

SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Thesis title

Author

SCHOOL OF COMPUTATION,
INFORMATION AND TECHNOLOGY —
INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Master's Thesis in Informatics

Thesis title

Titel der Abschlussarbeit

Author:	Author
Examiner:	Supervisor
Supervisor:	Advisor
Submission Date:	Submission date

I confirm that this master's thesis is my own work and I have documented all sources and material used.

Munich, Submission date

Author

Abstract

Contents

Abstract	iii
1 Introduction	2
1.1 Section	2
2 Bayesian Deep Learning	3
2.1 General Concepts	3
2.2 Priors	4
2.3 Inference	4
2.4 Properties of Bayesian Neural Network Posteriors	5
3 Flow Models	6
3.1 Overview	6
3.2 Flow Matching	6
3.2.1 The Objective	7
3.2.2 Training	8
3.2.3 Couplings and Conditional Paths	8
3.2.4 Connection to Diffusion (maybe remove)	9
3.2.5 Extensions (mini lit review of example uses)	9
3.3 Sampling / Integration	9
3.4 Computing Likelihoods	9
3.4.1 Faster Likelihoods Through Trace Estimation	10
4 Geometry of Neural Networks	11
4.1 Symmetries of Neural Networks	11
4.2 Canonical Representations of Neural Networks	11
5 Graph Neural Networks & MetaNets	12
6 Generative Models in Weight-Space	13
7 Method & Design Choices	14
8 Results	15

Contents

9 Discussion	16
10 Conclusion & Future Work	17
Abbreviations	18
List of Figures	19
List of Tables	20
Bibliography	21

Notation	Explanation
θ	Parameters of a neural network
$(x, y) \in \mathcal{D}$	Dataset with inputs x and labels/targets y
v_t	Time-dependent vector field
J_v	Jacobian of the vector field v
ϕ_t	Flow map from time 0 to time t

Table 1: Summary of notation used throughout the thesis.

1 Introduction

1.1 Section

2 Bayesian Deep Learning

A primary use case for a generative model over neural network weights is in Bayesian deep learning, where it can allow efficient inference by transporting the prior distribution to the posterior. Thus, to motivate the rest of the discussion, we first give an overview of concepts from Bayesian deep learning. Section 2.1 is a general introduction. It is followed by a review of inference methods (Laplace approximations, variational inference, MCMC-based methods) typically used for Bayesian neural networks, and we conclude with a review of literature around Bayesian deep learning particularly relevant for our work in Section 2.4.

2.1 General Concepts

In typical Bayesian fashion, Bayesian deep learning (refer to (MacKay, 1992; Neal, 1996) for foundational work and (Goan and Fookes, 2020; Arbel et al., 2023) for more recent reviews) aims to quantify the uncertainty in neural networks through probability distributions over their parameters, rather than obtaining a single solution by an SGD-like optimization method. Then, given the *posterior* distribution $p(\theta|\mathcal{D})$ over weights θ conditioned on the dataset \mathcal{D} , predictions are obtained via *Bayesian model averaging*:

$$p(y|x, \mathcal{D}) = \mathbb{E}_{\theta \sim p(\theta|\mathcal{D})} [p(y|x, \theta)] = \int p(y|x, \theta) p(\theta|\mathcal{D}) d\theta, \quad (2.1)$$

where the prediction is averaged over the posterior over the weights. Note that since the forward pass through the model $p(y|x, \theta)$, is deterministic, the uncertainty in predictions results solely from the uncertainty over parameters. To bring things together, Bayesian inference over neural network weights consists of three steps:

1. Specify prior $p(\theta)$.
2. Compute/sample posterior $p(\theta|\mathcal{D}) \propto p(\mathcal{D}|\theta)p(\theta)$.
3. Average predictions over the posterior.

The last step only requires forward passes through the model and thus is straightforward. The first two steps require deeper consideration.

2.2 Priors

Specifying a prior mainly consists of two choices: specifying an architecture, and specifying a probability distribution over the weights. The distribution is typically taken to be an isotropic Gaussian, which is an uninformative prior but straightforward to work with.

