

Kernelized Bandits Beyond Realizability

Davide Maran & Csaba Szepesvári

Setting

Input space $\mathcal{X} = [-1, 1]^m$, function $f : \mathcal{X} \rightarrow \mathbb{R}$ to optimize.

Setting

Input space $\mathcal{X} = [-1, 1]^m$, function $f : \mathcal{X} \rightarrow \mathbb{R}$ to optimize.

Kernel

The kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is symmetric and positive semi-definite, and

$$\sup_{x, y \in \mathcal{X}} |k(x, y)| \leq \kappa < +\infty.$$

Input space $\mathcal{X} = [-1, 1]^m$, function $f : \mathcal{X} \rightarrow \mathbb{R}$ to optimize.

Kernel

The kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is symmetric and positive semi-definite, and

$$\sup_{x, y \in \mathcal{X}} |k(x, y)| \leq \kappa < +\infty.$$

Interaction models

- (*Offline optimization*) Passive data:

$$(\mathbf{x}, \mathbf{y}) = ((x_1, y_1), \dots, (x_n, y_n)) \text{ s.t. } y_t = f(x_t) + \eta_t$$

for independent 1-subgaussian η_t .

Input space $\mathcal{X} = [-1, 1]^m$, function $f : \mathcal{X} \rightarrow \mathbb{R}$ to optimize.

Kernel

The kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is symmetric and positive semi-definite, and

$$\sup_{x, y \in \mathcal{X}} |k(x, y)| \leq \kappa < +\infty.$$

Interaction models

- (*Offline optimization*) Passive data:

$$(\mathbf{x}, \mathbf{y}) = ((x_1, y_1), \dots, (x_n, y_n)) \text{ s.t. } y_t = f(x_t) + \eta_t$$

for independent 1-subgaussian η_t .

- (*Online optimization*) Choose sequentially x_1, \dots, x_n and receive $y_t = f(x_t) + \eta_t$.

Target approximator

For $f_\star \in \mathcal{H} = \text{RKHS}(k; \mathcal{X})$,

$$\varepsilon := \sup_{x \in \mathcal{X}} |f(x) - f_\star(x)|.$$

Target approximator

For $f_\star \in \mathcal{H} = \text{RKHS}(k; \mathcal{X})$,

$$\varepsilon := \sup_{x \in \mathcal{X}} |f(x) - f_\star(x)|.$$

- 1 *Offline setting*: output \hat{x}_n , minimize

$$r_n = f(x^\star) - f(\hat{x}_n).$$

Target approximator

For $f_* \in \mathcal{H} = \text{RKHS}(k; \mathcal{X})$,

$$\varepsilon := \sup_{x \in \mathcal{X}} |f(x) - f_*(x)|.$$

- 1 *Offline setting:* output \hat{x}_n , minimize

$$r_n = f(x^*) - f(\hat{x}_n).$$

- 2 *Online setting:* minimize

$$R_n = \sum_{t=1}^n f(x^*) - f(x_t).$$

Target approximator

For $f_\star \in \mathcal{H} = \text{RKHS}(k; \mathcal{X})$,

$$\varepsilon := \sup_{x \in \mathcal{X}} |f(x) - f_\star(x)|.$$

- 1 *Offline setting:* output \hat{x}_n , minimize

$$r_n = f(x^\star) - f(\hat{x}_n).$$

- 2 *Online setting:* minimize

$$R_n = \sum_{t=1}^n f(x^\star) - f(x_t).$$

In terms of $n, \varepsilon, \|f_\star\|_{\mathcal{H}}$ and properties of \mathcal{H} .

