型时,可以利用基于解释的学习进行信任知识获取^[2]。但是,实际上领域模型通常是不确定和不完备的。本文提出了

一个基于诱导推理(abduction)和多元信任函数^[3]的知识获取范式,它能够处理不确定的领域模型,并自动编译出产生式规则^[4]。

二、知识表示与获取中的模型

依据知识获取的狭义含义,它是指把知识从已经存在的来源转移到另一种形式,这就涉及到所谓知识的模型,从一定的意义上说,模型在知识工程研究中占据着基础性的地位。我们可能希望知道专家系统的功能的极限是什么?从领域专家那儿能抽取到什么样的知识?如何有效地表示这种知识?对于这些问题的回答,都必须依据相应的模型做出,而这些问题是知识工程研究中经常遇到的。

例如,在 MYCIN 的作者看来,产生式规则是一种合理的模型,不仅将它作为一种表示专门知识的模型,即:使用产生式规则描述医学专家的细菌感染知识,而且也将它作为一种认知模型。……(产生式系统)所具有的“条件-行动”的一般形式或“若……则”的抽象形式适用于不同的内容和不同性质的问题,成为问题求解中对选择算子进行控制的一般机制。因此,在问题求解中,正确识别条件起着重要作用,也是正确应用算子的前提。从这个角度看,问题求解过程也就成为获得与应用一定的产生式系统的过程。产生式概念为问题求解研究,也为整个思维、学习的研究提出新的思路。^[5]在这样思想基础之上,产生式规则也就自然地成为知识获取的模型。事实上,既然专家知识全部由“若……则……”这样的规则组成,那么就有可能从专家那儿一次抽取出来。所

*) 攀登计划资助。杜劲松 博士。张尧庭 教授。

以,几乎不需要知识工程师的干预,就可以利用自动或半自动的工具把人类专家的点滴知识直接转移到专家系统中来。

产生式规则的倡导者认为这样一种上下文无关文法可以表示他们所需要的全部知识。这种假定只有在积木世界中才能被满足,因为由于规则的结构单一性,导致它所表示的知识只能是“肤浅的”,它是一种编译过的知识,它死板地指定问题与解之间的直接对应关系,虽然它在某些情况下所显示的性能良好,但它不能给出令用户信服的解释,只是给出被“激活”的规则链;另外,产生式系统缺乏求解问题的灵活性,是脆弱的:当所处的环境变化时,产生式系统的性能会急剧下降。由于缺乏关于问题领域的深层次知识,产生式系统被称为肤浅或浅层模型。

与此相对应,人们提出了深层模型^[6]。从深层模型的观点看,专家们在推理时确实使用启发式规则(事实上,这也是专家与非专家的一个主要区别之一)。但是,这些启发式知识只是构成了专门知识的浅层模型。另一方面,专家们拥有关于领域的深层的知识结构,即背景知识,包括领域中的对象、对象之间的层次与因果关系以及各种相关性,领域的基本原理,或称初始原理(first principles),正是深层结构与浅层的启发式知识的交互作用,才使得专家能够快速处理常见问题,而对于不常见的问题,专家也能给出正确答案。

大多数的研究者认为,在深层模型与浅层模型之间并没有一个很明确的界限,认为它们之间存在一种连续的刻度似乎更合适,如果把深层模型与浅层模型分别视作一个区间的上、下界,那么下界将是由简单的事实组成的查询表,而上界是因果模型及其它深层结构,编译的、启发式的知识落在中间。

通过表示系统的结构和功能而对系统进行推理,这是深层模型派的基本思想,使用深层模型,将获得更强的问题求解能力,因为在启发式规则力所不及的情形下,有可能依靠使用领域的深层表示从初始原理出发进行推理。其次,使用深层模型也可以得到更好的解释,第三,利用深层模型还可以依据领域的原理来评价系统的性能与行为。最后,深层模型也为知识的自动获取提供了可能(这正是下文将要讨论的)。

当然,深层模型也有它自身的局限性,最直接的一点是,使用深层模型时,搜索是在一个更大的空间中进行的。因此,与启发式模型相比,也许需要更长的时间才能得到一个解。

另一方面,对于深层的基于知识的系统的研究目前还处于早期阶段,尚缺乏成熟的技术。从上面的讨论我们看到,深层模型与浅层模型具有某种互补性。如果我们抛开认知的观点,仅仅希望专家系统能够达到专家求解问题能力的水平(即对于例行的问题能够快速解答,对于太复杂或不寻常的问题也能给出正确回答),那么把深层模型与浅层模型结合起来是一种可行的选择。例如,Sticklen and Chandrasekaran^[7]把编译层的推理与深层推理集成起来,编译层系统与深层系统之间的权衡实际上是计算效率与问题求解通用性之间的权衡。因此,只要安排得当,那么既能获得编译层问题求解的效率特性,又不失深层推理所具备的鲁棒性。

现在的问题是,应该如何建立领域的深层模型以及编译层模型呢?我们建设用多元信任函数来描述领域的深层知识,用诱导推理得到不确定推理规则(即编译规则),并且这种不确定性由信任函数描述,我们以自动编译诊断规则为例来阐述。