Different architectural decisions also result in different functions, even if the flattened weight vectors are identical, meaning that the choice of an architecture further specifies a prior in function space. As a simple example, keeping the depth and width of a neural network constant, even just changing the activation function from a ReLU to a sigmoid results in a different distribution of functions. The functional distribution can also be specified in a more deliberate way; e.g. a translation-invariant convolutional neural network, or a group-equivariant network (Cohen and Welling, 2016) puts probability mass only on functions satisfying certain equivariance constraints depending on the task at hand.

We keep this discussion short since the choice of a prior has tangential impact in the rest of the presentation, and we refer to recent reviews such as (Fortuin, 2022) for a more detailed treatment of Bayesian neural network priors.

2.3 Inference

- Variational methods
- MCMC-based methods
- Transformation-based methods (normalizing flow)

Laplace

VI

DE

MCMC introduction

HMC

Highlight distinction between chain-based and transformation based sampling, using normalizing flows and connections to optimal transport etc

– Stochastic –

SG-MCMC

SG-HMC

2.4 Properties of Bayesian Neural Network Posteriors

How they are like Boltzmann distributions, not arbitrary.

3 Flow Models

3.1 Overview

We train our flow using *flow matching* (Lipman et al., 2023; Albergo et al., 2023; Liu et al., 2022), which generalizes diffusion models with a more flexible design space. In this section we first formulate the flow matching objective (Sec. 3.2), explain the design choices it enables (Sec. 3.2.3), and describe in more detail how samples (Sec. 3.3) and likelihoods (Sec. 3.4) can be obtained using a flow model.

3.2 Flow Matching

Flow matching, first proposed in (Lipman et al., 2023; Albergo et al., 2023; Liu et al., 2022) aims to solve the problem of *dynamic transport*, i.e. finding a time-dependent vector field to transport the source (prior) distribution p_0 to the target (data) distribution p_1 . More formally, the vector field $u_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ leads to the ordinary differential equation (ODE)

$$dx = u_t(x)dt \quad (3.1)$$

and induces a *flow* $\phi : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ that gives the solution to the ODE at time t with starting point x_0 , such that

$$\frac{d}{dt}\phi_t(x_0) = u_t(\phi_t(x)) \quad (3.2)$$

$$\phi_0(x_0) = x_0. \quad (3.3)$$

Starting with p_0 , transformed distributions p_t can then be defined using this flow with the push-forward operation

$$p_t := [\phi_t]_{\#}(p_0) \quad (3.4)$$

and the instantaneous change in the density satisfies the *continuity equation*

$$\frac{\partial p}{\partial t} = -\nabla \cdot (p_t u_t) \quad (3.5)$$

which means that probability mass is conserved during the transformation. With these formulations, we say the vector field u_t *generates* the *probability path* (also called *interpolant*) p_t .

3.2.1 The Objective

The formulation above could also be applied to traditional continuous normalizing flows (CNFs) (Chen et al., 2018) as well, and flow matching is an instantiation of CNFs. However, continuous normalizing flows have in the past been trained using objectives which required solving and then backpropagating through the ODE, such as KL-divergence or other likelihood-based objectives, which made training costly. This problem was later addressed with diffusion models and their simpler regression objectives such as score matching and denoising (Sohl-Dickstein et al., 2015; Song et al., 2021; Ho et al., 2020) that proved to be very effective. As we will now demonstrate, the flow matching objective is also formulated as a simulation-free regression objective, and is more flexible than the diffusion objectives.

As explained in Section 3.2, the goal in flow matching is to learn a vector field $v_\theta : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ parametrized by a neural network.¹ If we know the ground truth vector field u and can sample from the intermediate p_t 's, we can directly optimize the flow matching objective

$$\mathcal{L}_{\text{FM}}(\theta) := \mathbb{E}_{t \sim \mathcal{U}(0,1), x_t \sim p_t(x)} \|v_\theta(t, x_t) - u_t(x_t)\|^2 \quad (3.6)$$

by first sampling a time point t and then $x_t \sim p_t$. However in practice, we neither have a closed form expression for u nor can sample from an arbitrary p_t without integrating the flow.

