

LEARNING GPs WITH BAYESIAN POSTERIOR OPTIMIZATION

Luiz F. O. Chamon, Santiago Paternain, and Alejandro Ribeiro

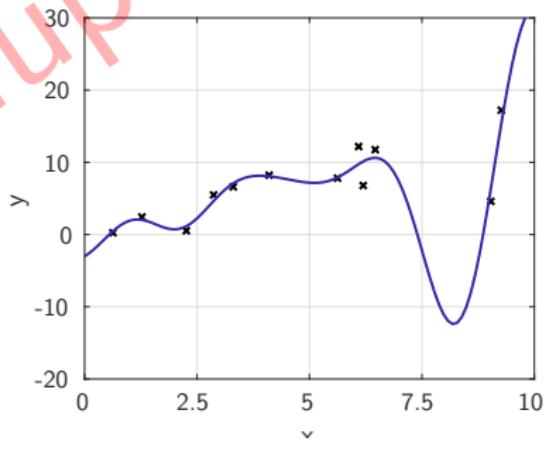
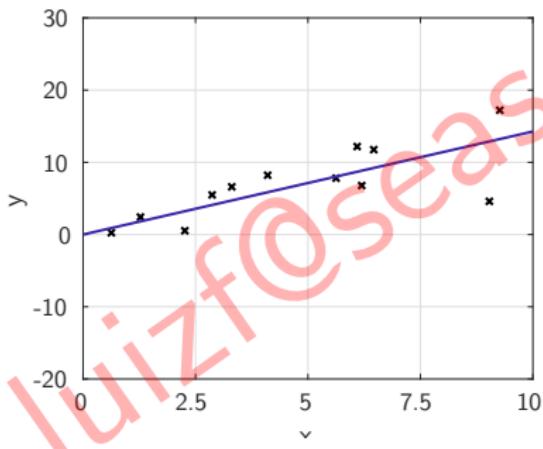
ASILOMAR 2019
November 4th, 2019

Dealing with complexity and uncertainty

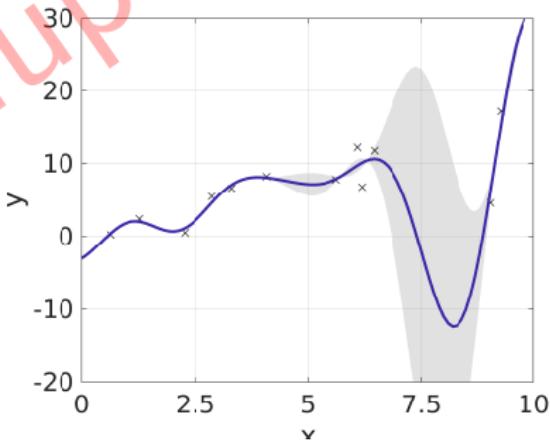
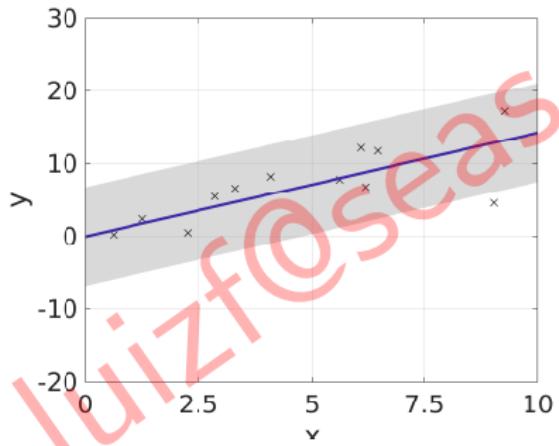


luizf@seas.upenn.edu

- ▶ Nonparametric methods



- ▶ Nonparametric methods
- ▶ Bayesian methods



- ▶ All of Bayesian inference

$$\underbrace{\mathbb{P}(\text{model})}_{\text{Prior}} + \underbrace{\mathbb{P}(\text{data} \mid \text{model})}_{\text{Likelihood}} \rightarrow \underbrace{\mathbb{P}(\text{model} \mid \text{data})}_{\text{Posterior}}$$

- ▶ All of Bayesian inference

$$\underbrace{\mathbb{P}(\text{model})}_{\text{Prior}} + \underbrace{\mathbb{P}(\text{data} \mid \text{model})}_{\text{Likelihood}} \rightarrow \underbrace{\mathbb{P}(\text{model} \mid \text{data})}_{\text{Posterior}}$$

