# Kronecker-factored Approximate Curvature for Physics-Informed Neural Networks

Felix Dangel\*, Johannes Müller\*, Marius Zeinhofer\*

NeurIPS 2024









We develop KFAC for, and apply it to, loss functions with differential operators.



#### We develop KFAC for, and apply it to, loss functions with differential operators.

- Goal: Learn PDE solution
  - $egin{aligned} \mathcal{L} m{u}(m{x}) &= m{f}(m{x}) & m{x} \in \Omega \ m{u}(m{x}) &= m{g}(m{x}) & m{x} \in \partial \Omega \end{aligned}$
- + Ansatz: Neural net  $u_{\theta}(\mathbf{x})$



#### We develop KFAC for, and apply it to, loss functions with differential operators.

Goal: Learn PDE solution

$$\mathcal{L}u(\mathbf{x}) = f(\mathbf{x}) \qquad \mathbf{x} \in \Omega$$
  
 $u(\mathbf{x}) = g(\mathbf{x}) \qquad \mathbf{x} \in \partial \Omega$ 

- + Ansatz: Neural net  $u_{\theta}(\mathbf{x})$
- + Sample  $\boldsymbol{x}_n \sim \Omega$ ,  $\boldsymbol{x}_n^{\mathrm{b}} \sim \partial \Omega$

$$\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \coloneqq \underbrace{\frac{1}{2N_{\Omega}} \sum_{n=1}^{N_{\Omega}} (\mathcal{L}\boldsymbol{u}_{\boldsymbol{\theta}}(\boldsymbol{x}_{n}) - f(\boldsymbol{x}_{n}))^{2}}_{L_{\Omega}(\boldsymbol{\theta})} + \underbrace{\frac{1}{2N_{\partial\Omega}} \sum_{n=1}^{N_{\partial\Omega}} \left(\boldsymbol{u}_{\boldsymbol{\theta}}(\boldsymbol{x}_{n}^{b}) - g(\boldsymbol{x}_{n}^{b})\right)^{2}}_{L_{\partial\Omega}(\boldsymbol{\theta})}}_{L_{\partial\Omega}(\boldsymbol{\theta})}$$



### We develop KFAC for, and apply it to, loss functions with differential operators.

Goal: Learn PDE solution +

 $\mathcal{L}u(\mathbf{x}) = f(\mathbf{x})$  $\mathbf{X} \in \Omega$  $u(\mathbf{x}) = q(\mathbf{x})$  $oldsymbol{x}\in\partial\Omega$ 

- + Ansatz: Neural net  $u_{\theta}(\mathbf{x})$
- + Sample  $\boldsymbol{x}_n \sim \Omega$ ,  $\boldsymbol{x}_n^{\rm b} \sim \partial \Omega$

#### Hard to train with popular methods 2d Poisson (D = 257)





### We develop KFAC for, and apply it to, loss functions with differential operators.

Goal: Learn PDF solution +

 $\mathcal{L}u(\mathbf{x}) = f(\mathbf{x})$  $oldsymbol{x}\in\Omega$  $u(\mathbf{x}) = q(\mathbf{x})$  $\mathbf{X}\in\partial\Omega$ 

 $L_{\Omega}(\boldsymbol{\theta})$ 

- + Ansatz: Neural net  $u_{\theta}(\mathbf{x})$
- + Sample  $\boldsymbol{x}_n \sim \Omega$ ,  $\boldsymbol{x}_n^{\rm b} \sim \partial \Omega$



 $L_{\partial\Omega}(\boldsymbol{\theta})$ 



### We develop KFAC for, and apply it to, loss functions with differential operators.

Goal: Learn PDF solution

 $\mathcal{L}u(\mathbf{x}) = f(\mathbf{x})$  $\mathbf{X} \in \Omega$  $u(\mathbf{x}) = q(\mathbf{x})$  $\mathbf{x} \in \partial \Omega$ 

- + Ansatz: Neural net  $u_{\theta}(\mathbf{x})$
- + Sample  $\mathbf{x}_n \sim \Omega$ ,  $\mathbf{x}_n^b \sim \partial \Omega$

