Bayesian optimization of time-varying functions

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

1/43

Petar M. Djurić

Department of Electrical Engineering Stony Brook University Stony Book, NY

(In collaboration with Fernando Llorente)

October 30, 2024

Some geography first





Outline of the presentation

Problem statement

- 2 Learning a black-box function with Gaussian processes
 - 1.1 Random-feature Gaussian processes for online learning
 - 1.2 Injecting dynamics via a Markov transition model
 - 1.2 Illustration of forgetting mechanisms
- Bayesian optimization
 - 2.1 Time-varying Bayesian optimization with dynamic random-feature Gaussian processes
 - 2.2 Time-varying acquisition function with windows of data
- 4 Experiments
- 5 Conclusions

Objective: Time-varying optimization (TVO), i.e., solving

$$\boldsymbol{x}_{t}^{\mathsf{opt}} = \arg \max_{\boldsymbol{x} \in \mathcal{X}_{t}} f_{t}(\boldsymbol{x}), \quad f_{t}(\boldsymbol{x}) : \mathcal{X}_{t} \subset \mathbb{R}^{d_{x}} \mapsto \mathbb{R}, \quad t \in \mathbb{N}$$

- Assumptions: $f_t(x)$ is unknown, $f_t(x)$ is continuous for every t, $f_t(x)$ is slowly varying over time.
- Allowance: The estimation of $\mathbf{x}_t^{\text{opt}}$ for t = 1, 2, ... is carried out by probing the function N times at each step.

Previous works of TVO and BO:

C. Cruz, J. R. González, and D. A. Pelta. Optimization in dynamic environments: a survey on problems, methods and measures. *Soft Computing*, 15:1427–1448, 2011.

provides review of non-probabilistic approaches for TVO.

 I. Bogunovic, J. Scarlett, and V. Cevher. Time-varying Gaussian process bandit optimization. *Artificial Intelligence and Statistics*, pp. 314–323. PMLR, 2016.

introduces algorithms for TV-BO, but does not use sparse, online approaches.

F. M. Nyikosa, M. A. Osborne, and S. J. Roberts. Bayesian optimization for dynamic problems. *arXiv preprint arXiv*:1803.03432, 2018.

This paper also presents algorithms for TV-BO, but these algorithms do not use sparse, online approaches.

Q. Lu, K. D. Polyzos, B. Li, and G. B. Giannakis. Surrogate modeling for Bayesian optimization beyond a single Gaussian process. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 2023.

uses sparse, online approaches for BO but *not* in the TVO setting.

S. Van Vaerenbergh, M. Lázaro-Gredilla, and I. Santamaría. Kernel recursive least-squares tracker for time-varying regression. *IEEE Transactions on Neural Networks and Learning Systems*, 23(8):1313–1326, 2012.

considers learning a time-varying function instead of TVO.

Our contributions:

- 1 We use a Bayesian approach for solving the TVO problem based on online Gaussian processes with forgetting.
- We study different forgetting mechanisms in this framework.
- 3 We propose the use of windows of past data to build the time-varying acquisition function.

Some real-world problems where we need to perform TVO

 robot/drone/self-driving vehicles learning in changing environments (adaptations needed because of changing obstacles, shifting weather conditions, or varying traffic. TV-BO helps optimize paths, control parameters, or sensor placements in real time)

smart grids and energy systems (needed for optimization of energy generation and consumption in real-time to balance supply and demand) *– financial markets* (adaptation to changing market conditions;
 TV-BO enables dynamic portfolio optimization)

 wireless networks and communications (used to adapt to changing network traffic, interference patterns, and user demand, resource allocation)

healthcare and medical imaging (adaptive radiation therapy, electrical stimulation of the brain)

aerospace engineering (flight path optimization, satellite trajectory control, rocket launch and re-entry optimization, UAV swarm coordination)

Gaussian processes (GPs) are a model for learning unknown functions

$$f: x \mapsto y$$

◆□> <圖> < E> < E> < E < のへで</p>

11/43

from observed input-outputs $(x_1, y_1), \ldots, (x_N, y_N)$.

GPs are the generalization of Gaussian distributions to functions (infinite dimensional):

 $f \sim GP(m,k)$

where m is the mean function and k is the covariance function.





