

# Bayesian Bandit:

- Bayesian measure: using the observation data to establish the belief of the hyperparameters of a distribution.

e.g.

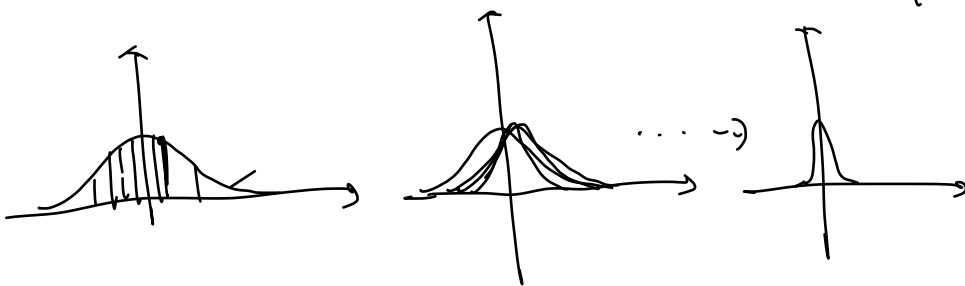
$$X \sim N(\mu, \sigma)$$

↑  
unknown.

observation:  $x_1, x_2, \dots, x_n$

Goal: Based on the observation, derive a belief for  $\mu$ .

↑  
 $\mu \sim p(\mu)$



the prob of obs is  $x$  for  $\mu$

$$p(\mu | x) = \frac{P(x|\mu)}{P(x)} p(\mu) = \frac{1}{Z} \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(\mu-x_i)^2}{2\sigma_i^2}} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mu-\mu)^2}{2\sigma^2}}$$

normalizing constant

a new prob measure for  $\mu$

$$p(x|\mu) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$p(\mu) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(\mu-\mu)^2}{2\sigma^2}}$$

At the beginning of round  $i$ ,  
we have prior measure

$$\mu \sim N(\mu_{i-1}, \sigma_{i-1}^2)$$

After one observation  $x_i$ ,

we update the posterior measure  $\mu \sim N(\mu_i, \sigma_i^2)$

$$\mu_i = \left( \frac{1}{\sigma^2} + \frac{1}{\sigma_{i-1}^2} \right)^{-1} \left[ \frac{x_i}{\sigma^2} + \frac{\mu_{i-1}}{\sigma_{i-1}^2} \right]$$

$$\sigma_i^2 = \left( \frac{1}{\sigma^2} + \frac{1}{\sigma_{i-1}^2} \right)^{-1}$$

or equivalently, before any observations, with prior measure  $N(\mu_0, \sigma_0^2)$   
 & after  $n$  observations  $X_1, \dots, X_n$ , we have empirical mean  $\bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$ ,  
 then the posterior measure of  $\mu \sim N(\mu_n, \sigma_n^2)$  will be

$$\mu_n = \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)^{-1} \left[ \frac{n}{\sigma^2} \bar{x} + \frac{1}{\sigma_0^2} \mu_0 \right]$$

$$\sigma_n^2 = \left( \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right)^{-1}$$

- Bayesian for Binomial distribution  $\begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p. \end{cases}$

$$p \sim \text{Beta}(a_0, b_0), \quad \rightarrow \mathbb{E}[p] = \frac{a_0}{a_0 + b_0}$$

prior

After  $n$  observations, we can update the posterior measure for  $p$ :  
 $p \sim \text{Beta}(a_n, b_n) \leftarrow \text{posterior.}$

$$\begin{cases} a = a_0 + \sum_{j=1}^n X_j \\ b = n - \sum_{j=1}^n X_j + b_0 \end{cases} \Rightarrow \mathbb{E}[p] = \frac{a_0 + \sum_{j=1}^n X_j}{n + a_0 + b_0}$$

or equivalently

$\rightarrow$  with prior  $p \sim \text{Beta}(a_{i-1}, b_{i-1})$  & observation  $X_i$

$$\text{posterior } p \sim \text{Beta}(a_{i-1} + X_i, 1 - X_i + b_{i-1})$$

Thompson Sampling → I-armed reward.

e.g.  $X_0 = \frac{1}{2}$

$X_t \sim \text{Binomial}(p)$       ↓ unknown

start with  $p \sim \text{Beta}(1, 1)$

At the beginning of round  $t$ , we have prior measure of  $p \sim \text{Beta}(a_{t-1}, b_{t-1}) \rightarrow p_t$

• If  $p_t > \frac{1}{2} \Rightarrow$  pull unknown arm

$\Rightarrow$  update the posterior measure of  $p$  based on

the observation  $x_t$ .

$$a_t = a_{t-1} + x_t$$

$$b_t = (-x_t + b_{t-1})$$

• If  $p_t < \frac{1}{2} \Rightarrow$  pull known arm

$$\Rightarrow \begin{cases} a_t = a_{t-1} \\ b_t = b_{t-1} \end{cases}$$

HW  $\rightarrow$  Try TS for Normal reward & Binomial reward with different  $\Delta$  & differe prior.

