# **Continuous Bayesian Bandits** & Open Problems

# Yuhua Zhu University of California - San Diego





## **Multi-armed Bandits**

#### $x \sim N(\mu, 1)$ , with unknown $\mu$ 100 Euro It costs 1 euro to play once



$$u_1$$

#### Q: What is the best strategy to get the maximum total rewards?

Method 1:

- pull  $a_1$  and  $a_2$  25 times each
- Get the empirical mean

 $\hat{\mu}_1, \hat{\mu}_2$ 

- Pull the arm has the larger empirical mean for the rest 50 times

Case 1:  $\mu_1 = 10$ ,  $\mu_2 = -10$ 

- -

Case 2:  $\mu_1 = 0.1$ ,  $\mu_2 = -0.1$ 

- worse arm

#### Tradeoff between "exploitation" and "exploration"



 $\mu_2$ 

- One does not need to pull 25 times to figure out  $a_1$  is better Waste too much time on the obviously worse arm



explore too much & exploit too little

- It is possible that  $\hat{\mu}_1 < \hat{\mu}_2$ , and the rest 50 times is pulling the

- Spend too little time to differentiate the two arms

explore too little & exploit too much





# Some better algorithms

#### **Tempered greedy:**

Pull 
$$a_1$$
 w.p.  $\pi_1 \propto e^{c\frac{s}{\mu}}$   
Pull  $a_2$  w.p.  $\pi_2 \propto e^{c\frac{s}{\mu}}$ 



Compare 
$$\frac{s_1}{q_1} + \sqrt{c \frac{\log(n)}{q_1}} = \frac{s_2}{q_2} + \sqrt{c \frac{\log(n)}{q_2}}$$



- $s_i$ : cumulative reward of  $a_i$
- $q_i$  : number of pulls of  $a_i$

K arms: 
$$\mathscr{A} = \{a_1, \cdots, a_K\}$$
, Horizon n: the total  
round i  $->$  pull  $A^i \in \mathscr{A} ->$  get reward  $X^i \in \mathbb{R}$   
policy  $\pi^i$ :  $H^i \to \Delta(\mathscr{A})$  Environment  $\nu$   
history  $H^i$ :  $(A^1, X^1, \cdots, A^{i-1}, X^{i-1})$   
 $(s_1^i, q_1^i), \cdots (s_K^i, q_K^i) \leftarrow$  Cumulative rewa

Goal: find

max

n: the total number of pulls

For ironment  $\nu \in \mathscr{C}$ : e.q.  $x_k \sim N(\nu_k, 1), \nu = (\nu_1, \dots, \nu_K)$  $\mathscr{C} = \mathbb{R}^K$ , each reality  $\langle - \rangle \nu \in \mathscr{C}$ 

ulative reward, number of pulls

d 
$$\{\pi^i\}_{i=1}^n$$
 to  

$$\mathbb{E}\sum_{i=1}^n X^i$$





- Advert Placement
- Recommendation Service
- Network Routing
- Dynamic Pricing
- Dynamic treatment

# amazon









#### Frequentist:

$$\min R(\nu) = \mathbb{E} \sum_{i=1}^{n} (\mu^{*}(\nu) - X^{i}) \longrightarrow \max \mathbb{E} \sum_{i=1}^{n} X^{i} = \mathbb{E} \sum_{i=1}^{n} \mu_{A^{i}}(\nu)$$
$$\mu^{*}(\nu) = \max\{\mu_{1}(\nu), \cdots, \mu_{K}(\nu)\}, \text{ where } \mu_{k} \text{ is the expectation of arm}$$

### It is not an optimization problem.

$$\max_{\pi} C = \mathbb{E} \sum_{i=1}^{n} \int_{\mathscr{C}} \mu_{A^{i}}(\nu) \rho^{i}(\nu) d\nu$$

- It is an optimization problem
- cumulative reward



• It is proved that a good policy can achieve  $R(\nu) = O(log(n))$  for  $\forall \nu \in \mathscr{E}$ 

 $\rho^{\prime}(\nu)$  is the posterior measure according to the bayesian rule

•The goal is to solve for the best policy that maximizes the expected



- [Bradt et al., 1956] developed the framework of bayesian optimal bandits
- Bayesian optimal bandit dominates the research of bandit problem from 1950 1980
- The biggest difficulty of bayesian framework is the computational cost: e.g. K arms: - exact solve the problem requires  $O(n^K)$  computational time
- One breakthrough is the Gitten's index [Gittins, 1979]: Computational cost:  $O(n^K) \rightarrow O(Kn^2)$ However, it can only be applied to infinite discounted cumulative reward:  $\mathbf{\Omega}$

$$\max C = \mathbb{E} \sum_{i=1}^{\infty} \int_{\mathscr{C}}$$

- For frequentist setting: [Agrawal, 1995], [Katehakis and Robbins, 1995] developed UCB - The computational cost is O(Kn)
- - Asymptotically achieve the optimal regret, but not necessarily optimal bayesian rewards

