

# 一种基于信念-期望-意图和效用的社会承诺机制\*

徐晋晖 石纯一

(清华大学计算机科学与技术系 北京 100084)

**摘要** 在多智能体系统中,为完成问题求解,智能体之间需建立起社会承诺.文章将信念-期望-意图和效用理论结合起来,提出了一种社会承诺机制,保证了智能体行为的逻辑理性和决策理性.该文的工作将 Rao & Georgeff 的信念-期望-意图理论和 S. Kraus 的激励承包思想有机地结合起来,改进了 Haddadi 的社会承诺机制,为 Castelfranchi 的社会承诺提供了实现支持.

**关键词** 多智能体系统,信念,期望,意图,效用,社会承诺.

**中图法分类号** TP18

在 MAS(multi-agent systems)中,当一个 Agent/组(如果不特殊说明是一组 Agent)根据目前的情景难以完成一个任务(目标)时,可以通过支付报酬委托其他 Agent 来完成全部或部分任务,这是一种承包现象.

可将承包中涉及的两个 Agent 称为委托者  $D$  和采纳者  $A$ ,所委托的任务称为项目  $P$ , $P$  可以是  $D$  的目标或规划中的一个子目标.目前的研究工作中<sup>[1~4]</sup>将  $A$  有意图去完成  $D$  委托的项目  $P$  称为  $A$  对  $D$  的社会承诺,其作用是共同联合行为的胶合剂,是个体承诺向集体承诺过渡的桥梁和组/队合作的基础.

Castelfranchi<sup>[2]</sup>对社会承诺进行了定性的分析,指出社会承诺中存在双方的社会约束,不足之处是没有考虑社会承诺的建立和解除过程,社会承诺的定义未能反映社会约束;Haddadi<sup>[1]</sup>讨论了组内 Agent 之间社会承诺的建立和解除的通信问题,但没有考虑社会约束,社会承诺局限于组内,不具有普遍性.

那么如何结合社会约束给出社会承诺的概念?如何在自主自利的 Agent 之间建立社会承诺?在环境动态变化的情况下,社会承诺如何变化?这些将是本文所要解决的问题.

本文通过合同结构来描述社会承诺,合同结构描述了双方的责任、禁止、允许等社会约束.社会承诺与 BDI (belief-desire-intention)之间有着密切的关系,甚至社会承诺可以通过 BDI 来进行定义,这样,社会承诺的建立和解除就可以通过 BDI 的演化来进行分析.对 Agent 的 BDI 可以从不同的角度来观察,我们从 Agent 角度考虑社会承诺的建立和解除.

Agent 是自主自利的,这样, $D$  委托给  $A$  一个  $P$ , $A$  可能采纳也可能不采纳,这是由  $A$  自主决定的;同样地,当  $D$  有多个候选委托者时,也要根据某种利益标准选择合适的  $A$ ,而当不存在候选者时,需激励 Agent 去采纳其项目,那么如何作出选择?如何激励?这是社会承诺建立过程中需要考虑的问题.本文以效用理论为基础对此进行讨论.

在 MAS 研究中,要求 Agent 的行为是理性的,从逻辑角度来说,Agent 的行为可以通过逻辑的方式推导出来;从决策的角度来说,Agent 的行为应该符合效用最大化.我们将 BDI 与效用理论相结合而建立的社会承诺机制满足了理性的这两种解释.

我们假定 Agent 的能力是有限的和不同的,Agent 是自主自利的,Agent 可以在同一个组内,通信渠道是正常的.

\* 本文研究得到国家自然科学基金和清华大学研究生院博士学位论文基金资助.作者徐晋晖,1966年生,博士生,主要研究领域为分布式人工智能应用基础.石纯一,1935年生,教授,博士生导师,主要研究领域为人工智能应用基础.

本文通讯联系人:徐晋晖,北京 100084,清华大学计算机科学与技术系

本文 1998-07-03 收到原稿,1998-09-11 收到修改稿

## 1 基本框架

### 1.1 效用概念

效用理论是 Agent 追求效用最大化的一种决策理论,有关的概念阐述如下.

