



バンディットの理論と応用

北海道大学情報科学研究科

中村篤祥



目次

1. バンディット問題とは
2. stableな確率的バンディット 問題の解法
 - (1) infinite horizon (geometric discount)の最適手法
 - (2) finite horizonの手法
3. Adversarial バンディット問題の解法
4. バンディットの拡張
 - (1) multiple play
 - (2) combinatorial bandit
 - (3) online linear optimization
 - (4) Gaussian Process Bandit
5. バンディット問題の応用



multi-armed bandit 問題とは [Robbins 1952]

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 $t (= 1, 2, \dots)$ 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。

$\gamma_1, \gamma_2, \dots$: discount sequence ($0 \leq \gamma_t \leq 1$)

目標: expected total discounted reward $\sum_{t=1}^{\infty} \gamma_t E(x_{i_t}(t))$ の最大化



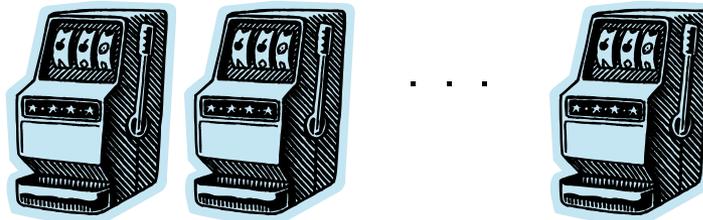
最も単純なmulti-armed bandit 問題

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Bernulli process

主な目標

試行回数	discount sequence	最大化
無限(infinite horizon)	geometric discount $\gamma_t = \gamma^{t-1} \quad (0 < \gamma < 1)$	$\sum_{t=1}^{\infty} \gamma^{t-1} E(x_{i_t}(t))$
T回(T-horizon)	T-horizontal uniform discount $\gamma_1 = \gamma_2 = \dots = \gamma_T = 1, \gamma_{T+1} = \gamma_{T+2} = \dots = 0$	$\sum_{t=1}^T E(x_{i_t}(t))$



知識の獲得と利用のトレードオフ

知識の獲得と利用のトレードオフ (exploration-exploitation trade-off)

選択回数の少なスロット
を選んで成功確率 θ_i の
推定値 $\hat{\theta}_i$ の信頼度を
上げる

知識の獲得(exploration)

