

部分観測問題におけるThompson抽出

土屋 平 / Taira Tsuchiya

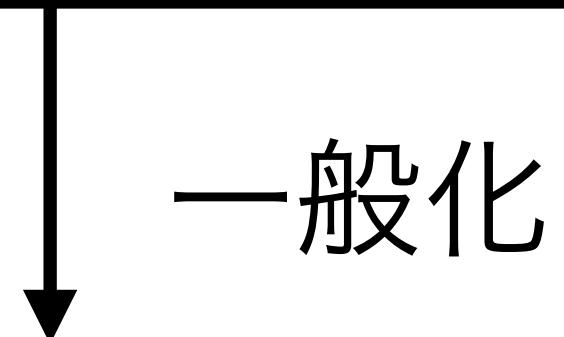


* 本発表は本多淳也先生、杉山将先生
との共同研究に基づきます。

Toshiba Symposium 2020, 2020年12月22日

講演の概要

前半: バンディット問題とThompson抽出



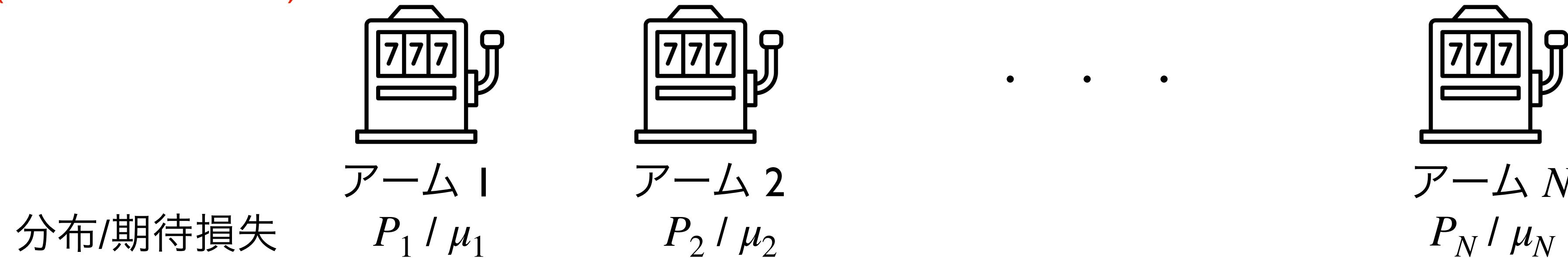
後半: 部分観測問題におけるThompson抽出の有効性に関する研究の紹介

[Tsuchiya, Honda & Sugiyama NeurIPS 2020]

T. Tsuchiya, J. Honda, & M. Sugiyama,
Analysis and Design of Thompson Sampling for Stochastic Partial Monitoring, In NeurIPS 2020.

多腕バンディット問題

- 未知の分布の集合 $(P_i)_{i=1,\dots,N}$ が備わった、逐次的意思決定モデル
 (= アーム、行動)



For round $t = 1, \dots, T$:

1. プレイヤーがアーム $i(t) \in \{1, \dots, N\}$ を選択する
2. アーム $i(t)$ の確率的な損失 ($\sim P_{i(t)}$) を観測する

引いたアーム $i(t)$ についてのみ
損失が観測される

目標: ラウンド全体の期待損失を最小化する

- 「探索」と「活用」のトレードオフに対処する必要がある

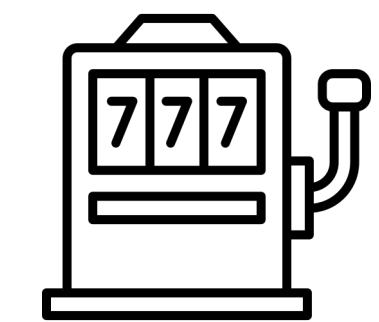
情報の少なく最適になり得る
アームを引く

現時点で最適に見える
アームを引く

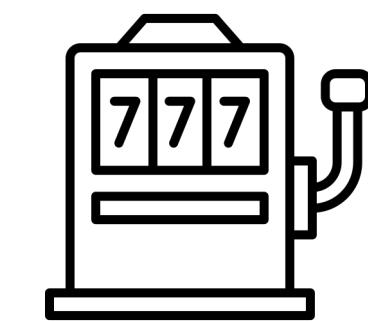
Thompson抽出 [Thompson 1933]

- あらゆる逐次的意思決定問題において経験的に最も有用な方策の1つ

例. 2腕Bernoulliバンディット



アーム1の損失
~ $\text{Ber}(\mu_1)$



アーム2の損失
~ $\text{Ber}(\mu_2)$

$$\text{Ber}(\mu) = \begin{cases} 1 & (\text{確率 } \mu) \\ 0 & (\text{確率 } 1 - \mu) \end{cases}$$

- 各ラウンド $t = 1, \dots, T$ で以下を行う

- 目的パラメータの事後分布を計算する

$\pi(\mu_1 | \text{時刻 } t \text{までの観測データ}), \pi(\mu_2 | \text{時刻 } t \text{までの観測データ})$ (目的パラメータ: μ_1, μ_2)