The *conditional flow matching* (CFM) framework first introduced in (Lipman et al., 2023) and then extended in (Tong et al., 2023) solves this problem by formulating the intermediate probability paths as mixtures of simpler paths,

$$p_t(x) = \int p_t(x | z) q(z) dz \quad (3.7)$$

where z is the conditioning variable and $q(z)$ a distribution over z (e.g. with $z := (x_0, x_1)$, $q(z) = p_0(x_0)p_1(x_1)$, $u_t(x | z) = x_1 - x_0$, and $p_t(x | z) = \mathcal{N}(x | (1-t)x_0 + tx_1, \sigma^2)$). Then similar to how p_t 's were generated by the vector field u_t , the conditional probability paths $p_t(x | z)$ are generated by conditional vector fields $u_t(x | z)$, and as shown in (Tong et al., 2023) u_t can be decomposed in terms of these conditional vector fields as

$$u_t(x) = \mathbb{E}_{z \sim q(z)} \frac{u_t(x | z) p_t(x | z)}{p_t(x)}. \quad (3.8)$$

Then similar to Equation 3.6, we have the conditional flow matching objective

$$\mathcal{L}_{\text{CFM}}(\theta) := \mathbb{E}_{t \sim \mathcal{U}(0,1), z \sim q(z), x_t \sim p_t(x|z)} \|v_\theta(t, x_t) - u_t(x_t | z)\|^2. \quad (3.9)$$

¹For conciseness, we interchangeably use the subscripts for vector fields to denote time ($u_t(x)$) and parameters ($v_\theta(t, x)$).

That is, we first sample a conditioning variable z , and then regress to the *conditional* vector field $u_t(x | z)$. Thus we obtain a tractable objective by defining sample-able conditional probability paths and a tractable conditional vector field. Moreover, as shown in (Tong et al., 2023), the FM and CFM objectives equivalent up to a constant and thus

$$\nabla_{\theta} \mathcal{L}_{\text{FM}}(\theta) = \nabla_{\theta} \mathcal{L}_{\text{CFM}}(\theta), \quad (3.10)$$

meaning we do not lose the expressive power of the FM objective by regressing only to the conditional vector fields. As we will show in Section 3.2.3, the choice of these conditional probability paths, vector fields, along with the conditioning variable itself, makes the flow matching approach particularly flexible.

3.2.2 Training

To sum up the discussion in the previous section, a step of training a flow model using CFM objective (Equation 3.9) proceeds as follows:

1. Sample $t \sim \mathcal{U}(0, 1)$, $z \sim q(z)$, and $x_t \sim p_t(x | z)$.
2. Compute $\mathcal{L}_{\text{CFM}}(\theta) = \|v_{\theta}(t, x_t) - u_t(x_t | z)\|^2$.
3. Update θ with $\nabla_{\theta} \mathcal{L}_{\text{CFM}}(\theta)$.

3.2.3 Couplings and Conditional Paths

With this framework established, the three main design choices for building a flow matching model are choosing the coupling $q(z)$, the conditional “ground truth” vector field $u_t(x | z)$, and the conditional probability paths $p_t(x | z)$. Starting with an arbitrary source distribution p_0 and target distribution p_1 , Tong et al. (2023) propose three different ways of constructing conditional paths from couplings between p_0 and p_1 , of which we focus on two (independent and optimal transport couplings). In all setups, the condition variable z corresponds to a pair (x_0, x_1) of source and target points.

Independent Coupling. The simplest way of obtaining is to sample independently from p_0 and p_1 ; i.e. $q(z) = p_0(x_0)p_1(x_1)$, with the conditional paths and the vector field defined as

$$p_t(x | z) = \mathcal{N}(x | (1 - t)x_0 + tx_1, \sigma^2) \quad (3.11)$$

$$u_t(x | z) = x_1 - x_0. \quad (3.12)$$

The conditional paths and the coupling defined this way are easily easy to sample from, but have undesirable properties such as crossing paths which which can lead to high

variance in the ground truth vector field for a specific point and time. Moreover in practice, independent couplings can lead to curved paths that incur higher integration errors, as there is no notion of straightness considered in this formulation.

Optimal Transport. To obtain straighter and shorter paths that are easier to integrate, Tong et al. (2023) propose to use the static 2-Wasserstein optimal transport map π as the coupling; i.e.

$$q(z) = \pi(x_0, x_1), \quad (3.13)$$

with the conditional paths and vector field defined as in Equations 3.11 and 3.12. The flow model thus obtained solves the dynamic optimal OT problem as $\sigma^2 \rightarrow 0$ (Proposition 3.4 in (Tong et al., 2023)). However, computing the exact OT map for the entire dataset is challenging, especially in high dimensions as in our problem. It can instead be approximated using mini-batches (Fatraş et al., 2021). This means at the end the OT problem is solved only to an approximation, but nevertheless results in straighter paths that cross less often, since intuitively an $x_0 \sim p_0$ is more likely to be coupled with $x_1 \sim p_1$ closer to it rather than an x_1 chosen uniformly random.