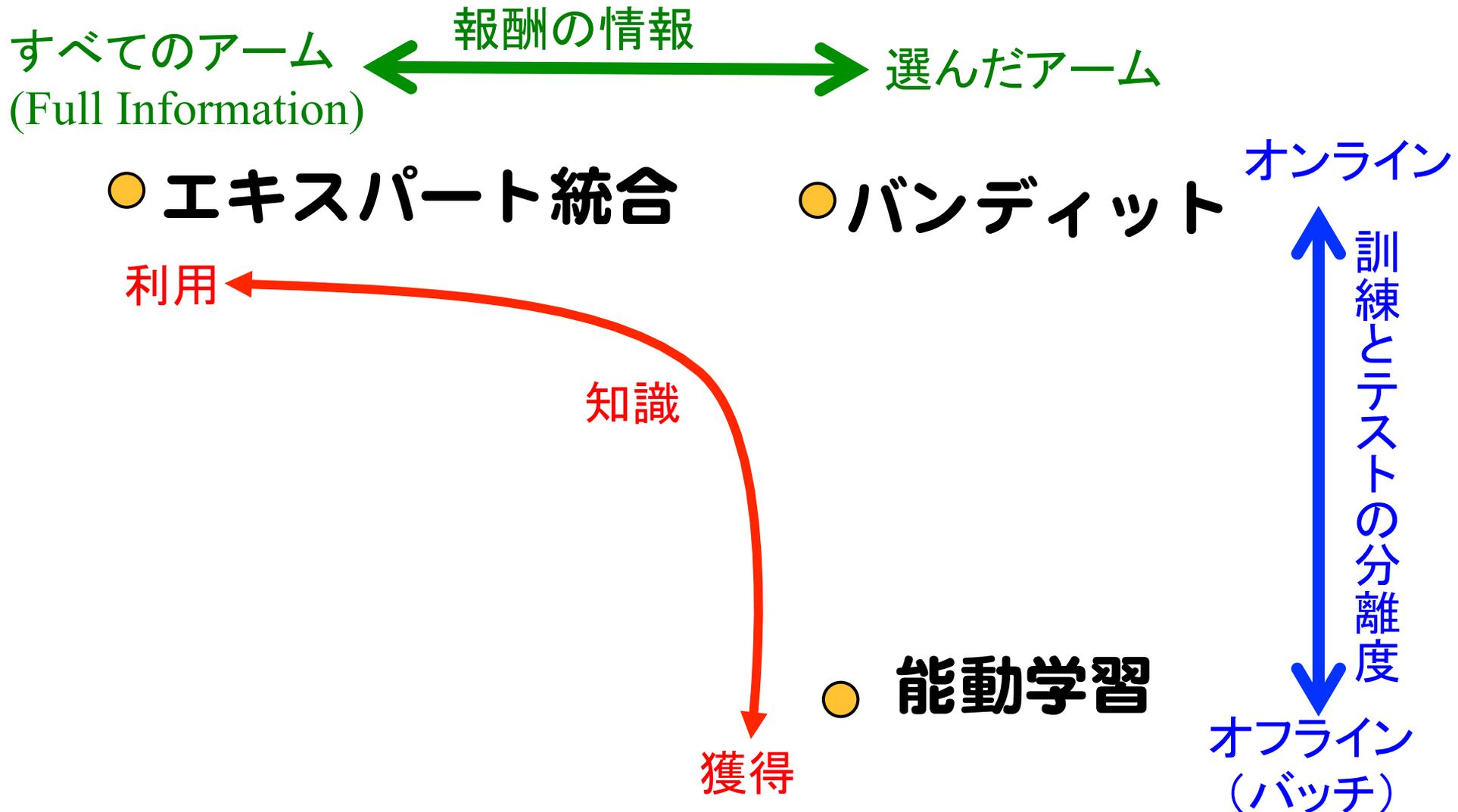


成功確率の推定値 $\hat{\theta}_i$
が最大のスロット i を
選択する

知識の利用(exploitation)



バンディット問題の位置づけ





バンディット問題の例

- クリニックトライアル
- ネットワークにおけるルート選択
- リコメンデーション・広告配信
- 商品価格の設定
- バンディット版オンライン最短路問題



目次

1. バンディット問題とは
2. **stableな確率的バンディット 問題の解法**
 - (1) **infinite horizon (geometric discount)の最適手法**
 - (2) **finite horizonの手法**
3. Adversarial バンディット問題の解法
4. バンディットの拡張
 - (1) multiple play
 - (2) combinatorial bandit
 - (3) online linear optimization
 - (4) Gaussian Process Bandit
5. バンディット問題の応用



stableなstochastic multi-armed bandit 問題

スロットマシンID

1

2

K



...



報酬分布

$$F_1(x|\theta_1) \quad F_2(x|\theta_2)$$

$$F_K(x|\theta_K)$$

}

}

}

報酬

$$x_1(t)$$

$$x_2(t)$$

$$x_K(t)$$

$$x_i(s), x_j(t) \quad (i \neq j, s, t = 1, 2, \dots)$$

は独立

期待報酬

$$\mu_1$$

$$\mu_2$$

$$\mu_K$$

(神様の)最適戦略: ずっと $i^* = \arg \max_i \mu_i$ を選び続ける

目標

(1) ベイズ最適な戦略

(2) 期待リグレット(神様の最適戦略による報酬との差)最小化

T-horizonの場合

$$\mu^* T - E\left(\sum_{t=1}^T x_{i_t}(t)\right) = \mu^* T - \sum_{i=1}^K \mu_i n_i$$