- ▶ Parametric models are finite dimensional

$$\mathbb{P}(\text{model}) = \mathbb{P}(\text{parameters})$$

- ▶ Nonparametric models are infinite dimensional

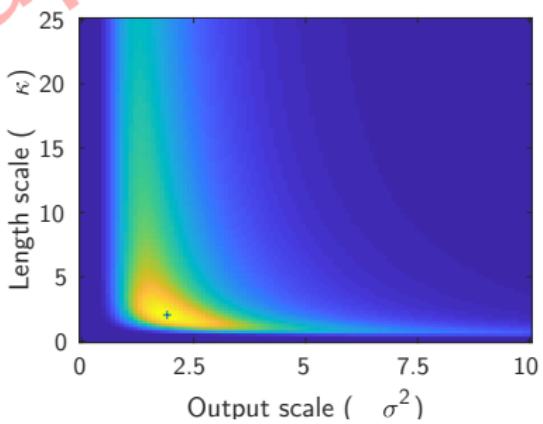
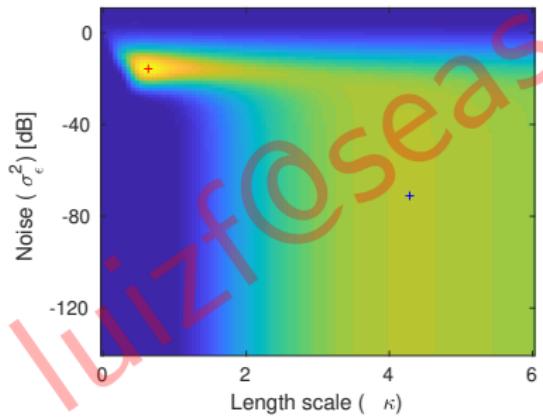
- ▶ GPs are priors on “smooth” functions
 - ✓ easy to specify: choose a covariance function (and hyperparameters)
 - ✓ flexible: wide variety of covariance functions (degree of smoothness, periodicity...)
 - ✓ tractable

- ▶ GPs are priors on “smooth” functions
 - ✓ easy to specify: choose a covariance function (and hyperparameters)
 - ✓ flexible: wide variety of covariance functions (degree of smoothness, periodicity...)
 - ✓ tractable
 - ▶ Still... which GP?
 - ✗ limited access to prior knowledge
 - ✗ hard to interpret hyperparameters
 - ✗ misspecifying GPs can be catastrophic
- [Bachoc'13, Beckers et al.'18, Zaytsev et al.'18]

Which GP?

- ▶ Maximize likelihood w.r.t. hyperparameters [Stein'99, RW'06]
 - Ambiguity: multimodal likelihood, local maxima
 - Indeterminacy: different parameters, same measure

$$\theta^* = \operatorname{argmax}_{\theta} \log \mathbb{P}(y | X, \theta)$$



- ▶ Maximize likelihood w.r.t. hyperparameters [Stein'99, RW'06]
 - Ambiguity: multimodal likelihood, local maxima
 - Indeterminacy: different parameters, same measure

$$\theta^* = \operatorname{argmax}_{\theta} \log \mathbb{P}(y | X, \theta)$$

- ▶ Hierarchical models [RG'02, RW'06, Gelman et al.'13]
 - Noninformative priors → improper posteriors
 - Hard to interpret, hard to set priors
 - Indeterminacy: setting one prior affects the others

$$\theta \sim \mathcal{P}$$

- ▶ Hybrid Bayesian–Optimization approach
 - **Bayesian:**
 - **Optimization:**

- ▶ Hybrid Bayesian–Optimization approach
 - **Bayesian:** obtain distribution over hyperparameters instead of point estimate
 - **Optimization:**

- ▶ Hybrid Bayesian–Optimization approach
 - **Bayesian:** obtain distribution over hyperparameters instead of point estimate
 - **Optimization:** minimize a risk measure instead of using Bayes rule

- ▶ Hybrid Bayesian–Optimization approach
 - **Bayesian:** obtain distribution over hyperparameters instead of point estimate
 - **Optimization:** minimize a risk measure instead of using Bayes rule
- ✓ Non-convex risk measures (0-1 loss, truncated MSE...)
- ✓ Incorporate complex structures in the prior:
maximum entropy, sparsity, moments...

