



Title	bandit問題の拡張について
Author(s)	中村, 篤祥
Citation	2010年度科学技術振興機構ERATO湊離散構造処理系プロジェクト講究録. p.295-297.
Issue Date	2011-06
Doc URL	http://hdl.handle.net/2115/48414
Type	conference presentation
Note	ERATO 湊離散構造処理系プロジェクト春のワークショップ (キックオフシンポジウム). 2010年5月28日 (金) ~ 29日 (土). ERATO湊プロジェクト研究室.
File Information	04.nakamura_06.pdf



[Instructions for use](#)

北海道大学 Hokkaido University

bandit 問題の拡張について

北海道大学情報科学研究科
中村篤祥

2010/05/08 ERATO湊離散構造処理系プロジェクトキックオフシンポジウム

北海道大学 Hokkaido University

multi-armed bandit 問題とは

スロットマシンID 1 2 ... K

時刻tにおける報酬(reward) $x_1(t)$ $x_2(t)$... $x_K(t)$ ← playerは知らない (選んだスロットマシンについてのみ知ることができる)

各時刻t(=1,2,...)においてplayerは以下のことを行う。

1. K台のスロットマシンから1台のスロットマシン i_t を選ぶ。
2. 選ばれたスロットマシン i_t から報酬 $X_{i_t}(t)$ を得る。

総利得 $\sum x_{i_t}(t)$ を最大化するためには各時刻においてどのようにスロットマシンを選べばよいか？

2010/05/08 ERATO湊離散構造処理系プロジェクトキックオフシンポジウム

北海道大学 Hokkaido University

確率的なmulti-armed bandit 問題

スロットマシンID 1 2 ... K

成功確率 θ_1 θ_2 ... θ_k

報酬 $x_1(t)$ $x_2(t)$... $x_k(t)$

← playerは知らない 時刻によらず一定

$$x_i(t) = \begin{cases} 1 & \text{if success} \\ 0 & \text{if fail} \end{cases}$$

- $x_1(t), x_2(t), \dots, x_k(t)$ は独立
- $x_i(1), x_i(2), \dots$ は未知の成功確率 θ_i のBernoulli process

目標1 expected total discounted reward $\sum_{t=1}^{\infty} \gamma^{t-1} E(x_{i_t}(t))$ の最大化 ($0 < \gamma < 1$)

目標2 与えられた時刻Tまでの総利得 $\sum_{t=1}^T x_{i_t}(t)$ の最大化

2010/05/08 ERATO湊離散構造処理系プロジェクトキックオフシンポジウム

北海道大学 Hokkaido University

確率的なmulti-armed bandit問題の戦略

- **Gittins Indexが最大のスロットマシンを選ぶ**
成功、失敗回数が α, β であるようなスロットの **Gittins Index $G(\alpha, \beta)$**
以下のような確率p: 成功確率がpであるとわかっているスロットマシン1と、(成功回数, 失敗回数)=(α, β)であるような成功確率未知のスロットマシン2があったとき、時刻t=1どちらのスロットマシンを選んで、optimal expected total discounted rewardが同じになるようなp
expected total discounted rewardを最大化する戦略 [Gittins and Jones 1974]
- **ϵ -greedy** [Sutton & Barto 1998]
1- ϵ の確率でそれまでの平均報酬が最大のマシンを選び、 ϵ の確率でランダムに選ぶ。
囲碁プログラムのモンテカルロ木探索に用いられて著名になった。
- **UCB1** [Auer, Cesa-Bianchi and Fisher 2002]
時刻tに $\bar{x}_j + \sqrt{\frac{2 \ln t}{n_j}}$ が最大のマシンjを選ぶ。
 \bar{x}_j : マシンjのそれまでの平均報酬
 n_j : マシンjをそれまでに選んだ回数

2010/05/08 ERATO湊離散構造処理系プロジェクトキックオフシンポジウム

北海道大学 Hokkaido University

adversarial bandit 問題 [Auer et al. 2002]

スロットマシンID 1 2 ... K

時刻tにおける報酬(reward) $x_1(t)$ $x_2(t)$... $x_k(t)$ ← 悪魔が選ぶ。 playerは知らない

各時刻t(=1,2,...,T)においてplayerは以下のことを行う。

1. K台のスロットマシンから1台のスロットマシン i_t を選ぶ。
2. 選ばれたスロットマシン i_t から報酬 $X_{i_t}(t)$ を得る。

$G_A(T) = \sum_{t=1}^T x_{i_t}(t)$: (乱択)アルゴリズムの時刻Tにおける総利得

期待総利得 $E(G_A(T))$ が大きなアルゴリズムAは？

2010/05/08 ERATO湊離散構造処理系プロジェクトキックオフシンポジウム

北海道大学 Hokkaido University

乱択アルゴリズムExp3 [Auer et al. 2002]

