

HW1: 博弈论与多臂老虎机算法基础

姓名： 学号： 日期：



1.1. 占优策略均衡与纳什均衡的关系

证明如下关于占优策略均衡与纳什均衡的关系的结论：

1. 如果每个参与人 i 都有一个占优于其它所有策略的策略 s_i^* ，那么 $s^* = (s_1^*, \dots, s_n^*)$ 是纳什均衡；
2. 如果每个参与人 i 都有一个严格占优于其它所有策略的策略 s_i^* ，那么 $s^* = (s_1^*, \dots, s_n^*)$ 是博弈的唯一纳什均衡。

Solution:

1. 我们要证明 $s^* = (s_1^*, s_2^*, \dots, s_n^*)$ 是一个纳什均衡。根据纳什均衡的定义，我们需要证明对于任何参与人 i ，他都无法通过单方面改变策略来提高收益。也就是要证明：

$$u_i(s_i^*, s_{-i}^*) \geq u_i(s_i, s_{-i}^*) \quad \forall s_i \in S_i$$

根据题设，对于每个参与人 i ，策略 s_i^* 是他所有策略中的一个（弱）占优策略。

根据占优策略的定义，对于任何参与人 i ，他的策略 s_i^* 满足：无论其他参与人选择什么策略组合 s_{-i} ，选择 s_i^* 的收益总是不会低于选择任何其他策略 s_i 的收益。即：

$$u_i(s_i^*, s_{-i}) \geq u_i(s_i, s_{-i}) \quad \forall s_i \in S_i, \forall s_{-i} \in S_{-i}$$

既然这个条件对所有的对手策略组合 s_{-i} 都成立，那么它自然也对 s_{-i}^* (即其他参与人选择其各自占优策略的组合) 成立。

因此，我们有：

$$u_i(s_i^*, s_{-i}^*) \geq u_i(s_i, s_{-i}^*) \quad \forall s_i \in S_i$$

这个不等式对所有参与人 $i \in N$ 都成立。这完全符合纳什均衡的定义。

因此，由每个参与人的占优策略构成的策略组合 s^* 是一个纳什均衡。

2. 这个证明包含两部分：a) 证明 s^* 是一个纳什均衡；b) 证明这个纳什均衡是唯一的。

a) 证明 s^* 是一个纳什均衡

根据严格占优策略的定义，对于任何参与人 i ，他的策略 s_i^* 满足：无论其他参与人选择什么策略组合 s_{-i} ，选择 s_i^* 的收益总是严格高于选择任何其他不等于 s_i^* 的策略 s_i 的收益。即：

$$u_i(s_i^*, s_{-i}) > u_i(s_i, s_{-i}) \quad \forall s_i \in S_{i(s_i \neq s_i^*)}, \forall s_{-i} \in S_{-i}$$

同样，这个条件对所有的对手策略组合 s_{-i} 都成立，自然也对 s_{-i}^* 成立。所以：

$$u_i(s_i^*, s_{-i}^*) > u_i(s_i, s_{-i}^*) \quad \forall s_i \in S_{i(s_i \neq s_i^*)}$$

这显然满足纳什均衡的条件 $u_i(s_i^*, s_{-i}^*) \geq u_i(s_i, s_{-i}^*)$ 。因此， s^* 是一个纳什均衡。

b) 证明这个纳什均衡是唯一的

我们使用反证法。

假设除了 $s^* = (s_1^*, \dots, s_n^*)$ 之外, 还存在另一个不同的纳什均衡, 我们称之为 $s' = (s'_1, \dots, s'_n)$ 。由于 $s' \neq s^*$, 那么必然存在至少一个参与人 j , 使得他的策略 s'_j 与他的严格占优策略 s_j^* 不同, 即 $s'_j \neq s_j^*$ 。

现在我们来考察参与人 j 的收益。根据严格占优策略 s_j^* 的定义, 对于任意的对手策略组合 s_{-j} , 选择 s_j^* 的收益都严格大于选择任何其他策略 s'_j 的收益。

$$u_j(s_j^*, s_{-j}) > u_j(s'_j, s_{-j}) \quad \forall s'_j \neq s_j^*, \forall s_{-j} \in S_{-j}$$

这个条件对于对手策略组合 s'_{-j} (即在 s' 均衡中, 除 j 以外的其他人的策略) 也成立。所以, 我们有:

$$u_j(s_j^*, s'_{-j}) > u_j(s'_j, s'_{-j})$$

这个不等式意味着: 在其他参与人都选择 s'_{-j} 的情况下, 参与人 j 如果将自己的策略从 s'_j 单方面变更为 s_j^* , 他的收益将会严格增加。

但这与我们最初的假设 " s' 是一个纳什均衡" 矛盾。因为根据纳什均衡的定义, 任何参与人都不可能通过单方面改变策略来获得更高的收益。

矛盾, 故假设是错误的。

因此, 不存在其他任何纳什均衡。由每个参与人的严格占优策略构成的策略组合 s' 是该博弈的唯一纳什均衡。

证毕!



1.2.N人古诺竞争

假设在古诺竞争中, 一共有 J 家企业。当市场中所有企业总产量为 q 时, 市场价格为 $p(q) = a - bq$ 。且每个企业生产单位产品的成本都是同一个常数 c , 即企业 i 的产量为 q_i 时该企业的成本为 $c_i(q_i) = c \cdot q_i$ 。假设 $a > c \geq 0, b > 0$ 。

1. 求纳什均衡下所有企业的总产量以及市场价格;
2. 讨论均衡价格随着 J 变化的情况, 你有什么启示?
3. 讨论 $J \rightarrow \infty$ 的均衡结果, 你有什么启示?

Solution:

1. 企业 i 的利润 π_i 取决于它自己的产量 q_i 和所有其它企业的总产量 $\sum_{j \neq i} q_j$ 。

$$\pi_i(q_i, q_{-i}) = p \cdot q_i - c \cdot q_i$$

将价格函数 $p(q) = a - b(q_i + \sum_{j \neq i} q_j)$ 代入:

$$\pi_i(q_i, q_{-i}) = \left[a - b \left(q_i + \sum_{j \neq i} q_j \right) \right] q_i - cq_i$$

$$\pi_i(q_i, q_{-i}) = aq_i - bq_i^2 - b \left(\sum_{j \neq i} q_j \right) q_i - cq_i$$

古诺模型的核心假设是, 每个企业在决定自己的产量时, 都将其他企业的产量视为给定值。因此, 企业 i 会选择 q_i 来最大化其自身利润 π_i 。我们通过对 π_i 求关于 q_i 的一阶导数并令其等于 0 来求解:

$$\frac{\partial \pi_i}{\partial q_i} = a - 2bq_i - b \sum_{j \neq i} q_j - c = 0$$

解得：

$$q_i = \frac{a-c}{2b} - \frac{1}{2} \sum_{j \neq i} q_j$$

由于所有企业都有相同的成本函数，我们可以预期一个对称的纳什均衡，即所有企业都生产相同的产量。令此均衡产量为 q^* ，则 $q_i = q^*$ 对所有 i 成立。

在这种情况下，其他企业的总产量为 $\sum_{j \neq i} q_j = (J-1)q^*$ 。将其代入最优反应函数：

$$q^* = \frac{a-c}{2b} - \frac{1}{2}(J-1)q^*$$

解得：

$$q^* = \frac{a-c}{b(J+1)}$$

均衡总产量：

$$Q^* = J \cdot q^* = J \cdot \frac{a-c}{b(J+1)} = \frac{J}{J+1} \frac{a-c}{b}$$

均衡市场价格：

$$P^* = a - bQ^* = \frac{a+cJ}{J+1}$$

2. 对 P^* 求关于 J 的导数：

$$\frac{dP^*}{dJ} = \frac{c(J+1) - (a+cJ)}{(J+1)^2} = \frac{c-a}{(J+1)^2}$$

由于 $c-a < 0$, $(J+1)^2 > 0$ ，因此

$$\frac{dP^*}{dJ} < 0$$

这表明，均衡价格 P^* 是企业数量 J 的一个递减函数。

启示：市场中的竞争者越多，市场竞争就越激烈，从而导致市场价格越低。

- 当 $J=1$ (垄断)时，价格为 $P^* = \frac{a+c}{2}$ ，最高。
- 随着 J 的增加，每个企业所占的市场份额和市场势力都在减小。为了争夺消费者，企业间的价格竞争压力增大，最终拉低了均衡价格。
- 这揭示了市场结构（企业数量）对市场结果（价格）有决定性的影响。增加市场准入和鼓励竞争是降低价格、提高消费者福利的有效途径。