ただし

$\mu^* = \max_i \mu_i$, n_i : $i_t = i$ の回数



One-armed bandit 問題

スロットマシンID

1



2



報酬分布

$$F_1(x|\theta_1)$$

$$F_2(x|\theta_2)$$

$\theta_1 \sim G(\theta)$:事前分布

報酬

$$x_1(t)$$

$$x_2(t)$$

期待報酬

unknown

p (known)

最適戦略: $i_t=1$ の最適戦略と $i_t=2$ の最適戦略で
discounted total reward が大きい方を選ぶ

$$\Leftrightarrow i_t = \begin{cases} 1 & \text{if } \Lambda(t) > p \\ 2 & \text{otherwise} \end{cases}$$

ただし

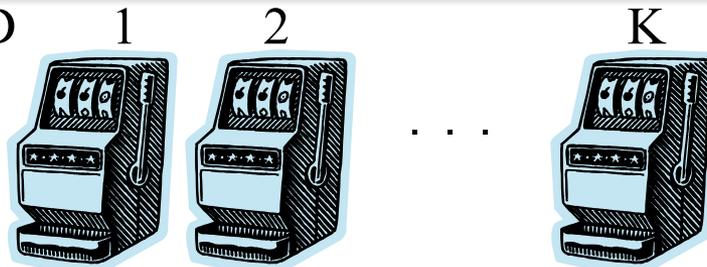
$$\Lambda(t) = \sup_N E\left(\sum_{s=1}^N \gamma^{s-1} x_1(t-1+s)\right) / \frac{1-\gamma^N}{1-\gamma}$$

: マシン1のGittins Index



K-armed bandit 問題の Gittins Index

スロットマシンID



$(\theta_1, \theta_2, \dots, \theta_K) \sim$
 $G(\theta_1, \theta_2, \dots, \theta_K)$: 事前分布

報酬分布	$F_1(x \theta_1)$	$F_2(x \theta_2)$...	$F_K(x \theta_K)$
	\downarrow	\downarrow		\downarrow
報酬	$x_1(t)$	$x_2(t)$		$x_K(t)$

最適戦略: $i_t = \arg \max_i \Lambda_i(t)$

ただし $\Lambda_i(t) = \sup_N E\left(\sum_{s=1}^N \gamma^{s-1} x_i(t-1+s)\right) / \frac{1-\gamma^N}{1-\gamma}$: マシン*i*の Gittins Index

定理 [Gittins and Jones 1974] 常に Gittins Index が最大の スロットマシンを選ぶ戦略は、expected total discounted reward $\sum_{t=1}^{\infty} \gamma^{t-1} E(x_{i_t}(t))$ を最大化する(ベイズ最適である)。



最も単純なK-armed bandit問題のGittins Index(1/2)

スロットマシンID

1

2

K



...



成功確率
報酬

θ_1

θ_2

$x_1(t)$

$x_2(t)$

θ_k

$x_K(t)$

playerは知らない
時刻によらず一定

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

成功回数が α 、失敗回数が β であるようなスロットの**Gittins Index $G(\alpha, \beta)$**

(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1と、
成功確率が p であるとわかっているスロットマシン2があったとき、
 その時刻にどちらのスロットマシンを選んでも、
 optimal expected total discounted rewardが同じになるような p



最も単純なK-armed bandit問題のGittins Index(2/2)

$R(\alpha, \beta, p)$: (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1と成功確率が p であるとわかっているスロットマシン2があるときの optimal expected total discounted reward

成功確率 θ_1 はパラメータ (α, β) のベータ分布従うとすれば

$$\frac{p}{1-\gamma} = \frac{\alpha+1}{\alpha+\beta+2} + \gamma \left(\frac{\alpha+1}{\alpha+\beta+2} R(\alpha+1, \beta, p) + \frac{\beta+1}{\alpha+\beta+2} R(\alpha, \beta+1, p) \right) \dots\dots\dots \textcircled{1}$$

①は十分大きな α, β に対する $R(\alpha, \beta, p)$ を $\frac{\alpha+1}{(1-\gamma)(\alpha+\beta+2)}$ で近似し、

