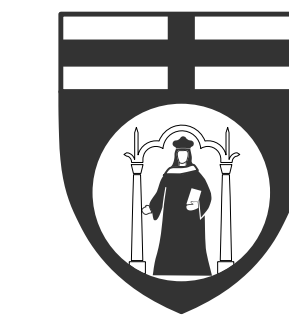


Gaussian Process Optimization with Adaptive Sketching: Scalable and No Regret

Daniele Calandriello*, Luigi Carratino†, Alessandro Lazaric‡, Michal Valko§, Lorenzo Rosasco*,†,¶

*LCSL - Istituto Italiano di Tecnologia, †Università degli Studi di Genova, ‡Facebook AI Research Paris,

§INRIA Lille - Nord Europe, ¶Massachusetts Institute of Technology



In a nutshell

Gaussian process-UCB (GP-UCB) is a popular Bayesian/bandit optimization alternative to grad student descent. However, it requires $\mathcal{O}(T^3)$ time and $\mathcal{O}(T^2)$ space and does not scale.

We introduce the first general GP optimization algorithm (BKB) that is **no regret** and **provably scalable**, with near-linear runtime $\mathcal{O}(Td_{\text{eff}}^2)$. It also maintains **valid posterior variance estimates** at all steps, while **previous approaches could under/over-estimate the confidence intervals** of the GP. BKB main ingredient is a novel **adaptive selection** of inducing point using **approximate posterior variance sampling**.

Gaussian process optimization

Arms $\mathcal{A} = \{\mathbf{x}_i\}_{i=1}^A$ with $\mathbf{x}_i \in \mathbb{R}^d$, similarity (kernel) $k(\cdot, \cdot)$ and RKHS \mathcal{H}

- For $t \in [1, \dots, T]$:
- (1) select $\mathbf{x}_{t+1} = \arg \max_{\mathbf{x}_i} u_t(\mathbf{x}_i)$
 - (2) Receive noisy feedback $y_{t+1} = f(\mathbf{x}_{t+1}) + \eta_{t+1}$
 - (3) Improve u_{t+1} for next time

Goal: minimize **regret** $R_T = \sum_{t=1}^T f(\mathbf{x}_*) - f(\mathbf{x}_t)$ vs. $\mathbf{x}_* = \arg \max_{\mathbf{x}_i} f(\mathbf{x}_i)$

Assumption: $f \in \mathcal{H}$ arbitrary but $\|f\| \leq F$ (frequentist/bandit regret)

Measuring the complexity of a GP

Maximum **information gain**: $\gamma_T \triangleq \max_{\mathcal{D} \subset \mathcal{A}, |\mathcal{D}|=T} \log \det(\mathbf{K}_{\mathcal{D}}/\lambda + \mathbf{I})$

Effective dimension/rank: $d_{\text{eff}} \triangleq \sum_{i=1}^T \sigma_i^2(\tilde{\mathbf{x}}_i)$

From γ_T to d_{eff} : $\log \det(\mathbf{K}_T + \mathbf{I}) \leq d_{\text{eff}} \log(T) \ll \gamma_T \log(T)$

GP-UCB and sparse GPs

GP-UCB:

$$\mu_t(\mathbf{x}) = \mathbf{k}_t(\mathbf{x})^\top (\mathbf{K}_t + \lambda \mathbf{I})^{-1} \mathbf{y}_t$$

$$u_t(\mathbf{x}) = \mu_t(\mathbf{x}) + \beta_t \sigma_t(\mathbf{x}) \quad \sigma_t^2(\mathbf{x}) = \frac{1}{\lambda} \left(\mathbf{k}(\mathbf{x}, \mathbf{x}) - \mathbf{k}_t(\mathbf{x})^\top (\mathbf{K}_t + \lambda \mathbf{I})^{-1} \mathbf{k}_t(\mathbf{x}) \right)$$

No regret $R_T^{\text{GP-UCB}} \leq \tilde{\mathcal{O}}(\sqrt{T}(\gamma_T + \sqrt{\gamma_T F}))$ but too slow $\mathcal{O}(At^2)$ per step