Can we do better than TS based on the posterior measure if we know how many rounds we left?

• Bayesian optimal policy

$$X_0 = \frac{1}{2}$$

$x_i \sim \text{Binomial}(p)$  ↙ unknown

Q: At the beginning of round  $n$ , I have prior measure for  $p \sim \text{Beta}(a_{n-1}, b_{n-1})$ , what is the best policy according to this posterior?

$$\mathbb{E}[X_n] = \mathbb{E}[p] = \frac{a_{n-1}}{b_{n-1} + a_{n-1}} \quad 0 \quad \frac{1}{2}$$

" > "      arm 1

" < "      arm 0 .

e.g.  $a_{n-1} = b_{n-1} = 1$ , then it is the same to pull unknown arm or known arm

Define  $V(a, b, i)$  as the optimal expected cumulative reward with  $n$  rounds left.

$$V(a, b, i) = \max \left\{ \frac{a_{n-1}}{b_{n-1} + a_{n-1}}, \frac{1}{2} \right\} .$$

Q: At the beginning of round  $n-1$ , I have prior measure for  $p \sim \text{Beta}(a_{n-2}, b_{n-2})$ , what is the best policy according to this posterior?

if arm 0 is pulled

$$V^2(a_{n-2}, b_{n-2}, 2) = \frac{1}{2} + V(a_{n-2}, b_{n-2}, 1)$$

if arm 1 is pulled:

$$V^1(a_{n-2}, b_{n-2}, 2) = \frac{a_{n-2}}{b_{n-2} + a_{n-2}} + \mathbb{E}[V(\underbrace{a_{n-1}, b_{n-1}}_{\substack{\uparrow \\ \text{the posterior measure of } p \text{ after} \\ \text{pulling arm 1 at round } n-1}}, 1)]$$

if arm 1 is pulled at round  $n-1$ :

$\left. \begin{array}{l} 1 \\ 0 \end{array} \right\}$  w.p.  $p \rightarrow$  posterior  $p \sim \text{Beta}(a_{n-2}+1, b_{n-2})$   
w.p.  $1-p \rightarrow$  posterior  $p \sim \text{Beta}(a_{n-2}, b_{n-2}+1)$

$$V^1(a_{n-2}, b_{n-2}, 2) = \frac{a_{n-2}}{b_{n-2} + a_{n-2}} + \frac{a_{n-2}}{b_{n-2} + a_{n-2}} \cdot V(a_{n-2}+1, b_{n-2}, 1) \\ + \frac{b_{n-2}}{b_{n-2} + a_{n-2}} \cdot V(a_{n-2}, b_{n-2}+1, 1)$$

$$\underline{V(a, b, 2) = \max \{ V^1, V^2 \}}$$

e.g. if  $a_{n-2} = b_{n-2} = 1$

$$\begin{aligned}V^1(a, b, 2) &= \frac{1}{2} + \frac{1}{2} V(2, 1, 1) + \frac{1}{2} V(1, 2, 1) \\&= \frac{1}{2} + \frac{1}{2} (\max\{\frac{2}{3}, \frac{1}{2}\} + \max\{\frac{1}{3}, \frac{1}{2}\}) \\&= \frac{1}{2} + \frac{1}{2} (\frac{2}{3} + \frac{1}{2}) = \frac{1}{2} + \frac{1}{2} \frac{4+3}{6} = \frac{1}{2} + \frac{7}{12} = \frac{13}{12}\end{aligned}$$

$$V^2(a, b, 2) = \frac{1}{2} + V(1, 1, 1) = \frac{1}{2} + \max\{\frac{1}{2}, \frac{1}{2}\} = 1$$

if 2 rounds left, even if the expectation is the same for 2 arms, it will automatically prefer to explore the unknown arm.

In general:  $V(a, b, n)$

$$= \max \left\{ \frac{1}{2} + V(a, b, n-1), \frac{a}{a+b} + \frac{a}{a+b} V(a+1, b, n-1) + \frac{b}{a+b} V(a, b+1, n-1) \right\}$$

with  $V(a, b, 0) = 0$

Q: What's the complexity to obtain the optimal Bayesian policy for MAB with horizon  $n$ ?

$$O(n^3)$$

If it is  $k$ -armed bandit prob,  $\Rightarrow O(n^{2k})$ .