Regret bound:

$$R_n \leq \tilde{O} \left(\gamma_n \sqrt{n} + \|f_\star\|_{\mathcal{H}} \sqrt{\gamma_n n} + \sqrt{\gamma_n} n \varepsilon \right).$$

Regret bound:

$$R_n \leq \tilde{O}(\gamma_n \sqrt{n} + \|f_\star\|_{\mathcal{H}} \sqrt{\gamma_n n} + \sqrt{\gamma_n n} \varepsilon).$$

Maximal information gain

$$\gamma_n := \max_{x_1, \dots, x_n \in \mathcal{X}} \log \det(I + \lambda^{-1} K), \text{ for } K_{ij} = k(x_i, x_j).$$

For k in the Matérn family [Vakili et al., 2021]:

$$\gamma_n \asymp n^{\frac{m}{2\nu+m}},$$

except for log –factors.

Regret bound:

$$R_n \leq \tilde{O}(\gamma_n \sqrt{n} + \|f_\star\|_{\mathcal{H}} \sqrt{\gamma_n n} + \sqrt{\gamma_n n} \varepsilon).$$

Maximal information gain

$$\gamma_n := \max_{x_1, \dots, x_n \in \mathcal{X}} \log \det(I + \lambda^{-1} K), \text{ for } K_{ij} = k(x_i, x_j).$$

For k in the Matérn family [Vakili et al., 2021]:

$$\gamma_n \asymp n^{\frac{m}{2\nu+m}},$$

except for log –factors.

Error amplification

If $\varepsilon \approx n^{-\frac{m}{2\nu+m}}$ this gives *linear regret*.

Notation

- $P_{\mathbf{x},\lambda}g$: KRR estimator data $(\mathbf{x}, g(\mathbf{x}))$
- $\sigma_{\mathbf{x},\lambda}$: posterior standard deviation

$$\sigma_{\mathbf{x},\lambda}(x) := \sqrt{k(x, x) - \Psi_{\mathbf{x}}(x)^{\top} (K + \lambda I)^{-1} \Psi_{\mathbf{x}}(x)}.$$

Notation

- $P_{\mathbf{x},\lambda}g$: KRR estimator data $(\mathbf{x}, g(\mathbf{x}))$
- $\sigma_{\mathbf{x},\lambda}$: posterior standard deviation

$$\sigma_{\mathbf{x},\lambda}(x) := \sqrt{k(x, x) - \Psi_{\mathbf{x}}(x)^{\top} (K + \lambda I)^{-1} \Psi_{\mathbf{x}}(x)}.$$

$$|P_{\mathbf{x},\lambda}f(x) - P_{\mathbf{x},\lambda}f_{\star}(x)| \leq \sqrt{\frac{n}{\lambda}} \sigma_{\mathbf{x},\lambda}(x) \varepsilon.$$

Notation

- $P_{\mathbf{x},\lambda}g$: KRR estimator data $(\mathbf{x}, g(\mathbf{x}))$
- $\sigma_{\mathbf{x},\lambda}$: posterior standard deviation

$$\sigma_{\mathbf{x},\lambda}(x) := \sqrt{k(x, x) - \Psi_{\mathbf{x}}(x)^{\top} (K + \lambda I)^{-1} \Psi_{\mathbf{x}}(x)}.$$

$$|P_{\mathbf{x},\lambda}f(x) - P_{\mathbf{x},\lambda}f_{\star}(x)| \leq \sqrt{\frac{n}{\lambda}} \sigma_{\mathbf{x},\lambda}(x) \varepsilon.$$

Can we refine this result?

$\mathbf{x} \stackrel{i.i.d.}{\sim} p$, $(\phi_i)_{i=1}^{\infty}$ orthogonal in $L^2(p)$ and

$$k(x, y) = \sum_{i=1}^{\infty} \mu_i \phi_i(x) \phi_i(y).$$

Population version

$\mathbf{x} \stackrel{i.i.d.}{\sim} p$, $(\phi_i)_{i=1}^{\infty}$ orthogonal in $L^2(p)$ and

$$k(x, y) = \sum_{i=1}^{\infty} \mu_i \phi_i(x) \phi_i(y).$$

Population operator

$$P_{\tau} f(x) := \sum_{i=1}^{\infty} \frac{\mu_i f_i \phi_i(x)}{\tau + \mu_i} \quad f_i = \langle f, \phi_i \rangle_{L^2(p)}.$$

For $\tau = \lambda/n$, $\|P_{\tau} f - P_{\mathbf{x}, \lambda} f\|_{\infty} = o(1)$.