 $\gamma^{i}\mu_{Ai}(\nu)\rho^{i}(\nu)d\nu$ 



### This talk: Continuous Bayesian bandits

- Computational cost:  $O(n^K) \rightarrow O(Kn)$
- In general, using numerical solution
  - the computational cost:  $O(n^K) \rightarrow O(N^K)$
  - Using DNN to break the curse of dimensionality

Goal: undiscounted bayesian optimal bandit problem  $\max_{\pi} C = \mathbb{E} \sum_{i=1}^{n} \int_{\mathscr{C}} \mu_{A^{i}}(\nu) \rho^{i}(\nu) d\nu$ 

- As  $n \to \infty$ , it converges to a Hamilton-Jacobi-Bellman equation

- Propose a regularized version of bayesian optimal bandits

- Analytic solution exists for some commonly visited cases Depends on mesh

### An illustrative example

 $a_1$  gives reward 1 w  $a_2$  gives determinis

#### **One-armed** bandit problem

- Initially, before round 1: assume  $\nu \sim \text{Beta}(\alpha, \beta)$

-  $a_1$  is chosen w.p.  $\pi^i = \pi(s^i, q^i)$ 

Goal: max  $\mathbb{E}\left[\sum_{i=1}^{n} \left(\int \mu\right)\right]$ 

w.p. 
$$\nu$$
 and 0 w.p.  $1 - \nu \longrightarrow \mu_1(\nu) = \nu$   
stic reward  $\mu_2$ 

- At round i: the posterior measure  $\rho^{i}(\nu)$  depends on -  $s^i$ : the cumulative reward from  $a_1$  up to round i - 1-  $q^i$ : the number of pulls from  $a_1$  up to round i-1-  $\rho^{i}(\nu) = \rho(\nu | s^{i}, q^{i}) \sim \text{Beta}(\alpha + s^{i}, \beta + q^{i} - s^{i})$ 

$$u_1(\nu)\rho^i(\nu)d\nu\right)\pi^i+\mu_2(1-\pi^i)$$

### **Discrete bayesian algorithm**

- If  $a_2$  is chosen at round i:  $w^i(s,q) = w^{i+1}(s,q)$
- If  $a_1$  is chosen at round i:  $w^i(s,q) = p(s,q) + p(s,q)w^{i+1}(s-1)$

where 
$$p(s,q) = \frac{\alpha + s}{\alpha + \beta + q}$$
 is the expense

- after pulling  $a_1$ ,

- 
$$q^{i+1} = q^i + 1$$
  
-  $\mathbb{P}(s^{i+1} = s^i + 1 | s^i) = p(s^i, q^i),$   
-  $\mathbb{P}(s^{i+1} = s^i | s^i) = 1 - p(s^i, q^i).$ 

 $w^{i}(s,q) = \max\left\{w^{i+1}(s,q) + \mu_{2}, \quad p(s,q) + p(s,q)w^{i+1}(s+1,q+1) + (1-p(s,q))w^{i+1}(s,q+1)\right\}$ 

With  $w^{n+1}(s,q) = 0$ , one can solve for  $w^i(s,q)$  backwards

- Let  $w^{i}(s,q)$  be the optimal expected cumulative reward starting from round i with  $(s^{i},q^{i}) = (s,q)$ 

$$(q) + \mu_2$$

$$(+1,q+1) + (1-p(s,q))w^{i+1}(s,q+1),$$

ected reward obtained at round i



### **Derivation to HJB equation**

$$p(s,q) + p(s,q)w^{i+1}(s+1,q+1) + (1-p(s,q))w^{i+1}(s,q+1) \} \qquad w^{n+1}(s,q) = U^{n+1}(s,q) = U$$