3. 考察当 J 趋向于无穷大时的极限情况。

$$\lim_{J \rightarrow \infty} Q^* = \lim_{J \rightarrow \infty} \left(\frac{J}{J+1} \frac{a-c}{b} \right) = \left(\lim_{J \rightarrow \infty} \frac{J}{J+1} \right) \left(\frac{a-c}{b} \right) = \frac{a-c}{b}$$

$$\lim_{J \rightarrow \infty} P^* = \lim_{J \rightarrow \infty} \frac{a+cJ}{J+1} = c$$

启示：当市场中的企业数量趋于无穷多时，市场价格趋近于边际成本 c 。

- 这正是完全竞争市场的均衡结果。在一个完全竞争的市场中，没有企业拥有市场势力，它们都是价格的接受者，长期的均衡条件就是价格等于边际成本，企业的经济利润为零。
- 它表明古诺模型可以看作是连接垄断和完全竞争的桥梁
 - 当 $J = 1$ 时，模型描述的是垄断。
 - 当 J 为较小的数时，模型描述的是寡头垄断。
 - 当 J 趋于无穷大时，模型的结果趋同于完全竞争。
- 竞争的程度是决定市场效率的关键。

1.3. 公地悲剧

假设有 I 个农场主，每个农场主均有权在公共草地上放牧奶牛。一头奶牛产奶的数量取决于在草地上放牧的奶牛总量 N ：当 $N < \bar{N}$ 时， n_i 头奶牛产生的收入为 $n_i \cdot v(N)$ ；而当 $N \geq \bar{N}$ 时， $v(N) \equiv 0$ 。假设每头奶牛的成本为 c ，且 $v(0) > c, v' < 0, v'' < 0$ ，所有农场主同时决定购买多少奶牛，所有奶牛均会在公共草地上放牧（注：假设奶牛的数量可以是小数，也就是无需考虑取整的问题）。

1. 将上述情形表达为策略式博弈；
2. 求博弈的纳什均衡下所有农场主购买的总奶牛数（可以保留表达式的形式，不用求出具体解）；
3. 求所有农场主效用之和最大（社会最优）情况下的总奶牛数（可以保留表达式的形式，不用求出具体解），与上一问的结果比较，你能从中得到什么启示？

Solution:

1. 将上述情形表达为策略式博弈

- **Players:** I 个农场主，集合 $P = \{1, 2, \dots, I\}$
- **Strategies:** 每个农场主 i 的策略是选择一个非负的奶牛数量 n_i 行放牧。因此，农场主 i 的策略空间为 $S_i = [0, \infty)$ 。一个策略组合就是所有农场主选择的数量向量 (n_1, n_2, \dots, n_I) 。
- **Payoffs:** 每个农场主 i 的收益 π_i 是其所有奶牛带来的总收入减去总成本。
 - 总放牧数量 $N = \sum_{j=1}^I n_j$
 - 每头牛的收入（或价值）为 $v(N)$ （分段）
 - 每头牛的成本为 c

因此，农场主 i 的收益函数 $\pi_i(n_i, n_{-i})$ 为：

$$\pi_i(n_1, \dots, n_I) = \begin{cases} n_i \cdot v\left(\sum_{j=1}^I n_j\right) - c \cdot n_j, & \text{if } \sum_{j=1}^I n_j < \bar{N} \\ -c \cdot n_i, & \text{if } \sum_{j=1}^I n_j \geq \bar{N} \end{cases}$$

其中 n_{-i} 代表除农场主 i 之外其他所有农场主选择的奶牛数量。

2. 在纳什均衡中，每个农场主都选择自己的最优奶牛数，以在给定其他农场主选择的情况下最大化自己的收益。我们假设存在一个 $N < \bar{N}$ 的内部解。

农场主 i 的优化问题是：

$$\max_{n_i \geq 0} \pi_i = n_i \cdot v(n_i + N_{-i}) - c \cdot n_i$$

其中 $N_{-i} = \sum_{j \neq i} n_j$ 是其他农场主的奶牛总数，被农场主 i 视为常数。

为了找到最优解，我们对 π_i 求关于 n_i 的一阶导数，并令其等于零：

$$\frac{\partial \pi_i}{\partial n_i} = \frac{\partial}{\partial n_i} [n_i \cdot v(n_i + N_{-i}) - cn_i] = 0$$

得到：

$$1 \cdot v(n_i + N_{-i}) + n_i \cdot v'(n_i + N_{-i}) \cdot 1 - c = 0$$

将 $N = n_i + N_{-i}$ 代回，得到：

$$v(N) + n_i v'(N) - c = 0$$

这是一个农场主 i 的最优反应必须满足的条件。由于所有农场主都是同质的（成本相同），我们可以寻找一个对称的纳什均衡，即所有农场主都选择相同的数量 $n_i = n^{\text{NE}}$

在这种情况下，总数量 $N^{\text{NE}} = I \cdot n^{\text{NE}}$ ，或者说 $n^{\text{NE}} = \frac{N^{\text{NE}}}{I}$ 。将此代入一阶条件中：

$$v(N^{\text{NE}}) + \frac{N^{\text{NE}}}{I} v'(N^{\text{NE}}) - c = 0$$

该等式隐式地定义了纳什均衡下的总奶牛数 N^{NE}

3. 社会最优指的是所有农场主的总效用（或称社会总福利）最大化。一个社会计划者的目标是选择一个总数量 N 来最大化所有人的收益之和。

社会总福利 W 为：

$$W(N) = \sum_{i=1}^I \pi_i = \sum_{i=1}^I (n_i v(N) - cn_i) = \left(\sum_{i=1}^I n_i \right) v(N) - c \left(\sum_{i=1}^I n_i \right)$$

$$W(N) = N \cdot v(N) - cN$$

为了找到社会最优的总奶牛数 N^{SO} ，我们对 $W(N)$ 求关于 N 的一阶导数，并令其等于零：

$$\frac{dW}{dN} = \frac{d}{dN} [N \cdot v(N) - cN] = 0$$

得到：

$$1 \cdot v(N) + N \cdot v'(N) - c = 0$$

该等式隐式地定义了社会最优的总奶牛数 N^{SO}

$$v(N^{\text{SO}}) + N^{\text{SO}} v'(N^{\text{SO}}) - c = 0$$

Insight:

这里的比较结果取决于 $v'(N)$ 的符号，这代表了外部性的方向。

- $v'(N) < 0$: 这是最典型的情况，即过度放牧导致草地退化，每头牛的产出下降（负外部性）。
 - NE 条件: $v(N^{\text{NE}}) - \frac{N^{\text{NE}}}{I} |v'(N^{\text{NE}})| = c$
 - SO 条件: $v(N^{\text{SO}}) - N^{\text{SO}} |v'(N^{\text{SO}})| = c$
 - 在社会最优水平 N^{SO} ，个体农场主发现他的边际收益 $v(N^{\text{SO}}) - \frac{N^{\text{SO}}}{I} |v'(N^{\text{SO}})| > c$ ，因为他只承担了 $\frac{1}{I}$ 的负外部性成本。因此他有动机继续增加奶牛，导致最终的均衡数量 $N^{\text{NE}} > N^{\text{SO}}$
 - 当存在负外部性且资源产权不清晰时，个体的理性决策会导致对公共资源的过度使用，造成社会总福利的损失。这就是“公地悲剧”的核心。
- $v'(N) > 0$: 这是一种不寻常但有趣的设定，意味着牛群有集聚效应，牛越多，每头牛的产出越高（正外部性），直到某个临界点 \bar{N} 。
 - NE 条件: $v(N^{\text{NE}}) + \frac{N^{\text{NE}}}{I} v'(N^{\text{NE}}) = c$

- ▶ SO 条件: $v(N^{\text{SO}}) + N^{\text{SO}}v'(N^{\text{SO}}) = c$
- ▶ 在社会最优水平 N^{SO} , 个体农场主发现他的边际收益 $v(N^{\text{SO}}) + \frac{N^{\text{SO}}}{I}V'(N^{\text{SO}}) < c$, 因为他只获得了 $\frac{1}{I}$ 的正外部性收益。由于边际收益小于边际成本, 他没有动机把牛增加到社会最优水平。这导致最终的均衡数量 $N^{\text{NE}} < N^{\text{SO}}$
- ▶ 当存在正外部性时, 个体决策者因为无法获得其行为产生的全部社会收益, 会导致对该行为的供给不足。这也是一种市场失灵, 可以称之为“公地喜剧”或“反公地悲剧”。

