

# Bayesian Neural Model Selection for Symmetry Learning

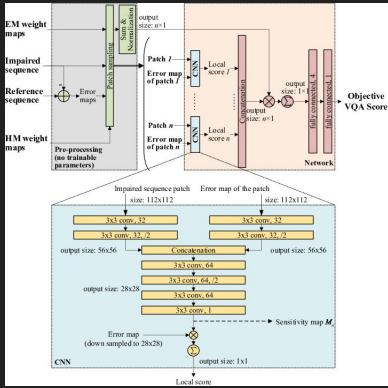
Tycho van der Ouderaa

@ JSM 2024 in Portland, US

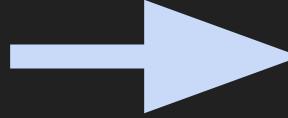
- I. What is Bayesian model selection?
- II. How to scale to deep neural networks
- III. Learning inductive bias and symmetries

# Future vision of ML

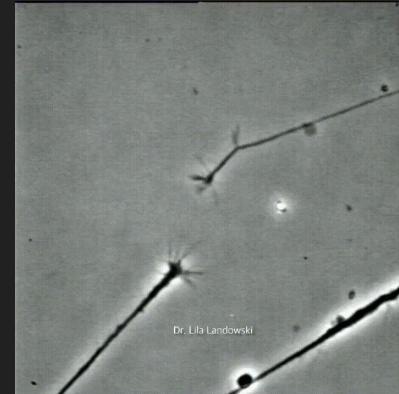
Architectures determined by  
**trial-and-error**



Bayesian model selection



Morphable architectures  
**automatic inductive bias**

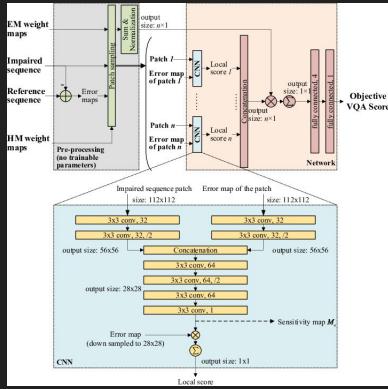


References:

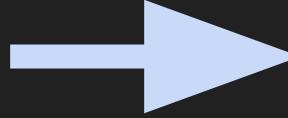
Left: *Bridge the Gap Between VQA and Human Behavior on Omnidirectional Video: A Large-Scale Dataset and a Deep Learning Model*, Chen Li, Mai Xu, Xinzhe Du, Zulin Wang  
Right video: "Neurons connecting to one another in a Petri dish - growth cones." by Dr Lila Landowski

# Future vision of ML

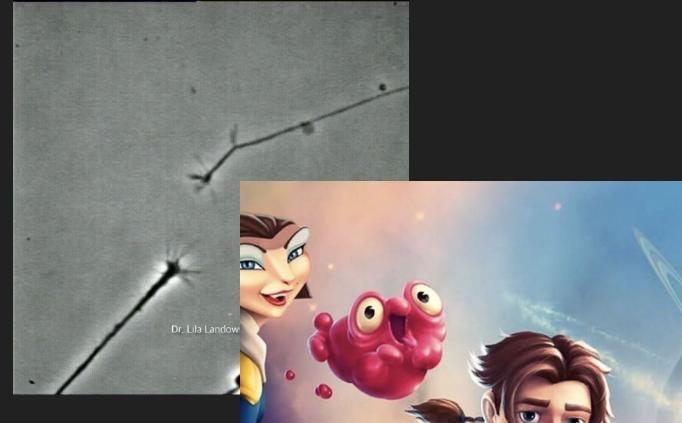
Architectures determined by  
**trial-and-error**



Bayesian model selection



Morphable architectures  
**automatic inductive bias**



References:

Left: *Bridge the Gap Between VQA and Human Behavior on Omnidirectional Video: A Large-Scale Dataset and a Deep Learning Model*, Chen Li, Mai Xu, Xinzhe Du, Zulin Wang

Right video: "Neurons connecting to one another in a Petri dish - growth cones." by Dr Lila Landowski

Right picture: Wells, R. (Producer), Clements, R., & Musker, J. (Directors). (2002). *Treasure Planet* [Film]. Walt Disney Pictures.

## Fitting a neural network

$$\text{NN}_{\theta} : \mathbb{R}^M \rightarrow \mathbb{R}^N$$

# Regular training (maximum likelihood / MAP)

Likelihood:  $p(\mathcal{D} \mid \theta)$

Prior:  $p(\theta)$

# Regular training (maximum likelihood / MAP)

Likelihood:  $p(\mathcal{D} \mid \boldsymbol{\theta})$

Prior:  $p(\boldsymbol{\theta})$

Posterior:  $p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})}$

# Regular training (maximum likelihood / MAP)

Likelihood:  $p(\mathcal{D} \mid \boldsymbol{\theta})$

Prior:  $p(\boldsymbol{\theta})$

Posterior:  $p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})}$

Optimal  $\boldsymbol{\theta}_* = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta} \mid \mathcal{D})$

$$= \operatorname{argmax}_{\boldsymbol{\theta}} \log \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathcal{D})} = \log p(\mathcal{D} \mid \boldsymbol{\theta}) + \log p(\boldsymbol{\theta})$$

$$= \operatorname{argmin}_{\boldsymbol{\theta}} \underbrace{-\log p(\mathcal{D} \mid \boldsymbol{\theta})}_{\text{negative log likelihood}} \underbrace{-\log p(\boldsymbol{\theta})}_{\text{regularizer}}$$

Regular loss

# Supervised training (maximum likelihood / MAP)

Likelihood:  $p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) = \prod_{n=1}^N \mathcal{N}(y_n \mid \text{NN}_{\boldsymbol{\theta}}(x_n), \sigma^2 \mathbf{I})$