GPs are attractive for many reasons:

BayesianImage: Constraint of the second second

However, they have drawbacks:

Cost © Kernel design © Given some data y, we aim to infer the unknown θ by computing

 $p(\boldsymbol{\theta}|\boldsymbol{y}) \propto p(\boldsymbol{y}|\boldsymbol{\theta}) \times p(\boldsymbol{\theta})$

posterior \propto likelihood \times prior

In GPs,

$$p(f|\boldsymbol{y}) \propto p(\boldsymbol{y}|f) \times p(f)$$

GP posterior \propto likelihood \times GP prior

Example of GPs

The distribution over plausible functions shrinks as we gather more data

The covariance between function values is expressed by a kernel function, e.g., RBF kernel:

$$\mathsf{Cov}(f(x), f(x')) = k(x, x'; \lambda) = \sigma_f^2 \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

where $\lambda = [\sigma_f, \ell]$ are hyperparameters.

- Covariance functions include strong assumptions about f(x).
- Hyperparameters allow to interpret the data.

Kernels











 $\ell = 0.5$

 $\ell = 0.01$



< □ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ ● ■ のへで 17/43 The performance of GPs depends a lot on the hyperparameters. 😳

Bayesian model selection: choose them by

$$\boldsymbol{\lambda}^* = \arg \max_{\boldsymbol{\lambda}} p(\boldsymbol{y}|\boldsymbol{\lambda}) = \arg \min_{\boldsymbol{\lambda}} - \log p(\boldsymbol{y}|\boldsymbol{\lambda})$$

The derivatives $\nabla_{\lambda} \log p(y|\lambda)$ are available so we can apply gradient descent algorithms. \odot

Learning of hyperparameters

GPs require the inversion of a $N \times N$ kernel matrix

$$\boldsymbol{K} = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_N) \\ \vdots & \ddots & \cdots \\ k(x_N, x_1) & \dots & k(x_N, x_N) \end{bmatrix}$$

The cost scales cubically with the number of datapoints N. Several successful approaches exist for reducing this cost, based on:

- inducing points
- random-feature expansions

Computation costs



Random-feature Gaussian process

When k(x, x') is a (normalized) stationary kernel,

$$k(\boldsymbol{x}, \boldsymbol{x}') = \int e^{i \mathbf{v}^{ op}(\boldsymbol{x} - \boldsymbol{x}')} p(\mathbf{v}) d\mathbf{v}$$

we can sample a set of random Fourier features $\{\mathbf{v}_j\}_{j=1}^J \sim p(\mathbf{v})$ from the power spectral density of the kernel to build the map

$$\phi(\boldsymbol{x}) = \frac{1}{\sqrt{J}} [\sin(\boldsymbol{x}^{\mathsf{T}} \mathbf{v}_1), \cos(\boldsymbol{x}^{\mathsf{T}} \mathbf{v}_1), \dots, \sin(\boldsymbol{x}^{\mathsf{T}} \mathbf{v}_J), \cos(\boldsymbol{x}^{\mathsf{T}} \mathbf{v}_J)]^{\mathsf{T}} \in \mathbb{R}^{2J}$$

and obtain a low-rank approximation of the kernel matrix

$$\mathbf{K} \approx \sigma^2 \mathbf{\Phi}^{\mathsf{T}} \mathbf{\Phi}$$

$$oldsymbol{\Phi}$$
 = $[oldsymbol{\phi}(oldsymbol{x}_1), \dots, oldsymbol{\phi}(oldsymbol{x}_N)] \in \mathbb{R}^{2J imes N}$

For (normalized) RBF kernel

$$k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{x}')^{\mathsf{T}} \Lambda^{-1}(\boldsymbol{x} - \boldsymbol{x}')\right)$$

the power spectral density is

$$p(\mathbf{v}) = \mathcal{N}(\mathbf{0}, \Lambda^{-1})$$

where

$$\Lambda = \mathsf{diag}\left(\left[\ell_1^2, \dots, \ell_d^2\right]\right)$$

◆□ > ◆□ > ◆三 > ◆三 > ・三 のへの

23/43

Random-feature Gaussian process

RF-GP is a Bayesian linear regression model with parameters $\boldsymbol{\theta} \in \mathbb{R}^{2J}$,