$$\max \left\{ w^{i+1}(s,q) + \mu_{2}, \quad p(s,q) + p(s,q)w^{i+1}(s+1,q+1) + (1-p(s,q))w^{i+1}(s,q+1) \right\} \qquad w^{n+1}(s,q) = \\ \text{Let } \hat{s} = \frac{1}{n}s, \hat{q} = \frac{1}{n}q \qquad \int v^{i}(\hat{s},\hat{q}) = \frac{1}{n}w^{i}(s,q) \\ \text{ax } \left\{ \frac{1}{n}\mu_{2} + v^{i+1}(\hat{s},\hat{q}), \frac{1}{n}\tilde{p}(\hat{s},\hat{q}) + \tilde{p}(\hat{s},\hat{q})v^{i+1}(\hat{s} + \frac{1}{n},\hat{q} + \frac{1}{n}) + (1-\tilde{p}(\hat{s},\hat{q}))v^{i+1}(\hat{s},\hat{q} + \frac{1}{n}) \right\}, \quad \text{where } \tilde{p}(s,q) = \frac{n^{-1}}{n^{-1}(\alpha)} \\ \text{By setting } \delta_{i} = \delta_{q} = \delta_{s} = \frac{1}{n}, \\ \frac{-v^{i}(\hat{s},\hat{q})}{\delta_{i}} + \max \left\{ \mu_{2}, \tilde{p}(\hat{s},\hat{q}) + \tilde{p}(\hat{s},\hat{q}) \frac{v^{i+1}(\hat{s} + \delta_{s},\hat{q} + \delta_{q}) - v^{i+1}(\hat{s},\hat{q} + \delta_{q})}{\delta_{s}} - \frac{v^{i+1}(\hat{s},\hat{q} + \delta_{q}) - v^{i+1}(\hat{s},\hat{q} + \delta_{q})}{\delta_{s}} + \frac{v^{i+1}(\hat{s},\hat{q} + \delta_{q}) - v^{i+1}(\hat{s},\hat{q})}{\delta_{q}} \right\} = \\ As n \to \infty \quad \int \delta_{i}, \delta_{q}, \delta_{s} \to 0 \\ \delta_{i}V(t, s, q) + \max \left\{ \mu_{2}, \hat{p}(s, q) + \hat{p}(s, q)\partial_{s}V(t, s, q) + \partial_{q}V(t, s, q) \right\} = 0, \quad V(1, s, q) = 0, \quad \text{where } \hat{p}(s, q) = \lim_{n \to \infty} \tilde{p}(s, q) \\ \downarrow$$

 $\partial_t V(t, s, q) + \max_{\pi \in [0,1]} \left\{ \hat{p}(s,q) + \hat{p}(s,q) \partial_s V(t,s,q) + \partial_q V(t,s,q) - \mu_2 \right\} \pi + \mu_2 = 0, \quad V(1,s,q) = 0,$ 







#### General Results for K-armed bayesian bandits

History at round i:

Policy at round i:  $\pi^i = \pi(\mathbf{s}^i, \mathbf{q}^i)$ , where  $\sum_k \pi_k(s, q) = 1$  $\int x^2 P_k^{\nu} \rho^i(\nu \mid s, q) \, dx \, d\nu$ ,  $E_k^p(s, q) = \int x^p P_k^{\nu} \rho^i(\nu \mid s, q) \, dx \, d\nu$ Posterior measure at round i:  $\rho^{i}(\nu) = \rho(\nu | \mathbf{s}^{i}, \mathbf{q}^{i})$ 

Let 
$$\bar{\mu}_k(s,q) = \int \mu_k(\nu) \rho^i(\nu \mid s,q) d\nu$$
,  $\bar{\sigma}_k^2(s,q) = \int d\nu$ 

Assume  $\hat{t} = \frac{t-1}{----1}$ 

Theorem [Z-Ying-Izzo, 22] If  $\lim_{n \to \infty} \frac{n}{f(n)} \bar{\mu}_k(f(n)\hat{s}, n\hat{p}) = \hat{\mu}_k(\hat{s}, \hat{p}), \lim_{n \to \infty} \frac{n}{f(n)^2} \bar{\sigma}_k^2(p)$ 

then 
$$V(t, \hat{s}, \hat{q}) = V(\frac{i-1}{n}, \frac{s}{f(n)}, \frac{q}{n}) = \frac{1}{f(n)}w^i(s, q)$$

$$\partial_t V + \max_{\sum_k \pi_k = 1} \left\{ \hat{\mu}_k \partial_{s_k} V + \partial_{q_k} V + \frac{1}{2} \hat{\sigma}_k^2 \partial_{s_k}^2 V + \hat{\mu}_k \right\} \pi_k = 0 \quad \text{with } V(1, \hat{s}, \hat{q}) = 0, \text{ for } \forall s, q$$

$$\mathbf{s}^{i} = (s_{k}^{i})_{k=1}^{K}, \quad \mathbf{q}^{i} = (q_{k}^{i})_{k=1}^{K}$$

$$\frac{1}{r}, \quad \hat{q} = \frac{1}{n}q, \quad \hat{s} = \frac{1}{f(n)}s$$

$$f(n)\hat{s}, n\hat{p}) = \hat{\sigma}_k^2(\hat{s}, \hat{p}), \lim_{n \to \infty} \frac{n}{f(n)^p} E_k^p(f(n)\hat{s}, n\hat{p}) = 0$$

satisfies the following HJB equation:

## **Continuous HJB and control problem**

$$\partial_{t}V + \max_{\sum_{k}\pi_{k}=1} \{\hat{\mu}_{k}\partial_{s_{k}}V + \partial_{q_{k}}V + \frac{1}{2}\hat{\sigma}_{k}^{2}\partial_{s_{k}}^{2}V + \hat{\mu}_{k}\}\pi_{k} = 0, \text{ with } V(1,\hat{s},\hat{q}) = 0, \text{ for } \forall s, q$$

$$\int_{V(t,\hat{s},\hat{q})} V(t,\hat{s},\hat{q}) = \max_{\pi} \mathbb{E}\int_{t}^{1} [\hat{\mu}_{k}(\hat{s}(\tau),\hat{q}(\tau))\pi_{k}(\hat{s}(\tau),\hat{q}(\tau))]d\tau$$

$$\max_{\pi} \mathbb{E} \int_{0}^{1} \sum_{k} \hat{\mu}_{k}(\hat{\mathbf{s}}, \hat{\mathbf{q}}) \pi_{k}(\hat{\mathbf{s}}, \hat{\mathbf{q}}) dt$$
  
s.t.  $d\hat{s}_{k}(t) = \hat{\mu}_{k}(\hat{\mathbf{s}}, \hat{\mathbf{q}}) \pi_{k}(\hat{\mathbf{s}}, \hat{\mathbf{q}}) dt + \hat{\sigma}_{k}(\hat{\mathbf{s}}, \hat{\mathbf{q}}) \sqrt{\pi_{k}(\hat{\mathbf{s}}, \hat{\mathbf{q}})} dB_{t}$   
 $d\hat{q}_{k}(t) = \pi_{k}(\hat{s}, \hat{q}) dt$   
with  $\hat{\mathbf{s}}(0) = \mathbf{0}, \quad \hat{\mathbf{q}}(0) = \mathbf{0}.$ 

# **Regularized HJB equation**

$$\max_{\pi} \quad \mathbb{E} \int_{0}^{1} \sum_{k} \hat{\mu}_{k}(\hat{s}, \hat{q}) \pi_{k}(\hat{s}, \hat{q}) - \lambda \pi_{k} \log(\pi_{k}) dt$$
  
s.t. 
$$d\hat{s}_{k}(t) = \hat{\mu}_{k}(\hat{s}, \hat{q}) \pi_{k}(\hat{s}, \hat{q}) dt + \hat{\sigma}(\hat{s}, \hat{q}) \sqrt{\pi_{k}(x)}$$
$$d\hat{q}_{k}(t) = \pi_{k}(\hat{s}, \hat{q}) dt$$
  
with 
$$\hat{s}(0) = 0, \quad \hat{q} = 0$$

$$\begin{split} \partial_t V + \lambda \log \left( \sum_k \exp\left[\frac{1}{\lambda} H_k\left(\partial_{s_k} V, \partial_{q_k} V, \partial_{s_k}^2 V\right)\right] \right) &= 0, \text{ with } V(1, \hat{s}, \hat{q}) = 0, \text{ for } \forall s, q \\ \pi_k \propto \exp\left[\frac{1}{\lambda} H_k\left(\partial_{s_k} V, \partial_{q_k} V, \partial_{s_k}^2 V\right)\right] \\ \end{split}$$

$$\end{split}$$

$$Where H_k\left(p, g, h\right) = \hat{\mu}_k(s, q)p + g + \frac{1}{2}\hat{\sigma}_k^2(s, q)h + \hat{\mu}_k(s, q)$$

- $\overline{(\hat{s},\hat{q})}dB_t$



Bernoulli rewards:

 $a_k$  gives reward  $\gamma$  w.  $\mathscr{E} = [0,1]^K$ Prior of  $\nu_k \sim \text{Beta}(\alpha$ 

For  $\gamma = O(n^{-p})$ , if  $\exists f(n) \leq O(n^{1/2-p})$ , s.

 $\longrightarrow \hat{\mu}_k(s,q) = -$ 

Normal rewards:

The reward of  $a_k$  follows  $\mathscr{E} = \mathbb{R}^{K}$ Prior of  $\nu_{k} \sim N(\alpha_{k}, \beta_{k})$ 

For  $\sigma = O(n^{-p})$ , if  $\exists f(n) \leq O(n^{1/2-p})$ , s.

