

Scalable Marginal Likelihood Estimation for Model Selection in Deep Learning

Maksim Tyurikov

MIPT, 2023

November 13, 2023

- 1 Motivation
- 2 Definition
- 3 Methods
- 4 Algorithm
- 5 Experiments
- 6 Literature

Motivation

Main idea

There are many approaches to selecting hyperparameters of models. Most of them rely on validation data, which may not be readily available. In this work, are presented a scalable marginal-likelihood estimation method to select both hyperparameters and network architectures, based on the training data alone.

Bayesian models

Bayesian models

We denote by $\mathbf{f}(\mathbf{x}, \boldsymbol{\theta})$ the C - dimensional real-valued output of a neural network with parameters $\boldsymbol{\theta} \in \mathbb{R}^P$, specified by a model \mathcal{M} which typically consists of a network architecture and hyperparameters.

Posterior distribution

$$p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{M}) \propto p(\mathcal{D}|\boldsymbol{\theta}, \mathcal{M})p(\boldsymbol{\theta}|\mathcal{M}) \quad (1)$$

A Bayesian model can then be defined using a likelihood and a prior, to get the posterior distribution.

Marginal likelihood

We assume that the data examples are sampled i.i.d. from $p(y_n|f(\mathbf{x}_n, \theta), \mathcal{M})$. The normalizing constant of the posterior, also known as the marginal likelihood, is given by the following expression:

$$p(\mathcal{D}|\mathcal{M}) = \int \prod_{i=1}^N p(\mathbf{y}_n|\mathbf{f}(\mathbf{x}_n, \theta), \mathcal{M})p(\theta|\mathcal{M})d\theta \quad (2)$$

The model \mathcal{M} might consist of the choice of network architecture (CNN, ResNet, etc.) and hyperparameters of the likelihood and prior, for example, observation noise and prior variances. Some of these are continuous parameters while others are discrete.

Empirical Bayes

A simple method is to pick the model that assigns the highest probability to the training data

$$\mathcal{M}_* = \arg \max_{\mathcal{M}} (p(\mathcal{D}|\mathcal{M})) \quad (3)$$

This procedure is also called type-II maximum likelihood estimation or empirical Bayes.

Laplace's method

Approximation

The method relies on a local quadratic approximation of $\log p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{M})$, around a maximum $\boldsymbol{\theta}_*$, resulting in a Gaussian approximation to $p(\boldsymbol{\theta}|\mathcal{D}, \mathcal{M})$, denoted by $q(\boldsymbol{\theta}|\mathcal{D}, \mathcal{M})$, and an approximation to the marginal likelihood, denoted by $q(\mathcal{D}|\mathcal{M})$ and shown below:

$$\log p(\mathcal{D}|\mathcal{M}) \simeq \log q(\mathcal{D}|\mathcal{M}) := \log p(\mathcal{D}, \boldsymbol{\theta}_*|\mathcal{M}) - \frac{1}{2} \log \left| \frac{1}{2\pi} \mathcal{H}_{\boldsymbol{\theta}_*} \right| \quad (4)$$

with Hessian: $\mathcal{H}_{\boldsymbol{\theta}_*} := \nabla_{\boldsymbol{\theta}\boldsymbol{\theta}}^2 \log p(\mathcal{D}, \boldsymbol{\theta}|\mathcal{M})$

Laplace's method

MargLik

However, computing \mathcal{H}_θ , a large matrix of size $P \times P$, and its determinant is infeasible in general. We refer to the log marginal likelihood as ***margLik***.

Hessian approximations

Generalized Gauss-Newton (GGN)

$$\mathcal{H}_\theta \simeq \mathcal{H}_\theta^{\text{GGN}} = \mathcal{J}_\theta^T \mathcal{L}_\theta \mathcal{J}_\theta + \mathcal{P}_\theta \quad (5)$$

The complexity is $\mathcal{O}(P^2 NC + PNC^2)$

Empirical Fisher (EF)

$$\mathcal{H}_\theta \simeq \mathcal{H}_\theta^{\text{EF}} = \mathcal{G}_\theta^T \mathcal{G}_\theta + \mathcal{P}_\theta \quad (6)$$