Algorithm Exp3 // exponential-weight algorithm
Parameter: $\gamma \in (0, 1]$ for exploration and exploitation
Initialization: $w_i(1) \leftarrow 1$ for $i=1, 2, \dots, K$

for t=1 to T

1. $p_i(t) \leftarrow (1-\gamma) \frac{w_i(t)}{\sum_{j=1}^K w_j(t)} + \frac{\gamma}{K}$ for $i=1, \dots, K$
2. i_t を $p_1(t), \dots, p_K(t)$ の分布に従ってランダムに選ぶ
3. 報酬 $x_{i_t}(t) \in [0, 1]$ を得る
4. for $j=1, \dots, k$

$$\hat{x}_j(t) \leftarrow \begin{cases} x_j(t)/p_j(t) & \text{if } j = i_t \\ 0 & \text{otherwise} \end{cases}$$

$$w_j(t+1) \leftarrow w_j(t) \exp(\gamma \hat{x}_j(t)/K)$$

2010/05/08 ERATO湊離散構造処理系プロジェクトキックオフシンポジウム

7 北海道大学 Hokkaido University
Exp3のexpected (weak) regret [Auer et al. 2002]

1つのスロットマシンを選び続けた場合の総利得の最大値 $G_{\max}(T) = \max_{i=1}^K \sum_{t=1}^T x_i(t)$
 と比べて平均的にどのくらい後悔するかを $G_{\max}(T) - E[G_{\text{Exp3}}(T)]$ で測ると

$$G_{\max}(T) - E[G_{\text{Exp3}}(T)] \leq 2.63\sqrt{TK \ln K} \quad \left(\gamma = \min\left\{1, \sqrt{\frac{K \ln K}{(e-1)T}}\right\} \text{のとき} \right)$$

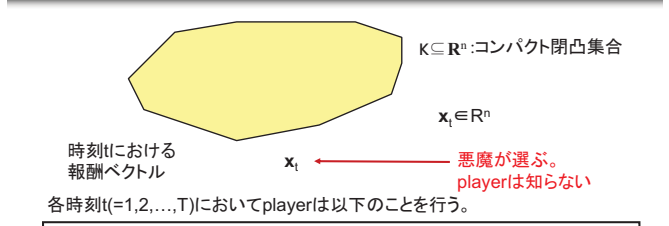
($x_i(t) \in [0, 1]$ なので、 $G_{\max}(T) - G_{\text{Exp3}}(T) \leq T$ が成り立つことに注意。)

また、ある報酬割合分布が存在し、どのような乱択アルゴリズム A に対しても

$$E[G_{\max}(T) - G_A(T)] \geq \frac{1}{20} \min\{\sqrt{TK}, T\}$$

が成り立つ。ただし、期待値は報酬割合とアルゴリズムの乱択の両方に関してとるものとする。

8 北海道大学 Hokkaido University
online linear optimization 問題 [McMahan and Blum 2004]



各時刻 $t(=1, 2, \dots, T)$ において player は以下を行う。

1. K から 1 点 q_t を選ぶ
2. 選ばれた点 q_t に対する利得 $x_t \cdot q_t$ を得る。

$G_A(T) = \sum_{t=1}^T x_t \cdot q_t$: (乱択)アルゴリズムの時刻 T における総利得
期待総利得 $E(G_A(T))$ が大きなアルゴリズム A は?

9 北海道大学 Hokkaido University
Abernethyらのアルゴリズム [Abernethy et al. 2008]

Algorithm BOLO (仮称) // Bandit Online Linear Optimization
 Input: $\eta > 0, \theta$ -self-concordant R
 Initialization: $q_1 \leftarrow \arg \min_{q \in K} R(q)$
 for $t=1$ to T