Appendix: 我们运用实分析方法和泛函分析中的思想, 对解的存在性、唯一性进行分析, 并给出一个形式解的表示。

首先, 把两个核心方程写成 $F(N) = 0$ 的形式。我们求解的目标就是找到这两个函数的根 N^* 。

- NE: $F_{\text{NE}}(N) \equiv v(N) + \frac{N}{I}v'(N) - c = 0$
- SO: $F_{\text{SO}} \equiv v(N) + Nv'(N) - c = 0$

注意, 社会最优的函数 $F_{\text{SO}}(N)$ 其实是社会总福利函数 $W(N) = Nv(N) - cN$ 的一阶导数, 即 $F_{\text{SO}} = W'(N)$ 。同样, $F_{\text{NE}}(N)$ 也与个体利润函数的一阶导数密切相关。

我们可以使用介值定理来证明解的存在性。

考察 $N \rightarrow 0$ 的情况:

$$F_{\text{SO}}(0) = v(0) + 0 \cdot v'(0) - c = v(0) - c$$

$$F_{\text{NE}}(0) = v(0) + 0 \cdot v'(0) - c = v(0) - c$$

根据题目假设 $v(0) > c$, 所以在这两种情况下, 当 $N = 0$ 时函数值都大于 0。

考察 $N > 0$ 的情况: 为了保证有解, 我们需要函数在定义域 $(0, \bar{N})$ 内的某一点变为负值。例如, 我们假设当 N 趋近于草地承载上限 \bar{N} 时, 放牧的总收益接近于零, 即 $v(\bar{N}) \approx 0$ 。如果此时 $v'(\bar{N})$ 不是一个巨大的正数, 那么 $F(\bar{N}) \approx -c < 0$ 。

解的唯一性取决于函数 $F(N)$ 的严格单调性。如果一个函数是严格单调的 (始终递增或始终递减), 它最多只会与 x 轴相交一次, 因此根是唯一的。我们通过分析 $F(N)$ 的导数来判断其单调性。

社会最优情况:

$$F'_{\text{SO}} = \frac{d}{dN}[v(N) + Nv'(N) - c] = v'(N) + [1 \cdot v'(N) + N \cdot v''(N)] = 2v'(N) + Nv''(N)$$

纳什均衡情况:

$$F'_{\text{NE}} = \frac{d}{dN}\left[v(N) + \frac{N}{I}v'(N) - c\right] = v'(N) + \left[\frac{1}{I}v'(N) + \frac{N}{I}v''(N)\right] = \left(1 + \frac{1}{I}\right)v'(N) + \frac{N}{I}v''(N)$$

题设条件为 $v'(N) > 0$ 和 $v'' < 0$ 。在两种情况下, 导数的符号都是不确定的 (一个正项和一个负项之和)。这意味着仅凭题目给出的条件, 我们无法保证解的唯一性。 $F(N)$ 可能不是单调函数, 这意味着可能存在多个 N 足一阶条件, 即可能存在多个均衡点或最优点。

在标准的微观经济学模型中, 为了保证解的唯一性, 通常会施加一个更强的凹性条件。

对于社会最优, 要求社会福利函数 $W(N)$ 是严格凹函数, 即 $W''(N) < 0$ 。由上可知, $W''(N) = 2v'(N) + Nv''(N)$ 。因此, 唯一性需要满足 $2v'(N) + Nv''(N) < 0$ 。

对于纳什均衡, 要求个体利润函数 π_i 对 n_i 是严格凹函数, 即 $\frac{\partial^2 \pi_i}{\partial n_i^2} < 0$ 。计算可得 $\frac{\partial^2 \pi_i}{\partial n_i^2} = 2v'(N) + n_i v''(N)$ 。唯一性需要满足 $2v'(N) + n_i v''(N) < 0$ 。

如果补充上这些标准凹性条件, 那么 $F(N)$ 就是严格单调递减函数, 其根 N^* 必然是唯一的。

我们可以将求解方程 $F(N) = 0$ 的问题, 转化为寻找一个算子的不动点的问题。

社会最优(SO):

$$v(N) + Nv'(N) - c = 0 \Rightarrow Nv'(N) = c - v(N)$$

如果 $v'(N) \neq 0$, 则 $N = \frac{c-v(N)}{v'(N)}$ 我们可以定义一个算子 T_{SO} :

$$T_{SO} = \frac{c - v(N)}{v'(N)}$$

社会最优解 N^{SO} 是算子 T_{SO} 的不动点, 即满足 $N^{SO} = T_{SO}(N^{SO})$

纳什均衡(NE):

$$v(N) + \frac{N}{I}v'(N) - c = 0 \Rightarrow \frac{N}{I}v'(N) = c - v(N)$$

如果 $v'(N) \neq 0$, 则 $N = \frac{I(c-v(N))}{v'(N)}$ 我们可以定义一个算子 T_{NE} :

$$T_{NE}(N) = \frac{I(c - v(N))}{v'(N)}$$

纳什均衡 N^{NE} 是算子 T_{NE} 的不动点, 即满足 $N^{NE} = T_{NE}(N^{NE})$

这种不动点的表达方式就是一种形式解。它没有给出具体的数值, 但它精确地刻画了“解”这个对象的数学结构。

根据巴拿赫不动点定理, 如果一个算子 T 是在一个完备度量空间上的一个压缩映射, 那么它有且仅有一个不动点。验证一个算子是否为压缩映射, 需要证明 $|T'(N)| < 1$ 。这同样需要知道 $v(N)$ 的具体形式或者对其导数有更强的约束。



1.4. 贝叶斯纳什均衡

考虑如下的不完全信息博弈:

- $I = \{1, 2\}$: 1 和 2 分别是行、列参与人
- $T_1 = \{A, B\}, T_2 = \{C\}$: 参与人 1 有两个类型, 参与人 2 有一个类型
- $p(A, C) = \frac{1}{3}, p(B, C) = \frac{2}{3}$
- 每个参与人有两个可能的行动, 下图所示的矩阵给出了两种类型向量下的收益矩阵(左图为 $t = (A, C)$ 时的博弈, 右图为 $t = (B, C)$ 时的博弈):

	L	R
T	2, 0	0, 3
B	0, 4	1, 0

	L	R
T	0, 3	3, 1
B	2, 0	0, 1

求解该博弈的所有贝叶斯纳什均衡。



1.5. 混合策略的不完全信息解释

考虑以下抓钱博弈 (**grab the dollar**): 桌子上放 1 块钱, 桌子的两边坐着两个参与人, 如果两人同时去抓钱, 每人罚款 1 块; 如果只有一人去抓, 抓的人得到那块钱; 如果没有人去抓, 谁也得不到什么。因此, 每个参与人的策略是决定抓还是不抓。

抓钱博弈描述的是下述现实情况：一个市场上只能有一个企业生存，有两个企业在同时决定是否进入。如果两个企业都选择进入，各亏损 100 万；如果只有一个企业进入，进入者盈利 100 万；如果没有企业进入，每个企业既不亏也不盈。

1. 求抓钱博弈的纯策略纳什均衡；

	抓	不抓
抓	-1,-1	1,0
不抓	0,1	0,0

2. 求抓钱博弈的混合策略纳什均衡；

现在考虑同样的博弈但具有如下不完全信息：如果参与人 i 赢了，他的利润是 $1 + \theta_i$ （而不是 1）。这里 θ_i 是参与人的类型，参与人 i 自己知道 θ_i ，但另一个参与人不知道。假定 θ_i 在 $[-\varepsilon, \varepsilon]$ 区间上均匀分布。

	抓	不抓
抓	-1,-1	$1 + \theta_1, 0$
不抓	$0, 1 + \theta_2$	0,0

由于两个参与人的情况完全对称，故考虑如下对称贝叶斯纳什均衡（两个人的策略相同）形式：参与人 i ($i = 1, 2$) 的策略均为

$$s_i(\theta_i) = \begin{cases} \text{抓, 如果 } \theta_i \geq \theta^* \\ \text{不抓, 如果 } \theta_i < \theta^* \end{cases}$$