Population version

$\mathbf{x} \stackrel{i.i.d.}{\sim} p$, $(\phi_i)_{i=1}^{\infty}$ orthogonal in $L^2(p)$ and

$$k(x, y) = \sum_{i=1}^{\infty} \mu_i \phi_i(x) \phi_i(y).$$

Population operator

$$P_{\tau} f(x) := \sum_{i=1}^{\infty} \frac{\mu_i f_i \phi_i(x)}{\tau + \mu_i} \quad f_i = \langle f, \phi_i \rangle_{L^2(p)}.$$

For $\tau = \lambda/n$, $\|P_{\tau} f - P_{\mathbf{x}, \lambda} f\|_{\infty} = o(1)$.

For bounded eigenfunctions

$$\|P_{\tau} f - P_{\tau} f_{\star}\|_{\infty} \leq \sqrt{d_{\text{eff}}(k|\tau)} \varepsilon.$$

Note: $\sqrt{\frac{n}{\lambda}} \sup_{x \in \mathcal{X}} \sigma_{\mathbf{x}, \lambda}(x) \propto \sqrt{d_{\text{eff}}(k|\tau)} + o(1)$.

Error decomposition

Idea:

$$r_n \leq 2 \cdot \sup_{x \in \mathcal{X}} |P_{\mathbf{x}, \lambda} \mathbf{y}(x) - f(x)|$$

Error decomposition

Idea:

$$r_n \leq 2 \cdot \sup_{x \in \mathcal{X}} |P_{\mathbf{x}, \lambda} \mathbf{y}(x) - f(x)|$$

Function decomposition

$$f(x) = f_*(x) + f_\varepsilon(x)$$

Error decomposition

Idea:

$$r_n \leq 2 \cdot \sup_{x \in \mathcal{X}} |P_{\mathbf{x}, \lambda} \mathbf{y}(x) - f(x)|$$

Function decomposition

$$f(x) = f_{\star}(x) + f_{\varepsilon}(x)$$

$$\begin{aligned} |P_{\mathbf{x}, \lambda} \mathbf{y}(x) - f(x)| &\leq \underbrace{|P_{\mathbf{x}, \lambda} f_{\star}(x) - f_{\star}(x)|}_{\text{stoch. y}} + \underbrace{|P_{\mathbf{x}, \lambda} f_{\varepsilon}(x) - P_{\tau} f_{\varepsilon}(x)|}_{\text{stoch. design}} \\ &\quad + \underbrace{|P_{\tau} f_{\varepsilon}(x) - f_{\varepsilon}(x)|}_{\text{misspecification}}. \end{aligned}$$

Error decomposition

Idea:

$$r_n \leq 2 \cdot \sup_{x \in \mathcal{X}} |P_{\mathbf{x}, \lambda} \mathbf{y}(x) - f(x)|$$

Function decomposition

$$f(x) = f_*(x) + f_\varepsilon(x)$$

$$\begin{aligned} |P_{\mathbf{x}, \lambda} \mathbf{y}(x) - f(x)| &\leq \underbrace{|P_{\mathbf{x}, \lambda} f_*(x) - f_*(x)|}_{\text{stoch. y}} + \underbrace{|P_{\mathbf{x}, \lambda} f_\varepsilon(x) - P_\tau f_\varepsilon(x)|}_{\text{stoch. design}} \\ &\quad + \underbrace{|P_\tau f_\varepsilon(x) - f_\varepsilon(x)|}_{\text{misspecification}}. \end{aligned}$$

- **stoch. y** $\lesssim \sqrt{d_{\text{eff}}(k|\tau)n^{-1/2}} \|f_*\|_{\mathcal{H}}$ [Hsu et al., 2012]
- **stoch. design** $\lesssim \kappa^2 n^{-1/2} (\varepsilon + n^{-1/2})$ [Smale and Zhou, 2005]
- **misspecification** $\lesssim \Lambda(P_\tau) \varepsilon$

