# Model Selection in Bandits

## Motivation

- MAB ($k$ arms)

- CAB (infinite action/arm set)

$$\begin{cases} \text{linear bandits} \\ \text{kernelized} \quad - \\ GP \quad \sim \\ \text{Lipschitz} \sim \end{cases}$$

Goal :

reward function

$f(x)$



$x_1 \qquad x_3 \quad x_2$

$$R(T) = \sum_{t=1}^{T} [ f(x^*) - f(x_t) ]$$

pseudo regret

$$\mathbb{E}[R(T)] = \mathbb{E}\left[ \sum_{t=1}^{T} f(x^*) - f(x_t) \right]$$

expected regret

Assumption : $f \in \mathcal{F}$

the environment (function space)
is exactly known

E.g. ① linear $f(x)$
$= \langle \theta, x \rangle$
we know. & $\theta \in \mathbb{R}^d$

② GP- bandit

$$f(x) \sim GP\left(0, \Sigma(\theta)\right)$$

GP- UCB : at each step $t$

$$x_t \in \arg\max_x \underbrace{\mu_{t-1}(x)}_{} + \sqrt{\beta}\, \sigma_t(x)$$

$$\sigma_t(x)$$

$$= K(x, x) - k_t(x)^\top \left(K_t + \sigma^2 I\right)^2$$

$$k_t(x)$$

Question :

What happens if these parameter(s)

is unknown?

$\begin{cases} ① \text{ can we do } \sim \text{ if params are known} \end{cases}$

② if not, what is the best we
can do (what algorithms)

---

Why is model selection "hard" for bandits

* exploration - exploitation

* lack of knowledge of parameters $\mathcal{F}$ (key of) creates a "hard" problem

---

Model selection in bandits:

Goal: select the best "suitable" algorithm among $M$ candidates : $B^* \in \{B_1 \cdots B_m\}$

"competitive" with $B^*$

---

Methods

A. "CORRAL"  ⎰ CORRAL
            ⎱ "smooth" CORRAL

RBBE  "regret bound balancing & elimination"
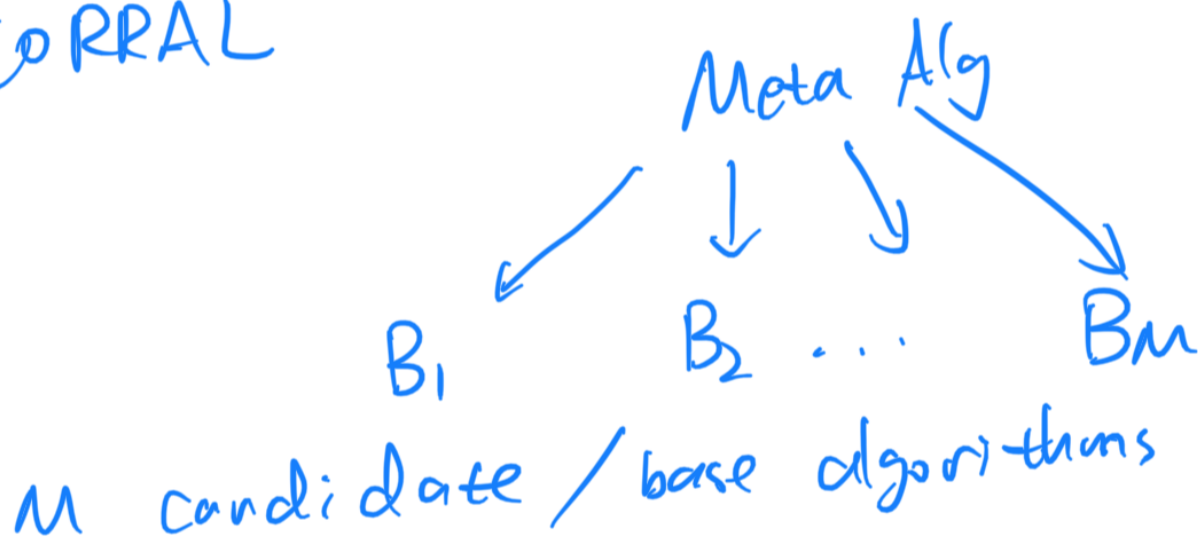* stochastic only
* better than Mod CB

B : "test-based"

- relies on statistical test
  to check whether a base algorithm
  is "misspecified"

$\left\{ \begin{array}{l} \text{Mod CB} \quad \text{* linear models} \\ \qquad\qquad \text{* not optimal} \\ \text{"model selection for linear contextual} \\ \quad \text{bandit" Foster et al. 2019} \\ \left[ d_1 \leq d_2 \leq d_3 < \ldots d_M \right] \end{array} \right.$

$$d*$$

---

CORRAL

Meta Alg



$B_1 \qquad B_2 \quad \ldots \qquad B_M$

$M$ candidate / base algorithms

Meta alg: adversarial bandit algorithm

---

Preliminary: adversarial bandits

Stochastic

$[P_1 \quad \ldots \quad P_K]$

| adversarial
no such assumptions on
static reward distributions

Compare to:
"the best arm in hindsight"

$$\left. \text{Reg}(T) \right|$$

$$= \mathbb{E}\left[ \max_{i \in [K]} \sum_{t=1}^{T} r_{t,i} - \sum r_{t,x_t} \right]$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{\text{randomness in the Alg}}$$

Exp 3: "Exponential weight Alg for
$$\underline{\phantom{xxxxxxxxxxxxxxxxxxxxxx}}$$
exploration - exploitation