3.2.4 Connection to Diffusion (maybe remove)

3.2.5 Extensions (mini lit review of example uses)

3.3 Sampling / Integration

3.4 Computing Likelihoods

In earlier normalizing flows that aim to learn a static mapping between the two distributions (Rezende and Mohamed, 2015), given source samples $x_0 \sim p_0$, likelihoods of the generated samples $z = f(x_0) \approx p_1$ can be computed exactly via the change of variables formula

$$\log p_1(z) = \log p_0(z) - \log \det |J_f(z)| \quad (3.14)$$

where J_f is the Jacobian of f . Thus, we can obtain exact likelihoods for the generated samples by taking the determinant of the Jacobian of the normalizing flow. Since Jacobian computations can be costly, this has motivated work on designing normalizing flows with easier to compute Jacobians, such as RealNVP (Dinh et al., 2017).

In a *continuous normalizing flow* on the other hand, the *instantaneous change of variables* formula (Chen et al., 2018) defines the change in probability mass through time. Given that the vector field v_t is continuous in t and uniformly Lipschitz continuous in \mathbb{R}^d , it

holds that

$$\frac{d \log p_t(\phi_t(x))}{dt} = -\operatorname{div}(v_t(\phi_t(x))) \quad (3.15)$$

$$= -\operatorname{Tr} \left(\frac{dv_t(\phi_t(x))}{dt} \right) \quad (3.16)$$

where $\frac{dv_t(\phi_t(x))}{dt} =: J_v(\phi_t(x))$ is the Jacobian of the vector field. Then we integrate over time to compute the full change in probability:

$$\log p_1(\phi_1(x)) = \log p_0(\phi_0(x)) - \int_0^1 \operatorname{Tr}(J_v(\phi_t(x))) dt. \quad (3.17)$$

Then we can integrate the Jacobian trace of the vector field through time (simultaneously with sampling) to obtain exact likelihoods for the generated samples.

3.4.1 Faster Likelihoods Through Trace Estimation

However, materializing the full Jacobian of the vector field can be prohibitively expensive, especially if the task is high dimensional (as in our case) since the log determinant computation has a time complexity of $O(d^3)$ (Grathwohl et al., 2018) without any restrictions on the structure of the Jacobian.

To alleviate this problem, (Grathwohl et al., 2018) propose to use the *Hutchinson trace estimator* (Hutchinson, 1990) for an unbiased estimate of the Jacobian trace of a square matrix:

$$\operatorname{Tr}(J_v) = \mathbb{E}_{p(\epsilon)} \left[\epsilon^T J_v \epsilon \right] \quad (3.18)$$

where $p(\epsilon)$ is chosen such that $\mathbb{E}[\epsilon] = 0$ and $\operatorname{Cov}(\epsilon) = I$, typically a Gaussian or a Rademacher distribution. Then, we can use this estimator in place of the explicit trace computation in Equation 3.17 and compute the likelihoods as

$$\log p_1(\phi_1(x)) = \log p_0(\phi_0(x)) - \int_0^1 \mathbb{E}_{p(\epsilon)} \left[\epsilon^T J_v(\phi_t(x)) \epsilon \right] dt. \quad (3.19)$$

The performance benefit of using the Hutchinson trace estimator results from the fact that the Jacobian-vector product $J_v \epsilon$ can be computed very efficiently by automatic differentiation (Baydin et al., 2018), giving the whole approach a time complexity of $O(d)$ only. Due to this significant performance improvement and being an unbiased estimate, the Hutchinson trace estimator has been widely used in the diffusion/flow model literature (Lipman et al., 2023; Song et al., 2021).

4 Geometry of Neural Networks

4.1 Symmetries of Neural Networks

Also introduce mode connectivity here, before canonical reps

Introduce permutation and scaling symmetries, and maybe other symmetries if we end up using them.