Q How to reduce the computational cost

1. Derive a limiting HJB eqn  $\Rightarrow$  independent of  $n$

"Continuous-in-time limit for Bayesian Bandits. 2-1220 - Ying JomLR 23"

$$t = \frac{i}{n}, \hat{a} = \frac{1}{n}a, \hat{b} = \frac{1}{n}b, \hat{V}(\hat{a}, \hat{b}, t) = \frac{1}{n}V(a, b, i), \quad a = n\hat{a}, \quad b = n\hat{b}$$

$$a+1 = n\hat{a}+1 = n(\hat{a} + \frac{1}{n})$$

$$\hat{V}(\hat{a}, \hat{b}, t) = \max \left\{ \frac{1}{2} \frac{1}{n} + \hat{V}(\hat{a}, \hat{b}, t - \frac{1}{n}), \frac{1}{n} \frac{\hat{a}}{\hat{a} + \hat{b}} + \frac{\hat{a}}{\hat{a} + \hat{b}} \hat{V}(\hat{a} + \frac{1}{n}, \hat{b}, t + \frac{1}{n}) + \frac{\hat{b}}{\hat{a} + \hat{b}} \hat{V}(\hat{a}, \hat{b} + \frac{1}{n}, t + \frac{1}{n}) \right\}$$

$$\frac{\hat{V}(\hat{a}, \hat{b}, t) - \hat{V}(\hat{a}, \hat{b}, t - \frac{1}{n})}{\frac{1}{n}} = \max \left\{ \frac{1}{2}, \frac{\hat{a}}{\hat{a} + \hat{b}} + \frac{1}{n} \frac{\hat{a}}{\hat{a} + \hat{b}} \left[ \hat{V}(\hat{a} + \frac{1}{n}, \hat{b}, t + \frac{1}{n}) - \hat{V}(\hat{a}, \hat{b}, t + \frac{1}{n}) \right] + \frac{1}{n} \frac{\hat{b}}{\hat{a} + \hat{b}} \left[ \hat{V}(\hat{a}, \hat{b} + \frac{1}{n}, t + \frac{1}{n}) - \hat{V}(\hat{a}, \hat{b}, t + \frac{1}{n}) \right] \right\}$$

$$\frac{\hat{V}(\hat{a}, \hat{b}, t) - \hat{V}(\hat{a}, \hat{b}, t - \frac{1}{n})}{\delta t} = \frac{1}{2} + \max \left\{ 0, \frac{\hat{a}}{\hat{a} + \hat{b}} - \frac{1}{2} + \frac{\hat{a}}{\hat{a} + \hat{b}} \frac{\hat{V}(\hat{a} + \frac{1}{n}, \hat{b}, t + \frac{1}{n}) - \hat{V}(\hat{a}, \hat{b}, t + \frac{1}{n})}{\delta a} + \frac{\hat{b}}{\hat{a} + \hat{b}} \frac{\hat{V}(\hat{a}, \hat{b} + \frac{1}{n}, t + \frac{1}{n}) - \hat{V}(\hat{a}, \hat{b}, t + \frac{1}{n})}{\delta b} \right\}$$

Define function  $\hat{V}(a, b, t)$  with  $t \in (0, 1)$ , then.  
 as  $\delta t, \delta a, \delta b \rightarrow 0$

$$\delta_t \hat{V} = \frac{1}{2} + \max_{\pi \in [0, 1]} \left\{ 0, p(\hat{a}, \hat{b}) - \frac{1}{2} + p(\hat{a}, \hat{b}) \delta_a V + (1 - p(\hat{a}, \hat{b})) \delta_b V \right\}$$

$$\delta_t \hat{V} = \frac{1}{2} + \max_{\pi \in [0, 1]} \left( p(\hat{a}, \hat{b}) - \frac{1}{2} + \hat{p}(\hat{a}, \hat{b}) \delta_a V + (1 - p(\hat{a}, \hat{b})) \delta_b V \right) \pi(\hat{a}, \hat{b}, t)$$

exact solution or numerical solution (indep of  $n$ ).

open prob:  $\textcircled{1}$  what theoretical guarantees can we have for this approximated soln?

