

Scalable Logit Gaussian Process Classification

Florian Wenzel^{1,3}, Théo Galy-Fajou², Christian Donner², Marius Kloft³ and Manfred Opper²

¹Humboldt-Universität zu Berlin, ²TU Berlin, ³TU Kaiserslautern

Summary

- We present a Gaussian process classification method building on **Pólya-Gamma data augmentation** and **inducing points**.
- We develop a fast **stochastic variational inference** algorithm based on efficient **natural gradient updates** which are given in closed-form.
- **Speedups of up to two orders of magnitude** while being competitive in terms of prediction performance.

Gaussian Process Classification

- Data: $X = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{d \times n}$ with labels $\mathbf{y} = (y_1, \dots, y_n) \in \{-1, 1\}^n$.

- Logit GP Classification Model:

$$p(y_i | \mathbf{f}, \mathbf{x}_i) = \sigma(y_i f(\mathbf{x}_i)) = (1 + \exp(-y_i f(\mathbf{x}_i)))^{-1}$$

$$\mathbf{f} \sim \text{GP}(0, K)$$

Pólya-Gamma Data Augmentation

Pólya-Gamma Distribution

- $\omega \sim \text{PG}(b, 0)$, $b > 0$ is defined by the moment generating function $\mathbb{E}_{\text{PG}(\omega|b,0)}[\exp(-\omega t)] = (\cosh^b(\sqrt{t/2}))^{-1}$.
- Idea: write **logistic function** in terms of Pólya-Gamma variables $\sigma(z_i) = (1 + \exp(-z_i))^{-1} = \frac{1}{2} \int \exp\left(\frac{z_i}{2} - \frac{z_i^2}{2} \omega_i\right) p(\omega_i) d\omega_i$
- Where $p(\omega_i) = \text{PG}(\omega_i|1, 0)$.

Pólya-Gamma Augmented Model:

$$p(\mathbf{y}, \boldsymbol{\omega}, \mathbf{f}) = p(\mathbf{y} | \mathbf{f}, \boldsymbol{\omega}) p(\mathbf{f}) p(\boldsymbol{\omega}) \propto \exp\left(\frac{1}{2} \mathbf{y}^\top \mathbf{f} - \frac{1}{2} \mathbf{f}^\top \Omega \mathbf{f}\right) p(\mathbf{f}) p(\boldsymbol{\omega})$$

Sparse Gaussian Process (Inducing Points)

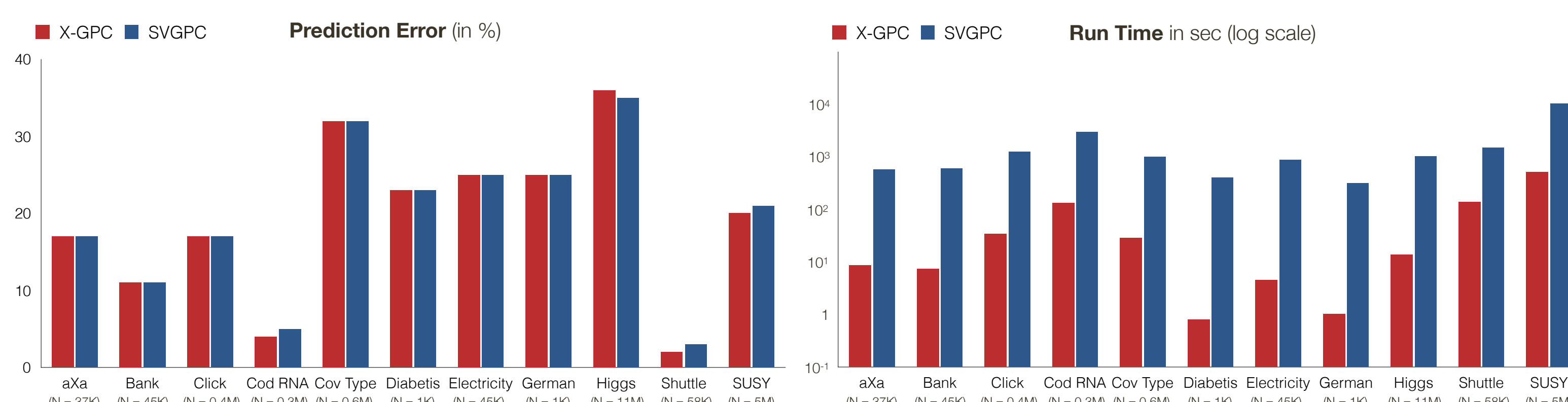
- Inference in GPs is typically $\mathcal{O}(n^3)$.
- Scalable approximation by using a sparse GP representation with m inducing points $(Z_1, \mathbf{u}_1), \dots, (Z_m, \mathbf{u}_m)$ (reduced complexity $\mathcal{O}(m^3)$)
- \mathbf{f} and the inducing variables $\mathbf{u} = (\mathbf{u}_1, \dots, \mathbf{u}_m)$ are connected via $p(\mathbf{f} | \mathbf{u}) = \mathcal{N}(\mathbf{f} | K_{nm} K_{mm}^{-1} \mathbf{u}, \tilde{K})$, $p(\mathbf{u}) = \mathcal{N}(\mathbf{u} | 0, K_{mm})$
- Where $\tilde{K} = K_{nn} - K_{nm} K_{mm}^{-1} K_{mn}$.

