

# Bandits

## Stochastic

- play action  $a_t$  (choice/decision)

observe reward  $r_t = r(a_t)$

$$E[r_t] = \mu(a_t) \leftarrow r_t = \mu(a_t) + \eta_t \leftarrow$$

$r, \mu$  - bounded  $[0,1]$  or subgaussian

concentration  $\hat{\mu} \rightarrow \mu$

hoeffding / Bernstein - finite, Lipschitz }  $\mu$   
 " martingale - linear, GB

WP  $\geq 1 - \delta$

$$R(T) = \sum_{t=1}^T (\mu(a^*) - \mu(a_t)) \leq \epsilon(\delta, T, \mathcal{A})$$

$\mathcal{A}$  action space

$$E[R(T)] = O(\sqrt{KT \log T})$$

$$\approx O(d \sqrt{T \log T}), \quad T^{\frac{d+1}{d+2}} (KT)^{\frac{1}{d+2}}$$

$K$  - finite  
 $d$  - linear, Lipschitz

$d, K, \nu$  - dim + kernel hyperparameters

## Algorithms:

Explore then exploit

$\epsilon$ -greedy

UCB (Upper Confidence Bound)  $|\mu - \hat{\mu}| \leq \epsilon$

Thompson

→ • Cumulative Regret  $\sum_{t=1}^T \mu(a^*) - \mu(a_t) = R(T)$

• Simple Regret  $\mu(a^*) - \mu(\hat{a}_T) = S(T)$   $\hat{a}_T$  - recommended action after  $T$  rounds

can get  $E[S(T)] = \frac{E[R(T)]}{T}$  at time  $T$  let  $\hat{a}_T = a$  w.p.  $\frac{n(a, T)}{T}$

$$E[R(T)] = E\left[\sum_t \mu(a^*) - \mu(a_t)\right] = \sum_a (\mu(a^*) - \mu(a)) E[n_a(T)]$$

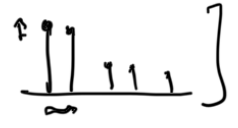
$$E[S(T)] = E[\mu(a^*) - \mu(\hat{a}_T)] = \sum_a (\mu(a^*) - \mu(a)) E\left[\frac{n(a, T)}{T}\right]$$

Conversely,

$$E[S(T)] \geq \min_{i^*} \Delta_i e^{-c E[R(T)]}$$

for all Bernoulli bands.  $\mu_i - \mu_{i^*}$

• Best arm identification  $P(\hat{a}_T \neq a^*) = e(T)$



→  $\min_{i \neq i^*} \Delta_i \cdot e(T) \leq E[S(T)] \leq e(T)$  if  $\mu \in [0,1]$

fixed confidence

given  $\delta$   $\min_{\tau} E[\tau]$  st.  $P(\hat{a}_\tau \neq a^*) \geq 1 - \delta$

↳ stopping time (random) "anytime" valid

Fixed budget

given  $T$ ,  $\min P(\hat{a}_T \neq a^*)$

• Reward estimation  $\hat{\mu}(i)$

• Average treatment effect (ATE)  $K=2$   
 $ATE = \mu(1) - \mu(0)$

Randomized control trials

Fixed allocation  $\pi = P(a=1)$

1) Inverse Probability Weighting (IPW) ✓  
 Horvitz-Thompson estimator

⇒  $\hat{\mu}(1) = \frac{1}{T} \sum_{t=1}^T r_t \frac{1_{a_t=1}}{\pi}$        $\hat{\mu}(0) = \frac{1}{T} \sum_{t=1}^T r_t \frac{1_{a_t=0}}{1-\pi}$

→  $ATE \hat{=} \frac{1}{T} \sum_{t=1}^T r_t \left( \frac{1_{a_t=1}}{\pi} - \frac{1_{a_t=0}}{1-\pi} \right)$

$E[\hat{\mu}(i)] = \mu(i)$  unbiased. ←

var  $\frac{1}{T} \left( \frac{E[r^2(1)]}{\pi} + \frac{E[r^2(0)]}{1-\pi} - \Delta^2 \right)$        $\Delta = \mu(1) - \mu(0)$

$\min_{\pi} \text{var} ? \equiv \min_{\pi} \text{MSE}$

$\pi$ 

Neyman policy

$$\pi_{\text{Ney}} = \frac{m_1}{m_1 + m_0}$$

$$= \frac{1}{2} \text{ iff } m_1 = m_0$$

$$\left. \frac{\hat{m}_1}{\hat{m}_1 + \hat{m}_0} \right\} \leftarrow$$

clipped  $0 < \hat{m}_0, \hat{m}_1 < 1$ 

② Direct method

$$\hat{\mu} \leftarrow \min_{f \in \mathcal{F}} \sum_i (f(i) - \mu(i))^2 \quad \text{-- misspecification}$$

③

AIPW (unbiased even if  $\hat{\mu}$  is biased / misspecified)

$$\frac{1}{T} \sum_{t=1}^T \left( \frac{1_{a_t=1}}{\pi} (r_t - \hat{\mu}_t(1)) + \hat{\mu}_T(1) - \frac{1_{a_t=0}}{1-\pi} (r_t - \hat{\mu}_t(0)) - \hat{\mu}_T(0) \right)$$

$E[\hat{\mu}] \neq \mu$

unbiased

$$\text{var} \quad \frac{\sigma^2(1)}{\pi} + \frac{\sigma^2(0)}{1-\pi} + E[(\mu(1) - \mu(0) - \Delta)^2]$$

$$\pi_{\text{AIPW}}^* = \frac{\sigma(1)}{\sigma(1) + \sigma(0)}$$

how to estimate in finite sample setting.  
dependence.Take away: so far considered bounded or equal var actions  
need to account for different variances (std dev).Adversarial setting- play action  $a_t$ observe reward  $r_t(a_t)$  bandit /  $\{r_t(a)\}_{a \in A}$  full feedbackreward can change at each step  $t$  for any  $a$ .

$$R(T) = \sum_t [r_t(a^*) - r_t(a_t)]$$

$$a^* = \arg \min_a \sum_t r_t(a)$$

Pseudo regret

$$a^* = \arg \min_a \sum_t E[r_t(a)]$$

Algos.:  $\left\{ \begin{array}{l} \text{Weighted Maj Vote (rewards binary)} \\ \text{Hedge} \\ \text{Exp(4)} \end{array} \right\}$  full feedback  
 - bandit

Sample arms/experts acc to  $w_t(a) / w_t(e)$   
 multiplicative/exp update  $w_{t+1}(a) = w_t(a) (1 - \epsilon)^{c_t(a)}$   
 $= e^{-\epsilon c_t(a)}$   
 $e^{-\epsilon} \approx 1 + \epsilon + \frac{\epsilon^2}{2} + \dots$

$$E[R(T)] \approx \sqrt{\underbrace{\sum_t c_t(a^*)}_{KT} \ln K}$$

$\downarrow$   
 $a^*$   
 $\downarrow$   
 $N$

Contextual bandits - middle ground.

env choose context  $x_t$  visible to player

play action  $a_t$

observe reward  $r_t = r(a_t, x_t)$   $E[r_t] = \mu(a_t | x_t)$

$$R(T) = \sum_{t=1}^T \left[ \underbrace{\max_a \mu(a | x_t)}_{\pi^*(x_t)} - \mu(a_t | x_t) \right]$$

$$E[R(T)] = O(\sqrt{d |X| T \log T})$$

$$= O(d \sqrt{T \log T})$$