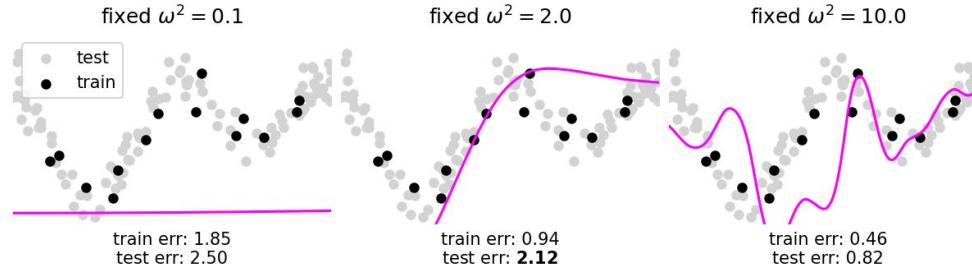
Prior:  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta} \mid \mathbf{0}, \sigma_{\text{prior}}^2 \mathbf{I})$

Posterior:  $p(\boldsymbol{\theta} \mid \mathbf{x}, \mathbf{y})$

Optimal  $\boldsymbol{\theta}_* = \operatorname{argmin}_{\boldsymbol{\theta}} \underbrace{\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \text{NN}_{\boldsymbol{\theta}}(x_n))^2}_{\text{mean squared error}} + \underbrace{\frac{1}{2\sigma_{\text{prior}}^2} \|\boldsymbol{\theta}\|_2^2}_{\text{weight-decay}}$

# Bayesian model selection

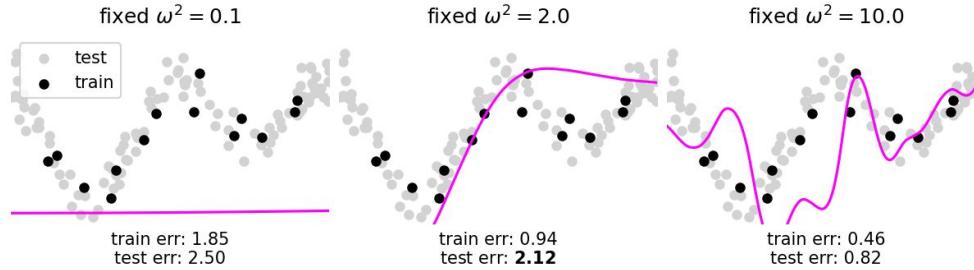
## Regular training (maximum likelihood)



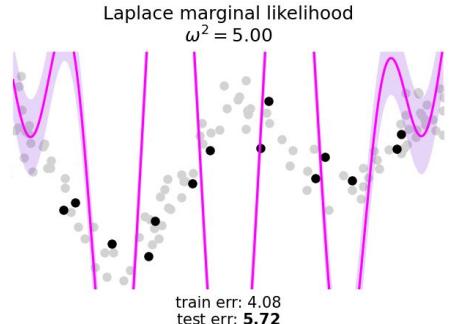
$$\text{NN}_{\theta}(x) = \mathbf{W}_3 \text{relu}(\mathbf{W}_2 \cos(\omega^2 \mathbf{W}_1 x))$$

# Bayesian model selection

## Regular training (maximum likelihood)



## Bayesian model selection (marginal likelihood)



Uses gradients.  
No expensive cross-validation

$$\text{NN}_{\theta}(\mathbf{x}) = \mathbf{W}_3 \text{relu}(\mathbf{W}_2 \cos(\omega^2 \mathbf{W}_1 \mathbf{x}))$$

Learn how to generalize

# Bayesian model selection (marginal likelihood)

$$p(\boldsymbol{\theta}, \boldsymbol{\eta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathcal{D} \mid \boldsymbol{\eta})} = \underbrace{\frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})}{p(\mathcal{D} \mid \boldsymbol{\eta})}}_{\text{usual posterior}} \underbrace{\frac{p(\mathcal{D} \mid \boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathcal{D})}}_{\text{hyper posterior}}$$

# Bayesian model selection (marginal likelihood)

$$p(\boldsymbol{\theta}, \boldsymbol{\eta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathcal{D} \mid \boldsymbol{\eta})} = \underbrace{\frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})}{p(\mathcal{D} \mid \boldsymbol{\eta})}}_{\text{usual posterior}} \underbrace{\frac{p(\mathcal{D} \mid \boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathcal{D})}}_{\text{hyper posterior}}$$

Marginal likelihood  $p(\mathcal{D} \mid \boldsymbol{\eta}) = \int_{\boldsymbol{\theta}} p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})d\boldsymbol{\theta}$

# Bayesian model selection (marginal likelihood)

$$p(\boldsymbol{\theta}, \boldsymbol{\eta} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathcal{D} \mid \boldsymbol{\eta})} = \underbrace{\frac{p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})}{p(\mathcal{D} \mid \boldsymbol{\eta})}}_{\text{usual posterior}} \underbrace{\frac{p(\mathcal{D} \mid \boldsymbol{\eta})p(\boldsymbol{\eta})}{p(\mathcal{D})}}_{\text{hyper posterior}}$$

Marginal likelihood  $p(\mathcal{D} \mid \boldsymbol{\eta}) = \int_{\boldsymbol{\theta}} p(\mathcal{D} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid \boldsymbol{\eta})d\boldsymbol{\theta}$

Optimise  $\boldsymbol{\eta}_* = \operatorname{argmax}_{\boldsymbol{\eta}} \log p(\mathcal{D} \mid \boldsymbol{\eta})$  [Type-II ML / Empirical Bayes]

Allows optimization of architecture/inductive bias with gradients!

How to integrate out neural network parameters?

# Gaussian approximation

Approximate posterior with a Gaussian

$$p(\boldsymbol{\theta} \mid \mathcal{D}, \boldsymbol{\eta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}})$$

How to find  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  ?

