



UNIVERSIDAD
DE GRANADA



Linearized Laplace Approximation for Modern Deep Learning (II)

Francisco Miguel Castro Macías
fcastro@ugr.es

Visual Information Processing Group (VIP)
Department of Computer Science and Artificial Intelligence (DECSAI)
University of Granada (UGR)

May 5, 2023



Overview

1. Introduction

- 1.1 Our goal
- 1.2 A short recap

2. The Linearized Laplace Approximation as a Gaussian Linear Model

3. Conjugate Gaussian Regression and the EM algorithm

- 3.1 Conjugate Gaussian Regression
- 3.2 EM algorithm

4. Evidence maximization using stochastic approximation

- 4.1 M step
- 4.2 E step

5. Sample-based Lin. Laplace inference

1. Introduction

- 1.1 Our goal
- 1.2 A short recap

2. The Linearized Laplace Approximation as a Gaussian Linear Model

3. Conjugate Gaussian Regression and the EM algorithm

- 3.1 Conjugate Gaussian Regression
- 3.2 EM algorithm

4. Evidence maximization using stochastic approximation

- 4.1 M step
- 4.2 E step

5. Sample-based Lin. Laplace inference



Our goal



Sampling-based inference for large linear models, with application to linearised Laplace

Today

Main references

- **Linearized Laplace Approximation as a Gaussian Linear Model.** Approximate Inference Turns Deep Networks into Gaussian Processes, Khan et al. (2019).
- **Evidence maximization using stochastic approximation.** Sampling-based inference for large linear models, with application to linearised Laplace, Antorán et al. (2022b).



Notation

Let $\mathbf{F} = (F_1, \dots, F_M): \mathbb{R}^N \rightarrow \mathbb{R}^M$ be differentiable and let $d\mathbf{F}_{\mathbf{x}}: \mathbb{R}^N \rightarrow \mathbb{R}^M$ be its differential application at $\mathbf{x} \in \mathbb{R}^N$ given by $d\mathbf{F}_{\mathbf{x}}(\mathbf{v}) = \partial_{\mathbf{x}}\mathbf{F}\mathbf{v}$ for each $\mathbf{v} \in \mathbb{R}^N$. Here, $\partial_{\mathbf{x}}\mathbf{F} \in \mathbb{R}^{M \times N}$ is the Jacobian matrix of \mathbf{F} at \mathbf{x} given by $(\partial_{\mathbf{x}}\mathbf{F})_{ij} = \frac{\partial F_j}{\partial x_i}(\mathbf{x})$. We denote $\nabla\mathbf{F}(\mathbf{x}) = (\partial_{\mathbf{x}}\mathbf{F})^{\top} \in \mathbb{R}^{N \times M}$.

When $M = 1$ and $\mathbf{F} = F$ is twice differentiable, the Hessian matrix of F at \mathbf{x} , denoted by $\nabla^2 F(\mathbf{x}) \in \mathbb{R}^{N \times N}$, is given by $(\nabla^2 F(\mathbf{x}))_{ij} = \frac{\partial^2 F}{\partial x_i \partial x_j}(\mathbf{x})$. If we consider the gradient as $\nabla F: \mathbb{R}^N \rightarrow \mathbb{R}^N$, then $\nabla^2 F(\mathbf{x}) = \partial_{\mathbf{x}}(\nabla F)^{\top}$.



Bayesian Deep Learning

Consider a dataset $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n) : n \in \{1, \dots, N\}\}$ where $\mathbf{x}_n \in \mathbb{R}^D$ and $\mathbf{y}_n \in \mathcal{Y}^C$ (\mathcal{Y} can be \mathbb{R} or $\{0, 1\}$). Let $\mathbf{f} = [f_1, \dots, f_C] : \mathbb{R}^D \times \mathbb{R}^P \rightarrow \mathbb{R}^C$ be a neural network that is trained to minimize the following loss function with respect to $\boldsymbol{\theta} \in \mathbb{R}^P$,

$$\mathcal{L}(\boldsymbol{\theta}, \mathbf{A}, \mathbf{f}) = \sum_n \ell(\mathbf{y}_n, \mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta})) + R(\boldsymbol{\theta}, \mathbf{A}), \quad \tilde{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{A}, \mathbf{f}),$$

where ℓ is a negative log-likelihood and \mathbf{A} denotes the regularization hyperparameters. Let us cast this under the Bayesian framework. Consider

$$p(\boldsymbol{\theta}, \mathcal{D}, \mathbf{A}; \mathbf{f}) = p(\mathbf{A})p(\boldsymbol{\theta} \mid \mathbf{A}) \prod_n p(\mathbf{y}_n \mid \mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta})) \prod_n p(\mathbf{x}_n),$$

$$p(\boldsymbol{\theta} \mid \mathbf{A}) = \exp(-R(\boldsymbol{\theta}, \mathbf{A}))/Z_R(\mathbf{A}), \quad p(\mathbf{y}_n \mid \mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta})) = \exp(-\ell(\mathbf{y}_n, \mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta}))/Z_\ell.$$

