

 **DTU Compute**
Department of Applied Mathematics and Computer Science

Approximate Inference for Bayesian Neural Networks

Silas Brack (s174433)

Kongens Lyngby 2022



DTU Compute

Department of Applied Mathematics and Computer Science

Technical University of Denmark

Matematiktorvet

Building 303B

2800 Kongens Lyngby, Denmark

Phone +45 4525 3031

compute@compute.dtu.dk

www.compute.dtu.dk

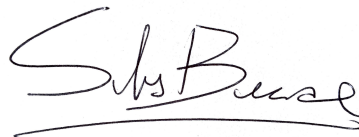
Summary

The first move is optional and contains general background information about your key research variable or variables. The second move is the statement of the problem indicating your research hypothesis/ question. The third move is the methodology move representing your participants, if any ; your materials/instruments, procedures, and data analysis. The fourth move portrays your findings. Finally, in the last optional move, you talk about the likely implications of the study. (Maybe 2x this in terms of size.)

Preface

This Master's thesis was prepared at the department of Applied Mathematics and Computer Science at the Technical University of Denmark in fulfillment of the requirements for acquiring a Master's degree in Mathematical Modelling and Computation.

Kongens Lyngby, August 4, 2022

A handwritten signature in black ink, reading "Silas Brack". The signature is written in a cursive style with a horizontal line underneath the name.

Silas Brack (s174433)

Acknowledgements

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Contents

Summary	i
Preface	iii
Acknowledgements	v
Contents	vii
List of Figures	ix
List of Tables	ix
1 Introduction	1
1.1 The problem	1
1.2 Current state of research in the area	1
1.3 My solution	1
1.4 Research objectives	1
2 Literature Review	3
3 Theory	5
3.1 Bayes' Theorem	5
3.2 Variational Inference	5
3.3 Methods	6
4 Method	11
5 Results	13
6 Experiments	15
7 Conclusion	17
7.1 Conclusion and Outlook	17
A An Appendix	19

Bibliography	21
---------------------	-----------

List of Figures

5.1	In-distribution calibration curves for various variational families with variational inference and Laplace approximation trained on MURA.	13
-----	---	----

List of Tables

5.1	ELBO and \hat{k} -statistic for variational inference	14
-----	---	----

x

CHAPTER 1

Introduction

1.1 The problem

1.1.1 What is the problem?

[Wid06]

1.1.2 Why is it interesting and important? (Motivation)

1.2 Current state of research in the area

1.2.1 Why is it hard? (E.g. why do naïve approaches fail?)

1.2.2 Why hasn't it been solved before? (Or, what's wrong with previous proposed solutions?)

1.3 My solution

1.3.1 How does my solution differ?

The ideas presented in this article essentially constitute a compilation of advice given by academics from multiple fields and institutions,

1.3.2 What are the key components of my approach and results? Are there any specific limitations?

1.4 Research objectives

1.4.1 What are the goals for this project?

This should be discussed in a bit more detail in a thesis, since there are certain objectives discussed in the project plan. In an article, a quick sentence which sum-

marises the rest of the introduction and states succinctly exactly which problem will be tackled and which approach will be made should suffice.

CHAPTER 2

Literature Review

CHAPTER 3

Theory

3.1 Bayes' Theorem

Bayes' theorem describes the probability of an event occurring with respect to prior knowledge of that event. Specifically for modelling, it can be used to relate latent variables with observed variables, allowing distributions on these latent variables to be inferred. We therefore want to find the posterior $p(\boldsymbol{\theta} | \mathbf{y})$ by According to Bayes' theorem, the posterior is given by

$$p(\boldsymbol{\theta} | \mathbf{y}) = \frac{p(\mathbf{y} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})} \quad (3.1)$$

where $p(\mathbf{y} | \boldsymbol{\theta})$ is the likelihood, parameterised by the model parameters $\boldsymbol{\theta}$, $p(\boldsymbol{\theta}) = p(\boldsymbol{\theta}; \boldsymbol{\alpha})$ is the prior, parameterised by its hyper-parameters $\boldsymbol{\alpha}$, and $p(\mathbf{y}) = p_{\boldsymbol{\theta}}(\mathbf{y})$ is the marginal likelihood (or evidence).

