

# Piecewise-stationary Bandit Problems with Side Observations

Jia Yuan Yu<sup>1</sup>   Shie Mannor<sup>1 2</sup>

<sup>1</sup>Department of Electrical and Computer Engineering  
McGill University

<sup>2</sup>Faculty of Electrical Engineering  
Technion

ICML, Montréal. June 2009

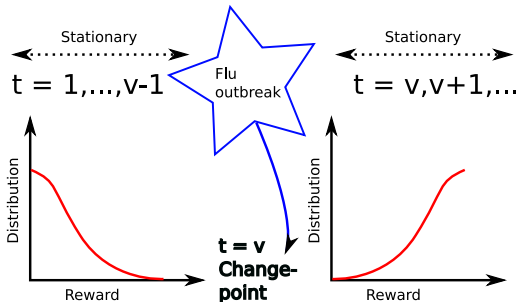
# Example: Investment with queries

- Investment experts  $1, \dots, n$ .



# Example: Investment with queries

- Each expert has a piecewise-stationary **reward**.
  - ▶ Stationary distribution, but with **abrupt** changes at **unknown** instants.



# The problem

- At each time instant:
  - ▶ choose and follow one expert,
  - ▶ query the reward of other experts,
  - ▶ receive reward of chosen expert,
  - ▶ observe rewards of queried experts,
  - ▶ pay cost-of-query.
- Goal: do as well as if reward distributions and change-points were known in advance.

# Why this model?

Stationary

vs.

Adversarial

Somewhere in between.

# Notation: Piecewise-stationary rewards

- Experts  $1, \dots, n$ .
- $b_t(i) \in [0, 1]$  is reward of  $i$ -th expert at time  $t$
- Rewards form a sequence of random vectors:

$$\underbrace{\begin{bmatrix} b_1(1) \\ b_1(2) \\ \vdots \\ b_1(n) \end{bmatrix}, \dots, \begin{bmatrix} b_{\nu_2}(1) \\ b_{\nu_2}(2) \\ \vdots \\ b_{\nu_2}(n) \end{bmatrix}}_{\text{distribution } f_1, \text{ mean } \beta_1}, \dots, \underbrace{\begin{bmatrix} b_{\nu_j}(1) \\ b_{\nu_j}(2) \\ \vdots \\ b_{\nu_j}(n) \end{bmatrix}, \dots, \begin{bmatrix} b_{\nu_{j+1}}(1) \\ b_{\nu_{j+1}}(2) \\ \vdots \\ b_{\nu_{j+1}}(n) \end{bmatrix}}_{\text{distribution } f_j, \text{ mean } \beta_j}, \dots$$

- **Unknown** change-points:  $\nu_2, \nu_3, \dots$
- Between change-points: fixed, but **unknown** distributions  $f_1, f_2, \dots$

# More notations

- At each time step  $t$ ,
  - ▶ pick one expert, say  $a_t$ ,
  - ▶ query the reward of a subset of additional experts, say  $\mathcal{S}_t$ ,
  - ▶ receive reward  $b_t(a_t)$ ,
  - ▶ observe  $\{b_t(j) \text{ for } j \in \mathcal{S}_t\}$ ,
  - ▶ pay query cost  $C_Q(|\mathcal{S}_t|)$ .

# Objective

- **Optimal** expected reward: when every **reward distribution is known in advance**.
- **Expected regret** at time  $T$ :

$$R_T \triangleq \sum_{t=1}^T \max_{i=1, \dots, n} \beta_t(i) - \sum_{t=1}^T \mathbb{E}[b_t(a_t)]. \quad (1)$$

- If we make  $\ell$  queries per step, the overall cost-per-step is

$$R_T/T + C_Q(\ell).$$



## Other models: Stochastic MAB

- Stochastic multi-armed bandit:
  - ▶ no change-point,
  - ▶ no queries.
- When there is **no** change-point, **few** queries are needed:
  - ▶ asymptotically optimal solution with logarithmic number of exploration steps.
- Expected regret of  $O(n \log(T))$ .