代价(cost)是 Agent 完成一个活动所需要的资源耗费,代价函数  $C$  是  $ACTIVITY(活动集) \rightarrow R(实数集)$  的映射.

收益(benefit)是 Agent 实现一个目标所带来的经济价值,收益函数  $B$  是  $GOAL(目标集) \rightarrow R$  的映射.

报酬(reward)是  $D$  对  $A$  完成  $P$  的奖励,报酬函数  $R$  是  $ACTIVITY * GOAL \rightarrow R$  的映射.报酬对  $D$  而言是代价,对  $A$  而言是收益.

委托者效用(utility)是  $D$  通过将  $P$  委托给  $A$  来完成所获得的利益,委托者效用函数  $U_D$  是  $BENEFIT * REWARD \rightarrow R$  的映射.

采纳者效用是  $A$  通过完成  $D$  委托的  $P$  所获得的利益,采纳者效用函数  $U_A$  是  $COST * REWARD \rightarrow R$  的映射.

### 1.2 结合效用的 BDI 理论

Agent 的 BDI 理论以 Rao & Georgeff<sup>[5,6]</sup>的工作为基础,本文考虑到 BDI 与效用的结合,作如下说明.

信念(belief)是 Agent 对于世界或环境的知识,用  $(BEL X \varphi)$  表示 Agent  $X$  相信  $\varphi$ .

目标(goal)是 Agent 希望达到的世界状态,用  $(GOAL X \varphi)$  表示  $X$  期望  $\varphi$  成立.

目标有多种来源,可以是用户或其他 Agent 提出的,也可以是用户自己产生的.目标的实现可以给 Agent 带来一定的收益  $b$ ,收益函数  $B$  依赖目标的来源,可以是用户或其他 Agent 说明的,也可以是根据信念和动机确定的或根据总目标计算的.以  $(GOAL X \varphi b)$  表示  $X$  有目标  $\varphi$  且收益是  $b$ ,有  $(GOAL X \varphi b) \equiv (GOAL X \varphi) \wedge (BEL X B(\varphi) = b)$ .

规划(plan)是一个有目标的活动,用  $(Has-Plan X PL \varphi)$  表示  $X$  有完成  $\varphi$  的规划  $PL$ .

Agent 执行规划必然带来资源的耗费  $c$ ,可以根据行为序列的构成由基本行为求得,也就是代价函数  $C$  依赖于基本行为的耗费和规划的构成形式.以  $(Has-Plan X PL \varphi c)$  表示  $X$  有完成  $\varphi$  的规划  $PL$  且代价是  $c$ ,有  $(Has-Plan X PL \varphi c) \equiv (Has-Plan X PL \varphi) \wedge (BEL X C(PL) = c)$ .

意图(intention)是 Agent 打算按相应的规划去实现目标,用  $(INT X \varphi)$  表示  $X$  有意图去实现  $\varphi$ .

Agent 对于这种意图的结果用效用  $u$  来判断, $u$  可以通过目标收益和规划代价来确定.以  $(INT X \varphi u)$  表示  $X$  有意图去实现  $\varphi$  且效用是  $u$ ,有  $(INT X \varphi u) \equiv (INT X \varphi) \wedge (BEL X U(b, c) = u)$ ,其中  $U$  是效用函数,为了简化,令  $u = b - c$ .

#### Rao & Georgeff 公理

$$(Has-Plan X PL \varphi) \wedge (GOAL X \varphi) \wedge (BEL X Pre(PL)) \wedge (Choose X PL \varphi) \Rightarrow (INT X Body(PL))$$

这里,将  $PL$  分为前提条件和实体,公理意指  $X$  有目标  $\varphi$ ,完成  $\varphi$  的规划,相信规划的前提条件存在,选择该规划作为实现  $\varphi$  的规划,可以推出  $X$  有意图去执行规划的实体部分.根据公理,本文引入效用得到下面的定理.