$$y = \boldsymbol{\phi}(\boldsymbol{x})^{\mathsf{T}} \boldsymbol{\theta} + \epsilon$$

where

$$p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{0}, \sigma_{\boldsymbol{\theta}}^{2}\mathbf{I}_{2J})$$

$$p(y|\boldsymbol{\theta}, \boldsymbol{x}) = \mathcal{N}(y|\boldsymbol{\phi}(\boldsymbol{x})^{\mathsf{T}}\boldsymbol{\theta}, \sigma_{\mathsf{obs}}^{2})$$

$$p(\boldsymbol{\theta}|\boldsymbol{y}, \boldsymbol{X}) \propto p(\boldsymbol{\theta}) \prod_{i=1}^{N} p(y_{i}|\boldsymbol{\theta}, \boldsymbol{x}_{i}) \propto \mathcal{N}(\boldsymbol{\theta}|\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Sigma}})$$

so

$$p(y_{\star}|\boldsymbol{y},\boldsymbol{x}_{\star}) = \int p(y_{\star}|\boldsymbol{\theta},\boldsymbol{x}_{\star}) p(\boldsymbol{\theta}|\boldsymbol{y},\boldsymbol{X}) d\boldsymbol{\theta} = \mathcal{N}(y_{\star}|\widehat{y}(\boldsymbol{x}_{\star}),\widehat{\sigma}^{2}(\boldsymbol{x}_{\star}))$$

Bayesian update: prediction + update

$$p(\theta|y_{1:t+1}, X_{1:t+1}) = \frac{p(y_{t+1}|\theta, x_{t+1})p(\theta|y_{1:t}, X_{1:t})}{p(y_{t+1}|y_{1:t}, X_{1:t+1})}$$

is possible in closed-form with constant per-step computation cost.

Prediction step. Let $\Phi_t = [\phi(x_1) \dots \phi(x_t)] \in \mathbb{R}^{2J \times t}$. The predictive density is

$$p(y_{t+1}|\mathcal{D}_t, \boldsymbol{x}_{t+1}) = \mathcal{N}(y_{t+1}|\widehat{y}_{t+1}, \widehat{\sigma}_{t+1}^2)$$
$$\widehat{y}_{t+1} = \boldsymbol{\phi}(\boldsymbol{x}_{t+1})^{\mathsf{T}} \widehat{\boldsymbol{\theta}}_t$$
$$\widehat{\sigma}_{t+1}^2 = \boldsymbol{\phi}(\boldsymbol{x}_{t+1})^{\mathsf{T}} \widehat{\boldsymbol{\Sigma}}_t \boldsymbol{\phi}(\boldsymbol{x}_{t+1}) + \sigma_{\mathsf{obs}}^2$$

where $\widehat{\Sigma}_{t} = \left(\frac{1}{\sigma_{obs}^{2}} \Phi_{t} \Phi_{t}^{T} + \frac{1}{\sigma_{\theta}^{2}} \mathbf{I}_{2J}\right)^{-1}$ and $\widehat{\theta}_{t} = \frac{1}{\sigma_{obs}^{2}} \widehat{\Sigma}_{t} \Phi_{t} \boldsymbol{y}_{t}$ are the moments of the posterior $p(\boldsymbol{\theta}|\mathcal{D}_{t}) = \mathcal{N}(\boldsymbol{\theta}|\widehat{\theta}_{t}, \widehat{\Sigma}_{t})$.

<ロト<部ト<E><E> E のQの 26/43 **Update step.** Upon receiving (x_{t+1}, y_{t+1}) , the update formulae are

$$\widehat{\boldsymbol{\theta}}_{t+1} = \widehat{\boldsymbol{\theta}}_t + \widehat{\sigma}_{t+1}^{-2} \widehat{\boldsymbol{\Sigma}}_t \boldsymbol{\phi}(\boldsymbol{x}_{t+1}) (\boldsymbol{y}_{t+1} - \widehat{\boldsymbol{y}}_{t+1})$$
$$\widehat{\boldsymbol{\Sigma}}_{t+1} = \widehat{\boldsymbol{\Sigma}}_t - \widehat{\sigma}_{t+1}^{-2} \widehat{\boldsymbol{\Sigma}}_t \boldsymbol{\phi}(\boldsymbol{x}_{t+1}) \boldsymbol{\phi}(\boldsymbol{x}_{t+1})^{\mathsf{T}} \widehat{\boldsymbol{\Sigma}}_t$$

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ _ 圖 _ 釣�?

