# Learning to Control Renewal Processes with Bandit Feedback

Semih Cayci<sup>1</sup>, Atilla Eryilmaz<sup>1</sup>, R. Srikant<sup>2</sup>

<sup>1</sup>The Ohio State University, ECE

<sup>2</sup> University of Illinois at Urbana-Champaign, CSL and ECE

## Multi-Armed Bandits

Prominent model for exploration and exploitation dilemma since 1933.



What are we investigating?

- Time dimension into multi-armed bandits as required by many applications
- Very general time and reward distributions including heavy-tailed distributions
- Novel dynamics + new dilemma

#### Multi-Armed Bandits: Applications

Traditional bandit models have a broad area of applications:



2

**Stochastic bandits** [Lai & Robbins, '85], [Auer, '02], [Agrawal, '11], [Bubeck, '12a]

Heavy-tailed reward [Bubeck, '12b], [Zhao, '11]







**Contextual bandits** [Dudik, '11], [Slivkins, '14] **Linear bandits** [Dani, '08], [Abbasi-Yadkori, '11]





**Budgeted bandits** [Gyorgy, '07], [Tran-Tranh, '10], [Badanidiyuru, '13], [Combes, '15], [Xia, '15]





- Cost and reward in [0, 1]
- No interrupt mechanism
- Limited control

# Recap: Classical Stochastic Multi-Armed Bandit Problem

In the classical stochastic bandit model:

Q: Is it always valid?

- Discrete time: Each arm pull takes a unit time for all arms.
- Control: Choose an arm I<sub>n</sub> among {1, 2, ..., K}
- Goal: Maximize cumulative reward in τ units of time.

$$\mathbb{E}\big[\mathsf{Rew}_{\pi}(\tau)\big] = \mathbb{E}\big[\sum_{\mathsf{n}=1}^{\tau}\mathsf{R}_{\mathsf{n}}^{(\mathsf{I}_{\mathsf{n}})}\big]$$

• Optimal policy: Arm with maximum expected reward.

$$\pi_n^{\texttt{opt}} = \arg\max_k \ \mathbb{E}[\mathsf{R}_n^{(k)}]$$

1<sup>st</sup>-order statistics

 Learning: Upper confidence bound (UCB), ε-greedy [Auer '02], Thompson sampling [Agrawal, '11], IDS [Russo, '14]



### Time Dimension in Bandits: Task Scheduling Example

Non-clairvoyant task scheduling with K task types and single server.

- Type-k task takes a random **completion time** X<sub>n</sub><sup>(k)</sup>.
- A random **reward** R<sub>n</sub><sup>(k)</sup> is obtained once the task is completed.
- For a **time budget**  $\tau$ , maximize total reward in  $[0, \tau]$
- Completion time might be **heavy-tailed** [Harchol-Balter, '99]
  - What if the completion time is "too long"? Should I wait?
- Optimal policy?

? 
$$\arg \max_{k} \mathbb{E}[\mathsf{R}_{\mathsf{n}}^{(k)}] \quad \arg \max_{k} \mathbb{E}\left[\frac{\mathsf{R}_{\mathsf{n}}^{(k)}}{\mathsf{X}_{\mathsf{n}}^{(k)}}\right] \quad \arg \max_{k} \frac{\mathbb{E}[\mathsf{R}_{1}^{(k)}]}{\mathbb{E}[\mathsf{X}_{1}^{(k)}]}$$

• Existing **MAB** models fall short for this application.





#### Time Dimension in Bandits: Applications

There are many other applications that have similar time-dependence. Here are some examples:



Adaptive routing in telecommunications:

- K parallel channels, transmit over one at a time, bandit feedback
- Transmission time of n<sup>th</sup> packet over channel k: X<sub>n</sub><sup>(k)</sup> (heavy-tailed [Asmussen, '10])
- Goal: Maximize throughput within [0,  $\tau$ ], i.e.,  $R_n^{(k)} = 1$





#### Bandits with Interrupts (BwI)



New dilemma: Should it interrupt an ongoing cycle and reject the reward for a more rewarding alternative, or wait and gather the reward?

# Optimal Policy with Known Statistics: Complexity

**BwI** problem is NP-hard even if all statistics are known.

Wanted: Approximation algorithms



To find an approximation algorithm:



Consider a **renewal reward process** (X<sub>n</sub>, R<sub>n</sub>).