Lebesgue constant

For an operator Π , we call

$$\Lambda(\Pi) := \sup_{g \in L^\infty(\mathcal{X})} \frac{\|\Pi g\|_\infty}{\|g\|_\infty}.$$

Lebesgue constant

For an operator Π , we call

$$\Lambda(\Pi) := \sup_{g \in L^\infty(\mathcal{X})} \frac{\|\Pi g\|_\infty}{\|g\|_\infty}.$$

$L^\infty - L^1$ duality

$$\Lambda(P_\tau) = \max_{x \in \mathcal{X}, f_i} \left\| \sum_{i=1}^{\infty} \frac{\mu_i f_i \phi_i(x)}{\tau + \mu_i} \phi_i(\cdot) \right\|_{L^1(\rho)}.$$

Intractable due to the generality of the ϕ 's

Spectral Lebesgue Growth

Many different kernels use the same features and different eigenvalues!

Many different kernels use the same features and different eigenvalues!

Bochner's theorem

Stationary periodic kernels admit a discrete Fourier basis as their exact eigenfunctions.

Many different kernels use the same features and different eigenvalues!

Bochner's theorem

Stationary periodic kernels admit a discrete Fourier basis as their exact eigenfunctions.

Fourier features \implies S.L.G.

S.L.G.

$$\sup_{x \in \mathcal{X}} \left\| \sum_{j=1}^i \phi_j(x) \phi_j(\cdot) \right\|_{L^1(p)} \lesssim \log(e + i).$$

The question turns on eigenvalue profiles

Call μ_i the eigenvalues sorted according to S.L.G.

The question turns on eigenvalue profiles

Call μ_i the eigenvalues sorted according to S.L.G.

Typical case = worst case

For a random eigenvalue profile $\mu_i \sim i^{-s} X_i$ with $X_i \stackrel{i.i.d.}{\sim} \text{Bern}(1/2)$,

$$\Lambda(P_\tau) \asymp \sqrt{d_{\text{eff}}(k|\tau)} \quad a.s.$$

The question turns on eigenvalue profiles

Call μ_i the eigenvalues sorted according to S.L.G.

Typical case = worst case

For a random eigenvalue profile $\mu_i \sim i^{-s} X_i$ with $X_i \stackrel{i.i.d.}{\sim} \text{Bern}(1/2)$,

$$\Lambda(P_\tau) \asymp \sqrt{d_{\text{eff}}(k|\tau)} \quad a.s.$$

Monotone spectrum

If $\mu_i \downarrow 0$ (monotonically),

$$\Lambda(P_\tau) \lesssim \log(e + \kappa/\tau).$$

No effective dimension, only log-dependence in κ and $1/\tau$.

The question turns on eigenvalue profiles

Call μ_i the eigenvalues sorted according to S.L.G.

Typical case = worst case

For a random eigenvalue profile $\mu_i \sim i^{-s} X_i$ with $X_i \stackrel{i.i.d.}{\sim} \text{Bern}(1/2)$,

$$\Lambda(P_\tau) \asymp \sqrt{d_{\text{eff}}(k|\tau)} \quad a.s.$$

Monotone spectrum

If $\mu_i \downarrow 0$ (monotonically),

$$\Lambda(P_\tau) \lesssim \log(e + \kappa/\tau).$$

No effective dimension, only log-dependence in κ and $1/\tau$.

Tensorized kernels

$k(x, y) = \prod_{j=1}^m k_j(x_j, y_j)$ where satisfies S.L.G. and $\mu_i \downarrow 0$.

$$\Lambda(P_\tau) \lesssim \log(e + \kappa^m/\tau)^{2m-1}.$$

Controlling the regret

$$R_n = \sum_{t=1}^n f(x^*) - f(x_t)$$

requires *exploration-exploitation trade-off*.

Controlling the regret

$$R_n = \sum_{t=1}^n f(x^*) - f(x_t)$$

requires *exploration-exploitation trade-off*.

