

# 報酬が過去の履歴に依存する場合のバンディットアルゴリズム

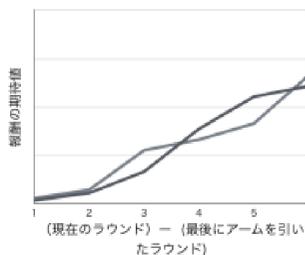
谷本 悠斗 総合研究大学院大学統計科学専攻博士課程(5年一貫制)5年

## 多腕バンディット問題

- 目標：リグレット最小化( $\Leftrightarrow$ 累積報酬最大化)
  - リグレット  $= \underbrace{E[\sum_{t=1}^T r_t^*]}_{\text{Oracle policy}} - \underbrace{E[\sum_{t=1}^T r_t]}_{\text{Learner}}$
- 探索と活用のトレードオフ
  - 探索：各アームの報酬の情報収集
  - 活用：探索での情報を基に累積報酬を最大化
- アームが定常（報酬の期待値がroundに関わらず一定）
  - 活用時に同じアームばかり引く
  - 満たさない例：推薦システム（同じ商品ばかり勧めてしまう）

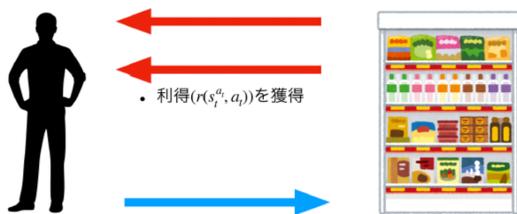
## 問題設定

- (非定常な) 報酬の構造
  - アーム*i*の報酬の期待値は最後にアーム*i*を引いた時点から経過したラウンド数が長いほど大きい
    - アームを引く  $\Rightarrow$  そのアームの期待値は低下
    - アームを引かない  $\Rightarrow$  そのアームの期待値は上昇
- 応用例：多様な商品の推薦
  - 商品を購入  $\rightarrow$  続けて同じ商品を勧められても購入しない
  - 推薦せずに時間が経過  $\rightarrow$  再び商品を購入する可能性が上昇



## 状態と利得の構造

- $K$ 次元の状態(state)  $s_t = (s_t^1, s_t^2, \dots, s_t^K)$  ( $K$ : アーム数)
  - 次期状態  $s_{t+1}^a, s_{t+1}^k = 1, s_{t+1}^k = \min\{s_t^k + 1, s_{max}\}$  (for  $k \neq a$ )
- $s_t^k$ : 最後にアーム*k*を引いた時点からの経過ラウンド数 ( $s_t^k = 1, 2, \dots, s_{max}$ ) ( $\forall k$ )
  - 利得  $(r(s_t^a, a))$  を獲得
- $s$ の組み合わせ:  $|s| = s_{max}^K$
- $r(s^k, a)$ : アーム*i*を引いた時の報酬
  - 報酬はアーム*k*の状態のみに依存
    - 報酬の単調性:  $r(s^a, a) \leq r(s^a + 1, a)$
  - 経過ラウンド数が一定以上 ( $s_{max}$ ) になるとそれ以上は利得は一定



## 各状態の次元ごとのQ関数の導入

- $s_{max}$  or  $K$ のいずれかが大きいと  $s$ の組み合わせ爆発によりQ学習ではうまくいかない
  - $\rightarrow$  Q関数の組み合わせを減らすアルゴリズムの必要性
- $k$ 次元目(アーム*k*)のみの状態のアーム*a*に関するQ関数:
 