Sparse GPs: given m inducing points $\mathcal{S} = \{\mathbf{x}_j\}_{j=1}^m$ (a.k.a. dictionary) replace $k(\mathbf{x}_i, \mathbf{x}_j)$ with $\tilde{k}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{k}_{\mathcal{S}}(\mathbf{x}_i)^\top \mathbf{K}_{\mathcal{S}}^{-1} \mathbf{k}_{\mathcal{S}}(\mathbf{x}_j)$

GP-UCB + DTC:

$$\tilde{\mu}_t(\mathbf{x}) = \tilde{\mathbf{k}}_t(\mathbf{x})^\top (\tilde{\mathbf{K}}_t + \lambda \mathbf{I})^{-1} \mathbf{y}_t$$

$$\tilde{u}_t(\mathbf{x}) = \tilde{\mu}_t(\mathbf{x}) + \tilde{\beta}_t \tilde{\sigma}_t(\mathbf{x}) \quad \tilde{\sigma}_t^2(\mathbf{x}) = \frac{1}{\lambda} \left(\mathbf{k}(\mathbf{x}, \mathbf{x}) - \tilde{\mathbf{k}}_t(\mathbf{x})^\top (\tilde{\mathbf{K}}_t + \lambda \mathbf{I})^{-1} \tilde{\mathbf{k}}_t(\mathbf{x}) \right)$$

Deterministic training conditional (DTC) a.k.a. **projected GP**

↳ $\mathcal{O}(Am^2 + m^3)$ per step since $\text{Rank}(\tilde{\mathbf{K}}_t) = m$ but is it **no regret**?

🤔 **Problem:** how to choose \mathcal{S} for good accuracy/regret? 🤔

Budgeted Kernelized Bandits (BKB)

\mathbf{X}_t changes over time

💡 \mathcal{S}_t must change with t

Accuracy-efficiency tradeoff of m

💡 adaptively resize \mathcal{S}_t

$\sigma_t^2(\cdot)$ captures informative arms

💡 include \mathbf{x}_i with large $\sigma_t^2(\mathbf{x}_i)$

Greedy inclusion hard to analyze

💡 random inclusion $p_{t,i} \propto \sigma_t^2(\cdot)$

for $t = \{1, \dots, T-1\}$ **do**

Compute $\tilde{\mu}_t(\mathbf{x}_i)$ and $\tilde{\sigma}_t^2(\mathbf{x}_i)$ for all \mathbf{x}_i ;

Select $\tilde{\mathbf{x}}_{t+1} \leftarrow \arg \max_{\mathbf{x}_i \in \mathcal{A}} \tilde{u}_t(\mathbf{x}_i)$;

for $i = \{1, \dots, t+1\}$ **do**

Set $\tilde{p}_{t+1,i} \leftarrow \bar{q} \cdot \tilde{\sigma}_t^2(\tilde{\mathbf{x}}_i)$;

Draw $q_{t+1,i} \sim \text{Bernoulli}(\tilde{p}_{t+1,i})$;

If $q_{t+1,i} = 1$ **then** include $\tilde{\mathbf{x}}_i$ in \mathcal{S}_{t+1} ;

end

end

Main result: BKB is scalable and no regret

Theorem: Let $\tilde{\beta}_t \triangleq 2\sqrt{(\sum_{s=1}^t \tilde{\sigma}_s^2(\tilde{\mathbf{x}}_s)) \log(t) + \log(1/\delta) + 3\sqrt{\lambda}F}$. Then w.p. $1 - \delta$, for all $t \in [T]$ and all $\mathbf{x} \in \mathcal{A}$, we have

$$\sigma_t^2(\mathbf{x})/2 \leq \tilde{\sigma}_t^2(\mathbf{x}) \leq 2\sigma_t^2(\mathbf{x}) \quad \text{and} \quad |\mathcal{S}_t| \leq \mathcal{O}(d_{\text{eff}} \log(t/\delta)),$$

and BKB suffers at most regret

$$R_T^{\text{BKB}} \leq 32\sqrt{T} \left(d_{\text{eff}} \log(T) + \sqrt{\lambda F^2 d_{\text{eff}} \log(T)} + \log(1/\delta) \right)$$

💡 $R_T^{\text{BKB}} \leq 16R_T^{\text{GP-UCB}} \log(T)$: no regret but only $\tilde{\mathcal{O}}(TAd_{\text{eff}}^2)$ time! 💡