#### but proposed methods don't scale 2d Poisson (D = 9873) $10^{-1}$ SGD $10^{-3}$





### We develop KFAC for, and apply it to, loss functions with differential operators.

Goal: Learn PDE solution

 $egin{aligned} \mathcal{L} m{u}(m{x}) &= m{f}(m{x}) & m{x} \in \Omega \ m{u}(m{x}) &= m{g}(m{x}) & m{x} \in \partial \Omega \end{aligned}$ 

- + Ansatz: Neural net  $u_{\theta}(\mathbf{x})$
- + Sample  $\boldsymbol{x}_n \sim \Omega$ ,  $\boldsymbol{x}_n^{\mathrm{b}} \sim \partial \Omega$



 $10^{1}$ 

 $10^{2}$ 

 $10^{3}$ 

KFAC\*

$$\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) \coloneqq \underbrace{\frac{1}{2N_{\Omega}} \sum_{n=1}^{N_{\Omega}} (\mathcal{L}\boldsymbol{u}_{\boldsymbol{\theta}}(\boldsymbol{x}_{n}) - f(\boldsymbol{x}_{n}))^{2}}_{L_{\Omega}(\boldsymbol{\theta})} + \underbrace{\frac{1}{2N_{\partial\Omega}} \sum_{n=1}^{N_{\partial\Omega}} \left( \boldsymbol{u}_{\boldsymbol{\theta}}(\boldsymbol{x}_{n}^{\mathrm{b}}) - g(\boldsymbol{x}_{n}^{\mathrm{b}}) \right)^{2}}_{L_{\partial\Omega}(\boldsymbol{\theta})}}_{L_{\partial\Omega}(\boldsymbol{\theta})}$$

 $10^{-5}$ 



+ Neural network ansatz  $u_{m{ heta}} = f^{(L)} \circ f^{(L-1)} \circ \ldots \circ f^{(1)}$ 



- + Neural network ansatz  $u_{m{ heta}} = f^{(L)} \circ f^{(L-1)} \circ \ldots \circ f^{(1)}$
- + Evaluating  $u_{\theta}(\mathbf{x})$

$$oldsymbol{z}^{(0)}\mapstooldsymbol{z}^{(1)}\mapstooldsymbol{z}^{(2)}\mapsto\cdots\mapstooldsymbol{z}^{(L)}=oldsymbol{u}_{oldsymbol{ heta}}(oldsymbol{x})$$

with vector-valued  $\mathbf{z}^{(l)} = f^{(l)}(\mathbf{z}^{(l-1)})$  and  $\mathbf{z}^{(0)} = \mathbf{x}$ 



- + Neural network ansatz  $u_{m{ heta}} = f^{(L)} \circ f^{(L-1)} \circ \ldots \circ f^{(1)}$
- + Evaluating  $u_{\theta}(\mathbf{x})$

$$\boldsymbol{z}^{(0)}\mapsto \boldsymbol{z}^{(1)}\mapsto \boldsymbol{z}^{(2)}\mapsto \dots\mapsto \boldsymbol{z}^{(L)}=\boldsymbol{u}_{\boldsymbol{ heta}}(\boldsymbol{x})$$

with vector-valued  $\boldsymbol{z}^{(l)} = f^{(l)}(\boldsymbol{z}^{(l-1)})$  and  $\boldsymbol{z}^{(0)} = \boldsymbol{x}$ 

+ Evaluating  $\mathcal{L}u_{\theta}(\mathbf{x})$ 

$$\mathbf{Z}^{(0)} \mapsto \mathbf{Z}^{(1)} \mapsto \mathbf{Z}^{(2)} \mapsto \ldots \mapsto \mathbf{Z}^{(L)} \mapsto \mathcal{L}u_{\theta}(\mathbf{x})$$

with **matrix-valued**  $Z^{(l)}$  and dependencies determined by Taylor-mode

## High-level Walkthrough



 $\begin{array}{l} \textbf{Computing differential operators yields networks with linear weight sharing layers.} \\ & \rightarrow \textbf{Apply KFAC for linear layers with weight sharing [Eschenhagen et al., 2023, NeurIPS]} \end{array}$ 

