

Learning Invariant Weights in Neural Networks

Oral presentation @ UAI 2022

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

Imperial College London

Open pdf in Adobe Acrobat to enable animations.

Invariances in Deep Learning

Symmetries in Deep Learning

Embedding symmetries into architectures leads to better models!



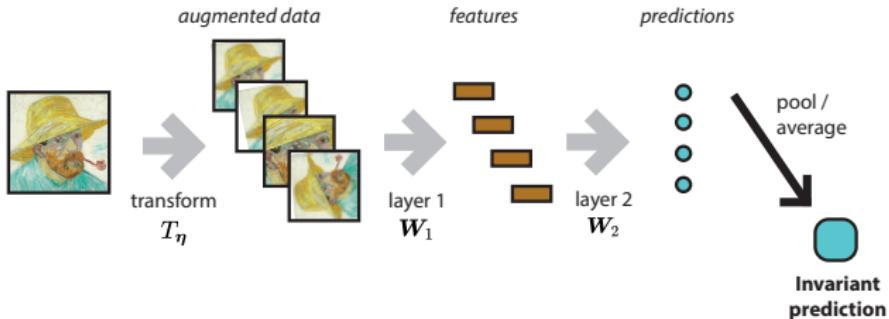
(a) Convolutions embed translation equivariance.

(b) Can be extended to other groups, such as rotation.

Symmetries need to be chosen or selected with cross-validation.

Can we *learn* the right invariances with gradients?

Constructing an invariant neural network



Consider a shallow network

$$g_{\theta}(T(x)) = \sigma(W_2 \circ \phi(W_1 \circ T \circ x))$$

integrated over a set of transformations

$$f_{\theta}(x; \eta) = \int g_{\theta}(T(x)) p_{\eta}(T) dT$$

We can apply transformations to the input or the weights first:

$$W_1 \circ (T \circ x) = (W_1 \circ T) \circ x$$

Parameterizing invariance

General parameterization of invariance

$$T = \exp \left(\sum_i \epsilon_i \eta_i \mathbf{G}_i \right), \quad \epsilon \sim U[-1, 1]^k$$

use set of affine generators

\mathbf{G}_1 : translation x

\mathbf{G}_2 : translation y

\mathbf{G}_3 : rotation

\mathbf{G}_4 : scale x

\mathbf{G}_5 : scale y

\mathbf{G}_6 : shear

Stochastic or deterministic sampling

The invariant predictor is described by an integral, which is intractable.

However, we can predict with an approximation from MC samples:

$$\hat{f}_{\theta}(x; \eta) = \frac{1}{S} \sum_{i=1}^S g_{\theta}(T_i(x))$$

to get an unbiased estimate

$$f_{\theta}(x; \eta) = \mathbb{E}_T [\hat{f}_{\theta}(x; \eta)]$$

Normally, we would use cross-validation to learn hyperparameters.

Finding optimal hyper-parameters η with Bayesian model selection

$$p(\theta, \eta | \mathcal{D}) = \frac{p(\mathcal{D}|\theta, \eta)p(\theta|\eta)p(\eta)}{p(\mathcal{D}|\eta)} \quad (\text{Full Bayes})$$

$$\hat{\eta} = \arg \max_{\eta} \int p(\mathcal{D}|\theta, \eta)p(\theta|\eta)d\theta \quad (\text{Empirical Bayes})$$

Well known procedure (Empirical Bayes, Type-II ML).

Used to learning invariances in GPs [van der Wilk; 2018].

Variational Inference

Marginal likelihood is intractable for neural networks.

We derive a lower bound using multi-sample Jensen's inequality:

$$\log p(\mathcal{D}) \geq \mathbb{E}_{q(\boldsymbol{\theta})} \left[\mathbb{E}_{\prod_{i=1}^N p_{\boldsymbol{\eta}}(T_i)} \left[\log p(y | \hat{f}_{\boldsymbol{\eta}}(\mathbf{x}, \boldsymbol{\eta})) \right] \right] - \text{KL}(q(\boldsymbol{\theta} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) || p(\boldsymbol{\theta}))$$

Similar to Nabarro et al., 2022 and Schwöbel et al., 2022.

Results

Learning affine invariances



Figure 1: Visualisation of filter banks trained on rotated CIFAR-10 data.



(a) Trained on regular MNIST.

(b) Trained on rotated MNIST.



(c) Trained on scaled MNIST.

(d) Trained on translated MNIST.

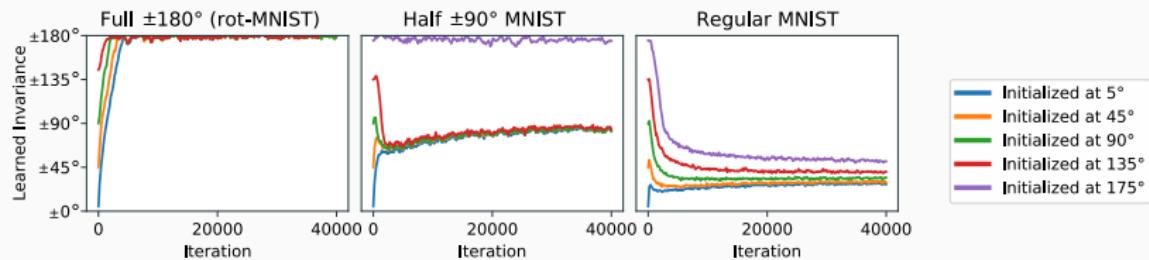
Figure 2: The same model capable of learning affine invariances learns filter banks with different invariances corresponding the data it was trained on.

Learned invariant filter banks

Recovering invariances

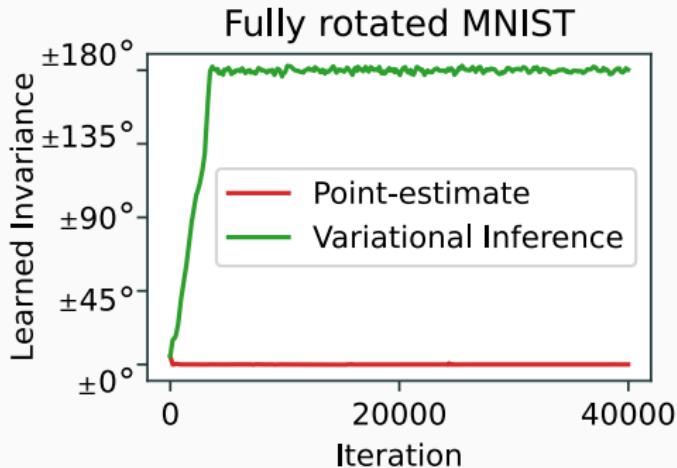
Relaxed invariance is not limited by the 'closure' axiom of groups.

We can learn the right amount of invariance from training data.



The necessity of the Bayesian approach

The marginal likelihood balances data fit and model complexity.



This proves to be very useful beyond predictive uncertainty.

Conclusion

Conclusion

We can learn invariant weights from data!

But, only in shallow networks

Exciting work in progress...

- Scaling the objective to deep networks

"Invariance Learning in Deep Neural Networks with Differentiable Laplace Approximations",

A Immer, TFA van der Ouderaa, V Fortuin, G Rätsch, M van der Wilk (2022)

- Parameterization of learnable layer-wise equivariance

"Relaxing Equivariance Constraints with Non-stationary Continuous Filters",

TFA van der Ouderaa, M van der Wilk (2022)

Very happy to engage in discussions on the topic!

Follow on  @tychovdo

Come chat at our poster ID# 419