- 事後分布から目的パラメータをサンプリングする

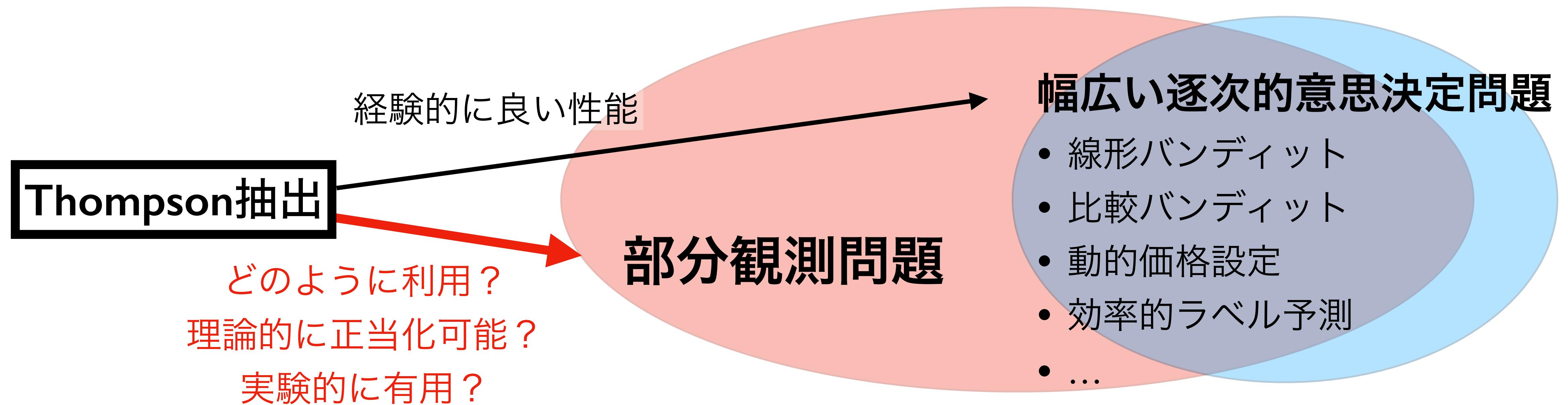
$\tilde{\mu}_1 \sim \pi(\mu_1 | \text{時刻 } t \text{までの観測データ}), \tilde{\mu}_2 \sim \pi(\mu_2 | \text{時刻 } t \text{までの観測データ})$

- サンプルしたパラメータをもとに最適な行動 (= 腕) を選択する

$\tilde{\mu}_1 \leq \tilde{\mu}_2$ ならばアーム1を引き, それ以外ならばアーム2を引く

研究課題

- 部分観測問題
 - ▶ 限られたフィードバックをもとに逐次的意思決定を行う問題の一般的な枠組み
- Thompson抽出
 - ▶ 幅広い逐次的意思決定問題に対して、経験的に最も有用な方策の1つ
 - ▶ 「探索」と「活用」のトレードオフを事後分布からのサンプリングによって扱う



部分観測問題の例: 動的価格設定

プレイヤー (= 販売者)

$t = 1$ 日目

ホテルのオーナー
(販売者)



ホテル1泊の宿泊料を決定
 $\{1000\text{円}, \dots, N\text{円}\}$

宿泊料 4,000円

$t = 2$ 日目

宿泊料 8,000円

$t = \dots$

敵対者

利用者の
内部状態 $j(t)$

(= 評価価格)



宿泊料 $\leq 9,000\text{円}$ なら利用



宿泊料 $\leq 5,000\text{円}$ なら利用



販売者は、フィードバック（宿泊する or 宿泊しない）のみ観測可能

フィードバックのみから、全体の損失を最小化する (= 全体の報酬を最大化する) ことを目指す

(機会) 損失

$$9,000 - 4,000 \\ = 5,000\text{円}$$

c 円 (定数)
($\because 5000 - 8000 < 0$)



フィード
バック

宿泊する

宿泊
しない



部分観測問題の定式化

- N 個の行動と M 個の内部状態からなる部分観測問題 $G = (L, H)$
- 損失行列 $L = (\ell_{i,j}) \in \mathbb{R}^{N \times M}$, フィードバック行列 $H = (h_{i,j}) \in \Sigma^{N \times M}$
(Σ : フィードバック記号の集合)

For round $t = 1, \dots, T$:

1. プレイヤーが行動 $i(t) \in \{1, \dots, N\}$ を選択する
2. 敵対者が内部状態を選択する $j(t) \stackrel{\text{i.i.d.}}{\sim} \text{Multi}(p^*)$ ($p^* \in \mathcal{P}_M$)
戦略 確率単体
3. プレイヤーが損失 $\ell_{i(t), j(t)}$ を被り, フィードバック $h_{i(t), j(t)}$ を観測する

- 目標: 擬リグレットの最小化 (= 全体の損失の最小化)

$$\text{Reg}(T) = \sum_{t=1}^T \left(\underline{L}_{i(t)}^\top p^* - \underline{L}_1^\top p^* \right)$$

実際に取った行動の期待損失 最適な行動の期待損失

行動1が最適であるとする

$L_i \in \mathbb{R}^M : L$ の i 列目

例 I. 動的価格設定 [Kleinberg & Leighton 2003]

N : (離散化した) 販売価格

M : (離散化した) 評価額

(*) 行: 販売価格, 列: 内部状態

行動

内部状態 $\sim p^*$ (戦略)

プレイヤー (=販売者)

$t = 1\text{日目}$

ホテルのオーナー
(販売者)

宿泊料 4,000円



ホテル1泊の宿泊料を決定
 $\{1000\text{円}, \dots, N\text{円}\}$

$t = 2\text{日目}$

宿泊料 8,000円

$t = \dots$

敵対者
内部状態 $j(t)$
(=評価価格)



宿泊料 $\leq 9,000\text{円}$ なら利用



宿泊料 $\leq 5,000\text{円}$ なら利用

(機会) 損失

$$9,000 - 4,000 = 5,000\text{円}$$

c 円 (定数)
($\because 5000 - 8000 < 0$)

フィード
バック

宿泊する

宿泊
しない

損失行列 (*)

$$\ell_{i,j} = \begin{cases} j - i & (j \geq i) \\ c & (\text{otherwise}) \end{cases}$$

					$j \geq i$
0	1	2	3	4	
c	0	1	2	3	
c	c				
c	c				
c	c	c	c	0	$j < i$

$$\Sigma = \{\text{宿泊する}(\bigcirc), \text{宿泊しない}(\times)\}$$

フィードバック行列 (*)

$$h_{i,j} = \begin{cases} \bigcirc & (j \geq i) \\ \times & (\text{otherwise}) \end{cases}$$

					$j \geq i$
\bigcirc	\bigcirc	\bigcirc	\bigcirc	\bigcirc	
\times	\bigcirc	\bigcirc	\bigcirc	\bigcirc	
\times	\times	\bigcirc	\bigcirc	\bigcirc	
\times	\times	\times	\bigcirc	\bigcirc	$j < i$
\times	\times	\times	\times	\bigcirc	

例2. 効率的ラベル予測 [Cesa-Bianchi+ 2005]