1. $\{e_1, e_2, \dots, e_n\} \leftarrow \nabla^2 R(q_t)$ の固有ベクトル
 $\{\lambda_1, \lambda_2, \dots, \lambda_n\} \leftarrow \nabla^2 R(q_t)$ の固有値
2. i_t を $\{1, 2, \dots, n\}$ からランダムに選ぶ。
 ϵ_t を $\{-1, 1\}$ からランダムに選ぶ。
3. $r_t \leftarrow q_t + \epsilon_t \lambda_{i_t}^{-1/2} e_{i_t}$
4. 報酬 $x_t \cdot r_t$ を得る。
5. 次のように更新する。

$$\hat{x}_t \leftarrow n(x_t \cdot r_t) \epsilon_t \lambda_{i_t}^{1/2} e_{i_t}$$

$$q_{t+1} \leftarrow \arg \min_{q \in K} \sum_{s=1}^t \hat{x}_s \cdot q + R(q)$$

10 北海道大学 Hokkaido University
BOLOのexpected regretの上界

$$\max_{q \in K} E\left(\sum_{t=1}^T x_t \cdot q\right) - E(G_{\text{BOLO}}(T)) = O\left(n\sqrt{\theta \ln T}\right)$$

$$\left(\eta = \frac{\sqrt{\theta \ln T}}{4n\sqrt{T}}, T \geq 8\theta \ln T \text{ のとき} \right)$$

11 北海道大学 Hokkaido University
インターネット広告はbandit問題



12 北海道大学 Hokkaido University
multiple play設定



Algorithm FExp3
 Parameter: $\gamma \in (0, 1]$
 Initialization: $w_i(1) \leftarrow 1$ for $i=1, 2, \dots, K$

- for $t=1$ to T
- if $\operatorname{argmax}_{j \in \{1, 2, \dots, K\}} w_j(t) \geq \left(\frac{1-\gamma}{k} - \frac{\gamma}{K}\right) \sum_{i=1}^K w_i(t) / (1-\gamma)$ then
 $\frac{\alpha_i}{\sum_{w_i(t) \geq \alpha_i} \alpha_i + \sum_{w_i(t) < \alpha_i} w_i(t)} = \left(\frac{1-\gamma}{k} - \frac{\gamma}{K}\right) / (1-\gamma)$ を満たす α_i をもとめる。
 $S_0(t) \leftarrow \{i : w_i(t) \geq \alpha_i\}$, $w'_i(t) \leftarrow \alpha_i$ for $i \in S_0(t)$
 else
 $S_0 \leftarrow \emptyset$
 - $w'_i(t) \leftarrow w_i(t)$ for $i \in \{1, 2, \dots, K\} - S_0(t)$
 - $p_i(t) \leftarrow k \left((1-\gamma) \frac{w'_i(t)}{\sum_{j=1}^K w'_j(t)} + \frac{\gamma}{K} \right)$ for $i = 1, \dots, K$

- $S(t) \leftarrow i$ が選ばれる確率が $p_i(t)$ になるように $\{1, 2, \dots, K\}$ から k 個選択
- 報酬 $x_i(t) \in [0, 1]$ for $i \in S(t)$ を得る
- for $j=1, \dots, k$

$$\tilde{x}_j(t) \leftarrow \begin{cases} x_j(t)/p_j(t) & \text{if } j \in S(t) \\ 0 & \text{otherwise} \end{cases}$$

$$w_j(t+1) \leftarrow \begin{cases} w_j(t) \exp(k \gamma \tilde{x}_j(t)/K) & \text{if } j \notin S_0(t) \\ w_j(t) & \text{otherwise} \end{cases}$$

同じ k 個のスロットマシンを選び続けた場合の総利得の最大値

$$G_{\max-k}(T) = \max_{S \subseteq \{1, 2, \dots, K\}, |S|=k} \sum_{t=1}^T \sum_{i \in S} x_i(t) \text{ と比べると}$$

$$G_{\max-k}(T) - E[G_{\text{FExp3}}(T)] \leq 2.63 \sqrt{kTK \ln \frac{K}{k}} \left(\gamma = \min \left\{ 1, \sqrt{\frac{k \ln(K/k)}{(e-1)kT}} \right\} \text{ のとき} \right)$$

また、ある報酬割当分布が存在し、どのような乱択アルゴリズム A に対しても

$$E[G_{\max-k}(T) - G_A(T)] \geq \min \left\{ \frac{1}{5} \left(\frac{K-k}{K} \right)^2 \sqrt{TK}, \frac{K-k}{8K} kT \right\}$$

が成り立つ。ただし、期待値は報酬割当とアルゴリズムの乱択の両方に関してとるものとする。

- 協調フィルタリングのbanditモデル
 ユーザの好みはいくつかの好みのタイプの混合として表現できると仮定して、最適な混合割合と比べた場合、リグレットが少ない方法を考案する。
- 情報検索のbanditモデル
 1回に k 個ずつ結果が提示されるとして T 回までみるとした場合に最も多く欲しいものが含まれている結果提示方法は？ 適合度フィードバックしながら適度に良い結果を返す。低次元の学習問題にどのように落とすか。