

Scalable Multi-Class Gaussian Process Classification via Data Augmentation

Théo Galy-Fajou¹, Florian Wenzel^{1,2}, Christian Donner¹
and Manfred Oppel¹

¹TU Berlin, Germany, ²TU Kaiserslautern, Germany

TL;DR

- We propose a **new scalable multi-class Gaussian process classification** approach building on a novel **modified softmax likelihood function**.
- This likelihood allows for a latent variable augmentation that leads to a **conditionally conjugate model** and enables **efficient variational inference** via **block coordinate ascent updates**.
- **Sparse Gaussian Processes** with independent inducing points and kernel parameters.

• • • • •

GP Multi-Class Problem

- Dataset with inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ and labels $\mathbf{y} = (y_1, \dots, y_N)$ where $y_i \in \{1, \dots, C\}$
- Classical softmax : $p(y_i = k | \mathbf{f}_i) = \frac{\exp(f_i^k)}{\sum_{c=1}^C \exp(f_i^c)}$
- With C functions $\mathbf{f} = (f^1, \dots, f^C)$ with $f_i^c = f^c(\mathbf{x}_i)$ independent priors : $f^c \sim \text{GP}(0, k^c)$

Logistic-softmax likelihood

$$p(y_i = k | \mathbf{f}_i) = \frac{\sigma(f_i^k)}{\sum_{c=1}^C \sigma(f_i^c)}$$

- $\sigma(x) = (1 + \exp(-x))^{-1}$ is the logistic function
- Softmax with transformation $h(f) = \log(\sigma(f))$
- Generalization of binary classification with $f^2 = -f^1$

• • • • •

Augmentations to conditional conjugacy

Augm. 1 : Gamma

$$\frac{1}{x} = \int_0^\infty \exp(-\lambda x) d\lambda \implies p(y_i = k | \mathbf{f}_i) = \sigma(f_i^k) \int_0^\infty \exp\left(-\lambda^i \sum_{c=1}^C \sigma(f_i^c)\right) d\lambda^i$$

$$\implies p(y_i = k, \lambda^i | \mathbf{f}_i) = \sigma(f_i^k) \prod_{c=1}^C \exp(-\lambda^i \sigma(f_i^c))$$

Augm. 2 : Poisson

- σ is bounded : $\sigma(x) = 1 - \sigma(-x)$

$$\implies p(y_i = k | \lambda^i, \mathbf{f}_i) = \sigma(f_i^k) \prod_{c=1}^C \exp(-\lambda^i (1 - \sigma(-f_i^c)))$$

$$\exp(x) = \sum_{n=1}^{\infty} x^n / n!$$

$$\implies = \sigma(f_i^k) \prod_{c=1}^C \sum_{n_c^i=1}^{\infty} \exp(-\lambda^i) \left(\lambda^i \sigma(-f_i^c)^{n_c^i} / n_c^i! \right)$$

$$\implies p(y_i = k, \{n_c^i\}_{c=1}^C | \lambda^i, \mathbf{f}_i) = \sigma(f_i^k) \prod_{c=1}^C \text{Po}(n_c^i | \lambda^i) \sigma(-f_i^c)^{n_c^i}$$

Augm. 3 : Pólya-Gamma

$$\sigma^n(x) = 2^{-n} \int_0^\infty \exp\left(\frac{nx}{2} - \frac{\omega x^2}{2}\right) \text{PG}(\omega | n, 0) d\omega, \quad \underbrace{\omega_1}_{\text{PG}(a,s)} + \underbrace{\omega_2}_{\text{PG}(b,s)} = \underbrace{\omega_3}_{\text{PG}(a+b,s)}$$

$$\implies \prod_{c=1}^C \int_0^\infty 2^{-(\delta_{kc} + n_i^c)} \exp\left(\frac{(\delta_{kc} - n_i^c) f_i^c}{2} - \frac{(f_i^c)^2}{2} \omega_i^c\right) \text{PG}(\omega_i^c | \delta_{kc} + n_i^c, 0) d\omega_i^c$$

Conjugate Likelihood

$$p(y_i = k | \{\omega_i^c\}_{c=1}^C, \{n_i^c\}_{c=1}^C, \lambda_i, \mathbf{f}_i)$$

$$\prod_{c=1}^C 2^{-(\delta_{kc} + n_i^c)} \exp\left(\frac{(\delta_{kc} - n_i^c) f_i^c}{2} - \frac{(f_i^c)^2}{2} \omega_i^c\right)$$

• • • • •

Inference

Inducing Points (VFE)

$$p(\mathbf{f}^c) = \int p(\mathbf{f}^c | \mathbf{u}^c) p(\mathbf{u}^c) d\mathbf{u}^c \quad \mathbf{u}^c \sim \text{GP}(0, k^c)$$

Mean Field Approximation

$$p(\mathbf{u}, \boldsymbol{\omega}, \mathbf{n}, \boldsymbol{\lambda} | \mathbf{y}) \approx \prod_{c=1}^C q(\mathbf{u}^c) \prod_{i=1}^N q(\lambda_i) q(n_i^c, \omega_i^c)$$

Block CAVI Updates

$$q^*(\theta) \propto \exp\left(\mathbb{E}_{q/\theta} [\log p(\theta, \boldsymbol{\Theta}_{/\theta}, \mathbf{y})]\right)$$

• • • • •

Experiments