その近似を使って動的計画法を行うことにより近似計算可能



T-horizon の場合の主な戦略

- Upper Confidence Indexによる方法 [Lai & Robbins 1985, Aggrawal 1995]

漸近的にベイズ最適な方法

$T \rightarrow \infty$

$E(n_i) \rightarrow \log T / D(X_i \parallel X_{i^*})$ where $n_i : i_t = i$ の回数,

KL情報量

$X_i \sim F_i(x|\theta_i)$, $i^* = \operatorname{argmax}_i \mu_i$

- ϵ -greedy [Sutton & Barto 1998]

$1 - \epsilon$ 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。

$\Rightarrow T \rightarrow \infty$ としてもリグレットは線形で増える

- **UCB1** [Auer, Cesa-Bianchi & 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

リグレット上界
$$\left[8 \sum_{i: \mu_i < \mu^*} \frac{1}{\mu^* - \mu_i} \right] \log T + \left(1 + \frac{\pi^2}{3} \right) \sum_{i=1}^K (\mu^* - \mu_i)$$



目次

1. バンディット問題とは
2. stableな確率的バンディット 問題の解法
 - (1) infinite horizon (geometric discount)の最適手法
 - (2) finite horizonの手法
3. Adversarial バンディット問題の解法
4. バンディットの拡張
 - (1) multiple play
 - (2) combinatorial bandit
 - (3) online linear optimization
 - (4) Gaussian Process Bandit
5. バンディット問題の応用



adversarial bandit 問題 [Auer et al. 2002]

スロットマシンID

1

2

K



...



$x_i(t) \in [0, 1]$ for all $i=1, 2, \dots, K$

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 $t(=1, 2, \dots, T)$ 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

期待リグレット $\max_i \sum_{t=1}^T x_i(t) - E(G_A(T))$ を小さくするアルゴリズム A は？



乱択アルゴリズムExp3 [Auer et al. 2002]

Algorithm Exp3

Parameter: $\gamma \in (0, 1]$

Initialization: $w_i(1) \leftarrow 1$ for $i=1, 2, \dots, K$

// **ex**ponential-weight algorithm

for **ex**ploration and **ex**ploitation

for $t=1$ to T

$$1. p_i(t) \leftarrow (1-\gamma) \frac{w_i(t)}{\sum_{j=1}^k w_j(t)} + \gamma \frac{1}{K} \quad \text{for } i=1, \dots, K$$

知識の利用

知識の獲得

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 $j=1, \dots, k$

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

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

乱択版Hedge

(full information設定)

との違い

1. 選択確率に一様分布を加える
2. 重みの更新に推定報酬を使う

不偏推定量になる!

$$E(\hat{x}_i(t)) = x_i(t)$$



Exp3の期待リグレット [Auer et al. 2002]

$G_{\max}(T) = \max_j \sum_{t=1}^T x_j(t)$: 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

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

(注意: $x_j(t) \in [0, 1]$ より、 $G_{\max}(T) - G_{\text{Exp3}}(T) \leq T$ が成立)

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

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

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。

	stable stochastic	adversarial
期待リグレット	$O(\log T)$	$\Theta(\sqrt{T})$



目次

1. バンディット問題とは
2. stableな確率的バンディット 問題の解法
 - (1) infinite horizon (geometric discount)の最適手法
 - (2) finite horizonの手法
3. Adversarial バンディット問題の解法
4. バンディットの拡張
 - (1) multiple play
 - (2) combinatorial bandit
 - (3) online linear optimization
 - (4) Gaussian Process Bandit
5. バンディット問題の応用