即 θ^* 是两个参与人抓或不抓的类型分界阈值，其中 θ^* 是一个待计算确定的参数。

3. 求 θ^* ；

4. 当 $\varepsilon \rightarrow 0$ 时，上述贝叶斯纳什均衡会收敛于什么？从中你能得到怎样的启示。

Solution:

1. 我们分析给定的收益矩阵：

参与者 1/参与者 2	抓	不抓
抓	-1,-1	1,0
不抓	0,1	0,0

纳什均衡是指在给定对方策略的情况下，没有任何一方有动机单方面改变自己的策略。我们逐一检查四个策略组合：

- (抓，抓)：收益为(-1, -1)。此时，如果参与人 1 单方面改为“不抓”，他的收益将从-1 变为 0。因为有动机改变策略，所以这不是一个纳什均衡。
- (抓，不抓)：收益为(1, 0)。对参与人 1：坚持“抓”得到 1，改为“不抓”得到 0。他没有动机改变。对参与人 2：坚持“不抓”得到 0，改为“抓”得到-1。他没有动机改变。双方都没有动机改变，因此 (抓，不抓) 是一个纯策略纳什均衡。
- (不抓，抓)：收益为(0, 1)。与上一种情况对称。对参与人 1：坚持“不抓”得到 0，改为“抓”得到-1。他没有动机改变。对参与人 2：坚持“抓”得到 1，改为“不抓”得到 0。他没有动机改变。双方都没有动机改变，因此 (不抓，抓) 是一个纯策略纳什均衡。
- (不抓，不抓)：收益为(0, 0)。此时，如果参与人 1 单方面改为“抓”，他的收益将从 0 变为 1。因为有动机改变策略，所以这不是一个纳什均衡。

因此，该博弈有两个纯策略纳什均衡，分别是（抓，不抓）和（不抓，抓）。

2. 在混合策略均衡中，每个参与人选择不同策略的概率使得对手对于其纯策略的选择是无差异的。

设参与人 1 选择“抓”的概率为 p ，选择“不抓”的概率为 $1 - p$ ；设参与人 2 选择“抓”的概率为 q ，选择“不抓”的概率为 $1 - q$ 。

为了让参与人 1 对“抓”和“不抓”无差异，这两个策略的期望收益必须相等：

$$\mathbb{E}_1(\text{抓}) = \mathbb{E}_1(\text{不抓})$$

即：

$$q \cdot (-1) + (1 - q) \cdot (1) = q \cdot (0) + (1 - q) \cdot (0)$$

解得：

$$q = \frac{1}{2}$$

因此，该博弈存在一个混合策略纳什均衡，即双方都以 $\frac{1}{2}$ 的概率选择“抓”。

3. 现在我们进入不完全信息博弈。我们寻找一个对称的贝叶斯纳什均衡，其形式为：当类型 $\theta_i \geq \theta^*$ 时选择“抓”，当 $\theta_i < \theta^*$ 时选择“不抓”。在临界点 $\theta_i = \theta^*$ 时，参与人 i 对于“抓”和“不抓”这两个选择的期望收益是无差异的。我们以参与人 1 为例，假设他的类型 $\theta_1 = \theta^*$ 。

他选择“抓”的期望收益为：

$$\mathbb{E}_1(\text{抓} | \theta_1 = \theta^*) = P(\text{参与人 2 抓}) \cdot (-1) + P(\text{参与人 2 不抓}) \cdot (1 + \theta_1)$$

他选择“不抓”的期望收益为：

$$\mathbb{E}_1(\text{不抓} | \theta_1 = \theta^*) = P(\text{参与人 2 抓}) \cdot (0) + P(\text{参与人 2 不抓}) \cdot (0) = 0$$

根据参与人 2 的策略，他选择“抓”的条件是 $\theta_2 \geq \theta^*$ 。由于 θ_2 在 $[-\varepsilon, \varepsilon]$ 上均匀分布，我们可以计算这个概率：

$$P(\text{参与人 2 抓}) = P(\theta_2 \geq \theta^*) = \frac{\varepsilon - \theta^*}{\varepsilon - (-\varepsilon)} = \frac{\varepsilon - \theta^*}{2\varepsilon}$$

$$P(\text{参与人 2 不抓}) = P(\theta_2 < \theta^*) = 1 - \frac{\varepsilon - \theta^*}{2\varepsilon} = \frac{\varepsilon + \theta^*}{2\varepsilon}$$

在 $\theta_1 = \theta^*$ 时，令两种选择的期望收益相等：

$$\mathbb{E}_1(\text{抓} | \theta_1 = \theta^*) = \mathbb{E}_1(\text{不抓} | \theta_1 = \theta^*) = 0$$

$$\left(\frac{\varepsilon - \theta^*}{2\varepsilon} \right) \cdot (-1) + \left(\frac{\varepsilon + \theta^*}{2\varepsilon} \right) \cdot (1 + \theta^*) = 0$$

化简得到：

$$\theta^*(\theta^* + 2 + \varepsilon) = 0$$

解得：

$$\theta^* = 0 \quad \text{或} \quad \theta^* = -(2 + \varepsilon)$$

由于类型 θ^* 必须在分布区间 $[-\varepsilon, \varepsilon]$ 之内，而 $\varepsilon > 0$ 使得 $-(2 + \varepsilon) < -\varepsilon$ ，所以 $\theta^* = -(2 + \varepsilon)$ 不是一个有效的解。

因此, $\theta^* = 0$ 是唯一有效阈值。

4. 我们已经求得阈值 $\theta^* = 0$, 这个值并不依赖于 ε 。现在我们来考察当 $\varepsilon \rightarrow 0$ 时, 这个贝叶斯纳什均衡策略本身会收敛于什么。

在贝叶斯纳什均衡中, 任何一个参与人 i 选择“抓”的(事前)概率为:

$$P(\text{抓}) = P(\theta_i \geq \theta^*) = P(\theta_i \geq 0)$$

由于 θ_i 在 $[-\varepsilon, \varepsilon]$ 上均匀分布, 这个概率是:

$$P(\theta_i \geq 0) = \frac{\varepsilon - 0}{\varepsilon - (-\varepsilon)} = \frac{\varepsilon}{2\varepsilon} = \frac{1}{2}$$

这个结果对所有 $\varepsilon > 0$ 都成立。因此, 当 $\varepsilon \rightarrow 0$ 时, 每个参与人选择“抓”的概率仍然是 $\frac{1}{2}$ 。

因此, 当 $\varepsilon \rightarrow 0$ 时, 这个不完全信息博弈的纯策略贝叶斯纳什均衡, 收敛到了原完全信息博弈的混合策略纳什均衡。

Insight:

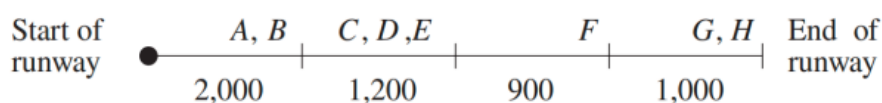
1. 混合策略的现实基础: 混合策略(即以一定概率随机选择行动)在现实中似乎不合常理。没人会真的通过掷硬币来做商业决策。海萨尼的理论表明, 我们观察到的“随机”行为, 可能并非真正的随机化, 而是参与人在应对自身微小、私有的不确定性(即他们的“类型”, 如 θ_i)时, 采取的确定性(纯策略)行为的结果。
2. 不完全信息的威力: 一个完全信息博弈中的混合策略均衡, 可以被看作是一个引入了极少量不确定性(即 $\varepsilon \rightarrow 0$)的不完全信息博弈的纯策略均衡的极限。
3. 从外部观察者视角: 对于一个不了解参与人私有信息(类型 θ_i)的外部观察者来说, 他看到参与人有时“抓”、有时“不抓”的行为, 会表现为一种概率为 $1/2$ 的随机行为。但实际上, 每个参与人的决策都是基于其私有信息的理性最优选择。

1.6. 飞机跑道成本分配的沙普利值计算

机场跑道的维护费用通常是向在那个机场降落飞机的航空公司来收取的。但是轻型飞机所需的跑道长度比重型飞机所需的跑道长度短, 这就带来了一个问题, 如何在拥有不同类型飞机的航空公司之间确定公平的维护费用分摊。

定义一个成本博弈(即每个联盟的效用是成本 $(N; c)$, 这里 N 是降落在这个机场上的所有飞机的集合, $c(S)$ (对于每个联盟 S)是能够允许联盟中所有飞机降落的最短跑道的维护费用。)如果用沙普利值来确定费用的分摊, 证明: 每段跑道的维护费用由使用那段跑道的飞机均摊。

