

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

Information Disclosure in Dynamic Innovation Contests

Jussi Keppo

NUS Business School and Institute of Operations Research and Analytics
National University of Singapore, Singapore
keppo@nus.edu.sg

Linsheng Zhuang

Institute of Operations Research and Analytics
National University of Singapore, Singapore
linsheng.z@u.nus.edu

...

Key words: Digital Economy, Data Protection Regulation, Innovation Contest

1. Introduction

- Kaggle¹ ...
- Meta-kaggle dataset [Risdal and Bozsolik \(2022\)](#).

1.1. Literature Review

This paper focuses on the two players innovation contest with a continuous time where the players' relative position is public information throughout the game. This is closely related to tug-of-war contest, which, to our knowledge, was first formally given by [Harris and Vickers \(1987\)](#) as a one-dimensional simplification of the multi-stage R&D race. The output processes are model by Brownian motions drifted with effort inputs, which is followed by [Budd et al. \(1993\)](#) who model the state of a dynamic competition of two innovative duopoly firms by a Brownian motion drifted by the effort gap, and solve the equilibrium approximately. Furthermore, [Moscarini and Smith \(2007\)](#) model the tug-of-war state as the gap of the two outputs directly, and draw an analytical equilibrium of the pure strategies.

¹ <https://www.kaggle.com>

...

Information disclosure in contest - [Bimpikis et al. \(2019\)](#).

...

Closest paper - [Ryvkin \(2022\)](#).

...

2. The Model

Two players, i and j , compete for a prize $\theta > 0$ in a contest. Winner gets the prize and loser gets nothing. The contest starts at time zero. At every time $t \geq 0$, the representative player i chooses an effort level $q_{i,t}$ and burdens a quadratic cost $C_i(q_{i,t}) = c_i q_{i,t}^2 / 2$, with a lower c_i corresponding to higher ability. Denoted by y_t the *output gap* of player i and j at time t , driven by

$$dy_t = (q_{i,t} - q_{j,t})dt + \sigma dW_t \quad (1)$$

where W_t is a Brownian motion and $\sigma > 0$ measures the innovation risk.

The contest is equipped with a submission system that allows participants to upload their algorithms at any time and receive immediate feedback. For simplicity, we further assume that agents submit their intermediate results whenever they make progress. This setup enables the contest organizers to monitor all players' progress $x_{i,t}$ and $x_{j,t}$ in real time. Moreover, the true output level, evaluated by the system, is only known by the game designer but not the two players. At any time $t > 0$, the contest designer emits a *public* signal of the real output gap y_t . The signal is ambiguous and the game holder controls the ambiguity. The dynamic of signal is

$$dZ_t = y_t dt + \frac{dB_t}{\sqrt{\lambda}} \quad (2)$$

where B_t is standard Brownian motion independent with $(W_{i,t})$ and $(W_{j,t})$, and the parameter λ is set by the game holder to control the precision of signal. The larger the λ , the more accurate the signal would be.

The information set of both players at time $t \geq 0$ is $I_t \equiv \{Z_s : 0 \leq s \leq t\}$. Player i estimates the unknown output gap y_t based on the information set I_t . Let $\tilde{y}_t \equiv E(y_t | I_t)$ be the estimated output gap and $S_t \equiv E[(\tilde{y}_t - y_t)^2 | I_t]$ be the estimation variance. According to Chapter 1.2 of [Bensoussan \(1992\)](#), *Kalman-Bucy filter* gives the dynamics of \tilde{y}_t and S_t ,

$$d\tilde{y}_t = (q_{i,t} - q_{j,t})dt + \lambda S_t (dZ_t - \tilde{y}_t dt) \quad (3)$$