$$Q_\pi(s^k, a) := E_\pi[\sum_{t=0}^T \gamma^t r(s_t^a, a) | s_0^k = s^k, a_0 = a] = E_\pi[r(s^a, a) + \gamma V_\pi((s^k)) | s_0^k = s^k, a_0 = a]$$
  - $\pi$ : policy,  $\gamma \in (0, 1)$ : 割引率,  $T$ : 期間,  $V_\pi(s^k) = E_\pi[\sum_{t=0}^T \gamma^t r(s_t^a, a) | s_0^k = s^k]$  ( $k$ 次元目の状態の価値関数)
  - Q関数の組み合わせ:  $K^2 \times s_{max}$
  - 通常の( $K$ 次元の状態をもつ)Q関数( $Q_\pi(s, a) := E_\pi[\sum_{t=0}^T \gamma^t r(s_t^a, a) | s_0 = s, a_0 = a]$ )は  $K$ 個の  $Q_\pi(s^k, a)$  ( $k = 1, 2, \dots, K$ )の平均と等しい (i. e.  $Q_\pi(s, a) = \frac{1}{K} \sum_{k=1}^K E_\pi[Q_\pi(s^k, a) | s_0 = s, a_0 = a]$ )

## アルゴリズム

### Algorithm 1 State-separated Q-learning

```

1: Input  $T$ :Horizon,  $E$ :Number of Exploration,  $\gamma$ :Discount Rate
2: Initialize  $Q(s^k, a) \leftarrow 0 \forall (s, k, a), s^k \leftarrow s_{max} \forall k$ 
3: for  $t = 1, 2, \dots, T$  do
4:   choose  $a \leftarrow \begin{cases} \text{randomly} & (t \leq E) \\ \text{argmax}_{a'} \frac{1}{K} \sum_{k=1}^K Q(s^k, a') & (t \geq E) \end{cases}$ 
5:    $r \leftarrow r(s^a, a)$ 
6:    $Q(s^k, a) \leftarrow Q(s^k, a) + \alpha(r(s^a, a) + \gamma \max_{a'} Q((s^k)', a) - Q(s^k, a)) \forall k$ 
7:    $s^k \leftarrow (s^k)'$   $\forall k$ 
8: end for
    
```

- アームの選択
  - 探索：指定された回数( $E$ )ランダムにアームを選択
  - 活用：各アームの  $K$ 個のQ関数の平均を比較して最も大きいアームを選択
- 選んだアームに関する  $K$ 個のQ関数を更新

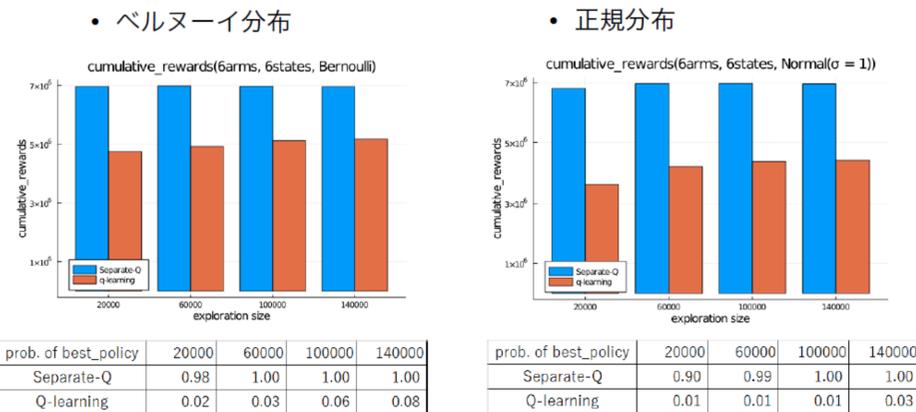
## シミュレーション

- 利得の分布：ベルヌーイ分布、正規分布 (分散 = 1)
- ラウンド数:  $10^7$ 、シミュレーション回数: 100
- 2つの実験のセッティング
  - 状態が1つ上がるたびにランダムに利得の期待値が上昇
  - $s_{max}$ でのアームの利得のみ以下のように期待値を固定

	(アーム数, 状態数)	状態が一つ上がるたびに上がる利得の期待値	$S_{max}$ の利得の期待値
実験1	(6, 6)	$0.07 \times \text{Uni}(0, 1)$	全部のアームで0.7
実験2	(10, 5)	$0.015 \times \text{Uni}(0, 1)$	5つのアームのみ0.7

- 評価の指標 (通常のQ学習と比較)
  - 累積利得の平均(cumulative rewards), 最適政策の割合(prob. of best policy)

## 実験1 (アーム数 = 6, 状態数 = 6)



## 実験2 (アーム数 = 10, 状態数 = 5)

