

纽科姆难题与囚徒困境的关系新探

丁雨辰

(南京大学 哲学学院,江苏 南京 210023)

摘要:刘易斯曾通过引入“完美复制人”策略,论证纽科姆难题与囚徒困境是同一类问题的不同表现形式,但学界对此存在颇多争议。运用逻辑悖论构成的“三要素”理论重新考察纽科姆难题与囚徒困境,可以进一步澄清两者能够依据相同的“交互信念”要素而构成逻辑悖论,并在逻辑结构上具有同构性。据此揭示传统决策理论与博弈论处理囚徒困境时所隐含的理论预设,有助于正确认识两种理论的功能与局限,进而通过引入情境量化域转变思想,可为化解此类悖论提供统一思路。

关键词:纽科姆难题;囚徒困境;逻辑悖论;情境量化域转变

中图分类号:B81

文献标志码:A

文章编号:1672-7835(2025)04-0045-11

纽科姆难题(newcomb problem)和囚徒困境(prisoners' dilemma)^①是传统决策理论、博弈论与博弈逻辑研究的两大关键疑难。刘易斯(D. Lewis)曾通过引入“完美复制人”策略,试图论证纽科姆难题和囚徒困境是同一个问题的不同表现形式^②,但此后学界仍就两者是否为同一问题展开了持续争论。反对者认为刘易斯的论证改变了纽科姆难题的基本结构,而囚徒困境与原初的纽科姆难题具有本质上的不同;支持者则认为反对者对刘易斯论证的反驳并不成立,并为“两问题本质相同”的观点辩护。本文试图借助逻辑悖论构成的“三要素”理论,从悖论研究的视角讨论问题之症结,为刘易斯的基本观点做出新的辩护,并探讨这种讨论的重要意义。

一、关于纽科姆难题和囚徒困境关系的论争

纽科姆难题由诺奇克(R. Nozick)于1969年提出。诺奇克称其源自物理学家纽科姆(W. Newcomb)在思考囚徒困境时的一个构造,该问题的最初版本可以描述如下^③。

一个超级生物可以正确预测你的选择结果,同时你也知道其具有这种正确预测能力。假设你处在如下情境中:现在有两个盒子,Box1 和 Box2,此时你有两种选择:(1)只拿 Box2;(2)同时拿 Box1 和 Box2。其中 Box1 中一直放有1 000 美元;当超级生物预测到你只拿 Box2 时,其会在你拿取之前,在 Box2 中放入100 万美元。那么,什么样的行动才是更合理的?支持只拿 Box2 的人认为,如果你要选择同时拿 Box1 和 Box2,那么超级生物会预测到,并不在 Box2 中放钱,因此收益是1 000 美元。而只拿 Box2 的收益是1 000 000 美元。根据做期望效用更大的选择的 MEU(maximum expected utility)原则,此时应只拿 Box2。而支持同时拿 Box1 和 Box2 的人认为,因为超级生物的预测与放钱都是在你实际进行选择之前,因此 Box2 带来的收益是确定的,则对你来说,只拿 Box2 的收益是 M 美元,而同时拿 Box1 和

收稿日期:2025-03-14

基金项目:国家社会科学基金重大项目(18ZDA031)

作者简介:丁雨辰(1992—),男,宁夏石嘴山人,博士生,主要从事逻辑哲学研究。

①“纽科姆难题”在英文中主要有“newcomb problem”“newcomb's problem”两种的不同写法,“囚徒困境”则主要有“prisoners' dilemma”“prisoner's dilemma”“prisoner dilemma”三种的不同写法,本文采用刘易斯1979年论文标题中的用法。

②Lewis D. “Prisoners' Dilemma is a Newcomb Problem”, *Philosophy & Public Affairs*, 1979, 8(3):235-240.

③Andson A R, Benacerraf P, Grünbaum, etc. *Essays in Honor of Carl G. Hempel*. Dordrecht: D. Reidel Publishing Company, 1969, pp.114-146.

Box2 的收益是 $M+1\ 000$ 美元。根据在不同情况下都做占优选择的 DP (dominant principle) 原则,此时应只拿 Box2。这样,两种不同的合理行动原则在此情境下出现了矛盾。随着贝叶斯框架被引入决策理论,诺奇克又在新理论框架内采用了证据决策理论 (evidential decision theory) 和因果决策理论 (causal decision theory),来替代 MEU 原则和 DP 原则作为合理行动原则发挥的作用。最终,诺奇克本人采取了一种折中主义的立场,认为人们会综合证据决策理论和因果决策理论,根据自身经验在实际选择中为两个原则分配权重。

囚徒困境则是在 20 世纪 50 年代即已提出的一个关于非零和又非合作博弈的经典难题^①:两个合伙作案的嫌疑人被警察逮捕,警察将他们分开审讯,两人都有不招供和招供两种选择,并有如下情形成立:(1)若两人都招供,则两人都被判坐 5 年牢;(2)若两人都不招供,则因证据不足,两人都坐 1 年牢;(3)若两人中只有一人招供而另一人不招供,则招供者因立功表现被释放,不招供者坐 10 年牢。上述信息均为两个嫌疑人所知晓,其以坐牢年份衡量的收益矩阵如表 1 所示。

表 1 经典囚徒困境收益矩阵

		嫌疑人 B	
		不招供	招供
嫌疑人 A	不招供	(-1, -1)	(-10, 0)
	招供	(0, -10)	(-5, -5)

如果 B 选择不招供,那么对于 A 而言,不招供的收益是 -1,招供的收益是 0,此时理性的选择是招供。如果 B 选择招供,那么 A 不招供的收益是 -10,招供的收益是 -5,此时也同样应该选择招供。显然,对于 A 来说“招供”是一个占优策略,按照占优(DP)原则,他应该选择“招供”。对于 B 也可由类似过程得出“招供”是一个占优策略。但这种个体理性的选择最终导致 (-5, -5) 的集体非最优结果,违背了最大期望效用(MEU)原则。要实现集体最优的(-1, -1)结果,必须通过事前的“合作”承诺。

表面上看,纽科姆难题是单个主体在特定条件下的选择困境,而囚徒困境涉及两个主体选择的相互影响,似乎有根本性差异。刘易斯的论证策略是,将纽科姆难题中引起诸多质疑的“始终能够正确预测的超级生物”消除,替换为一个局中人 A 的“完美复制人”B,以其代替超级生物的“预测”在纽科姆难题中的作用,同时其也可以取代局中人 B 在囚徒困境中的位置。这样,只需要考察 A 的完美复制人的选择,就可以对 A 的选择进行提前预测^②。而将表 1 中的坐牢年份替换为纽科姆难题中的金钱收益,更容易地看出其间的关系(见表 2)。

