

Bayesian bandits

finite arms/actions

Mean reward vector $\mu = \begin{bmatrix} \mu(1) \\ \vdots \\ \mu(K) \end{bmatrix} \sim P$ (prior)

Thompson Sampling

Observe history $H_{t-1} = H = \{(a_1, r_1), \dots, (a_{t-1}, r_{t-1})\}$
Sample mean vector from posterior $\mu_t \sim P_H(\mu)$ ← posterior
Choose best arm acc to μ_t , $a_t = \underset{a}{\operatorname{argmax}} \mu_t(a)$
Add a_t, r_t to $H_{t-1} \rightarrow H_t$.

Bayesian Regret $O(\sqrt{KT \log T})$

a_t and a^* have same dist given H_{t-1} ✓

$E[\text{UCB}(a_t, H_{t-1}) | H_{t-1}] = E[\text{UCB}(a^*, H_{t-1}) | H_{t-1}]$ ✓

$a^* : \mu \sim P \quad a^* = \underset{a}{\operatorname{argmax}} \mu(a) \quad \underline{\Pr(a^* | H_{t-1})}$
 $a_t : P_{H_{t-1}}(\mu) \quad \mu_t \sim P_{H_{t-1}}(\mu) \quad a_t = \underset{a}{\operatorname{argmax}} \mu_t(a) \quad \underline{\Pr(a_t | H_{t-1})}$ }

$$\begin{aligned} \Pr(a_t | H_{t-1}) &= \Pr(\mu_t(a_t) \geq \mu_t(a) \forall a | H_{t-1}) \\ &= \Pr(\mu(a_t) \geq \mu(a) \forall a | H_{t-1}) \\ &= \Pr(a^* | H_{t-1}) \end{aligned}$$

Continuous arms (correlated) - Gaussian Process prior

reward $\mu(x) \sim \text{GP}(m(x), k(x, x')) \quad \forall x \in \mathcal{X}$
↳ continuous

covariance function = $\begin{cases} x^T x' & \text{linear (Bayesian)} \\ \exp\{-\frac{\|x-x'\|^2}{2r^2}\} & \text{RBF or squared kernel} \\ \text{Matern kernel} & \end{cases}$

$$z = \mu(x) + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\mu \sim \text{GP}(0, K)$$



Observations $(x_1, y_1) \dots (x_T, y_T)$

Posterior is also $\text{GP}(m_T(x), k_T(x, x'))$

$$m_T(x) = K_T^{-1}(x) (K_T + \sigma^2 I)^{-1} \vec{y}_T$$

$$\vec{y}_T = [y_1 \dots y_T]^T$$

$$K_T = [k(x_i, x_j)]_{x_i, x_j \in \text{observations}}$$

$$k_T(x) = [k(x, x_1) \dots k(x, x_T)]^T$$

Posterior Covariance $k_T(x, x') = k(x, x') - k_T^T(x) (K_T + \sigma^2 I)^{-1} k_T(x')$

$$\sigma_T^2(x) = k_T(x, x) \quad \text{posterior variance} \equiv \text{uncertainty}$$

Experimental Design (Bayesian with GP prior)

Information gain $I(z_A, \mu) = \frac{1}{2} \log |I + \sigma^2 K_A|$

A - subset of sample from X $I(z_A, \mu_A) \quad K_A = [k(x, x')]_{x, x' \in A}$

$\max_{|A| \leq T} I(z_A, \mu) \leftarrow$ is NP-hard (even for GP)
 $\underbrace{\hspace{10em}}_{F(A)}$

$F(A)$ is a submodular function Krause-Guestron '05

$X, Y, X \subseteq Y$. for every $x \in Y$

$$F(X \cup x) - F(X) \geq F(Y \cup x) - F(Y)$$

Diminishing returns

monotonically \uparrow Any submodular function can be approximately (upto a constant) optimized using an efficient greedy procedure. Nemhauer '978

$$F(A_T) \geq (1 - \frac{1}{e}) \max_{|A| \leq T} F(A)$$

Greedy algorithm:

$$x_t = \operatorname{argmax}_{x \in X} F(A_{t-1} \cup x)$$

$$A_{t-1} = \{x_1, \dots, x_{t-1}\}$$

$$= \operatorname{argmax}_x \sigma_{t-1}(x)$$

GP Bandits (exp design doesn't exploit, so doesn't optimize cum regret)

