

一种结合效用的 Agent 思维状态模型*

徐晋晖, 石纯一

(清华大学 计算机科学与技术系, 北京 100084)

E-mail: xujinhui@263.net

http://www.tsinghua.edu.cn

摘要: 建立 Agent 思维状态模型是 Agent 理论研究的一个重要课题, 结合效用提出一种 Agent 思维状态模型 BDICU (belief desire intention combined with utility), 使 Agent 的行为具有逻辑理性和决策理性, 为副作用问题提供了解决方法. 同时, 给出目标和意图的生成和更新规则. BDICU 模型改进和扩充了 Rao 和 Georgeff 的信念-期望-意图理论, 为逻辑和效用理性 Agent 系统提供了实现支持.

关键词: Agent; 思维状态; 信念; 目标; 意图; 效用; 可能世界

中图法分类号: TP18 **文献标识码:** A

Agent 的理论研究普遍基于一种“意向姿态(intentional stance)”的观点, 把人类行为过程中的信念、能力、意图、承诺等心智成分引入 Agent 的研究中, 来刻画和描述 Agent 的概念和特征, 改变了经典 AI 研究中所采用的单纯模拟人类的知识表示与推理能力的研究模式^[1~3].

Cohen 和 Levesque^[4]首先进行了思维状态模型的研究, 提出基本模态算子 Bel, Goal, Happens, Done 的模态逻辑系统, 意图作为一个导出概念加以定义, 分析了 Agent 的信念、目标、规划、意图、承诺间理性平衡以及意图的作用. 进而, Rao 和 Georgeff^[5,6]给出 BDI 模型, 建立了 BDI 解释器, 讨论了承诺性质. Konolige 和 Pollack^[7]针对副作用问题, 以非正规模态逻辑作为描述工具, 建立了一个关于意图的形式化系统. Linder^[8]试图给出统一的形式化框架来研究动机类别的 Agent 认识状态, 并提出 Agent 效用理性研究的想法. 但是, 上述对思维状态模型的研究还存在 3 个问题: (1) 理论与实现的脱离问题, 原因是可能世界过于抽象, 难以在现实模型中找到对应的映射, 以致于实现的 Agent 系统不能符合所给的理论^[9]; (2) 思维状态模型没有考虑决策理性, 而把决策的过程体现在具体的实现中, 使得逻辑与实现脱节^[8]; (3) 副作用问题, 这是由正规模态逻辑的逻辑全知问题所导致的, 虽然有的模型给出了部分解决方法^[7], 但尚未满意地加以解决.

本文针对以上问题给出了一种结合效用的思维状态模型 BDICU (belief desire intention combined with utility), 在语言上给出了决策所涉及的参数的刻画, 在语义解释模型中引入了两个可能世界之间的变化值和实现变化代价度量函数, 使得逻辑理性和效用理性结合起来. 概括 Rao 和 Georgeff 所给出的公理分析了 BDICU, 建立了带效用限制的 KU 公理, 使副作用问题可以得到避免. 另外, 从实现角度分析了目标和意图的生成以及更新问题.

1 非形式化分析

理性 Agent 要不断地进行决策行为, 必然涉及效用问题^[9,10]. 效用可以通过行为的代价和行为所导致的 Agent 状态改变的收益来求得, 这样, 依代价和收益就可以刻画 Agent 的决策过程.

* 收稿日期: 1999-01-07; 修改日期: 1999-09-03

基金项目: 国家自然科学基金资助项目(69773026; 69733020); 教育部高等学校重点实验室访问学者基金资助项目; 清华大学博士学位论文基金资助项目

作者简介: 徐晋晖(1966—), 男, 山西人, 博士, 主要研究领域为多 Agent 系统; 石纯一(1935—), 男, 河北人, 教授, 博士生导师, 主要研究领域为人工智能应用基础.

作为 Agent 的思维状态成分信念、目标和意图等,应该有对于代价和收益的心理刻画。

例 1:我相信我将成为一个科学家。

我想成为一个科学家。

我打算成为一个科学家。

我作为一个 Agent 相信、想和打算自己成为一个科学家,同时对于成为一个科学家的代价和收益是有一个心理评判的,并且这种评判指导我的行为。

例 2:我执行成为一个科学家的行为。

对于这样的一个行为,如何去评价其理性应从两个角度来进行:(1) 这种行为是逻辑理性的,即我相信是可能的同时自己有这样的目标,就是现有思维状态模型刻画的理性;(2) 这种行为是决策理性的,即这种行为在所有可能的行为中(因为可能存在多个逻辑理性的行为,如:成为一个政治家,成为一个企业家等)效用是极大的。