😊 $\tilde{\beta}_t$ **computable** in $\tilde{\mathcal{O}}(Ad_{\text{eff}}^2)$ time replacing worst-case bounds on γ_T

😊 **No assumptions** on k (e.g., not only stationary k)

😞 **No free lunch**: worst-case falls back to GP-UCB

😞 **Not incremental**: have to recompute \mathcal{S}_t at each step

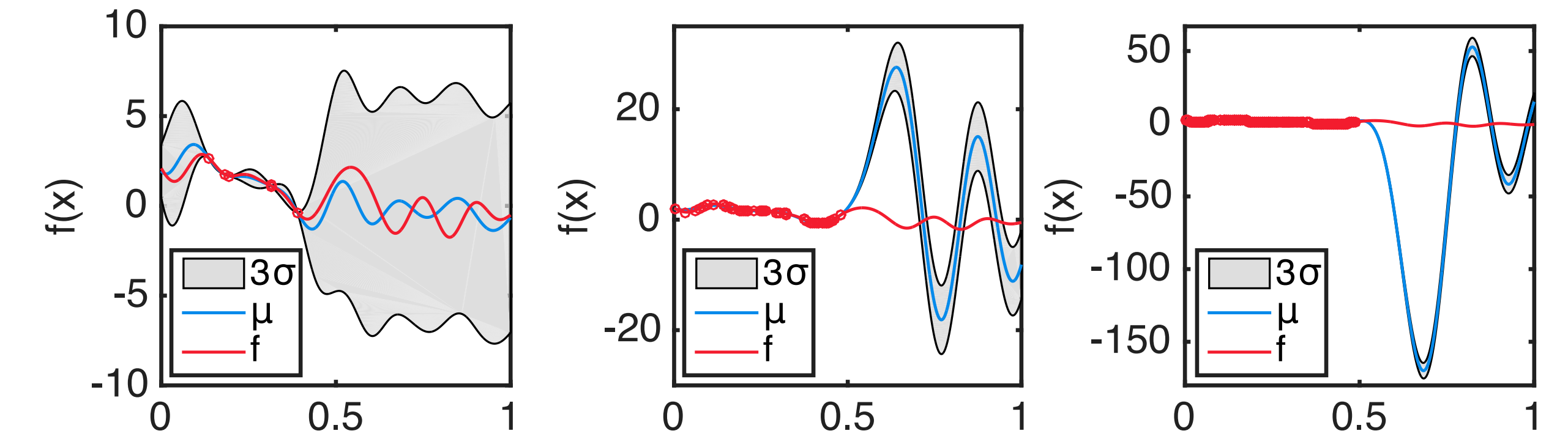
Alg.	$k(\cdot, \cdot)$	m	$R_T/R_T^{\text{GP-UCB}}$
TS-QFF [2]	stationary	$2^d \gamma_T$	16
SOFUL [3]	linear	k	$1 + \sum_{i=k+1}^T \lambda_i(\mathbf{K}_T)$
BKB	any	d_{eff}	$16 \log(T)$

😞 DTC is **not a GP** (not consistent), but now a **justified** heuristic

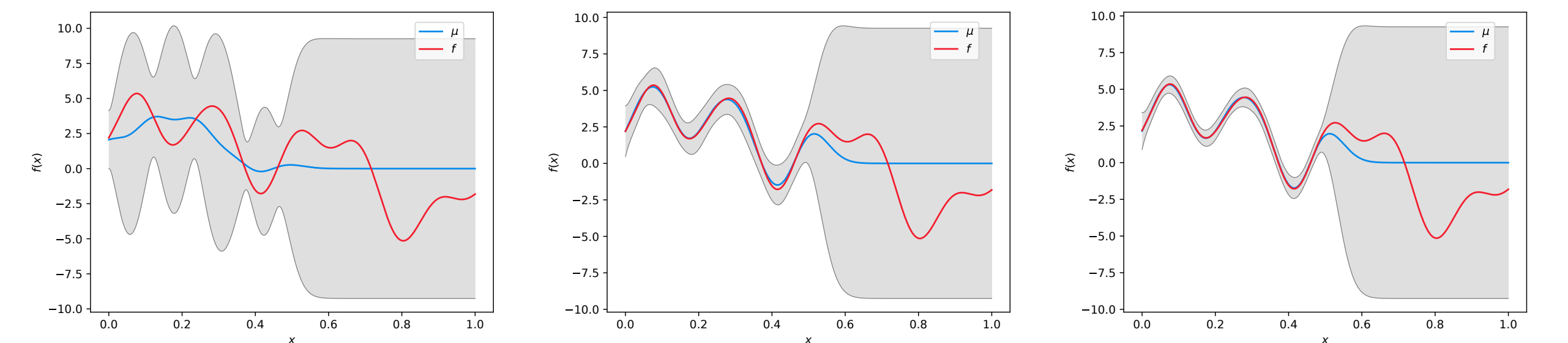
😊 Easy **extension to infinite** \mathcal{A} , but **how to optimize posterior**?