Final Augmented Model

$$p(\mathbf{y}, \boldsymbol{\omega}, \mathbf{f}, \mathbf{u}) = p(\mathbf{y} | \boldsymbol{\omega}, \mathbf{f}) p(\boldsymbol{\omega}) p(\mathbf{f} | \mathbf{u}) p(\mathbf{u}).$$

- \mathbf{y} - labels
- \mathbf{f} - latent decision function (modeled as GP)
- \mathbf{u} - inducing points
- $\boldsymbol{\omega}$ - Pólya-Gamma variables

Performance and run time on several datasets



Stochastic Variational Inference

- We propose a fast and scalable inference algorithm.
- The lower bound is given in closed-form which enables efficient optimization.
- Efficient second-order optimization based on natural gradient updates.

Variational Approximation

- Apply VI to marginal joint distribution $p(\mathbf{y}, \boldsymbol{\omega}, \mathbf{u}) = p(\mathbf{y} | \boldsymbol{\omega}, \mathbf{u}) p(\boldsymbol{\omega}) p(\mathbf{u})$.
- Variational distribution: $q(\mathbf{u}, \boldsymbol{\omega}) = q(\mathbf{u}) q(\boldsymbol{\omega})$.
- With $q(\omega_i) = \text{PG}(\omega_i | 1, c_i)$ and $q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}, \Sigma)$.

Variational Lower Bound

$$\log p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u}, \boldsymbol{\omega})}[\log p(\mathbf{y} | \mathbf{u}, \boldsymbol{\omega})] - \text{KL}(q(\mathbf{u}, \boldsymbol{\omega}) || p(\mathbf{u}, \boldsymbol{\omega}))$$

$$\geq \mathbb{E}_{p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) q(\boldsymbol{\omega})}[\log p(\mathbf{y} | \boldsymbol{\omega}, \mathbf{f})] - \text{KL}(q(\mathbf{u}, \boldsymbol{\omega}) || p(\mathbf{u}, \boldsymbol{\omega}))$$

$$=: \mathcal{L}$$

- Is given in closed-form (no sampling needed).

Parameter Updates

(based on mini-batch \mathcal{S} of size s)

- Pólya-Gamma parameters (local):

$$c_i = \sqrt{\tilde{K}_{ii} + \boldsymbol{\kappa}_i \Sigma \boldsymbol{\kappa}_i^\top + \boldsymbol{\mu}^\top \boldsymbol{\kappa}_i \boldsymbol{\kappa}_i^\top \boldsymbol{\mu}}$$

- GP parameters in natural parameterization (global):