Then $p(\boldsymbol{\theta} \mid \mathcal{D}, \mathbf{A}; \mathbf{f}) \propto \exp(-\mathcal{L}(\boldsymbol{\theta}, \mathbf{A}, \mathbf{f}))$. Therefore, minimizing $\mathcal{L}(\boldsymbol{\theta}, \mathbf{A}, \mathbf{f})$ is equivalent to maximizing $\log p(\boldsymbol{\theta} \mid \mathcal{D}, \mathbf{A}; \mathbf{f})$.



Posterior predictive: Laplace Approximation

We make (probabilistic) predictions for a new input \mathbf{x}^* using the posterior predictive

$$p(\mathbf{y}^* \mid \mathbf{x}^*, \mathcal{D}; \mathbf{f}) = \int_{\mathbb{R}^P} p(\mathbf{y}^* \mid \mathbf{f}(\mathbf{x}^*, \boldsymbol{\theta})) p(\boldsymbol{\theta} \mid \mathcal{D}; \mathbf{f}) d\boldsymbol{\theta} = \mathbb{E}_{p(\boldsymbol{\theta} \mid \mathcal{D}; \mathbf{f})} [p(\mathbf{y}^* \mid \mathbf{f}(\mathbf{x}^*, \boldsymbol{\theta}))].$$

Laplace Approximation (LA).

$$p(\boldsymbol{\theta} \mid \mathcal{D}; \mathbf{f}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \tilde{\boldsymbol{\theta}}, \boldsymbol{\Sigma}),$$
$$\boldsymbol{\Sigma}^{-1} = \nabla_{\boldsymbol{\theta}}^2 R(\tilde{\boldsymbol{\theta}}) + \sum_n \nabla_{\boldsymbol{\theta}}^2 \ell(\mathbf{y}_n, \mathbf{f}(\mathbf{x}_n, \tilde{\boldsymbol{\theta}})).$$



The Linearized Laplace Approximation

We consider the *local linearization* of \mathbf{f} around $\tilde{\boldsymbol{\theta}}$,

$$\mathbf{f}^{\text{lin}}(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{f}(\mathbf{x}, \tilde{\boldsymbol{\theta}}) + \partial_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}, \tilde{\boldsymbol{\theta}})(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})$$

Linearized Laplace Approximation (LLA)

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathcal{D}; \mathbf{f}) &\approx p(\boldsymbol{\theta} \mid \mathcal{D}; \mathbf{f}^{\text{lin}}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \tilde{\boldsymbol{\theta}}, \Sigma_{\text{GGN}}), \\ \Sigma_{\text{GGN}}^{-1} &= \nabla_{\boldsymbol{\theta}}^2 R(\tilde{\boldsymbol{\theta}}) + \sum_n \partial_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_n, \tilde{\boldsymbol{\theta}})^\top \Lambda_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}_n, \mathbf{y}_n) \partial_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}_n, \tilde{\boldsymbol{\theta}}), \\ \Lambda_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{y}_n) &= \nabla_{\mathbf{f}}^2 \ell(\mathbf{y}_n, \mathbf{f})|_{\mathbf{f}=\mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta})} \end{aligned}$$

Alternative adaptation of LLA, Antorán et al. (2022a)



Objective. Obtain a pair (θ^*, \mathbf{A}^*) that satisfies

$$\theta^* \in \arg \min_{\theta} \mathcal{L}(\theta, \mathbf{A}^*, \mathbf{f}^{\text{lin}}), \quad \mathbf{A}^* \in \arg \max_{\mathbf{A}} \mathcal{M}(\mathbf{A}, \theta^*).$$