$$\frac{dS_t}{dt} = \sigma^2 - \lambda S_t^2 \quad (4)$$

Hence, the conditional distribution $y_t|I_t \sim \mathcal{N}(\tilde{y}_t, S_t|I_t)$ is fully captured by the mean \tilde{y}_t and variance S_t . If $\lambda = 0$, we have $S_t = S_0 + \sigma^2 t$, i.e., the estimation variance is increasing in time linearly. If $\lambda > 0$, the solution of (4) is

$$S_t = \begin{cases} \bar{S} \cdot \tanh \left\{ t \cdot \sigma \sqrt{\lambda} + \tanh^{-1} (S_0/\bar{S}) \right\} & \text{if } S_0 < \bar{S} \\ \bar{S} & \text{if } S_0 = \bar{S} \\ \bar{S} \cdot \coth \left\{ t \cdot \sigma \sqrt{\lambda} + \coth^{-1} (S_0/\bar{S}) \right\} & \text{if } S_0 > \bar{S} \end{cases} \quad (5)$$

Specifically, $\bar{S} = \sigma/\sqrt{\lambda}$ when $\lambda > 0$ and $\bar{S} = \infty$ when $\lambda = 0$. Please refer to Appendix A for the derivations. Figure 1 shows the evolution of S_t in time: estimation variance S_t converges to *steady state* \bar{S} as time goes by regardless of the starting estimation variance. For simplicity, we henceforth assume that $S_0 = \bar{S}$, hence $S_t \equiv \bar{S}$.



Figure 1 The evolution of S_t in time t , given that $\lambda = 1$ and $\sigma = 1$.

Following Ryvkin (2022), let's consider a dynamic contest with a fixed deadline. Suppose the contest is terminated when time $t = T > 0$. Since the steady state estimation variance \bar{S} is fixed as displayed above, the state of the game is fully characterized by a tuple (\tilde{y}_t, t) . At any time $0 \leq t < T$, player i optimizes her effort level $q_{i,\tau}$ in the remaining contest period $\tau \in [t, T)$ according to the following optimization problem,

$$V^i(\tilde{y}_t, t; q_{j,t}, \Theta_i) = \max_{\{q_{i,\tau}\}_{\tau=t}^T} \mathbb{E} \left(\theta \cdot 1_{\tilde{y}_T > 0} - \int_t^T C_i(q_{i,\tau}) d\tau \middle| I_t \right) \quad (6)$$

where $\Theta_i \equiv \{\theta, \lambda, \sigma, c_i\}$, subject to constraints (3), (4) and $q_{i,\tau} \geq 0$ for all $\tau \in [t, T)$. The optimization problem for player j is just symmetric to that of player i as $V^j(\tilde{y}_t, t) = V^i(-\tilde{y}_t, t)$. The corresponding Hamilton-Jacobi-Bellman (HJB) equation for player i is

$$0 = \max_{q_{i,t} \geq 0} \left[-\frac{c_i q_i^2}{2} + V_y^i \cdot (q_{i,t} - q_{j,t}) + V_t^i + \frac{V_{yy}^i}{2} \lambda \bar{S}^2 \right]$$

By definition, we have $\lambda \bar{S}^2 = \sigma^2$. Under the assumption of inner solution, we plug into the first order conditions $q_{i,t} = V_y^i/c_i$ and $q_{j,t} = -V_y^j/c_j$, we have the system of equations

$$\begin{aligned} \frac{1}{2c_i}(V_y^i)^2 + \frac{1}{c_j}V_y^iV_y^j + V_t^i + V_{yy}^i \frac{\sigma^2}{2} &= 0 \\ \frac{1}{2c_j}(V_y^j)^2 + \frac{1}{c_i}V_y^jV_y^i + V_t^j + V_{yy}^j \frac{\sigma^2}{2} &= 0 \end{aligned}$$

subject to boundary conditions $V^i(-\infty, t) = 0$, $V^i(+\infty, t) = \theta$, $V^j(-\infty, t) = \theta$, $V^j(+\infty, t) = 0$, $V^i(\tilde{y}_T, T) = \theta \cdot 1_{\tilde{y}_T > 0}$ and $V^j(\tilde{y}_T, T) = \theta \cdot 1_{\tilde{y}_T < 0}$.