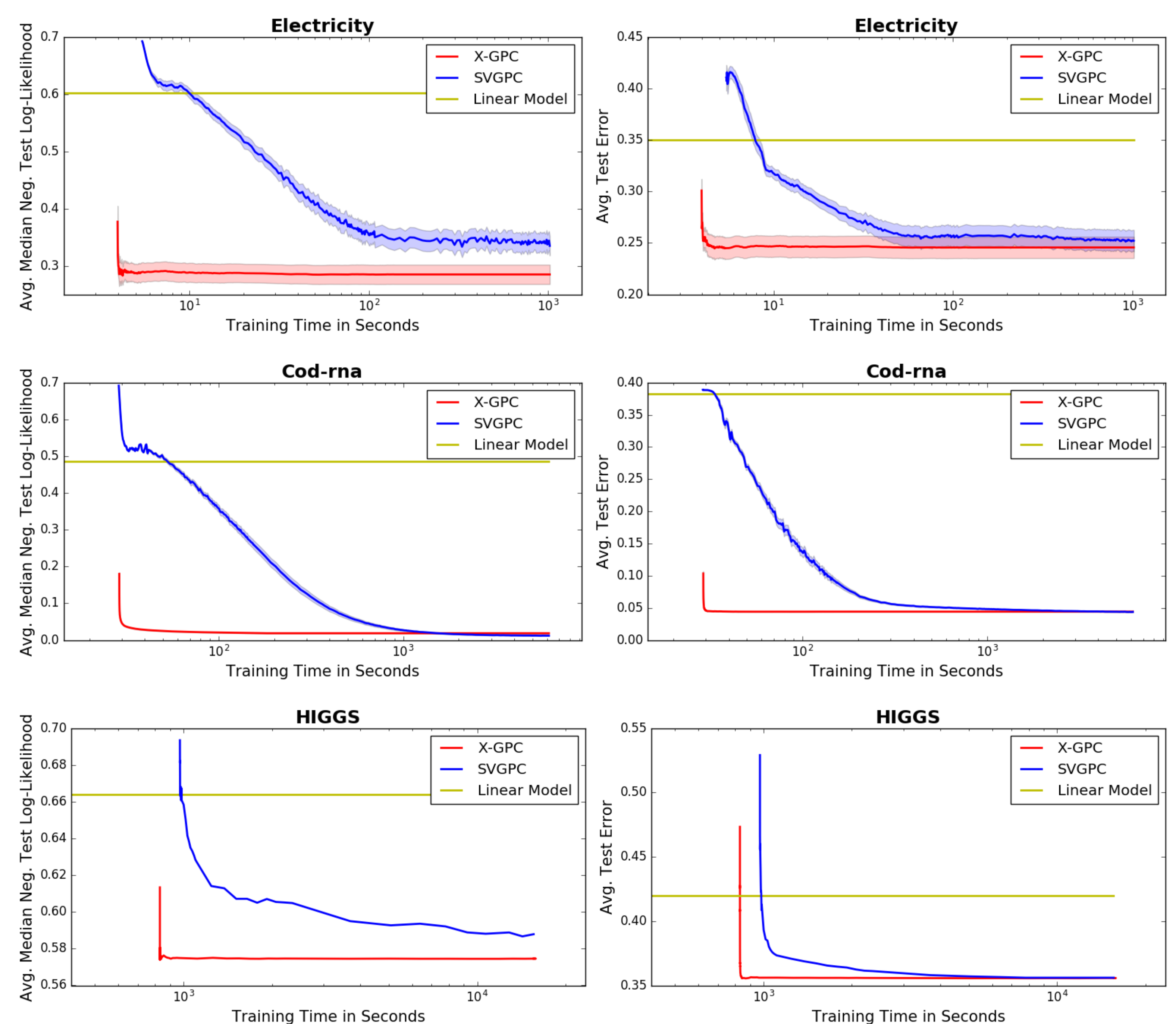
$$\tilde{\nabla}_{\boldsymbol{\eta}_1} \mathcal{L}_{\mathcal{S}} = \frac{n}{2s} \boldsymbol{\kappa}_{\mathcal{S}}^\top \mathbf{y}_{\mathcal{S}} - \boldsymbol{\eta}_1$$

$$\tilde{\nabla}_{\boldsymbol{\eta}_2} \mathcal{L}_{\mathcal{S}} = -\frac{1}{2} \left(K_{mm}^{-1} + \frac{n}{s} \boldsymbol{\kappa}_{\mathcal{S}}^\top \Theta_{\mathcal{S}} \boldsymbol{\kappa}_{\mathcal{S}} \right) - \boldsymbol{\eta}_2$$

- Where $\boldsymbol{\kappa}_i = K_{im} K_{mm}^{-1}$, $\Theta = \text{diag}(\boldsymbol{\theta})$ and $\theta_i = \frac{1}{4c_i} \tanh\left(\frac{c_i}{2}\right)$.

Experiments

Prediction performance as function of time



Contact:

- Mail: wenzelfl@hu-berlin.de
- Web: www.florian-wenzel.de