





### Is mean-field a bad variational approximation?

- People often assume that the commonly used mean-field approximation (independent weights) is bad in variational inference for Bayesian Neural Networks.
- We challenge this. Instead, in **deeper models the** mean-field approximation is fine.

fine in deep models

"diagonal approximation is no good because of the weak strong posterior correlations in the parameters." at least one mode of – David MacKay, 1992

This challenges a foundational assumption that motivates much research (e.g., Louizos and Welling 2016; Sun et al. 2017; Oh et al. 2019).

## **Implications: Solve Difficulties of** Scaling MFVI. Don't Be Fancy.

- In bigger networks, MFVI gets better. So we must resolve problems scaling MFVI. E.g., gradient variance, computational cost.
- We may not need complex posterior approximations.
- Architecture matters for Bayesian approximations. Can't evaluate methods meant for deep models in shallow models.

the University of Oxford.

# Liberty or Depth: Deep Bayesian Neural Nets Do Not **Need Complex Weight Posterior Approximations**

Sebastian Farquhar, Lewis Smith, Yarin Gal {first.last@cs.ox.ac.uk}

### **Empirically: Mean-field Assumption Is Better In Deep Models**



# **Theoretically: Depth 'Simulates' Correlation and True Posterior Has At** Least One 'Approximately Mean-field' Mode

We introduce a new tool: 'local product matrices' for analysis piecewise linear models via linear ones. In linear case, we can 'flatten' a deep model with matrix multiplication and look at correlations in the 'product matrix'. This lets us analyse the **function** output distribution induced by weights. Diagonal weights induce complex covariance in product.

 $\mathbf{o} = \begin{pmatrix} W_{11}^{(2)} & W_{12}^{(2)} \\ W_{21}^{(2)} & W_{22}^{(2)} \end{pmatrix} \mathbf{x}$  $= \begin{pmatrix} B_{11}A_{11} + B_{12}A_{21} & B_{11}A_{12} + B_{12}A_{22} \\ B_{21}A_{11} + B_{22}A_{21} & B_{21}A_{12} + B_{22}A_{22} \end{pmatrix} \mathbf{x}$ 

**Prop 1:** 3+ layers gives off-diagonal covariance anywhere in product matrix. **Prop 2:** MVG/K-FAC is special case of 3 meanfield layers.

Fig 4. Visualizing covariance of product matrix of K meanfield Gaussian layers trained on FashionMNIST. Starts diagonal. Develops covariance.





(a) 1-layer.

This research was supported by The Alan Turing Institute and the EPSRC grant number EP/P00881X/1 at the Centre for Doctoral Training in Autonomous Intelligent Machines and Systems and Cyber Security at



	Method	Covariance	Acc.	NLL	ECE
	VOGN <sup>‡</sup>	Diagonal	67.4%	1.37	0.029
	Noisy K-FAC <sup>††</sup>	MVG	66.4%	1.44	0.080
et	SWAG -Diag <sup>†</sup>	Diagonal	78.6%	0.86	0.046
et	SWAG <sup>†</sup>	Low-rank	78.6%	0.83	0.020
	SWAG -Diag <sup>†</sup>	Diagonal	80.0%	0.86	0.057
	$SWAG^{\dagger}$	Low-rank	79.1%	0.82	0.028

Fig 3. On Imagenet, in deep models there is no clear advantage to complex covariance approximations over mean-field (diagonal).

(b) 5-layers.



(c) 10-layers.

Prop 3: Extends Prop 1 to 'local product matrix' for piecewise non-linearities like Leaky ReLU. **Prop 4:** A mean-field network can approximate true posterior density arbitrarily closely.