By the **key renewal theorem**, we have:

$$\mathbb{E}\Big[\sum_{n=1}^{\mathsf{N}(\tau)}\mathsf{R}_n\Big] = \frac{\mathbb{E}[\mathsf{R}_n]}{\mathbb{E}[\mathsf{X}_n]} \cdot \tau + \mathsf{o}(\tau)$$



For an action (k, b), the **renewal reward rate** is defined as follows:



**Interpretation:** Reward per unit time if (k, b) is persistently chosen.

(See [Asmussen, '08] and [Gut, '09])

# Optimal Policy with Known Statistics: A Good Static Approximation

**Approximation algorithm:** Pull (k, b) with the largest reward rate persistently until the time budget is depleted.

 $\pi_n^* = \arg \max_{(k,b)} r^{(k)}(b)$ 

for all n.

How well does  $\pi^*$  approximate the optimal policy  $\pi^{opt}$ ?





Low complexity, time-invariant.



The optimality gap with  $\pi^{opt}$  is **bounded** for all  $\tau > 0$ .



**Challenge:** r<sup>(k)</sup>(b) depends on the tails, not only mean.

#### **Proposition 1. (Optimality Gap for** $\pi^*$ **)**

**For any**  $\tau > 0$ , the following inequality holds for the static policy  $\pi^*$ :

$$\mathbb{E}[\mathsf{Rew}_{\pi^{\mathsf{opt}}}(\tau)] - \mathbb{E}[\mathsf{Rew}_{\pi^*}(\tau)] = \mathsf{O}(1)$$

Consequently,  $\pi^*$  is asymptotically optimal as  $\tau \rightarrow \infty$ .



Appropriate for low-regret learning algorithms.

## Optimal Policy with Known Statistics: Optimal Interrupt Time

**Renewal reward rate** for (k, b): 
$$r^{(k)}(b) = \frac{\mathbb{E}[R_n^{(k)}\mathbb{I}_{\{X_n^{(k)} \le b\}}]}{\mathbb{E}[X_n^{(k)} \land b]}$$



Does interruption improve reward rate?

$$b_k^* = \max_{b \in \mathcal{B}} r^{(k)}(b)$$

#### **Proposition 2. (Optimal Interrupt Time)**

Interruption at a finite time is optimal, i.e.,  $b_k^* \neq \infty$ , if and only if the following holds for some  $b < \infty$ :

 $\mathbb{E}[X_1^{(k)} - b | X_1^{(k)} > b] > \mathbb{E}[X_1^{(k)}]$ 

Interruption is optimal if the **mean residual life** at some time **b** is larger than the mean completion time of a fresh new cycle.

For **correlated**  $(X_n^{(k)}, R_n^{(k)})$ , the condition becomes the following:

# $\frac{\mathbb{E}[\mathsf{R}_1^{(k)} - \mathsf{b}|\mathsf{X}_1^{(k)} > \mathsf{b}]}{\mathbb{E}[\mathsf{X}_1^{(k)} - \mathsf{b}|\mathsf{X}_1^{(k)} > \mathsf{b}]} < \frac{\mathbb{E}[\mathsf{R}_1^{(k)}]}{\mathbb{E}[\mathsf{X}_1^{(k)}]}$



All **heavy-tailed** and some light-tailed completion time distributions lead to a finite optimal interrupt time.



Exponential distribution acts as a barrier case: **memoryless property** leads to indifference to interruption.



Most light-tailed distributions (folded Gaussian, uniform, logistic, gamma) have **decreasing MRL**, thus no interrupt is optimal.



All have the same mean → Tails are important

## Optimal Policy with Known Statistics: Optimal Interrupt Time

**P** Renewal reward rate for (k, b): 
$$r^{(k)}(b) = \frac{\mathbb{E}[R_n^{(k)}\mathbb{I}_{\{X_n^{(k)} \leq b\}}]}{\mathbb{E}[X_n^{(k)} \wedge b]}$$



Does interruption improve reward rate?

$$b_k^* = \max_{b \in \mathcal{B}} r^{(k)}(b)$$

#### **Proposition 2. (Optimal Interrupt Time)**

Interruption at a finite time is optimal, i.e.,  $b_k^* \neq \infty$ , if and only if the following holds for some  $b < \infty$ :

 $\mathbb{E}[X_1^{(k)} - b | X_1^{(k)} > b] > \mathbb{E}[X_1^{(k)}]$ 

Interruption is optimal if the **mean residual life** at some time **b** is larger than the mean completion time of a fresh new cycle.