The Nash equilibrium is summarized in the following lemma. We include a simplified version of the proof in the appendix:

LEMMA 1 (Ryvkin 2022). *In the Markov perfect equilibrium, the players' efforts in state $m_{i(j)}(\tilde{y}_t, t) : \mathbb{R} \times [0, T) \rightarrow \mathbb{R}_+$ are given by*

$$m_{i(j)}(\tilde{y}_t, t) = \frac{e^{-z^2/2}}{\sqrt{2\pi\sigma^2(T-t)}} \cdot \frac{\sigma^2}{2} [\gamma(\rho_i) + \gamma(\rho_j)] [1 - \rho(z)^2] [1 \pm \rho(z)] \quad (7)$$

where $z = \tilde{y}_t/(\sigma\sqrt{T-t})$, $\rho(z) = \gamma^{-1}(\Phi(z)[\gamma(\rho_i) + \gamma(\rho_j)] - \gamma(\rho_j))$ and

$$\begin{aligned} \gamma(u) &= \frac{u}{1-u^2} + \frac{1}{2} \ln \frac{1+u}{1-u}, \quad u \in (-1, 1) \\ \rho_i &= \frac{e^{w_i} + e^{-w_j} - 2}{e^{w_i} - e^{-w_j}}, \quad \rho_j = \frac{e^{w_j} + e^{-w_i} - 2}{e^{w_j} - e^{-w_i}}, \quad w_{i(j)} = \frac{\theta}{\sigma^2 c_{i(j)}}. \end{aligned}$$

The variables $w_{i(j)}$ represent the abilities of two players, while $\rho_{i(j)}$ normalizes $w_{i(j)}$ into the interval $(-1, 1)$. It is not hard to see that $\gamma(\cdot)$ is strictly increasing on $(-1, 1)$, ranging from $-\infty$ to $+\infty$. Moreover, the equilibrium effort $m_{i(j)}$ can be represented to the product of $\phi(y; 0, \sigma^2(T-t))$, the probability density of normal distribution with mean zero and variance $\sigma^2(T-t)$ at the state y and an amplitude factor that only depends on the composite variable z .

3. Estimation

In this section, we describe the estimation procedure. We first outline the data generation process, establishing the connection between the empirical data and the theoretical model discussed previously. Then, we introduce a structural estimation method using Bayesian framework.

3.1. Data Generating Process

For each contest on Kaggle, the observable information can be classified into three primary components:

The first component consists of essential contest details, including the contest duration, prize structure, information disclosure policy, and other governing rules. Contrary to the assumptions

of our model, a typical contest usually involves multiple teams rather than just two. We index the participating teams of the contest by $i \in \{1, 2, \dots, n\}$. To fully leverage the potential of the data and establish a connection with our theoretical model, let's assume that each participant perceives a competitor they are playing against at every moment, denoted as j . This perceived competitor is typically understood as the most prominent individual on the leaderboard, i.e., the person ranked first. When team i themselves hold the top position, their perceived competitor is the individual who poses the greatest threat, namely the person ranked second. Furthermore, the contests may feature intricate prize structures, such as the provision of multiple awards, rather than adhering to a winner-takes-all format. The issue will be discussed in Section 5.

The second component captures the submission events of each player i to the system, denoted by the sequence $\{\hat{t}_k^i\}_{k=1}^{N_i}$. Here, N_i represents for the total number of submissions by player i , and t represents for the time of each submission. We understand the submission events of player i and j as two conditional independent inhomogeneous Poisson processes, driven by the intensity functions $\tau_i(t)$ and $\tau_j(t)$.² Then, during any time interval \mathcal{S} of the contest duration \mathcal{T} , the Poisson arrival rate of submissions of the representative player i is given by $\int_{s \in \mathcal{S}} \tau_i(s) ds$. We assume the submission intensity $\tau_i(t)$ is proportional to the effort level $m_i(\tilde{y}_t, t)$. More specifically,

$$\tau_i(t) = r \cdot m_i(\tilde{y}_t, t) \quad (8)$$

where $r > 0$ is the common ratio of submission intensity and effort.