## Other models: Stochastic MAB

- Stochastic multi-armed bandit:
  - ▶ no change-point,
  - ▶ no queries.
- When there is **no** change-point, **few** queries are needed:
  - ▶ asymptotically optimal solution with logarithmic number of exploration steps.
- Expected regret of  $O(n \log(T))$ .

## Other models: Expert-prediction

- Prediction with expert-advice:
  - ▶ arbitrarily many change-points,
  - ▶ free queries of every expert at every time step.
- Adversarial multi-armed bandit:
  - ▶ arbitrarily many change-points,
  - ▶ no queries.
- Different notion of **adversarial** regret:

$$\max_{i=1,\dots,n} \sum_{t=1}^T \beta_t(i) - \sum_{t=1}^T \mathbb{E}[b_t(a_t)]. \quad (2)$$

- When there are **many** change-points, queries have **small** effect:
  - ▶ with free queries:  $O(\sqrt{T \log(n)})$ ;
  - ▶ without any query:  $O(\sqrt{Tn \log(n)})$  (Auer et al., 2002b).

## Other models: Expert-prediction

- Prediction with expert-advice:
  - ▶ arbitrarily many change-points,
  - ▶ free queries of every expert at every time step.
- Adversarial multi-armed bandit:
  - ▶ arbitrarily many change-points,
  - ▶ no queries.
- Different notion of **adversarial** regret:

$$\max_{i=1,\dots,n} \sum_{t=1}^T \beta_t(i) - \sum_{t=1}^T \mathbb{E}[b_t(a_t)]. \quad (2)$$

- When there are **many** change-points, queries have **small** effect:
  - ▶ with free queries:  $O(\sqrt{T \log(n)})$ ;
  - ▶ without any query:  $O(\sqrt{Tn \log(n)})$  (Auer et al., 2002b).

# Known results for piecewise-stationary bandit without queries

- (Hartland et al., 2006) provide a partial solution that detects only distribution changes in the current best expert.

## Definition

$k \triangleq$  number of changes up to time  $T$

- If  $k$  is known in advance:
  - ▶ Fixed-share algorithm (Herbster & Warmuth, 1998) gives a regret of  $O(\sqrt{nkT \log(\bar{T})})$ .
  - ▶ Discounted- & Sliding-window-UCB algorithm (Garivier & Moulines, 2008) give a regret of  $O(n\sqrt{kT \log(T)})$ .
  - ▶ Lower-bound of  $\Omega(\sqrt{T})$  without queries.

## Our result for piecewise-stationary rewards

- Queries reduce regret from  $\Omega(\sqrt{T})$  to  $O(nk \log(T))$ ,
- **without** prior knowledge of  $k$ .

# Our approach

- Run standard algorithm for stochastic MAB,
  - ▶ e.g., UCB algorithm of (Auer et al., 2002a).
- Reset when we detect a change-point.
- How to detect change-points?
  - ▶ Naive solution: detect changes in the distribution directly using change-detection algorithms (CUSUM, etc.).
  - ▶ Simpler solution: detect changes in the empirical mean over windows of appropriate length.

# Our approach

- Run standard algorithm for stochastic MAB,
  - ▶ e.g., UCB algorithm of (Auer et al., 2002a).
- Reset when we detect a change-point.
- How to detect change-points?
  - ▶ Naive solution: detect **changes in the distribution** directly using change-detection algorithms (CUSUM, etc.).
  - ▶ Simpler solution: detect **changes in the empirical mean** over windows of appropriate length.



# Our approach

- Run standard algorithm for stochastic MAB,
  - ▶ e.g., UCB algorithm of (Auer et al., 2002a).
- Reset when we detect a change-point.
- How to detect change-points?
  - ▶ Naive solution: detect **changes in the distribution** directly using change-detection algorithms (CUSUM, etc.).
  - ▶ Simpler solution: detect **changes in the empirical mean** over windows of appropriate length.