- ▶ Hybrid Bayesian–Optimization approach
 - **Bayesian:** obtain distribution over hyperparameters instead of point estimate
 - **Optimization:** minimize a risk measure instead of using Bayes rule
- ✓ Non-convex risk measures (0-1 loss, truncated MSE...)
- ✓ Incorporate complex structures in the prior:
maximum entropy, sparsity, moments...
- ✗ Non-convex, infinite dimensional optimization problem

- ▶ Hybrid Bayesian–Optimization approach
 - **Bayesian:** obtain distribution over hyperparameters instead of point estimate
 - **Optimization:** minimize a risk measure instead of using Bayes rule

- ✓ Non-convex risk measures (0-1 loss, truncated MSE...)
- ✓ Incorporate complex structures in the prior:
maximum entropy, sparsity, moments...
- ✓ Non-convex, infinite dimensional optimization problem
⇒ *simple, efficient solution using duality*

The Bayesian part

The optimization part

Solving Bayesian posterior optimization problems

luizf@seas.upenn.edu

- ▶ **Data:** (\mathbf{x}_i, y_i) with $y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma_\epsilon^2)$ for an unknown f
- ▶ **Goal:** determine $(f \mid \mathbf{X}, \mathbf{y})$
- ▶ **How?** Bayes' rule and GP prior

- ▶ A GP is a stochastic process whose finite dimensional marginals are jointly Gaussian
- ▶ Formally, $\mathbb{GP}(m, k)$ is a distribution over functions g such that $[g(\mathbf{x}_1) \ \cdots \ g(\mathbf{x}_n)] \sim \mathcal{N}(\mathbf{m}, \mathbf{K})$ for all $n \in \mathbb{N}$

$$\mathbf{m} = \begin{bmatrix} m(\mathbf{x}_1) \\ \vdots \\ m(\mathbf{x}_n) \end{bmatrix} \text{ and } \mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

- ▶ Typically, $m \equiv 0$

- ▶ **Data:** (\mathbf{x}_i, y_i) with $y_i \sim \mathcal{N}(f(\mathbf{x}_i), \sigma_\epsilon^2)$ for an unknown f
- ▶ **Goal:** determine $(f \mid \mathbf{X}, \mathbf{y})$
- ▶ **How?** Bayes' rule and GP prior

$$(f(\bar{\mathbf{x}}) \mid \mathbf{X}, \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\mu} = \bar{\mathbf{k}}^T \mathbf{K}^{-1} \mathbf{y}, \quad \boldsymbol{\Sigma} = k(\bar{\mathbf{x}}, \bar{\mathbf{x}}) - \bar{\mathbf{k}}^T \mathbf{K}^{-1} \bar{\mathbf{k}}$$

$$\bar{\mathbf{k}} = \begin{bmatrix} k(\bar{\mathbf{x}}, \mathbf{x}_1) \\ \vdots \\ k(\bar{\mathbf{x}}, \mathbf{x}_n) \end{bmatrix} \text{ and } \mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \cdots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \cdots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

- ▶ First level: unknown function f

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, f) \mathbb{P}(f \mid \mathbf{X})$$

- ▶ First level: unknown function f

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, f, \boldsymbol{\theta}) \mathbb{P}(f \mid \mathbf{X}, \boldsymbol{\theta})$$

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$

- ▶ First level: unknown function f

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, f, \boldsymbol{\theta}) \mathbb{P}(f \mid \mathbf{X}, \boldsymbol{\theta})$$

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \underbrace{\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})}_{\theta\text{-posterior}} d\boldsymbol{\theta}$$

- ▶ First level: unknown function f

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, f, \boldsymbol{\theta}) \mathbb{P}(f \mid \mathbf{X}, \boldsymbol{\theta})$$

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \underbrace{\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})}_{\theta\text{-posterior}} d\boldsymbol{\theta}$$

- ▶ Second level: hyperparameters $\boldsymbol{\theta}$

$$\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) \mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X})$$

- ▶ First level: unknown function f

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, f, \boldsymbol{\theta}) \mathbb{P}(f \mid \mathbf{X}, \boldsymbol{\theta})$$

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \underbrace{\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})}_{\boldsymbol{\theta}\text{-posterior}} d\boldsymbol{\theta}$$