For **correlated**  $(X_n^{(k)}, R_n^{(k)})$ , the condition becomes the following:

# $\frac{\mathbb{E}[\mathsf{R}_1^{(k)} - \mathsf{b}|\mathsf{X}_1^{(k)} > \mathsf{b}]}{\mathbb{E}[\mathsf{X}_1^{(k)} - \mathsf{b}|\mathsf{X}_1^{(k)} > \mathsf{b}]} < \frac{\mathbb{E}[\mathsf{R}_1^{(k)}]}{\mathbb{E}[\mathsf{X}_1^{(k)}]}$



All **heavy-tailed** and some light-tailed completion time distributions lead to a finite optimal interrupt time.



Exponential distribution acts as a barrier case: **memoryless property** leads to indifference to interruption.



Most light-tailed distributions (folded Gaussian, uniform, logistic, gamma) have **decreasing MRL**, thus no interrupt is optimal.



All have the same mean  $\rightarrow$  **Tails are important** 

## Optimal Policy with Known Statistics: Optimal Interrupt Time

**P** Renewal reward rate for (k, b): 
$$r^{(k)}(b) = \frac{\mathbb{E}[R_n^{(k)}\mathbb{I}_{\{X_n^{(k)} \leq b\}}]}{\mathbb{E}[X_n^{(k)} \wedge b]}$$



Does interruption improve reward rate?

$$b_k^* = \max_{b \in \mathcal{B}} r^{(k)}(b)$$

#### **Proposition 2. (Optimal Interrupt Time)**

Interruption at a finite time is optimal, i.e.,  $b_k^* \neq \infty$ , if and only if the following holds for some  $b < \infty$ :

 $\mathbb{E}[X_1^{(k)} - b | X_1^{(k)} > b] > \mathbb{E}[X_1^{(k)}]$ 

Interruption is optimal if the **mean residual life** at some time **b** is larger than the mean completion time of a fresh new cycle.

For **correlated**  $(X_n^{(k)}, R_n^{(k)})$ , the condition becomes the following:

 $\frac{\mathbb{E}[\mathsf{R}_1^{(k)} - \mathsf{b}|\mathsf{X}_1^{(k)} > \mathsf{b}]}{\mathbb{E}[\mathsf{X}_1^{(k)} - \mathsf{b}|\mathsf{X}_1^{(k)} > \mathsf{b}]} < \frac{\mathbb{E}[\mathsf{R}_1^{(k)}]}{\mathbb{E}[\mathsf{X}_1^{(k)}]}$ 



All **heavy-tailed** and some light-tailed completion time distributions lead to a finite optimal interrupt time.



Exponential distribution acts as a barrier case: **memoryless property** leads to indifference to interruption.



Most light-tailed distributions (folded Gaussian, uniform, logistic, gamma) have **decreasing MRL**, thus no interrupt is optimal.



All have the same mean → Tails are important

#### Algorithm Design: Preliminaries

**Objective:** For all (k, b) actions, learn the reward rate:

$$\mathsf{r}^{(k)}(\mathsf{b}) = \frac{\mathbb{E}[\mathsf{R}_n^{(k)}\mathbb{I}_{\{X_n^{(k)} \leq \mathsf{b}\}}]}{\mathbb{E}[\mathsf{X}_n^{(k)} \wedge \mathsf{b}]}$$

Far beyond the first-order statistics unlike classical MAB.

Assumption: Finite but arbitrary set of interrupt times

$$\mathcal{B} = \left\{ b_1, b_2, \ldots, b_L = \infty \right\}$$

Imposed by the nature: CPU scheduling, 1/3/7/14 days in free trials, etc.

Strategy: Use upper confidence bound as a surrogate for r<sup>(k)</sup>(b<sub>1</sub>)

Find concentration inequalities for reward rate.

**Need:** Estimators that provide exponential convergence rate despite the heavy tails.



#### Exploit the information structure

imposed by the temporal dynamics.

There is an information sharing between different (k, b) actions based on observability.



3

Propose **UCB-BwI** for the online learning problem.

Based on the concentration inequalities, exploit the information structure for low-regret.

