

“十二五”国家重点图书
经济学术译丛·当代制度分析前沿系列

博弈论与社会契约 (第2卷·下) 公正博弈

GAME THEORY AND THE SOCIAL CONTRACT (VOLUME 2)

— JUST PLAYING

肯·宾默尔
(Ken Binmore)

著

潘春阳
陈雅静
陈琳

译

 上海财经大学出版社

“十二五”国家重点图书

经济学
译丛
当代制度分析前沿系列

博弈论与社会契约

(第2卷·下)

公正博弈

肯·宾默尔
(Ken Binmore) 著

潘春阳
陈雅静 译

陈琳

 上海财经大学出版社

本书由上海文化发展基金会图书出版专项基金资助出版

图书在版编目(CIP)数据

博弈论与社会契约(第2卷)·公正博弈/(英)宾默尔(Binmore, K.)著;潘春阳等译. —上海:上海财经大学出版社,2016.1

(经济学术译丛·当代制度分析前沿系列)

书名原文:Game Theory and the Social Contract II:Just Playing
ISBN 978-7-5642-1875-1/F·1875

I. ①博… II. ①宾… ②潘… III. ①博弈论—研究 ②社会契约—研究 IV. ①F224.32 ②F246

中国版本图书馆CIP数据核字(2014)第073362号

□责任编辑 温涌

□封面设计 周卫民

BOYILUN YU SHEHUI QIYUE

博弈论与社会契约

(第2卷)

公正博弈

肯·宾默尔 著
(Ken Binmore)

潘春阳 陈雅静 陈琳 译

上海财经大学出版社出版发行
(上海市武东路321号乙 邮编200434)

网 址: <http://www.sufep.com>

电子邮箱: webmaster@sufep.com

全国新华书店经销

上海叶大印务发展有限公司印刷装订

2016年1月第1版 2016年1月第1次印刷

890mm×1240mm 1/32 22.375印张 621千字

印数:0 001—3 000 定价:80.00元(上下册)

目 录

“当代制度分析前沿系列”总序/1

译者序/1

致歉/1

系列前言/1

阅读指南/1

导论 设定场景/3

0.1 走向何方? /3

0.2 妥协的艺术/4

0.2.1 高跷上的胡言乱语/5

0.2.2 社会契约/7

0.2.3 改革/9

0.2.4 原初状态/12

0.2.5 讨价还价/17

0.3 道德哲学/20

0.3.1 传统的哲学分类/21

0.3.2 世纪末/23

0.4 非合作博弈理论/26

0.4.1 最后通牒博弈/27

0.4.2 反常现象? /36

0.5 合作博弈理论/46

0.5.1 联盟形式下的博弈/46

0.6 纳什规划/51

0.7 实施/59

第1章 磋商探微/71

1.1 现实的讨价还价模型/71

1.2 讨价还价难题/72

1.2.1 收益区域/73

1.2.2 纳什讨价还价难题/77

1.2.3 讨价还价集/80

1.2.4 分钱博弈/82

1.2.5 埃奇沃思箱形图/88

1.3 讨价还价解/92

1.3.1 纳什讨价还价解/93

1.3.2 卡莱-斯莫洛定斯基解/97

1.3.3 存在人际比较的讨价还价/99

1.4 讨价还价解的特征化/103

1.4.1 卡莱-斯莫洛定斯基公理/105

1.4.2 纳什公理/110

1.4.3 再磋商公理/113

1.5 有承诺的讨价还价/118

1.5.1 纳什要价博弈/118

1.5.2 固定的和可变的威胁/122

1.6 无信任交易/126

1.6.1 重复博弈/127

1.6.2 过渡协议/131

1.7 无承诺的讨价还价/139

1.7.1 交替出价博弈/141

1.7.2 鲁宾斯坦模型有多现实呢? /150

2

1.8 研究讨价还价的其他方法/156

- 1.8.1 科斯定理/157
- 1.8.2 对社会契约讨价还价的高德观点/160

第2章 伊甸园里的演化/167

- 2.1 善、正当与适当/167
- 2.2 功利主义/169
 - 2.2.1 至善/172
 - 2.2.2 道德伪善主义者/173
 - 2.2.3 理想观察者/175
 - 2.2.4 哲学王/179
 - 2.2.5 社会契约方法/183
 - 2.2.6 是规则功利主义还是行为功利主义? /185
 - 2.2.7 整体情况/188
- 2.3 虚幻假定? /192
 - 2.3.1 效用的人际比较/193
- 2.4 演化伦理学/204
- 2.5 演化与正义/210
 - 2.5.1 互惠/213
 - 2.5.2 亲缘关系/214
 - 2.5.3 均衡选择/233
 - 2.5.4 移情与公平/241
 - 2.5.5 长期、短期和中期/257
- 2.6 非目的论功利主义/259
 - 2.6.1 伊甸园里的承诺/260
 - 2.6.2 中期的人际比较/274
 - 2.6.3 重述罗尔斯的故事/280
- 2.7 道德作为一种短期现象/283
 - 2.7.1 公主和豌豆/284
 - 2.7.2 正义是如何运作的/286