表 2 刘易斯型囚徒困境收益矩阵

		嫌疑人 B	
		不招供	招供
嫌疑人 A	不招供	(100 万, 100 万)	(0, 1 001 000)
	招供	(1 001 000, 0)	(1 000, 1 000)

刘易斯认为使用“完美复制人”策略并没有改变纽科姆难题本身,若承认在使用“完美复制人”策略后的类囚徒困境中应选择“不招供”,则在原本的纽科姆难题中就应当选择“只拿 Box2”;反之,若承认在使用“完美复制人”策略后的类囚徒困境中应选择“招供”,则在原本的纽科姆难题中就应当选择“同时拿 Box1 和 Box2”。鉴于经典囚徒困境的纳什均衡解是(招供,招供),按照相同的思路,使用“完美复制人”策略后的纽科姆难题的对应解也应当是(拿 Box1 和 Box2,拿 Box1 和 Box2)。

由此可见,使用“完美复制人”策略后纽科姆难题的收益矩阵,与经典囚徒困境收益矩阵在结构上相同。刘易斯认为原本的纽科姆难题、使用“完美复制人”策略后的纽科姆难题与经典囚徒困境都具有

^①Straffin P D. *Game Theory and Strategy* (vol.36). Washington D.C.: Mathematical Association of America, 1993, pp.73~80.

^②Lewis D. “Prisoners’ Dilemma is a Newcomb Problem”, *Philosophy & Public Affairs*, 1979, 8 (3): 235~240.

如下相同的前提:

- (1) 选择主体具有选择任意行动的自由;
- (2) 选择主体有可能获得最大期望效益,但这与主体的选择没有因果关系。

但其不同之处在于,在原本的纽科姆难题中:

- (3') 局中人 A 会获得 100 万,当且仅当局中人 A 被预测没有拿 1 000。

在表 2 以金钱收益衡量的刘易斯型囚徒困境和如表 3 使用“完美复制人”策略后纽科姆难题中则是:

- (3'') 局中人 A 会获得 100 万,当且仅当局中人 B 没拿 1 000。

于是问题就变成了“局中人 A 被预测没拿 1 000”和“局中人 B 没拿 1 000”是否等同。刘易斯认为可以通过“完美复制人”替换策略说明,在纽科姆难题中“局中人 A 被预测没有拿 1 000”就是“局中人 A 的完美复制人没拿 1 000”。因而刘易斯论证得出,纽科姆难题和囚徒困境是一个问题的不同表述。在使用“完美复制人”策略后的纽科姆难题具有如下类似于囚徒困境的收益矩阵(见表 3)。

表 3 使用“完美复制人”策略后纽科姆难题的收益矩阵

		局中人 B(局中人 A 的完美复制人)	
		只拿 Box2	拿 Box1 和 Box2
局中人 A	只拿 Box2	(100 万, 100 万)	(0, 1 001 000)
	拿 Box1 和 Box2	(1 001 000, 0)	(1 000, 1 000)

索贝尔(J.H. Sobel)承认如刘易斯所言,某些具体的囚徒困境和纽科姆难题例子可以满足将两个问题视为一个的看法,但他对所有囚徒困境都是纽科姆难题这一论述予以质疑^①。他认为,依据使用“完美复制人策略”的纽科姆难题的条件,要想确定某种特定选择,则需要满足:

(3_{sn}) 在局中人 A 选择只拿 Box2 的条件下局中人 A 的完美复制人选择只拿 Box2 的概率是 1,且在局中人 A 选择拿 Box1 和 Box2 的条件下局中人 A 的完美复制人选择拿 Box1 和 Box2 的概率是 1。

而在刘易斯使用“完美复制人”策略后的纽科姆难题中,则需要满足:

(3_{sp}) 局中人 A 选择只拿 Box2 并且局中人 A 的完美复制人选择只拿 Box2 或局中人 A 选择拿 Box1 和 Box2 并且局中人 A 的完美复制人选择拿 Box1 和 Box2 的概率是 1。

索贝尔指出,存在(3_{sp})为真但(3_{sn})为假的情况。他给出了双胞胎例子:一对双胞胎做出相同选择的概率非常高,其中任何一位会认为自己和对方都会做出某种特定选择,且知道自己和对方做出这种选择的理由是不同的。假设双胞胎都只会拿 Box2,在囚徒困境式选择中可以体现为两位双胞胎选择的结果是相同的,即(3_{sp})为真。但由于不同理由,在双胞胎 A 选择拿 Box1 和 Box2 的条件下,双胞胎 B 选择拿 Box1 和 Box2 的概率可以是不确定的。因此,索贝尔认为由于双胞胎做出选择的理由不同,在满足(3_{sp})的情况下(3_{sn})不一定能得到满足。这种情况的产生源于两位局中人存在认知不对称(epistemic asymmetry)。由此他得出结论,当(3_{sp})和(3_{sn})同真时,囚徒困境和使用“完美复制人”策略后的纽科姆难题相同,当(3_{sp})为真(3_{sn})为假时,此时的囚徒困境不再是纽科姆难题。

与索贝尔承认有一部分囚徒困境与纽科姆难题同构不同,贝穆德斯(J.L. Bermúdez)则从主体认知情形出发,认为引入“完美复制人”方案会致使新构建出的纽科姆难题与囚徒困境中主体所面临认知情形不一致,致使囚徒困境不会成为纽科姆难题^②。

贝穆德斯指出,在刘易斯的论证中即便(3')和(3'')同真,但由于存在被忽略的第四个条件,导致仍会存在需要探讨的情况。他将此缺少的条件补充为:

^①Campbell R, Sowden L. *Paradoxes of Rationality and Cooperation: Prisoner's Dilemma and Newcomb's Problem*. Vancouver: University of British Columbia Press, 1985, pp.263–274.

^② 参见 Cf. Bermúdez J L. “Prisoner's Dilemma and Newcomb's Problem: Why Lewis's Argument Fails”, *Analysis (Oxford)*, 2013, 73 (3): 423–429; Peterson M. *The Prisoner's Dilemma*. Cambridge: Cambridge University Press, 2015, pp.115–132.