# WMD algorithm

## Windowed Mean-shift Detection

Break time horizon into intervals of length  $\tau$ , compute empirical mean in each interval:

$$\underbrace{b_1, b_2, \dots, b_\tau}_{\hat{b}_1}, \underbrace{b_{\tau+1}, \dots, b_{2\tau}}_{\hat{b}_2}, \dots, \underbrace{b_{(m-1)\tau+1}, \dots, b_{m\tau}}_{\hat{b}_m} \dots$$

At each time step  $t$ :

- 1 Follow UCB algorithm:
  - ▶ Play the expert with highest upper-confidence index:

$$\hat{b}(i) + \sqrt{2 \log(T) / \#(i)}.$$

- 2 Query experts with equal frequency.
- 3 Detect changes: If  $\|\hat{b}_m - \hat{b}_r\|_\infty > \epsilon$ , reset UCB sub-algorithm.

# Guarantee

## Piecewise-stationary bandit with queries

### Theorem

Suppose that at every change-point, the mean reward of some expert changes by at least  $2\epsilon$ , i.e.,  $|\beta_{\nu_j}(i) - \beta_{\nu_{j+1}}(i)| > 2\epsilon$ . The WMD algorithm with windows of length  $\tau = \lfloor \frac{n}{\ell} \rfloor \cdot \lfloor \frac{\log(T)}{2\epsilon^2} \rfloor$  achieves a regret of

$$R_T \leq \frac{7}{\epsilon^2} \frac{kn}{\ell} \log(T) + \frac{C}{\Delta^2} kn \log(T) + \frac{6C}{\Delta^2} n^2, \quad (3)$$

for every sequence of change-points  $\nu_1, \nu_2, \dots$  and every choice of post-change distributions  $f_{\nu_1}, f_{\nu_2}, \dots$

- This regret is  $O(nk \log(T))$  **without** prior knowledge of  $k$ , but **with** queries.
- Recall: lower-bound of  $\Omega(\sqrt{T})$  **with** prior knowledge of  $k$ , but **without** queries.

# Proof ideas

- $L \triangleq$  Expected number of intervals between change-point and its detection.
- $N(T) \triangleq$  Expected number of false detections up to  $T$ .
- Two components of the regret:
  - ▶ expected number of resets is  $k + N(T)$ ; hence, regret between resets is

$$(Cn/\Delta^2)(k + N(T)) \log(T).$$

- ▶ regret due to delay in detection:

$$k(L + 1)\tau.$$

- Bound  $N(T)$  and  $L$  using Hoeffding's Inequality on empirical mean of i.i.d. rewards in each window.

## A (partial) lower bound

- Consider class of algorithms that **detect-and-react**.
  - ▶ Constraint: as many switches between experts as distribution changes;
  - ▶ or a vanishing frequency of false detection.
- Optimal delay to change-detection is  $\Omega(\log(T))$  (Lorden, 1971).
- For such algorithms, regret is lower-bounded by  $\Omega(k \log(T))$ .

# Experiments

## Piecewise-stationary rewards with Bernoulli distributions

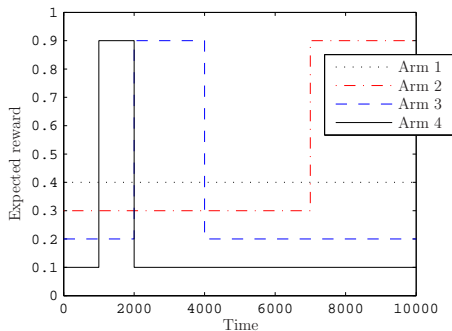
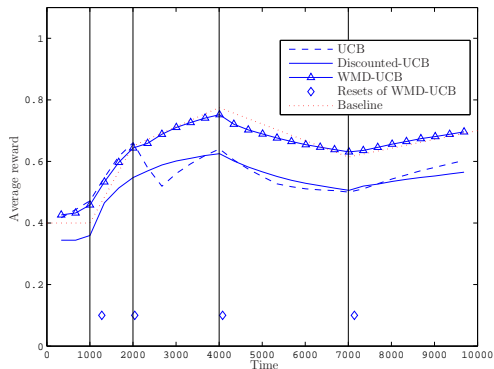


Figure: Four experts with Bernoulli distributed piecewise-stationary rewards.