- ▶ Second level: hyperparameters $\boldsymbol{\theta}$

$$\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) \mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X})$$

- ▶ **Issue:** choosing $\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X})$
(interpretation, indeterminacy, informativeness)

- ▶ First level: unknown function f

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, f, \boldsymbol{\theta}) \mathbb{P}(f \mid \mathbf{X}, \boldsymbol{\theta})$$

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \underbrace{\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})}_{\theta\text{-posterior}} d\boldsymbol{\theta}$$

- ▶ Second level: hyperparameters $\boldsymbol{\theta}$

$$\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) \propto \mathbb{P}(\mathbf{y} \mid \mathbf{X}, \boldsymbol{\theta}) \mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X})$$

- ▶ **Issue:** choosing $\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X})$
(interpretation, indeterminacy, informativeness)

The Bayesian part

The optimization part

Solving Bayesian posterior optimization problems

luizf@seas.upenn.edu

► Statistical learning:

$$\phi^* = \operatorname{argmin}_{\phi \in \mathcal{F}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(\phi(\mathbf{x}), y)] + R(\phi)$$

- \mathcal{D} is an *unknown* probability distribution over pairs (\mathbf{x}, y)
- $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}_+$ is a loss function
- R is a regularizer
- \mathcal{F} is a space of functions $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$

► **Statistical learning:**

$$\phi^* = \operatorname{argmin}_{\phi \in \mathcal{F}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(\phi(\mathbf{x}), y)] + R(\phi)$$

► **Empirical risk minimization:**

$$\hat{\phi}^* = \operatorname{argmin}_{\phi \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell(\phi(\mathbf{x}_i), y_i) + R(\phi)$$

- \mathcal{D} is an *unknown* probability distribution over pairs (\mathbf{x}, y)
- $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}_+$ is a loss function
- R is a regularizer
- \mathcal{F} is a space of functions $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$
- Data: $(\mathbf{x}_i, y_i) \sim \mathcal{D}$

► **Statistical GP learning:**

$$\Gamma^* = \operatorname{argmin}_{\Gamma \in \mathcal{GP}} \mathbb{E}_{(x,y) \sim \mathcal{D}, f \sim \Gamma} [\ell(f(\boldsymbol{x}), y)] + R(\gamma)$$

► **Empirical GP-risk minimization:**

$$\hat{\Gamma}^* = \operatorname{argmin}_{\Gamma \in \mathcal{GP}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{f \sim \Gamma} [\ell(f(\boldsymbol{x}_i), y_i)] + R(\gamma)$$

- \mathcal{D} is an *unknown* probability distribution over pairs (\boldsymbol{x}, y)
- $\ell : \mathbb{R}^2 \rightarrow \mathbb{R}_+$ is a loss function
- R is a regularizer
- \mathcal{GP} is the “space of Gaussian processes”
- Data: $(\boldsymbol{x}_i, y_i) \sim \mathcal{D}$

- ▶ **Empirical GP-risk minimization:**

$$\hat{\Gamma}^* = \operatorname{argmin}_{\Gamma \in \mathcal{GP}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{f \sim \Gamma} [\ell(f(\mathbf{x}_i), y_i)] + R(\gamma)$$

- ▶ Challenge: optimizing over \mathcal{GP}
(isomorphic to the space of positive semi-definite functions)

▶ **Empirical GP-risk minimization:**

$$\hat{\Gamma}^* = \operatorname{argmin}_{\Gamma \in \mathcal{GP}} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{f \sim \Gamma} [\ell(f(\mathbf{x}_i), y_i)] + R(\gamma)$$

- ▶ Challenge: optimizing over \mathcal{GP}
(isomorphic to the space of positive semi-definite functions)
- ▶ Leverage the first level of the hierarchical model

$$\mathbb{P}(f \mid \mathbf{X}, \mathbf{y}) = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) \mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y}) d\boldsymbol{\theta}$$

- ▶ "Parameterize" (a subset of) \mathcal{GP} using $\mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$

- ▶ Bayesian posterior optimization:

$$p^* = \operatorname{argmin}_{p \in \mathcal{P}} \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \text{GP}(0, k_{\theta})} [\ell(f(\mathbf{x}_i), y_i)] p(\boldsymbol{\theta}) d\boldsymbol{\theta} + R(p)$$
$$\hat{\Gamma}_p^* = \int \mathbb{P}(f \mid \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}) p^*(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