我们在 Agent 的思维状态模型 BDI 中引入一个二元组 bc , 来描述对应的代价和收益,以 $BEL(A, \varphi, bc)$, $GOAL(A, \varphi, bc)$, $INTEND(A, \varphi, bc)$ 表示。

BDI 对应的语义解释是可能世界模型,在可能世界中存在一种可达关系。从一个可能世界到另一个可能世界需要执行相应的动作,而这种动作需要 Agent 对应的资源消耗,同时,Agent 进入不同的可能世界会产生状态上的改变。

例 3:目前的 Agent 可能世界是学生,下一个可达的可能世界是科学家。

从学生世界到科学家世界,Agent 要作出各种努力,这就是可能世界变化的代价,处于科学家和学生世界对应的价值是不同的。为了给出 $BEL(A, \varphi, bc)$, $GOAL(A, \varphi, bc)$, $INTEND(A, \varphi, bc)$ 中 bc 的语义解释,我们将可能世界之间存在的可能世界变化的代价和价值引入到解释模型中。

对于副作用问题,我们以一个研究 BDI 常引用的例子来加以说明。

例 4:知道拔牙会导致痛苦,这样,如果有拔牙的意图,就会有痛苦的意图。

这是明显不符合 Agent 理性的,如果一个模型纯粹考虑逻辑的理性,这个问题则是不可避免的,因为副作用本质上是一个符合逻辑理性的公式,但并不符合效用理性。如果结合效用理性的观点,可以得知痛苦的效用是很低的,这样,即使采纳了拔牙的意图,也不会采纳痛苦的意图。可见,要想解决副作用问题,有必要引入效用理性的考虑。

2 形式化模型

BDICU 是在 Rao 和 Georgeff 的 BDI 模型的基础上结合效用的一种扩充,下面的描述延用了有关符号,我们仅对扩充之处加以说明。

2.1 语言

定义 1. BC 是一个二元组集合。元素 $bc = (b, c)$, 其中 b 表示收益, c 表示代价。效用函数 U 是 $BC \rightarrow R$ 的映射。

效用函数 U 通常满足 $\partial U / \partial b > 0$ 和 $\partial U / \partial c < 0$ 约束,也就是说,效用随收益的增大而增大,随代价的变小而增大。新增加的状态公式如定义 2(其他的公式与已有的公式一致)。

定义 2. 如果 φ 是一个状态公式且 $bc \in BC$, 那么 $BEL(A, \varphi, bc)$, $GOAL(A, \varphi, bc)$, $INTEND(A, \varphi, bc)$ 是状态公式。

2.2 语义

定义 3. 一个解释模型 $M = \langle W, B, D, I, \Phi, RB, RC \rangle$ 。

其中 W 是可能世界集合, B, D, I 分别是信念、愿望和意图可达关系, 函数 $RB: W \times W \rightarrow R$ 度量可能世界的价值变化量; 函数 $RC: W \times W \rightarrow R$ 度量可能世界之间可达的代价, Φ 是原子命题的真值函数。

对应于 $BEL(A, \varphi, bc)$, $GOAL(A, \varphi, bc)$, $INTEND(A, \varphi, bc)$ 的语义解释如下:

$(M, v, w) \models BEL(A, \varphi, bc)$ 当且仅当 $\forall w' \in B(w), (M, v, w') \models \varphi$ 且 $RB(w', w) \geq b$ 且 $RC(w', w) \leq c$ 。

$\langle M, v, w \rangle \models \text{GOAL}(A, \varphi, bc)$ 当且仅当 $\forall w' \in D(w), \langle M, v, w' \rangle \models \varphi$ 且 $RB(w', w) \geq b$ 且 $RC(w', w) \leq c$.

$\langle M, v, w \rangle \models \text{INTEND}(A, \varphi, bc)$ 当且仅当 $\forall w' \in I(w), \langle M, v, w' \rangle \models \varphi$ 且 $RB(w', w) \geq b$ 且 $RC(w', w) \leq c$.

b, c 分别是对应的 w 的信念(目标、意图)可达世界集合 $B(w)(G(w), I(w))$ 中状态改变的最小值和改变所花费代价的最大值. 如果在给定 b, c 的情况下, 则对于所有 w 信念(目标、意图)可达的世界 w' 可能有 $\langle M, v, w' \rangle \models \varphi$, 但是不满足其他两个条件, 我们可以将全部满足的那些可能世界称为效用信念(目标、意图)可达的世界, 分别记为 BU, GU, IU , 这样便有下面的性质 1.