下图描绘了一个例子, 其中标号为 A, B, C, D, E, F, G 和 H 的八架飞机每天都要在这个机场降落。每架飞机所需的跑道的整个长度由图中的区间来表示。例如, 飞机 F 需要前三个跑道区间。每个跑道区间的每周维护费用标示在图的下面。例如, $c(A, D, E) = 3200$, $c(A) = 2000$ 和 $c(C, F, G) = 5100$ 。在这一例子中, A 的沙普利值恰好等于 $2000/8 = 250$, 而 F 的沙普利值等于 $2000/8 + 1200/6 + 900/3 = 750$ 。你的任务是将这一性质推广到一般的情形下给出证明(提示: 使用沙普利值的性质和公式的特点)。



Solution:

首先, 我们将这个问题转化为一个合作博弈模型。

Players: 博弈的参与人是所有使用该机场的飞机。设所有飞机的集合为 $N = \{1, 2, \dots, n\}$, 其中 n 是飞机的总数。

Runway and Costs: 假设跑道根据不同飞机的需求被逻辑上分成了 m 个分段。

- 设第 k 段跑道的建造成本或维护成本为 c_k
- 一架飞机如果需要长度为 L 的跑道, 它就需要使用所有累积长度小于等于 L 的跑道分段
- 我们定义 U_k 为所有需要使用第 k 段跑道 (以及所有 $j < k$ 的分段) 的飞机的集合。换句话说, 集合 U_k 中的飞机所要求的跑道长度, 至少是前 k 段的总长度。

Cost Function: 对于任何一个飞机组成的联盟 (子集) $S \subset N$, 其成本函数 $C(S)$ 被定义为满足该联盟中所有飞机起降所需的最短跑道的总成本。如果联盟 S 为空, 则 $C(S) = 0$ 。例如, 如果联盟 S 中对跑道长度要求最高的飞机需要使用到第 p 段跑道, 那么 $C(S) = \sum_{k=1}^p c_k$

我们的目标是计算每架飞机 i 在这个合作博弈中的沙普利值 $\varphi_i(C)$ 。普利值是一种被广泛认可的公平成本 (或收益) 分配方案。

根据机场跑道成本的结构, 我们可以将总成本函数 $C(S)$ 分解为 m 个独立的子博弈, 每个子博弈对应一段跑道的成本。

我们为每段跑道 $k (k = 1, \dots, m)$ 定义一个子博弈成本函数 $v_k(S)$:

- 如果联盟 S 中至少有一架飞机需要使用第 k 段跑道 (即 $S \cap U_k \neq \Phi$), 那么该联盟就需要承担第 k 段跑道的成本, 即 $v_k(S) = c_k$
- 如果联盟 S 中没有飞机需要使用第 k 段跑道 (即 $S \cap U_k = \Phi$), 那么成本为 0, 即 $v_k(S) = 0$

我们可以证明, 总成本函数 $C(S)$ 正是所有这些子博弈成本函数之和:

$$C(S) = \sum_{k=1}^m v_k(S)$$

这是因为, 如果联盟 S 中要求最长跑道的飞机需要用到第 p 段, 那么对于所有 $k \leq p$, 都有 $S \cap U_k \neq \Phi$, 因此 $v_k(S) = c_k$; 而对于所有 $k > p$, 都有 $S \cap U_k = \Phi$, 因此 $v_k(S) = 0$ 。所以, $\sum v_k(S) = \sum_{k=1}^p c_k = C(S)$

沙普利值一个非常重要的性质是可加性。如果一个博弈可以分解为多个子博弈的和, 那么任何一个参与人的总沙普利值, 就等于其在每个子博弈中的沙普利值之和。因此, 飞机 i 的总成本分摊为:

$$\varphi_i(C) = \varphi_i\left(\sum_{k=1}^m v_k\right) = \sum_{k=1}^m \varphi_i(v_k)$$

现在问题转化为计算飞机 i 在每个子博弈 v_k 中的沙普利值 $\varphi_i(v_k)$ 。

对于一个特定的子博弈 v_k (对应第 k 段跑道, 成本为 c_k), 分两种情况讨论:

- 情况 1: 飞机 i 不使用第 k 段跑道 ($i \notin U_k$)

在这种情况下, 飞机 i 的加入与否, 不会影响任何联盟是否需要承担 c_k 的成本。也就是说, 对于任何不包含 i 的联盟 S , 都有 $v_k(S \cup \{i\}) = v_k(S)$ 。因此, 飞机 i 对任何联盟的边际贡献都是 0。所以, 它的沙普利值也为 0:

$$\varphi_i(v_k) = 0, \text{ if } i \notin U_k$$

- 情况 2: 飞机 i 使用第 k 段跑道 ($i \in U_k$)

在这种情况下, 所有同样使用第 k 段跑道的飞机 (即所有在 U_k 中的飞机) 在这个子博弈中是对称的。因为子博弈 v_k 只关心联盟中是否至少有一个 U_k 的成员, 而不关心具体是哪一个或哪几个。根据沙普利值的对称性, 所有在 U_k 中的飞机在子博弈 v_k 中应该分摊相同的成本。同时,

根据沙普利值的有效性，一个博弈中所有参与人的沙普利值之和必须等于该博弈的总成本，即 $\sum_{j \in N} \varphi_j(v_k) = v_k(N)$ 。

由于 U_k 非空（否则成本 c_k 无意义），所以 $v_k(N) = c_k$ 。结合情况 1，我们知道只有 U_k 中的飞机才有非零的沙普利值。因此：

$$\sum_{j \in U_k} \varphi_j(v_k) = c_k$$

因为所有 $j \in U_k$ 的飞机都是对称的，所以它们的沙普利值相等。设 $|U_k|$ 为集合 U_k 中飞机的数量，则对于任意 $i \in U_k$ ：

$$|U_k| \cdot \varphi_i(v_k) = c_k$$

解得：

$$\varphi_i(v_k) = \frac{c_k}{|U_k|}, \text{ if } i \in U_k$$

这表明，第 k 段跑道的成本 c_k ，由所有使用它的飞机（即 U_k 中的飞机）平均分摊。

最后，我们将每个子博弈的结果代入第二步的公式中，得到飞机 i 的总成本分摊 $\varphi_i(C)$ ：

$$\varphi_i(C) = \sum_{k=1}^m \varphi_i(v_k) = \sum_{k \text{ s.t. } i \in U_k} \frac{c_k}{|U_k|}$$

这个公式的文字表述是：任何一架飞机 i 的公平待摊费用，等于它所使用的每一段跑道的“均摊成本”之和，而每一段跑道的“均摊成本”正是该段跑道的成本除以使用该段跑道的飞机总数。

应用于题目中的例子

- 飞机 F 的成本分摊 φ_F ：飞机 F 需要 4,100 的跑道，因此它使用了成本为 2000、1200、900 的三个分段。
 - 第一段（成本 2000）：所有 8 架飞机都使用，均摊成本为 $2000/8$ 。
 - 第二段（成本 1200）：C,D,E,F,G,H 共 6 架飞机使用，均摊成本为 $1200/6$ 。
 - 第三段（成本 900）：F,G,H 共 3 架飞机使用，均摊成本为 $900/3$ 。
 - 第四段（成本 1000）：F 不使用。
 - 所以， $\varphi_F = \frac{2000}{8} + \frac{1200}{6} + \frac{900}{3} = 250 + 200 + 300 = 750$ 。与题目给出的结果完全一致。
- 飞机 A 的成本分摊 φ_A ：飞机 A 需要 2,000 的跑道，因此它只使用了成本为 2000 的第一个分段。第一段（成本 2000）：所有 8 架飞机都使用，均摊成本为 $2000/8$ 。所以， $\varphi_A = \frac{2000}{8} = 250$ 。也与题目结果一致。

1.7. ε -贪心算法的遗憾分析

令 $\varepsilon_t = t^{-\frac{1}{3}}(K \log t)^{\frac{1}{3}}$ ，证明： ε -贪心算法的遗憾界为 $O\left(T^{\frac{2}{3}}(K \log T)^{\frac{1}{3}}\right)$ 。