**定理.**  $(Has-Plan X PL \varphi c) \wedge (GOAL X \varphi b) \wedge (BEL X Pre(PL)) \wedge (BEL X U(b, c) > 0) \wedge (Max(U(b, c))) \Rightarrow (INT X Body(PL) u)$ .

证明提示:由于  $(BEL X U(b, c) > 0) \wedge (Max(U(b, c)))$  可以保证  $(Choose X PL \varphi)$ .

意指  $X$  有收益是  $b$  的目标  $\varphi$ ,完成  $\varphi$  的代价为  $c$  的规划  $PL$ ,相信  $PL$  规划的前提条件存在,效用大于 0 且在所有的可选规划中效用最大,可以推出  $X$  有意图去执行带来效用  $u$  的规划的实体部分.

定理说明了结合效用的 BDI 理论,可以保证两种理性的解释.

### 1.3 合同结构和社会承诺

关于社会承诺的定义<sup>[2,4]</sup>,只是从 BDI 的角度给出的,忽略了  $D$  和  $A$  的责任、禁止和允许,如  $(S-COMM X Y \varphi) \equiv (INT X \varphi) \wedge (GOAL Y (DONE X \varphi))$ .在人类承包现象中,通常用合同来描述委托者和采纳者双方的责任

等约束,于是我们引入合同结构.

令 合同结构  $CS=(\text{合同号}, \text{对象组}, \text{约束}, \text{期限})$

其中,对象组:  $= (D, A)$ , 对象:  $= D|A$ ,

约束:  $= (\text{对象}, \text{约束说明})$ , 约束说明:  $= \text{责任}; \text{禁止}; \text{允许}$ ,

责任:  $= \text{责任描述表}$ , 责任描述:  $= (\text{责任}; \text{条件}, \text{必须执行行为}, \text{报酬})$ ,

禁止:  $= \text{禁止描述表}$ , 禁止描述:  $= (\text{禁止}; \text{条件}, \text{不能执行行为}, \text{罚金})$ ,

允许:  $= \text{允许描述表}$ , 允许描述:  $= (\text{允许}; \text{条件}, \text{行为})$ ,

期限:  $= (\text{有效期}; \text{开始时间到结束时间})$ .

从外部观察者的角度,可以通过合同结构来建立社会承诺,

$$(S-COMM \ X Y \ cs) \equiv (MBEL \ X Y \ cs) \wedge (EINT \ X Y \ cs) \wedge (MBEL \ X Y \ (EINT \ X Y \ cs)).$$

这里,  $cs$  是一个具体的合同,  $MBEL$  是相互信念,  $EINT$  是每一个 Agent 都有意图,  $EINT(X Y cs) \equiv INT(X cs) \wedge INT(Y cs)$ .

从 Agent 角度来定义,

$$(S-COMM_D \ A \ cs) \equiv (BEL \ D \ cs) \wedge (INT \ D \ cs \ u_D) \wedge (BEL \ D \ (INT \ A \ cs)).$$

$$(S-COMM_A \ D \ cs) \equiv (BEL \ A \ cs) \wedge (INT \ A \ cs \ u_A) \wedge (BEL \ A \ (INT \ D \ cs)).$$

其中  $(INT \ X \ cs)$  是指  $X$  有意图去履行合同  $cs$ ;  $u_D$  和  $u_A$  分别是  $D$  和  $A$  履行合同的效用. 根据定理对于  $D$  而言有

$$(GOAL \ D \ \varphi \ b) \wedge (Has-Plan \ D \ cs \ \varphi \ c) \wedge (BEL \ D \ U(b, c) > 0 \wedge (Max \ (U(b, c)))) \Rightarrow (INT \ D \ cs \ u).$$

这里,  $D$  的目标是  $\varphi$ ,  $cs$  是  $D$  实现该目标的规划,  $c$  是  $D$  付给  $A$  的报酬.