# Approximate Bayesian inference

Laplace approximation

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{MAP}}, \mathbf{H}^{-1})$$

Variational Inference

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

# Approximate Bayesian inference

## Laplace approximation

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \underline{\boldsymbol{\theta}_{\text{MAP}}}, \underline{\mathbf{H}^{-1}})$$

(local)

Pretrained network + Hessian

## Variational Inference

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}})$$

(global)

Optimize ELBO objective

# Only correlate rows and columns

Independence between layers

$$\mathbb{E}[\mathbf{H}] = \bigoplus_{l=1}^L \mathbf{H}_l$$

Independence between inputs and outputs

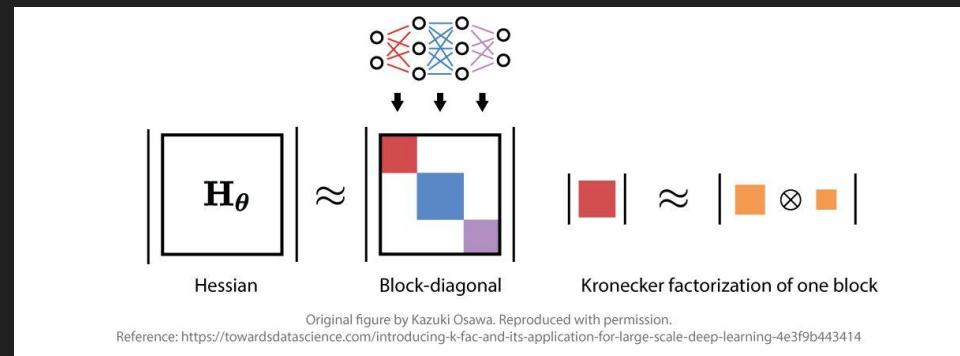
$$\mathbb{E}[\mathbf{H}] = \mathbb{E}[\mathbf{g}_l \mathbf{g}_l^T \otimes \mathbf{x}_l \mathbf{x}_l^T] \approx \mathbb{E}[\mathbf{g}_l \mathbf{g}_l^T] \otimes \mathbb{E}[\mathbf{x}_l \mathbf{x}_l^T]$$

$$\mathbb{E}[\mathbf{g}_l \mathbf{g}_l^T] \approx \frac{1}{N} \sum_{n=1}^N \mathbf{g}_l^{(n)} (\mathbf{g}_l^{(n)})^T = \mathbf{S}_l$$

$$\mathbb{E}[\mathbf{x}_l \mathbf{x}_l^T] \approx \frac{1}{N} \sum_{n=1}^N \mathbf{x}_l^{(n)} (\mathbf{x}_l^{(n)})^T = \mathbf{A}_l$$

KFAC Hessian approximation:

$$\mathbb{E}[\mathbf{H}] = \bigoplus_{l=1}^L (\mathbf{S}_l \otimes \mathbf{A}_l)$$



Original figure by Kazuki Osawa. Reproduced with permission.  
Reference: <https://towardsdatascience.com/introducing-k-fac-and-its-application-for-large-scale-deep-learning-4e3f9b443414>

# Approximate Bayesian inference

## Laplace approximation

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \underline{\boldsymbol{\theta}_{\text{MAP}}, \mathbf{H}^{-1}})$$

$$\mathbf{H} = \bigoplus_{l=1}^L (\mathbf{S}_l \otimes \mathbf{A}_l) \quad [1, 2]$$

## Variational Inference

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \underline{\boldsymbol{\mu}, \Sigma})$$

$$\Sigma^{-1} = \bigoplus_{l=1}^L (\mathbf{S}_l \otimes \mathbf{A}_l) \quad [3]$$

[1] Optimizing Neural Networks with Kronecker-factored Approximate Curvature

James Martens and Roger Grosse

In ICML 2016

[2] Scalable Marginal Likelihood Estimation for Model Selection in Deep Learning

Alexander Immer, Matthias Bauer, Vincent Fortuin, Gunnar Ratsch, Mohammad Emtiyaz Khan

In ICML 2021

[3] Structured and efficient variational deep learning with matrix gaussian posteriors

Christos Louizos and Max Welling

In ICML 2016

# Approximate Bayesian inference

## Laplace approximation

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{MAP}}, \mathbf{H}^{-1})$$

$$\mathbf{H} = \bigoplus_{l=1}^L (\mathbf{S}_l \otimes \mathbf{A}_l)$$

## Variational Inference

$$q(\boldsymbol{\theta}) \approx \mathcal{N}(\boldsymbol{\theta} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\Sigma}^{-1} = \bigoplus_{l=1}^L (\mathbf{S}_l \otimes \mathbf{A}_l)$$

[4] Stochastic Marginal Likelihood Gradients using Neural Tangent Kernels

Alexander Immer, Tycho F.A. van der Ouderaa, Mark van der Wilk, Gunnar Ratsch, Bernhard Schölkopf

In ICML 2023

[5] Variational Inference Failures Under Model Symmetries: Permutation Invariant Posteriors for Bayesian Neural Networks

Yoav Gelberg, Tycho F. A. van der Ouderaa, Mark van der Wilk, Yarin Gal

In ICML 2024, GRaM workshop (best paper)