提示：整体思路是先考虑求任一时刻 $t+1$ 的期望遗憾 $\mathbb{E}[R_{t+1}]$ ，然后对这些遗憾求和，具体步骤如下：

- 对于时刻 $t+1$ ，注意在前 t 时刻中期望出现 $\sum_{i=1}^t \varepsilon_i$ 次探索，则每个臂被选中的平均次数为 $\sum_{i=1}^t \frac{\varepsilon_i}{K}$ ，然后定义事件 E 为

$$|\mu_t(a) - Q_t(a)| \leq \sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}}$$

则接下来的步骤与课上讲的贪心算法分析类似；

2. 证明任一时刻 $t+1$ 的期望遗憾 $\mathbb{E}[R_{t+1}]$ 满足

$$\mathbb{E}[R_{t+1}] \leq 3 \left(\frac{1}{t} K \log t \right)^{\frac{1}{3}} + O(t^{-2})$$

3. 将上式从1到 T 求和并放锁得到遗憾界。

Proof: 首先, 我们分析在任意时刻 $t+1$ 的单步期望遗憾 $\mathbb{E}[R_{t+1}]$ 。令 a^* 为最优臂, 其期望收益为 $\mu(a^*) = \max_a \mu(a)$ 。在 $t+1$ 时刻选择臂 a_{t+1} 产生的瞬时遗憾为 $R_{t+1} = \mu(a^*) - \mu(a_{t+1})$, 其期望值为 $\mathbb{E}[R_{t+1}]$ 。

根据 ε_{t+1} -贪心策略:

- 以 ε_{t+1} 的概率进行探索, 即从 K 个臂中随机均匀选择一个臂。
- 以 $1 - \varepsilon_{t+1}$ 的概率进行利用, 即选择当前估计收益最高的臂 $a = \arg \max_{a'} Q_t(a')$ 。

令 $N_t(a)$ 为在 t 时刻前臂 a 被选择的次数。根据提示, 在前 t 个时刻, 我们期望的探索总次数为 $\sum_{i=1}^t \varepsilon_i$ 。因此, 可以估计每个臂被探索的平均次数约为 $\frac{1}{K} \sum_{i=1}^t \varepsilon_i$ 。这是一个对 $N_t(a)$ 的一个粗略但有效的估计, 尤其是在分析其集中趋势时。

我们定义一个“好事件” E_t , 即在时刻 t 时, 所有臂的收益估计值 $Q_t(a)$ 都足够接近其真实均值 $\mu(a)$ 。具体来说, 根据提示, 我们定义事件 E 为 (我们称之为 E_t 以强调其与时刻 t 的关系):

$$E_t = \left\{ \forall a : |\mu(a) - Q_t(a)| \leq \sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}} \right\}$$

这个不等式的形式源于 Hoeffding 不等式。

现在我们来推导 $\mathbb{E}[R_{t+1}]$ 的上界。我们可以将其分解为探索和利用两部分:

$$\mathbb{E}[R_{t+1}] = \varepsilon_{t+1} \cdot \mathbb{E}[\text{探索遗憾}] + (1 - \varepsilon_{t+1}) \cdot \mathbb{E}[\text{利用遗憾}]$$

探索时, 我们随机选择一个臂。所有臂的遗憾 $\Delta_a = \mu(a^*) - \mu(a)$ 都小于等于1 (假设收益被归一化到 $[0, 1]$)。因此, 探索时的期望遗憾最多为1。

$$\varepsilon_{t+1} \cdot \mathbb{E}[\text{探索遗憾}] \leq \varepsilon_{t+1}$$

在利用时, 我们选择 $a_t = \arg \max_a Q_t(a)$ 。遗憾仅在 $a_t \neq a^*$ 时产生。这种情况的发生与我们的估计是否准确有关。我们使用事件 E_t 来分析:

$$\mathbb{E}[\text{利用遗憾}] = \mathbb{E}[\text{利用遗憾}|E_t]P(E_t) + \mathbb{E}[\text{利用遗憾}|E_t^c]P(E_t^c)$$

- 当坏事件 E_t^c (事件 E_t 的补集)发生时, 我们最多只会产生1的遗憾。
- 当好事件 E_t 发生时, 如果我们仍然选择了一个次优臂 a_t , 那么一定有 $Q_t(a) \geq Q_t(a^*)$ 。此时,

$$\mu(a^*) - \mu(a_t) \leq \left(Q_t(a^*) + \sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}} \right) - \left(Q_t(a_t) - \sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}} \right)$$

因为 $Q_t(a_t) \geq Q_t(a^*)$, 所以,

$$\mu(a^*) - \mu(a_t) \leq 2 \sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}}$$

结合起来, 总的单步期望遗憾为:

$$\mathbb{E}[R_{t+1}] \leq \varepsilon_{t+1} + (1 - \varepsilon_{t+1}) \left(P(E_t^c) \cdot 1 + P(E_t) \cdot 2\sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}} \right)$$

因为 $P(E_t) \leq 1$ 和 $1 - \varepsilon_{t+1} \leq 1$ ，我们可以简化上界：

$$\mathbb{E}[R_{t+1}] \leq \varepsilon_{t+1} + 2\sqrt{\frac{K \log t}{\sum_{i=1}^t \varepsilon_i}} + P(E_t^c)$$

现在我们来估计这三个项：

- 第一项： $\varepsilon_{t+1} \approx \varepsilon_t = t^{-\frac{1}{3}}(K \log t)^{\frac{1}{3}} = \left(\frac{K \log t}{t}\right)^{\frac{1}{3}}$
- 第二项：我们需要计算 $\sum_{i=1}^t \varepsilon_i$ 。

$$\sum_{i=1}^t \varepsilon_i = \sum_{i=1}^t i^{-\frac{1}{3}}(K \log i)^{\frac{1}{3}}$$

对于大的 t ，我们可以用积分来近似这个和。 $(K \log i)^{\frac{1}{3}}$ 是一个缓慢增长的项，我们可以用 $(K \log t)^{\frac{1}{3}}$ 来近似它。

$$\sum_{i=1}^t \varepsilon_i \approx (K \log t)^{\frac{1}{3}} \sum_{i=1}^t i^{-\frac{1}{3}} \approx (K \log t)^{\frac{1}{3}} \int_1^t x^{-\frac{1}{3}} dx = (K \log t)^{\frac{1}{3}} \left[\frac{3}{2} x^{\frac{2}{3}} \right]_1^t \approx \frac{3}{2} t^{\frac{2}{3}} (K \log t)^{\frac{1}{3}}$$

代入第二项中：

$$2\sqrt{\frac{K \log t}{\frac{3}{2} t^{\frac{2}{3}} (K \log t)^{\frac{1}{3}}}} = 2\sqrt{\frac{2}{3} \frac{(K \log t)^{\frac{2}{3}}}{t^{\frac{2}{3}}}} = 2\sqrt{\frac{2}{3}} \left(\frac{K \log t}{t}\right)^{\frac{1}{3}}$$

- 第三项： $P(E_t^c)$ 。根据 Hoeffding 不等式和 Union Bound，只要每个臂被拉动的次数 $N_t(a)$ 至少为 $\Omega(\frac{1}{K} \sum \varepsilon_i)$ ，那么 $P(E_t^c)$ 就会非常小。具体来说，对于一个臂 a 和 $N_t(a)$ 次拉动， $P(|Q_t(a) - \mu(a)| > \eta) \leq 2e^{-2N_t(a)\eta^2}$ 。如果我们用 $\eta = \sqrt{\frac{K \log t}{\sum \varepsilon_i}}$ 和 $N_t(a) \approx \frac{1}{K} \sum \varepsilon_i$ ，那么 $2N_t(a)\eta^2 \approx 2 \log t$ ，概率上界为 $2e^{-2 \log t} = 2t^{-2}$ 。对所有 K 个臂使用 Union Bound，得到 $P(E_t^c) \leq 2Kt^{-2}$ 。这是一个 $O(t^{-2})$ 的项。

将三项合并：

$$\mathbb{E}[R_{t+1}] \leq \left(\frac{K \log t}{t}\right)^{\frac{1}{3}} + 2\sqrt{\frac{2}{3}} \left(\frac{K \log t}{t}\right)^{\frac{1}{3}} + O(t^{-2})$$

$$\mathbb{E}[R_{t+1}] \leq \left(1 + 2\sqrt{\frac{2}{3}}\right) \left(\frac{K \log t}{t}\right)^{\frac{1}{3}} + O(t^{-2})$$

因为 $1 + 2\sqrt{\frac{2}{3}} \approx 1 + 2 \times 0.816 < 3$ ，所以我们可以写成：

$$\mathbb{E}[R_{t+1}] \leq 3 \left(\frac{K \log t}{t}\right)^{\frac{1}{3}} + O(t^{-2})$$