 $\longrightarrow \hat{\mu}_k(s,q) =$ 

#### **Two special cases**

p. 
$$\nu_k$$
 and  $-\gamma$  w.p.  $1 - \nu_k$ , ( $\gamma$  is known)  
 $\frac{\alpha_k, \beta_k}{p_k}$   
t.  $\lim_{n \to \infty} \frac{\gamma(\alpha_k - \beta_k)}{f(n)} = \hat{\alpha}_k$ ,  $\lim_{n \to \infty} \frac{\alpha_k + \beta_k}{n} = \hat{\beta}_k$ ,  
 $\frac{s + \hat{\alpha}_k}{q + \hat{\beta}_k}$ ,  $\hat{\sigma}_k \equiv \frac{n^{1/2}}{f(n)}\gamma$ .  
Dow  $N(\nu_k, \sigma^2)$  with known  $\sigma$   
 $\frac{\alpha_k^2}{p_k^2}$   
t.  $\lim_{n \to \infty} \frac{\alpha_k \beta_k^{-2}}{f(n)} = \hat{\alpha}_k$ ,  $\lim_{n \to \infty} \frac{\sigma^2 \beta_k^{-2}}{n} = \hat{\beta}_k$ ,  
 $\frac{s + \hat{\alpha}_k}{q + \hat{\beta}_k}$ ,  $\hat{\sigma}_k \equiv \frac{n^{1/2}}{f(n)}\sigma$ .

Theorem [Z-Ying-Izzo, 22]  
when 
$$\hat{\mu}(s, q) = \frac{s + \hat{\alpha}_k}{q + \hat{\beta}_k}$$
, then the optimal policy for the unreg  
 $\pi_k^*(s, q) = \begin{cases} 1, & k = \operatorname{argmax}_k \left\{ \mu(s_k, q_k) \right\} \\ 0, & o . w . \end{cases}$ 

For finite horizon one-armed bandits problem:

When 
$$\mu\left(\frac{s}{f(n)}, \frac{q}{n}\right) > \hat{\mu}_2, \pi = 1$$

gularized HJB is



### Analytic solution for the regularized HJB

Theorem [Z-Ying-Izzo, 22]  
when 
$$\hat{\mu}_k(s,q) = \frac{s + \hat{\alpha}_k}{q + \hat{\beta}_k}$$
, then the optimal policy  $\pi_k^*(s,q) \propto e^{\mu_k(s_k,q_k)}$ 

For finite horizon one-armed bandits problem:

When 
$$\mu\left(\frac{s}{f(n)}, \frac{q}{n}\right) > \hat{\mu}_2, \pi = 1 \iff \pi \propto \exp\left(\frac{n}{f(n)}\frac{s_i + f(n)\hat{\alpha}}{q_i + n\hat{\beta}}\right)$$
  
Compare with tempered greedy algo:  $\pi \propto \exp\left(c\frac{s_i + a}{q_i + b}\right)$ 

/ for the regularized HJB is



#### $a_k$ gives rewa

#### K = 3, $\nu_1 = 1/2$ , $\nu_2 = \nu_3 = \nu \in [0,1]$



## **Numerical Experiments**

ard 
$$\begin{bmatrix} 1 \text{ w.p. } \nu_k; \\ -1 \text{ w.p. } 1 - \nu_k \end{bmatrix}$$

# K = 10, $\nu_1 = 1/2$ , $\nu_2 = \dots = \nu_{10} = \nu \in [0,1]$







# K = 3, $\nu_1 = 0$ , $\nu_2 = \nu_3 = \nu \in [-1,1]$



## **Numerical Experiments**

eward ~ 
$$N(\nu_k, 1)$$

## K = 10, $\nu_1 = 0$ , $\nu_2 = \dots = \nu_{10} = \nu \in [-1,1]$





- efficiently for large K. (DNN?)
- Convergence rate to the HJB equation:  $\frac{1}{n} w^{i}(s,q) \to v(t,\hat{s},\hat{q})$

- Discrete: 
$$h(n, s, q) > \mu_2 \rightarrow \pi =$$
  
-  $h(n, s, q) \rightarrow g\left(\frac{s}{f(n)}, \frac{q}{n}\right)$ 

If  $\mu(s,q) = \frac{s+\alpha}{a}$ , is the HJB equation well-defined?

- Will the HJB equation converge to some "mean-field limit" as  $K \to \infty$
- Does there exist a PDE for UCB algorithm as  $n \to \infty$ ?



When the exact solution is hard to calculate, how to compute the numerical solution

# = 1; Continuous: $g(\hat{s}, \hat{q}) > \mu_2 \rightarrow \pi = 1$