# Simulation: Average expected regret



- UCB algorithm (Auer et al., 2002a);
- Discounted-UCB algorithm of (Kocsis & Szepesvári, 2006) with prior knowledge of the number  $k$  of changes;
- WMD-UCB algorithm with 1 query per step.

# Query-regret trade-off

- Overall expected cost-per-step at time  $T$ :

$$C_Q(\ell) + R_T/T.$$

- If  $C_Q(\ell) = c_q \times \ell$ , then optimize cost-per-step with respect to  $\ell$ .
- Optimal query rate is

$$\ell^* = \sqrt{(7kn/c_q) \log(T)/T}.$$



# Open questions

- Non-uniform querying (confidence bound-type algorithms for mean-shift detection).
- Probabilistic queries (with a failure probability).
- Combine query and decision: we query two experts and then receive the highest reward of the two.
- Restless bandit problems.

# References

- Auer, P., Cesa-Bianchi, N., & Fischer, P. (2002a). Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47, 235–256.
- Auer, P., Cesa-Bianchi, N., Freund, Y., & Schapire, R. E. (2002b). The nonstochastic multiarmed bandit problem. *SIAM J. Computing*, 32, 48–77.
- Garivier, A., & Moulines, E. (2008). On upper-confidence bound policies for non-stationary bandit problems. Preprint. <http://arxiv.org/abs/0805.3415>.
- Hartland, C., Gelly, S., Baskiotis, N., Teytaud, O., & Sebag, M. (2006). Multi-armed bandit, dynamic environments and meta-bandits. Preprint. <http://hal.archives-ouvertes.fr/hal-00113668/en/>.
- Herbster, M., & Warmuth, M. K. (1998). Tracking the best expert. *Machine Learning*, 32, 151–178.
- Kocsis, L., & Szepesvári, C. (2006). Discounted-UCB. 2nd PASCAL Challenges Workshop.
- Lorden, G. (1971). Procedures for reacting to a change in distribution. *Ann. Math. Statist.*, 42, 1897–1908.
- Pollak, M. (1985). Optimal detection of a change in distribution. *Ann. Statist.*, 13, 206–227.

# Shiryayev-Roberts scheme

Raise alarm ( $f_0 \rightarrow f_\theta$ ) at time  $t$  if

$$\sum_{k=1}^t \frac{f_\theta(b_k) \dots f_\theta(b_t)}{f_0(b_k) \dots f_0(b_t)} \geq A. \quad (4)$$

## Theorem (Average run-lengths (Pollak, 1985))

- 1 If there is no change, then  $\mathbb{E}_\infty[\text{Alarm Time}] \geq A$ .
- 2 If a change occurs at time 1, then

$$\mathbb{E}_1[\text{Alarm Time}] = [\log A + \log \log A + O(1)] / D(f_\theta; f_0).$$

This is *optimal*.

# Detect-and-act algorithms

Class of algorithms where:

- # switches between experts  $\leq$  # distribution changes + 1

Consequence:

- Expected # false-detections  $\leq 1$ .

## Theorem (Regret lower-bound)

*For every fixed algorithm of our class (fixed  $\ell$  and query mechanism), there exists a piecewise-stationary source such that detect-and-react algorithm has a regret of at least*

$$R_T \geq k \frac{n}{\ell} \log T / D(f_\theta; f_0).$$

Proof:

- Bernoulli-distributed sources:
  - ▶ Detecting change in distribution requires detecting change in mean.
- Average run-length theorem.