- プレイヤーが商品のラベル（正か負か）をオンラインに決定する
- 3通りの行動が存在:
 - 正とラベル付け (P)
 - 負とラベル付 (N)
 - 専門家に聞いて、正しいラベルを教えてもらう

$$L = \begin{pmatrix} 0 & c_{N \rightarrow P} \\ c_{P \rightarrow N} & 0 \\ q & q \end{pmatrix}$$

$c_{N \rightarrow P} > 0$: N を P と間違えるコスト

$c_{P \rightarrow N} > 0$: P を N と間違えるコスト

$q > 0$: 専門家に聞くコスト

$$H = \begin{pmatrix} \text{None} & \text{None} \\ \text{None} & \text{None} \\ P & N \end{pmatrix}$$

部分観測問題におけるThompson抽出の適用

未知パラメータ: 敵対者の戦略 p^*

Thompson抽出の素朴な一般化:

- I. 目的パラメータの事後分布を計算する

$$\pi(p \mid \text{時刻 } t \text{ までの観測データ}) \propto \pi(p) \prod_{i=1}^N \exp \left(-n_i \mathcal{D}_{\text{KL}} \left(q_i^{(t)} \| S_i p \right) \right)$$

2. 事後分布から目的パラメータをサンプリングする

$$\tilde{p}_t \sim \pi(p \mid \text{時刻 } t \text{ までの観測データ})$$

3. サンプルしたパラメータをもとに最適な行動を選択する

$$\text{行動 } i(t) := \arg \min_{i \in [N]} \underline{L_i^\top \tilde{p}_t} \text{ を取る}$$

行動 i の期待損失

複雑な事後分布

→サンプリングが困難 😞

n_i : 時刻 t までに行動 i を取った回数

$q_i(t)$: 時刻 t における行動 i を取ったときの経験フィードバック分布

S_i : 行動 i の信号行列

Bayes-update PM for TS (BPM-TS) [Vanchinathan+ 2014]

- 戦略 p^* の推定値をガウス分布の事前分布を用意してベイズ的更新
- 仮定: 内部状態は共分散行列 I_M と平均未知のガウス分布に従う
(実際には $\text{Multi}(p^*)$ に従う)

😊 高速な計算

😊 経験的に最も性能の良いアルゴリズムの一つ

😢 真の事後分布との乖離

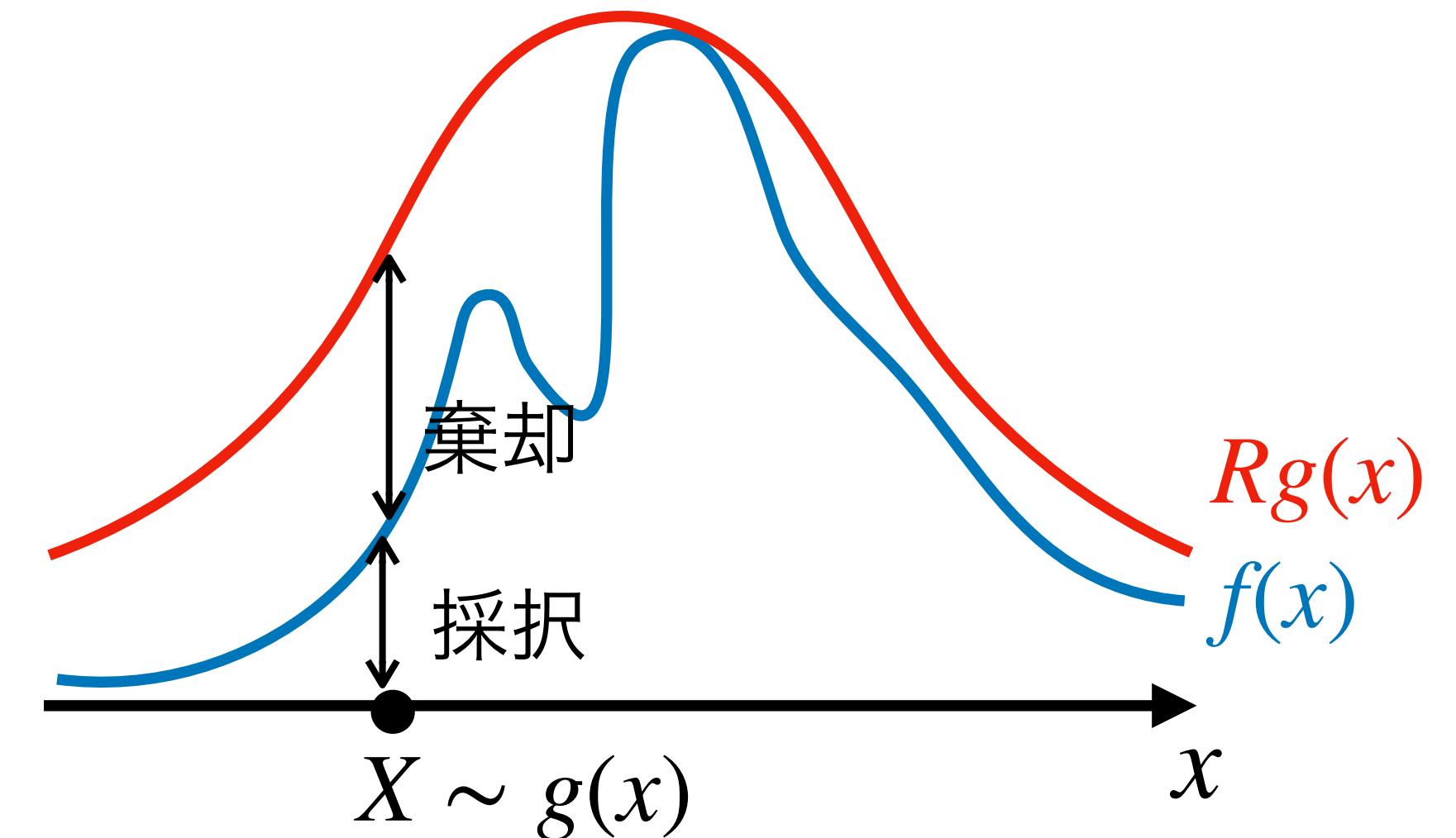
$$\mathcal{N}(\text{時刻 } t \text{ におけるパラメータ}) \longleftrightarrow_{\text{乖離}} \pi(p) \prod_{i=1}^N \exp \left(-n_i \mathcal{D}_{\text{KL}} \left(q_i^{(t)} \| S_i p \right) \right)$$