Furthermore, the posterior predictive probability, i.e., the probability of observing some new data, is given by

$$p(\mathbf{y}_{\text{new}} | \mathbf{y}) = \int p(\mathbf{y}_{\text{new}} | \boldsymbol{\theta})p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}. \quad (3.2)$$

3.2 Variational Inference

In variational inference (VI), we minimise the KL divergence between the posterior distribution $p(\boldsymbol{\theta} | \mathbf{y})$ and a variational distribution $q_{\boldsymbol{\phi}}(\boldsymbol{\theta})$ parameterised by its variational parameters $\boldsymbol{\phi}$ by optimising with respect to these parameters. The KL divergence is defined as

$$\begin{aligned} D_{\text{KL}} [q(\boldsymbol{\theta}) \parallel p(\mathbf{y} | \boldsymbol{\theta})] &\equiv \mathbb{E}_q \left[\log \frac{q(\boldsymbol{\theta})}{p(\mathbf{y} | \boldsymbol{\theta})} \right] \\ &= \mathbb{E}_q [\log q(\boldsymbol{\theta})] - \mathbb{E}_q [\log p(\mathbf{y} | \boldsymbol{\theta})] \end{aligned} \quad (3.3)$$

However, the KL divergence contains the evidence term $\log p(\mathbf{y})$, which is intractable. Instead, we define a lower bound on this marginal likelihood term that is surrogate to the KL divergence (known as the ELBO). Since they are surrogate, ELBO fulfils

the property $\mathcal{L}[q_\phi] = \log p(\mathbf{y}) - D_{\text{KL}}[q(\boldsymbol{\theta}) \parallel p(\mathbf{y} | \boldsymbol{\theta})]$, minimising the KL divergence is equivalent to maximising the ELBO, which is defined as

$$\begin{aligned} \mathcal{L}[q_\phi] &\equiv \mathbb{E}_q[\log p(\mathbf{y}, \boldsymbol{\theta})] - \mathbb{E}_q[\log q(\boldsymbol{\theta})] \\ &= \mathbb{E}_q[\log p(\mathbf{y} | \boldsymbol{\theta})] + \mathbb{E}_q[\log p(\boldsymbol{\theta})] - \mathbb{E}_q[\log q(\boldsymbol{\theta})] \\ &= \mathbb{E}_q[\log p(\mathbf{y} | \boldsymbol{\theta})] - \underbrace{\mathbb{E}_q[\log p(\boldsymbol{\theta})]}_{\text{Cross-entropy}} + \underbrace{-\mathbb{E}_q[\log q(\boldsymbol{\theta})]}_{\text{Entropy}} \end{aligned} \quad (3.4)$$

where the first term $\mathbb{E}_q[\log p(\mathbf{y} | \boldsymbol{\theta})]$ is the data (or likelihood) term, $-\mathbb{E}_q[\log p(\boldsymbol{\theta})]$ is the cross-entropy of the prior with respect to the variational approximation, and $-\mathbb{E}_q[\log q(\boldsymbol{\theta})]$ is the entropy term of the variational approximation. The last two terms can be interpreted as the regularising KL divergence (or relative entropy) from the prior to the variational approximation, $D_{\text{KL}}[q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta})] = \mathbb{E}_q[\log q(\boldsymbol{\theta})] - \mathbb{E}_q[\log p(\boldsymbol{\theta})]$. Furthermore, the likelihood term in the ELBO can be decomposed further by assuming independence in the observations, calculated as the term given by $\mathbb{E}_q[\log p(\mathbf{y} | \boldsymbol{\theta})] = \frac{1}{N} \sum_i^N \mathbb{E}_q[\log p(y_i | \boldsymbol{\theta})]$. Overall, BNNs have been found to be relatively ineffective unless the number of observations is greater than the number of model parameters.

3.3 Methods

Multiple approximate inference methods were implemented. Specifically, *maximum a posteriori* estimation, Laplace approximation, VI with mean-field, full-rank, low-rank and radial variational families, deep ensembles, and MultiSWAG. These methods will be described in the following sections.

3.3.1 Maximum a Posteriori Estimation

Maximum a posteriori estimation finds a point estimate of the posterior given by its maximum. It can therefore be interpreted as a Delta distribution estimate of the posterior $p(\boldsymbol{\theta} | \mathbf{y}) \approx \delta(\boldsymbol{\theta}_{\text{MAP}})$, where $\boldsymbol{\theta}_{\text{MAP}}$ is defined as

$$\begin{aligned} \boldsymbol{\theta}_{\text{MAP}} &= \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{y}) \\ &= \max_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta} | \mathbf{y}) \\ &= \max_{\boldsymbol{\theta}} [\log p(\mathbf{y} | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})] \\ &= \max_{\boldsymbol{\theta}} \left[\sum_{i=1}^N \log p(y_i | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \right]. \end{aligned} \quad (3.5)$$

