

# Backpropagation Beyond the Gradient

## Dissertation

der Mathematisch-Naturwissenschaftlichen Fakultät  
der Eberhard Karls Universität Tübingen  
zur Erlangung des Grades eines  
Doktors der Naturwissenschaften  
(Dr. rer. nat.)

vorgelegt von  
**Felix Julius Dangel, M. Sc.**  
aus Stuttgart

Tübingen  
2022

Gedruckt mit Genehmigung der Mathematisch-Naturwissenschaftlichen Fakultät der Eberhard Karls Universität Tübingen.

Tag der mündlichen Qualifikation: 31. Januar 2023

Dekan: Prof. Dr. Thilo Stehle  
1. Berichterstatter: Prof. Dr. Philipp Hennig  
2. Berichterstatter: Prof. Dr. Hendrik Lensch  
3. Berichterstatterin: Prof. Dr. Asja Fischer

**Disclaimer:** this thesis uses Federico Marotta's kaobook template based on Ken Arroyo Ohori's doctoral thesis. I am grateful to Frank Schneider who added many stylistic features and would like to thank him, Katharina Ott, Frederik Kunstner, Julia Grosse, Jonathan Wenger, Matthias Werner, Marvin Pförtner, Runa Eschenhagen, and Nathanel Bosch for their help to improve the text.



# Acknowledgments

I want to express my deepest gratitude to Philipp Hennig for his guidance that allowed me to seamlessly change fields for the Ph.D. I am inspired by your passion for research, and grateful for your advice and everything I have learned from you. Thank you for allowing me to travel the world, for spontaneously stepping in on poster presentations, and for taking care of rebellious coffee machines.

Also, I am grateful to Asja Fischer and Hendrik Lensch who invested their time and effort into reviewing this thesis, as well as Zeynep Akata and Georg Martius for serving on my thesis advisory committee and forming my examination committee.

I want to thank Franziska Weiler (Methods of Machine Learning (MoML) group), Leila Masri & Sara Sorce (International Max Planck Research School for Intelligent Systems) for their effective administrative support.

This work was made possible by great collaborations. I want to thank Stefan Harmeling, Frederik Kunstner, Frank Schneider, and Lukas Tatzel for these experiences and their patience with me. I very much appreciated our regular coffee breaks, coding sessions, “grammar police” discussions, and “random questions of today”. This also applies more broadly to the entire MoML group that provided a productive and friendly work environment. I consider myself lucky to have met all of you (apologies for the regular mess in the office kitchen). I also enjoyed supervising Christian Meier, Paul Fischer, Shrisha Bharadwaj, Tim Schäfer, Jessica Bader, and Elisabeth Knigge. Thank you for contributing to this work; hopefully you learned as much as I did.

During my time in Tübingen, I was fortunate enough to be surrounded by amazing people. I want to express my gratitude to Jack Brady, Nicolò Zottino, Ori Press, Emilia Magnani, Laura Iacovissi, Georg Hafner, Massimiliano Mancini, Ilya Petrov, Faezeh Zakeri, Wafaa Mohammed, Martina Contisciani, Alessandro Lonardi, Nicolò Ruggeri, Diego Theuerkauf, Ahmad Gazar, and everyone I forgot to list here. You made this a really great time.

Thank you to Jan Koschke, Marcel Bauer, and Jens Hollmann for being great friends, and your support during—and long before—my studies. Please do not ever call me Doctor D. (you know what I mean...).

I am grateful to my family, especially my parents, my sister, and my grandparents (happy birthday, Oma!). They sparked my curiosity about science, trusted in me, and provided me with everything I could hope for throughout the years. Eu também mando um abraço forte à minha segunda família brasileira. Dona Cícera, Senhor João Bosco, Gabriela, fico muito feliz por ter vocês na minha vida.

Finally, I cannot thank Daniela Leite enough for her support during this journey, and for being such a wonderful person. Thank you for showing me the world from perspectives I would never have known otherwise. I hope the future holds many more adventures and moments we share.

Thank you!

*Felix Dangel*  
Tübingen, August 31, 2022



# Abstract

Automatic differentiation is a key enabler of deep learning: previously, practitioners were limited to models for which they could manually compute derivatives. Now, they can create sophisticated models with almost no restrictions and train them using first-order, *i.e.* gradient, information. Popular libraries like PyTorch [126] and TensorFlow [1] compute this gradient efficiently, automatically, and conveniently with a single line of code. Under the hood, reverse-mode automatic differentiation, or gradient backpropagation, powers the gradient computation in these libraries. Their entire design centers around gradient backpropagation.

These frameworks are specialized around one specific task—computing the average gradient in a mini-batch. This specialization often complicates the extraction of other information like higher-order statistical moments of the gradient, or higher-order derivatives like the Hessian. It limits practitioners and researchers to methods that rely on the gradient. Arguably, this hampers the field from exploring the potential of higher-order information and there is evidence that focusing solely on the gradient has not lead to significant recent advances in deep learning optimization [145].

To advance algorithmic research and inspire novel ideas, information beyond the batch-averaged gradient must be made available at the same level of computational efficiency, automation, and convenience.

This thesis presents approaches to simplify experimentation with rich information beyond the gradient by making it more readily accessible. We present an implementation of these ideas as an extension to the backpropagation procedure in PyTorch. Using this newly accessible information, we demonstrate possible use cases by (i) showing how it can inform our understanding of neural network training by building a diagnostic tool, and (ii) enabling novel methods to efficiently compute and approximate curvature information.

First, we extend gradient backpropagation for sequential feedforward models to Hessian backpropagation which enables computing approximate per-layer curvature. This perspective unifies recently proposed block-diagonal curvature approximations. Like gradient backpropagation, the computation of these second-order derivatives is modular, and therefore simple to automate and extend to new operations.

Based on the insight that rich information beyond the gradient can be computed efficiently and at the same time, we extend the backpropagation in PyTorch with the BackPACK library. It provides efficient and convenient access to statistical moments of the gradient and approximate curvature information, often at a small overhead compared to computing just the gradient.

Next, we showcase the utility of such information to better understand neural network training. We build the Cockpit library that visualizes what is happening inside the model during training through various instruments that rely on BackPACK’s statistics. We show how Cockpit provides a meaningful statistical summary report to the deep learning engineer to identify bugs in their machine learning pipeline, guide hyperparameter tuning, and study deep learning phenomena.

Finally, we use BackPACK’s extended automatic differentiation functionality to develop ViViT, an approach to efficiently compute curvature information, in particular curvature noise. It uses the low-rank structure of the generalized Gauss-Newton approximation to the Hessian and addresses shortcomings in existing curvature approximations. Through monitoring curvature noise, we demonstrate how ViViT’s information helps in understanding challenges to make second-order optimization methods work in practice.

This work develops new tools to experiment more easily with higher-order information in complex deep learning models. These tools have impacted works on Bayesian applications with Laplace approximations [40], out-of-distribution generalization [64, 130], differential privacy [179], and the design of automatic differentiation systems. They constitute one important step towards developing and establishing more efficient deep learning algorithms.



# Zusammenfassung

Automatisches Differenzieren stellt eine wesentliche Komponente für Deep Learning dar: Zuvor waren Anwender auf Modelle beschränkt, deren Ableitungen sich manual berechnen ließen. Jetzt können sie komplexe Modelle mit fast beliebiger Struktur entwerfen und diese mit Gradienteninformation trainieren. Software-Bibliotheken wie PyTorch [126] und TensorFlow [1] berechnen den Gradienten effizient und automatisch in einer Codezeile. Im Hintergrund geschieht dies per automatischer Differenzierung im Rückwärtsmodus, genannt Backpropagation. Das Design dieser Bibliotheken basiert auf Backpropagation.

Solche Bibliotheken sind besonders auf eine spezielle Funktion – das Berechnen des gemittelten Gradienten über einen Mini-Batch – ausgerichtet. Diese Spezialisierung erschwert oft die Berechnung anderer Größen, wie höhere statistische Momente des Gradienten oder Ableitungen höherer Ordnung, etwa der Hesse-Matrix. Sie beschränkt Anwender und Forscher auf gradientenbasierte Methoden und hindert die Erforschung des Potenzials höherer Ableitungen und statistischer Momente. Es gibt Anzeichen, dass dieser Fokus auf den Gradienten keine signifikanten Fortschritte in der Optimierung neuronaler Netze erlaubt hat [145].

Um die algorithmische Forschung voranzutreiben und neue Ideen zu inspirieren müssen andere Größen, die über den gemittelten Gradienten eines Mini-Batches hinausgehen, genauso einfach – das heißt automatisiert und effizient – zugänglich gemacht werden.

Diese Arbeit stellt Ansätze vor, die das Experimentieren mit diversen Größen jenseits des Gradienten erleichtern, indem sie diese leichter zugänglich machen. Wir implementieren diese Ideen durch Erweiterung der bestehenden Backpropagation in PyTorch. Mit diesen neu zugänglichen Größen demonstrieren wir Anwendungsszenarien indem wir zeigen wie diese (i) anhand eines Diagnosetools zum besseren Verständnis des Trainings neuronaler Netze führen und (ii) neue Methoden für effiziente und approximative Berechnung von Krümmungsinformation ermöglichen.

Zunächst erweitern wir Gradient Backpropagation für sequenzielle neuronale Netze auf Backpropagation von Hesse-Matrizen, welche die Berechnung von Krümmung innerhalb einer Schicht ermöglicht. Diese Formulierung vereinheitlicht kürzlich vorgeschlagene Näherungsverfahren für Krümmung durch block-diagonale Matrizen. Wie Gradient Backpropagation ist dieses Verfahren für Ableitungen zweiter Ordnung zwischen Schichten entkoppelt und daher einfach zu automatisieren und zu erweitern.

Basierend auf der Erkenntnis, dass reichhaltige Information jenseits des Gradienten gleichzeitig mit diesem berechnet werden kann, erweitern wir die Standard-Backpropagation von PyTorch mit der BackPACK-Bibliothek. Diese bietet effizienten Zugriff auf statistische Gradientenmomente und approximative Krümmungsinformation. Im Vergleich zur Gradientenberechnung ist der Mehraufwand oft gering.