- (4')局中人A对(3')具有高度信心；
 (4'')局中人A对(3'')具有高度信心。

在贝穆德斯看来,想要说明(4')和(4'')等价,需要借助如下条件才能予以证明:

- (5)局中人A对预测“局中人A没拿1000,当且仅当局中人B没拿1000”具有高度信心。

然而,在囚徒困境中,(5)为真会导致局中人的选择不再满足(1)。比如回到表1来看,就只应该存在选择(招供,招供)和(不招供,不招供)两种结果。因此,纽科姆难题中的主体认知情形会产生一种有别于囚徒困境的认知情形,实际表现为(5)和(1)产生的矛盾,因此使用“完美复制人”策略后的纽科姆难题不是囚徒困境,任何试图将两者等同起来的做法从一开始就注定要失败。

为了维护自己的论点,贝穆德斯后续引入了参数刻画方案和策略刻画方案,来对使用“完美复制人”策略后的纽科姆难题和囚徒困境进行区分。他认为适用于参数刻画方案需要满足仅有主体的选择一个自由变量,而其他参数则由环境设定的情况,因而使用“完美复制人”策略后的纽科姆难题适用于此方案。而策略刻画方案则是用来刻画囚徒困境这种至少涉及两名局中人的情形,而局中人的选择是否理性,取决于“其他局中人的选择——而这些选择的理性程度又部分取决于该主体的选择”等因素^①。

沃克(M.T. Walker)认为,即便贝穆德斯的上述分析很有意义,但其提出的反驳并不充分^②。他指出,虽然选择结果只剩下了(招供,招供)和(不招供,不招供)两个选择,但按照贝穆德斯的主张,在思考之前,局中人仍然应该将四种结果都视为可能的结果。而贝穆德斯的结论只能表明,在思考的过程中可能存在“去自由化(de-liberation)”,即局中人在思考过程中,发现自己可能理性地自由做出的选择范围逐渐缩小。因而,贝穆德斯并没有对刘易斯论证本身做出真正的反驳。

进而,沃克区分了两种条件句:决策行动条件句(d-conditionals)和附加条件句(a-conditionals),他认为决策行动条件句本身就是关于决策行动的表征,相较于附加条件句在因果链条上的回溯性更强。在原本的纽科姆难题中,主要使用的是决策行动条件句,而在囚徒困境中则使用附加条件句。沃克据此指出,刘易斯观点的主要问题在于混淆了这两种条件句,把并不具有强回溯能力的附加条件句当成了决策行动条件句使用,这才是刘易斯把两者视为同一个问题的关键。同时,他也为常用于博弈论中的纳什均衡为什么没有在纽科姆难题中被使用进行了辩护。他认为一旦决定采取某种行动,纳什均衡的说服力就很容易被低估,而纽科姆难题中的局中人采取的都是决策行动条件句,从而导致纳什均衡不会被采纳。

支持刘易斯观点并为之做出新的辩护的有巴蒂奇(E. Badici)和张智皓等学者。巴蒂奇认为贝穆德斯提出的反驳基于对刘易斯论证的误解^③,张智皓则比巴蒂奇的工作更进一步,为刘易斯的观点做了更为细致的辩护。他说明,贝穆德斯提出的(5)并不会如其所言,最终导致囚徒困境中的两位局中人都只会做出相同的选择^④。根据刘易斯的论证,(5)实际应表示为:

- (5')局中人A对“局中人A的完美复制人没拿1000,当且仅当局中人B没拿1000”具有高度信心。

而在使用“完美复制人”策略后的纽科姆难题中,局中人B就是指“局中人A的完美复制人”,所以这是一个永真的命题。而且这是A对其“完美复制人行动”的信念,与A是否真的采取行动无关。进而,在刘易斯的方案中,局中人的选择不会是相同的,仍能保证(1)成立。张智皓认为,贝穆德斯的论证想要成立,还需补充如下条件:

^①参见 Bermúdez J L. “Strategic vs. parametric choice in Newcomb’s problem and the prisoner’s dilemma: Reply to Walker”, *Philosophia*, 2015, 43(3): 787–794; Bermúdez J L. “Does Newcomb’s problem really exist”, *Newcomb’s Problem*. Ahmed A. (Ed), Cambridge: Cambridge University Press, 2018, pp.19–41.

^②Walker M T. “The Real Reason Why the Prisoner’s Dilemma is Not a Newcomb Problem”, *Philosophia*, 2014, 42(3): 841–859.

^③Badici E. “Prisoner’s Dilemma and Newcomb’s Problem: Two Problems Or One”, *Philosophia*, 2023, 51(5): 2543–2557.

^④张智皓:《纽康难题是另一种囚犯困境吗?》,《逻辑学研究》2017年第1期。

(6)局中人 A 对“局中人 A 没拿 1 000,当且仅当局中人 B 没拿 1 000”具有高度信心。

这在刘易斯的“完美复制人”语境中就是:

(6')局中人 A 对“局中人 A 没拿 1 000,当且仅当局中人 A 的完美复制人没拿 1 000”具有高度信心。

根据“完美复制人”策略本身的含义,其在纽科姆难题和囚徒困境中都为真。而张智皓又进一步表明,(6)成立实际建立在 EDT 的基础上,而在 CDT 中(6)实际是与决策无关的。在采取 CDT 的观点时,使用“完美复制人”策略后的纽科姆难题和囚徒困境则仍面对相同的结构。在采取 EDT 的立场后,使用“完美复制人”策略后的纽科姆难题和囚徒困境都会消除两种选择结果,也具有相同的结构。

由以上评述可以见得,参与有关争论的学者多从传统决策理论与博弈论的视角展开分析,以各自的理论前提为基础,建构对不同理论框架的辩护。既往争论的焦点在于对“完美复制人”策略与“行动-结果”关系的思考。但综合来看,他们的思考更多采取的是各自预设的理论视角,而没有从逻辑悖论研究的视角做出系统分析。本文试图表明,系统使用逻辑悖论的“三要素”理论,有助于我们对有关问题进行深入分析,发现其中蕴藏的问题。