Exp3のmultiple play版

[Uchiya, Nakamura & Kudo 2010, Kale, Reyzin & Shapire 2010],

- 各アーム*i*の選択確率が p_i になるように*k*本選択
 - メリット: 報酬 $x_j(t)$ の不偏推定量となる推定式がそのまま使える

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

- 次の条件を満たさなければならない

$$0 \leq p_j \leq 1, \sum_{j=1}^K p_j = k$$

Exp3のアーム選択確率 $\{p_j\}$ に対しては $\sum_j p_j = 1$

→各アーム*i*を kp_j で選択するように変更

→ $p_j \leq 1/k$ を満たす必要がある

→ p_j は重みの割合 $w_i / \sum_j w_j$ で決まるので

capping法[Warmuth & Kuzmin 2008]を用いて大きな重みをカット

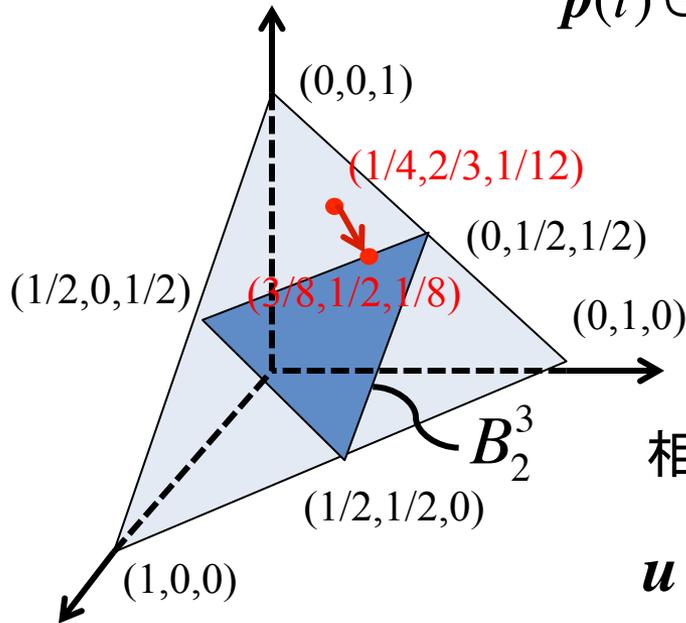
- $O(K)$ で効率よく選択できる方法が存在する

Dependent Rounding 法[Gandhi et al. 2006]



Capping法 [Warmuth & Kuzmin 08]

$K=3, k=2$ のとき



$p(t) \in B_k^K = \left\{ y \in [0, 1/k]^K : \sum_{i=1}^K y_i = 1 \right\}$ を満たさなければならない

$p(t)$ が凸集合 B_k^K の外に出てしまった場合、 $p(t)$ を B_k^K に射影したベクトルを $p(t)$ として用いる。

相対エントロピー $d(u, v) = \sum_{i=1}^K u_i \ln \frac{u_i}{v_i}$ を用いた場合

$u \notin B_k^K$ は $\operatorname{argmin}_{v \in B_k^K} d(u, v)$ へ射影される。

バンディットの場合は重みのベクトル $(w_1/\sum_j w_j, \dots, w_K/\sum_j w_j)$ を

$$B_{k,\gamma}^K = \left\{ y \in \left[0, \frac{1/k - \gamma/K}{1-\gamma} \right]^K : \sum_{i=1}^K y_i = 1 \right\} \text{ へ射影}$$



multiple play 設定のリグレット上界

BCH: Bandit版Capped Hedge Algorithm

Exp3.M: Bandit版Modified Capped Hedge Algorithm

[Uchiya, Nakamura & Kudo 10]

$$G_{\max-k} - E[G_{BCH}], G_{\max-k} - E[G_{Exp3.M}] \leq 2\sqrt{(e-1)TK \ln \frac{K}{k}}$$

ただし $G_{\max-k} \stackrel{def}{=} \max_{S \subseteq \{1, \dots, K\}, |S|=k} \sum_{i \in S} \sum_{t=1}^T x_i(t)$

$$\sqrt{\frac{(e-1)K}{(e-2)k}} \text{ 倍}$$

Cf. full Informationの場合は

$$2k\sqrt{(e-2)T \ln \frac{K}{k}}$$

損失版は[Warmuth & Kuzmin 08]



Combinatorial Bandit問題 [Cesa-Bianchi & Lugosi 2009]