Danach demonstrieren wir den Nutzen solcher Information zum besseren Verständnis des Trainings neuronaler Netze. Wir entwickeln die Cockpit-Bibliothek, die während des Trainings anhand verschiedener Instrumente – basierend auf BackPACKs Größen – visualisiert, was innerhalb des Modells geschieht. Wir zeigen, wie Cockpit Ingenieuren hilft, Fehler in deren Pipeline zu identifizieren, Hyperparameter zu wählen und Deep Learning Phänomene zu untersuchen.

Zuletzt verwenden wir BackPACKs Funktionalität zur Entwicklung von ViViT, einem Verfahren zur effizienten Berechnung von Krümmungsinformation, insbesondere deren Rauschen. ViViT basiert auf dem äußeren Produkt in der verallgemeinerten Gauß-Newton-Matrix und behebt Probleme bestehender Methoden. Durch Beobachtung von Krümmungsrauschen zeigen wir, wie ViViTs Größen dabei helfen, Herausforderungen für die erfolgreiche Realisierung krümmungsbasierter Optimierer zu verstehen.

Diese Arbeit entwickelt neue Tools zum vereinfachten Experimentieren mit Information höherer Ordnung in komplizierten tiefen Netzen. Diese Tools haben Arbeiten zu bayesschen Anwendungen mit Laplace-Approximationen [40], Out-of-Distribution Generalisierung [64, 130], Differential Privacy [179], sowie das Design von Deep Learning Bibliotheken beeinflusst. Sie stellen einen wichtigen Schritt zur Entwicklung und Etablierung effizienterer Algorithmen für Deep Learning dar.



# Table of Contents

<b>Acknowledgments</b>	<b>v</b>
<b>Abstract</b>	<b>vii</b>
<b>Zusammenfassung</b>	<b>ix</b>
<b>Table of Contents</b>	<b>xi</b>
<b>Notation</b>	<b>xiii</b>
<b>1. Overview</b>	<b>1</b>
1.1. Introduction . . . . .	1
1.2. Outline . . . . .	3
 <b>I. BACKGROUND &amp; MOTIVATION</b>	<b>7</b>
<b>2. Deep Learning Components</b>	<b>9</b>
2.1. Empirical Risk Minimization . . . . .	9
2.2. Neural Networks . . . . .	18
2.3. Automatic Differentiation . . . . .	22
<b>3. Higher-order Information</b>	<b>31</b>
3.1. Popular Deep Learning Optimizers . . . . .	32
3.2. Second-order Optimization . . . . .	35
3.3. Other Applications . . . . .	43
 <b>II. BACKPROPAGATION BEYOND THE GRADIENT</b>	<b>49</b>
<b>4. Modular Block-diagonal Curvature Approximations for Feedforward Architectures</b>	<b>51</b>
4.1. Introduction . . . . .	51
4.2. Notation . . . . .	53
4.3. Modular Hessian Backpropagation . . . . .	53
4.4. Experiments & Implementation . . . . .	59
4.5. Conclusion . . . . .	62
<b>5. BackPACK: Packing More into Backprop</b>	<b>63</b>
5.1. Introduction . . . . .	63
5.2. Theory & Implementation . . . . .	65
5.3. Evaluation & Benchmarks . . . . .	70
5.4. Experiments . . . . .	70
5.5. Conclusion . . . . .	73
<b>6. Cockpit: A Practical Debugging Tool for the Training of Deep Neural Networks</b>	<b>75</b>
6.1. Introduction & Motivation . . . . .	75
6.2. Cockpit's Instruments . . . . .	77
6.3. Experiments . . . . .	81
6.4. Showcase . . . . .	84

6.5. Benchmark . . . . .	84
6.6. Conclusion . . . . .	86
<b>7. ViViT: Curvature Access Through the GGN’s Low-rank Structure</b>	<b>89</b>
7.1. Introduction & Motivation . . . . .	89
7.2. Notation & Method . . . . .	91
7.3. Experiments . . . . .	95
7.4. Related Work . . . . .	101
7.5. Use Cases . . . . .	101
7.6. Conclusion . . . . .	102
<b>III. CONCLUSION &amp; FUTURE DIRECTIONS</b>	<b>103</b>
<b>8. Conclusion &amp; Future Directions</b>	<b>105</b>
8.1. Summary & Impact . . . . .	107
8.2. Future Work . . . . .	108
<b>IV. APPENDIX</b>	<b>111</b>
<b>A. Additional Material for Chapter 4</b>	<b>113</b>
A.1. Matrix Derivatives . . . . .	113
A.2. HBP for FCNNs . . . . .	115
A.3. HBP for Loss Functions . . . . .	118
A.4. HBP for CNNs . . . . .	119
A.5. Experimental Details . . . . .	124
<b>B. Additional Material for Chapter 5</b>	<b>127</b>
B.1. BackPACK Extensions . . . . .	127
B.2. Benchmark Details . . . . .	134
B.3. Experimental Details . . . . .	135
B.4. BackPACK Cheat Sheet . . . . .	139
<b>C. Additional Material for Chapter 6</b>	<b>141</b>
C.1. Code Example . . . . .	141
C.2. Instrument Overview . . . . .	142
C.3. Mathematical Details . . . . .	143
C.4. Additional Experiments . . . . .	155
C.5. Implementation Details & Additional Benchmarks . . . . .	157
C.6. Cockpit View of Convex Stochastic Problems . . . . .	166
<b>D. Additional Material for Chapter 7</b>	<b>167</b>
D.1. Mathematical Details . . . . .	167
D.2. Experimental Details . . . . .	168
D.3. Implementation Details . . . . .	193
<b>Bibliography</b>	<b>197</b>

# Notation

The notation is influenced by Goodfellow et al. [60].

## Tensors, Matrices, Vectors, Numbers

$a$	A scalar
$\mathbf{a}$	A column vector
$A$	A matrix
$\mathbf{A}$	A tensor
$a_i$ or $[\mathbf{a}]_i$	The $i$ th entry of the vector $\mathbf{a}$
$A_{i,j}$ or $[A]_{i,j}$	The $(i, j)$ th entry of the matrix $A$ (row $i$ , column $j$ )
$[A]_{i,:}$ (or $[A]_{:,j}$ )	The $i$ th row (or $j$ th column) of the matrix $A$
$A_{i,j,k}$ or $[\mathbf{A}]_{i,j,k}$	The $(i, j, k)$ th entry of the tensor $\mathbf{A}$
$\text{vec}(A), \text{vec}(\mathbf{A})$	Matrix/tensor flattened into a vector; convention implies $\text{vec}(ABC) = (\mathbf{C}^\top \otimes A) \text{vec}(B)$
$\text{diag}(\mathbf{a})$	The square matrix with vector $\mathbf{a}$ on the diagonal and zeros elsewhere
$\text{diag}(A)$	The vector containing the diagonal elements of the matrix $A$
$\text{diag}(A_1, \dots, A_L)$	A block-diagonal matrix with diagonal blocks given by square matrices $A_1, \dots, A_L$
$\text{Tr}(A), \det(A)$	Trace and determinant of a matrix $A$
$\ \mathbf{a}\ _2$	$L_2$ norm of vector $\mathbf{a}$ , i.e. $\ \mathbf{a}\ _2^2 = \mathbf{a}^\top \mathbf{a}$
$\text{eig}(A) := \{(\lambda_k, \mathbf{e}_k)\}_k$	Eigendecomposition of the matrix $A$ , eigenpairs $(\lambda_k, \mathbf{e}_k)$ satisfy $A\mathbf{e}_k = \lambda_k \mathbf{e}_k$
$(\lambda_k(A), \mathbf{e}_k(A))$	$k$ th eigenpair (eigenvalue, eigenvector) of matrix $A$
$A \otimes B$	Kronecker product of two matrices, For two vectors $\mathbf{a}, \mathbf{b}$ , one has $\mathbf{a} \otimes \mathbf{b}^\top = \mathbf{a}\mathbf{b}^\top$
$\mathbf{A} \odot \mathbf{B}, A \odot B, \mathbf{a} \odot b$	Elementwise multiplication (Hadamard product) of two tensors, matrices, vectors
$\mathbf{A} \oslash \mathbf{B}, A \oslash B, \mathbf{a} \oslash b$	Elementwise division (Hadamard division) of two tensors, matrices, vectors
$\mathbf{A}^{\odot 2}, A^{\odot 2}, \mathbf{a}^{\odot 2}$	Elementwise square of a tensor, matrix, vector
$\mathbf{A}^{\odot 1/2}, A^{\odot 1/2}, \mathbf{a}^{\odot 1/2}$	Elementwise square root of a tensor, matrix, vector

## Empirical Risk Minimization

A datum is usually indicated by a subscript  $n$ .

$(x, y)$	Labeled datum with input features $x$ and target $y$
$(x_n, y_n)$	Datum $n$ from a dataset
$D$	Total number of parameters in a model
$\theta \in \Theta := \mathbb{R}^D$	Parameter vector of a model
$f := f_\theta(x)$	Prediction of a model $f_\theta$ for input features $x$
$\ell$ or $\ell(f, y)$	Loss function to compare prediction and target; convex in $f$
$\mathbb{D} := \{(x_n, y_n)\}_{n=1}^{ \mathbb{D} }$	A dataset containing instances of labeled data $(x_n, y_n)$ indexed by $n$
$\mathbb{B}$	A mini-batch $\mathbb{B} \subseteq \mathbb{D}$
$N$	Number of data in a mini-batch or a dataset, depending on the context
$f_n := f_\theta(x_n)$	Model prediction for datum $n$
$\ell_n$ or $\ell(f_n, y_n)$	Loss of datum $n$
$p_{\mathbb{D}}(x, y)$	Empirical distribution of a dataset $\mathbb{D}$
$\mathcal{L}_{\mathbb{D}}(\theta)$	Empirical risk implied by the empirical distribution of a dataset $\mathbb{D}$
$\mathcal{L}_{\mathbb{D}_{\text{train}}}(\theta), \mathcal{L}_{\mathbb{B}}(\theta), \text{etc.}$	Training loss, mini-batch loss, etc.