2.8 为什么不是功利主义? /291

第3章 理性互惠/299

3.1 相互帮助/299

3.2 在适当理论中的权利/311

3.2.1 权利如同策略吗? /313

3.2.2 维持均衡的规则/316

3.2.3 道德责任/318

3.2.4 自由意志/321

3.2.5 永不绝望! /332

3.3 无名氏定理/334

3.3.1 文化基因/335

3.3.2 有限自动机/336

3.3.3 计算收益/339

3.3.4 互惠分成/340

3.3.5 罪与罚/347

3.3.6 相互监管的监管者/353

3.3.7 针锋相对? /356

3.3.8 合作是如何演化出来的? /363

3.4 大型社会中的社会契约/374

3.4.1 社会转移/375

3.4.2 友谊和联盟/378

3.4.3 警察力量/379

3.4.4 惩罚无辜者/380

3.4.5 领导与权威/383

3.5 情感的作用/387

3.5.1 显见之事/387

3.5.2 管窥之见/388

3.6 程序正义/393

- 3.6.1 从无政府到国家/395
- 3.6.2 自然均衡/397
- 3.7 再磋商/399
 - 3.7.1 从这里到那里/401
 - 3.7.2 在原初状态中的再磋商/404
 - 3.7.3 使惩罚与罪行相符/406
 - 3.7.4 再磋商防御/407
- 3.8 道德观念是什么? /409
 - 3.8.1 混淆喜好和观念/409
 - 3.8.2 论亚当·斯密难题/420
 - 3.8.3 后福利主义/423
 - 3.8.4 道德相对主义不是什么? /425

第4章 向往乌托邦/431

- 4.1 导论/431
- 4.2 嫉妒/433
- 4.3 经济学中的公平/438
 - 4.3.1 免嫉妒/439
 - 4.3.2 福利主义/449
- 4.4 心理学中的公平/455
- 4.5 人类学中的公平/460
 - 4.5.1 分享和关心/462
 - 4.5.2 觅食社会中的执行/466
 - 4.5.3 史前的无政府状态? /471
 - 4.5.4 小型群体中的亲属关系/474
- 4.6 道德博弈/485
 - 4.6.1 公平的社会契约/485
 - 4.6.2 失乐园/487
 - 4.6.3 对原初状态建模/489

- 4.6.4 正义何时被分配? /494
- 4.6.5 罗尔斯是无辜的! /501
- 4.6.6 中期的人际比较/504
- 4.6.7 共识与环境/513
- 4.6.8 短期中的道德/515
- 4.6.9 平均主义 vs. 功利主义/517
- 4.6.10 复乐园? /520
- 4.7 价值和权力/521
 - 4.7.1 权力意志? /522
 - 4.7.2 比较静态/526
 - 4.7.3 按需分配? /528
 - 4.7.4 劳动使人自由? /532
 - 4.7.5 各尽所能/535
 - 4.7.6 高尚和卑贱/536
 - 4.7.7 社会主义和资本主义/538
- 4.8 市场与长期/540
 - 4.8.1 瓦尔拉斯讨价还价解/546
 - 4.8.2 错误表达私人偏好/553
 - 4.8.3 公平价格的概念/563
 - 4.8.4 时间毁了一切/565
- 4.9 未完成的事业/566
 - 4.9.1 大的社会和同盟/567
 - 4.9.2 不完全的信息和机制设计/568
 - 4.9.3 不断变化的生存博弈/569
- 4.10 完美的理想国度? /571
 - 4.10.1 辉格党原则是什么? /571
 - 4.10.2 哪里有辉格党原则? /575
- 4.11 休谟主义与人道主义/578

附录 A 真的如此! /585

- A.1 自然主义/585
 - A.1.1 因果倒置/587
- A.2 把“人”模型化/587
 - A.2.1 身体的力量/590
 - A.2.2 理性/591
 - A.2.3 激情/593
 - A.2.4 经验/594