- 各時刻 t において、 K 本のアームの集合 $X = \{1, 2, \dots, K\}$ において、あらかじめ決められた部分集合族 $S \subseteq 2^X$ からアームの集合 S_t を選択

multiple playの場合 : $S = \{U \subseteq X : |U| = k\}$

- 選ばれたアームに対する報酬の合計 $\sum_{i \in S_t} x_i(t)$ のみ知らされる
- 各アームの報酬は、同時に選択された他のアームの報酬に影響されない



ComBandアルゴリズム [Cesa-Bianchi & Lougosi 2009]

- 集合 S を K 次元ベクトル $1_S \in \{0,1\}^K$ で表現
- 集合 S 毎に重み q_S をもつ
- 集合族 S 上の確率分布 $\{p_S\}$ に従って S_t を選択
ただし、

$$p_S = (1-\gamma)q_S + \gamma/K$$

- アドバーサリーが決めた報酬 $x(t) = (x_1(t), \dots, x_K(t))$ に対し、
プレイヤーが選んだ S_t に属するアームに対する報酬の合計
 $x(t)^T 1_{S_t}$ が与えられる
- $x(t)$ の不偏推定量となる推定値 $\hat{x}(t)$ を次式で求める

$$\hat{x}(t) = (x(t)^T 1_{S_t}) P_t^{-1} 1_{S_t} \quad \text{where } P_t = E_{\{p_S\}} [1_S 1_S^T]$$

- S に対するロスの推定値 $\sum_{i \in S} \hat{x}_i(t)$ を使って、exponentialに q_S を更新



Combinatorial Bandit問題のリグレット上界

$$E(L_{\text{ComBand}}) - \min_{S \in \mathcal{S}} \sum_{t=1}^T x(t)^T \mathbf{1}_S \leq \left(2 + \frac{B}{K \lambda_{\min}(M)} \right) B \sqrt{TK \ln N}$$

ただし、

$$M = \frac{1}{N} \sum_{S \in \mathcal{S}} \mathbf{1}_S \mathbf{1}_S^T$$

λ_{\min} : M の非ゼロ最小固有値

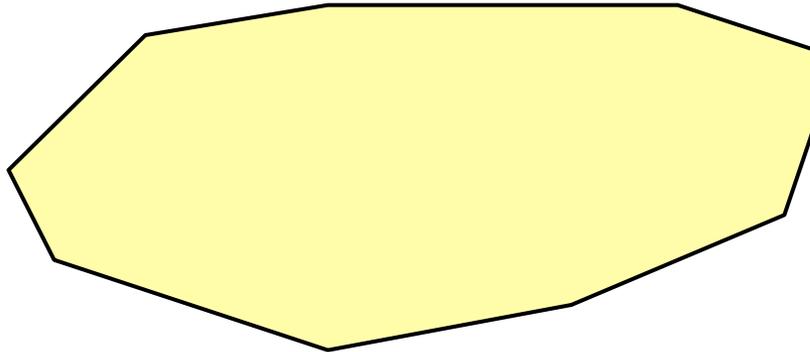
$$B = \max_{S \in \mathcal{S}} |S|, \quad N = |\mathcal{S}|$$

Cf. 上の上界をmultiple play設定に適用すると

$O\left(\sqrt{k^3 TK \ln K}\right)$ ←Capping法をよるアルゴリズムより $O(k)$ 倍悪い



online linear optimization 問題 [McMahan and Blum 2004]



$Q \subseteq \mathbb{R}^n$: コンパクト閉凸集合

$x_t \in \mathbb{R}^n$

時刻 t における
報酬ベクトル

x_t

悪魔が選ぶ。
playerは知らない

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

1. Q 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

期待リグレット $\max_{q \in K} \sum_{t=1}^T x_t \cdot q - E(G_A(T))$ を小さくするアルゴリズム A は？



Abernethyらのアルゴリズム [Abernethy et al. 2008]