## Neural Networks

The layer number is indicated by parenthesized superscripts  $(l)$ .

$L$	Total number of layers
$d^{(l)}$	Number of parameters in layer $l$ ; total number of parameters is $D = \sum_{l=1}^L d^{(l)}$
$\theta^{(l)} \in \mathbb{R}^{d^{(l)}}$	Parameter vector of layer $l$ , potentially empty for parameter-free layers like activations
$h^{(l-1)}$	Number of (hidden) inputs fed into layer $l$
$M := h^{(0)}, C := h^{(L)}$	Input feature dimension, output dimension (number of classes for classification)
$z^{(l-1)} \in \mathbb{R}^{h^{(l-1)}}$	(Hidden) features fed into layer $l$ (output of layer $l - 1$ )
$x := z^{(0)}, f := z^{(L)}$	Input to the neural network, and its prediction for input $x$
$f_{\theta^{(l)}}^{(l)}$	Layer $l$ parameterized by $\theta^{(l)}$ , mapping input $z^{(l-1)}$ to output $z^{(l)}$
$f_{\theta} := f_{\theta^{(L)}}^{(L)} \circ \dots \circ f_{\theta^{(1)}}^{(1)}$	Sequential feedforward neural network parameterized by $\theta$ , maps input $x$ to output $f$
$\theta$	Parameter vector, concatenation of parameters over layers $\theta := (\theta^{(1)\top}, \dots, \theta^{(L)\top})^\top$

## Derivatives

$\nabla, J$  and  $\nabla^2$  denote the gradient, Jacobian, and Hessian, respectively.

$J_a b$	Jacobian matrix of a vector $b$ w.r.t. a vector $a$ , $[J_a b]_{i,j} = \partial[b]_i / \partial[a]_j$
$J_{\mathbf{A}} \mathbf{B}$	Generalized Jacobian matrix for tensor variables, $[J_{\mathbf{A}} \mathbf{B}]_{i,j} = \partial[\text{vec } \mathbf{B}]_i / \partial[\text{vec } \mathbf{A}]_j$
$\nabla_a b := (J_a b)^\top$	Gradient vector of a scalar $b$ w.r.t. a vector $a$ , $[\nabla_a b]_i = \partial b / \partial a_i$
$\nabla_a^2 b$	Hessian matrix of a scalar $b$ w.r.t. a vector $a$ , $[\nabla_a^2 b]_{i,j} = \partial^2 b / \partial[a]_i \partial[a]_j$ (symmetric)
$\nabla_a^2 b, \nabla_{\mathbf{A}}^2 \mathbf{B}$	Generalized Hessian matrix (in general not quadratic, hence not symmetric) of a vector $b$ w.r.t. a vector $a$ , or more general tensor variables
$g_n(\theta) := \nabla_{\theta} \ell_n(\theta)$	Gradient of the loss implied by sample $n$
$g_{\mathbb{D}}(\theta) := \nabla_{\theta} \mathcal{L}_{\mathbb{D}}(\theta)$	Gradient of the empirical risk implied by a dataset $\mathbb{D}$
$g_{\mathbb{B}}(\theta) := \nabla_{\theta} \mathcal{L}_{\mathbb{B}}(\theta)$	Mini-batch gradient
$H_n(\theta) := \nabla_{\theta}^2 \ell_n(\theta)$	Hessian of the loss implied by sample $n$
$H_{\mathbb{D}}(\theta) := \nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}}(\theta)$	Hessian of the empirical risk implied by a dataset $\mathbb{D}$
$H_{\mathbb{B}}(\theta) := \nabla_{\theta}^2 \mathcal{L}_{\mathbb{B}}(\theta)$	Mini-batch Hessian
$H^{(l)}(\theta^{(l)})$ or $H(\theta^{(l)})$	The block in the Hessian corresponding to layer $l$
$G_{\mathbb{D}}(\theta)$	Generalized Gauss-Newton matrix on a dataset $\mathbb{D}$
$G^{(l)}(\theta^{(l)})$ or $G(\theta^{(l)})$	The block in the generalized Gauss-Newton matrix corresponding to layer $l$

## Statistics

$\mathcal{U}(\{1, \dots, N\})$	Uniform distribution over $\{1, \dots, N\}$
$\mathcal{N}(x   \mu, \sigma^2)$	Uni-variate normal/Gaussian distribution of random variable $x$ , with mean $\mu$ , positive variance $\sigma^2$ , and density $\mathcal{N}(x   \mu, \sigma^2) = 1/\sigma\sqrt{2\pi} \exp[-1/2((x-\mu)/\sigma)^2]$
$\mathcal{N}(\mathbf{x}   \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Multi-variate normal/Gaussian distribution of random vector $\mathbf{x}$ with mean vector $\boldsymbol{\mu}$ , PSD covariance matrix $\boldsymbol{\Sigma}$ , and density $\mathcal{N}(\mathbf{x}   \boldsymbol{\mu}, \boldsymbol{\Sigma}) = 1/(\sqrt{2\pi \det \boldsymbol{\Sigma}}) \exp[-1/2(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})]$
$\text{Cat}(c   \mathbf{p})$	Multinomial/Categorical distribution with probabilities $\mathbf{p}$ for categories $c$

## Miscellaneous

$\log$	The natural logarithm (base e, i.e. $\log(e) = 1$ )
--------	---

$\text{onehot}(c)$	One-hot vector of class $c$ with $\text{onehot}(c) = \delta_{i,c}$
$\text{softmax}(\mathbf{a})$	Softmax probabilities of the logits $\mathbf{a}$ , $[\text{softmax}(\mathbf{a})]_c = \exp(a_c)/\sum_{i=1} \exp(a_i)$ .
$\delta_{i,j}, \delta(\mathbf{x} - \mathbf{a})$	Kronecker delta ( $\delta_{i,i} = 1$ and $\delta_{i,j \neq i} = 0$ ), Dirac delta distribution
$(\mathbb{X} \rightarrow \mathbb{Y})$	Signature of a function that maps between $\mathbb{X}$ and $\mathbb{Y}$
$\{\mathbf{x}_n\}$ or $\{\mathbf{x}_n\}_n$	A set/collection of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ over the index set implied by $n$
$\hat{\mathbf{e}}_i$	Unit vector in direction $i$ , <i>i.e.</i> $\hat{\mathbf{e}}_i = \text{onehot}(i)$
$\mathbf{1}_m$	An $m$ -dimensional vector containing ones everywhere
$\log(\mathbf{a}), \exp(\mathbf{a})$	Elementwise natural logarithm and exponential function of a vector
$m_{\theta_t}(\theta)$	Local approximation of the loss in around $\theta_t$

## Acronyms & Abbreviations

<i>E.g.</i> or <i>e.g.</i>	For example ( <i>exempli gratia</i> )
<i>Etc.</i> or <i>etc.</i>	And so on ( <i>et cetera</i> )
<i>I.e.</i> or <i>i.e.</i>	That is ( <i>id est</i> )
<i>I.i.d.</i> or <i>i.i.d.</i>	Independent and identically distributed
<i>W.r.t.</i> or <i>w.r.t.</i>	With respect to
AD	Automatic differentiation
API	Application Programming Interface
BDA	Block diagonal approximation
CG	Conjugate gradients
CNN	Convolutional neural network
CPU	Central processing unit
DNN	Deep neural network
DP	Differential privacy
FCNN	Fully-connected neural network
GGN	Generalized Gauss-Newton (matrix)
GN	Gauss-Newton (matrix)
GPU	Graphics processing unit
HBP	Hessian backpropagation
JMP	Jacobian-matrix product
JVP	Jacobian-vector product
KFAC	Kronecker-factored curvature
KFC	Kronecker factors for convolution
KFLR	Kronecker-factored low rank
KFRA	Kronecker-factored recursive approximation
MAP	Maximum a posteriori (estimation)
MC	Monte Carlo
MJP	Matrix-Jacobian product
ML	Machine learning
MLE	Maximum likelihood estimation
MLP	Multi-layer perceptron
NGD	Natural gradient descent
PCH	Positive-curvature Hessian
PD	Positive definite
PSD	Positive semi-definite
ResNet	Residual (neural) network
SNR	Signal-to-noise ratio
TPU	Tensor processing unit
VJP	Vector-Jacobian product



# Overview 1.

## 1.1 Introduction

Deep learning has achieved significant performance gains over traditional methods on various tasks like image classification [41, 68, 91, 153, 159, 181], image generation [35, 58, 84, 131], machine translation [7, 102, 172], and game play [11, 151, 152]. These applications are powered by machine learning (ML) frameworks [1, 32, 126, 165] that tremendously reduce the added complexity for practitioners of using highly-efficient implementations on hardware accelerators like GPUs [88] and TPUs [83]. The convenience introduced by these libraries makes deep learning more accessible and is one main factor for its popularity and success.

It is surprisingly easy to build, train, and deploy a model, despite deep learning being computationally extremely demanding. Typically, neural networks have millions [68, 91, 153, 159, 181], billions [24, 129], even trillions [49] of parameters, and are trained on large data sets [41, 98] that can only be processed in smaller batches. The parameters are adjusted during training, which is phrased as optimization of a performance measure called “the loss”.