4.2 Canonical Representations of Neural Networks

Mainly based on the deep toroids paper

5 Graph Neural Networks & MetaNets

6 Generative Models in Weight-Space

7 Method & Design Choices

8 Results

9 Discussion

10 Conclusion & Future Work

Abbreviations

List of Figures

List of Tables

1	Summary of notation used throughout the thesis.	1
---	---	---

Bibliography

- Albergo, Michael S. et al. (2023). *Stochastic Interpolants: A Unifying Framework for Flows and Diffusions*. DOI: [10.48550/arXiv.2303.08797](#). arXiv: [2303.08797 \[cond-mat\]](#).
- Arbel, Julyan et al. (2023). *A Primer on Bayesian Neural Networks: Review and Debates*. DOI: [10.48550/arXiv.2309.16314](#). arXiv: [2309.16314 \[cs, math, stat\]](#).
- Baydin, Atilim Gunes et al. (2018). “Automatic Differentiation in Machine Learning: A Survey.” In: *Journal of Machine Learning Research*.
- Chen, Ricky T. Q. et al. (2018). “Neural Ordinary Differential Equations.” In: *Advances in Neural Information Processing Systems*. Vol. 31. Curran Associates, Inc.
- Cohen, Taco and Max Welling (2016). “Group Equivariant Convolutional Networks.” In: *Proceedings of The 33rd International Conference on Machine Learning*. PMLR, pp. 2990–2999.
- Dinh, Laurent et al. (2017). “Density Estimation Using Real NVP.” In: *International Conference on Learning Representations*.
- Fatras, Kilian et al. (2021). *Minibatch Optimal Transport Distances; Analysis and Applications*. arXiv: [2101.01792 \[cs, stat\]](#).
- Fortuin, Vincent (2022). “Priors in Bayesian Deep Learning: A Review.” In: *International Statistical Review* 90.3, pp. 563–591. ISSN: 0306-7734, 1751-5823. DOI: [10.1111/insr.12502](#).
- Goan, Ethan and Clinton Fookes (2020). “Bayesian Neural Networks: An Introduction and Survey.” In: vol. 2259, pp. 45–87. DOI: [10.1007/978-3-030-42553-1_3](#). arXiv: [2006.12024 \[cs, stat\]](#).
- Grathwohl, Will et al. (2018). *FFJORD: Free-form Continuous Dynamics for Scalable Reversible Generative Models*. DOI: [10.48550/arXiv.1810.01367](#). arXiv: [1810.01367](#).
- Ho, Jonathan et al. (2020). “Denoising Diffusion Probabilistic Models.” In: *Advances in Neural Information Processing Systems*. Vol. 33. Curran Associates, Inc., pp. 6840–6851.
- Hutchinson, M.F. (1990). “A Stochastic Estimator of the Trace of the Influence Matrix for Laplacian Smoothing Splines.” In: *Communications in Statistics - Simulation and Computation* 19.2, pp. 433–450. ISSN: 0361-0918. DOI: [10.1080/03610919008812866](#).
- Lipman, Yaron et al. (2023). *Flow Matching for Generative Modeling*. DOI: [10.48550/arXiv.2210.02747](#). arXiv: [2210.02747 \[cs, stat\]](#).
- Liu, Xingchao et al. (2022). *Flow Straight and Fast: Learning to Generate and Transfer Data with Rectified Flow*. DOI: [10.48550/arXiv.2209.03003](#). arXiv: [2209.03003](#).

- MacKay, David J C (1992). “Bayesian Methods for Adaptive Models.” In.
- Neal, Radford M. (1996). *Bayesian Learning for Neural Networks*.
- Rezende, Danilo and Shakir Mohamed (2015). “Variational Inference with Normalizing Flows.” In: *Proceedings of the 32nd International Conference on Machine Learning*. PMLR, pp. 1530–1538.
- Sohl-Dickstein, Jascha et al. (2015). “Deep Unsupervised Learning Using Nonequilibrium Thermodynamics.” In.
- Song, Yang et al. (2021). *Score-Based Generative Modeling through Stochastic Differential Equations*. arXiv: [2011.13456 \[cs, stat\]](#).
- Tong, Alexander et al. (2023). *Improving and Generalizing Flow-Based Generative Models with Minibatch Optimal Transport*. DOI: [10.48550/arXiv.2302.00482](#). arXiv: [2302.00482 \[cs\]](#).