GP-UCB

$$x_t = \operatorname{argmax}_x m_{t-1}(x) + \sqrt{\beta_t} \sigma_{t-1}(x)$$

\uparrow exploit \uparrow explore

$$K(\dots) \leq 1$$

$$\beta_t = 2 \log \frac{|X| t^2}{6\delta} \quad |X| \text{ finite}$$

Theorem: Regret $R_T = \sum_{t=1}^T \lambda(x^*) - \lambda(x_t) = O(\sqrt{T \beta_T \gamma_T}) \quad \forall T \quad \text{w.p.} \geq 1 - \delta$

max info gain $\gamma_T = \max_{|A|=T} \mathcal{I}(\lambda_A, \mu)$

Regret	Linear	RBF	Matern (ν)
$\sqrt{T \beta_T \gamma_T}$		continuous	finite
β_T	—————	$d \log \frac{T}{\delta}$ or $\log \frac{ X T}{\delta}$	—————
γ_T	$d \log T$	$(\log T)^{d+1}$	$T \frac{d(d+1)}{2\nu+d(d+1)} \log T$

Proof sketch: finite X

$$\text{w.p.} \geq 1 - \delta \quad |\mu(x) - m_{t-1}(x)| \leq \sqrt{\beta_t} \sigma_{t-1}(x) \quad \forall x \quad \mu \sim \text{GP}(m_{t-1}(x), \sigma_{t-1}(x))$$

conditioned on H_{t-1}

$$R_t = \mu(x^*) - \mu(x_t) \leq \underbrace{m_{t-1}(x^*) + \sqrt{\beta_t} \sigma_{t-1}(x^*)}_{\text{conc}} - \mu(x_t)$$

$$\sum_{t=1}^T m_{t+1}(x_t) + \sqrt{\beta_t} \sigma_{t+1}(x_t) - \mu(x_t)$$

$$\stackrel{\text{conc.}}{\leq} 2\sqrt{\beta_t} \sigma_{t+1}(x_t)$$

$$\sum_{t=1}^T R_t \leq 2 \sum_{t=1}^T \sqrt{\beta_t} \sigma_{t+1}(x_t) = O(\sqrt{T \beta_T} \sigma_T) \quad \text{w.p. } 1-\delta$$

$$\leq 2\sqrt{\beta_T} \sum_{t=1}^T \sigma_{t+1}(x_t) \cdot 1$$

$$\leq 2\sqrt{\beta_T} \sqrt{T} \sqrt{\sum_{t=1}^T \sigma_{t+1}^2(x_t)} \quad \text{IF}$$

$$\text{IF}(R_A, \mu) = \frac{1}{2} \sum_{t=1}^T \log\left(1 + \frac{\sigma_{t+1}^2(x)}{\sigma^2}\right) \rightarrow \sigma^{-2} K_A$$

$$\sum_{t=1}^T \sigma_{t+1}^2(x) = \sigma^2 \sum_{t=1}^T \frac{\sigma_{t+1}^2(x)}{\sigma^2} = O\left(\sigma^2 \sum_{t=1}^T \log\left(1 + \frac{\sigma_{t+1}^2(x)}{\sigma^2}\right)\right)$$

$$z \leq O(\log(1+z))$$

$$e^z \leq 1+z$$

Continuous setting

$X \subseteq \mathbb{R}^d$ compact, convex assumption on kernel (\equiv Lipschitz)

follows discretization

RKHS version

μ smoothness controlled by kernel

$$\|\mu\|_k \leq B$$

- low RKHS norm induced by kernel k .

$\mu \in \mathcal{H}_k(X)$ - RKHS(k) - Reproducing Kernel Hilbert Space

$$\langle \cdot, \cdot \rangle_k \text{ obeys } \langle \mu, k(x) \rangle_k = \mu(x)$$

$$\text{RKHS norm } \langle \mu, \mu \rangle_k = \|\mu\|_k^2$$

Then: $\mu \in \text{RKHS}(k) \quad \|\mu\|_k^2 \leq B$

$$\beta_t \approx 2B + C \gamma_t \log \gamma_t \quad \text{w.p. } 1-\delta \quad R_T = O(\sqrt{T \beta_T} \sigma_T)$$

... .. on GP

... .. A.

Doesn't assume $\mu = \dots$

uniform for all $\mu \in \mathcal{V}_k$
& $\|\mu\|_k^2 \leq B$

Proof sketch: Freedman's inequality

$$Z_t = \|m_t - \mu\|_{K_T}^2$$

- 3 terms

same as HW
linear problem

GP - Thompson sampling.

$$\mu_t \sim \text{GP}(m_t, K_t)$$

$$x_t = \arg \max_{x \in \mathcal{X}} \mu_t(x)$$

$$\mu_t(x) + \frac{1}{\sqrt{\beta_t}} \sqrt{K_t(x, x)}$$