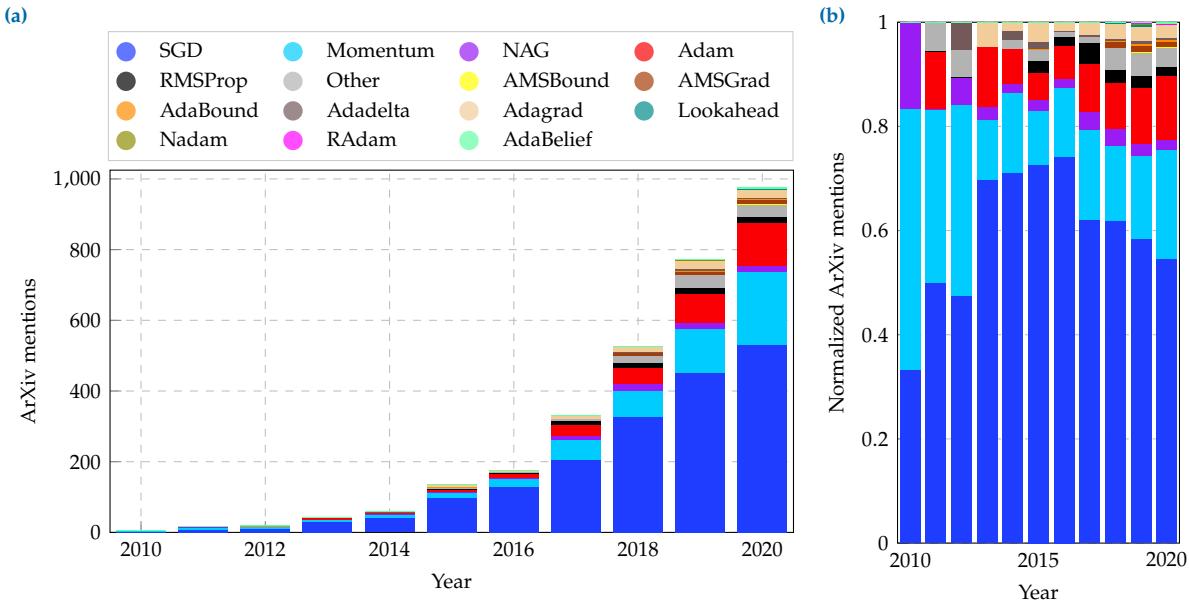
Training, *i.e.* optimization, proceeds in an iterative fashion: to update the model, a mini-batch of data is fed through the model (forward pass) to evaluate its current performance—the mini-batch loss. This loss is then differentiated *w.r.t.* the model’s parameters to obtain the mini-batch gradient (backward pass). Finally, an optimizer incorporates this gradient to update the model parameters. Popular frameworks like PyTorch [126] and TensorFlow [1] realize this procedure in code that looks similar to the following pseudocode:

```
1 dataset = ...      # Learning task examples
2
3 model = ...        # Practitioner's choice
4 loss_func = ...    # Practitioner's choice
5
6 optimizer = ...    # First-order method
7
8 while not_converged: # Standard training loop
9     features, targets = dataset.next_minibatch()
10
11     # Forward pass: Compute the loss
12     predictions = model(features)
13     loss = loss_func(predictions, targets)
14
15     # Backward pass: Compute the gradient
16     loss.backward()
17
18     # Update model parameters using the gradient
19     optimizer.step()
20     optimizer.zero_grad()
```

1.1 Introduction . . . . .	1
1.2 Outline . . . . .	3

**Procedure 1.1: Canonical deep learning training loop.** After setting up the data, model, loss function, and optimizer, iterate over batches: in each iteration, compute the mini-batch loss in a forward pass, and its gradient with a backward pass. Then use the gradient as learning signal to update the model parameters.

This framework frees practitioners from implementation details and lets them focus on, *e.g.*, specifying the neural network (Line 3) and performance criterion (Line 4). No matter how complicated a model



**Figure 1.1:** All popular deep learning optimizers are first-order methods that rely on the gradient. Reproduced from Schmidt et al. [145] with permission from the authors. (a) Number of ArXiv mentions in titles and abstracts of deep learning optimization methods, with other methods grouped into a single category. (b) The same plot with mentions normalized to the unit interval.

a practitioner may come up with: as long as it is differentiable, it will be compatible with the above training loop, and can be automatically differentiated and trained. This allows practitioners to focus on modeling aspects of their problem. Thanks to automatic differentiation built into ML libraries, they do not have to worry about low-level details of how to compute the learning signal. They obtain the gradient with a single function call that performs the backward pass (Line 16).

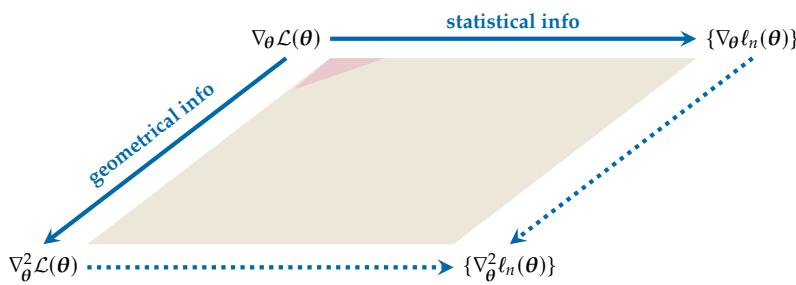
However, by abstracting these details, [Procedure 1.1](#) limits practitioners and researchers to gradient access only. This is not restrictive for the most popular optimizers to train DNNs as they only rely on the gradient ([Figure 1.1](#)). But the development of novel methods focuses mainly on inventing alternative update rules involving the gradient. Currently, there exist more than one hundred such methods; see e.g. [145] for an overview. This leads to the question whether one of these update rules performs significantly better than others. A broad comparison of the methods in [Figure 1.1](#), however, raises concerns that those methods do not seem to be specialized to problems or leverage different problem properties:

*Despite efforts by the community, there is currently no method that clearly dominates the competition. [...] tuning helps about as much as trying other optimizers.*

Schmidt et al. [145]

One reason for this may be that despite the diversity of update rules, they all rely on the same gradient information. Research and the development of novel methods should therefore use more than just the gradient.

Such information can be of statistical or geometrical nature ([Figure 1.2](#)). The loss/gradient is the empirical mean of the distribution of per-sample losses/gradients. Higher moments of the gradients, such as their variance, could enable more robust estimation of the mean. Higher-order derivatives, e.g. curvature in form of the Hessian, encode the loss



**Figure 1.2: Beyond the average gradient.** Higher-order statistical (noise) and geometrical (curvature) information provide a richer local approximation of the loss in deep learning, which relies mostly on the gradient highlighted in ●. The goal of this thesis is making higher-order information ● more readily available (figure inspired by [22]).  $\nabla_{\theta} \mathcal{L}$ ,  $\nabla_{\theta}^2 \mathcal{L}$  denote average gradient and Hessian;  $\{\nabla_{\theta} \ell_n\}$ ,  $\{\nabla_{\theta}^2 \ell_n\}$  are the distributions of per-sample gradients and Hessians.

landscape's geometry. These quantities are useful to build approximations of the loss that take its stochastic nature and curvature into account.

Optimization methods that use more than the gradient are the default in other applications (*e.g.* convex optimization, generalized linear models) and have also been considered in deep learning [5, 14, 53, 69, 106, 109, 178, 183]. However, they are often so complicated to implement that practitioners barely try them out: efficiently computing with higher-order information is challenging because its explicit representation is often too large to be stored in memory. Therefore, practical methods rely on implicit schemes, such as matrix-vector products, or light-weight structured approximations. Understanding and efficiently implementing them requires expert knowledge, which complicates experimentation.

The adoption of, and research on, algorithms depend on ease-of-use within frameworks. Existing software has closely evolved with the popularity of first-order methods that only demand efficient access to the gradient. This leads to the question on how to make higher-order information more readily available within such frameworks to expand the available toolkit. As backpropagation [137] forms the computational backbone of deep learning, this thesis describes approaches to extend the backpropagation algorithm, and answers the research questions:

- (Q1) *Which information beyond the gradient is efficiently accessible?*
- (Q2) *How to compute this information—conveniently, automatically, and efficiently—re-using the existing backpropagation implementation of ML frameworks?*
- (Q3) *How to use this information to advance gradient-based deep learning?*

## 1.2 Outline

This thesis contains three main parts: (I) an introduction to the setting, related concepts, and relevant notation along with motivation; (II) the achieved scientific contributions; and (III) a discussion of their impact and future directions. Throughout the text, additional details (examples, auxiliary calculations, *etc.*) are provided in margins. Readers that are familiar with the subject should feel encouraged to read the main text without interruptions of this additional material. An extensive appendix (IV) complements the main parts. It collects lengthier mathematical discussions and details about experiments and implementations, along with additional empirical results.

Part I provides background material and motivates the goal of this work. Chapter 2 introduces the components of supervised deep learning

through gradient-based empirical risk minimization (Procedure 1.1). The goal is to highlight how each component introduces structure in the empirical risk that can be leveraged for convenient and efficient computation. This structure is helpful to identify new classes of interesting quantities and understand the computational pipeline that will be extended in this work. Specifically, we inspect the following components:

- ▶ **The loss:** Section 2.1 presents supervised learning through empirical risk minimization and outlines connections of loss-based learning to approximating the unknown data distribution through samples via maximum likelihood estimation (MLE) or maximum a posteriori (MAP) estimation. Further, it highlights the finite-sum structure in the loss. This allows for massive speed-ups on parallel hardware accelerators like GPUs [88] when evaluating the loss, and stochastic approximation via sub-sampling with mini-batches, which introduces noise into the computed quantities.
- ▶ **The model:** Section 2.2 outlines structure in the optimization space given by the parameter space of a deep neural network. DNNs are usually highly over-parameterized, containing millions [68, 91, 153, 159, 181], billions [24, 129], even trillions [49] of parameters. Yet, they are built from relatively simple functional units—layers, or modules—glued together through function composition. We give a selected overview of layers, and present sequential feedforward neural networks, which are the model class this work focuses on.
- ▶ **The gradient:** Section 2.3 outlines the importance of automatic differentiation (AD) for ML. It starts with fundamentals on how to automate the chain rule for functions represented by a computation graph that tracks the dependencies between input, output, and intermediate variables. One important property is the modularity of AD, which allows for an elegant and extensible implementation: because the chain rule relates the derivative of a function composition to the derivatives of its composites, only these composites need to implement derivatives. New operations that cannot be composed from primitives can be added by specifying their derivatives. The discussion concludes by presenting gradient backpropagation [137], the most prominent evaluation scheme to compute gradients of the loss in deep learning which will be extended in this work.