## Tricks of the trade:

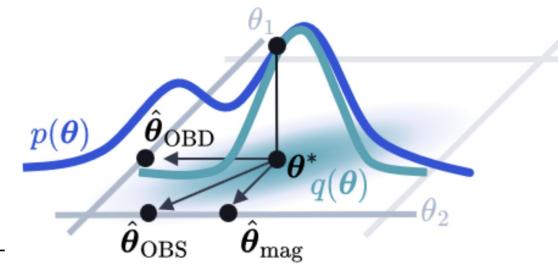
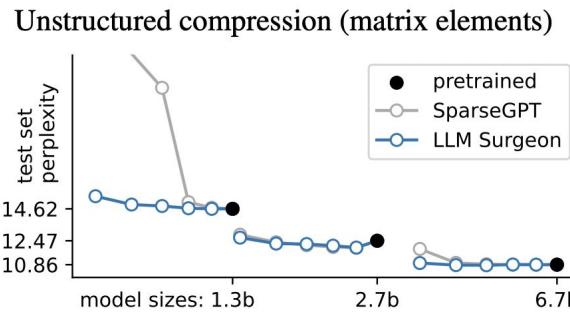
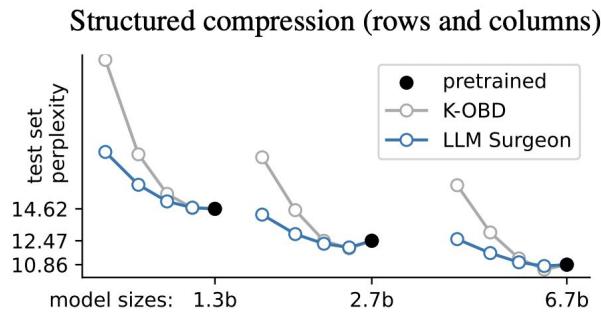
- Easy inverse / determinant (Kronecker identities in Laplace [1,2], Cholesky parameterization in VI as App. D in [13])
- Use mini-batches (easy in VI, requires care in Laplace [4])
- Account for weight-symmetries (with correction in KL term [5])
- Optimize prior variance and output variance empirically ([2], closed-form VI updates in App. D of [13])

# Laplace approximations at scale

# Laplace of networks with billions of parameters

Laplace can be applied to modern architectures [6]

State-of-the-art in structured LLM pruning [7]:



[6] Kronecker-Factored Approximate Curvature for Modern Neural Network Architectures

Runa Eschenhagen, Alexander Immer, Richard E Turner, Frank Schneider, Philipp Hennig

In NeurIPS 2023

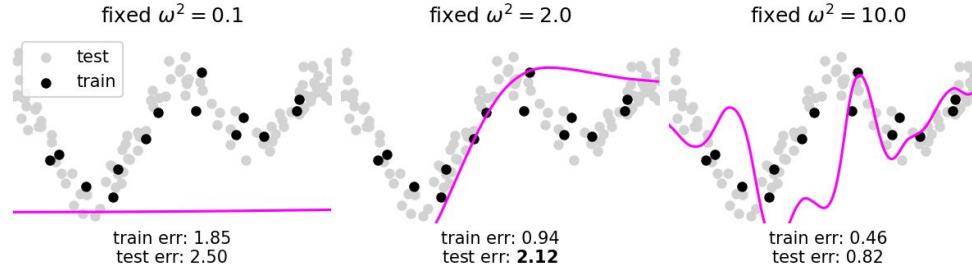
[7] The LLM Surgeon

Tycho F.A. van der Ouderaa, Markus Nagel, Mart van Baalen, Yuki M. Asano, Tijmen Blankevoort

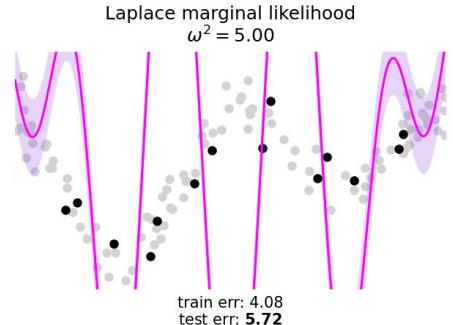
In ICLR 2024

# Bayesian model selection

## Regular training (maximum likelihood)



## Bayesian model selection (marginal likelihood)

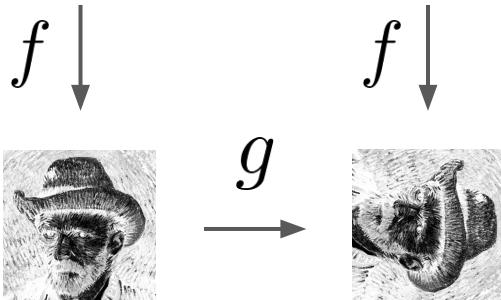
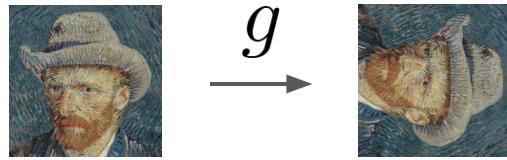


$$\text{NN}_{\theta}(x) = \mathbf{W}_3 \text{relu}(\mathbf{W}_2 \cos(\omega^2 \mathbf{W}_1 x))$$

# Inductive bias in deep learning

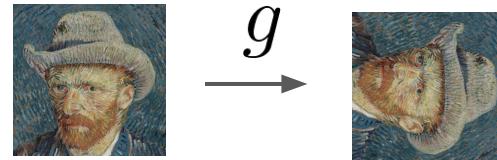
# Equivariance

$$f(g \cdot x) = g \cdot f(x)$$



# Invariance

$$f(g \cdot x) = f(x)$$

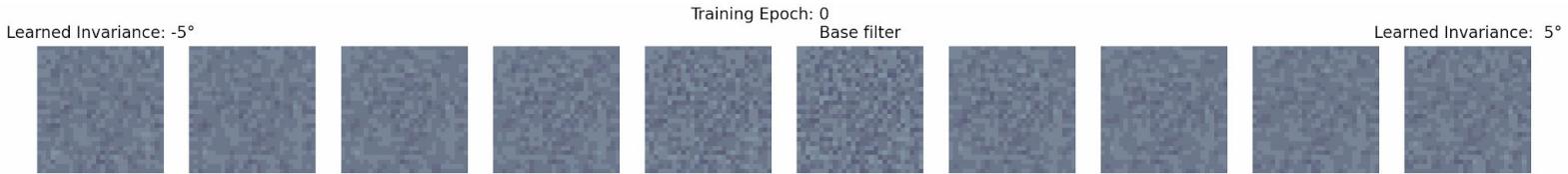


Example image:

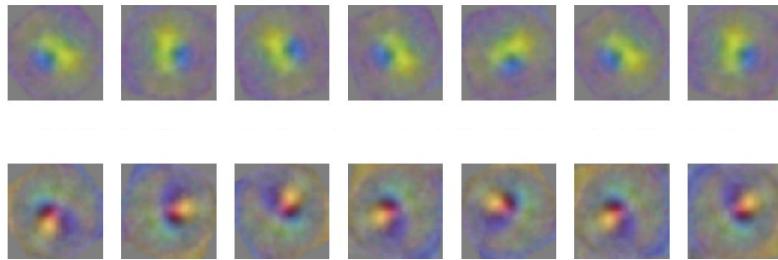
Vincent van Gogh, Self-Portrait with Grey Felt Hat, 1887  
Van Gogh Museum Amsterdam

Can we learn symmetries?

## Learned filter bank on rotated MNIST:



## Learned filter bank on rotated CIFAR-10:



arXiv:2202.1349v1 [stat.ML] 25 Feb 2022

**Learning Invariant Weights in Neural Networks**

Tycho F.A. van der Ouderaa<sup>1</sup> Mark van der Wilk<sup>1</sup>  
<sup>1</sup>Imperial College London, UK

**Abstract**

Assumptions about invariances or symmetries in data are often made to produce more parsimonious models. Many common models use it as machine learning as constraint to project correctly onto a lower dimensional space. This paper experiments in convolutional neural networks, and introduces a new way to learn invariance in neural networks. It has been shown that marginal likelihood often provides a better way to learn invariance than the maximum Likelihood Process. We propose a weight-space regularizer that encourages weights to be invariant in the marginal likelihood to both invariance in scale and rotation resulting in naturally higher performing models.

**2 RELATED WORK**

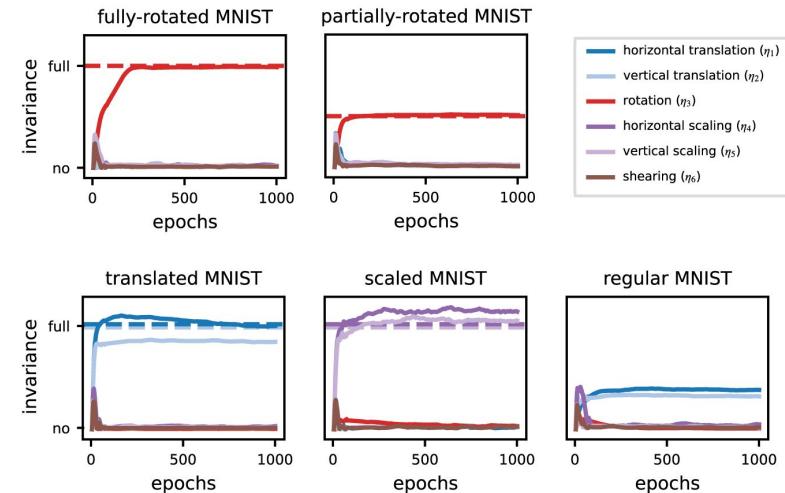
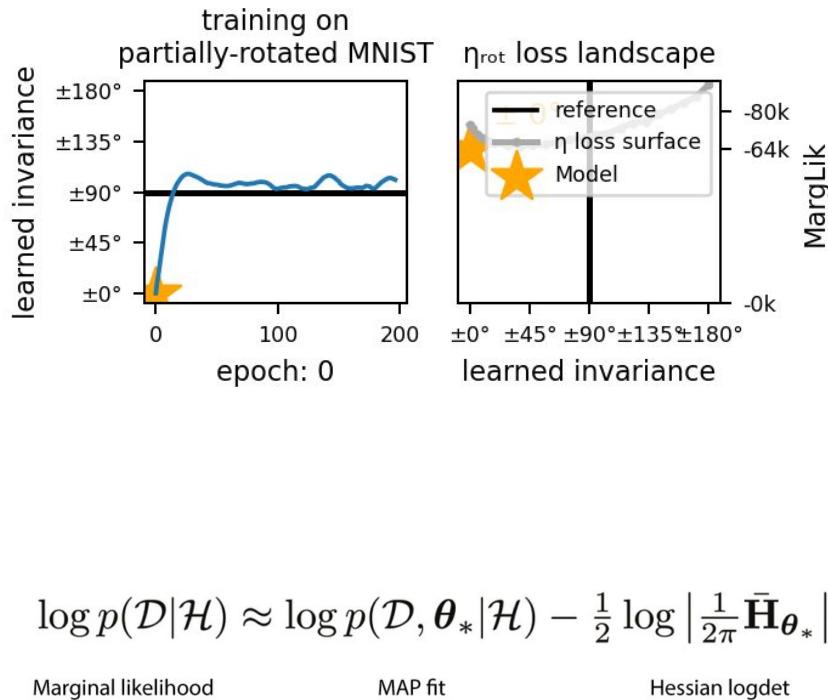
Convolutional neural networks (CNNs) have been successfully applied to many computer vision tasks since the rise in the waves of Deep Learning (LeCun et al., 2015). It is commonly understood that the translational invariance that arises from the stride of the convolutional layers is key for its outstanding performance on many tasks.

In contrast to the success of CNNs, learning invariance in linear models has been less successful. One reason is that linear models are typically not able to learn invariance in the same way as CNNs. Many studies have proposed ways to incorporate other types of invariance in linear models. For example, invariance in the weights of neural networks (Welling et al., 2014; Matus et al., 2017; Ermon et al., 2017; Weller and Welling, 2016) and invariance in the input data (Huang et al., 2017; Gómez et al., 2017) have been proposed. These efforts allow practical operation in neural networks for learning invariance. However, these methods for learning invariance are typically hard and must be known and specified by hand.