- ▶ Optimization variable: $p(\boldsymbol{\theta}) = \mathbb{P}(\boldsymbol{\theta} \mid \mathbf{X}, \mathbf{y})$
- ▶ Alternative interpretation: mixture of GPs with weights $p(\boldsymbol{\theta})$
- ✖ Non-convex, infinite dimensional optimization problem

The Bayesian part

The optimization part

Solving Bayesian posterior optimization problems

luizf@seas.upenn.edu

- ▶ Bayesian posterior optimization:

$$\underset{p \in \mathcal{P}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \text{GP}(0, k_{\theta})} [\ell(f(x_i), y_i)] p(\theta) d\theta + R(p)$$

- ▶ Bayesian posterior optimization:

$$\underset{p \in \mathcal{P}}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \text{GP}(0, k_{\theta})} [\ell(f(x_i), y_i)] p(\theta) d\theta + R(p)$$

- ▶ Measure p is non-atomic and absolutely continuous

- ▶ Bayesian posterior optimization:

$$\underset{p \in L_1^+}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \mathbb{GP}(0, k_{\theta})} [\ell(f(\mathbf{x}_i), y_i)] p(\boldsymbol{\theta}) d\boldsymbol{\theta} + R(p)$$

subject to $\int p(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1$

- ▶ Measure p is non-atomic and absolutely continuous

- ▶ Bayesian posterior optimization:

$$\begin{aligned} & \underset{p \in L_1^+}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \mathbb{GP}(0, k_{\theta})} [\ell(f(\mathbf{x}_i), y_i)] p(\boldsymbol{\theta}) d\boldsymbol{\theta} + R(p) \\ & \text{subject to} \quad \int p(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1 \end{aligned}$$

- ▶ Measure p is non-atomic and absolutely continuous
- ▶ R is a separable functional

- ▶ Bayesian posterior optimization:

$$\begin{aligned} \underset{p \in L_1^+}{\text{minimize}} \quad & \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \text{GP}(0, k_\theta)} [\ell(f(\mathbf{x}_i), y_i)] p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ & + \sum_{j=1}^m \lambda_j \int h_j [p(\boldsymbol{\theta}), \boldsymbol{\theta}] d\boldsymbol{\theta} \end{aligned}$$

subject to $\int p(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1$

- ▶ Measure p is non-atomic and absolutely continuous
- ▶ R is a separable functional

► Bayesian posterior optimization:

$$\begin{aligned} \underset{p \in L_1^+}{\text{minimize}} \quad & \frac{1}{n} \sum_{i=1}^n \int \mathbb{E}_{f \sim \text{GP}(0, k_\theta)} [\ell(f(\mathbf{x}_i), y_i)] p(\boldsymbol{\theta}) d\boldsymbol{\theta} \\ & + \sum_{j=1}^m \lambda_j \int h_j [p(\boldsymbol{\theta}), \boldsymbol{\theta}] d\boldsymbol{\theta} \end{aligned}$$

subject to $\int p(\boldsymbol{\theta}) d\boldsymbol{\theta} = 1$

- Measure p is non-atomic and absolutely continuous
- R is a separable functional
- ℓ and h_j are (possibly non-convex) normal integrands

- ▶ Strong duality
(BPO problem is an SFP [Chamon'19])

- ▶ Strong duality
(BPO problem is an SFP [Chamon'19])
- ▶ Exchangeability of the infimum and integral operators
(normal integrand + separability of L_1^+ [Rockafellar'76])

- ▶ Strong duality
(BPO problem is an SFP [Chamon'19])
- ▶ Exchangeability of the infimum and integral operators
(normal integrand + separability of L_1^+ [Rockafellar'76])
- ▶ Lagrangian can be minimized efficiently, often in closed-form
(separability + Gaussian integrals)

Optimizing posteriors

$$1) \bar{\ell}(\boldsymbol{\theta}) = \int \left[\frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i), y_i) \right] \underbrace{\mathbb{P}(f(\bar{\mathbf{x}}) \mid \boldsymbol{\theta})}_{\mathcal{N}(f(\bar{\mathbf{x}}) \mid \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})} df$$