The third component of the observed data is an open leaderboard that records the real-time rankings and scores of each participant, denoted by \hat{x}_t^i . We interpret the difference in scores between i and j displayed on the leaderboard as the signal Z_t (defined in (2)) released by the contest organizer. Specifically, let's denote $\hat{y}_t = \hat{x}_t^i - \hat{x}_t^j$ the gap between displayed scores, and interpret it as the signal intentionally released by the contest organizer:

$$dZ_t = \hat{y}_t dt \quad (9)$$

As indicated in (2), we assume that \hat{y}_t represents a noisy signal, with its precision governed by λ . Moreover, according to (1), the observed real-time gap \hat{y}_t on leaderboard evolves as

$$d\hat{y}_t = [m_i(\tilde{y}_t, t) - m_j(\tilde{y}_t, t)] dt + \sigma dW_t + \frac{dB_t}{\sqrt{\lambda}} \quad (10)$$

where the term σdW_t captures the innovation shock, and $dB_t/\sqrt{\lambda}$ represents the signal noise. It is important to recognize that the generation of \hat{y}_t is inherently tied to the players' strategic interactions, as it depends on their estimates of the underlying state, \tilde{y}_t . In turn, \tilde{y}_t evolves dynamically

² That is, given the intensity functions $\tau_i(t)$ and $\tau_j(t)$, the submission events $\{\hat{t}_k^i\}_{k=1}^{N_i}$ and $\{\hat{t}_k^j\}_{k=1}^{N_j}$ are mutually independent.

based on \hat{y}_t , since players continually update their beliefs in response to observed data. As a result, the generation of \hat{y}_t and \tilde{y}_t proceeds jointly. By (3) and (9), the estimated gap \tilde{y}_t by the two players is jointly generated by

$$d\tilde{y}_t = [m_i(\tilde{y}_t, t) - m_j(\tilde{y}_t, t)] dt + \sqrt{\lambda} \sigma (\hat{y}_t - \tilde{y}_t) dt \quad (11)$$

In practice, the leaderboard signals are inherently noisy, as organizers deliberately disclose only a subset of the full dataset to participants to mitigate the risk of overfitting. The proportion of the released data is generally known to all participants. In addition to the public leaderboard, most competitions hosted on Kaggle also maintain a private leaderboard, where organizers evaluate the true predictive performance of participants' models using the full dataset.

Beyond the uncertainty introduced by partial data disclosure, a second source of signal noise arises from the timing of leaderboard updates: rankings are refreshed only after model submissions, creating a delay relative to the true standings. It should be noted that such lag would bias our model estimates only if participants strategically timed their submissions. Although such strategic behaviour is theoretically possible, we abstract from it in this study. We assume that each submission follows a period of substantive effort, allowing the leaderboard rankings to broadly reflect the relative performance of algorithms based on the publicly available data.

3.2. Bayesian Framework

By definition, the likelihood function of a realization of the representative player i 's submission events $\{t_k^i\}_{k=1}^{N_i}$ is given by

$$p\left(\{\hat{t}_k^i\}_{k=1}^{N_i} | \tau_i\right) = \exp\left\{-\int_{s \in \mathcal{T}} \tau_i(s) ds\right\} \prod_{k=1}^{N_i} \tau_i(\hat{t}_k^i) \quad (12)$$

where τ_i is defined in (8). Similar to the other player j . Next, suppose we sample the leaderboard gap \hat{y}_t at time points (t_1, t_2, \dots, t_N) , obtaining observations $\{\hat{y}_{t_k}\}_{k=1}^N$. Let $t_0 = 0$ and suppose the initial gap y_0 follows $\mathcal{N}(\mu_0, \sigma_0)$. Then, according to the assumed data-generating process of \hat{y}_t in (10), the corresponding likelihood function is

