

Laplace Approximations for Deep Learning

A Bayesian Odyssey in Uncertainty: from Theoretical Foundations to Real-World Applications

Alexander Immer





Deep Learning

Data $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$, weights $\theta \in \mathbb{R}^P$, and model hyperparameters \mathcal{M}



 \rightarrow Successful at learning complex patterns from large data sets



Probabilistic View





Deep Learning Workflow





Part I: Bayesian Model Selection with Laplace



Part II: Predictive Uncertainties with Laplace





The Bayesian Approach



Advantage 1) Bayesian Model Selection

Optimize marginal likelihood: $p(\mathcal{D}|\mathcal{M}) = \int p(\mathcal{D}|\theta, \mathcal{M}) p(\theta|\mathcal{M}) d\theta$



\rightarrow Occam's razor: balance data-fit and model complexity 1,2

¹Blumer, Ehrenfeucht, Haussler, Warmuth. *Occam's Razor*. Information processing letters 1987. ²Rasmussen, Ghahramani: *Occam's Razor*. NIPS 2001.



Advantage 2) Predictive Uncertainty

Posterior predictive: $p(y_*|x_*, D) = \int p(y_*|x_*, \theta) p(\theta|D) d\theta$ Posterior



→ Additionally provide epistemic uncertainty

Aleatoric vs. Epistemic Uncertainty

Potential definition of uncertainties as variance decomposition:

$$\mathbb{V}_{p(\mathbf{y}|\mathbf{x}_{*},\mathcal{D})}[\mathbf{y}] = \int \mathbb{V}_{p(\mathbf{y}|\mathbf{x}_{*},\theta)}[\mathbf{y}]p(\theta|\mathcal{D}) \,\mathrm{d}\theta + \int \left(\mathbb{E}_{p(\mathbf{y}|\mathbf{x}_{*},\theta)}[\mathbf{y}] - \mathbb{E}_{p(\mathbf{y}|\mathbf{x}_{*},\mathcal{D})}[\mathbf{y}]\right)^{2} p(\theta|\mathcal{D}) \,\mathrm{d}\theta$$
aleatoric

Variance in labels for same input \rightarrow aleatoric uncertainty

No similar input to $x_* \rightarrow$ epistemic uncertainty





Laplace Approximations



Laplace Approximation



- 1) MargLik: $\log q(\mathcal{D}|\mathcal{M}) = \log p(\mathcal{D}, \hat{\theta}|\mathcal{M}) \frac{1}{2} \log |\frac{1}{2\pi} \mathbf{H}_{\hat{\theta}}(\mathcal{M})|$
- 2) Posterior: $q(\theta) = \mathcal{N}(\theta; \hat{\theta}, \mathbf{H}_{\hat{\theta}}^{-1})$
- 3) Predictive: $q(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) = \int p(\mathbf{y}_*|\mathbf{x}_*, \theta) \mathcal{N}(\theta; \hat{\theta}, \mathbf{H}_{\hat{\theta}}^{-1}) d\theta$

But: Hessian in $O(P^2)$ not tractable for Deep Learning

¹Laplace. *Mémoire sur les probabilités*. Mémoires de l'Académie royale des sciences de Paris 1778. ²MacKay. *A practical Bayesian framework for backpropagation networks*. Neural Computation 1992.





\rightarrow Tractable for extremely large models

¹Ritter, Botev, Barber. A Scalable Laplace Approximation for Neural Networks. ICLR 2018. ²Martens, Grosse. Optimizing Neural Networks with Kronecker-factored Approximate Curvature. ICML 2015.

1 Generalized Gauss-Newton (GGN) instead of Hessian:



→ Positive-semidefinite and reduced cost

2 Decompositions and sampling for Jacobian-vector products:

$$\begin{aligned} \mathbf{J}_{\theta}(\mathbf{x})^{\top} \nabla_{\mathbf{f}}^{2} \log p(\mathbf{y}|\mathbf{f}) \mathbf{J}_{\theta}(\mathbf{x}) \\ &= \sum_{k} \mathbf{v}_{k} \mathbf{v}_{k}^{T} & \mathsf{GGN} & \mathcal{O}(\mathcal{PK}) \\ &= \mathbb{E}_{p(\mathbf{y}|\mathbf{f})} [\nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})^{2}] & \mathsf{Fisher} & \mathcal{O}(\mathcal{PS}) \ \mathcal{S} \ge 1 \\ &\approx \nabla_{\mathbf{f}} \log p(\mathbf{y}|\mathbf{f})^{2} & \mathsf{Empirical Fisher} & \mathcal{O}(\mathcal{P}) \end{aligned}$$

 \rightarrow Reduce cost from O(PK) to O(P) and not require Jacobian

3 KFAC¹ is a layer-wise block-diagonal structured approximation:

$$\begin{bmatrix} \sum_{n} \mathbf{J}_{\theta}(\mathbf{x}_{n})^{\top} \mathbf{\Lambda}(\mathbf{f}_{n}) \mathbf{J}_{\theta}(\mathbf{x}_{n}) \end{bmatrix}_{l} = \sum_{n} \begin{bmatrix} \mathbf{a}_{l,n} \otimes \mathbf{g}_{l,n} \end{bmatrix} \mathbf{\Lambda}(\mathbf{f}_{n}) \begin{bmatrix} \mathbf{a}_{l,n} \otimes \mathbf{g}_{l,n} \end{bmatrix}^{\top}$$