对于  $A$  而言,有

$$(GOAL \ A \ cs \ r) \wedge (Has-Plan \ A \ PL \ cs \ c) \wedge (BEL \ A \ U(r, c) > 0 \wedge (Max \ (U(r, c)))) \Rightarrow (INT \ A \ cs \ u).$$

这里,  $A$  的目标是  $cs$ , 并且有实现该目标的规划,  $r$  是执行合同的报酬.

这样,通过合同使  $D$  和  $A$  之间建立起社会承诺,与现有的定义相比,一是确立了双方的约束,二是强调了社会承诺的相互性.

## 2 社会承诺过程

社会承诺的建立过程是在  $D$  和  $A$  之间,通过协商达到一个共同同意的合同,并且作出意图去履行合同的决定;社会承诺的解除过程是在合同履行中,责任、禁止和允许所描述的条件出现,导致  $D$  和  $A$  需要重新考虑合同,作出有关的决定和行为.

### 2.1 $D$ 的社会承诺建立

(1) 发现自己不能完成的目标,  $(GOAL \ D \ \varphi \ b) \wedge \neg (Has-plan \ D \ PL \ \varphi)$ ;

(2) 产生委托其他 Agent 完成  $\varphi$  的目标,  $(GOAL \ D \ (Contract \ \varphi))$ , 由于存在实现委托目标的规划,执行下面的有关步骤;

(3) 通知其他 Agent 信息  $((Contract \ D \ \varphi), r_D, \text{响应时间})$ , 这里,  $r_D$  是  $D$  提出的初始报酬,通知可以有针对性地发送;

(4) 在响应时间内无  $(Can-Adopt \ A \ \varphi \ r_A \ c)$  消息接收,可以认为委托规划失败,结束下面的执行,  $r_A$  是  $A$  要求的初始报酬,  $c$  是  $A$  实现  $\varphi$  的代价,如果  $A$  与  $D$  不在同一个合作组内,该信息可以为 0 ( $A$  对  $D$  保密其实际的投入);

(5) 从多个候选者中选择一个  $A$ ,

① 计算由每一个候选者完成  $\varphi$  带来的效用; IF  $A$  和  $D$  在同一组内 THEN  $u_D = b - c$  ELSE  $u_D = b - r_A$ , 这说明  $A$  和  $D$  在一个组内追求整体效用 (因为  $u_D = b - c = b - r_A + r_A - c$ ), 否则考虑  $D$  的自身效用;

② 将  $u_D$  从大到小排队;

③ 从中取出第 1 个,就是预选的  $A$ ;

(6) 如果  $u_D < 0$ , 需要重新设计对  $A$  的报酬, 直到  $u_D > 0$  满意为止, 这时有 ( $Want D A \varphi r$ ), 将该信息通知对应的  $A$ ;

(7) 在规定的时间内没有收到 ( $Willing A \varphi$ ), 这时可以认为  $A$  不满意, 需要重新考虑报酬或其他候选者;

(8) 建立合同,

① 如果  $D$  收到  $A$  的合同  $cs$  信息, 可以根据情况确定是否满意, 如果满意, 则形成履行合同的意图 ( $INT D cs u_D$ ) 并通知  $A$ , 否则, 对合同进行修改, 并将修改的  $cs$  通知  $A$ ;

② 如果  $D$  没收到有到  $A$  的合同  $cs$  信息,  $D$  可以生成  $cs$  并通知  $A$ ;

③ 如果  $D$  收到  $A$  的 ( $INT A cs$ ), 认为  $A$  同意  $cs$ ;

直到相互满意或达不成一致意见为止, 如果是后者, 则需考虑其他候选者.

(9) 社会承诺建立,

满足以下两个条件,  $D$  就建立起对应的社会承诺:

① 当  $D$  发出  $cs$  信息, 且收到 ( $INT A cs$ ), 这时有 ( $INT D cs$ );

② 当  $D$  收到  $cs$  信息, 且作出 ( $INT D cs$ ), 并收到 ( $INT A cs$ ).

对于  $A$  来说, 基本上是对  $D$  社会承诺建立的反过程.