[Gauss-Hermite quadrature]

$$2) p_d(\boldsymbol{\theta}, \mu) = \operatorname{argmin}_{p \geq 0} (\bar{\ell}(\boldsymbol{\theta}) + \mu)p + \sum_{j=1}^m \lambda_j h_j(p, \boldsymbol{\theta})$$

[(often) closed-form]

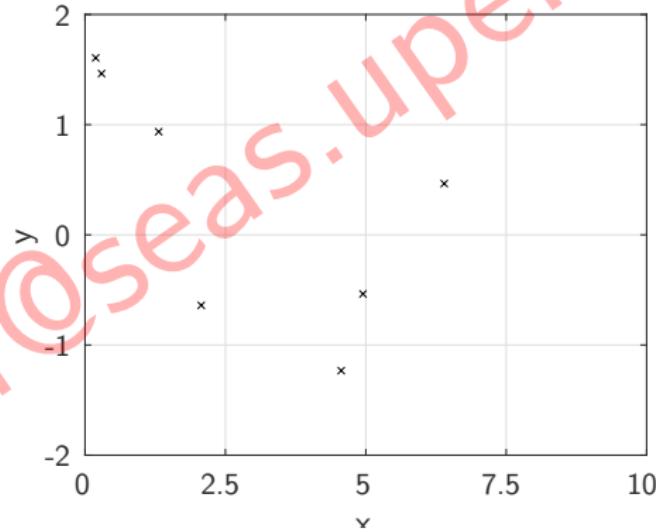
$$3) \mu^* = \operatorname{argmax}_{\mu \in \mathbb{R}} \mathcal{L}(p_d(\boldsymbol{\theta}, \mu), \mu)$$

[SGD or PBA]

$$4) p^*(\boldsymbol{\theta}) = p_d(\mu^*, \boldsymbol{\theta})$$

- GP (RBF): $\sigma^2 = 1$, $\kappa = 1$, and $\sigma_\epsilon^2 = 10^{-1}$

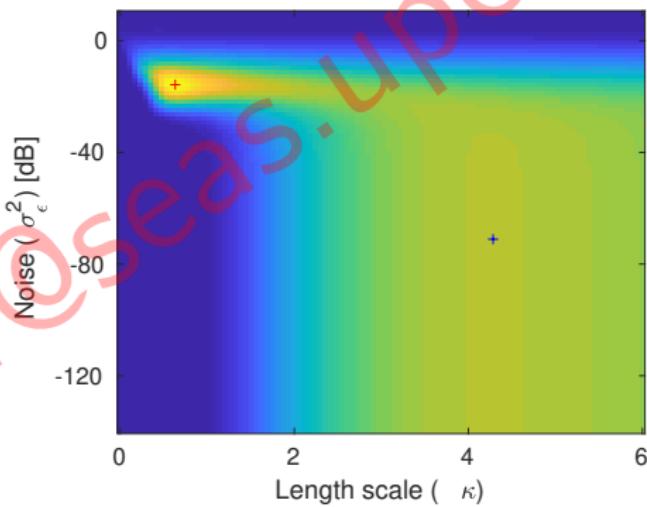
$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left[-\kappa \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2} \right] + \sigma_\epsilon^2 \mathbf{I}$$



Numerical examples

- GP (RBF): $\sigma^2 = 1$, $\kappa = 1$, and $\sigma_\epsilon^2 = 10^{-1}$

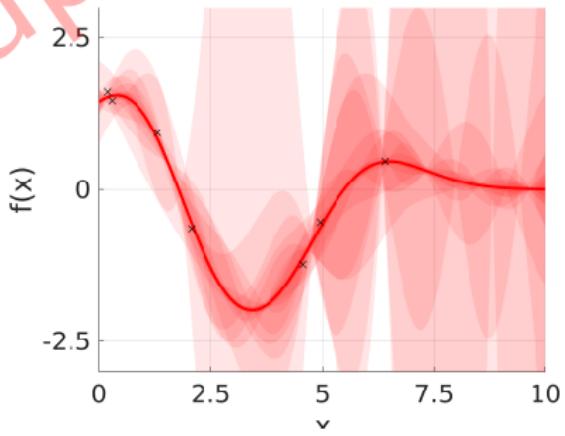
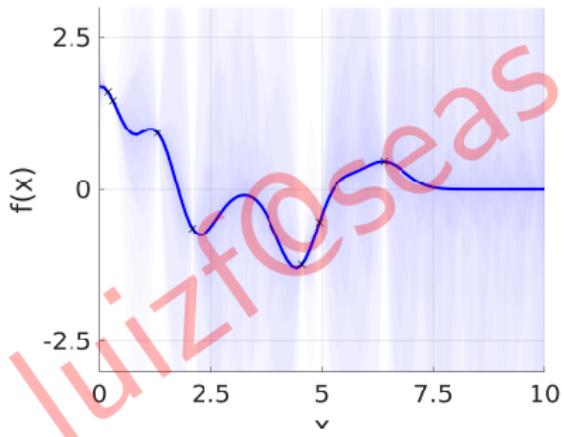
$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left[-\kappa \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2} \right] + \sigma_\epsilon^2 \mathbf{I}$$