性质 1. $BU(w) \subseteq B(w), DU(w) \subseteq D(w), IU(w) \subseteq I(w)$.

$\neg \text{BEL}(\varphi, bc), \neg \text{GOAL}(\varphi, bc)$ 和 $\neg \text{INTEND}(\varphi, bc)$ 为真, 当且仅当所对应的解释中, 3 个与条件里有 1 个不成立.

3 有关公理

3.1 模型性质公理

根据通常的约定 BEL 满足 KD45 系统公理, GOAL 和 INTEND 满足 KD 系统公理, 在 BDICU 中, BEL, 不含 bc 的 GOAL 和不含 bc 的 INTEND 公式符合以上的公理. 但是, 对于含有 bc 的 GOAL 和 INTEND 公式满足带效用阈值限制的公理 KU.

定义 4. 效用阈值 AU 是 Agent 作出决策时的最低效用值.

通常地, 不同的 Agent 有不同的 AU, 并且一个 Agent 在不同的条件下会有不同的 AU.

公理 1(KU 公理).

$\text{GOAL}(A, P \rightarrow Q, (U(bc_1) \geq AU) \wedge (U(bc_2) \geq AU)) \rightarrow (\text{GOAL}(A, P, bc_1) \rightarrow \text{GOAL}(A, Q, bc_2)),$ (1)

$\text{INTEND}(A, P \rightarrow Q, (U(bc_1) \geq AU) \wedge (U(bc_2) \geq AU)) \rightarrow (\text{INTEND}(A, P, bc_1) \rightarrow \text{INTEND}(A, Q, bc_2)).$ (2)

bc_1 和 bc_2 分别是 P 和 Q 的效用描述, KU 公理是指 Agent 如果以 P 作为目标(意图), 那么 Q 也成为 Agent 的目标(意图)是有条件的, 也就是说, 对应的效用不能小于该 Agent 的效用阈值 AU.

利用该公理可以避免副作用问题. 用例 4 来说明, 有 $\text{GOAL}(A, \text{拔牙} \rightarrow \text{牙痛}, (U(\text{拔牙的 } bc_1) \geq AU \wedge (U(\text{牙痛的 } bc_2) < AU)))$, 这样, 即使有 $\text{GOAL}(A, \text{拔牙}, bc_1)$, 但是不能有 $\text{GOAL}(A, \text{牙痛}, bc_2)$. 对于 INTEND 的解释类似.

推论 1. 由 $\text{GOAL}(A, P \rightarrow Q)$ 不能推出 $(\text{GOAL}(A, P, bc_1) \rightarrow \text{GOAL}(A, Q, bc_2))$, 由 $\text{INTEND}(A, P \rightarrow Q)$ 不能推出 $(\text{INTEND}(A, P, bc_1) \rightarrow \text{INTEND}(A, Q, bc_2))$.

3.2 BEL, GOAL 和 INTEND 关系公理

根据 Rao 和 Georgeff 给出的有关公理, 我们可以得到以下类似的公理.

公理 2(信念-目标相容公理). $\text{GOAL}(A, \varphi, bc) \rightarrow \text{BEL}(A, \varphi, bc)$.

公理 3(目标-意图相容公理). $\text{INTEND}(A, \varphi, bc) \Rightarrow \text{GOAL}(A, \varphi, bc)$.

公理 4(关于意图的信念公理). $\text{INTEND}(A, \varphi, bc) \Rightarrow \text{BEL}(\text{INTEND}(A, \varphi, bc))$.

公理 5(关于意图的目标公理). $\text{INTEND}(A, \varphi, bc) \Rightarrow \text{GOAL}(\text{INTEND}(A, \varphi, bc))$.

公理 6(意图放弃公理). $\text{INTEND}(A, \varphi, bc) \Rightarrow \text{inevitabl} \diamond \neg (\text{INTEND}(A, \varphi, bc))$.

4 模型实现分析

根据 Rao 和 Georgeff 给出的 BDI 解释器:

BDI Interpreter

```
initialize_state();
do
  options := option_generator(event_queue, B, G, I);
  selected_options := deliberate(options, B, G, I);
  update_intentions(selected_options, I);
```

```

execute(I);
get_new_external_events();
drop_successful_attitudes(B, G, I);
drop_impossible_attitudes(B, G, I);
until quit

```

可以发现,在一个思维状态模型的具体实现过程中,核心问题是:信念、目标和意图的生成与更新.对于这些问题,Rao 和 Georgeff 在有关逻辑理性的基础研究中给出了分析,这里,我们结合逻辑理性和效用理性来加以描述.关于信念的处理可以参考有关文献,本文着重于目标和意图方面.