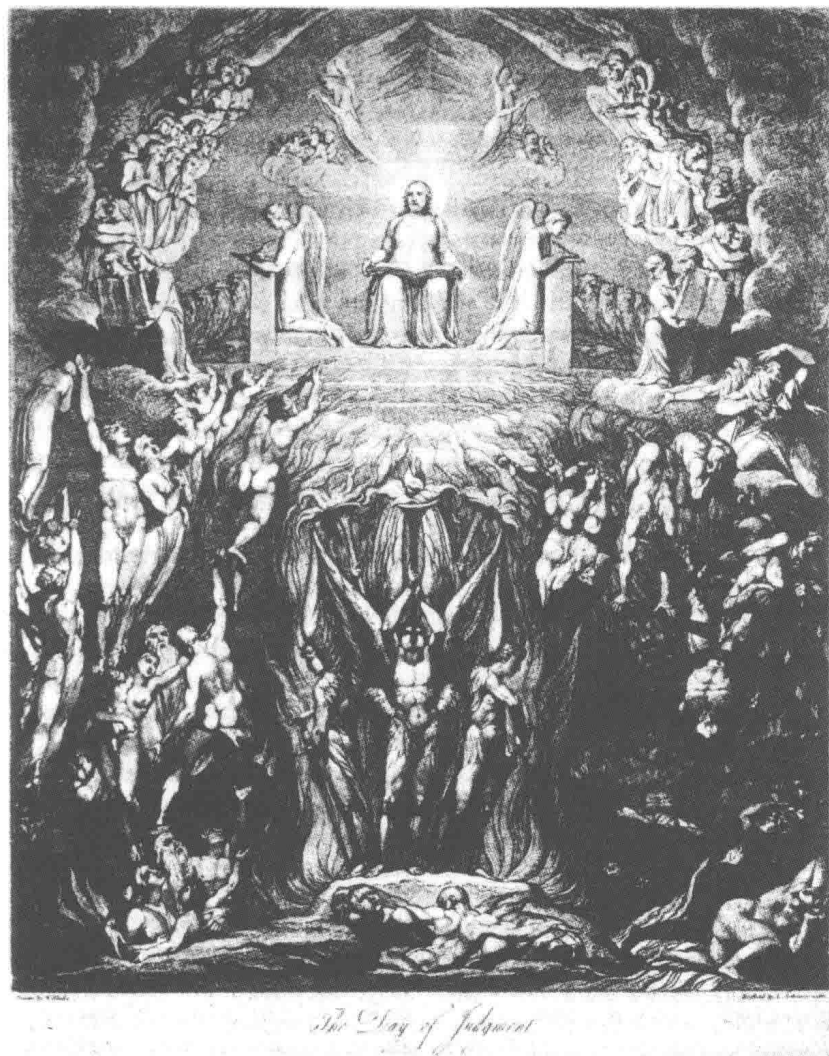
附录 B 哈萨尼学派/599

- B.1 导言/599
- B.2 目的论的功利主义/600
- B.3 非目的论的功利主义/606

附录 C 讨价还价理论/611

- C.1 简介/611
- C.2 轮流出价博弈/611
- C.3 偏好/613
- C.4 静态子博弈完美均衡/615
- C.5 不稳定均衡/618
- C.6 广义纳什讨价还价解/622
- C.7 纳什规划/623

参考文献/627



那些决定我们社会地位的人的语焉不详的反对,也经常足以使我们维系在由一个次优社会契约所决定的均衡路径之上。如果我们可以劝导人们更多地关注生存博弈在人死后将继续进行这一思想,而不仅是获得人们的口头支持,那么威廉·布莱克所面临的惩罚就能使最优的社会契约在当下实现。

第 3 章

理性互惠

让我总是做这个,我就会总是让你做这个,你就会让别人总是做这个。你在一个商业国家中承担了人的全部责任。

——潘克斯先生(Mr. Pancks),引自狄更斯(Dickens),《小杜丽》(*Little Dorrit*)

3.1 相互帮助



前一章以亚当和夏娃被驱逐出伊甸园而告终。在本书的剩余部分,再也不存在哲学王或者其他外部执行机构以维系某个他们一起制定的契约。唯一存在的“执法者”将是亚当和夏娃自己。

如果他们的生存博弈是如图 3.1(a)所示的单次囚徒困境,那么结局将变得十分悲惨。在存在一个哲学王以实施协议的情况下,在如图 3.1(b)阴影部分所示的合作收益区域内的任何一个收益对,都可以成为一个社会契约(参考第 1.2.1 节)。在不存在哲学王的情况下,亚当和夏娃将被局限于运行他们生存博弈中的均衡状态。但是,单次囚徒困境博弈只有一个均衡点——在这个均衡点中,每个局中人都采用“鹰”这一强占