27/43

- We consider a general Markov transition model where, before receiving the t + 1 datapoint, θ_{t+1} is a linear combination of the parameter at time t, a drift term, and random noise.
- Letting $\theta_0 \sim p_0(\theta_0)$. For t = 0, 1, ..., the parameter vector evolves as

$$\boldsymbol{\theta}_{t+1} = a_1 \boldsymbol{\theta}_t + a_2 \mathbf{m}_t + a_3 \mathbf{u}_t$$

where $\theta_t \sim p(\theta_t | D_t)$, \mathbf{m}_t is a time-varying vector, and \mathbf{u}_t is a noise process.

We can employ popular "forgetting" mechanisms

B2P:
$$\boldsymbol{\theta}_{t+1} = \sqrt{\lambda} \, \boldsymbol{\theta}_t + \sqrt{1-\lambda} \, \mathbf{u}_t, \qquad \mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \sigma_{\theta}^2 \mathbf{I}_{2J}), \quad \lambda \in [0, 1]$$

UI: $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \sqrt{\frac{1}{\gamma} - 1} \mathbf{u}_t, \qquad \mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \widehat{\boldsymbol{\Sigma}}_t), \quad \gamma \in (0, 1]$

"Back to Prior" forgetting

"Uncertainty Injection" forgetting

(Time-invariant) Bayesian optimization solves

$$oldsymbol{x}_{\star}$$
 = arg $\max_{oldsymbol{x}\in\mathcal{X}}f(oldsymbol{x})$

using a probabilistic surrogate p(f|D) and an acquisition function $\alpha(x; D)$.

BO sequentially acquires points that are "useful" by repeatedly doing:

1 $\boldsymbol{x}_{t+1} = \arg \max_{\boldsymbol{x} \in \mathcal{X}} \alpha(\boldsymbol{x}; \mathcal{D}_t)$

$$2 \quad \mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup (\boldsymbol{x}_{t+1}, y_{t+1})$$



Vanilla BO with UCB acquisition function

Time-invariant BO on a time-varying function

TVBO with UI

Bayesian optimization with dynamic RF-GP

The model

$$\theta_t \sim p(\theta_t | \mathcal{D}_t)$$

$$\theta_{t+1} = a_1 \theta_t + a_2 \mathbf{m}_t + a_3 \mathbf{u}_t$$

$$y_{t+1} = \phi^{\mathsf{T}}(\boldsymbol{x}) \theta_{t+1} + \epsilon_{\mathsf{obs}}$$

 For instance, the GP-UCB acquisition function with B2P forgetting is

$$\begin{aligned} \alpha_{t+1}^{UCB}(\boldsymbol{x}) &= \sqrt{\lambda} \boldsymbol{\phi}(\boldsymbol{x})^{\mathsf{T}} \widehat{\boldsymbol{\theta}}_{t} \\ &+ \beta_{t+1} \sqrt{\lambda \boldsymbol{\phi}(\boldsymbol{x})^{\mathsf{T}} \widehat{\boldsymbol{\Sigma}}_{t} \boldsymbol{\phi}(\boldsymbol{x}) + (1-\lambda) \sigma_{\theta}^{2} \boldsymbol{\phi}(\boldsymbol{x})^{\mathsf{T}} \boldsymbol{\phi}(\boldsymbol{x})} \end{aligned}$$

Bayesian optimization with dynamic RF-GP

Setting
$$a_1 = \sqrt{\lambda}$$
, $a_2 = 1 - \sqrt{\lambda}$, and $a_3 = \sqrt{1 - \lambda}$ with $\lambda \in (0, 1)$, we get

$$\alpha_{t+1}^{UCB}(\boldsymbol{x}) = \sqrt{\lambda}\phi(\boldsymbol{x})^{\mathsf{T}}\widehat{\boldsymbol{\theta}}_{t} + (1 - \sqrt{\lambda})\phi(\boldsymbol{x})^{\mathsf{T}}\mathbf{m}_{t} \\ + \beta_{t+1}\sqrt{\lambda\phi(\boldsymbol{x})^{\mathsf{T}}\widehat{\boldsymbol{\Sigma}}_{t}\phi(\boldsymbol{x}) + (1 - \lambda)\phi(\boldsymbol{x})^{\mathsf{T}}\mathbf{C}_{t}\phi(\boldsymbol{x})}.$$