Some studies have proposed to learn invariance through the use of Gaussian Processes (GP) (Rasmussen and Williams, 2006). By this do we embed symmetries in the weights of a neural network. For example, the work of van der Wilk et al. (2021) does parameter sharing of layer weights but requires a meta-learning procedure that also relies on a validation loss. Some of the work in Bayesian

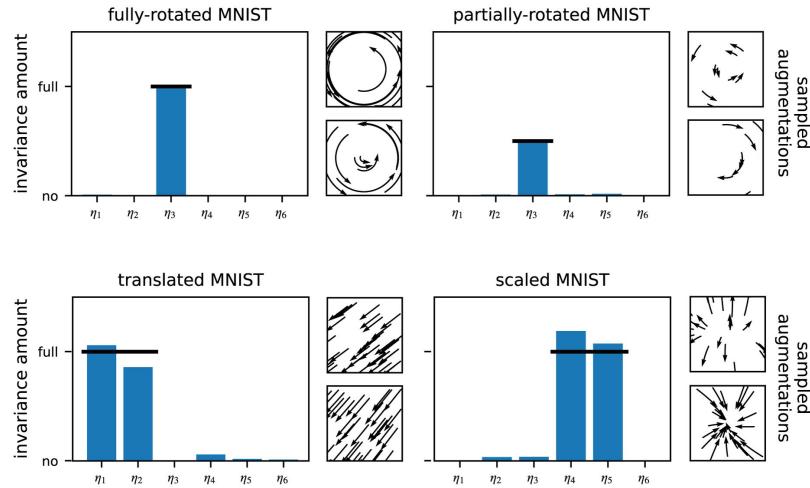
[8] **Learning Invariant Weights in Neural Networks**  
Tycho F.A. van der Ouderaa and Mark van der Wilk  
In ICML 2021 (UDL Workshop) and UAI 2022 (Oral)

# Invariance Learning in Deep Nets



[9] Invariance Learning in Deep Neural Networks with Differentiable Laplace Approximations  
Alexander Immer\*, Tycho F.A. van der Ouderaa\*,  
Vincent Fortuin, Gunnar Rätsch and Mark van der Wilk  
In NeurIPS 2022

# Invariance Learning in Deep Nets



**Invariance Learning in Deep Neural Networks  
with Differentiable Laplace Approximations**

Alexander Immer<sup>1,2</sup> Tycho F.A. van der Ouderaa<sup>3\*</sup> Vincent Fortuin<sup>1</sup> Gunnar Rätsch<sup>1,2</sup> Mark van der Wilk<sup>3</sup>

---

**Abstract**

Data augmentation is commonly applied to improve performance of deep learning by increasing the knowledge that invariant features can be in the input and the output. Currently, the common data augmentation is chosen by human effort and costly cross-validation, which makes it cumbersome to do so. We propose a more efficient and convenient gradient-based method for selecting the data augmentation. Our approach relies on phrasing data augmentation as an invariance in the prior distribution of the learned neural network model selection, which has been shown to work in Gaussian processes, but not yet for deep neural networks. We compare our method to a Kullback–Leibler–Fisher Laplace approximation to the marginal likelihood as our objective, which can be optimised without human supervision or validation data. The proposed gradient-based method can recover invariances present in the data, and that this improves generalisation on image datasets.

**1. Introduction**

Data augmentation is a commonly used machine learning technique that is essential to high-performing deep learning and computer vision systems. It is a process that adds noise to a set of distributions of transformations, by fitting a model with inputs that are transformed in a way that is known to leave the output class unchanged. This procedure can be regarded as artificially creating more data, which is vital for training deep neural networks. Yet, choosing the right transformation is an expensive and task-specific process that relies on domain knowledge and human effort. In addition, it is often necessary to validate this choice. This quickly becomes cumbersome when many parameters are considered, particularly if they are continuous.

<sup>1</sup>Equal contribution, order decided by coin flip. <sup>2</sup>Department of Computer Science, ETH Zurich, Switzerland. <sup>3</sup>Max Planck Institute for Intelligent Systems, University of Tübingen, Germany, and Department of Computing, Imperial College London, UK. Correspondence to: Alexander Immer (alexander.immer@inf.ethz.ch), Tycho F.A. van der Ouderaa (tycho.vanderouderaa@imperial.ac.uk).

arXiv:2202.10638v1 [stat.ML] 22 Feb 2022

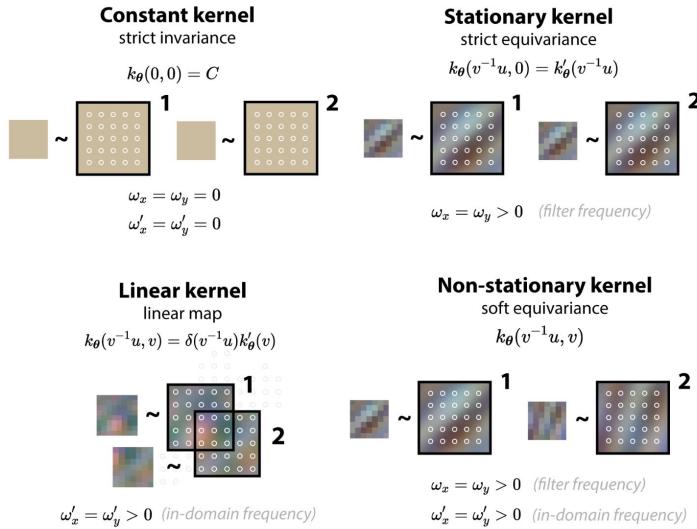
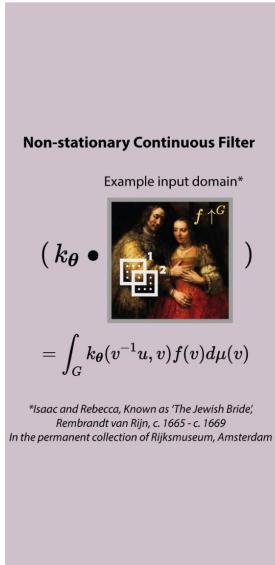
Figure 1: Optimising the marginal invariance of a neural network by minimising the Laplace approximation to the log marginal likelihood. The method recovers the true invariance with which the data was generated, during training, without supervision.