4.1 目标的生成与更新

目标一般是指当前未成立,并且相信将来会成立的命题.这样,则有:

目标生成规则.如果一个 Agent 对于 φ 知道目前不成立,相信将来成立,并且 φ 的实现所带来的效用大于等于其效用阈值,那么 Agent 可以将 φ 作为其目标.

$$BEL(A, \neg \varphi) \wedge BEL(A, \text{inevitabl } \varphi, bc) \wedge (U(bc) \geq AU) \Rightarrow GOAL(A, \varphi, bc).$$

这样,Agent 不是将所有当前未成立,并且将来会成立的命题都作为目标,而是要满足效用阈值.

目标冲突消解规则.如果一个 Agent 有两个冲突的目标,那么选择效用大的目标,删除效用小的目标.

$$GOAL(A, \varphi, bc_1) \wedge GOAL(A, \varphi, bc_2) \wedge (U(bc_1) > U(bc_2)) \Rightarrow GOAL(A, \varphi, bc_1).$$

目标冲突有显式和隐式冲突两种,显式冲突即 $\varphi = \neg \varphi$;隐式冲突是指两个目标不能同时执行,或者一个目标的实现导致另一个目标不能实现.

目标持续规则.Agent 有一个目标直到相信目标成立,或该目标必然不能实现,或即使可以实现但是效用值低于效用阈值.

$$GOAL(A, \varphi, bc) \Rightarrow \text{inevitabl } [GOAL(A, \varphi, bc) \text{ Until } (BEL(A, \varphi) \vee (BEL(A, \text{inevitabl } \neg \varphi) \vee (BEL(A, \varphi, bc) \wedge (U(bc) < AU)))]].$$

4.2 意图的生成与更新

意图是对于目标的承诺,意图的生成依赖于是否有实现目标的规划.一个规划是一棵由行为或子目标节点构成的树.一个 Agent 有一个实现目标 φ 的规划 PL 可以用 $\text{Has_PLAN}(A, PL, \varphi, c)$ 来描述,其中 c 是对应规划的代价,关于其对应的语义解释见文献[11].

意图生成规则.如果一个 Agent 有目标,并且有对应的规划,且规划的代价小于等于目标预期的代价.

$$GOAL(A, \varphi, bc) \wedge BEL(A, \text{Has_PLAN}(A, PL, \varphi, c_1)) \wedge (c_1 \leq c) \Rightarrow \text{INTEND}(A, \varphi, bc).$$

规划选择规则.如果 Agent 有两个可以实现目标 φ 的规划,那么选择代价小的规划.

$$\text{Has_PLAN}(A, PL_1, \varphi, c_1) \wedge \text{Has_PLAN}(A, PL_2, \varphi, c_2) \wedge (c_1 \leq c_2) \Rightarrow \text{Has_PLAN}(A, PL_1, \varphi, c_1).$$

意图持续是指 Agent 何时重新考虑承诺的目标,Rao 和 Georgeff 给出 3 种不同的承诺策略,粗略地描述了不同性格的 Agent.本文结合效用给出意图持续规则.

意图持续规则.Agent 有一个意图直到相信其实现,或对应的目标不存在,或没有可以完成该目标的规划.

$$\text{INTEND}(A, \varphi, bc) \Rightarrow \text{inevitabl } [\text{INTEND}(A, \varphi, bc) \text{ Until } (BEL(A, \varphi) \vee \neg GOAL(A, \varphi, bc) \vee \neg \text{Has_Plan}(A, PL, \varphi, c))].$$

目标不存在通过目标持续规则确定,没有规划是指规划库中无规划,或者有规划但是代价大于意图的代价.

5 结 语

Rao 和 Georgeff^[5,6]建立的思维状态模型,可以刻画 Agent 的逻辑理性,但是难以刻画效用理性,也难以解决模态逻辑固有的副作用问题,所给出的实现模型忽略了效用在目标和意图产生中的作用.Linder^[8]提出应建立逻辑理性和效用理性的思维状态模型,但是没有给出具体的工作.S. Kraus^[10]只是从效用理性的观点出发分析 Agent 的行为过程,没有结合思维状态研究逻辑理性问题.Cohen 和 Levesque^[4],Konolige 和 Pollack^[7]虽然建

立了对应的思维状态模型,但是没有给出实现的考虑.