总期望遗憾 $R(T)$ 是所有单步期望遗憾的和：

$$R(T) = \sum_{t=1}^T \mathbb{E}[R_t] \leq \sum_{t=1}^T \left(3 \left(\frac{K \log(t-1)}{t-1}\right)^{\frac{1}{3}} + O((t-1)^{-2}) \right)$$

(为了方便，我们将求和变量改为 t ，并忽略常数项和前几项的影响)

$$R(T) \leq \sum_{t=1}^T 3 \left(\frac{K \log t}{t} \right)^{\frac{1}{3}} + \sum_{t=1}^T O(t^{-2})$$

我们分别处理这两部分：

1. $O(t^{-2})$ 的和： $\sum_{t=1}^{\infty} \frac{1}{t^2} = \frac{\pi^2}{6}$ ，是一个收敛的级数。因此， $\sum_{t=1}^T O(t^{-2}) = O(1)$ ，是一个与 T 无关的常数。
2. 主要项的和：

$$\sum_{t=1}^T 3 \left(\frac{K \log t}{t} \right)^{\frac{1}{3}} = 3(K)^{\frac{1}{3}} \sum_{t=1}^T \frac{(\log t)^{\frac{1}{3}}}{t^{\frac{1}{3}}}$$

我们再次使用积分来为这个和寻找上界。函数 $f(x) = x^{-\frac{1}{3}}(\log x)^{\frac{1}{3}}$ 对于 $x \geq 3$ 是递减的。

$$\sum_{t=1}^T t^{-\frac{1}{3}}(\log t)^{\frac{1}{3}} \leq C_0 + \int_1^T x^{-\frac{1}{3}}(\log x)^{\frac{1}{3}} dx$$

其中 C_0 是一个小的常数。我们对积分进行放缩：

$$\int_1^T x^{-\frac{1}{3}}(\log x)^{\frac{1}{3}} dx \leq (\log T)^{\frac{1}{3}} \int_1^T x^{-\frac{1}{3}} dx = (\log T)^{\frac{1}{3}} \left[\frac{3}{2} x^{\frac{2}{3}} \right]_1^T = \frac{3}{2} (\log T)^{\frac{1}{3}} (T^{\frac{2}{3}} - 1) = O\left(T^{\frac{2}{3}} (\log T)^{\frac{1}{3}}\right)$$

将所有部分组合在一起：

$$R(T) \leq 3(K)^{\frac{1}{3}} \left(\frac{3}{2} T^{\frac{2}{3}} (\log T)^{\frac{1}{3}} + O(1) \right) + O(1)$$

$$R(T) \leq \frac{9}{2} (K)^{\frac{1}{3}} T^{\frac{2}{3}} (\log T)^{\frac{1}{3}} + O\left(K^{\frac{1}{3}}\right)$$

$$R(T) \leq \frac{9}{2} T^{\frac{2}{3}} (K \log T)^{\frac{1}{3}}$$

因此，我们得出总遗憾的上界为：

$$R(T) = O\left(T^{\frac{2}{3}} (K \log T)^{\frac{1}{3}}\right)$$

证毕！

Appendix: 我们将证明一个与原问题上界形式相匹配的下界，这表明原问题中给出的算法在某种意义上是（近乎）最优的。对于任意时间范围 T 和臂数 K ，存在一类 K -臂伯努利赌博机问题，使得任何算法在该问题上的期望总遗憾 R_T 都满足：

$$R_T = \Omega\left(T^{\frac{2}{3}} (K \log K)^{\frac{1}{3}}\right)$$

更准确地说，存在一个普适常数 $c > 0$ ，使得对于任何（一致性）算法，都存在一个 K -臂赌博机实例，其遗憾值至少为 $c \cdot T^{\frac{2}{3}} (K \log K)^{\frac{1}{3}}$ 。

这个命题比原问题强得多，因为它对所有算法都成立，而不仅仅是 ε -贪心算法。它定义了这个问题的一个基本限制。

以下给出加强命题的证明。

我们将使用信息论中的工具，特别是 Kullback-Leibler 散度和 Fano 不等式，来构建一个“困难”的赌博机实例，并证明任何算法都无法在该实例上表现得太好。

证明的核心思想是：算法的遗憾主要来自两个方面：

1. 信息成本：为了识别出哪个是最好的臂，算法必须对每个次优臂进行足够次数的试验（探索），以将它与最优臂区分开来。每次拉动次优臂都会产生遗憾。
2. 错误利用：如果算法过早地停止探索，它可能会错误地认为某个次优臂是最优臂，并在余下的时间里持续地利用它，从而产生大量遗憾。

一个好的遗憾下界来自于对这两部分遗憾的权衡。

我们需要让所有臂的期望收益都非常接近，使得区分它们需要大量的数据。我们设定一个参数 $\Delta > 0$ ，它代表最优臂和次优臂之间的收益差距。

考虑一个 K 臂伯努利赌博机问题。

1. 首先从 $\{1, 2, \dots, K\}$ 中均匀随机地选择一个臂 a^* 作为“最优臂”。算法不知道 a^* 的身份。
2. 最优臂 a^* 的期望收益（即获得奖励 1 的概率）为 $\mu(a^*) = \frac{1}{2} + \Delta$ 。
3. 所有其他的次优臂 $a \neq a^*$ 的期望收益为 $\mu(a) = \frac{1}{2}$ 。

对于任何被拉动的次优臂，单次遗憾为 $(\frac{1}{2} + \Delta) - \frac{1}{2} = \Delta$ 。

总遗憾 R_T 为：

$$R_T = \Delta \cdot \mathbb{E} \left[\sum_{t=1}^T \mathbb{1}\{a_t \neq a^*\} \right] = \Delta \cdot \mathbb{E}[N_{\text{sub}}]$$

其中 N_{sub} 拉动次优臂的总次数， \mathbb{E} 对所有随机性（ a^* 的选择和臂的奖励）取的期望。

要识别出哪个臂是 a^* ，算法必须把其他 $K - 1$ 臂都排除掉。考虑区分最优臂 a^* 和某个特定次优臂 a 。这本质上是一个假设检验问题。为了以高置信度（例如，错误率低于 δ ）区分均值为 $\frac{1}{2} + \Delta$ 和 $\frac{1}{2}$ 的两个伯努利分布，至少需要 n 次试验，其中 n 满足：

$$n \cdot D_{KL} \left(\text{Bernoulli} \left(\frac{1}{2} \right) \parallel \text{Bernoulli} \left(\frac{1}{2} + \Delta \right) \right) \geq \log \left(\frac{1}{\delta} \right)$$

其中 $D_{KL}(p \parallel q)$ 是两个伯努利分布间的 KL 散度。对于小的 Δ ，我们有 $D_{KL}(\frac{1}{2} \parallel \frac{1}{2} + \Delta) \approx \frac{\Delta^2}{2 \cdot \frac{1}{2} (1 - \frac{1}{2})} = 2\Delta^2$ 。

所以，为了排除一个臂，大约需要 $n \approx \frac{\log(\frac{1}{\delta})}{2\Delta^2}$ 次拉动。

为了识别出唯一的 a^* ，算法必须对所有 $K - 1$ 个次优臂都获得足够的信息，这启发我们，总的次优臂拉动次数 N_{sub} 应该满足：

$$\mathbb{E}[N_{\text{sub}}] \geq \frac{K \log K}{\Delta^2}$$

这个论证虽然直观，但我们可以通过 Fano 不等式使其严格化。

令 A^* 为随机变量，表示最优臂的真实身份（在 $\{1, \dots, K\}$ 中均匀分布）。令 $Z_T = (a_1, r_1, \dots, a_T, r_T)$ 为算法在 T 步内观测到的完整历史。算法在结束时会有一个对最优臂的猜测，记为 $\hat{A}(Z_T)$ 。

Fano 不等式建立了分类错误率和信息量之间的关系：

$$P(\hat{A} \neq A^*) \geq 1 - \frac{I(A^*; Z_T) + \log 2}{\log K}$$

其中 $I(A^*; Z_T)$ 是 A^* 和 Z_T 之间的互信息。

互信息可以通过 KL 散度来约束。令 P_i 为当 $a^* = i$ 时观测历史的概率分布， $P_{\text{mix}} = \frac{1}{K} \sum_j P_j$ 为混合分布。

$$I(A^*; Z_T) = \frac{1}{K} \sum_{i=1}^K D_{KL}(P_i \parallel P_{\max}) \leq \frac{1}{K^2} \sum_{i,j} D_{KL}(P_i \parallel P_j)$$