Numerical examples

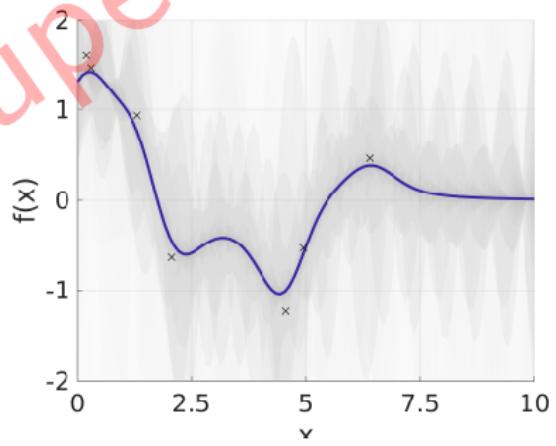
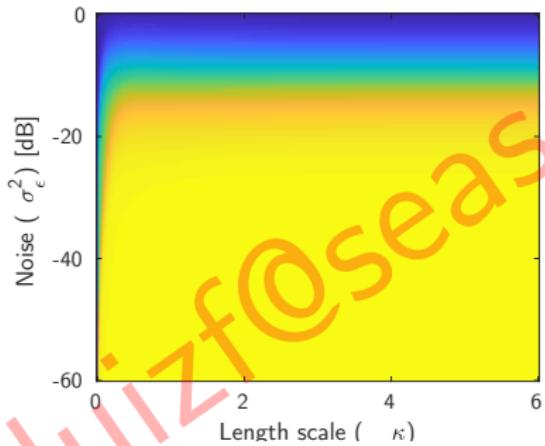
- GP (RBF): $\sigma^2 = 1$, $\kappa = 1$, and $\sigma_\epsilon^2 = 10^{-1}$

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp \left[-\kappa \frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2} \right] + \sigma_\epsilon^2 \mathbf{I}$$



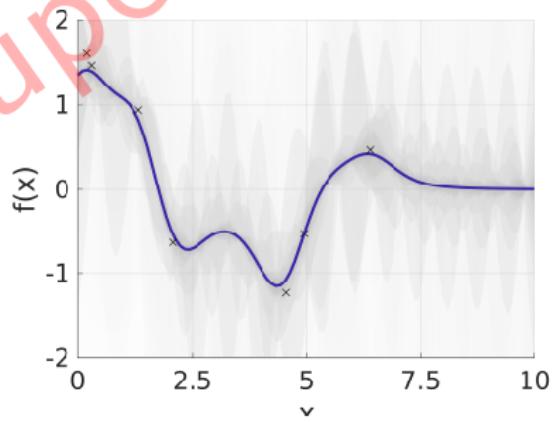
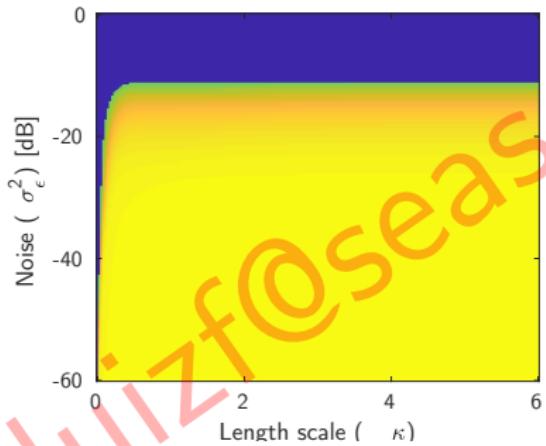
Numerical examples