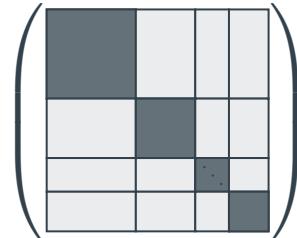
Chapter 3 motivates why focusing only on the gradient might be problematic, and motivates the importance of higher-order information, *e.g.* in form of noise (stochasticity) and curvature, to further enrich deep learning algorithms. First, Section 3.1 exemplifies the gradient’s dominance in optimization for deep learning where state-of-the-art methods are first-order methods. Next, Section 3.2 introduces relevant higher-order information for second-order optimization methods that rely on curvature. It discusses popular curvature matrices, such as the Hessian, generalized Gauss-Newton (GGN), Fisher information, and gradient covariance matrix, as well as their connections. Section 3.3 further motivates the utility of such information through concrete examples.

This concludes the background material. We identified and appreciated the strengths of deep learning libraries, their potential shortcomings that might complicate advances in the field, and found higher-order information as a promising direction for their extension.

**Part II** tackles this manuscript's goal to *extend the functionality of machine learning libraries to efficiently, automatically, and conveniently provide access to higher-order information*. It describes doing so by extending the gradient backpropagation algorithm of ML libraries to go beyond the gradient.

**Chapter 4** presents Hessian backpropagation (HBP), an extension of gradient backpropagation to compute layer-wise curvature in sequential feedforward neural networks. Just like gradient backpropagation recovers the gradient vectors in blocks that correspond to layers, local curvature, *i.e.* second-order partial derivatives of the loss *w.r.t.* parameters in a layer can be evaluated by backpropagating Hessians. Its computation is disentangled to the modular level, which allows for an elegant and extensible implementation in analogy to gradient backpropagation. We describe the backpropagation operations at the per-layer level, resulting in an algorithm that computes local curvature in an automated fashion, and at the same time as the gradient. Adaptations of the exact procedure recover positive semi-definite block diagonal approximations (BDAs)—*e.g.* of the GGN—and recently proposed Kronecker-factored curvature approximations [21, 31, 109] of the Hessian, unifying their computation.

**Chapter 4:** HBP: block-diagonal curvature via Hessian backpropagation.



[github.com/f-dangel/hbp](https://github.com/f-dangel/hbp)

**Disclaimer 1.1** **Chapter 4** is based on the peer-reviewed conference paper with the following co-author contributions:

F. Dangel, S. Harmeling, and P. Hennig. "Modular Block-diagonal Curvature Approximations for Feedforward Architectures". *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2020 [37]

	Ideas	Experiments	Analysis	Writing
<b>F. Dangel</b>	70 %	80 %	70 %	65 %
S. Harmeling	10 %	10 %	10 %	10 %
P. Hennig	20 %	10 %	20 %	25 %

**Chapter 5** further generalizes this idea and presents BackPACK, an efficient framework that extends the gradient backpropagation of PyTorch. The library provides access to higher-order statistical information about the gradient distribution, like individual gradients or an estimate of their variance, and structured curvature information, like the Hessian/GGN diagonal and block-diagonal Kronecker-factored curvature approximations. This is achieved with slightly more flexible implementations of AD functionality and by backpropagating additional information through the graph. Importantly, most quantities add small overhead to the gradient, making their exploration for research more attractive.

**Chapter 5:** BackPACK: an efficient framework built on top of PyTorch that extends the backpropagation algorithm.



[github.com/f-dangel/backpack](https://github.com/f-dangel/backpack)

**Disclaimer 1.2** **Chapter 5** is based on the peer-reviewed conference publication with the following co-author contributions:

F. Dangel, F. Kunstner, and P. Hennig. "BackPACK: Packing more into Backprop". *International Conference on Learning Representations (ICLR)*. 2020 [38]

	Ideas	Experiments	Analysis	Writing
<b>F. Dangel</b>	33 %	55 %	45 %	35 %
F. Kunstner	33 %	45 %	45 %	45 %
P. Hennig	33 %	0 %	10 %	20 %

**Chapter 6:** Cockpit: a debugging tool for the training of deep neural networks.



[github.com/f-dangel/cockpit](https://github.com/f-dangel/cockpit)

**Chapter 6** introduces Cockpit, a live-monitoring debugging tool that consists of various instruments which leverage higher-order information, efficiently provided by BackPACK. By providing a deeper look into the inner workings of neural networks through the lens of this information, Cockpit can help identify bugs in the ML pipeline, while keeping the computational overhead acceptable. This demonstrates the utility of higher-order information to assist deep learning practitioners and researchers to better understand their problems and conduct research.

**Disclaimer 1.3** **Chapter 6** is based on the peer-reviewed conference paper with the following co-author contributions:

F. Schneider, F. Dangel, and P. Hennig. "Cockpit: A Practical Debugging Tool for Training Deep Neural Networks". *Advances in Neural Information Processing Systems (NeurIPS)*. 2021 [147]

	Ideas	Experiments	Analysis	Writing
F. Schneider	45 %	40 %	40 %	45 %
<b>F. Dangel</b>	40 %	50 %	40 %	40 %
P. Hennig	15 %	10 %	20 %	15 %

**Chapter 7:** ViViT: efficient computation with the GGN's low-rank structure.



[github.com/f-dangel/vivit](https://github.com/f-dangel/vivit)

**Chapter 7** presents ViViT, a method that leverages the low-rank structure in the GGN to efficiently extract eigenvalues, eigenvectors, per-sample first- and second-order directional derivatives, and Newton steps. In contrast to other popular curvature approximations, ViViT is capable of tracking off-diagonal curvature blocks, offers principled approximations to trade off cost and accuracy, and allows studying noise in the curvature. Under the hood, ViViT's quantities are efficiently computed during backpropagation, building on BackPACK's advanced AD functionality. This demonstrates how such functionality enables investigations of unexplored structure in higher-order information which may be used to develop novel algorithms and better understand challenges to make them work in practice, *e.g.* noise.

**Disclaimer 1.4** **Chapter 7** is based on the peer-reviewed journal paper with the following co-author contributions:

F. Dangel, L. Tatzel, and P. Hennig. "ViViT: Curvature Access Through The Generalized Gauss-Newton's Low-Rank Structure". *Transactions on Machine Learning Research (TMLR)* (2022) [39]

	Ideas	Experiments	Analysis	Writing
<b>F. Dangel</b>	50 %	40 %	40 %	40 %
L. Tatzel	35 %	50 %	40 %	45 %
P. Hennig	15 %	10 %	20 %	15 %

**Part III** summarizes the findings of this manuscript *w.r.t.* the initial questions from [Page 3](#), as well as their impact and relation to the latest developments in the field. Finally, the manuscript identifies future research directions for the presented extended AD eco-system around PyTorch, and more broadly for the development of future ML libraries.

# **Part I.**

## Background & Motivation



# 2.

## Deep Learning Components

Broadly speaking, ML seeks to find a well-performing algorithm for a given task from “experience”. This thesis considers supervised deep learning, where “experience” is given through annotated examples in form of a dataset. A datum  $(x, y)$  consists of *input features*  $x$  and *targets*, or *labels*,  $y$ . The algorithm, or *model*, is a deep neural network ([Section 2.2](#)) that tries to predict  $y$  from  $x$  and is selected by minimizing a performance criterion on the available data (the empirical risk, [Section 2.1](#)) using optimization methods which rely on automatically computed derivatives ([Section 2.3](#)). This chapter reviews these components, highlighting their structure *w.r.t.* implementation in ML libraries. For a broader introduction to deep learning, see *e.g.* [60].

<a href="#">2.1 Empirical Risk Minimization</a>	9
<a href="#">2.2 Neural Networks</a>	18
<a href="#">2.3 Automatic Differentiation</a>	22

### 2.1 Empirical Risk Minimization

“Learning” is connected to optimization by the risk minimization paradigm. The idea is to define a performance metric that assesses the quality of the model’s prediction. Then, learning happens by maximizing that performance metric. Conversely, one can specify a metric for the prediction’s error (also referred to as *risk*), and minimize the latter.

For example, an intuitive way to assess performance for a classification task is *accuracy*, the ratio of correct and total predictions (the error would be the ratio of incorrect and total predictions). However, such direct performance measures are hard to optimize with derivative-based methods<sup>1</sup> and must be substituted by a surrogate function that approximates the original performance measure, but is easier to optimize. In the deep learning terminology, these surrogates are commonly referred to as *loss functions*. This section expands on risk minimization, its characteristic properties in deep learning, and its probabilistic interpretation.

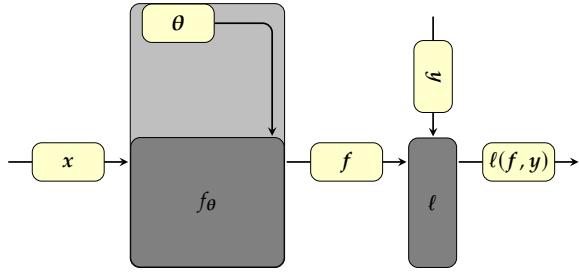
1: *e.g.* the accuracy on one datum is either 0 or 1, and hence its derivative vanishes everywhere it is defined.

#### 2.1.1 Notation & Mathematical Details

##### Risk & Empirical Risk

Consider supervised learning with the goal to learn the functional relation between inputs  $x \in \mathbb{X}$  and targets  $y \in \mathbb{Y}$  ([Figure 2.1](#)). The mapping is described by a model  $f_\theta : \mathbb{X} \rightarrow \mathbb{F}$  with adjustable parameters  $\theta \in \Theta$  that produces predictions  $f := f_\theta(x) \in \mathbb{F}$  for an input  $x$ . A prediction’s error is assessed through a convex loss function  $\ell : \mathbb{F} \times \mathbb{Y} \rightarrow \mathbb{R}$ .