😢 理論解析が与えられていない

棄却サンプリング

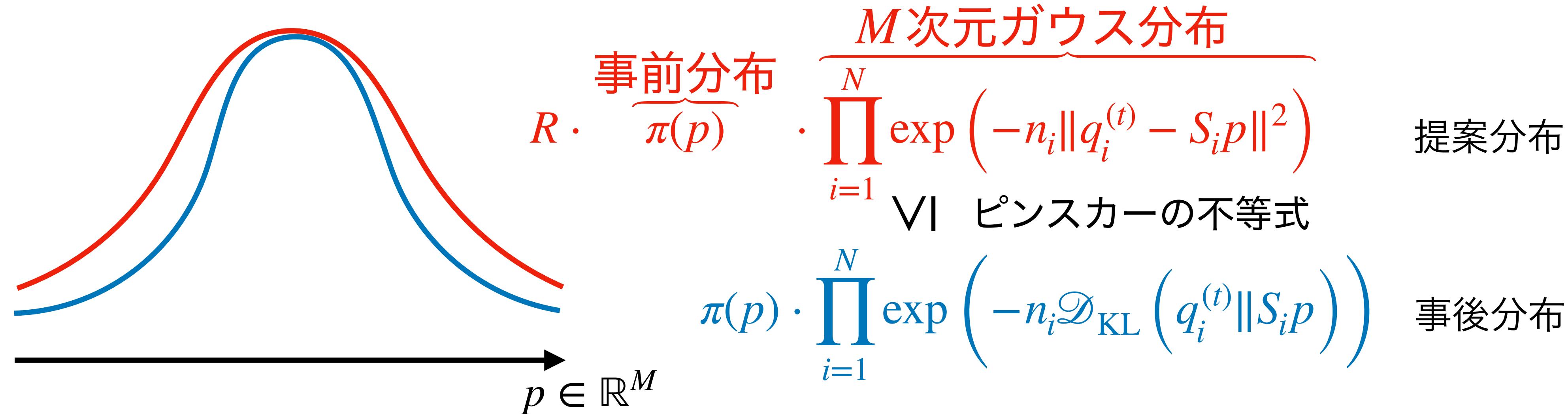
- 複雑な分布 $f(x)$ から独立同分布に従う標本を得る手法
- 提案分布 $g(x)$ を用意し、以下を行う：

容易にサンプルを得られる分布
- 1. サンプル $X \sim g(x)$ を得る
- 2. 確率 $f(X)/Rg(X)$ で X を採択する $R = \sup_x f(x)/g(x)$
- 3. 採択されるまで繰り返す
- タイトな提案分布を用意する必要がある



提案法 (TSPM)