$\textcircled{2}$  for  $\epsilon$ -armed bandit with no exact solution for HJB eqn, any method to find the soln efficiently?

$$p = \frac{a + b_0}{b + b_0}$$





## Infinite horizon with discount factor $\gamma$

$$V(a,b) = \max \left\{ \lambda + \gamma V(a,b), \frac{a}{a+b} + \gamma \left[ \frac{a}{a+b} V(a+1,b) + \frac{b}{a+b} V(a,b+1) \right] \right\} \quad (\infty)$$

$$S \in \mathbb{Z}_+^2; A = \{1,0\}; S_1^{A=1} = \begin{cases} \binom{\alpha+1}{\beta} & \text{w.p. } \frac{\alpha}{\alpha+\beta} = P(S_1=1 | S_0=A=1) \\ \binom{\alpha}{\beta+1} & \text{w.p. } \frac{\beta}{\alpha+\beta} = P(S_1=0 | S_0=1) \end{cases}, \quad r(1,S) = \frac{\alpha}{\alpha+\beta}$$

$$S_1^{A=0} = S_0, \quad r(0,S) = \lambda$$

$$V(S) = \max_a \left\{ r(0,S) + \gamma V(S), r(1,S) + \gamma \left[ P(S_1=1 | \dots) V(\dots) + P(S_1=0 | \dots) V(\dots) \right] \right\}$$

$$V(S) = \max_a \left\{ r(S,a) + \gamma \mathbb{E}[V(S_1) | S_0=S, A_0=a] \right\}$$

$$= \max_{\substack{\pi_a \in [0,1] \\ \sum_a \pi_a = 1}} \left[ \sum_{a \in A} (r(S,a) + \gamma \mathbb{E}[V(S_1) | S_0=S, A_0=a]) \pi_a \right]$$

$$= \max_{\pi} \mathbb{E}_{\substack{a \sim \pi(a|S) \\ S_1 \sim P(S_1|a_0, S_0)}} \left[ r(S,a) + \gamma V(S_1) \mid S_0=S \right]$$

for fixed  $S$ , a prob distribution for  $a$ .

Naively, one can cut the infinite horizon to finite horizon & then

derive the soln backwards, but as  $\gamma \rightarrow 1$ , the horizon will be larger.

open prob: (3) Can we derive a similar limiting PDE by rescaling the parameters

k-armed prob :

$$V(\alpha_1, \beta_1, \dots, \alpha_k, \beta_k) = \max_j \left\{ \frac{\alpha_j}{\alpha_j + \beta_j} + V(\dots, \alpha_{j+1}, \beta_j, \dots) \frac{\alpha_j}{\beta_j + \alpha_j} + V(\dots, \alpha_j, \beta_{j+1}, \dots) \frac{\beta_j}{\alpha_j + \beta_j} \right\}$$

$|r| < 1 \Rightarrow \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \leq \frac{\gamma^T}{1-\gamma} \leq \epsilon \Rightarrow T \sim O\left(\frac{1}{1-\gamma}\right)$   
 $O\left(\left(\frac{1}{1-\gamma}\right)^{2k}\right)$

Gittins index :

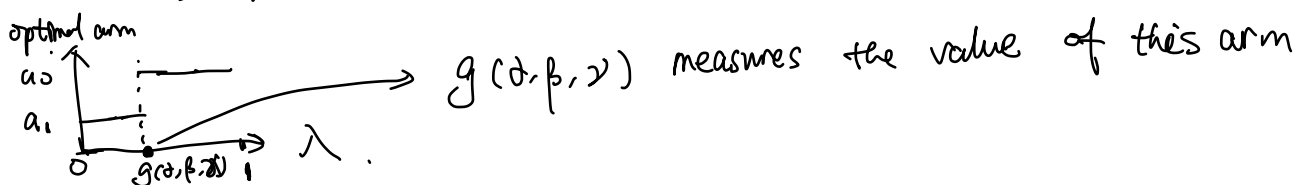
- when  $S \in \mathbb{Z}_+^{2k}$ , instead of viewing it as a coupled system, one decouple it by giving each arm an equivalent value at their current state  
 Gittins index

Let  $g(\alpha, \beta, \gamma)$  be the Gittins index for Binomial reward with prior measure Beta  $(\alpha, \beta)$  & discount factor  $\gamma$ ,

$$\text{optimal policy} = \underset{k \in [K]}{\text{argmax}} \left\{ g(\alpha_k, \beta_k, \gamma) \right\}$$

- View k-armed Bandit as k one-armed bandit prob.

For one-arm bandit prob with prior measure Beta  $(\alpha, \beta)$   
 U.S. known arm with deterministic reward  $\lambda$ .



The Gittins index is the deterministic arm reward such that it is equivalent to pull the known arm & unknown arm

One algo to calculate  $g(s)$

Initialization: set  $\bar{g}$  &  $g$ .

while  $\bar{g} - g > \epsilon \Rightarrow$  set  $g = \frac{\bar{g} + g}{2}$

open prob ①  
may reduce the  
cost of this step.

solve the optimal policy for one-armed bandit with an unknown arm with state  $s$  & known arm with  $r = g$

If the optimal policy is to pull the unknown arm

$g = g$  — — — — — known arm  
 $\bar{g} = g$

② open prob ② For a couple HTB, can we decouple it with similar idea?