二、作为逻辑悖论的囚徒困境及其化解路径

随着关于悖论研究的一般认识论与方法论研究的兴起,国内外学界已基本上就逻辑悖论构成的“三要素”达成了共识。人们已不再将悖论归结为单纯的悖论性语句或简单的矛盾等价式,而将逻辑悖论视为由“三要素”系统构成。目前国际学界使用较多的塞恩斯伯里(R.M. Sainsbury)定义:“从明显可接受的前提,经过明显可接受的推理,得出明显不可接受的结论”^①,可视为“三要素”的明确标识。但该定义的广泛可接受性源自“明显可(不可)接受”的高度模糊性,其他类似的“三要素”定义可视为其在不同程度上的精致化^②。本文所使用的张建军的逻辑悖论定义:从“公认正确的背景知识”出发,经过“严密无误的逻辑推导”,可以“建立矛盾等价式”^③,即可视为塞恩斯伯里定义的一种精致化。其中,“公认正确的背景知识”表明了逻辑悖论的基本的相对性特征,即相对于将悖论的“前提”要素作为公共信念的认知共同体而言,亦呈现出悖论的语用特征。由此引出“矛盾等价式”必须经过“严密无误的逻辑推导”,即有效演绎,而不只是“明显可接受”,这表明悖论的严格刻画必须明确建构矛盾等价式的所有前提,以及其所依赖的所有推理法则,而这些推理法则也在“公认正确的背景知识”或“认知共同体的公共信念”之中,即为认知共同体成员普遍接受。近来孔斯又在合理行动悖论研究中强调了博弈论中关于 common belief 与 mutual belief 这两种公共信念的差别,认为悖论的第一要素应为后者(包括认知共同体成员都在使用但并不一定明确认知的预设),而不一定是更强的前者(认知共同体成员都具有明确共识的信念)^④。我们认为,明确这一点对于揭示形成悖论的隐含预设至关重要。因而本文后续将第一要素称为认知公共体的“交互信念”。

囚徒困境之所以不被部分博弈论研究者视为悖论,是因为他们认为局中人应当基于独立决策原则行动,不应考虑其他外部因素。作为典型的“非合作博弈”模型,囚徒困境中由帕累托改进所引发的个体理性与集体理性冲突,在标准博弈论框架下被认为不属于讨论范畴。因此,囚徒困境本身不会产生逻辑矛盾——其纳什均衡解正是基于这种严格的“非合作”前提得出的理论结果。然而,这种观点实际上混淆了问题本身与解决方案两个不同层面。即便我们暂且不考量囚徒困境经过帕累托改进后显现的个

^①R. M. 塞恩斯伯里:《悖论》,刘叶涛、雒自新、冯立荣译,中国人民大学出版社 2020 年版,第 1 页。

^②就笔者视域所及,关于“三要素”(three-elements)这个名称的使用,西方学界最早出现在 2003 年出版的“哲学中心问题”丛书中由英国女逻辑学家奥林执笔的《悖论》一书之中(Olin D. Central Problems of Philosophy: Paradox. Bucks: Acumen Publishing Limited, 2003, p. 6);我国学界则由张建军 1992 年在《悖论的逻辑与方法论问题》一文中首次出现(载于张建军,黄展骥:《矛盾与悖论研究》,香港:黄河文化出版社 1992 版,第 58 页),而“三要素”界说所体现的悖论的语用学性质(认知共同体相对性),则是张建军在 2001 发表的《论作为语用学概念的逻辑悖论——兼复马佩先生》一文(载于《江海学刊》2001 年第 6 期)才得以明确指认的。

^③张建军:《逻辑悖论研究引论(修订本)》,人民出版社 2014 版,第 4—10 页。

^④罗伯特·C.孔斯:《信念悖论与策略合理性》,张建军等译,中国人民大学出版社 2020 版,第 199 页。

体与集体理性冲突,同样可以在囚徒困境中通过揭示其中的“交互信念”而构建出矛盾等价式。

(一) 囚徒困境何以成为逻辑悖论

基于普莱尔(A.N. Prior)刻画意外考试悖论的思路^①,将时间点用(思考的)行动阶段进行替换,以行动阶段为两轮的囚徒困境为例:

设 K_x 表示“A在(自己思考的)第x行动阶段知道”; H_y 表示“A知道B在(A思考的)第y行动阶段确定选择”; $x-1, y-1$ 表示“在x,y的前一行动阶段”。其中x的定义域(量化域)是{行动未开始的阶段(z), 第一行动阶段(m), 第二行动阶段(t)}; y的定义域(量化域)是{第一行动阶段(m), 第二行动阶段(t)}。

行动阶段定义:A行动需要建立在其知道B采取什么样的行动,A产生一个(思考)行动选择结果则定义为一阶段。如,假设A知道B招供是第一行动阶段,假设A知道B的选择后最终决定招供则是第二阶段。

在如上定义下,形成囚徒困境的DP原则与MEU原则依据的前提(交互信念)可刻画如下:

$(N_1) \sim H_m \rightarrow H_t$ (如果A知道B没有在第一行动阶段确定选择,那么A知道B会在第二行动阶段确定选择。)

$(N_2) H_t \rightarrow \sim H_m$ (如果A知道B在第二行动阶段确定选择,那么A知道B没有在第一行动阶段确定选择。)

$(N_3) \sim H_m \rightarrow K_m(\sim H_m)$ (如果A知道B在第一行动阶段没有确定选择,那么A知道自己知道B在第一行动阶段没有确定选择。)

$(N_4) H_y \rightarrow \sim K_{y-1}(H_y)$ (如果A知道B在行动阶段y后确定选择,那么A在行动阶段y的前一阶段不知道自己知道B在行动阶段y后确定选择。)

除经典逻辑中的推理规则外,再加入如下两个规则:

(L_1) 如果“ $\alpha \rightarrow \beta$ ”是系统定理,那么 $(\alpha \rightarrow \beta) \rightarrow (K_x \alpha \rightarrow K_x \beta)$;

(L_2) 如果“ α ”是系统定理,那么“ $\alpha \rightarrow K_z(\alpha)$ ”。

根据上述前提(“交互信念”),我们使用经典逻辑法则做如下推理:

- | | |
|---|-----------------|
| (1) $(H_t \rightarrow \sim H_m) \rightarrow ((\sim H_m \rightarrow K_m(\sim H_m)) \rightarrow (H_t \rightarrow K_m(\sim H_m)))$ | PC 定理、US |
| (2) $(\sim H_m \rightarrow K_m(\sim H_m)) \rightarrow (H_t \rightarrow K_m(\sim H_m))$ | (1)、 N_2 、MP |
| (3) $H_t \rightarrow K_m(\sim H_m)$ | (2)、 N_3 、MP |
| (4) $(\sim H_m \rightarrow H_t) \rightarrow (K_x(\sim H_m) \rightarrow K_x(H_t))$ | N_1, L_1 |
| (5) $K_x(\sim H_m) \rightarrow K_x(H_t)$ | $N_1, (4)$ 、MP |
| (6) $(H_t \rightarrow K_m(\sim H_m)) \rightarrow ((K_m(\sim H_m) \rightarrow K_m(H_t)) \rightarrow (H_t \rightarrow K_m(H_t)))$ | PC 定理、US |
| (7) $(K_m(\sim H_m) \rightarrow K_m(H_t)) \rightarrow (H_t \rightarrow K_m(H_t))$ | (3)、(6)、MP |
| (8) $K_m(\sim H_m) \rightarrow K_m(H_t)$ | (5)、行动阶段定义 |
| (9) $H_t \rightarrow K_m(H_t)$ | (7)、(8)、MP |
| (10) $H_t \rightarrow K_{t-1}(H_t)$ | (9)、行动阶段定义 |
| (11) $(H_y \rightarrow \sim K_{y-1}(H_y)) \rightarrow (\sim \sim K_{y-1}(H_y) \rightarrow \sim H_y)$ | PC 定理、US |
| (12) $(H_y \rightarrow \sim K_{y-1}(H_y)) \rightarrow (K_{y-1}(H_y) \rightarrow \sim H_y)$ | (11)、PC 定理 |
| (13) $K_{y-1}(H_y) \rightarrow \sim H_y$ | (12)、 N_4 、MP |
| (14) $((H_t \rightarrow K_m(H_t)) \rightarrow ((K_m(H_t) \rightarrow \sim H_t) \rightarrow (H_t \rightarrow \sim H_t)))$ | PC 定理、US |
| (15) $(K_m(H_t) \rightarrow \sim H_t) \rightarrow (H_t \rightarrow \sim H_t)$ | (9)、(14)、MP |
| (16) $K_m(H_t) \rightarrow \sim H_t$ | (13)、行动阶段定义 |
| (17) $H_t \rightarrow \sim H_t$ | (15)、(16)、MP |

^①Prior A N. "The Paradox of the Prisoner in Logical Form", *Synthese*, 2012, 188(3):411–416.

(18) $(H_t \rightarrow \sim H_t) \rightarrow \sim H_t$	PC 定理、US
(19) $\sim H_t$	(17)、(18)、MP
(20) $(\sim H_m \rightarrow H_t) \rightarrow (\sim H_t \rightarrow \sim \sim H_m)$	PC 定理、US
(21) $(\sim H_m \rightarrow H_t) \rightarrow (\sim H_t \rightarrow H_m)$	(20)、PC 定理
(22) $\sim H_t \rightarrow H_m$	N_1 、(21)、MP
(23) H_m	(19)、(22)、MP
(24) $H_m \rightarrow K_z(H_m)$	L_2
(25) $K_z(H_m)$	(23)、(24)、MP
(26) $K_{m-1}(H_m) \rightarrow \sim H_m$	(13)、行动阶段定义
(27) $K_z(H_m) \rightarrow \sim H_m$	(26)、行动阶段定义
(28) $\sim H_m$	(25)、(27)、MP

(23) 和 (28) 矛盾,从而可以在经典逻辑法则下建构矛盾等价式。依据“三要素”理论所揭示的悖论的相对性,相对于所有以前提 N_1-N_4 、 L_1 和 L_2 规则以及经典逻辑法则为“交互信念”的认知共同体而言,上述囚徒困境就构成一个符合“三要素”定义的逻辑悖论。

(二)作为逻辑悖论的囚徒困境的化解

在讨论纽科姆难题与囚徒困境的同构关联之前,我们可以先依据上述对囚徒困境悖论的严格刻画,探讨该悖论的一种化解路径。如果这种化解路径同样也可化解纽科姆难题,即可作为二者“本质上是同一个问题”的一种佐证。

通过对上述囚徒困境作为悖论推导过程的回溯分析,我们首先可发现规则 L_2 会导致过强的逻辑全能问题(logical omniscience problem)。这一强规则要求认知主体能够知道所有系统定理,显然不符合现实认知能力的有限性特征。借鉴达克(H.N. Duc)的认知悖论修正方案^①,我们可对 L_2 进行弱化处理以贴近日常认知直觉,得到改进后的规则 L'_2 :如果“ α ”是系统定理,那么“ $\alpha \rightarrow \Diamond K_z(\alpha)$ ”。这种改进既保留了原系统的核心推导能力,又避免了过度理想化的认知假设,使模型更贴合实际决策情境中有限理性主体的认知特征。原推导(24)、(25)变为:

(24') $H_m \rightarrow \Diamond K_z(H_m)$	L'_2
(25') $\Diamond K_z(H_m)$	(23)、(24)、MP

通过改进规则 L'_2 ,使得矛盾初步在现实世界 w_0 中得到化解。 $\Diamond K_z(H_m)$ 为真这一前提,实际是占优策略原则(DP)和期望效用最大化原则(MEU)得以适用的基本认知条件。其理论依据在于:这两个决策原则都要求“招供”和“不招供”这两种纯策略本身具有认知可能性,也就是刘易斯提出两个问题的共同前提(1)。具体而言,对局中人来说,在每个行动阶段都必须满足 $\Diamond K_z(H_m)$ 这一条件,即存在认知上可能知道对方选择的可能世界。 $\Diamond K_z(H_m)$ 成立,则表明嫌疑人完全可以找到一个现实世界可通达的可能世界 w_1 ,在 w_1 中 $K_z(H_m)$ 成立,这是假设性思考之所以可行的条件。但在这个可能世界中,我们就会发现对局中人来说,其完全可能再次得到矛盾 $H_m \wedge \sim H_m$ 。

再借助学界最近在讨论“实质蕴涵怪论”过程中提出的“情境量化域转变”的思想^②,可以更好地说明矛盾再次产生的原因。我们可以在 N_1-N_4 及相应规则基础上,设定一个可能世界模型 \mathcal{M} 为一个四元组 $\langle W, R, V, S_c \rangle$,来更好地表示囚徒困境得以构成的“交互信念”。其中 W 表示可能世界集合 $\{\dots, w_i, w_{i+1}, \dots\}$, R 表示可能世界间可及关系, V 为赋值, S_c 为主体 c 面临的情境量化域集合 $\{\dots, s_j, s_{j+1}, \dots\}$,并特别设定囚徒困境实际所处的现实世界为 w_0 。

问题的关键在于认识到,局中人的(思考)行动会导致量化域的转变,这种转变本质上是主体思维

^①Duc H N. “Reasoning about Rational, but not logically Omniscient Agents”, *Journal of Computation*, 1997, 7(5): 633–648.