三、一个知识获取范式

所谓范式,是指一个形式化的抽象表示体系或框架,而不是提供实现用的实际的数据结构。我们将知识获取理解为从已存在的来源中转移,那么一个知识获取范式就应该为以下两方面提供一个统一的框架:(1)关于知识源与目标知识的描述机制;(2)从知识源到目标知识的转移机制。

在我们建议的知识获取范式中,用多元信任函数(以及与它相关联的图形信任模型)来描述知识源,用产生式规则描述目标知识,这种规则的不确定程度由信任函数给出,利用诱导推理以及信任函数的运算来实现从知识源到目标知识的转移。下面,我们结合一个例子来阐明这一范式。

3.1 用多元信任函数表示深层知识

这里深层知识(即知识源)包括关于领域的对象以及对象之间相互关系的知识。试看一棵故障树,图1^[8]。

在该图中,M表示某一征兆, X_1, \dots, X_n 表示初始原因,其余结点为中间变量。为简单起见,假定所有变量均为二值的,有时,我们也用相应小写字母表示变量的状态。例如: \bar{x}_i 表示 $X_i=0$, x_i 表示 $X_i=1$,其中“1”表示故障状态,“0”表示正常状态。表示“或门”,例如, $M=1$ iff F和 X_i 中有一个为1;表示“与门”,例如, $D=1$ iff A,B均为1。这样一棵故障树构

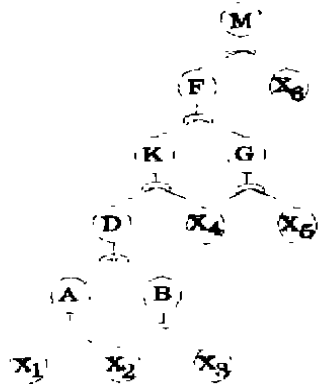


图 1

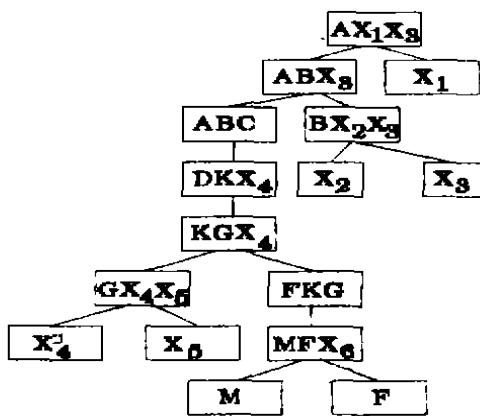


图 2

成了系统 M 的一因果模型。

信任函数是古典概率论的扩充,可以利用它描述关于变量的概率知识。例如,关于结点 X_1 , 可以有一个信任函数如下(用对应的基本概率赋值表示): $m\{1\}=0.05, m\{0\}=0.9, m\{0,1\}=0.05$ 。可以从上下概率的观点来理解信任函数,这样,结点 X_1 有故障的概率在区间 $[0.05, 0.1]$, 而它正常的概率在区间 $[0.9, 0.95]$ 。

当我们考虑变量的乘积空间时,我们得到多元信任函数。利用多元信任函数不仅可以表达概率知识,而且可以表达变量之间的逻辑关系。例如,由 F、K、G 这三个变量组成的或门可以用乘积空间 $F \times K \times G$ 上的信任函数表示,它具有以下基本概率赋值: $m\{(1,1,0), (1,0,1), (1,1,1), (0,0,0)\}=1$ 。于是,故障树可以利用若干个多元信任函数来表示。换句话说,若干个多元信任函数的组合就表示了关于整

个故障树的知识。把这些多元信任函数的识别骨架视为超边,那么这些超边就形成了一个图形信任模型,每条超边所表示的信息称为局部信息。

在这样的深层模型上,通常需要作如下的推理:
(1)系统处于故障状态的概率是多少?
(2)如果系统处于故障状态,那么哪些原因最有可能发生?
(3)如果又得到关于模型的进一步信息,那么系统发生故障的概率是多少?……一般而言,为了完成这些推理,要首先恢复全体变量的联合信任函数。但是,对一类特殊的图形模型(例如 Markov 树模型),可以仅从局部信息就可以解决:计算每个结点的边缘信任函数只需通过交换相邻结点的信息而完成,这将大大降低计算复杂性。图 2 给出了一棵 Markov 树,它满足:若 N_1 和 N_2 是树的任两个结点,那么 $N_1 \cap N_2$ 包含在 N_1 与 N_2 之间路径上的每个结点中。

3.2 用诱导推理确定故障原因

诱导推理是寻找一个数据集合的最好解释。以故障树为例,我们用诱导推理找出导致 M 故障的全部原因,假定因果理论由 Horn 子句组成,文[4]给出了递归算法。

每个诱导解释都满足最小性:即它的任何真子集都不能解释给定的数据集合。经过计算, $M=1$ 的全部诱导解释为: $\{X_1=1, X_3=1\}, \{X_2=2\}, \{X_1=1\}, \{X_5=1\}, \{X_6=1\}$ 。从故障树可以分析出,仅仅 X_1 或 X_2 故障而其余为正常时,是不能导致 M 故障的。所以 $\{X_1=1, X_3=1\}$ 满足最小性。