优策略并获得零回报。

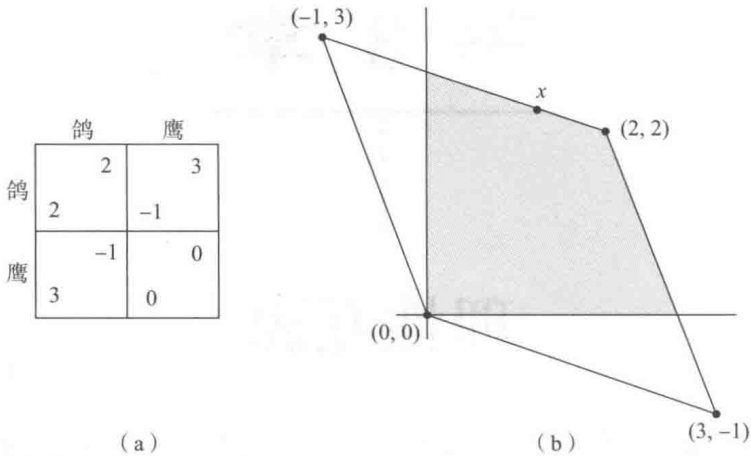


图 3.1 囚徒困境

然而,正如第 1 卷所反复强调的,遵循某些政治哲学家们的做法是错误的,他们致力于将我们的生存博弈构建为囚徒困境——或者是斗鸡博弈、性别战,抑或是其他单次博弈。如果我们的生存博弈的确是单次囚徒困境,那么我们将永远不可能进化为社会动物。尽管仍然包含了大量近于夸张的简化,但将生存博弈构建为一个无限期重复博弈显然更有意义。正如我在第 1 卷第 2.2.6 节所解释的,这样的建模方法解决了一个难题,那就是如何找到“正确”的博弈理论框架以刻画人类困境。真正重要的是,存在某个博弈是重复的,而我们选择重复哪个博弈在很大程度上是次要的。

在一个无限期重复博弈中,如果局中人仍然选择单次博弈中的一个纳什均衡,那么这仍然是一个均衡状态。在无限期囚徒困境博弈中,如果亚当和夏娃总是计划采用鹰策略,那么这也是一个均衡状态。这一均衡状态通常与霍布斯[251]提出的残酷的自然状态是一致的,是唯一一个能够替代独裁社会契约的选择。当然,还存在很多其他的纳什均衡点。在图 3.1(b)阴影区域中的每一个收益对都代表了无限期重复囚徒困境博弈的一个纳什均衡结局(第 1 卷第 2.2.6 节)。特别地,亚当和夏娃可以

选择任何帕累托有效的收益对作为他们的社会契约,且这一收益对相比霍布斯自然状态是一种帕累托改进。

正如我在第1卷第2.2.6节所解释的,在一个无限期重复博弈中,支撑这种合作的纳什均衡的机制正是互惠性(reciprocity)——艾克斯罗德(Axelrod)[31]提出的著名的“针锋相对”(TIT-FOR-TAT)策略便是一个良好的例证。

合作的演化 这一策略要求一个局中人一开始便采用鸽策略,之后总是选择对手上一期采用的策略(第1卷第2.2.6节)。两个“针锋相对”的策略便构成了无限期囚徒困境博弈的一个纳什均衡。由于无论是背叛还是合作都是相互的,因而对“针锋相对”策略的最优反应,便是在所有时期都选择合作(假设贴现率并不是特别高)。因为采用“针锋相对”策略的局中人永远不会成为第一个背叛者,所以两个都选择了“针锋相对”策略的局中人总是能够保持合作状态。任何一个局中人都会对对手的策略选择做出最优反应。

人们认识到互惠性是人类社会性的主要动力,至少可以追溯到公元前7世纪,^[1]但是,关于其内在机制可以辨识的现代表述,则不得不等待大卫·休谟的出现。正如休谟[267, p. 521]所解释的:

我学会了对别人进行服务,虽然我并没有感怀他的任何恩德,因为我可以预料,他会报答我的服务,并期望得到同样的另一次服务,并且也为了与我或与其他人维持同样的互助往来关系。而且相应地,在我为他服务之后,他会从我的行为中获得利益,他会被诱导以履行他的义务,因为他预见到他的拒绝可能产生的后果。

我怀疑,那些在50年代重新创立这一理念的博弈理论家们,可能完

[1] 施与那些施与别人的人,不施与那些不施与别人的人。——赫西奥德(Hesiod)
[246]