The complexity is $\mathcal{O}(P^2 N)$

MAP

To optimize the network parameters θ we perform regular neural network training on the maximum a posteriori (MAP).

$$\log p(\mathcal{D}, \theta | \mathcal{M}) = \sum_{n=1}^N p(\mathbf{y}_n | \mathbf{f}(\mathbf{x}_n, \theta)) + \log p(\theta) \quad (7)$$

Continuous hyperparameters

To optimize the continuous hyperparameters \mathcal{M} , we perform gradient ascent on the marginal likelihood estimate

$$\mathcal{M}^\partial \leftarrow \mathcal{M}^\partial + \gamma \nabla_{\mathcal{M}^\partial} \log q(\mathcal{D} | \mathcal{M}) \quad (8)$$

Algorithm

Algorithm 1 Marginal likelihood based training

- 1: **Input:** dataset \mathcal{D} , likelihood $p(\mathcal{D}|\theta, \mathcal{M})$, prior $p(\theta|\mathcal{M})$. Initial model \mathcal{M} , step size γ , steps K , burn-in epochs B , marglik frequency F . GGN or EF.
- 2: initialize θ of the neural network
- 3: **for** each *epoch* **do**
- 4: $\theta \leftarrow \text{trainEpoch}(\text{objective in Eq. 6})$
- 5: **if** $\text{epoch} > B$ and $\text{epoch} \bmod F = 0$ **then**
- 6: compute $\log q(\mathcal{D}|\mathcal{M})$ (Eq. 3) with $\mathbf{H}_{\theta}^{\text{GGN}}$ or $\mathbf{H}_{\theta}^{\text{EF}}$
- 7: **for** K steps **do**
- 8: $\mathcal{M}^{\partial} \leftarrow \mathcal{M}^{\partial} + \gamma \nabla_{\mathcal{M}^{\partial}} \log q(\mathcal{D}|\mathcal{M})$
- 9: update \mathcal{M} and $\log q(\mathcal{D}|\mathcal{M})$
- 10: **end for**
- 11: **end if**
- 12: **end for**
- 13: Return marginal likelihood $\log q(\mathcal{D}|\mathcal{M})$ and optionally posterior approximation $q(\theta|\mathcal{D}, \mathcal{M})$.

Example approximation

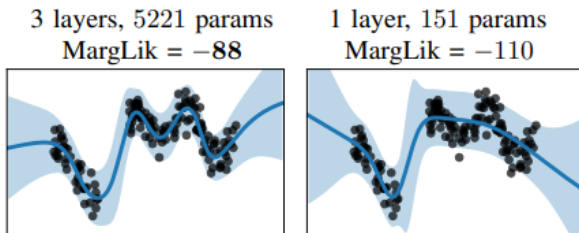


Figure: A deeper 3-layer network (left) has a better marginal likelihood compared to a 1-layer network (right). This agrees with the fit where the deeper network appears to explain the ‘sinusoidal’ trend better.

The connection between accuracy and MargLik

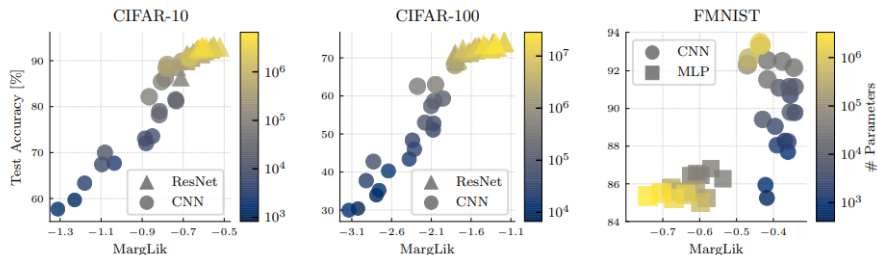


Figure: Each dot above shows a model of different size and/or architecture (around 40 models per plot of varying widths and depths). Models with higher training marginal-likelihood tend to have higher test accuracy. For similar performance, smaller models tend to have a higher marginal-likelihood as desired. Marker size and color changes with the number of parameters.

- 1 **Main article** Scalable Marginal Likelihood Estimation for Model Selection in Deep Learning.