1. Start with an initial guess of \mathbf{A}^* and minimize $\mathcal{L}(\theta, \mathbf{A}^*, \mathbf{f})$ to obtain $\tilde{\theta}$.
 2. Minimize $\mathcal{L}(\theta, \mathbf{A}^*, \mathbf{f}^{\text{lin}})$ obtaining θ^* . In most situations, $\theta \mapsto \mathcal{L}(\theta, \mathbf{A}, \mathbf{f}^{\text{lin}})$ will be convex, yielding a minimizer (convergence guaranteed).
 3. Maximize $\mathcal{M}(\mathbf{A}, \theta^*)$ obtaining \mathbf{A}^* . For every θ , $\mathbf{A} \mapsto \mathcal{M}(\mathbf{A}, \theta)$ is concave, yielding a maximizer (convergence guaranteed).
 4. If (θ^*, \mathbf{A}^*) is not a stationary point, move to step 2.
- Finally, use the pair (θ^*, \mathbf{A}^*) to compute the posterior approximation.

1. Introduction

- 1.1 Our goal
- 1.2 A short recap

2. The Linearized Laplace Approximation as a Gaussian Linear Model

3. Conjugate Gaussian Regression and the EM algorithm

- 3.1 Conjugate Gaussian Regression
- 3.2 EM algorithm

4. Evidence maximization using stochastic approximation

- 4.1 M step
- 4.2 E step

5. Sample-based Lin. Laplace inference





The LLA as a GLM

We seek a Gaussian Linear Model (GLM) whose posterior distribution is equal to the LLA approximation. We define the following *pseudo-inputs*,

$$\tilde{\mathbf{y}}_n = \phi(\mathbf{x}_n)\tilde{\boldsymbol{\theta}} - \boldsymbol{\Lambda}_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{y}_n)^{-1} \mathbf{r}_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}_n, \mathbf{y}_n),$$

where $\mathbf{r}_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{y}_n) = \nabla_{\mathbf{f}} \ell(\mathbf{y}, \mathbf{f})|_{\mathbf{f}=\mathbf{f}(\mathbf{x}_n, \boldsymbol{\theta})}$ and $\phi(\mathbf{x}) = \partial_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}, \tilde{\boldsymbol{\theta}})$. We consider the *transformed dataset*

$$\tilde{\mathcal{D}} = \{(\mathbf{x}_1, \tilde{\mathbf{y}}_1), \dots, (\mathbf{x}_N, \tilde{\mathbf{y}}_N)\}.$$

Next, we consider the linear model

$$\tilde{\mathbf{y}}_n = \phi(\mathbf{x}_n)\boldsymbol{\theta} + \boldsymbol{\eta}_n, \quad \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \quad \boldsymbol{\eta}_n \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}_n, \mathbf{y}_n)^{-1}) \quad (1)$$

The LLA as a GLM (Khan et al., 2019, Theorem 1).

The Linearized Laplace Approximation $\mathcal{N}(\boldsymbol{\theta} \mid \tilde{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{\text{GGN}})$ is equal to the posterior distribution $p(\boldsymbol{\theta} \mid \tilde{\mathcal{D}})$ of the GLM (1).

1. Introduction

- 1.1 Our goal
- 1.2 A short recap

2. The Linearized Laplace Approximation as a Gaussian Linear Model

3. Conjugate Gaussian Regression and the EM algorithm

- 3.1 Conjugate Gaussian Regression
- 3.2 EM algorithm

4. Evidence maximization using stochastic approximation

- 4.1 M step
- 4.2 E step

5. Sample-based Lin. Laplace inference



Conjugate Gaussian Regression



We study Bayesian Conjugate Gaussian Regression (CGR) with multidimensional outputs,

$$\begin{aligned}\mathcal{D} &= \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\} \subset \mathbb{R}^D \times \mathbb{R}^C, \\ \phi: \mathbb{R}^D &\rightarrow \mathbb{R}^{C \times P}, \quad \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \quad \boldsymbol{\eta}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{B}_n^{-1}), \\ \mathbf{y}_n &= \phi(\mathbf{x}_n)\boldsymbol{\theta} + \boldsymbol{\eta}_n.\end{aligned}$$

This can be written

$$\begin{aligned}\mathbf{Y} &= \Phi\boldsymbol{\theta} + \boldsymbol{\eta}, \quad \boldsymbol{\theta} \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1}), \quad \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{B}^{-1}), \\ \mathbf{Y} &= [\mathbf{y}_1^\top; \dots; \mathbf{y}_N^\top]^\top \in \mathbb{R}^{NC}, \quad \Phi = [\phi(\mathbf{x}_1)^\top; \dots; \phi(\mathbf{x}_N)^\top]^\top \in \mathbb{R}^{NC \times P}, \\ \mathbf{B} &= \begin{bmatrix} \mathbf{B}_1 & & \\ & \ddots & \\ & & \mathbf{B}_N \end{bmatrix} \in \mathbb{R}^{NC \times NC}\end{aligned}$$



EM algorithm

Goal. Infer the posterior distribution for the parameters θ given our observations, under the setting of $\mathbf{A} = \alpha \mathbf{I}$ most likely to have generated the observed data.