#### Concentration Inequalities for Renewal Processes

Wanted: UCB for 
$$r^{(k)}(b) = \frac{\mathbb{E}[\mathsf{R}_n^{(k)}\mathbb{I}_{\{X_n^{(k)} \le b\}}]}{\mathbb{E}[X_n^{(k)} \land b]}$$

First candidate: Empirical reward rate [Asmussen, '08], [Karlin, '83]  $\hat{r}^{(k)}(b_l) = \frac{\sum_{i=1}^{s} R_i \mathbb{I}_{\{X_i^{(k)} \le b_l\}}}{\sum_{i=1}^{s} (X_i \land b_l)} \rightarrow r^{(k)}(b_l) \text{ a.s.}$ 

Problem: Heavy-tailed completion time and reward

Convergence rate is **polynomial**, not exponential (Chebyshev)

$$\mathbb{P}\big(\hat{\mathsf{r}}_{\mathsf{s}}(\mathsf{b}_{\mathsf{L}}) \leq \mathsf{r}(\mathsf{b}_{\mathsf{L}}) - \Delta_{\mathsf{0}}(\epsilon)\big) = \mathsf{O}\Big(\frac{1}{\mathsf{s}^{\gamma} \epsilon^{1+\gamma}}$$

for  $\Delta_0(\epsilon) = \frac{(1+r)\epsilon}{\mu+\epsilon}$ . Bound is tight [Catoni, '10, Bubeck '12].

**Cause:** Outliers due to the heavy tails pull the sample mean away from the ensemble mean.

Fix: Robust median-of-means estimator [Nemirovski, '82], [Bubeck, '12]

Given {U<sub>i</sub>: i = 1, 2, ..., s}, w = 
$$\lfloor 8 \log(e^{\frac{1}{8}}\delta^{-1}) \wedge \frac{s}{2} \rfloor$$
 and m =  $\lfloor \frac{s}{w} \rfloor$   

$$M(U_{1:s}) \triangleq med \left\{ \frac{1}{m} \sum_{i=1}^{m} U_i, \dots, \frac{1}{m} \sum_{i=(w-1)m+1}^{wm} U_i \right\}$$

For each (k, b<sub>l</sub>), the median-of-means estimator for the reward rate  $r^{(k)}(b_l)$ :

$$\overline{r}_{s}(b) = \frac{\mathsf{M}\big(\mathsf{R}_{i}\mathbb{I}_{\{X_{i} \leq b\}}: i \leq s}{\mathsf{M}\big(X_{i} \wedge b: i \leq s\big)}$$

#### Proposition 3. (Conc. Ineq. for Renewal Reward Processes)

For any decision (k, b) and  $\delta \in (0, 1)$ ,

$$\mathbb{P}\Big(\overline{\mathbf{r}}_{\mathbf{s}}^{(\mathsf{k})}(\mathsf{b}) \leq \mathbf{r}^{(\mathsf{k})}(\mathsf{b}) - \Delta\big(\epsilon(\delta)\big)\Big) \leq \delta$$
  
ere  $\epsilon(\delta) = \beta\Big(\frac{\log(\nu\delta^{-1})}{\mathsf{s}}\Big)^{\frac{\gamma}{1+\gamma}}$ 

Exponential convergence rate

wh

Sub-Gaussian (optimal) accuracy-confidence tradeoff if  $\gamma = 1$  [Bubeck, '12b]

#### UCB-BwI Algorithm

#### **Information structure:**



#### UCB-BwI Algorithm

At epoch (n+1), UCB-Bwl Algorithm makes a decision as follows:

$$\left(\mathsf{I}_{n+1},\mathsf{B}_{n+1}^{(\mathsf{I}_{n+1})}\right) \in \underset{(\mathsf{k},\mathsf{b}_{l})}{\operatorname{arg\,max}} \left\{\overline{\mathsf{r}}_{\mathsf{n},\overline{\mathsf{T}}_{l}^{(\mathsf{k})}(\mathsf{n})}(\mathsf{b}_{l}) + \frac{(1+\mathsf{r})\epsilon_{\mathsf{n},\overline{\mathsf{T}}_{l}^{(\mathsf{k})}(\mathsf{n})}}{\mu + \epsilon_{\mathsf{n},\overline{\mathsf{T}}_{l}^{(\mathsf{k})}(\mathsf{n})}}\right\}$$

where

$$\epsilon_{\mathsf{n},\mathsf{s}} = \beta \Big[ \frac{\log \left( 2\mathsf{e}^{\frac{1}{8}} (\mathsf{n}+1)^4 \right)}{\mathsf{s}} \Big]^{\frac{\gamma}{1+\gamma}}$$