Explore-then-commit

- 1 Sample uniformly the first $n_0 = n^\alpha$ queries
- 2 Put the remaining $n - n^\alpha$ on \hat{x}_{n_0}

Controlling the regret

$$R_n = \sum_{t=1}^n f(x^*) - f(x_t)$$

requires *exploration-exploitation trade-off*.

Explore-then-commit

- 1 Sample uniformly the first $n_0 = n^\alpha$ queries
- 2 Put the remaining $n - n^\alpha$ on \hat{x}_{n_0}
- 3 **Sub-optimal**

Controlling the regret

$$R_n = \sum_{t=1}^n f(x^*) - f(x_t)$$

requires *exploration-exploitation trade-off*.

Explore-then-commit

- 1 Sample uniformly the first $n_0 = n^\alpha$ queries
- 2 Put the remaining $n - n^\alpha$ on \hat{x}_{n_0}
- 3 **Sub-optimal**

Smooth level curves

- 1 Assume level sets $\{f(\cdot) \geq M\}$ are smooth

Controlling the regret

$$R_n = \sum_{t=1}^n f(x^*) - f(x_t)$$

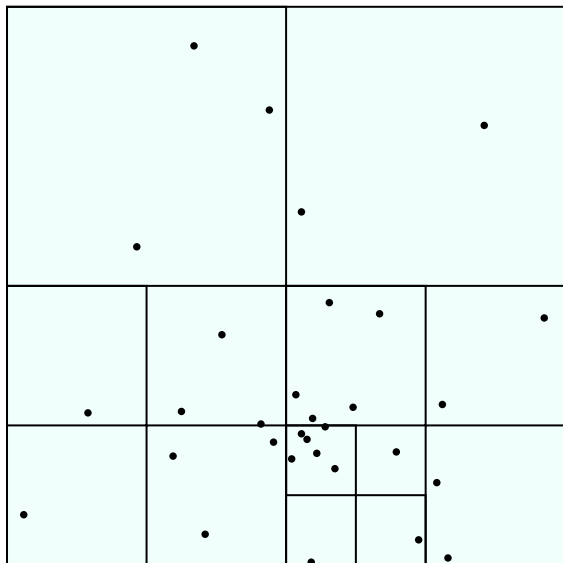
requires *exploration-exploitation trade-off*.

Explore-then-commit

- 1 Sample uniformly the first $n_0 = n^\alpha$ queries
- 2 Put the remaining $n - n^\alpha$ on \hat{x}_{n_0}
- 3 **Sub-optimal**

Smooth level curves

- 1 Assume level sets $\{f(\cdot) \geq M\}$ are smooth
- 2 **Very restrictive**



Split any region A with too many point into 2^m sub-regions. Always keeps that

$$\gamma_{|A \cap \mathbf{x}|} \leq \text{polylog}(n).$$

Domain splitting

Let $\Pi : L^\infty(\mathcal{X}) \rightarrow L^\infty(\mathcal{X})$ such that

$$\Pi[f; \mathcal{X}](x) = \underbrace{1\{x \in \mathcal{X}_0\} \Pi[f; \mathcal{X}_0](x)}_{\Pi_0} + \underbrace{1\{x \in \mathcal{X}_0^c\} \Pi[f; \mathcal{X}_0^c](x)}_{\Pi_1},$$

Domain splitting

Let $\Pi : L^\infty(\mathcal{X}) \rightarrow L^\infty(\mathcal{X})$ such that

$$\Pi[f; \mathcal{X}](x) = \underbrace{1\{x \in \mathcal{X}_0\} \Pi[f; \mathcal{X}_0](x)}_{\Pi_0} + \underbrace{1\{x \in \mathcal{X}_0^c\} \Pi[f; \mathcal{X}_0^c](x)}_{\Pi_1},$$

Then, $\Lambda(\Pi) = \max \{ \Lambda(\Pi_0), \Lambda(\Pi_1) \}$.