Algorithm BOLO (仮称) // Bandit Online Linear Optimization

Input: $\eta > 0$, θ -self-concordant $\Gamma(\cdot)$

Initialization: $\mathbf{q}_1 \leftarrow \operatorname{argmin}_{\mathbf{q} \in Q} \Gamma(\mathbf{q})$

for $t=1$ to T

1. $\{e_1, e_2, \dots, e_n\} \leftarrow \nabla^2 \Gamma(\mathbf{q}_t)$ の固有ベクトル
 $\{\lambda_1, \lambda_2, \dots, \lambda_n\} \leftarrow \nabla^2 \Gamma(\mathbf{q}_t)$ 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. $\mathbf{r}_t \leftarrow \mathbf{q}_t + \varepsilon_t \lambda_{i_t}^{-1/2} e_{i_t}$
4. 報酬 $\mathbf{x}_t \cdot \mathbf{r}_t$ を得る。
5. 次のように更新する。

$$\hat{\mathbf{x}}_t \leftarrow n(\mathbf{x}_t \cdot \mathbf{r}_t) \varepsilon_t \lambda_{i_t}^{1/2} e_{i_t}$$

$$\mathbf{q}_{t+1} \leftarrow \operatorname{argmin}_{\mathbf{q} \in Q} - \sum_{s=1}^t \hat{\mathbf{x}}_s \cdot \mathbf{q} + \Gamma(\mathbf{q})$$

期待リグレット

$$\max_{\mathbf{q} \in Q} E\left(\sum_{t=1}^T \mathbf{x}_t \cdot \mathbf{q}\right) - E(G_{BOLO}(T))$$

$$= O\left(n\sqrt{\theta T \ln T}\right)$$

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



Gaussian Process バンディットアルゴリズム(1/2)

[Dorand, Glowacka & Showe-Taylor 2009]

X : アーム集合

κ : X 上のカーネル(アーム間の報酬の共分散を表す)

各アーム x の報酬 y は次式により与えられるとする。

$$y = f(x) + \varepsilon \quad \text{where } \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$f \sim \text{GP}(0, \kappa(x, x'))$ と仮定(事前分布)

$$\left(\begin{array}{l} E(f(x)) = 0 \text{ for all } x \in X \\ E[f(x)f(x')] = \kappa(x, x') \text{ for all } x, x' \in X \end{array} \right)$$

$(x_1, y_1), \dots, (x_t, y_t)$ を得た後の $f(x)$ の事後分布は $\mathcal{N}(\mu_t(x), \sigma_t^2(x))$

ただし、 $(C_t)_{i,j} = \kappa(x_i, x_j) + \sigma^2 \delta_{i,j}$

$$(\mathbf{k}_t(x))_i = \kappa(x, x_i)$$

$$\mu_t(x) = \mathbf{k}_t(x)^T C_t^{-1} \mathbf{y}_t$$

$$\sigma_t^2(x) = \kappa(x, x) - \mathbf{k}_t(x)^T C_t^{-1} \mathbf{k}_t(x)$$



Gaussian Processバンディットアルゴリズム(2/2)

UCBタイプのアームの選択を行う

$$x_{t+1} = \operatorname{argmax}_{x \in X} \{f_t(x) = \mu_t(x) + B(t)\sigma_t(x)\}$$

where $B(t)$: 知識の利用と獲得のバランスをとる関数

$$\left[\text{UCB1では } \sqrt{\log t} \text{ に比例する量} \right]$$



目次

1. バンディット問題とは
2. stableな確率的バンディット 問題の解法
 - (1) infinite horizon (geometric discount)の最適手法
 - (2) finite horizonの手法
3. Adversarial バンディット問題の解法
4. バンディットの拡張
 - (1) multiple play
 - (2) combinatorial bandit
 - (3) online linear optimization
 - (4) Gaussian Process Bandit
5. バンディット問題の応用



バンディットの応用(最近の話題)

- ゲームの探索木における準最適手の探索

UCT[Kosis & Szepesvari 2006]

- ニュース記事のリコメンド

LinUCBを使ったcontextual バンディット手法

[Li, Chu, Langford & Schapire 2010]

ログデータを使ってオフラインでバンディット手法を評価する方法

[Li, Chu, Langford & Wang 2011]