$$p\left(\{\hat{y}_{t_k}\}_{k=1}^N | m_i, m_j\right) = \phi(y_0 | \mu_0, \sigma_0) \times \prod_{k=0}^{N-1} \phi\left(\hat{y}_{t_{k+1}} - \hat{y}_{t_k} \mid \int_{t_k}^{t_{k+1}} m_i(\tilde{y}_s, s) - m_j(\tilde{y}_s, s) ds, \left(\sigma^2 + \frac{1}{\lambda}\right)(t_{k+1} - t_k)\right) \quad (13)$$

where $\hat{y}_{t_0} = y_0$. Here, we assume that time is discretized into uniform intervals of length Δ , so that the integrals in the above equations can be approximated by summations.

To evaluate the likelihood functions (12) and (13), we must first compute the equilibrium trajectories of both the perceived output gap (\tilde{y}_t) and the effort levels $(m_i(\tilde{y}_t, t), m_j(\tilde{y}_t, t))$, as implied

by equations (7) and (11). Importantly, the construction of \tilde{y} and the effort functions m_i and m_j depends on the underlying (unobserved) parameters c_i , c_j , σ , λ , r , μ_0 and σ_0 , which are themselves subject to estimation. We denote these underlying parameters collectively as \mathcal{P} .

Once these parameters are specified, the equilibrium paths of the perceived output gap (\tilde{y}_t) and the corresponding effort levels ($m_i(\tilde{y}_t, t), m_j(\tilde{y}_t, t)$) can be deterministically computed. However, evaluating $m_{i(j)}(\tilde{y}_t, t)$ via equation (7) requires numerically approximating the inverse function γ^{-1} , which poses challenges for the use of gradient-based Markov Chain Monte Carlo (MCMC) sampling methods, such as Hamiltonian Monte Carlo (HMC) and the No-U-Turn Sampler (NUTS). Hence, we approximate this inverse function with an analytical form³:

$$\gamma^{-1}(x) \approx \frac{2}{\pi} \arctan(0.856 \cdot x) \quad (14)$$

Figure 2 compares the equilibrium effort function $m_i(\tilde{y}_t, t)$ derived from the approximate analytical form of γ^{-1} in equation (14) with that obtained from the numerically accurate solution. As illustrated, the approximation closely replicates the true function.