*Regression* ([Example 2.1](#)) and *classification* ([Example 2.2](#)) are two common tasks that will be used frequently in later chapters. The remainder of this text sets  $\mathbb{X} = \mathbb{R}^M$  and  $\mathbb{F} = \mathbb{R}^C$ , with input and prediction space dimensions  $M, C$ , respectively ( $C$  corresponds to the number of classes for classification). Because the discussion focuses on neural networks as model, it sets  $\Theta = \mathbb{R}^D$  in what follows.



**Figure 2.1: Components of supervised learning:** The goal is to infer parameters  $\theta$  of a model  $f_\theta$ , that relates inputs  $x$  to predictions  $f$ , by minimizing a loss function  $\ell$  between the prediction  $f$  and label  $y$ .

Assume that the population of input-target pairs, *i.e.* the *data-generating process*, follows a distribution  $p_{\text{data}}(x, y)$ . A model's expected risk under this distribution is defined as

$$\begin{aligned}\mathcal{L}_{p_{\text{data}}}(\theta) &:= \mathbb{E}_{p_{\text{data}}(x, y)} [\ell(f_\theta(x), y)] \\ &= \iint_{\mathbb{X}, \mathbb{Y}} \ell(f_\theta(x), y) p_{\text{data}}(x, y) dx dy.\end{aligned}\quad (2.3a)$$

The incentive for a model to perform well is to achieve a small expected risk. Therefore, training a model is minimizing [Equation \(2.3a\)](#),

$$\underset{\theta}{\text{minimize}} \mathcal{L}_{p_{\text{data}}}(\theta). \quad (2.3b)$$

But in practice, [Equation \(2.3b\)](#) is inaccessible because the data-generating process  $p_{\text{data}}(x, y)$  is unknown. Instead, the problem is approximated through an *i.i.d.* dataset  $\mathbb{D} = \{(x_n, y_n) \in \mathbb{X} \times \mathbb{Y}\}_n$  of labeled data collected from  $p_{\text{data}}(x, y)$ .  $p_{\text{data}}(x, y)$  is approximated by the empirical distribution

$$p_{\mathbb{D}}(x, y) = \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \delta(x - x_n) \delta(y - y_n) \quad (2.4a)$$

implied by  $\mathbb{D}$ . The model's empirical risk on  $\mathbb{D}$  follows from substituting  $p_{\text{data}}$  with  $p_{\mathbb{D}}$  in [Equation \(2.3a\)](#), which yields

$$\mathcal{L}_{\mathbb{D}}(\theta) := \mathcal{L}_{p_{\mathbb{D}}}(\theta) = \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \ell(f_\theta(x_n), y_n). \quad (2.4b)$$

**Example 2.1 (Least squares regression & square loss)** Regression associates features in  $\mathbb{X} = \mathbb{R}^M$  with targets in  $\mathbb{Y} = \mathbb{R}^C$ . A prediction in  $\mathbb{F} = \mathbb{R}^C$  compares to its ground truth via the mean squared error<sup>2</sup>

$$\ell(f, y) = \frac{1}{C} \sum_{c=1}^C (y_c - f_c)^2 = \frac{1}{C} \|y - f\|_2^2 \quad (2.1)$$

2: There exist different conventions for the normalization factor. This text adapts the implementation of `MSELoss` (with `reduction="mean"` mode) in PyTorch for consistency with the code presented in later chapters. Normalizing by  $1/C$  is also close to what the name, mean squared error, suggests.

3: Sometimes, the softmax is considered part of the model rather than the loss function. This text assigns it to the loss function, in line with the PyTorch implementation of `CrossEntropyLoss`, that combines softmax and cross-entropy, which is numerically more stable [60, Chapter 4].

**Example 2.2 (C-class classification & softmax cross-entropy loss)** Classification assigns features in  $\mathbb{X} = \mathbb{R}^M$  to classes in  $\mathbb{Y} = \{1, \dots, C\}$  using a model to  $\mathbb{F} = \mathbb{R}^C$ . The softmax cross-entropy loss<sup>3</sup> maps the model's prediction to a probability distribution over classes, then uses cross-entropy to compare it with the ground truth,

$$\ell(f, y) = -\log([p(f)]_y) = -\log p(f)^\top \text{onehot}(y) \quad (2.2)$$

where  $p(f) = \text{softmax}(f)$  and  $[\text{onehot}(y)]_c = \delta_{c,y}$ .

In practice, learning happens by minimizing an empirical risk,

$$\underset{\theta}{\text{minimize}} \mathcal{L}_{\mathbb{D}}(\theta). \quad (2.4c)$$

### Challenges for Optimization

Training a model relies on optimization algorithms which are initialized at some  $\theta_0 \in \Theta$  and seek to iteratively improve the solution to [Equation \(2.4c\)](#). At an iteration  $t$ , the optimizer is given access to information about the objective at  $\theta_t$ , like derivatives ([Section 2.3](#)), to deduct a step  $\theta_{t+1} \leftarrow \theta_t$ , potentially using additional information from past observations. In general, the more local information is available, the more powerful a single update step can be. But richer information is often more costly to compute. The large-scale nature of deep learning poses challenges on the accessible information:

- ▶ **Big data:** deep learning datasets are often large because more data gives better approximations of the true distribution  $p_{\text{data}}$  via  $p_{\mathbb{D}}$ , and thus better task performance. The simultaneous computation of all per-datum losses for  $\mathcal{L}_{\mathbb{D}}(\theta)$  does not fit into memory (for many contemporary tasks, it is even infeasible to hold  $\mathbb{D}$  in memory<sup>4</sup>). Still, to obtain  $\mathcal{L}_{\mathbb{D}}$ , one could sequentially compute and reduce its summands on data chunks of manageable size. In practice, it is more common though to approximate  $\mathcal{L}_{\mathbb{D}}$  on a randomly drawn small subset, a mini-batch,  $\mathbb{B} \subseteq \mathbb{D}$  with  $|\mathbb{B}| \ll |\mathbb{D}|$ . While this avoids a computationally expensive full sweep over the data, subsampling introduces *noise* in the loss ([Section 2.1.2](#)). This noise is inherited by the computed information supplied to the optimizer. Algorithms must therefore take into account the stochastic nature of their observations.
- ▶ **Very large models:** the parameter space dimension  $D$  of DNNs usually exceeds the (already large) amount of data  $|\mathbb{D}|$ , i.e.  $D \gg |\mathbb{D}|$ . It affects the complexity to store and compute information, such as derivatives of the loss *w.r.t.*  $\theta$ , and makes it more challenging to efficiently work with higher-order information<sup>5</sup>. Therefore, working with higher-order information relies on implicit schemes (e.g. matrix-free Hessian-vector products [[127](#)]) or light-weight structured approximations (e.g. through Kronecker products [[109](#)]). Such approaches are usually technical and challenging to implement. They also add significant computational work that must be compensated for by an improved update step quality in optimizers.

In addition to these computational aspects, there exist other challenges:

- ▶ **Non-convexity:** although the loss function  $\ell$  is convex *w.r.t.* the model's output  $f$ , convexity does not carry through to the model's parameters<sup>6</sup>. This is because DNNs  $f_{\theta}$  are highly non-linear, and therefore generally non-convex, in  $\theta$ . Local minima of convex functions are global minima, meaning that local improvement gets us closer to a global solution. But non-convex problems like [Equation \(2.4c\)](#) have multiple local minima that need not be global. Through local improvements, an optimizer can arrive at one of these local minima, but with no path of improvement to a global

4: E.g. ImageNet [[41](#)] consists of 1,281,167 training images. The Kaggle download requires roughly 166 GB of memory.

5: For example, interesting quantities for an optimizer are the loss landscape's local slope (gradient) and curvature (Hessian, [Definition 3.1](#)). While the gradient is cheap to compute and store ( $D$  elements), holding the  $D \times D$  Hessian in memory is infeasible.

6: Convexity of  $\ell$  in  $f$  is still useful to construct PSD approximations to the Hessian, see [Section 3.2.3](#).

solution. This depends on various aspects, such as the update rule, hyperparameters, initialization, *etc.*

- ▶ **Generalization:** learning is not pure optimization. While minimizing the empirical risk [Equation \(2.4c\)](#) improves the model's performance on the collected data, the actual goal is to obtain good performance on new unseen data, *i.e.* generalize well. To achieve generalization, it is crucial to prevent optimization to overfit specifics in the data. It is common practice to split the data into three disjoint sets  $\mathbb{D}_{\text{train}}$ ,  $\mathbb{D}_{\text{valid}}$ , and  $\mathbb{D}_{\text{test}}$ . The train set's empirical risk  $\mathcal{L}_{\mathbb{D}_{\text{train}}}(\theta)$  is minimized, and the validation loss  $\mathcal{L}_{\mathbb{D}_{\text{valid}}}(\theta)$  serves to identify hyperparameters that lead to generalization on  $\mathbb{D}_{\text{valid}}$ . The held-out examples in the test set  $\mathbb{D}_{\text{test}}$  are used to assess generalization to new data. Another way to improve generalization is to use more data during training. Data augmentation [150] allows for cheap generation of new examples without collecting new data. Sometimes, it may be desirable to penalize model properties by adding a regularization term to the objective in [Equation \(2.4c\)](#).

### 2.1.2 Batching & Noise

Due to the large-scale nature of  $\mathbb{D}$  and  $f_\theta$  in the empirical risk, [Equation \(2.4b\)](#) is usually stochastically approximated through a mini-batch  $\mathbb{B} \subseteq \mathbb{D}$ , and assessed through a mini-batch loss

$$\mathcal{L}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(x_n, y_n) \in \mathbb{B}} \ell(f_\theta(x_n), y_n). \quad (2.5)$$

In the following, per-sample predictions and losses will often be abbreviated as  $f_n := f_\theta(x_n)$  and  $\ell_n := \ell(f_\theta(x_n), y_n)$ .