(

 $\bigcap \overline{T}_{I}^{(k)}(n)$  is the effective sample size  $\rightarrow$  information structure

- **Feedback** Is  $X_n^{(k)}$  smaller than b?
  - If so, values of  $X_n^{(k)}$  and  $R_n^{(k)}$

 $\mathbf{Q}$  These questions can be answered for all b'  $\leq$  b by this feedback [Observability]

Formally,  $Y(k, b') \in \sigma(Y(k, b)), \forall b' \leq b$ 

 $\bigcirc$  Action (k, b) increases the sample size for all (k, b') such that b'  $\leq$  b.

## Performance Analysis: Regret Bounds for UCB-BwI

2

3

Let

 $d_{min}^{(k)} \leq r^{(k)}(b_k^*) - \max_{1 \leq l \leq L} \ r^{(k)}(b_l) \quad \ \text{(Suboptimality in k)}$ 

 $d^{(k)} = r^* - \max_{1 \leq l \leq L} \ r^{(k)}(b_l) \quad \mbox{(Suboptimality across k)}$ 

#### Theorem 1. (Regret Upper Bound for UCB-BwI)

Regret under UCB-Bwl satisfies the following bound for all  $\tau$  > 0:

$$\overline{\mathsf{Reg}}_{\pi}(\tau) \leq \sum_{\mathsf{k}:\mathsf{d}^{(\mathsf{k})}>0} \mathsf{C}_{\mathsf{k}} \log\Big(\frac{\tau}{\min_{\mathsf{k}} \mathbb{E}[\mathsf{X}_{1}^{(\mathsf{k})}]}\Big) + \mathsf{O}(\mathsf{K} \times \mathsf{L})$$

where

$$\mathsf{C}_{\mathsf{k}} = \mathcal{O}\Big(\big(\frac{1}{\mathsf{d}_{\mathsf{min}}^{(\mathsf{k})}}\big)^{\frac{1}{\gamma}} + \big(\frac{1}{\mathsf{d}^{(\mathsf{k})}}\big)^{\frac{1}{\gamma}}\Big)$$

Logarithmic regret in  $\tau$ .

Action space: O(K x L) but  $C_k$  is O(K) as a result of info structure.

Proof Sketch: Asynchronous decisions + random number of trials

Dealing with asynchronous decisions: regret rate  

$$\overline{\text{Reg}}_{\pi}(\tau) \leq \left[\sum_{n=1}^{N_{\pi}(\tau)} \sum_{k,l} \mathbb{I}_{\{\pi_{n} = (k,b_{l})\}} \left(r^{*} - r^{(k)}(b_{l})\right) \mu_{\text{max}}\right] + \text{Res}$$
Dealing with randomness of N<sub>\pi</sub>(\theta): High prob. upper bounds for N<sub>\pi</sub>(\theta)
$$\overline{\text{Reg}}_{\pi}(\tau) \leq \sum_{k,l} \mathbb{E}[\mathsf{T}_{l}^{(k)}(\bar{n})] \left(r^{*} - r^{(k)}(b_{l})\right) \mu_{\text{max}}$$

$$+ \mathsf{K} \times \mathsf{L} \times r^{*} \mu_{\text{max}} \sum_{n > \bar{n}} \mathbb{P}(\mathsf{N}_{\pi}(\tau) > \mathsf{n}) + \text{Res}$$

$$O(1) \text{ if } \bar{\mathsf{n}} = \frac{2\tau}{\mu}$$

- Number of suboptimal decisions in  $2\tau/\mu$  epochs: Bandit analysis + information structure
- Matching lower bound of order  $\Omega(K \log \tau) \rightarrow order optimality$

### Performance Analysis: Numerical Results

#### Task scheduling with heterogeneous statistics



#### Conclusions



Presented BwI framework to incorporate time dimension into sequential learning.



General completion time and reward distributions: Interruption as a new control



Non-parametric learning algorithm UCB-BwI that achieves order optimality in all parameters  $\tau$ , K, L