Time-varying BO with windows of data

We propose to use a sliding window of w samples $\mathcal{D}_t^w = \{(x_\tau, y_\tau)\}_{\tau=t-w+1}^t$ to compute the parameters $\{\mathbf{m}_t, \mathbf{C}_t\}$,

$$\mathbf{m}_{t} = \left(\mathbf{\Phi}_{t}^{w} (\mathbf{\Phi}_{t}^{w})^{\mathsf{T}} + \frac{\sigma_{\mathsf{obs}}^{2}}{\sigma_{\theta}^{2}} \mathbf{I}_{2J}\right)^{-1} \mathbf{\Phi}_{t}^{w} \boldsymbol{y}_{t}^{w} + \mathbf{C}_{t} = \left(\frac{1}{\sigma_{\mathsf{obs}}^{2}} \mathbf{\Phi}_{t}^{w} (\mathbf{\Phi}_{t}^{w})^{\mathsf{T}} + \frac{1}{\sigma_{\theta}^{2}} \mathbf{I}_{2J}\right)^{-1},$$

where $\mathbf{\Phi}_t^w$ and \mathbf{y}_t^w are built using \mathcal{D}_t^w .

Note that:

when $w = 1 \implies \mathbf{m}_t = 0$ and $\mathbf{C}_t = \sigma_{\theta}^2 \mathbf{I}_{2J} \implies$ we recover BO with prior forgetting.

when $w = t \Longrightarrow \mathbf{m}_t = \widehat{\boldsymbol{\theta}}_t$ and $\mathbf{C}_t = \widehat{\boldsymbol{\Sigma}}_t \Longrightarrow$ we recover time-invariant BO.

Time-varying BO with windows of data

TVBO with "Back 2 Window"

Let $x \in [0, 1]$

$$f_t(x) = \sin[2\pi(x+0.15-0.01t)]$$

where t goes from t = 1 to t = 85.

We compare the average regret,

$$\bar{R}_T = \frac{1}{T} \sum_{t=1}^T \left| f_t^* - f_t(\widehat{x}_t) \right|$$

- where $f_t^* = \max_x f_t(x) = 1$ for all t.
- Goal: study window size w and forgetting strength λ.



The optimal value of w increased as λ decreased: when we allow the current model to forget more and more (smaller λ), it is beneficial that the forgetting factor carries an increasing amount of past data (bigger w). However, for all λ , increasing w eventually hampers performance.

Localization example

■ Target Position *x* ∈ ℝ². The network has *L* = 3 sensors. We observed *M* measurements at every sensor,

$$z_{m,l} \sim \mathcal{N}\left(A\log\left(\parallel \boldsymbol{x} - \mathbf{s}_l \parallel\right), \sigma_z^2\right),$$

where $l = 1, 2, 3, m = 1, \dots, M$, and A is a constant.

 Goal: track target over time by estimating the MAP of the sequence of posterior distributions,

$$\boldsymbol{x}_{\mathsf{map},t} = \arg \max_{\boldsymbol{x}} \bar{\pi}(\boldsymbol{x}|\mathbf{Z}_t),$$

where \mathbf{Z}_t are measurements obtained at $t = 0, \dots, T$.



- We have presented a framework for time-varying optimization based on dynamic random-feature Gaussian processes.
- Random-feature Gaussian processes have several advantages such as scalability and capacity for online learning. Moreover, time dynamics can be injected by considering a Markov transition model.
- Combining random-feature Gaussian processes with a general Markov time evolution, we have proposed a framework for time-varying Bayesian optimization and studied the resulting time-varying acquisition function.
- More specifically, we proposed a novel time-varying Bayesian optimization algorithm that uses windows of past data to build the additional exploration and exploitation terms. The time-invariant setting, as well as time-varying Bayesian optimization with prior forgetting, are special cases.