3.3 诊断规则的获取

在实际的专家系统中,使用产生式规则来描述诊断知识。以图 1 为例,规则形如:“如果 M 发生故障,那么可能是由 X_1, \dots, X_1 引起的”。其中 $\{X_1, \dots, X_1\}$ 是 M 的一个诱导解释。利用 Markov 树模型上的局部传播与组合算法,可以确定各诱导解释的信任测度与似真测度。

使用这样一种编译型的知识至少有两个好处:
(1)可以实现长推理链的捷径,加速诊断推理速度,因为不必每次都从因果域理论的基本原理出发进行推理;
(2)把信任测度与诱导解释相关联,当选择不同的初始原因进行检查和修复时可以作出某种优良性的决策。

3.4 范式的鲁棒性

把信任函数理论引入到知识获取过程中,能够处理不确定的情形,多元信任函数既具有统计途径的优点,又具有图形模型(因果网络)的优点;以概率

(下转第 63 页)

户进程运行在客户机上,服务进程运行在服务器上。采用这种体系结构,不但可以方便用户,即用户不必关心共享问题,而且可以提高系统的效率。

(2)多缓冲区技术。LNDBMS 的服务进程包括两个进程:接收进程和处理进程。这两个进程共用三个缓冲区:FullQueue、EmptyQueue、NullQueue。其中 FullQueue 用于存放客户进程发来的请求,EmptyQueue 用于存放空闲的 NCB 块(NCB 块即 Network Control Block,它用于 NETBIOS 调用),NullQueue 用于存放 NCB 外壳(即 NCB 块的空指针),接收进程采用异步接收的方式,它从 EmptyQueue 中取出一个 NCB 块用于接收以后,把 NCB 外壳放在 NullQueue 中,当实际收到客户进程的请求时,其中断例程从 NullQueue 中取出一个 NCB 外壳,将之组装成一个 NCB 块,然后放到 FullQueue 中,它们的关系如图1。

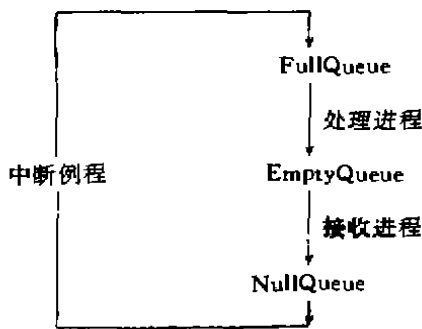


图1

(3)点对点通信, LNDBMS 采用点对点通信技术,其客户进程和服务进程间的通信采用 NETBIOS 协议,这自然也会提高效率。

(4)它是一个专用的数据库系统,只具备数据库系统中最基本的功能,如 CREATE TABLE、CREATE INDEX、INSERT、SELECT、UPDATE 等,它还具有专用的索引技术,如 HASH,因此与 ORACLE 等通用数据库系统相比,它的效率自然快得多。

5. 结束语

本文介绍了我们开发的面向数据采集的专用数据库系统—LNDBMS,它具有很高的效率和较好的安全性,有关其设计和实现技术及其作为联邦数据库系统中一个成员数据库系统的技术,将另文发表。

参考文献

- [1] ORACLE RDBMS Database Administrators Guide, Version 7, 1992
- [2] ORACLE RDBMS Server, Version 7, 1992
- [3] W. Kim, et al., Architecture of the Orion Next-Generation Database System, IEEE Trans. on Knowledge and Data Engineering, Vol. 2 No. 1, 1990
- [4] A. Delis, et al., Performance and Scalability of Client/Server Database Architectures, Proc. of the 18th International Conf. on VLDB, 1992
- [5] W. David Schwaderer, NetBIOS C 程序员指南, 上海科学普及出版社, 1990

(上接第47页)

为基础,知识的更新就反映为条件概率的变化;图形模型又允许知识以模块为单位增长,所以,这样一种知识获取范式能够适应噪声、不确定性以及变化的环境,即具有鲁棒性的特点,在开发能够自动获取、维护诊断专家系统的启发式知识库的工具时,这种鲁棒性正是一个值得追求的目标。

3.5 范式的非单调性

我们以一棵简单的故障树模型说明了基于模型的知识获取范式,实际的因果模型远比故障树复杂,但是可以认为故障树模型是建立因果模型的基本部件。由此,我们得到的模型有可能是不完备的,而这种不完备性决定了推理的非单调性。概率与诱导推

理都是可以与非单调推理结合起来的,参见文[4, 9, 10]。

四、结论

正如我们在引言中所指出的,本文所讨论的知识获取是指从已存在的来源中转移,在我们所提出的范式中,把用多元信任函数表示的领域知识转换成产生式规则,值得指出的是,“人工智能的核心是知识的表示与推理”这一观点已为大多数研究者所接受,我们提出的范式旨在为具有不确定性的知识的表示与推理提供严格的框架,把数理统计中的结果引入到人工智能中来。(参考文献共10篇略)