Figure 1: Optimising the marginal invariance of a neural network by minimising the Laplace approximation to the log marginal likelihood. The method recovers the true invariance with which the data was generated, during training, without supervision.

[9] **Invariance Learning in Deep Neural Networks with Differentiable Laplace Approximations**  
 Alexander Immer\*, Tycho F.A. van der Ouderaa\*,  
 Vincent Fortuin, Gunnar Rätsch and Mark van der Wilk  
 In NeurIPS 2022

What about learning equivariances?

# Relaxing equivariance constraints



---

Relaxing Equivariance Constraints with  
Non-stationary Continuous Filters

---

Tycho F.A. van der Ouderaa  
Imperial College London  
United Kingdom

David W. Romero  
Vrije Universiteit Amsterdam  
The Netherlands

Mark van der Wilk  
Imperial College London  
United Kingdom

---

**Abstract**

Equivariances can provide useful inductive biases in neural network modeling, with translation equivariance of convolutional neural networks as canonical example. Equivariances can be embedded into architectures through weight-sharing and many recent theoretical results show that this can be done in a representation-invariant way. The type of symmetry is typically fixed and must be chosen in advance. Although some tasks are inherently equivariant, many tasks do not strictly obey such symmetries in which case equivariant models can be overly restrictive. In this work, we propose a parameter-efficient relaxation of equivariance that has soft-equivariant linear map, equivariance and strict-equivariance as special case, between which can be interpolated in a parameter-efficient manner. The proposed parameterization can be easily integrated into existing layers to learn equivariant symmetry structure in neural networks. Compared to non-equivariant or strict-equivariant baselines, we experimentally verify that soft equivariance can lead to improved performance in terms of test accuracy on CIFAR-10 and CIFAR-100 image classification tasks.

---

**1 Introduction**

Symmetries, such as equivariances and invariances, can be embedded into neural network architectures and provide an inductive bias that leads to better generalization, data-efficiency and higher overall performance. Convolutional layers are known to provide translational equivariance in simple euclidean spaces, and many recent works have extended the construction to more complex groups and manifolds. When choosing correctly, equivariant layers can bring large gains in model performance. Symmetries are typically fixed, must be specified in advance, and can not be adjusted.

Symmetries that are embedded in the network architecture enforce a hard constraint on the functions the network can represent. This is used when a problem is inherently symmetric, but can become problematic when the problem does not obey such symmetry constraints. For example, 1's and 2's do not encode absolute positional information due to translation equivariance, and '0's and '9's become indistinguishable under strict rotational invariance. Relaxed symmetry constraints can mitigate such potential problems by allowing the network to learn equivariant representations on-the-fly.

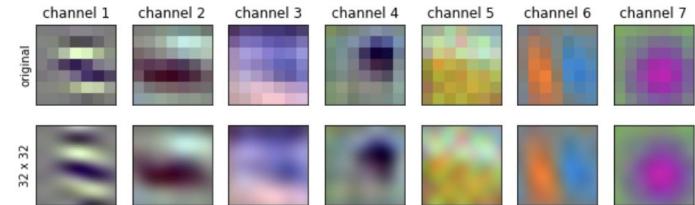
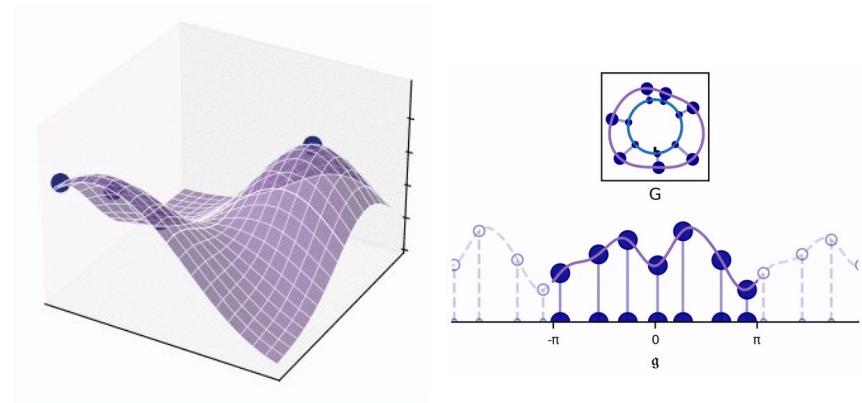
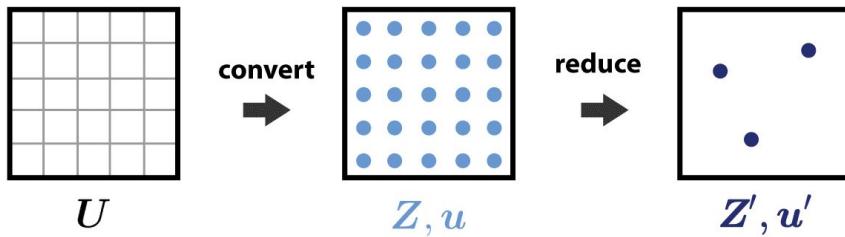
There have been attempts to automatically learn symmetry structures in neural networks from data. Existing approaches commonly rely on parameterizations that relax symmetry constraints, but often focus on invariance, which are easier to parameterize than equivariances. Allowing parameterizations that are equivariant in their own right can make a useful contribution to research in such symmetry discovery methods to enable layer-by-layer learning of equivariant symmetry structures.