As an optimisation problem, maximising the posterior corresponds to minimising the regularised loss,

$$\boldsymbol{\theta}_{\text{MAP}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{y}; \boldsymbol{\theta}) \quad (3.6)$$

$$= \arg \min_{\boldsymbol{\theta}} \left[- \sum_{i=1}^N \log p(y_i | \boldsymbol{\theta}) - \log p(\boldsymbol{\theta}) \right]. \quad (3.7)$$

In this formulation, the term $-\sum_{i=1}^N \log p(y_i | \boldsymbol{\theta})$ is known as the empirical loss or reconstruction loss and the term $\log p(\boldsymbol{\theta})$ is the regulariser.

3.3.2 Laplace Approximation

In the Laplace approximation (LA) [Dax+21], the posterior is approximated by a Gaussian, similarly to mean-field VI. However, instead of finding the optimal Gaussian distribution locations and scales by maximising the ELBO, the LA finds the location by computing the MAP solution $\boldsymbol{\theta}_{\text{MAP}}$ and the scale by approximating the log posterior with a second degree Taylor expansion around this solution ($\boldsymbol{\theta}_0 = \boldsymbol{\theta}_{\text{MAP}}$) and determining the curvature via its Hessian matrix $\mathbf{H} = \nabla_{\boldsymbol{\theta}}^2 \log p(\boldsymbol{\theta} | \mathbf{y})|_{\boldsymbol{\theta}_{\text{MAP}}}$. Since the Taylor expansion is performed around the MAP solution, the first order derivative is zero, and the expansion is simply given by

$$\ln p(\boldsymbol{\theta} | \mathbf{y}) \approx \ln p(\boldsymbol{\theta}_0 | \mathbf{y}) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \quad (3.8)$$

$$\tilde{p}(\boldsymbol{\theta} | \mathbf{y}) \approx p(\boldsymbol{\theta}_0 | \mathbf{y}) \exp \left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \right).$$

Normalising this unnormalised posterior yields the Laplace approximation of the posterior

$$\begin{aligned} p(\boldsymbol{\theta} | \mathbf{y}) &\approx \sqrt{\frac{\det \mathbf{H}}{(2\pi)^D}} \exp \left(-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) \right) \\ &\approx \mathcal{N}(\boldsymbol{\theta}_0, \mathbf{H}^{-1}) = \mathcal{N}(\boldsymbol{\theta}_{\text{MAP}}, \mathbf{H}^{-1}). \end{aligned} \quad (3.9)$$

3.3.3 Mean-Field Variational Approximation

One common variational family used to approximate the real posterior is a product of independent Gaussian distributions (the mean-field approximation) such that each model parameter is sampled from a normal distribution and is independent of all other model parameters, yielding a Gaussian with a diagonal covariance matrix. In mean-field VI, model weights are sampled from an approximate posterior $q(\boldsymbol{\theta}) \approx p(\boldsymbol{\theta} | \mathbf{y})$ as

$$\begin{aligned} q(\boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}) \Leftrightarrow \\ \boldsymbol{\theta} &= \boldsymbol{\mu} + \boldsymbol{\sigma} \circ \boldsymbol{\epsilon} \end{aligned} \quad (3.10)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For this type of variational approximation, the entropy term is given by $\mathbb{H}[q(\theta)] = -\sum_i \log \sigma_i$ and the cross-entropy of the prior relative to the variational approximation is calculated via Monte Carlo simulation by taking the prior log probability with respect to mean-field posterior samples as

$$\begin{aligned} \mathbb{H}[q(\theta), p(\theta)] &= - \int q(\theta) \log(p(\theta)) d\theta \\ &\approx \frac{1}{S} \sum_{s=1}^S \log p(\theta^{(s)}) \end{aligned} \quad (3.11)$$

where $\theta^{(s)} \sim \mathcal{N}(\mu, \sigma)$.

The use of the mean-field approximation for BNNs has been found to be unreliable [Wu+18] and, for increasingly wide neural networks, converges towards the prior [CPD21]. However, for *deep* BNNs, the mean-field assumption may be reasonable [FSG20].

3.3.4 Full-Rank Variational Approximation

As an alternative to the mean-field approximation, a full-rank approximation doesn't assume independence of the model parameters, instead being sampled as

$$q(\theta) = \mathcal{N}(\mu, \Sigma_{\text{FR}}). \quad (3.12)$$

For D model parameters, while a mean-field approximation has scale parameters σ^D , the full-rank approximation has scale parameters $\Sigma_{\text{FR}}^{D \times D}$. The variational parameters scale quadratically with the number of model parameters, which constitutes a significant limitation of this type of approximation. Specifically, since NNs are typically overparameterised, possessing from thousands to millions to billions of parameters, the quadratic number of variational parameters makes this type of VI infeasible for BNNs. As such, full-rank VI is not applied in any of the experiments in this paper.