本文在 Rao 和 Georgeff 工作的基础上,提出了结合效用的思维状态模型 BDICU,实现了逻辑理性和效用理性的结合,保证了 Agent 行为的逻辑理性和决策理性.为副作用问题的解决提供了方法.结合模型的实现,给出了符合逻辑和效用的目标和意图的生成以及更新规则. BDICU 模型对 Rao 和 Georgeff 的信念-期望-意图理论进行了改进和扩充,解决了副作用问题和单纯的逻辑理性的片面性,为逻辑和效用理性 Agent 系统提供了实现支持.以 BDICU 模型作为支持,所提出的 Agent 之间的社会承诺建立机制^[2]验证了该模型的优点.

References:

- [1] Wooldridge, M., Jennings, N. R. Intelligent agents: theory and practice. Knowledge Engineering Review, 1995,10(2):115~152.
- [2] Shoham, Y. Agent-Oriented programming. Artificial Intelligence, 1993,60(1):51~92.
- [3] Ma, Guang-wei, Xu, Jin-hui, Shi, Chun-yi. About the mental state model for Agent. Journal of Software, 1999,10(4):342~348 (in Chinese).
- [4] Cohen, P. R., Levesque, H. J. Intention is choice with commitment. Artificial Intelligence, 1990,43(2-3):213~261.
- [5] Rao, Georgeff. Modeling rational Agents within a BDI-architecture. In: Allen, J., Fikes, R., Sandewall, W eds. Principles of Knowledge Representation and Reasoning. San Mateo, CA: Morgan Kaufmann Publishers, Inc., 1991. 473~484.
- [6] Rao, Georgeff. An abstract architecture for rational Agents. In: Nebel, B., Rich, C., Swartout, W eds. Principles of Knowledge Representation and Reasoning. San Mateo, CA: Morgan Kaufmann Publishers, Inc., 1992. 439~449.
- [7] Konolige, K., Polack, M. E. A representationalist theory of intention. In: Bajcsy, R ed. Proceedings of the 13th International Joint Conference on Artificial Intelligence. San Mateo, CA: Morgan Kaufmann Publishers, Inc., 1993.390~395.
- [8] Linder, B. V., der Hock, W. V., Ch. Meyer, J. J. Formalising motivational attitudes of Agents on preference, goals and commitments. In: Wooldridge, M., Muller, J. P., Tambe, M eds. Intelligent Agent II: Agent Theories, Architectures, and Languages. Berlin: Springer-Verlag, 1996. 17~32.
- [9] Wooldridge, M. This is MYWORLD: the logic of an agent-oriented DAI testbed. In: Wooldridge, M., Jennings, N. R eds. Intelligent Agents—Proceedings of the ECAI'94 Workshop on Agent Theories, Architectures, and Languages. Berlin: Springer-Verlag, 1995. 160~178.
- [10] Kraus, S. An overview of incentive contracting. Artificial Intelligence, 1996,83(3):297~346.
- [11] Haddadi, A. S. Communication and Cooperation in Agent Systems. Berlin: Springer-Verlag, 1996. 1~134.
- [12] Xu, Jin-hui, Shi, Chun-yi. One mechanism of social commitment based on belief-desire-intention and utility. Journal of Software, 1999,10(8):829~834 (in Chinese).

附中文参考文献:

- [3] 马光伟,徐晋晖,石纯一. Agent 思维状态模型. 软件学报,1999,10(4):342~348.
- [12] 徐晋晖,石纯一. 一种基于信念-期望-意图和效用的社会承诺机制. 软件学报,1999,10(8):829~834.

A Model of Mental States Combined with Utility for Agents

XU Jin-hui, SHI Chun-yi

(Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China)

E-mail: xujinhui@263.net

http://www.tsinghua.edu.cn

Received January 7, 1999; accepted September 3, 1999

Abstract: It is an important subject to build the model of mental states for Agents in the theoretic research of agent. A model BDICU (belief desire intention combined with utility) of mental states combined with utility is presented in this paper, which guarantees agent's action with rational logic and decision, the problem of side effect with a solution is provided, and the rules for producing and updating goal and intention are also presented. This model modifies and extends Rao and Georgeff's Belief-Desire-Intention theory, and provides rational agent systems with implementing support on logic and utility.

Key words: agent; mental state; belief; goal; intention; utility; possible world