In this work, we propose to relax equivariance constraints by generalizing the convolution operator with non-stationary filters that also depend on the absolute input or output group element. Applied to deep learning, this results in a novel parameter-efficient layer able to interpolate between (i) non-equivariant dot products akin to a fully-connected layer, (ii) strict group equivariant convolutions and (iii) strict group invariant constraints (Fig. 1).

Preprint. Under review.

[10] Relaxing Equivariance Constraints with Non-stationary Continuous Filters  
Tycho F.A. van der Ouderaa, David Romero and Mark van der Wilk  
In NeurIPS 2022

# Continuous kernels



## [11] Sparse Convolutions on Lie Groups

Tycho F.A. van der Ouderaa and Mark van der Wilk

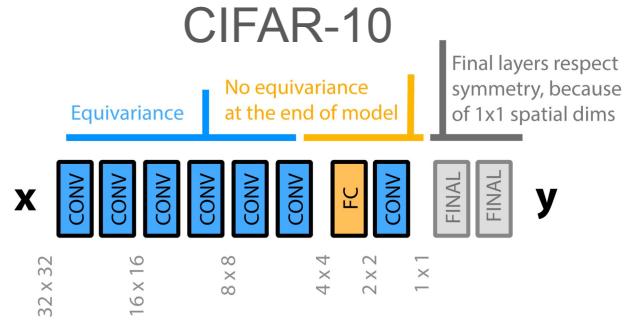
In NeurIPS 2022 (Workshop on Symmetry and Geometry in Neural Representations)

# Learning layer-wise equivariances

MNIST

Example:  
  
(3, upper left)

Symmetry	Prediction task	MAP		Diff. Laplace		
		FC 112.7 M	CONV 0.4 M	FC 112.7 M	CONV 0.4 M	Learned (ours) 1.8 M
Strict symmetry	(digit)	95.73	99.38	97.04	99.23	98.86
Partial symmetry	(digit, quadrant)	95.10	24.68	95.46	24.50	99.00



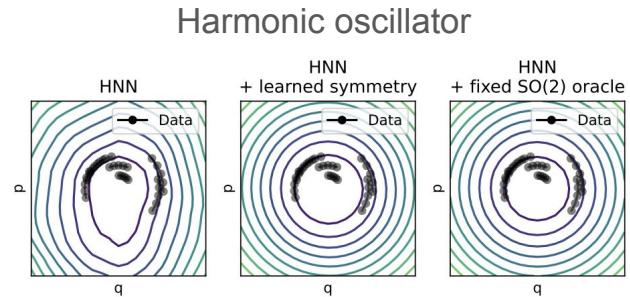
[12] Learning Layer-wise Equivariances Automatically using Gradients

Tycho F.A. van der Ouderaa, Alexander Immer, Mark van der Wilk  
In NeurIPS 2023 (Awarded with spotlight)

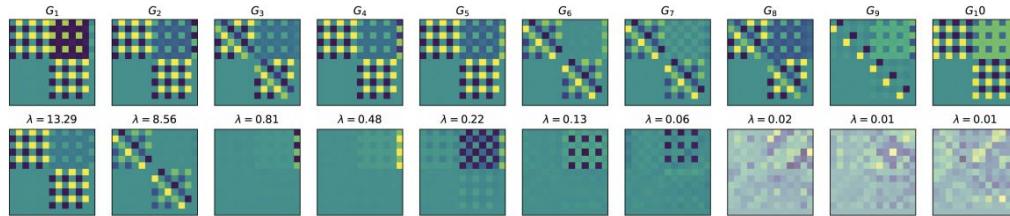
# Symmetries in Dynamical Systems

# Noether's razor

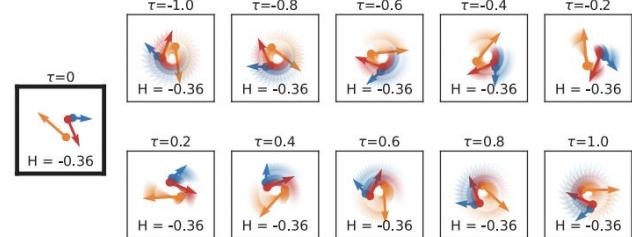
## Automatic symmetry discovery in dynamical systems



Generators associated to learned conserved quantities and their singular value decomposition.



Example of rotational symmetry associated to generator  $G_9$



Learned dynamics: <b>2d 3-body system</b>	Train data				Test data		Test data (moved)		Test data (wider)	
	Train MSE	NLL/N	KL/N	-ELBO/N ( $\downarrow$ )	Test MSE ( $\downarrow$ )					
HNN	0.0028	-13.87	13.34	-9.52	0.0016	0.0035	0.0016	0.0016	0.0016	
HNN + learned symmetry (ours)	0.0017	-20.09	7.28	<b>-12.81</b>	<b>0.0006</b>	<b>0.0004</b>	<b>0.0006</b>	<b>0.0006</b>	<b>0.0006</b>	
HNN + fixed SE(2) (reference oracle)	0.0019	-19.27	7.96	<b>-11.32</b>	<b>0.0006</b>	<b>0.0006</b>	<b>0.0006</b>	<b>0.0006</b>	<b>0.0006</b>	

[13] Noether's razor

Tycho F.A. van der Ouderaa, Mark van der Wilk, Pim de Haan  
In submission.

# Bayesian Neural Model Selection for Symmetry Learning



Tycho van der Ouderaa

 @tychovdo

 tychovdo@gmail.com

Supervised by Dr. Mark van der Wilk

Postgraduate student (PhD) at:



Visiting researcher at:

