### Gaussian Processes (GPs)

- task.

 $f_X \sim N(\mu_X, K_{XX})$  where  $f_X \in \mathbb{R}^n$ ;  $(f_X)_i = f(x_i)$  $y_X \in \mathbb{R}^n$ ;  $(\mathbf{y}_X)_i = \mathbf{y}_i$  $\mu_X \in \mathbb{R}^n$ ;  $(\mu_X)_i = \mu(x_i)$  $K_{XX} \in \mathbb{R}^{n imes n}$ ;

Write  $K_{XX}$  as K when unambiguous.





### **GP Regression with Derivatives**

Model function values and derivatives by a multi-output GP:

$$\begin{bmatrix} f_X \\ \nabla f_X \end{bmatrix} \sim \mathcal{N}(\mu_X^{\nabla}, \mathcal{K}_{XX}^{\nabla}), \quad \mu^{\nabla}(x) = \begin{bmatrix} \mu(x) \\ \nabla \mu(x) \end{bmatrix}, \quad k^{\nabla}(x, x') = \begin{bmatrix} k(x, x') \\ \nabla_x k(x, x') \end{bmatrix}$$

- ► Results in a larger kernel matrix  $K_{XX}^{\nabla} \in \mathbb{R}^{n(d+1) \times n(d+1)}$ .



### Kernel Learning with Derivatives

$$\mathscr{L}(\mathbf{y}^{\mathbf{v}}|\boldsymbol{\theta}) = \mathscr{L}_{\mathbf{y}^{\nabla}} + \mathscr{L}_{|\mathbf{K}^{\nabla}|} - \frac{\mathbf{H}(\mathbf{u}^{\mathbf{v}}+\mathbf{v})}{2}$$

$$\mathscr{L}_{y\nabla} = -\frac{1}{2} (y^{\nabla} - \mu_X^{\nabla})^T c, \qquad \qquad \frac{\partial \mathscr{L}_{y\nabla}}{\partial \theta_i} = \frac{1}{2} c^T \left( \frac{\partial \widetilde{K}^{\nabla}}{\partial \theta_i} \right)$$
$$\mathscr{L}_{|K\nabla|} = -\frac{1}{2} \log \det \widetilde{K}^{\nabla}, \qquad \qquad \frac{\partial \mathscr{L}_{|K\nabla|}}{\partial \theta_i} = -\frac{1}{2} \operatorname{tr} \left( \widetilde{K}^{\nabla} \right)$$

- ► Naive approach: Compute Cholesky factorization of  $K^{\nabla}$ .
- Challenges:

### Main Ideas

# Scaling Gaussian Process Regression with Derivatives Kun Dong<sup>1</sup> David Eriksson<sup>1</sup> Eric Hans Lee<sup>2</sup> David Bindel<sup>2</sup> Andrew Gordon Wilson<sup>3</sup>



Applied Math<sup>1</sup>, CS<sup>2</sup>, ORIE<sup>3</sup>

$$K_{UU} \begin{bmatrix} W \\ \partial W \end{bmatrix}$$

## **Dimensionality Reduction via Active Subspace Learning**

- Gradients allow us to uncover low-dimensional structure.

$$C = \int$$

- ► Fit GP with gradient information in the active subspace.

## **Recover Implicit Surface with D-SKI**

- Noisy Stanford bunny: 25K points and noisy normals.
- Fit an implicit GP surface:  $f(x_i) = 0$ ,  $\nabla f(x_i) = n_i$



Figure: (Left) Original surface (Middle) Noisy surface (Right) D-SKI reconstruction from noisy surface

# **Bayesian Optimization with Derivatives and Active Subspace Learning**

- while Budget not exhausted do Sample point  $x_{n+1}$ , value  $f_{n+1}$ , and gradient  $\nabla f_{n+1}$ Update data  $\mathscr{D}_{i+1} = \mathscr{D}_i \cup \{x_{n+1}, f_{n+1}, \nabla f_{n+1}\}$



# Discussion

- ► We achieve large-scale Bayesian optimization with derivatives.
- Implementation available at:

https://github.com/ericlee0803/GP\_Derivatives.

# Cornell University

Many high-dimensional problems have low-dimensional structure.  $\triangleright$   $\lambda_i$  is the average change in f given a perturbation along  $q_i$ :

 $\int_{\Omega} \nabla f(x) \nabla f(x)^{\mathsf{T}} dx = Q \wedge Q^{\mathsf{T}}.$ 

 $\blacktriangleright$  Active subspace: Leading  $\tilde{d}$  eigenvectors describe most of the change in f.

Calculate active subspace projection  $P \in \mathbb{R}^{d \times d}$  using sampled gradients Fit GP with gradient information defined by kernel  $k^{\nabla}(P^T x, P^T x')$ Optimize acquisition function,  $u_{n+1} = \arg \max \mathscr{A}(u)$  with  $x_{n+1} = Pu_{n+1}$ 

 $\blacktriangleright$  We test on 5D Ackley embedded in  $[-10, 15]^{50}$  and 5D Rastrigin in  $[-4, 5]^{50}$ . ► We apply D-SKI in a random 2D subspace of the estimated active subspace.

Gradient information is valuable for GP regression, but scalability is a problem. Our approach: Fast MVMs, iterative methods, and dimensionality reduction.