Domain splitting

Let $\Pi : L^\infty(\mathcal{X}) \rightarrow L^\infty(\mathcal{X})$ such that

$$\Pi[f; \mathcal{X}](x) = \underbrace{1\{x \in \mathcal{X}_0\} \Pi[f; \mathcal{X}_0](x)}_{\Pi_0} + \underbrace{1\{x \in \mathcal{X}_0^c\} \Pi[f; \mathcal{X}_0^c](x)}_{\Pi_1},$$

Then, $\Lambda(\Pi) = \max\{\Lambda(\Pi_0), \Lambda(\Pi_1)\}$.

Lebesgue constant and information gain

For KRR, $\Lambda(P_{A \cap \mathbf{x}, \lambda}) \lesssim \sqrt{\gamma |A \cap \mathbf{x}|}$.

Domain splitting

Let $\Pi : L^\infty(\mathcal{X}) \rightarrow L^\infty(\mathcal{X})$ such that

$$\Pi[f; \mathcal{X}](x) = \underbrace{1\{x \in \mathcal{X}_0\} \Pi[f; \mathcal{X}_0](x)}_{\Pi_0} + \underbrace{1\{x \in \mathcal{X}_0^c\} \Pi[f; \mathcal{X}_0^c](x)}_{\Pi_1},$$

Then, $\Lambda(\Pi) = \max\{\Lambda(\Pi_0), \Lambda(\Pi_1)\}$.

Lebesgue constant and information gain

For KRR, $\Lambda(\mathbb{P}_{A \cap \mathbf{x}, \lambda}) \lesssim \sqrt{\gamma_{|A \cap \mathbf{x}|}}$. With π -GP-UCB,

$$\gamma_{|A \cap \mathbf{x}|} \leq \text{polylog}(n),$$

on **all regions**.

Eigendecay on subdomains

Restricting k to an hypercube of side $\rho > 0$ gives $\mu_{\rho,i} \leq \rho^\alpha \mu_i$ and bounded $\phi_{\rho,i}$.

Eigendecay on subdomains

Restricting k to an hypercube of side $\rho > 0$ gives $\mu_{\rho,i} \leq \rho^\alpha \mu_i$ and bounded $\phi_{\rho,i}$.

$k \in \text{Matérn family} \implies \alpha = 2\nu$.

Eigendecay on subdomains

Restricting k to an hypercube of side $\rho > 0$ gives $\mu_{\rho,i} \leq \rho^\alpha \mu_i$ and bounded $\phi_{\rho,i}$.

$k \in \text{Matérn family} \implies \alpha = 2\nu$.

$$R_n \leq \tilde{O}(\|f_\star\|_{\mathcal{H}} \sqrt{\gamma_n n} + n\varepsilon).$$

- Is there some useful kernel such that the $\sqrt{d_{\text{eff}}(k|\tau)}\varepsilon$ is inevitable?

- Is there some useful kernel such that the $\sqrt{d_{\text{eff}}(k|\tau)}\varepsilon$ is inevitable?
- Offline method \approx localization of the spectrum, online method \approx localization of the domain



Bogunovic, I. and Krause, A. (2021).
Misspecified gaussian process bandit optimization.
Advances in neural information processing systems, 34:3004–3015.



Hsu, D., Kakade, S. M., and Zhang, T. (2012).
Random design analysis of ridge regression.
In *Conference on learning theory*, pages 9.1–9.24. JMLR Workshop and Conference Proceedings.



Janz, D., Burt, D., and González, J. (2020).
Bandit optimisation of functions in the matérn kernel rkhs.
In *International Conference on Artificial Intelligence and Statistics*, pages 2486–2495. PMLR.



Smale, S. and Zhou, D.-X. (2005).
Shannon sampling ii: Connections to learning theory.
Applied and Computational Harmonic Analysis, 19(3):285–302.



Vakili, S., Khezeli, K., and Picheny, V. (2021).
On information gain and regret bounds in gaussian process bandits.
In *International Conference on Artificial Intelligence and Statistics*, pages 82–90. PMLR.