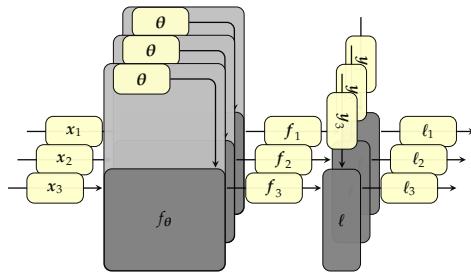
#### Batched Computations

Evaluation of the mini-batch loss in [Equation \(2.5\)](#) can be parallelized ([Figure 2.2a](#)). All mini-batch features  $\{x_n\}$  are mapped to predictions  $\{f_n\}$  by the *same instructions*  $f_\theta$ , and compared in parallel with the labels  $\{y_n\}$ , resulting in the per-sample losses  $\{\ell_n\}$ . Hardware accelerators can use this structure to achieve significant speed-up of evaluating the loss, or its derivatives ([Section 2.3](#)).

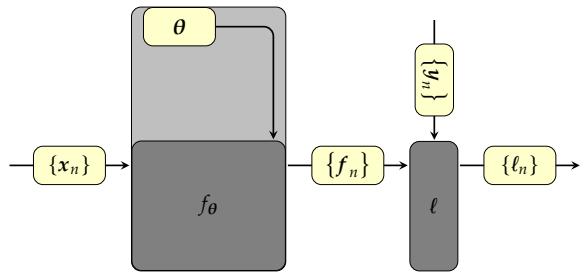
This single-instruction-multiple-data structure of the loss is often baked into ML libraries. Many of their operations natively support batched behavior, *i.e.* accept stacked inputs and assume one, usually the first, axis to correspond to a batch axis. The operation is then applied to all slices along the batch axis ([Figure 2.2b](#)).

The concept of `map` in functional programming [73] formalizes applying a function to a collection of inputs, which can be seen as a function transformation. Batched operations can be understood as transformations of the original operation. Recently developed ML libraries [23, 72] make batching explicit by providing a `vmap` interface to automatically vectorize, *i.e.* parallelize, the application of `map`.

(a) Same instructions, multiple data



(b) Batched operations



**Figure 2.2:** Different visualizations of empirical risk computation graphs. (a), The empirical risk is a weighted sum of per-sample losses that are computed independently and with the same instructions, allowing for efficient parallelization. (b) Many ML libraries natively support batched operations for efficiency. In code, this batching is often assumed, but not explicitly expressed. The vmap concept allows to make batching explicit:  $f_\theta$  and  $\ell$  in (b) correspond to vectorized versions  $\text{vmap}(f)$  and  $\text{vmap}(\ell)$  from (a).

**Definition 2.1 (vmap)** Let  $f : \mathbb{X} \rightarrow \mathbb{Y}, x \mapsto f(x)$  denote a function. The vectorized  $\text{map}$ ,  $\text{vmap}(f)$ , of  $f$  w.r.t. its argument  $x$  is a function that accepts a collection of inputs and maps each item by  $f$ , resulting in a collection of outputs,

$$\text{vmap} : (\mathbb{X} \rightarrow \mathbb{Y}) \rightarrow (\mathbb{X}^N \rightarrow \mathbb{Y}^N) \quad \forall N \in \mathbb{N}$$

such that

$$\text{vmap}(f)(\{x_1, x_2, \dots, x_N\}) = \{f(x_1), f(x_2), \dots, f(x_N)\}.$$

Collections are often abbreviated as  $\{x_n\} := \{x_1, x_2, \dots, x_N\}$  and  $\{f(x_n)\} := \{f(x_1), f(x_2), \dots, f(x_N)\}$ . Multi-variate functions can be mapped w.r.t. to a subset of arguments.

Making batching explicit with  $\text{map}$ , the mini-batch loss computation uses a vectorized model  $\text{vmap}(f_\theta)$  w.r.t. the input features  $x$ , such that

$$\text{vmap}(f_\theta)(\{x_n\}) = \{f_n(\theta)\}. \quad (2.6a)$$

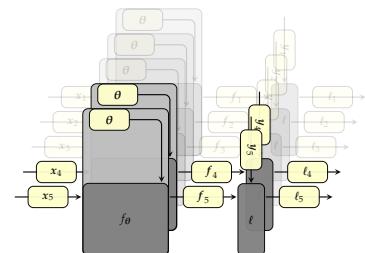
Per-sample losses, which are reduced into the mini-batch loss, are

$$\{\ell_n(\theta)\} = \text{vmap}(\ell)(\{f_n(\theta)\}, \{y_n\}). \quad (2.6b)$$

Numerically, collections of vectors like  $\{x_n\}$  etc. are represented as matrices, and more generally, collections of  $r$ -dimensional arrays are stacked into a  $(r+1)$ -dimensional arrays with an additional batch axis (e.g.,  $\{x_n\}$  is a  $|\mathbb{B}| \times M$  matrix). These arrays can then be efficiently processed in hardware accelerators.

## Noise

While the sum structure in the loss Equation (2.4b) can be efficiently parallelized, it can also be used for stochastic approximation via subsampling (Figure 2.3). The mini-batch loss  $\mathcal{L}_{\mathbb{B}}$  in Equation (2.5) is an estimator of the empirical risk  $\mathcal{L}_{\mathbb{D}}$  implied by the stochastic sampling procedure of  $\mathbb{B}$ . Increasing the batch size makes this estimator more precise but more costly to compute.



**Figure 2.3:** Illustration of stochastic subsampling: The computation graph of the empirical risk (with five data in this example) is only evaluated on a subset of data (two in this example) to save computations. While this preserves parallelization, the transparent parts are not evaluated, which introduces noise.

To see this cost-accuracy trade-off, consider the loss of a single datum  $\ell_n$  where  $n$  is uniformly drawn from  $\{1, \dots, |\mathbb{D}|\}$ . Then,  $\ell_n$  is a random variable, implied by the sampling distribution  $n \sim p(n) = \mathcal{U}(\{1, \dots, |\mathbb{D}|\})$ , and an unbiased estimator of the empirical risk,

$$\mathbb{E}_{p(n)}[\ell_n(\boldsymbol{\theta})] = \sum_{n=1}^{|\mathbb{D}|} p(n) \ell_n(\boldsymbol{\theta}) = \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}), \quad (2.7a)$$

with variance

$$\begin{aligned} \sigma^2 := \text{Var}_{p(n)}[\ell_n(\boldsymbol{\theta})] &= \mathbb{E}_{p(n)}[(\ell_n(\boldsymbol{\theta}) - \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}))^2] \\ &= \mathbb{E}_{p(n)}[\ell_n(\boldsymbol{\theta})^2] - \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta})^2. \end{aligned} \quad (2.7b)$$

Next, consider a batch with  $|\mathbb{B}|$  randomly drawn samples  $\mathbf{n} = (n_1, \dots, n_{|\mathbb{B}|})$  from a joint distribution  $p(\mathbf{n})$ , *i.e.* the mean of  $\ell_{n_1}, \dots, \ell_{n_{|\mathbb{B}|}}$ ,

$$\mathcal{L}_{\mathbf{n}}(\boldsymbol{\theta}) = \frac{1}{|\mathbb{B}|} \sum_{i=1}^{|\mathbb{B}|} \ell_{n_i}(\boldsymbol{\theta}). \quad (2.8a)$$

If samples are drawn uniformly *i.i.d.* ( $p(\mathbf{n}) = \prod_{i=1}^{|\mathbb{B}|} p(n_i)$  with  $p(n_i) = \mathcal{U}(\{1, \dots, |\mathbb{D}|\})$ ), this estimator is also unbiased,

$$\mathbb{E}_{p(\mathbf{n})}[\mathcal{L}_{\mathbf{n}}(\boldsymbol{\theta})] = \frac{1}{|\mathbb{B}|} \sum_{i=1}^{|\mathbb{B}|} \mathbb{E}_{p(n_i)}[\ell_{n_i}(\boldsymbol{\theta})] = \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}), \quad (2.8b)$$

but has smaller variance than the single-sample estimator:

$$\text{Var}_{p(\mathbf{n})}[\mathcal{L}_{\mathbf{n}}(\boldsymbol{\theta})] = \mathbb{E}_{p(\mathbf{n})}[\mathcal{L}_{\mathbf{n}}(\boldsymbol{\theta})^2] - \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta})^2 = \frac{\sigma^2}{|\mathbb{B}|}. \quad (2.8c)$$

Using more samples, *i.e.* a larger  $|\mathbb{B}|$ , decreases the variance, and thereby reduces noise in the mini-batch loss. The central limit theorem [50] connects the mini-batch estimator to a normal distribution. As the mini-batch size  $|\mathbb{B}|$  approaches infinity,  $\mathcal{L}_{\mathbf{n}}$  converges to a normal distribution with mean and variance from Equations (2.8b) and (2.8c)

$$\lim_{|\mathbb{B}| \rightarrow \infty} : \mathcal{L}_{\mathbf{n}} \sim \mathcal{N}(\mathcal{L}_{\mathbb{D}} \mid \mathcal{L}_{\mathbb{D}}, \sigma^2/|\mathbb{B}|). \quad (2.9)$$

While applications rely on finite batch sizes, it is sometimes useful to assume that this Gaussian distribution holds approximately.

Noise in the mini-batch loss propagates to other quantities, like the mini-batch gradient  $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta})$ , that are used for applications like training. Therefore, it represents a fundamental challenge for deep learning methods. It can be assessed through higher-order statistical moments such as the variance (the centered second moment, see Section 3.2.5 and Appendix C.3 for examples). Motivated by the central limit theorem Equation (2.9)—if the Gaussian approximation holds sufficiently well—only second-order statistical moments are required.

Like higher-order derivatives of multi-variate functions, higher-order statistical moments of multi-variate random variables such as the mini-batch gradient  $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta}) \in \mathbb{R}^D$  (first moment) scale exponentially with dimension  $D$  (*e.g.* the  $D \times D$  gradient covariance matrix, Section 3.2.5). Therefore, they are challenging to work with using naive approaches.