3.3.5 Low-Rank Variational Approximation

As a compromise between the last two methods, a low-rank variational approximation uses a low-rank approximation of the covariance matrix Σ_{LR} .

$$q(\theta) = \mathcal{N}(\mu, \Sigma_{\text{LR}}). \quad (3.13)$$

This approximation is obtained by reconstructing a covariance matrix parameterised by the covariance factor $\tilde{\Sigma}$ and diagonal element vector σ_{diag} as

$$\begin{aligned} \hat{\Sigma}_{\text{LR}} &= \tilde{\Sigma} \tilde{\Sigma}^T + \sigma_{\text{diag}} \\ \Sigma_{\text{LR}} &= \min_{\tilde{\Sigma}_{\text{LR}}} \left\| \Sigma - \hat{\Sigma}_{\text{LR}} \right\|_F \end{aligned} \quad (3.14)$$

3.3.6 Radial Variational Approximation

Farquhar, Osborne, and Gal [FOG20] introduces the Radial approximate posterior, For each NN layer, this posterior is sampled similarly to the mean-field ((3.10)), but by normalising the ϵ term (projecting it onto a hypersphere) yielding a direction term $\epsilon/\|\epsilon\|$, and scaling it by the distance r . Therefore, for each layer, the radial posterior is sampled as

$$\begin{aligned} q(\theta) &= \text{Radial}(\mu, \sigma) \\ \theta &= \mu + \sigma \circ \frac{\epsilon}{\|\epsilon\|} r \end{aligned} \quad (3.15)$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $r \sim \mathcal{N}(0, 1)$. The entropy term of the Radial approximate posterior is given $\mathbb{H}[q(\theta)] = -\sum_i \log \sigma_i + \text{const}$ (and is therefore approximately equal to the entropy of the mean-field approximation up to a constant) and the cross-entropy term is calculated via Monte Carlo simulation as for mean-field VI (Section 3.3.3), but by sampling from a radial posterior $\theta^{(s)} \sim \text{Radial}(\mu, \sigma)$ as in (3.15).

3.3.7 Deep Ensembles

In deep ensembles [LPB17], M neural networks are trained with different initialisations. In this way, multiple local MAP solutions $\theta^{(m)}$ are obtained and equally considered, such that $p(\theta = \theta^{(m)} | \mathbf{y}) = 1/M$ for every $\theta^{(m)} \in \{\theta^{(1)}, \dots, \theta^{(M)}\}$ (and 0 elsewhere). To make predictions with a deep ensemble, we use the posterior predictive given by

$$\begin{aligned} p(\mathbf{y}_{\text{new}} | \mathbf{y}) &= \int p(\mathbf{y}_{\text{new}} | \theta) p(\theta | \mathbf{y}) d\theta \\ &= \frac{1}{M} \sum_{m=1}^M p(\mathbf{y}_{\text{new}} | \theta^{(m)}) \end{aligned} \quad (3.16)$$

where $p(\mathbf{y}_{\text{new}} | \theta^{(m)})$ are the normalised logits, or predicted probabilities, from model m in the ensemble.

3.3.8 MultiSWAG

Stochastic weight averaging estimates the final point estimate of the weights of a neural network θ_{SWA} as the average of the model parameters θ at the end of each epoch obtained via gradient descent after convergence has been reached.

For SWAG [Mad+19], the standard deviation of the model parameters σ_{SWA} is also calculated. From the SWA estimates of the model parameter mean and standard deviations, the posterior is then estimated as

$$p(\theta | \mathbf{y}) \approx \mathcal{N}(\theta_{\text{SWA}}, \sigma_{\text{SWA}}). \quad (3.17)$$

For MultiSWAG [WI20], M SWAG estimates are obtained from different initialisations, as in deep ensembles. Then, a Gaussian mixture model (GMM) is formed as a combination of each Gaussian SWAG posterior estimate $\mathcal{N}(\boldsymbol{\theta}_{\text{SWA}}^{(m)}, \boldsymbol{\sigma}_{\text{SWA}}^{(m)})$ where each component is equally weighed, as

$$p(\boldsymbol{\theta} | \mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^M \mathcal{N}(\boldsymbol{\theta}_{\text{SWA}}^{(m)}, \boldsymbol{\sigma}_{\text{SWA}}^{(m)}). \quad (3.18)$$

CHAPTER 4

Method

CHAPTER 5

Results

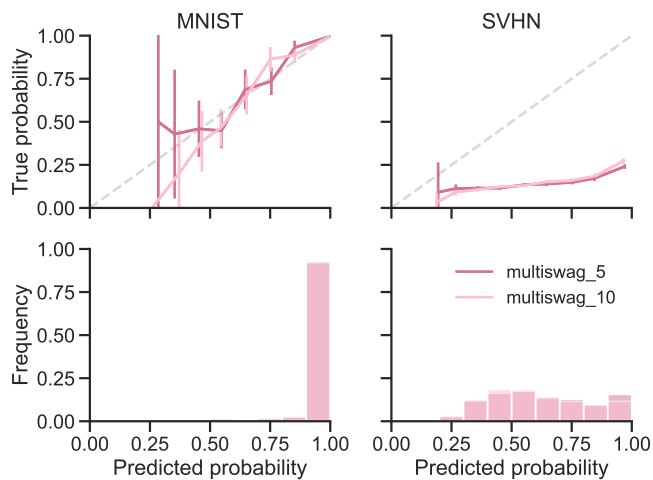


Figure 5.1: In-distribution calibration curves for various variational families with variational inference and Laplace approximation trained on MURA..

Table 5.1: ELBO and \hat{k} -statistic for variational inference using different variational families trained on the eight schools test problem, with three replicates for normalizing flows..

Type	Flows	$\mathcal{L}[q_\phi]$		\hat{k}		Params
		μ	σ	μ	σ	
Mean-field	—	-33.41	—	0.87	—	20
Full-rank	—	-32.60	—	0.82	—	110
Planar	4	-33.31	0.02	0.84	0.01	85
	8	-32.93	0.15	0.80	1e-3	169
	16	-32.25	0.11	0.73	0.03	337
	32	-31.84	0.06	0.65	0.05	673
Radial	4	-37.73	0.01	0.90	0.01	58
	8	-37.66	0.02	0.89	0.01	106
	16	-34.46	1.72	0.84	0.07	202
	32	-32.36	0.23	0.74	0.07	394
IAF	4	-31.53	0.09	0.57	0.03	63114
	8	-31.53	0.01	0.59	0.05	125918
	16	-31.61	0.06	0.54	0.01	251526
	32	-31.67	0.03	0.50	0.02	502742

CHAPTER 6

Experiments

CHAPTER 7

Conclusion

7.1 Conclusion and Outlook

7.1.1 Summary and key results

Remember to actually include some results and a quantitative look at the results that were obtained.

7.1.2 Outlook and future developments

Which avenues are most promising for future research? Either in order to solve the problem this paper tables or to solve the next problem.

APPENDIX A

An Appendix

Lorem ipsum dolor sit amet, consectetur adipisicing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Bibliography

- [CPD21] Beau Coker, Weiwei Pan, and Finale Doshi-Velez. “Wide Mean-Field Variational Bayesian Neural Networks Ignore the Data.” In: *International Conference on Machine Learning* (2021).
- [Dax+21] Erik Daxberger et al. “Laplace Redux-Effortless Bayesian Deep Learning.” In: *Advances in Neural Information Processing Systems* (2021).
- [FOG20] Sebastian Farquhar, Michael A Osborne, and Yarin Gal. “Radial bayesian neural networks: Beyond discrete support in large-scale bayesian deep learning.” In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2020.
- [FSG20] Sebastian Farquhar, Lewis Smith, and Yarin Gal. “Liberty or depth: Deep bayesian neural nets do not need complex weight posterior approximations.” In: *Advances in Neural Information Processing Systems* (2020).
- [LPB17] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. “Simple and scalable predictive uncertainty estimation using deep ensembles.” In: *Advances in neural information processing systems* (2017).
- [Mad+19] Wesley J Maddox et al. “A simple baseline for bayesian uncertainty in deep learning.” In: *Advances in Neural Information Processing Systems* 32 (2019).
- [WI20] Andrew G Wilson and Pavel Izmailov. “Bayesian deep learning and a probabilistic perspective of generalization.” In: *Advances in neural information processing systems* (2020).
- [Wid06] Jennifer Widom. *Tips for Writing Technical Papers*. 2006.
- [Wu+18] Anqi Wu et al. “Deterministic variational inference for robust bayesian neural networks.” In: *arXiv preprint arXiv:1810.03958* (2018).