### 2.2 D 的社会承诺解除

情形 1. ( $BEL D \neg (GOAL D \varphi b)$ ),  $D$  原有的实现  $\varphi$  的目标不存在.

(1) 如果  $A$  和  $D$  在一个组内, 或  $cs$  不在合同期限内, 那么  $D$  可以解除承诺, 同时告知  $A$  及原应 (所有承诺解除原则上要求通知对方, 并附上原应);

(2) 如果  $cs$  在期限内, 且解除承诺带来的罚金  $\leq$  报酬, 那么  $D$  可以解除承诺, 同时支付相应的罚金, 否则, 罚金  $>$  报酬  $D$ , 需根据问题领域来作决定.

情形 2. ( $BEL D (GOAL D \varphi bn) \wedge (bn < b)$ ),  $D$  原有目标存在, 但是收益要比起初相信的小.

(1) 如果  $A$  和  $D$  在一个组内, 且  $bn < c$ , 那么解除承诺, 否则  $bn > c$ , 可以继续;

(2) 如果  $A$  和  $D$  不在一个组内,  $cs$  不在期限内, 且  $bn < r$ , 那么  $D$  有两种选择: 解除承诺或减少  $r$ , 这需要与  $A$  协商, 否则  $bn > r$  可以继续;

(3) 如果  $A$  和  $D$  不在一个组内,  $cs$  在期限内,  $bn < r$ , 且  $bn + \text{罚金} < r$ , 那么  $D$  也有类似的两种选择, 否则  $bn + \text{罚金} > r$  可以继续.

情形 3. 接收  $A$  发来的协商, 要么是解除承诺, 要么是增加报酬支付.

计算解除承诺和增加支付的效用, 选择效用大者, 必要时修改  $cs$ .

情形 4. 接收  $A$  发来的解除承诺信息,  $D$  相应解除.

类似地, 有对  $A$  的社会承诺解除.

### 3 例子

我们以卡车运输调度问题为例来说明社会承诺的建立和解除过程.

一个车场 Agent  $D$  接收到用户的任务  $P$ , 对应的约束包括: 责任 (最迟完成时间  $T$ )、禁止 (打开货物)、允许 (分批运输), 完成任务的报酬是 1800 元. 根据目前的情况,  $D$  难以完成  $P$ , 产生 ( $GOAL D (Contract P)$ ),  $D$  知道存在可能的车场  $A1$  (属于同一公司) 和  $A2$  (不属于同一公司), 这时向他们发送 ( $Contract D P$ ), 1500 元, 10 分钟).

10 分钟内收到: ( $Can-Adopt A1 P 1500$  元 代价: 1700 元) 和 ( $Can-Adopt A2 P 1600$  元 代价: 0 元),  $D$  计算得: 委托  $A1$  的效用是 100 元, 委托  $A2$  的效用是 200 元, 这时选择  $A2$  作为委托对象.

经过协商产生如下的合同  $C1$ .

(合同号: 1, ( $D, A2$ ),

( $D$ , (责任:  $A2$  完成  $P$ , 支付报酬 1600 元, 报酬 0 元)

(禁止:非特殊情况,禁止终止合同,罚金 200 元)

(允许:意外情况,与  $A_2$  通信)...)

( $A_2$ , (责任:条件正常,完成  $P$ ,报酬 1600 元)

(禁止:运输中禁止打开货物,罚金 200 元)

(禁止:非特殊情况禁止终止合同,罚金 200 元)

(允许:意外情况,与  $D$  通信)...)