No variance starvation

😞 Sparse GPs become over/underconfident when $d_{\text{eff}} \gg m$ [5] 😞



😊 BKB adaptively matches GP rank m and d_{eff} 😊



Alg.	$\tilde{\sigma}_t(\cdot)$	$\tilde{k}(\cdot, \cdot)$	worst case
SoR	$\tilde{\mathbf{k}}(\mathbf{x}, \mathbf{x}) - \tilde{\mathbf{k}}_t(\mathbf{x})^\top (\tilde{\mathbf{K}}_t + \lambda \mathbf{I})^{-1} \tilde{\mathbf{k}}_t(\mathbf{x})$	fixed	0
DTC	$\mathbf{k}(\mathbf{x}, \mathbf{x}) - \tilde{\mathbf{k}}_t(\mathbf{x})^\top (\tilde{\mathbf{K}}_t + \lambda \mathbf{I})^{-1} \tilde{\mathbf{k}}_t(\mathbf{x})$	fixed	$k(\mathbf{x}, \mathbf{x})$
BKB	$\mathbf{k}(\mathbf{x}, \mathbf{x}) - \tilde{\mathbf{k}}_t(\mathbf{x})^\top (\tilde{\mathbf{K}}_t + \lambda \mathbf{I})^{-1} \tilde{\mathbf{k}}_t(\mathbf{x})$	adaptive	$(1 \pm 1/2)\sigma_t(\cdot)$