① sample $X_t \sim P_{\cdot t}$

② receive reward $r_{t, x_t}$

③ estimates reward for all arms
based on $r_{t, x_t}$

④ Updates $P_{t,i} \to P_{t+1,i} \quad \forall i \in [K]$

$$\hat{r}_{t,i} = \frac{\mathbb{I}\{X_t = i\}}{P_{t,i}} \, r_{t, x_t}$$

$$\hat{S}_{t,i} = \sum_{s=1}^{t} \hat{r}_{s,i} \quad \forall i \in [K] \qquad \text{* exponential weights}$$

Then: $\boxed{P_{t+1,i} = \dfrac{\exp\left(\eta \, \hat{S}_{t,i}\right)}{\sum_{j} \exp\left( - \right)}}$

Regret: $\mathbb{E}[\text{Reg}(T)] = O\left(\sqrt{KT \log(K)}\right)$

$$\left(\eta \propto \sqrt{\log(K)/(TK)}\right)$$

---

Problem with Exp 3: (as Meta Alg)

"exponential" weight can shrink too small

and "starve" a base algorithm.

---

CORRAL (different from Exp 3)

linear weights instead of exponential

"less extreme"

$$P_{t+1, i} \propto \left(-\eta \sum_{s=1}^{t} \hat{r}_{s,i} + \underbrace{Z}_{}\right)^{-1}$$

normalization

- - - - - - - - - - - - - - - - -

Input: $M$ base Algo $\{B_1, \ldots B_M\}$

initial $[\Gamma \quad \eta]$ (adaptive $\eta$)

· initialize all base algorithms

· $\gamma = 1/T$ (lower bound on sampling prob)

$P_{1, j} = 1/M \quad \forall j = 1, \ldots M$

$$\overline{P}_{1,j} = P_{1,j}$$

For $t = 1, \dots, T$ do:

- sample $j_t \sim \overline{P}_t$

- $B_{j_t}$ for one step

- receive reward
$$\underline{r_t} = f(\underbrace{X_t \sim B_{j_t}}_{\text{context}}) + \underbrace{\xi_t}_{\leftarrow \text{noise}}$$

- send feedback
$$\frac{r_t}{\overline{P}_{t,j_t}} \mathbb{1}\{j = j_t\} \text{ to all } B_j \quad j \in M \qquad |$$

$\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim\sim$

- Update $P_t \to P_{t+1}$ via OMD
$$(y_t, r_t, \overline{P}_t, j_t)$$

- smooth $P_t \to \overline{P}_t$

"Smooth" CORRAL  Pacchiano et al. 2020

* base algos only upates when selected

* original reward $r_{t,j_t}$ back to base

* only works in stochastic settings

* memory : a nightmare

Smooth wrapper

Base $B_j$

if selected at $t$:

$x_t \sim B_j$

$r_t = f(x_t) + \mathcal{P}_t$

$r_t^{(2)}(i^*)$ is bounded w.h.p

Smoothed $B_{\hat{j}}$
(internal state $s$)

if selected at $t$:

step 1:

$x_t^{(1)} \sim B_j$

$r_t^{(1)} = f(x_t^{(1)}) + \mathcal{P}_t$

step 2:

$q \sim \text{Uniform}(1, 2 \dots s)$

$x_t^{(2)} \sim B_j, q$

$r_t^{(2)} = f \cdots \cdots$

update $s = s+1$

(expected)

Regret decomposition

$$R(T) = \mathbb{E}\left[ \sum_{t=1}^{T} f(x^*) - f(x_t) \right]$$

$$= \mathbb{E}\left[ \sum^{T} f(x^*) - f(x_{i_t^*}) \right]$$

$$\overbrace{\phantom{xxxxxxxxx}}^{t-1}_{\#\mathrm{I}}$$

$$+ \mathbb{E}\Big[\underbrace{\sum_{t=1}^{T} f(X_{i^*,t}) - f(X_t)}_{\#\mathrm{II}}\Big]$$

#II: "model selection" cost

$$\propto \text{poly}(M)$$

#I: regret of base $i^*$

     w.r.t optimal action/policy

* has high-probability bounds!

- Assumes all $B_i$, $i \in [M]$

   has high-prob bound

$$U_i(t, \delta) \qquad \text{w.p.} \cdot 1-\delta$$

       if not misspecified

- $R(T) = \tilde{O}\big(\text{poly}(M) \, \underline{U_{i^*}(t, \delta)}\big)$

     w.p. $1-\delta$

---

Linear bandit results (model-selection)

       is known

$O(d_* \sqrt{T})$    if ...

& $\mathcal{X}$ is infinite
(action space)

smooth CORRAL

$$\frac{\tilde{O}(d_*^2 \sqrt{T})}{d_{1,2...M} = [1, 2, 2^2, .. 2^{\lceil \mu CD \rceil}]} \qquad d_* \subseteq D$$

Mod CB:
$$\tilde{O}\left( (k^{\frac{1}{4}} T^{\frac{3}{4}}) + \sqrt{kT} d_* \right)$$

K arms

RBBE : $\tilde{O}(d_*^2 \sqrt{T})$

lower bound :

$$\underline{\Omega(\sqrt{dT})} \quad \text{for} \quad d \leq \delta T$$

finite (k) armed setting

More complex settings for model selection

⌐ Lipschitz constant $L$
⌐ smoothness,

kernel parameters $\{$ lengthscale, etc..

for GP / kernelized bandits

...