Iterative procedure.

(E step) Given \mathbf{A} , the posterior for θ is computed exactly as

$$\Pi = \mathcal{N}(\bar{\theta}, \mathbf{H}^{-1}), \quad \mathbf{H} = \Phi^{\top} \mathbf{B} \Phi + \mathbf{A}, \quad \bar{\theta} = \mathbf{H}^{-1} \Phi^{\top} \mathbf{B} \mathbf{Y}.$$

(M step) Given Π , the evidence for the model is optimized using

$$\log p(\mathbf{Y}) \geq \underbrace{-\frac{1}{2} \left[\bar{\theta}^{\top} \mathbf{A} \bar{\theta} + \log \det \left(\mathbf{I} + \mathbf{A}^{-1} \Phi^{\top} \mathbf{B} \Phi \right) \right]}_{\mathcal{M}(\mathbf{A})} + \text{const.}$$

Problem. The E step requires inverting a $P \times P$ matrix. The M step requires evaluating its log-determinant. Both are cubic operations in P . The dual form has cubic complexity in NC .

1. Introduction

- 1.1 Our goal
- 1.2 A short recap

2. The Linearized Laplace Approximation as a Gaussian Linear Model

3. Conjugate Gaussian Regression and the EM algorithm

- 3.1 Conjugate Gaussian Regression
- 3.2 EM algorithm

4. Evidence maximization using stochastic approximation

- 4.1 M step
- 4.2 E step

5. Sample-based Lin. Laplace inference





M step

(M step) Given Π , the evidence for the model is optimized using

$$\log p(\mathbf{Y}) \geq \underbrace{-\frac{1}{2} \left[\bar{\boldsymbol{\theta}}^\top \mathbf{A} \bar{\boldsymbol{\theta}} + \log \det \left(\mathbf{I} + \mathbf{A}^{-1} \boldsymbol{\Phi}^\top \mathbf{B} \boldsymbol{\Phi} \right) \right]}_{\mathcal{M}(\mathbf{A})} + \text{const.}$$

The function $\mathbf{A} \mapsto \mathcal{M}(\mathbf{A})$ is concave. Its maximizer is given by $\nabla_{\mathbf{A}} \mathcal{M}(\mathbf{A}) = \mathbf{0}$, which leads to

$$\bar{\boldsymbol{\theta}}^\top \bar{\boldsymbol{\theta}} \mathbf{A} = \underbrace{\boldsymbol{\Phi}^\top \mathbf{B} \boldsymbol{\Phi} \left(\mathbf{A} + \boldsymbol{\Phi}^\top \mathbf{B} \boldsymbol{\Phi} \right)^{-1}}_{\text{Tr}(\cdot)} \underbrace{\Rightarrow \bar{\boldsymbol{\theta}}^\top \mathbf{A} \bar{\boldsymbol{\theta}} = \text{Tr} \left(\underbrace{\mathbf{H}^{-1} \boldsymbol{\Phi}^\top \mathbf{B} \boldsymbol{\Phi}}_{\gamma} \right)}_{\gamma}.$$

The quantity γ is the *effective dimension of the regression problem*. Setting $\mathbf{A} = \alpha \mathbf{I}$ leads to $\alpha = \gamma / \|\bar{\boldsymbol{\theta}}\|^2$.



M step: computing γ

Using Hutchinson (1989)'s trick,

$$\gamma = \text{Tr}(\mathbf{H}^{-1} \mathbf{\Phi}^\top \mathbf{B} \mathbf{\Phi}) = \mathbb{E}_{\mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})} [\boldsymbol{\xi}^\top \mathbf{\Phi}^\top \mathbf{B} \mathbf{\Phi} \boldsymbol{\xi}] \approx \frac{1}{K} \sum_{k=1}^K \hat{\boldsymbol{\xi}}_k^\top \mathbf{\Phi}^\top \mathbf{B} \mathbf{\Phi} \hat{\boldsymbol{\xi}}_k, \quad \hat{\boldsymbol{\xi}}_k \sim \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$$

Remarks.