Step 1. タイトかつサンプリングが容易な提案分布を用意する

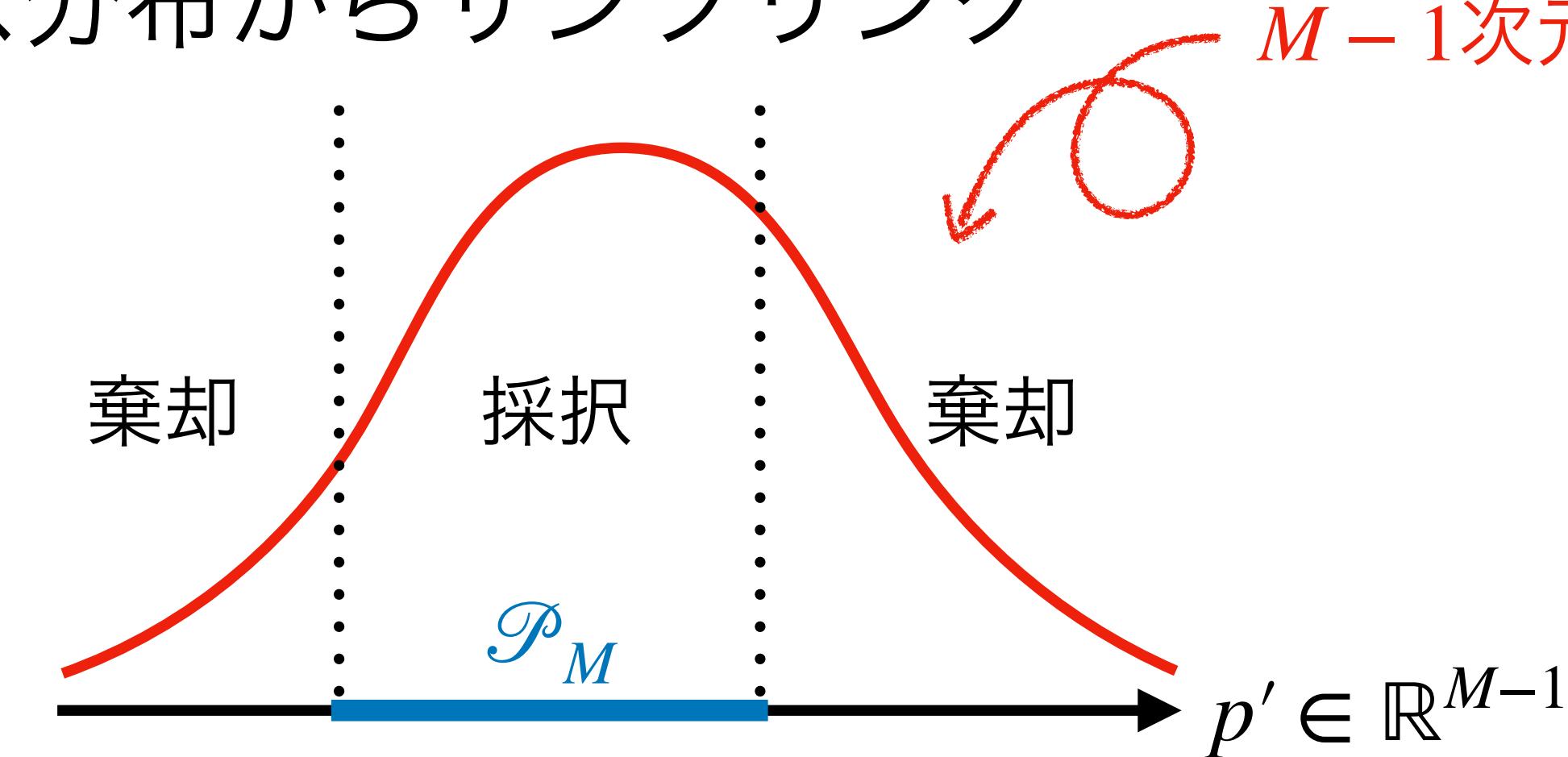


Step 2. 確率単体 \mathcal{P}_M 上に制限されたガウス分布からサンプリング

ガウス分布

$$\frac{\pi(p)}{\prod_{i=1}^N \exp(-n_i \|q_i^{(t)} - S_i p\|^2)}$$

確率単体 \mathcal{P}_M に制限



リグレット上界

定理（簡略版）.

任意の局所的観測可能な線形部分観測問題に対して、

TSPM-Gaussianの期待リグレットが以下で抑えられる:

$$\mathcal{O} \left(\max \left\{ \frac{AN \sum_{i \in [N]} \Delta_i}{\Lambda^2}, \frac{AN^3 \max_{i \in [N]} \Delta_i}{\Lambda^2} \right\} \log T \right).$$

問題依存の定数

ラウンド数 T への対数依存

- 部分観測問題に対する, Thompson抽出の初めての問題依存対数リグレット
- 線形バンディット問題に対する, Thompson抽出の初めての対数リグレット

A, N : フィードバックと行動の選択肢数

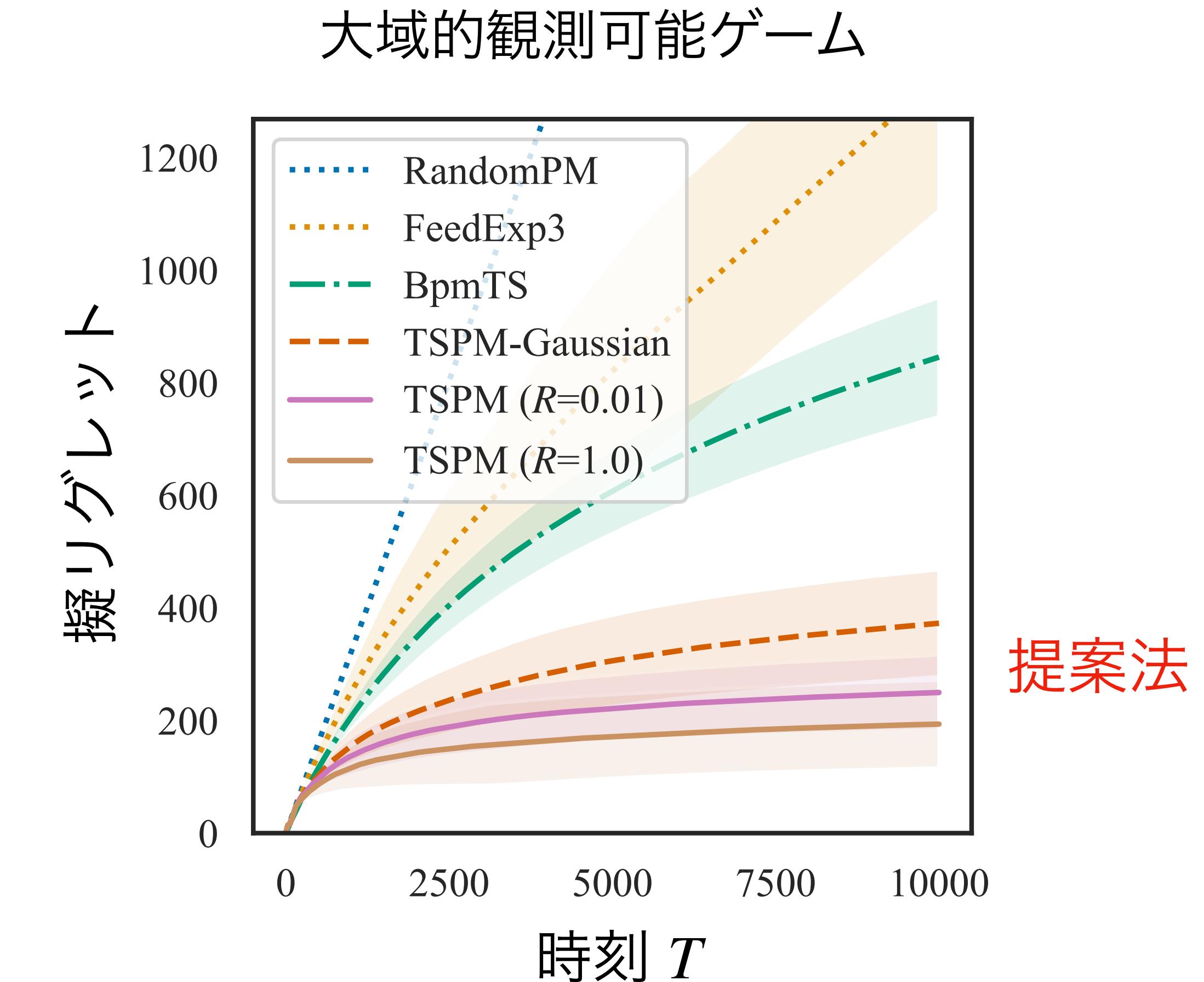
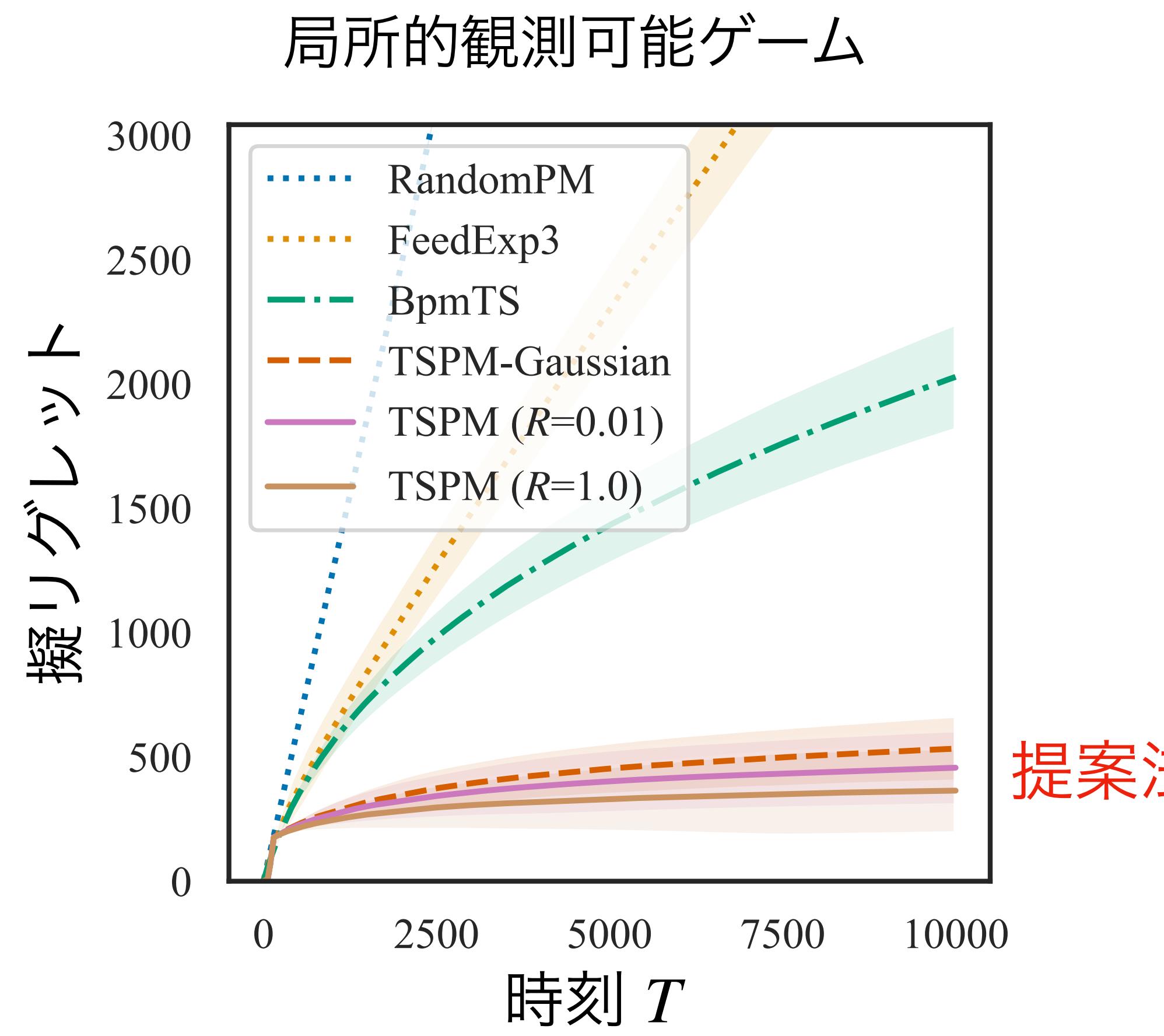
Δ_i : 行動 i の最適行動との期待損失の差