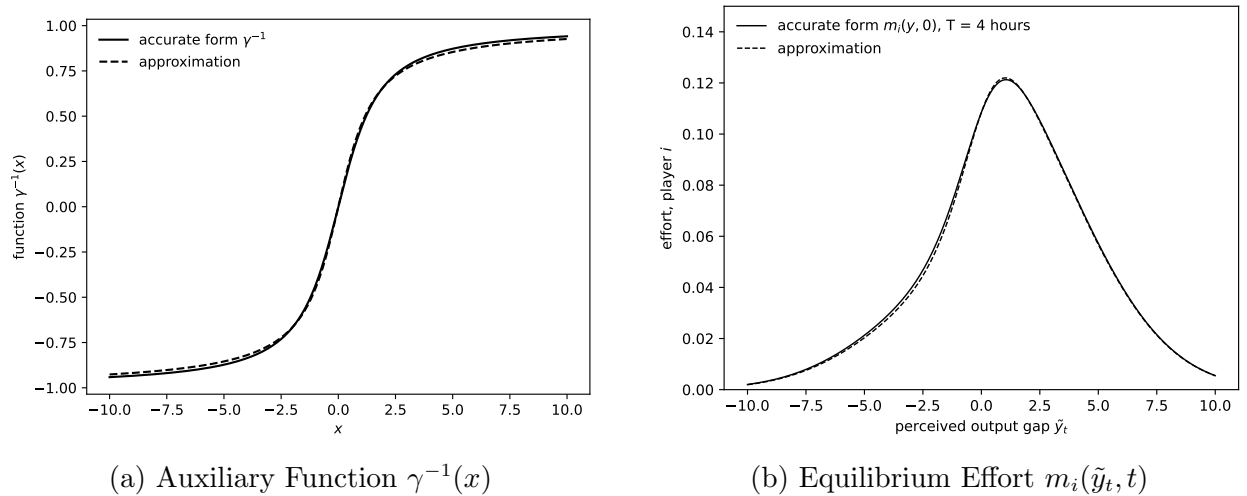


Figure 2 Comparison of the Accurate and Approximate Forms ($\theta = 1$, $c_i = c_j = 1$, $\sigma = 1$, $\Delta = 1/24$)

Unknown parameters to be estimated:

- σ (Total innovation risk): prior - Gamma
- c_i and c_j (Capacities): prior - Gamma

4. Synthetic Data

Before applying our estimation procedure to real-world contest data, we evaluate its potential on synthetically generated data.

³ The parameter $a = 0.856$ is obtained by minimizing the infinity norm of the difference between the numerical inverse of $\gamma(\cdot)$ and the approximation $\frac{2}{\pi} \arctan(ax)$. The resulting maximum approximation error, measured in the infinity norm, is approximately 0.019.

4.1. Data Generation

A central challenge in generating synthetic data lies in dynamically constructing a point process that conforms to an inhomogeneous Poisson process. We first generate candidate submission events according to a homogeneous Poisson process over the whole contest duration, using a fixed high intensity $\tau^* \geq \sup_t \tau_{i(j)}(t)$. Then, we apply the classical thinning procedure (Lewis and Shedler 1979) to determine whether each candidate event is accepted.

Algorithm 1 Synthetic Data Simulation

Input: $\Delta, T, c_i, c_j, \sigma, \lambda, r, \tau^*, y_0, \mu_0$

Sample $\{s_k^{i(j)}\}_{k=1}^{N_{i(j)}}$ from homogeneous Poisson process $(\tau_{i(j)}^*)$ on $[0, T)$

Sample $\{u_k^{i(j)}\}_{k=1}^{N_{i(j)}}$ from uniform distribution on $[0, 1]$

Sample series of Brownian motions (W_t) and (B_t)

Initialize $\ell^{i(j)} = 0, \hat{y}_0 = 0, \tilde{y}_t = \mu_0$

for $t = 0$ to T **do**

$m_{i(j)}(\tilde{y}_t, t) \leftarrow (7), \tau_{i(j)}^*(t) \leftarrow (8)$ {Use $\theta, c_{i(j)}, \sigma, r, T$ }

$y_{t+\Delta} \leftarrow (1); \tilde{y}_{t+\Delta} \leftarrow (11)$ {Use $\sigma, \lambda, \Delta, W_{t+\Delta}$ }

for $s_k^{i(j)} \in [t, t + \Delta)$ **do**

if $u_k^{i(j)} < \tau_{i(j)}(s_k^{i(j)})/\tau_{i(j)}^*$ **then**

$\hat{t}_\ell^{i(j)} \leftarrow s_k^{i(j)}; \ell^{i(j)} \leftarrow \ell^{i(j)} + 1$ {Accept the submission event}

$\hat{y}_{t+\Delta} \leftarrow (10)$ {Use $\sigma, \lambda, B_{t+\Delta}$ }

end if

end for

end for

Fill in all missing entries in (\hat{y}_t) by propagating the last observed value.

Output: $(\tilde{y}_t), (\hat{y}_t), \{\hat{t}_k^{i(j)}\}_{k=1}^{N_{i(j)}}, (m_{i(j)}(\tilde{y}_t, t))$

4.2. Bayesian Inference

Based on the following assumptions, we will next generate the submission data for two players in a data analysis competition:

- Total uncertainty: $\sigma = 5$
- Capacities: $c_i = 1.5$ and $c_j = 2$
- Contest duration: from 2025-01-01 to 2025-03-31
- Starting point $\tilde{y}_0 = 0$

4.3. Discussions

...