^②张建军,张顺:《条件句的语义排歧与假设性思考的量化机制——五论从形式蕴涵看“实质蕴涵怪论”》,《湖南科技大学学报(社会科学版)》2021年第6期。

过程的外在表征,而非决定因素。基于此,我们可以将 $\diamond K_z(H_m)$ 中的行动序列进行外延化处理,使其成为可量化的对象。当采用情境语义学的表述方式时,可以表述为:在主体c的量化域S中存在一个可量化的境况 s_0 ,在这个境况中 $\diamond K(H_m)$ 为真,即 $(s_0 \in S_c, (s_0, (\diamond K(H_m),)1))$ 。如果借助可能世界语义学,那么也可以写成“ $\diamond \diamond K(H_m)$ ”为真。在传统决策理论与博弈论分析框架中,“ $L_3: \diamond \diamond \alpha \rightarrow \diamond \alpha$ ”亦是隐含于其中的“交互信念”。即使将行动阶段从两轮拓展到n轮,也总有 $\diamond K_{n-2}(H_{n-1})$ 成立。也就是说,主体在最后一轮行动序列的两轮之前,可能知道对方在下一轮行动序列中会确定选择。以此类推,可以得到主体在每轮行动顺序时,都可能知道对方在下一轮行动序列中会确定选择,将每次思考行动都外延化处理即“ $\diamond \cdots \diamond K(H_m)$ ”。则根据 (L_3) 仍会得到“ $\diamond K(H_m)$ ”为真,依然可以在此类分析中得到矛盾。由此可见,具有以上全部“交互信念”的主体还会陷入“二难选择”的悖论,对理论本身来说,则仍未脱出导致悖论这一理论事实。

而 N_1-N_4, L_1, L'_2, L_3 则分别是除经典逻辑系统外需要承认的公理和规则。在推导步骤(23)之前的论证过程中,我们可以将其理解为:在囚徒困境的现实世界 w_0 中,当情境量化域 S_{c1} 给定时对特定命题进行的逻辑推导过程。需要特别注意的是,在这一过程中矛盾信念的产生机制的关键不在于具体情境 s_i 的变化,而在于命题函数所对应的整体情境量化域 S_c 的转变。如果认为仅仅是具体情境 s_i 发生了变化,那么仍旧仅能表述同一框架下的命题关系。但通过上面的分析,我们可以清楚知悉由于主体c的思考行动,使得情境量化域 S_c 发生了转变。通过把握形式蕴涵起到表达共变元命题函数普遍联系的作用,则(24)可表述为“对于任意 $s \in S_{c1}, \langle s, H_m \rightarrow K_z(H_m) \rangle$ 为真”,再结合(23)就可以得到“对于任意 $s \in S_{c1}, \langle s, K_z(H_m) \rangle$ 为真”。而(25)本意表述为“对于任意 $s' \in S_{c2}, \langle s', K_z(H_m) \rangle$ 成立”,其中 S_{c1} 的要求是一致的情境量化域集合, S_{c2} 则是 $K_z(H_m)$ 为真的情境量化域,两者的不同正是出于主体c采取了思考行动使需要达成的目的变化而导致的。显然,使 $K_z(H_m)$ 为真的情境量化域不一定是一致的情境量化域集合,如前文所述 w_1 所对应的情境量化域会得到矛盾。

只有在两个命题函数对应的两个情境量化域 S_{c1} 和 S_{c2} 相同时,才能完成原来的推导。在传统决策理论与博弈论分析框架中隐含的 L_3 规则却忽视了这一变化,致使 S_{c1} 和 S_{c2} 被当作是相同的量化域而进行了处理。通过上述分析可以呈现情境量化域转变视角在解决此类问题中的方法论价值。总结来看,囚徒困境问题本身产生矛盾等价式的原因,正是把在情境量化域集合 S_{c2} 中得到的结论,当成了在无矛盾的现实世界 w_0 以及无矛盾情境量化域集合 S_{c1} 中可以直接适用的结论。显而易见,相较于“ $\diamond K(H_m)$ ”,“ $\diamond \diamond K(H_m)$ ”中的两个“ \diamond ”所对应的量化域的集合是不同的,我们应对其加上下标以进行区分。 \diamond_s 所对应的情境量化域是用外延化的主体行动序列来表示其转变, \diamond_w 则对应原来命题成立的可能世界对应的情境量化域,因此“ $\diamond \diamond K(H_m)$ ”实际应写为“ $\diamond_s \diamond_w K(H_m)$ ”。至此囚徒困境以 N_1-N_4 以及相应规则作为“交互信念”所形成的悖论,得以通过“情境量化域转变”的方案获得化解。

正如李莉曾经指出:“‘由于囚徒之间的理性程度是相近的,因此行为选择也基本相同’为真”,即当囚徒困境中的嫌疑人能够以较高概率确知对方的效用函数时,该困境才构成真正的“二难选择”^①。张智皓的研究也与这一观点相通,他通过关于“玩家对于‘我没有拿我的一千元,当且仅当,我没拿我的一千元你没有拿你的一千元’有高度的信心”的分析,揭示了类似的认知前提。沃克虽然讨论了选择行动所产生的结果,却忽略了思考行动也会对选择结果产生影响。然而,如前文所论证的,即便不考虑这种强认知假设——局中人对彼此效用函数的确定性认知,单就囚徒困境本身而言,只要有上述“交互信念”,我们就已经可以建构出一个严格的逻辑悖论。而解悖的关键就在于诉诸“情境量化域”的转变。

三、作为逻辑悖论的纽科姆难题与囚徒困境之同构

本文对刘易斯观点的新辩护的关键环节在于,借助与囚徒困境相同的 N_1-N_4 以及相应规则作为

^①李莉:《信念与行动的矛盾:单次囚徒困境博弈的重构与消解》,《湖北大学学报(哲学社会科学版)》2020年第5期。

^②张智皓:《纽康难题是另一种囚犯困境吗?》,《逻辑学研究》2017年第1期。

“交互信念”,可以用相同的方式将纽科姆难题建构为一个符合“三要素”定义的逻辑悖论,继而通过与囚徒困境相同的“情境量化域转变”方案对其进行化解,并对前述相关论争的核心问题给出回应。

(一) 原本的纽科姆难题何以成为悖论

将 N_1-N_4 以及相应规则作为“交互信念”,以行动阶段保持两轮的纽科姆难题为例,将囚徒困境中的嫌疑人替换为纽科姆难题中的局中人(及其完美复制人,即局中人 B)就可以得到纽科姆难题成悖所依据的“交互信念”:

设 K_x 表示“局中人 A 在(自己的思想)第 x 行动阶段知道……”; H_y 表示“A 知道 B 在(A 的思想的)第 y 行动阶段确定选择”; $x-1, y-1$ 表示“在 x, y 前一行动阶段”。其中 x 的定义域(量化域)是{行动未开始的阶段(z), 第一行动阶段(m), 第二行动阶段(t)}, y 的定义域(量化域)是{第一行动阶段(m), 第二行动阶段(t)}。