$\Lambda = \min_{j \neq k} \Delta_{j,k} / \|z_{j,k}\|$ ($\Delta_{j,k}$: 行動 j と k の期待損失の差,

$z_{j,k} \in \mathbb{R}^{2A}$: 損失とフィードバックを関連させるベクトル)

動的価格設定におけるリグレット比較実験

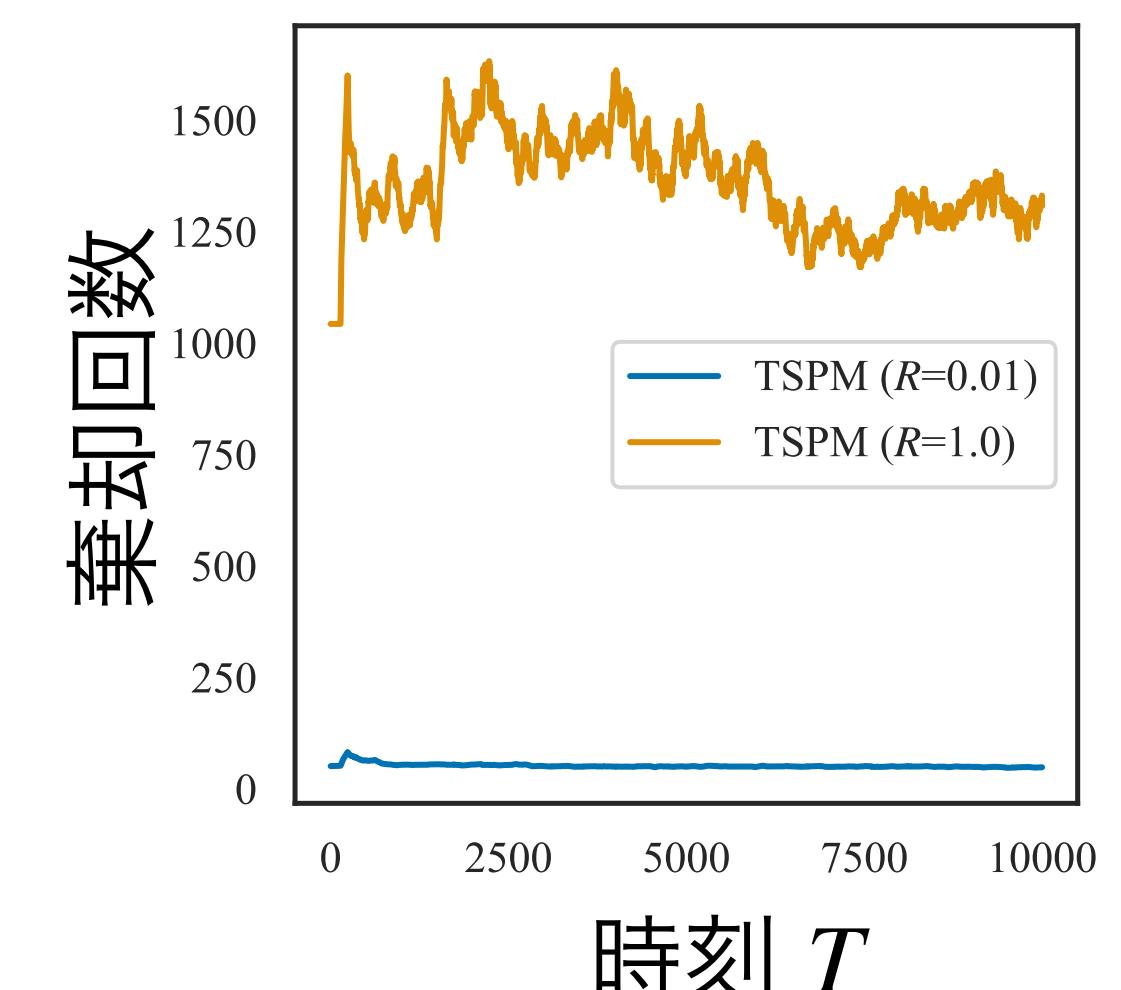
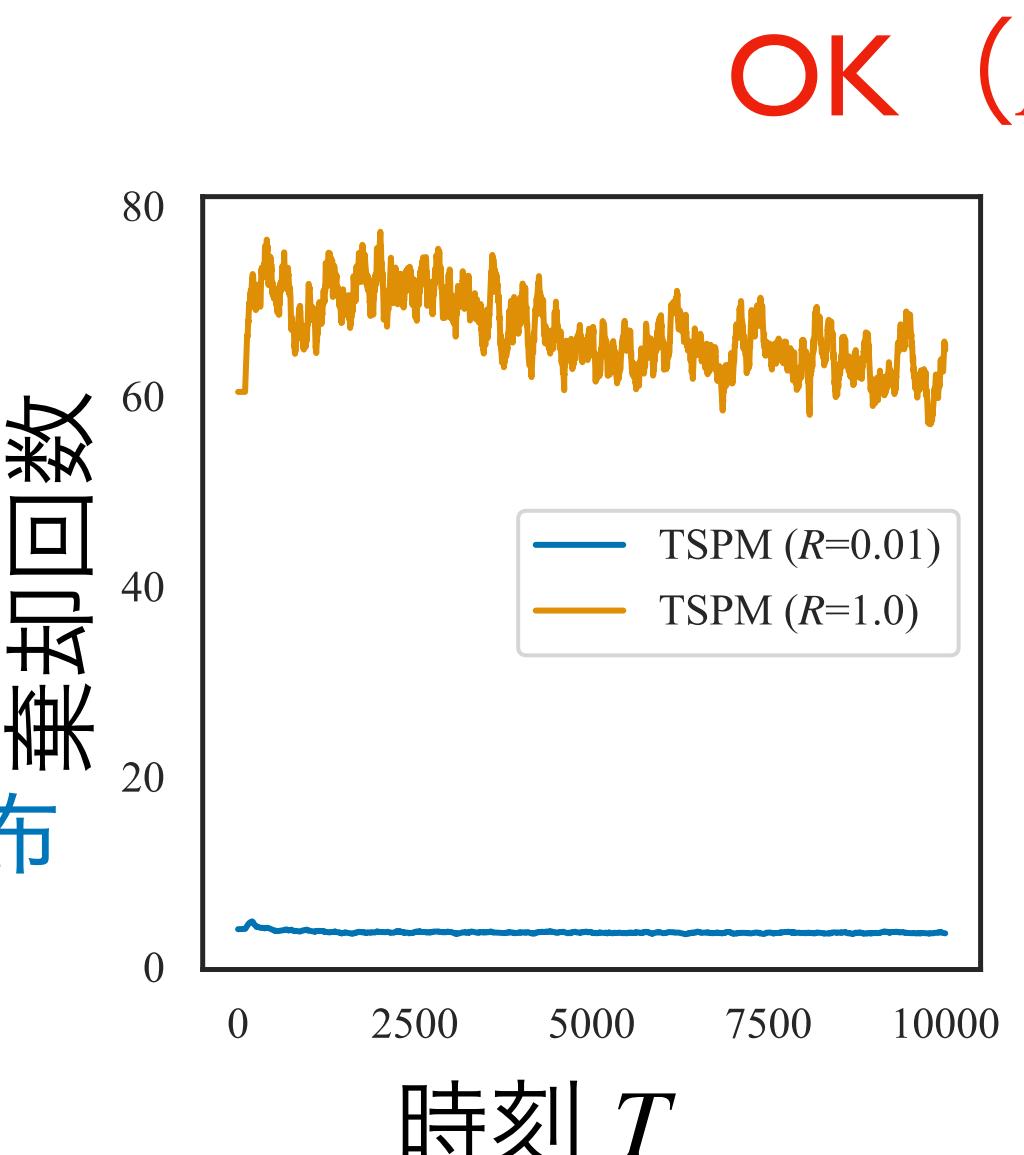
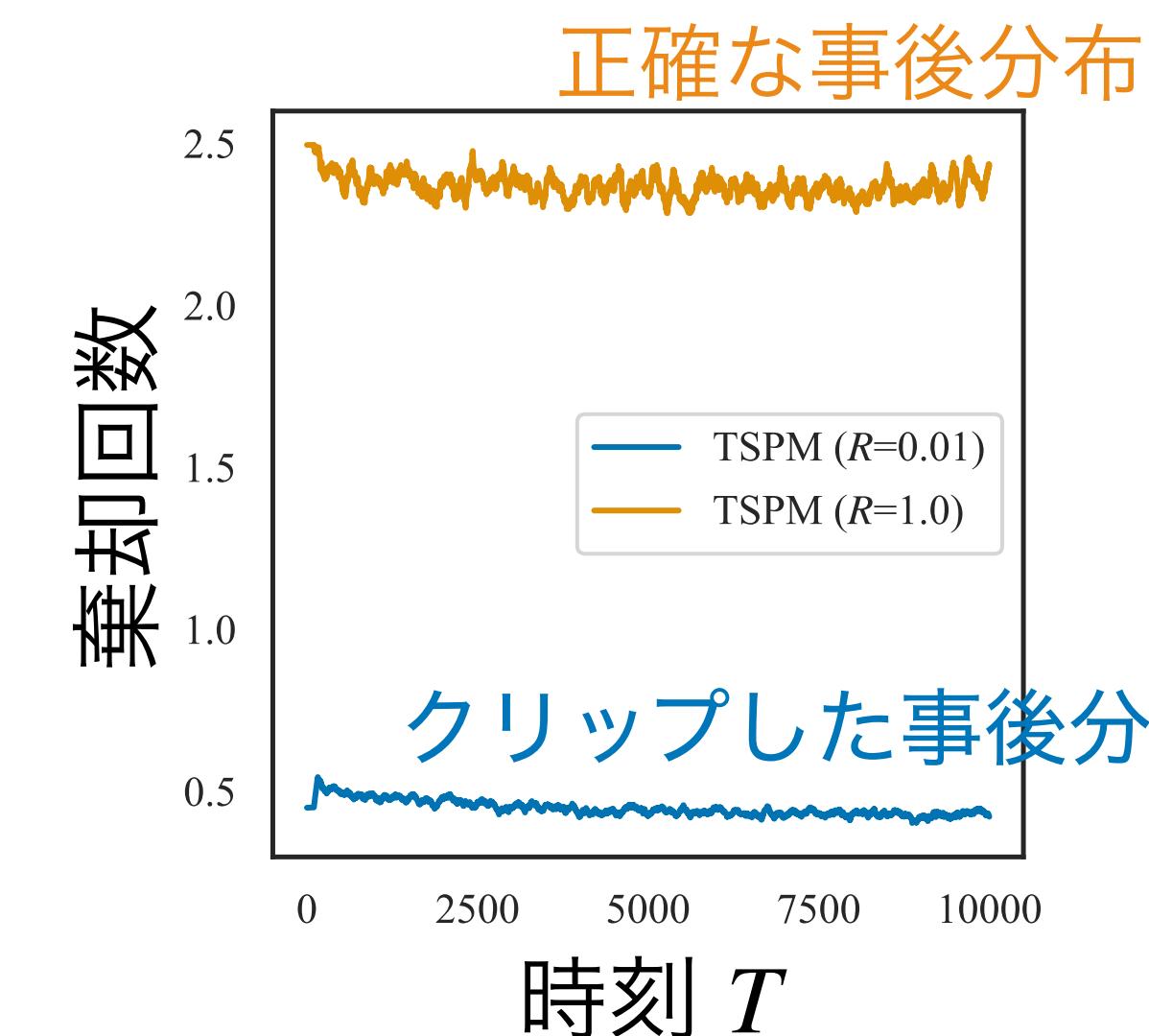
- 既存手法の性能を大きく改善