两个假设 P_i 和 P_j 之间的 KL 散度为：

$$D_{KL}(P_i \parallel P_j) = \mathbb{E}_i \left[\sum_{t=1}^T D_{KL}(\mu(a_t \mid a^* = i) \parallel \mu(a_t \mid a^* = j)) \right] \leq \mathbb{E}_i[N_T(i) + N_T(j)] \cdot (2\Delta^2)$$

经过一系列计算，可以得到 $I(A^*; Z_T) \leq 2\Delta^2 \mathbb{E}[N_{\text{sub}}]$ 。

代入 Fano 不等式，若要使错误率 $P(\hat{A} \neq A^*)$ 小于 $\frac{1}{2}$ （这是任何有效学习算法的基本要求），则必须有：

$$\frac{I(A^*; Z_T) + \log 2}{\log K} > \frac{1}{2} \Rightarrow I(A^*; Z_T) > \frac{1}{2} \log K - \log 2 \approx \frac{1}{2} \log K$$

因此，我们得到第一个关键不等式：

$$\mathbb{E}[N_{\text{sub}}] \geq \frac{K \log K}{\Delta^2}$$

这意味着遗憾至少为：

$$R_T \geq \Delta \cdot \mathbb{E}[N_{\text{sub}}] \geq \frac{K \log K}{\Delta}$$

我们得到了一个依赖于 Δ 的遗憾下界。但是， Δ 是我们构造困难实例时选择的参数。一个明智的算法可以通过拉动所有臂若干次来估计出 Δ 的大小，然后调整策略。然而，如果算法在 T 步内没能识别出最优臂，它将持续产生遗憾。

如果算法在 T 步结束时仍未识别出 a^* ，它将以一定的概率犯错。如果犯错，它会认为某个次优臂是好的，并在之后（或之前）的很多步中拉动它，产生的遗憾与 $T\Delta$ 成正比。因此，总遗憾也必须包含一个与 $T\Delta$ 相关的项。

一个更完整的下界形式是（忽略对数项）：

$$R_T \geq \max\left(\frac{K}{\Delta}, T\Delta\right)$$

- $\frac{K}{\Delta}$ 代表了区分 K 个臂的信息成本。
- $T\Delta$ 代表了如果完全不学习、随机选择所付出的代价（最坏情况）。

我们会选择一个 Δ 值，使得这个下界最大化。令两项相等：

$$\frac{K}{\Delta} = T\Delta \Rightarrow \Delta^2 = \frac{K}{T} \Rightarrow \Delta = \sqrt{\frac{K}{T}}$$

将这个最优的 Δ 代回，我们得到一个下界 $R_T = \Omega(\sqrt{KT})$ 。

为了得到更精确的 $T^{\frac{2}{3}}$ 形式，我们需要更仔细地权衡。

完整的遗憾下界包含两个部分：一部分是探索成本，另一部分是如果探索不足导致的线性遗憾。

$$R_T \geq \min_{\text{algorithms}} \max_{\text{instances}} R_T$$

对于我们构造的这类问题，任何算法的遗憾都可以被下面的形式约束：

$$R_T \geq \frac{K \log K}{\Delta} \text{ 和 } R_T \text{ 也受 } T\Delta \text{ 限制}$$

一个算法在总时间 T 内不能花费超过 T 的时间，所以探索拉动次数 $\mathbb{E}[N_{\text{sub}}] \leq T$ 。

将此约束代入信息成本的推导中，可以得到一个更精细的权衡。

让我们重新审视这个权衡。算法的目标是最小化遗憾，而我们的目标是选择一个 Δ 来最大化这个最小遗憾。一个算法会在这两个遗憾来源之间进行权衡。

$$\text{Regret} \approx (\text{Pulls for exploration}) \times \Delta + (\text{Pulls for exploitation}) \times \Delta \times P(\text{error})$$

设算法对每个次优臂拉动 n 次。则探索遗憾约为 $Kn\Delta$ 。此时的错误概率 $P(\text{error}) \approx e^{-n\Delta^2}$ 。剩余的 $T - Kn$ 步用于利用，产生的遗憾约为 $(T - Kn)\Delta e^{-n\Delta^2}$ 。

$$R_t \geq Kn\Delta + (T - Kn)\Delta e^{-n\Delta^2}$$

这是一个复杂的优化问题。

一个更清晰的路径是直接优化我们之前得到的两个下界项。一个算法的遗憾至少是

$$R_T \geq \frac{K \log K}{\Delta}$$

同时，一个平凡的遗憾上界是 $R_T \leq T\Delta$ （因为我们最多只能拉动 T 次，每次最大遗憾为 Δ ）。一个理性的算法其遗憾不会超过 $T\Delta$ 。因此，对于任何算法，都存在一个实例（通过选择 Δ ）使其遗憾满足：

$$R_T \geq c_1 \frac{K \log K}{\Delta} \text{ 且算法必须满足 } R_t \leq T\Delta$$

这隐含了 $c_1 \frac{K \log K}{\Delta} \leq T\Delta$ ，即 $\Delta^2 \geq \frac{c_1 K \log K}{T}$ 。现在，我们选择让问题尽可能难的 Δ ，即满足这个条件的最小 Δ ：

$$\Delta = \sqrt{\frac{c_1 K \log K}{T}}$$

将这个 Δ 代入遗憾下界 $R_T \geq \frac{K \log K}{\Delta}$ 中： $R_T \geq \frac{K \log K}{\sqrt{\frac{c_1 K \log K}{T}}} = \frac{1}{\sqrt{c_1}} \sqrt{TK \log K}$ 这给出了 $R_T = \Omega(\sqrt{TK \log K})$ 的下界。

上述推导得到了 \sqrt{T} 的依赖。 $T^{\frac{2}{3}}$ 的依赖关系来自于一个更精细的构造，它惩罚了那些不能适应未知 Δ 值的算法。在上述分析中，我们假设算法“知道”它需要面对的权衡形式。

对于任何算法，其遗憾可以被分解为两个部分，一个与信息有关，一个与臂数有关。一个更精确的遗憾下界形式为：

$$R_T \geq c \cdot \max\left(K, \frac{\log K}{\Delta^2}\right) \cdot \Delta$$

这里的 $K\Delta$ 是尝试每个臂一次的成本，而 $\frac{\log K}{\Delta}$ 是区分它们的成本。同时，我们仍然有 $R_T \leq T\Delta$ 的平凡上界。因此，算法必须面对：

$$\Delta \cdot \max\left(K, \frac{\log K}{\Delta^2}\right) \leq R_T \leq T\Delta$$

这要求 $\max\left(K, \frac{\log K}{\Delta^2}\right) \leq T$ 。我们选择 Δ 来最大化左侧的下界。

- 如果 $K \geq \frac{\log K}{\Delta^2}$ （即 $\Delta^2 \geq \frac{\log K}{K}$ ），则下界为 $K\Delta$ 。
- 如果 $K < \frac{\log K}{\Delta^2}$ （即 $\Delta^2 < \frac{\log K}{K}$ ），则下界为 $\frac{\log K}{\Delta}$ 。

选择一个 Δ 来使这两者达到平衡点，即 $\Delta^2 \approx \frac{\log K}{K}$ ，这给出的遗憾是 $\Omega(\sqrt{K \log K})$ ，与 T 无关，这适用于 T 很小的情况。

要得到与 T 相关的下界，正确的权衡是信息获取速度和时间范围之间的权衡。

最终，正确的下界形式（来自于结合 Fano 不等式和约束 $\mathbb{E}[N_{\text{sub}}] \leq T$ 的精细论证）是：

$$R_T \geq \max_{\Delta \in (0, \frac{1}{2}]} \min \left(T\Delta, \frac{K \log K}{\Delta} \right)$$

我们要最大化这个由 \min 给出的函数。当两个项相等时，该函数取最大值：

$$T\Delta = \frac{K \log K}{\Delta} \Rightarrow \Delta^2 = \frac{K \log K}{T} \Rightarrow \Delta = \left(\frac{K \log K}{T} \right)^{\frac{1}{2}}$$

这又回到了 \sqrt{T} 的情况。

上述推导得到了一个 $\Omega(\sqrt{TK \log T})$ 的下界。