5. Case Study

...

To Do:

regression: λ and the proportion of data release

6. Conclusion

Appendix A: Solve S_t in Equation (4)

If S is in steady state $dS/dt = 0 \Leftrightarrow S = \bar{S} \equiv \sigma/\sqrt{\lambda}$. If S is not in steady state, i.e. $S \neq \bar{S}$, we first isolate the two variables and get

$$\frac{dS}{\sigma - \lambda S^2} = dt$$

Then, we take the integral on both sides

$$t = \int \frac{dS}{\sigma - \lambda S^2} = \frac{1}{\sigma\sqrt{\lambda}} \int \frac{dS\sqrt{\lambda}/\sigma}{1 - (S\sqrt{\lambda}/\sigma)^2} \equiv \frac{1}{\sigma\sqrt{\lambda}} \int \frac{du}{1 - u^2}$$

where $u = S\sqrt{\lambda}/\sigma = S/\bar{S}$. Hence,

$$\sigma\sqrt{\lambda} \cdot t = \begin{cases} \tanh^{-1}(u) - K_1, & \text{if } |u| < 1 \\ \coth^{-1}(u) - K_2, & \text{if } |u| > 1 \end{cases} = \begin{cases} \tanh^{-1}(S/\bar{S}) - K_1, & \text{if } S < \bar{S} \\ \coth^{-1}(S/\bar{S}) - K_2, & \text{if } S > \bar{S} \end{cases}$$

Thus, we conclude the non-steady state case that

$$S = \begin{cases} \bar{S} \cdot \tanh(\sigma\sqrt{\lambda} \cdot t + K_1), & \text{if } S < \bar{S} \\ \bar{S} \cdot \coth(\sigma\sqrt{\lambda} \cdot t + K_2), & \text{if } S > \bar{S} \end{cases}$$

Finally, we determine the constants K_1, K_2 by the initial condition S_0 and have

$$K_1 = \tanh^{-1}(S_0/\bar{S})$$

$$K_2 = \coth^{-1}(S_0/\bar{S})$$

References

- Bensoussan A (1992) Stochastic Control of Partially Observed Systems (Cambridge University Press), 1 edition, ISBN 9780511526503, URL <http://dx.doi.org/10.1017/CB09780511526503>.
- Bimpikis K, Ehsani S, Mostagir M (2019) Designing dynamic contests. Operations Research 67(2):339–356, URL <http://dx.doi.org/10.1287/opre.2018.1823>.
- Budd C, Harris C, Vickers J (1993) A model of the evolution of duopoly: Does the asymmetry between firms tend to increase or decrease? The Review of Economic Studies 60(3):543–573.
- Harris C, Vickers J (1987) Racing with uncertainty. The Review of Economic Studies 54(1):1–21, ISSN 00346527, 1467937X, URL <http://www.jstor.org/stable/2297442>.
- Lewis PAW, Shedler GS (1979) Simulation of nonhomogeneous poisson processes by thinning. Naval Research Logistics Quarterly 26(3):403–413, ISSN 0028-1441, 1931-9193, URL <http://dx.doi.org/10.1002/nav.3800260304>.
- Moscarini G, Smith L (2007) Optimal dynamic contests. Working Paper 1–18, science (June).
- Risdal M, Bozsolik T (2022) Meta kaggle. URL <http://dx.doi.org/10.34740/KAGGLE/DS/9>.
- Ryvkin D (2022) To fight or to give up? dynamic contests with a deadline. Management Science 68(11):8144–8165, ISSN 0025-1909, 1526-5501, URL <http://dx.doi.org/10.1287/mnsc.2021.4206>.