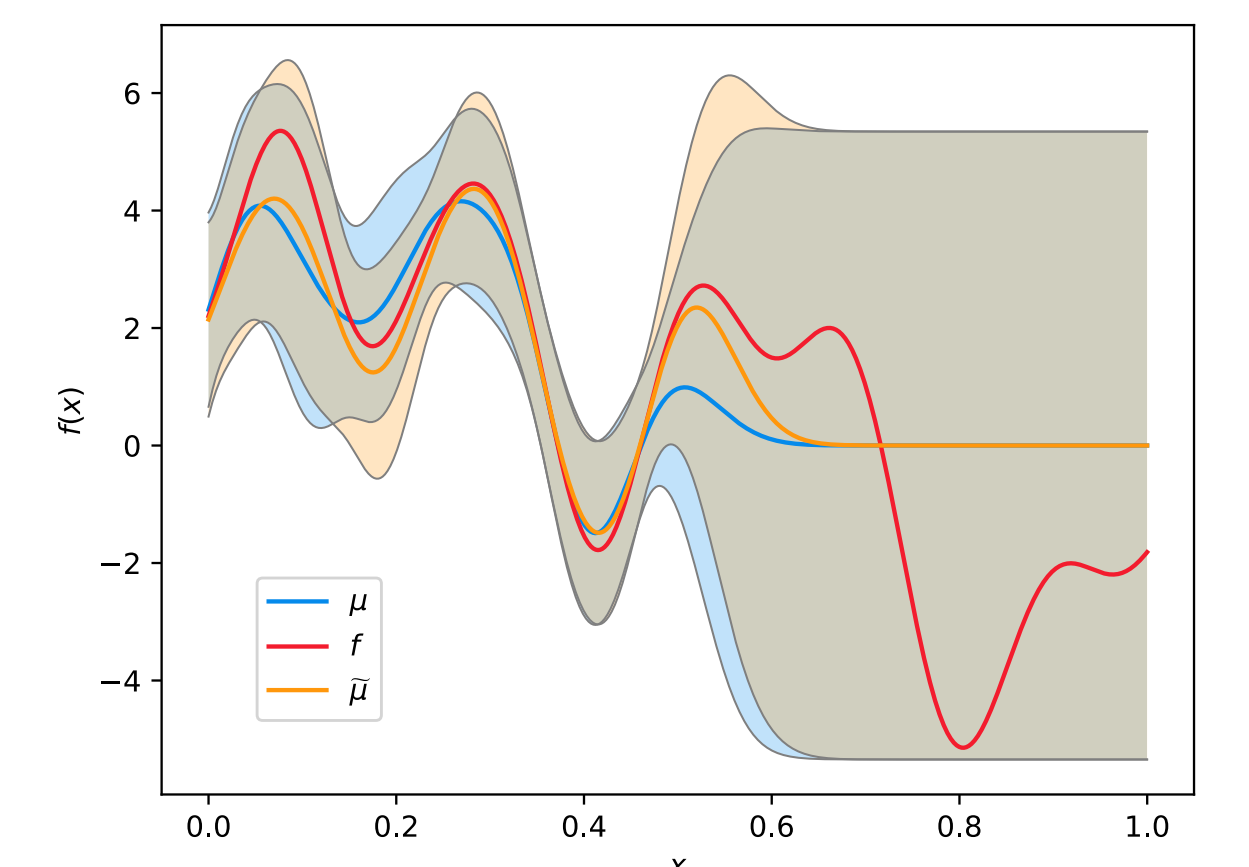
Not simply an approximate GP-UCB

Confidence intervals

$$C_t = [\mu_t(\mathbf{x}) \pm \beta_t \sigma_t(\mathbf{x})],$$

$$\tilde{C}_t = [\tilde{\mu}_t(\mathbf{x}) \pm \tilde{\beta}_t \tilde{\sigma}_t(\mathbf{x})]$$

$$C_t \not\subset \tilde{C}_t, \quad \tilde{C}_t \not\subset C_t$$



Analysis: OFUL [4] + **regularization** with projection \mathbf{P}_t on $\text{Im}(\mathcal{S}_t)$

$$\tilde{\beta}_t \leq \underbrace{\frac{\text{regular. bias}}{\lambda \|f\|}} + \underbrace{\eta_t^\top \mathbf{P}_t \mathbf{K}_t \mathbf{P}_t (\mathbf{P}_t \mathbf{K}_t \mathbf{P}_t + \lambda \mathbf{I})^{-1} \eta_t}_{\text{reduced variance}} + \underbrace{\|(\mathbf{I} - \mathbf{P}_t) \mathbf{K}_t^{1/2} f\|}_{\text{extra approximation bias}}$$

Lemma: $\mathbf{I} - \mathbf{P}_t \preceq (1 + \varepsilon)\lambda(\mathbf{K}_t + \lambda \mathbf{I})^{-1}$ when \mathcal{S}_t sampled $\sigma_t(\cdot)$

$$\|(\mathbf{I} - \mathbf{P}_t) \mathbf{K}_t^{1/2} f\|^2 \leq (1 + \varepsilon)\lambda \|(\mathbf{K}_t + \lambda \mathbf{I})^{-1/2} \mathbf{K}_t^{1/2} f\|^2 \leq (1 + \varepsilon)\lambda \|f\|^2$$

Non-uniform error: **self-normalized bias** focuses on essential parts of \mathcal{A}

😊 no need for uniform bounds, ε -grids, and $\exp\{d\}$ dependencies

[1] Srinivas et al. Gaussian process optimization in the bandit setting: No regret and experimental design. ICML 2010 [2] Mutný et al. Efficient high-dimensional Bayesian optimization with additivity and quadrature Fourier features. NeurIPS 2018

[3] Kuzborskij et al. Efficient linear bandits through matrix sketching. AISTATS 2019 [4] Abbasi-Yadkori et al. Improved algorithms for linear stochastic bandits. NeurIPS 2011 [5] Wang et al. Batched Large-scale Bayesian Optimization in High-dimensional Spaces. AISTATS 2018