Factors vary with layer type²
$$\sum_{n} \begin{bmatrix} \mathbf{a}_{l,n} \mathbf{a}_{l,n}^{\top} \end{bmatrix} \otimes \begin{bmatrix} \mathbf{g}_{l,n} \mathbf{\Lambda}(\mathbf{f}_{n}) \mathbf{g}_{l,n}^{\top} \end{bmatrix} \approx \frac{1}{N} \begin{bmatrix} \sum_{n} \mathbf{a}_{l,n} \mathbf{a}_{l,n}^{\top} \end{bmatrix} \otimes \begin{bmatrix} \sum_{n} \mathbf{g}_{l,n} \mathbf{\Lambda}(\mathbf{f}_{n}) \mathbf{g}_{l,n}^{\top} \end{bmatrix}$$

 $\mathcal{O}(H^2)$

$$\rightarrow$$
 Reduce storage $O(P^2)$ to $O(P)$ despite having off-diagonal entries!

O(H⁴)





Part I: Bayesian Model Selection for Deep Learning

Structured Laplace Approximations

Structured Hessian approximations yield lower bounds to the Laplace marglik¹:





 \rightarrow Theoretical justification for approximations



Practical Bayesian Model Selection with Laplace

Online Laplace approximation every F steps during training¹:

1) Gradient-based optimization of differentiable ${\cal M}$

 $\mathcal{M} \leftarrow \mathcal{M} + \gamma \nabla_{\mathcal{M}} \log q(\mathcal{D}|\mathcal{M})$

Gradient-based Selection

 \rightarrow Can optimize high-dimensional hyperparameters without re-training

2) Comparison of models and checkpoints

 $\log q(\mathcal{D}|\mathcal{M})$ Discrete Selection

¹Immer, Bauer, Fortuin, Rätsch, Khan. Scalable Bayesian Model Selection for Deep Learning. ICML 2021.



Example: Gradient-Based Model Selection

Optimize prior precision (regularization) with parameters in one training run



 \rightarrow Validation-based selection requires many iterations

CIFAR-10 Classification



 \mathcal{M} = regularization per layer

Compare train performance: Standard vs MargLik

Compare test performance:

 \rightarrow MargLik generalizes better

 \rightarrow Optimize thousands of regularization parameters in only 2x runtime



Marginal Likelihood for Architecture Comparison





 \rightarrow Discrete selection possible after optimizing the prior

Examples of Bayesian Model Selection



Immer, Bauer, Fortuin, Rätsch, Khan. Scalable Bayesian Model Selection for Deep Learning. ICML 2021. Antoran, Barbano, ..., Jin. Uncertainty Estimation for Computed Tomography with a Linearised Deep Image Prior. TMLR 2023. Dhahri, Immer, Charpentier, Günnemann, Fortuin. Shaving Weights with Occam's Razor. Preprint 2024.

• Selecting architecture, structure, variables

Zhou*, Yang*, Wang, Pan. *BayesNAS: A Bayesian Approach for Neural Architecture Search.* CVPR 2019. van der Ouderaa, Immer, van der Wilk. *Learning Layer-wise Equivariances Automatically using Gradients.* NeurIPS 2023. Bouchiat, Immer, Yeche, Rätsch, Fortuin. *Improving Neural Additive Models with Bayesian Principles.* ICML 2024.

Learning invariances/data augmentation

van der Wilk, Bauer, John, Hensman. *Learning Invariances using the Marginal Likelihood*. NeurIPS 2018. Immer*, van der Ouderaa*, Fortuin, Rätsch, van der Wilk. *Invariance Learning in Deep Neural Networks […]*. NeurIPS 2022.





Limitations of Bayesian Model Selection

• Required probabilistic model

• Can still require manual tuning or validation

• No benefit on saturated benchmarks





Part 2: Predictive Uncertainties



Laplace Posterior Predictive Underfitting

$q(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) = \int p(\mathbf{y}_*|\mathbf{x}_*, \theta) \mathcal{N}(\theta; \hat{\theta}, \mathbf{H}_{\hat{\theta}}^{-1}) d\theta$



Problem: Laplace posterior predictive underfits¹



Linearized Laplace

Hessian approximations linearize implicitly²

$$f_{\hat{\theta}}^{\text{lin}}(\mathbf{x}, \theta) = f(\mathbf{x}, \hat{\theta}) + \mathbf{J}_{\hat{\theta}}(\mathbf{x})(\theta - \hat{\theta})$$

We should keep this in mind for predictions:

$$q(\mathbf{y}_*|\mathbf{x}_*, \mathcal{D}) \approx \int p(\mathbf{y}_*|\mathbf{f}_{\hat{\theta}}^{\text{lin}}(\mathbf{x}_*, \theta)) \mathcal{N}(\theta; \hat{\theta}, \mathbf{H}_{\hat{\theta}}^{-1}) \, \mathrm{d}\theta$$

\rightarrow Stabilizes and improves Laplace predictive





Linearized Model as a Gaussian Process



\rightarrow Enables different posterior approximation techniques

Immer, Korzepa, Bauer. Improving predictions of Bayesian neural nets via local linearization. AISTATS 2021. Deng, Zhou, Zhu. Accelerated Linearized Laplace Approximation for Bayesian Deep Learning. NeurIPS 2022.