- ▶ Loss: ℓ_2 -norm
- ▶ Regularization: negative entropy



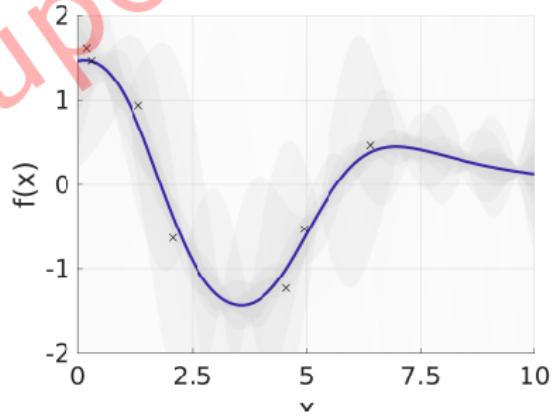
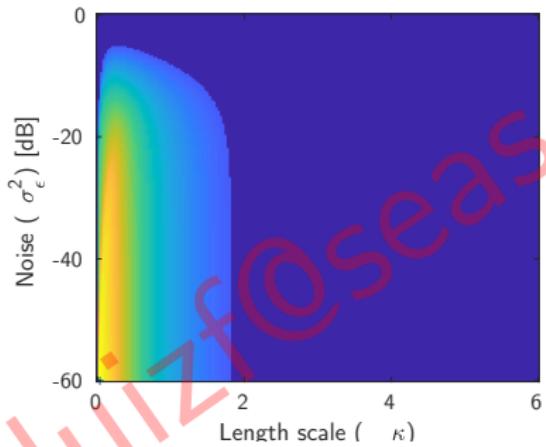
Numerical examples

- ▶ Loss: ℓ_2 -norm
- ▶ Regularization: negative entropy + L_0



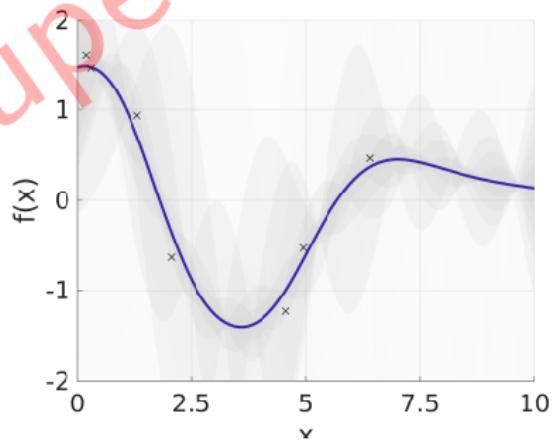
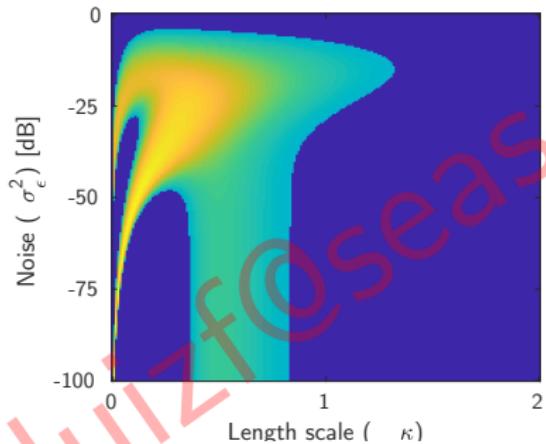
Numerical examples

- ▶ Loss: ℓ_2 -norm
- ▶ Regularization: negative entropy + L_0 + $\mathbb{E}[\kappa]$



Numerical examples

- ▶ Loss: Leave-one-out ℓ_2 -norm
- ▶ Regularization: negative entropy + L_0



- ▶ Priors for nonparametric Bayesian methods are hard to specify and learning them from data is challenging
- ▶ *Bayesian posterior optimization*: replace the prior by a statistical optimization problem
- ▶ Despite the non-convexity and infinite dimensionality, posterior optimization problems can be solved efficiently

LEARNING GPs WITH BAYESIAN POSTERIOR OPTIMIZATION

Luiz F. O. Chamon, Santiago Paternain, and Alejandro Ribeiro

ASILOMAR 2019
November 4th, 2019