- New hyperparameter: the number of samples, K .
- How to sample efficiently from $\mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$?



E step: sampling from the posterior using SGD

Sample-then-optimize (Papandreou and Yuille, 2010), (Antorán et al., 2022b, Appendix C.1).

Let $\mathbf{E} = [\epsilon_1^\top; \dots; \epsilon_N^\top]^\top \in \mathbb{R}^{NC}$, with $\epsilon_n \sim \mathcal{N}(\mathbf{0}, \mathbf{B}_n^{-1})$, and $\theta^0 \sim \mathcal{N}(\mathbf{0}, \mathbf{A}^{-1})$. Then, the minimizer of the following loss is a random variable $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$,

$$L(\mathbf{z}) = \underbrace{\frac{1}{2} (\Phi \mathbf{z} - \mathbf{E})^\top \mathbf{B} (\Phi \mathbf{z} - \mathbf{E})}_{T_1} + \underbrace{\frac{1}{2} (\mathbf{z} - \theta^0)^\top \mathbf{A} (\mathbf{z} - \theta^0)}_{T_2}$$

- $\nabla_{\mathbf{z}} T_1 = \Phi^\top \mathbf{B} (\Phi \mathbf{z} - \mathbf{E}) = \sum_n \phi(\mathbf{x}_n)^\top \mathbf{B}_n (\phi(\mathbf{x}_n) \mathbf{z} - \epsilon_n) \rightarrow$ It requires stochastic approximation for large datasets.
- $\nabla_{\mathbf{z}} T_2 = \mathbf{A} (\mathbf{z} - \theta^0) \rightarrow$ It does not require stochastic approximation.



E step: sampling from the posterior using SGD

An alternative loss is proposed,

$$L'(\mathbf{z}) = \underbrace{\frac{1}{2} \mathbf{z}^\top \Phi^\top \mathbf{B} \Phi \mathbf{z}}_{T'_1} + \underbrace{\frac{1}{2} (\mathbf{z} - \boldsymbol{\theta}^n)^\top \mathbf{A} (\mathbf{z} - \boldsymbol{\theta}^n)}_{T'_2}, \quad \boldsymbol{\theta}^n = \boldsymbol{\theta}^0 + \mathbf{A}^{-1} \Phi^\top \mathbf{B} \mathbf{E}$$

- $\nabla_{\mathbf{z}} T'_1 = \Phi^\top \mathbf{B} \Phi \mathbf{z} = \sum_n \phi(\mathbf{x}_n)^\top \mathbf{B}_n \phi(\mathbf{x}_n) \mathbf{z} \rightarrow$ It requires stochastic approximation for large datasets.
- $\nabla_{\mathbf{z}} T'_2 = \mathbf{A} (\mathbf{z} - \boldsymbol{\theta}^n) \rightarrow$ It does not require stochastic approximation.



E step: sampling from the posterior using SGD

Using L' instead of L may lead to a lower minibatch gradient variance^{1,2}:

$$\begin{aligned}\widehat{\mathbf{g}}_n &= \phi(\mathbf{x}_n)^\top \mathbf{B}_n (\phi(\mathbf{x}_n)\mathbf{z} - \epsilon_n), \quad \widehat{\mathbf{g}}'_n = \phi(\mathbf{x}_n)^\top \mathbf{B}_n \phi(\mathbf{x}_n)\mathbf{z}, \\ \text{Var}(\widehat{\mathbf{g}}_n) - \text{Var}(\widehat{\mathbf{g}}'_n) &= \underbrace{\text{Var}\left(\Phi^\top \mathbf{B} \mathbf{E}\right) - 2 \text{Cov}\left(\Phi^\top \mathbf{B} \Phi \mathbf{z}, \Phi^\top \mathbf{B} \mathbf{E}\right)}_{\Delta}\end{aligned}$$

Criterion: L' is preferred over L if $\text{Tr } \Delta > 0$,

- At initialization, when \mathbf{z} is independent of \mathbf{E} , $\text{Tr } \Delta = \text{Tr}\left(\Phi^\top \mathbf{B} \Phi\right) > 0$.
- At convergence, when $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1})$, $\text{Tr } \Delta > 0$ if $2\alpha\gamma > \text{Tr}\left(\Phi^\top \mathbf{B} \Phi\right)$.

¹ $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}]) (\mathbf{y} - \mathbb{E}[\mathbf{y}])^\top]$, $\text{Var}(\mathbf{x}) = \text{Cov}(\mathbf{x}, \mathbf{x})$.

²Note that \mathbf{z} is a random variable.

1. Introduction

- 1.1 Our goal
- 1.2 A short recap

2. The Linearized Laplace Approximation as a Gaussian Linear Model

3. Conjugate Gaussian Regression and the EM algorithm

- 3.1 Conjugate Gaussian Regression
- 3.2 EM algorithm

4. Evidence maximization using stochastic approximation

- 4.1 M step
- 4.2 E step

5. Sample-based Lin. Laplace inference





Sample-based predictions

Recall that³

$$\mathbf{f}^{\text{lin}}(\cdot, \boldsymbol{\theta}) \sim \mathcal{N} \left(\mathbf{f}(\cdot, \tilde{\boldsymbol{\theta}}), \partial_{\boldsymbol{\theta}} \mathbf{f}(\cdot, \tilde{\boldsymbol{\theta}})^{\top} \boldsymbol{\Sigma}_{\text{GGN}} \partial_{\boldsymbol{\theta}} \mathbf{f}(\cdot, \tilde{\boldsymbol{\theta}}) \right).$$

Therefore, for $r: \mathbb{R}^C \rightarrow \mathcal{Y}^C$,

$$\mathbb{E}_{\mathbf{f} \sim \mathbf{f}^{\text{lin}}(\mathbf{x}^*, \boldsymbol{\theta})} [r(\mathbf{f})] = \frac{1}{K} \sum_{k=1}^K r \left(\mathbf{f}(\mathbf{x}^*, \tilde{\boldsymbol{\theta}}) + \partial_{\boldsymbol{\theta}} \mathbf{f}(\mathbf{x}^*, \tilde{\boldsymbol{\theta}}) \boldsymbol{\xi}_k \right), \quad \boldsymbol{\xi}_k \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\text{GGN}}).$$

³ $\phi(\cdot) = \partial_{\boldsymbol{\theta}} \mathbf{f}(\cdot, \tilde{\boldsymbol{\theta}}), \mathbf{H}^{-1} = \boldsymbol{\Sigma}_{\text{GGN}}$



Sample-based Lin. Laplace inference

Algorithm 1: Sampling-based linearised Laplace hyperparameter learning and inference

Inputs: initial $\alpha > 0$; $k, k' \in \mathbb{N}$, number of samples for stochastic EM and prediction, respectively.

Compute g-prior scaling vector s as in (17)

Sample random regularisers $\theta_1^n, \dots, \theta_k^n$ per (7)

while α has not converged **do**

 Find posterior mode $\bar{\theta}$ by optimising linear model loss $\mathcal{L}(h(\theta, \cdot))$, given in (13)

 Draw posterior samples $\zeta_1 \dots \zeta_k$ by optimising objective L' with $\theta_1^n, \dots, \theta_k^n$

 Estimate effective dimension $\hat{\gamma}$, per (5), using samples $\zeta_1 \dots \zeta_k$

 Update prior precision $\alpha \leftarrow \hat{\gamma} / \|\bar{\theta}\|_2^2$

Sample k' random regularisers $\theta_1^{n'}, \dots, \theta_{k'}^{n'}$ using optimised α

Draw corresponding posterior samples $\zeta'_1, \dots, \zeta'_{k'}$ using loss L'

Output: posterior samples $\zeta'_1, \dots, \zeta'_k$

Thank you!



References I



- Javier Antorán, David Janz, James U Allingham, Erik Daxberger, Riccardo Rb Barbano, Eric Nalisnick, and José Miguel Hernández-Lobato. Adapting the linearised laplace model evidence for modern deep learning. In *International Conference on Machine Learning*, pages 796–821. PMLR, 2022a.
- Javier Antorán, Shreyas Padhy, Riccardo Barbano, Eric Nalisnick, David Janz, and José Miguel Hernández-Lobato. Sampling-based inference for large linear models, with application to linearised laplace. *arXiv preprint arXiv:2210.04994*, 2022b.
- Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 18(3):1059–1076, 1989.

References II



- Mohammad Emtiyaz E Khan, Alexander Immer, Ehsan Abedi, and Maciej Korzepa. Approximate inference turns deep networks into gaussian processes. *Advances in neural information processing systems*, 32, 2019.
- George Papandreou and Alan L Yuille. Gaussian sampling by local perturbations. *Advances in Neural Information Processing Systems*, 23, 2010.