行动阶段定义:局中人 A 行动需要建立在其知道 B 采取什么样的行动,A 产生一个(思想)行动选择结果则定义为一阶段。如,假设 A 知道 B 选 Box2 是第一行动阶段,假设 A 知道 B 的选择后最终决定同时拿 Box1 和 Box2 则是第二阶段。

则使用与囚徒困境悖论相同的前提与规则,经过与上一章相同的推理,同样可以得到局中人 A 在 w_1 中知道预测者在第一行动阶段确定选择并且局中人 A 知道预测者在第二行动阶段确定选择的矛盾。

假设第一行动阶段 Box1 和 Box2 都拿,如 A 知道超级生物在第一行动阶段确定选择,那么才会导致收益变为 1 000 美元,那么进入第二行动阶段最终选 Box2 才是合理的。如 A 知道超级生物在第二行动阶段确定选择,那么超级生物知道 A 会选 Box2,则同时拿 Box1 和 Box2 才是合理的。假设第一行动阶段只拿 Box2,如 A 知道超级生物在第一行动阶段确定选择,那么同时拿 Box1 和 Box2 才是合理的。如 A 知道超级生物在第二行动阶段确定选择那么超级生物知道 A 会同时拿 Box1 和 Box2,则只选 Box2 才是合理的。综上可以得到,无论先思考采取哪种方案,最终由于推理得到的矛盾结论,都会使得选 Box2 和同时拿 Box1 和 Box2 两个矛盾行为都是合理的。纽科姆难题从与囚徒困境相同的“交互信念”出发,经过“严密无误的逻辑推导”也最终得到了“矛盾等价式”,这也是纽科姆难题成为一个逻辑悖论的原因。

(二) 原本的纽科姆难题的化解

在引入“情境量化域转变”的思想后,经过类似的推导,纽科姆难题会最终得到 $\diamond_s \diamond_w K(H_m)$ 为真,在纽科姆难题中两种理性原则产生矛盾的原因就在于此。在两轮的纽科姆难题中,合理选择行动的 DP 原则和 MEU 原则(或是 CDT 和 EDT 原则)都要求有相同的前提: $\diamond_w K(H_m)$ 为真,只有如此才会产生两种方案。但通过上述分析我们只能得到 $\diamond_s \diamond_w K(H_m)$ 为真。而此时对应的“情境量化域”并不一定同时满足 DP 原则和 MEU 原则(或是 CDT 和 EDT 原则)适用的要求。通常来说, $\diamond \diamond K(H_m)$ 只有叠置可能算子消去才能得到 $\diamond K(H_m)$,这要求起码满足规则 L_3 。如此一来,其对定义在 $W \times W$ 上的可及关系 R 就有了专门的要求,这只能在传统决策理论与博弈论理论框架中静态地分析问题时才生效,不能刻画主体实际面对的由思考行动本身所带来的动态过程。而在传统决策理论与博弈论中,将情境量化域集合 S_e 与可能世界集合 W 进行混淆,误解了这两个发挥不同作用的层面,由规则 L_3 最终得到 $\diamond K(H_m)$,仍会导致纽科姆难题成为一个悖论。

显然,纽科姆难题与囚徒困境问题本身,实际上都可以依据上述相同的“交互信念”得以通过逻辑推导得出矛盾等价式,因此通过逻辑悖论的三个构成要素来看,两者确实是同一个问题的不同表述。可以对此种悖论做出如下界定:从 N_1-N_4 以及相应规则出发,借助可能世界四元组 $\langle W, R, V, S_e \rangle$,承认规则 L_3 ,即 $\diamond \diamond \alpha \rightarrow \diamond \alpha$ 对可及关系 R 产生的限制,如果未正确区分表征行动的情境量化域集合 S_e 与可能世界集合 W ,即可通过有效逻辑推导得到矛盾等价式。又因为这种悖论本质地涉及行动合理性辩护,所以两者确实是同一个合理行动悖论。

从现实认知的角度来看,上述要求的可及关系 R 显然不能一直得到满足,如果 R 满足正规模态系统 S_5 所要求的“自反”“对称”“传递”等关系,那么如埃尔斯特(J. Elster)所言,对现实的有限理性人而

言,这样的系统是“平庸的”^①。刻画现实中人的思考不能使用这样的平庸系统。在静态分析系统中,由于未能充分考虑行动的动态影响,研究者常常会将情境量化域集合 S_c 与可能世界集合 W 相混淆。而在传统决策理论与博弈论中,又恰好同时满足了这种混淆和规则 L_3 ,因此囚徒困境在决策理论与博弈论作为前提的分析中,会得到矛盾的结果是题中应有之义。

在增添“情境量化域转变”的思想作为“交互信念”后,纽科姆难题中反应的命题函数对应的量化域分别为 S_{c1} 和 S_{c2} ,两个合理原则得到在一致量化域中适用的辩护,进而可以对孔斯的纽科姆难题解悖方案^②进行修正,就会得到如下两个不会得到矛盾的公式,悖论可以得到化解:

$$\sim \langle s, \langle J, p, w_0, 1 \rangle \rangle \rightarrow p, \text{ 其中 } s \in \{S_{c1} \mid s_0, s_1, \dots, s_n\}^{\circledR},$$

$$\langle s', \langle J, (p \rightarrow (\sim \langle s, \langle J, p, w_0, 1 \rangle \rangle)), 1 \rangle \rangle, \text{ 其中 } s \in \{S_{c1} \mid s_0, \dots, s_n\}, s' \in \{S_{c2} \mid s_a, \dots, s_m\}, \text{ 且 } s \neq s'.$$

相较于孔斯原本将情境、主体分开的“奥斯汀式”命题,如此构建更有助于表明主体行动与量化域转变。在囚徒困境的传统解释语境中,传统决策理论与博弈论也都面临着对“交互信念”中规则 L_3 :“ $\Diamond \Diamond \alpha \rightarrow \Diamond \alpha$ ”预设的忽视。而即便承认 L_3 这一交互信念,对表征行动的情境量化域集合 S_c 与可能世界集合 W 的混淆,仍会导致矛盾的出现。因此,刘易斯将使用“完美复制人”策略后的纽科姆难题和囚徒困境视为同一问题是正确的,只是“完美复制人”策略看起来过强,引起了诸多误解。本文明确地将 N_1-N_4 以及相应规则作为“交互信念”之后,也可以得到原本纽科姆难题与囚徒困境是同一种逻辑悖论的答案。