(有效期: $T_1$  到  $T_2$ )

这时有( $INT\ D\ C_1\ 200$  元)和( $INT\ A_2\ C_1\ 300$  元),并且作出对应的社会承诺。

当  $D$  发现原有的报酬是 1500 元时,根据  $1500 + 200 > 1600$  保持合同的继续执行。

应用 Haddadi 的社会承诺机制则难以反映该例中的社会约束和效用决策。

#### 4 结 语

Castelfranchi<sup>[2]</sup>提出了社会承诺的概念,讨论了社会承诺的性质,但是没有给出社会承诺的建立和解除机制,也没有考虑效用对社会承诺的影响。虽然也意识到了社会承诺中双方的责任等方面的约束,但是没给出合同结构的概念,利用合同结构来描述双方在社会承诺中的约束。

S. Kraus<sup>[7]</sup>提出了在 MAS 中激励承包的思想和在不同情况下计算报酬的方法,从效用理论的观点出发,考虑如何计算合适的报酬,才能使得其他 Agent 采纳其委托的目标,但是没有考虑社会承诺的概念以及与 BDI 理论的结合问题。

A. Haddadi<sup>[1]</sup>基于 Rao & Georgeff 的 BDI 理论,研究了组内 Agent 之间的社会承诺建立机制,提出了合作可能、预承诺、承诺的概念以及相应的通信策略问题。但他们的工作只局限于  $D$  和  $A$  在一个合作组内。在这种情况下,Agent 是非自利的,不需要效用理论来调节,可是  $D$  可以将  $P$  委托给非组内的  $A$  来完成,这样,他的工作就具有局限性。另外,他们没有考虑双方的责任等约束。本文通过合同的概念,将  $D$  和  $A$  连接起来,则更符合人类承包现象。

本文提出了合同结构的概念,用来联系社会承诺的双方,并给出了社会承诺的定义——社会承诺是双方的,并且意图去履行合同;对 BDI 和效用理论的结合进行了尝试;从 Agent 的角度来分析社会承诺机制,所给描述类似于算法形式,因而具有一定的可操作性。

**致谢** 本文的研究工作得到国家自然科学基金资助,项目编号为 69773026 和 69733020。另外,本文得到上海铁道大学施鸿宝教授的指导,在此表示感谢。

#### 参 考 文 献

- 1 Haddadi A S. Communication and Cooperation in Agent Systems. Berlin: Springer-Verlag, 1996. 1~134
- 2 Castelfranchi C. Commitments: from individual intentions to groups and organizations. In: Victor Lesser ed. Proceedings of the 1st International Conference on Multi-agent Systems. Cambridge, MA: AAAI Press/MIT Press, 1995. 41~48
- 3 Castelfranchi C. Modeling social action for AI agents. In: Pollack M E ed. Proceedings of the 15th International Joint Conference on Artificial Intelligence. San Mateo, CA: Morgan Kaufmann Publishers, 1997. 1567~1576
- 4 Dunin-Keplicz B, Verbrugge R. Collective commitments. In: Durfee F ed. Proceedings of the 2nd International Conference on Multi-agent Systems. Menlo Park, California: AAAI Press, 1996. 56~63
- 5 Rao, Georgeff. Modeling rational agents within a BDI-architecture. In: Llen J A, Fikes R, Sandewall W eds. Proceedings of the Principles of Knowledge Representation and Reasoning. San Mateo, CA: Morgan Kaufmann Publishers, 1991. 473~484
- 6 Rao, Georgeff. An abstract architecture for rational agents. In: Nebel B, Rich C, Swartout W eds. Proceedings of the Principles of Knowledge Representation and Reasoning. San Mateo, CA: Morgan Kaufmann Publishers, 1992. 439~449
- 7 Kraus S. An overview of incentive contracting. Artificial Intelligence, 1996, 83(2): 297~346

## One Mechanism of Social Commitment Based on Belief-Desire-Intention and Utility

XU Jin-hui SHI Chun-yi

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

**Abstract** The social commitment between agents needs to be built for solving problem in the multi-agent systems. A mechanism of social commitment by combining belief-desire-intention and utility theory is presented in this paper, which guarantees agent's rational action on logic and decision. The authors combine with Rao & Georgeff's Belief-Desire-Intention theory and S. Kraus's idea of incentive contract, improve Haddadi's mechanism of social commitment, and provide Castelfranchi's concept of social commitment with implementing support is provided in this paper.

**Key words** Multi-agent systems, belief, desire, intention, utility, social cc © 中国科学院软件研究所 <http://www.jos.org.cn>