棄却サンプリングにおける棄却頻度

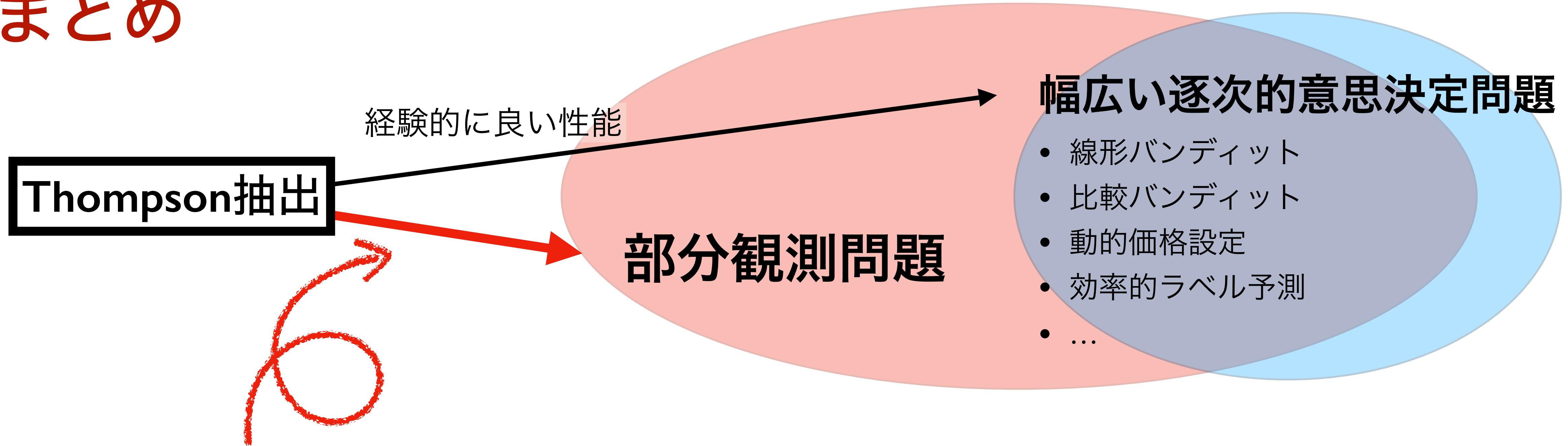
- 棄却サンプリングにおいて望ましい性質
 - I. ラウンド数が進むにつれ、棄却頻度が増えない **OK**
 2. 分布のサポートの次元 $M - 1$ が増えるにつれ、棄却頻度が増えない

- 局所的観測可能ゲーム
- 動的価格設定



定数 · $\frac{\text{事後分布の密度関数}}{R \cdot \text{提案分布の密度関数}}$ の確率で採択

まとめ



I. どのように利用？

タイトな提案分布からのサンプリングによる新しいThompson抽出アルゴリズム

2. 理論的に正当化可能？

大域的観測可能性を満たす連續緩和した問題に対して可能

3. 実験的に有用？

Yes! 理論的性質を示していない大域的観測可能ゲームに対しても有用