(三)“情境量化域转变”方案对论争的回应

针对既往研究中对“完美复制人”策略与“行动-结果”关系的思考,“情境量化域转变”方案也可以给予良好回应。刘易斯在考虑“完美复制人”策略时,过多考虑了预测的“还原性”,直接选取了局中人 A 的完美复制人作为讨论对象,但是仅从正确“预测”其选择结果的角度来看,我们只需要挖掘其推理所依据的“交互信念”和过程即可。而且对于论证原本纽科姆难题和囚徒困境是同一个逻辑悖论来说,“完美复制人”策略过强,并不是必需的。因此,不仅仅是引入“完美复制人”策略的纽科姆难题与囚徒困境是同一个问题,只要引入承认依据 N_1-N_4 以及相应规则作为“交互信念”的认知共同体而导致矛盾等价式的情形,都可以得出其与原本的纽科姆难题是同一种逻辑悖论的结论。

巴蒂奇、张智皓对索贝尔、贝穆德斯的批评也基于其对“完美复制人”策略的“错误解读”。从“情境量化域转变”方案出发,只选取具有相应“交互信念”的认知共同体,则可以避免贝穆德斯在认知情形上的直接责难。但不可否认的是,认知情形角度确实展现出了问题之所在:主体在采取不同认知情形时的思考行动,可能会导致其所面临的情境量化域产生转变。在对纽科姆难题的讨论中贝穆德斯特别强调了这点,并以此为基准,批评将纽科姆难题与囚徒困境视为同一个问题的看法。但以“情境量化域转变”方案作囚徒困境的分析时,可以清晰地看到,囚徒困境本身也存在主体在采取不同认知情形时的思考行动,也可能会导致其所面临的情境量化域产生转变的问题。但索贝尔、贝穆德斯却都忽略了这一点,认为囚徒困境并没有遭受此类问题,进而得出两者非同一种问题的结论。“情境量化域转变”方案展示出,之所以囚徒困境不被认为如此,是因为传统决策理论和博弈论分析框架并没有对此给出良好的解释。张智皓所给出的采取特定视角后使用“完美复制人”策略的纽科姆难题与囚徒困境的同步结构变化,正是这类问题的一种特殊情况。

沃克对两种条件句的区分颇具创见。但如同贝穆德斯一样,其没有考虑到将囚徒困境视为使用“附加条件句”并不是囚徒困境问题本身所要求的,而是作为讨论背景的传统决策理论和博弈论所隐含的“交互信念”的限制。这也能表明为什么纳什均衡并不只在纽科姆难题中不被使用,甚至在囚徒困境

^①乔恩·埃尔斯特:《逻辑与社会:矛盾与可能世界》,贾国恒、张建军译,南京大学出版社2015年版,第23页。

^②罗伯特·C·孔斯:《信念悖论与策略合理性》,张建军等译,中国人民大学出版社2020版,第133页。

^③ $\langle s_j, \langle J, p, w_0, 1 \rangle \rangle$,其中 $s_j \in \{S_{c1} \mid s_0, s_1, \dots, s_n\}$,表示在现实世界中主体 c 的量化域 S_1 中的一个情境 s_j 中,命题 p 是可证例的。

中也不应该被直接使用,而只是在以传统决策理论和博弈论去分析使用“完美复制人”策略后的纽科姆难题与囚徒困境中才会被使用。因此,传统决策理论与博弈论分析框架必须予以合理修正,而对于“行动-结果”关系的讨论,也不能仅仅局限于传统决策理论与博弈论分析框架中。在探讨以纽科姆难题与囚徒困境为代表的“合理行动悖论”时,需要重点考察为以往研究所忽视的思考行动所带来的量化域转变问题,重新对过去分析“行动-结果”关系所依据的“交互信念”本身是否可靠进行考察。

结语

本文分析表明,纽科姆难题与囚徒困境在以传统决策理论和博弈论框架的分析中,都具备相同的“交互信念”作为前提,经过“严密无误的逻辑推导”可以得到“矛盾等价式”,因而是同一种合理行动悖论,故刘易斯说它们本质上是同样的问题是正确的。但通过本文考察,“完美复制人”这一概念过强,诉诸 N_1-N_4 以及相应规则作为认知共同体的“交互信念”,将原本的纽科姆难题和囚徒困境都构建为“三要素”齐全的悖论,可为刘易斯的结论做出新的辩护。而使用“完美复制人”策略后的纽科姆难题和囚徒困境之所以不被某些学者视为同一个问题,主要原因在于对其进行分析时所采用的传统决策理论和博弈论分析框架隐含的预设。在决策理论和博弈论蓬勃发展、广泛应用的今天,正确认识这种隐含的预设,有助于我们更好地看待、理解借助这类分析系统得出的论证。在进一步理解行动带来的量化域转变后,亦有助于对其进行系统性反思和修正。

总之,知识本身多形成于不同的情境量化域之下,而形成于不同情境量化域的知识放在一起,却可能会存在隐含的矛盾。从悖论的视角考察纽科姆难题和囚徒困境,有助于我们发现问题的核心,重新评估对“交互信念”的认识,反思我们思考过程中的深层因素,消除其中的隐含矛盾。将促使量化域转变的行动引入纽科姆难题与囚徒困境的研究,既能为讨论合理行动悖论群落的解悖方案提供新的思路,也可以为进一步挖掘合理行动悖论研究在当代行动哲学与行动逻辑中的作用提供助力。

A New Exploration of the Relationship Between the Newcomb Problem and the Prisoners' Dilemma

DING Yuchen

(School of Philosophy, Nanjing University, Nanjing 210023, China)

Abstract: David Lewis once demonstrated that the Newcomb Problem and the Prisoners' Dilemma are different manifestations of the same problem by introducing the “Perfect Replicant” strategy. However, there is considerable controversy in the academic community over this. Reexamining the Newcomb Problem and the Prisoners' Dilemma with the “three-element” theory composed of logical paradoxes can further clarify why they form logical paradoxes based on the same elements of “mutual belief” and have isomorphic logical structures. This reveals the theoretical presuppositions implicit in traditional decision theory and game theory when dealing with problems such as the Prisoners' Dilemma, which is conducive to a correct understanding of the functions and limitations of the two theories. And through the idea of transformation of quantitative domain, it provides a unified approach to resolving such paradoxes.

Key words: newcomb problem; prisoners' dilemma; logical paradox; transformation of quantitative situational domain

(责任编辑 曾小明)