- + Neural network ansatz  $u_{m{ heta}} = f^{(L)} \circ f^{(L-1)} \circ \ldots \circ f^{(1)}$
- + Evaluating  $u_{\theta}(\mathbf{x})$  Linear layer:  $\mathbf{z}^{(l)} = \mathbf{W}\mathbf{z}^{(l-1)}$

$$\boldsymbol{z}^{(0)}\mapsto \boldsymbol{z}^{(1)}\mapsto \boldsymbol{z}^{(2)}\mapsto \dots\mapsto \boldsymbol{z}^{(L)}=\boldsymbol{u}_{\boldsymbol{ heta}}(\boldsymbol{x})$$

with vector-valued  $\boldsymbol{z}^{(l)} = f^{(l)}(\boldsymbol{z}^{(l-1)})$  and  $\boldsymbol{z}^{(0)} = \boldsymbol{x}$ 

+ Evaluating  $\mathcal{L}u_{\theta}(\mathbf{x})$  Linear layer:  $\mathbf{Z}^{(l)} = \mathbf{W}\mathbf{Z}^{(l-1)}$ 

$$\mathbf{Z}^{(0)} \mapsto \mathbf{Z}^{(1)} \mapsto \mathbf{Z}^{(2)} \mapsto \ldots \mapsto \mathbf{Z}^{(L)} \mapsto \mathcal{L}u_{\theta}(\mathbf{x})$$

with **matrix-valued**  $Z^{(l)}$  and dependencies determined by Taylor-mode



- $L(\boldsymbol{ heta}) = L_{\Omega}(\boldsymbol{ heta}) + L_{\partial\Omega}(\boldsymbol{ heta})$
- + **Goal:** Approximate Gauss-Newton matrix

$${f G}({m heta})={f G}_{\Omega}({m heta})+{f G}_{\partial\Omega}({m heta})$$

$$\mathbf{V}$$

 $L(oldsymbol{ heta}) = L_\Omega(oldsymbol{ heta}) + L_{\partial\Omega}(oldsymbol{ heta})$ 

+ **Goal:** Approximate Gauss-Newton matrix

$${f G}({m heta})={f G}_{\Omega}({m heta})+{f G}_{\partial\Omega}({m heta})$$

#### + KFAC

 $G_{\Omega}(W) \approx A_{\Omega} \otimes B_{\Omega}$ (generalized to diff. ops)  $G_{\partial\Omega}(W) \approx A_{\partial\Omega} \otimes B_{\partial\Omega}$ 

[Martens and Grosse, 2015, ICML]

## Algorithm, Experiments & Conclusion

# $\mathbf{V}$

100d Poisson.  $D \approx 10^6$ 

 $10^{4}$ 

100

 $10^{-1}$ 

 $10^{-2}$ 

Our KFAC-based optimizer outperforms popular methods, scales to large neural networks, and performs similarly to Hessian-free [Martens, 2010, ICML].

4+1d heat.  $D \approx 10^5$ 

 $10^{(}$ 

 $10^{-2}$ 

 $10^{-4}$ 

- $L(oldsymbol{ heta}) = L_{\Omega}(oldsymbol{ heta}) + L_{\partial\Omega}(oldsymbol{ heta})$
- + **Goal:** Approximate Gauss-Newton matrix

$$oldsymbol{\mathsf{G}}(oldsymbol{ heta}) = oldsymbol{\mathsf{G}}_\Omega(oldsymbol{ heta}) + oldsymbol{\mathsf{G}}_{\partial\Omega}(oldsymbol{ heta})$$

+ KFAC

 $G_{\Omega}(W) \approx A_{\Omega} \otimes B_{\Omega}$ (generalized to diff. ops)  $G_{\partial\Omega}(W) \approx A_{\partial\Omega} \otimes B_{\partial\Omega}$ [Martens and Grosse, 2015, ICML]





- Runa Eschenhagen, Alexander Immer, Richard E. Turner, Frank Schneider, and Philipp Hennig. Kronecker-Factored Approximate Curvature for Modern Neural Network Architectures. In Advances in Neural Information Processing Systems (NeurIPS), 2023.
- James Martens. Deep learning via Hessian-free optimization. In International Conference on Machine Learning (ICML), 2010.
- James Martens and Roger Grosse. Optimizing Neural Networks with Kronecker-factored Approximate Curvature. In International Conference on Machine Learning (ICML), 2015.
- Johannes Müller and Marius Zeinhofer. Achieving high accuracy with PINNs via energy natural gradient descent. In **International Conference on Machine Learning (ICML)**, 2023.