Laplace Posterior Predictive





• Last-layer is very cheap and effective baseline

Kristiadi, Hein, Hennig. *Being Bayesian, even just a bit, fixes overconfidence in relu networks.* ICML 2020. Ober, Rasmussen. *Benchmarking the neural linear model for regression.* AABI 2019.

• Laplace can improve ensembling further

Eschenhagen, Daxberger, Hennig, Kristiadi. Mixtures of Laplace Approximations [...]. BDL@NeurlPS 2021.

Linearized predictive can have a closed-form (approximation)

Immer, Korzepa, Bauer. Improving predictions of Bayesian neural nets via local linearization. AISTATS 2021. Deng, Zhou, Zhu. Accelerated Linearized Laplace Approximation for Bayesian Deep Learning. NeurIPS 2022.

Limitations



- Scalability issues in some cases¹
- Linearization can be expensive²

• Post-hoc cannot fix badly trained network

laplace-torch¹ on github.com/AlexImmer/Laplace



양 Fork 70 → ☆ Star 447

Contributors 12







Thank you!

Co-Authors

Matthias Bauer Kouroche Bouchiat Peter Bühlmann **Bertrand Charpentier** Ryan Cotterell Erik Daxberger Rayen Dhahri Runa Eschenhagen Vincent Fortuin Stephan Günnemann Philipp Hennig Lucas Torroba Hennigen Francis Jacob Mohammad Emtiyaz Khan

Maciej Korzepa Agustinus Kristiadi Kjong-Van Lehmann Alexander Marx Emanuele Palumbo Gunnar Rätsch Frank Schneider Bernhard Schölkopf **Christoph Schultheiss** Stefan Stark Richard Turner Tycho van der Ouderaa Mark van der Wilk Julia Vogt Hugo Yèche



Estimating Heteroscedastic Aleatoric Uncertainty



Model mean $\mu(\mathbf{x})$ and variance $\sigma^2(\mathbf{x})$ of a Gaussian likelihood:

$$\operatorname{og} p(\mathbf{y}|\mathbf{x}, \mu, \sigma^2) \propto -\frac{1}{2} \log \sigma^2(\mathbf{x}) - \frac{(\mathbf{y} - \mu(\mathbf{x}))^2}{2\sigma^2(\mathbf{x})}$$

Objective balances mean and variance \rightarrow Hard to optimize and regularize^{2,3}

¹Wessels, Mendez-Mancilla, ..., Sanjana. *Massively parallel Cas13 screens reveal principles for guide RNA design*. Nature Biotechnology 2020. ²Stirn, Wessels, ..., Knowles. *Faithful Heteroscedastic Regression with Neural Networks*. AISTATS 2023. ³Seitzer, Tavakoli, Antic, Martius. *On the Pitfalls of Heteroscedastic Uncertainty Estimation with Probabilistic Neural Networks*. ICLR 2022.



Problem with Mean-Variance Parameterization

Mean-variance parameterization does not yield concave log likelihood:

$$\frac{\partial^2 \log \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}, \sigma^2)}{\partial \sigma^2 \partial \sigma^2} = \frac{\sigma^2 - 2(\boldsymbol{\mu} - \mathbf{y})^2}{2\sigma^6} \qquad \qquad \text{Can be positive or negative.} \\ \rightarrow \text{ potentially indefinite Hessian} \\ \rightarrow \text{ not concave} \end{cases}$$

 \rightarrow Cannot apply Hessian or Laplace approximations naively



Revisiting the Natural Parameterization

Natural parameterization always yields a concave log likelihood^{1,2} with

$$\eta_1 = \frac{\mu}{\sigma^2}$$
 and $\eta_2 = -\frac{1}{2\sigma^2}$

 \rightarrow Hessian and Laplace approximations straightforward

 \rightarrow Neural network models both natural parameters

¹Le, Smola, Canu. Heteroscedastic Gaussian process regression. ICML 2005. ²Martens. New insights and perspectives on the natural gradient method. JMLR 2020.



Heteroscedastic Regression Illustration





Epistemic Uncertainty in Heteroscedastic Regression

Problem: predict gene knockdown efficacy of guides in CRISPR system²



→ Improvements due to both natural parameterization and Laplace

¹Seitzer, Tavakoli, Antic, Martius. On the Pitfalls of Heteroscedastic Uncertainty Estimation with Probabilistic Neural Networks. ICLR 2022. ²Stirn, Wessels, ..., Knowles. Faithful Heteroscedastic Regression with Neural Networks. AISTATS 2023.