全没有意识到休谟的思想。因为对他们而言,显而易见的是,单次博弈的应用范围十分有限,他们的工作仅仅是将纳什[381]当时新创造的均衡概念拓展至重复博弈。通常的情况便是,当某种思想出现的时机到来时,许多学者几乎同时想到了重复博弈理论的基本结论的不同版本[奥曼和马斯库勒(Aumann and Maschler)[28]]。之所以现在我们把这一结论称为无名氏定理(folk theorem),是因为没有人知道这一定理地发现应该归功于谁比较合适。该定理认为,单次博弈合作收益区域中的所有令人感兴趣的结局,也可以成为该博弈的无限期重复版本的均衡状态。因此,在重复博弈的情势下,我们并不需要一个哲学王来执行契约;也就是说,我们自己都做不了的事,外部执法机构也必将无能为力。

我不知道当时的博弈理论家们对这一发现的广泛意义的理解达到了什么程度,但是他们显然没有成功地将这一福音传播至生物群落,结果是,特里弗斯(Trivers)[524]在大约15年之后重新发现了这一基本思想,他成为一个更为成功的福音传道者。因此,无名氏定理背后的作用机制被认为是互惠利他主义(reciprocal altruism)的,尽管在这一机制运作时,并没有真正的利他主义包含其中。

在梅纳德·史密斯(Maynard Smith)[347]广受赞赏的《演化与博弈论》(*Evolution and the Theory of Games*),以及艾克斯罗德(Axelrod)[31]更为成功的《合作的演化》这两部名著出版十多年以后,严肃的博弈理论几乎仍然是演化生物学家和政治科学家们的未知领域。例如,杜格纳托夫(Dugnatov)[163, p. 11]的《动物间的合作》(*Cooperation among Animals*)告诉我们,冯·诺依曼和摩根斯坦对博弈论的贡献在于提出了囚徒困境!因此,我决定在本章详尽地讨论无名氏定理及其内涵,当然,首先我们有必要对这一定理在社会契约理论中的地位进行拓展回顾。

自利 自由主义者想当然地认为,那些使各种大家族团结起来的相互尊重和兄弟情义的共同纽带,也必然提供了防止大型社会分崩离析的黏合剂。事实上,对许多学者来说,自利动机能够提供足够的向心力以维

持社会运转的思想令人十分厌恶。我同意利他主义者相比利己主义者更能够成为合得来的伙伴,但我认为,人们需要成为利他主义者,才能使社会运转这一信念与人类历史的事实并不一致。陀思妥耶夫斯基(Dostoyevsky) [162, p. 35] 在其自传《死屋手记》(*House of the Dead*) [1]中描述他在俄国沙皇集中营的囚徒经历时,非常清晰地阐明了这一点:

那里的绝大多数人都堕落得无可救药。诽谤和流言从未停息;这就是一个地狱、一个灵魂的黑夜。但是,没有人敢于反抗这一内在规律,而都选择接受监狱的规则;每一个人都屈服于它们。当做得过头的人来到监狱时,恐怖的气氛便笼罩了整个村镇。然而,新来的罪犯潜移默化地变得十分顺从,并融入了整体的环境。

总之,我们并不需要遵循埃尔斯特(Elster)[172]强调的爱与责任是连接社会的黏合剂的观点。更为现实的做法是,与詹姆斯·麦迪逊(James Madison)[215, p. 285]一起在大卫·休谟中寻找一个恰当的隐喻。〔2〕稳定的社会无须黏合剂以使其连接在一起,因为它们建立在与干石墙相同的基础之上,这些干石墙正是牧羊人用以为分清他们的财产而建造的。正如一只狒狒给另一只狒狒梳毛刷洗以期获得一次类似的服务作为回报那样,在干石墙中的每一块石头都被其邻居的重量固定在位,同时也贡献自己的力量将其邻居固定在位。

吉基尔博士(Dr. Jekylls)的社会相比海兹先生(Mr. Hydes)的社会无

〔1〕 道德哲学家们经常引用文学著作作为证据,以支持其所辩护的某种心理学理论的现象早已司空见惯。特别是荷马(Homer)和亨利·詹姆斯(Henry James)的作品似乎还广泛地成为灵感的来源。但确定无疑的是,人们从虚构作品中学习到的都是作者自身对人性的理解。然而,陀思妥耶夫斯基的《死屋手记》并不是一部虚构的作品,而是当时他亲眼所见并为之震惊的事实记录。

〔2〕 人类的幸福与繁荣……或许可以与拱顶建筑相类比,其中每一块石块,如果单独出来,都会掉落地面。但是通过各部分的相互支撑和组合,整个建筑结构便得以支撑。——大卫·休谟[266]