### 2.1.3 Probabilistic Interpretation

Risk-based supervised learning is often connected to learning the unknown true distribution  $p_{\text{data}}(x, y)$  through a model distribution  $p_\theta(x, y)$  and data  $\mathbb{D}$ , using estimation techniques to find a good set of parameters that minimizes a measure of dissimilarity between  $p_{\text{data}}$  and  $p_\theta$ . The following section links the loss function  $\ell$  and the model  $f_\theta$  to probabilistic objects. This will be helpful to identify additional structure in the risk minimization problem and use it for structural approximation of higher-order information (e.g. the Fisher information matrix, [Section 3.2.4](#)), and to motivate probabilistic applications (e.g. Laplace approximations, [Section 3.3.1](#)).

#### Connections to Maximum Likelihood Estimation (MLE)

The KL-divergence between the true and the model distribution,

$$\begin{aligned} D_{\text{KL}}(p_{\text{data}}(x, y) \parallel p_\theta(x, y)) \\ = \mathbb{E}_{p_{\text{data}}(x, y)} [\log p_{\text{data}}(x, y) - \log p_\theta(x, y)], \end{aligned} \quad (2.10a)$$

can be used to measure their dissimilarity. Minimizing the above expression over  $\theta$ , and dropping parameter-independent terms leads to

$$\begin{aligned} & \underset{\theta}{\text{minimize}} D_{\text{KL}}(p_{\text{data}}(x, y) \parallel p_\theta(x, y)) \\ \Leftrightarrow & \underset{\theta}{\text{minimize}} \mathbb{E}_{p_{\text{data}}(x, y)} [-\log p_\theta(x, y)]. \end{aligned} \quad (2.10b)$$

$p_{\text{data}}(x, y)$  is inaccessible and therefore empirically approximated through data  $\mathbb{D} = \{(x_n, y_n)\}_n$ . Under the *i.i.d.* assumption in the data,  $p_{\text{data}}$  is replaced with the empirical distribution  $p_{\mathbb{D}}$  from [Equation \(2.4a\)](#) and yields the accessible optimization task

$$\underset{\theta}{\text{minimize}} \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} -\log p_\theta(x_n, y_n) \quad (2.10c)$$

This expression resembles the empirical risk minimization [Equation \(2.4c\)](#) with a specific loss function  $\ell$  that produces the negative log-probability of a datum  $(x_n, y_n)$ .

Supervised learning only processes features  $x$  to predict labels  $y$ . The probabilistic model  $p_\theta(x, y)$  thus has a more special form, in that it only parameterizes the likelihood  $y | x$ ,

$$p_\theta(x, y) = p_\theta(y | x)p(x). \quad (2.11a)$$

Since only the likelihood contains parameters, [Equation \(2.10b\)](#) simplifies to minimizing the expected negative log-likelihood

$$\begin{aligned} & \underset{\theta}{\text{minimize}} \mathbb{E}_{p_{\text{data}}(x, y)} [-\log p_\theta(y | x) - \log p(x)] \\ \Leftrightarrow & \underset{\theta}{\text{minimize}} \mathbb{E}_{p_{\text{data}}(x, y)} [-\log p_\theta(y | x)] \end{aligned}$$

and with empirical approximation through *i.i.d.* data as [Equation \(2.10c\)](#),

$$\underset{\theta}{\text{minimize}} \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} -\log p_{\theta}(y_n | x_n). \quad (2.11b)$$

7: Finding the negative log-likelihood's minimum, [Equation \(2.11b\)](#), is equivalent to finding the maximum of the *i.i.d.* data's likelihood

$$p(\mathbb{D} | \theta) = \prod_{(x_n, y_n) \in \mathbb{D}} p_{\theta}(y_n | x_n) p(x_n).$$

The MLE satisfies

$$\begin{aligned} \theta_{\text{MLE}} &= \arg \max p(\mathbb{D} | \theta) \\ &= \arg \max \log p(\mathbb{D} | \theta) \\ &= \arg \min -\log p(\mathbb{D} | \theta). \end{aligned}$$

Inserting  $p(\theta | \mathbb{D})$  and dropping parameter-independent terms recovers the problem [Equation \(2.11b\)](#).

8: Inserting mean and covariance into the negative log-probability yields,

$$\begin{aligned} &-\log \mathcal{N}(y; \mu, \Sigma) \\ &= 1/2(y - \mu)^T \Sigma^{-1} (y - \mu) \\ &\quad + 1/2 [\log \det \Sigma + C \log 2\pi], \end{aligned}$$

*i.e.* the square loss [Equation \(2.1\)](#) up to a  $\theta$ -independent term which does not affect optimization.

9: With  $p_c$  denoting the probability to observe class  $y = c$ ,

$$-\log \text{Cat}(y; p) = -\log p_y,$$

which is the softmax cross-entropy loss [Equation \(2.2\)](#).

This minimization problem corresponds to MLE<sup>7</sup> with a statistical model  $p_{\theta}(y | x)$ . This is a specific form of empirical risk minimization where the model's prediction  $f_{\theta}(x)$  parameterizes a likelihood  $q$  for  $y | f$ , and a negative log-likelihood loss function  $\ell$ , *i.e.*

$$p_{\theta}(y | x) = q(y | f_{\theta}(x)), \quad (2.12a)$$

$$\ell(f, y) = -\log q(y | f). \quad (2.12b)$$

Both the square loss and softmax cross-entropy loss ([Examples 2.1](#) and [2.2](#)) have such a probabilistic interpretation, see [Examples 2.3](#) and [2.4](#).

**Example 2.3 (Probabilistic interpretation of square loss)** The square loss [Equation \(2.1\)](#) is the negative log-likelihood of a Gaussian centered around the model's prediction with diagonal constant covariance,

$$\ell(f, y) = -\log q(y | f)$$

$$\text{with } q(y | f_{\theta}(x)) = \mathcal{N}(y | \mu, \Sigma)$$

where<sup>8</sup>  $\mu = f_{\theta}(x)$  and  $\Sigma = C/2I$ .

**Example 2.4 (Probabilistic interpretation of softmax cross-entropy loss)** The softmax cross-entropy loss [Equation \(2.2\)](#) is the negative log-likelihood of a multinomial distribution parameterized by the softmax probabilities,

$$\ell(f, y) = -\log q(y | f)$$

$$\text{with } q(y | f_{\theta}(x)) = \text{Cat}(y; p)$$

where<sup>9</sup>  $p = \text{softmax}(f_{\theta}(x))$ .

Following the maximum likelihood principle results in the MLE parameter  $\theta_{\text{MLE}}$  which satisfies [Equation \(2.11b\)](#) and gives rise to the distribution  $p_{\theta_{\text{MLE}}}(y | x) = q(y | f_{\theta_{\text{MLE}}}(x))$  as an approximation to  $p_{\text{data}}(y | x)$ .

## Connections to Maximum A Posteriori (MAP) Estimation

MLE maximizes the likelihood  $p(\mathbb{D} | \theta)$ . In a probabilistic formulation with a prior  $p(\theta)$  over the parameters and evidence  $p(\mathbb{D}) = \int_{\Theta} p(\mathbb{D} | \theta)p(\theta) d\theta$  for the data, one can instead consider the posterior

$$p(\theta | \mathbb{D}) = \frac{p(\mathbb{D} | \theta)p(\theta)}{p(\mathbb{D})} \quad (2.13)$$

which corrects the prior with data observations. Such a posterior is useful to form probabilistic beliefs over predictions  $y_*$  for new inputs  $x_*$ ,

$$\begin{aligned} p(y_* | x_*, \mathbb{D}) &= \int_{\Theta} p(y_* | x_*, \mathbb{D}, \theta) p(\theta | x_*, \mathbb{D}) d\theta \\ &= \int_{\Theta} p(y_* | x_*, \theta) p(\theta | \mathbb{D}) d\theta. \end{aligned} \quad (2.14)$$

However, this requires integration over the posterior, which itself is almost always intractable and thus needs to be approximated.

The MAP principle approximates the posterior with a delta distribution around the posterior mode, *i.e.* the point of maximum posterior density,

$$p(\theta | \mathbb{D}) \approx \delta(\theta - \theta_{\text{MAP}}) \quad \text{where } \theta_{\text{MAP}} = \arg \max_{\theta} p(\theta | \mathbb{D}). \quad (2.15)$$

It is connected to empirical risk minimization with a regularization term that results from the prior  $p(\theta)$ . To see this, reformulate [Equation \(2.15\)](#) to minimize the negative log-posterior, expand Baye's rule ([Equation \(2.13\)](#)) and neglect the parameter-independent evidence term. Then, apply the same assumptions as for MLE<sup>10</sup>, *i.e.* *i.i.d.* data and  $\theta$  only parameterizing the likelihood  $y | x$ . This yields

$$\begin{aligned} \theta_{\text{MAP}} &= \arg \min_{\theta} -\log p(\theta | \mathbb{D}) \\ &= \arg \min_{\theta} -\log p(\mathbb{D} | \theta) - \log p(\theta) \\ &= \arg \min_{\theta} \sum_{(x_n, y_n) \in \mathbb{D}} -\log p(y_n | x_n, \theta) - \log p(\theta) \\ &= \arg \min_{\theta} \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} -\log p(y_n | x_n, \theta) - \frac{\log p(\theta)}{|\mathbb{D}|}. \end{aligned} \quad (2.16)$$

In analogy to [Equation \(2.12\)](#), the first term is an empirical risk ([Equation \(2.11b\)](#)) with negative log-likelihood loss of a distribution  $q$  for targets given model predictions,

$$p(y | x, \theta) = q(y | f_\theta(x)) \quad \text{and} \quad \ell(f, y) = -\log q(y | f). \quad (2.17)$$

However, this risk is extended by a regularization term from the prior,

$$\theta_{\text{MAP}} = \arg \min_{\theta} \mathcal{L}_{\mathbb{D}}(\theta) + r(\theta) \quad \text{where } r(\theta) = -\frac{\log p(\theta)}{|\mathbb{D}|}. \quad (2.18)$$

[Equations \(2.13\), \(2.16\)](#) and [\(2.18\)](#) connect the posterior with the loss via  $\log p(\mathbb{D} | \theta) = -|\mathbb{D}| \mathcal{L}_{\mathbb{D}}(\theta)$  and  $\log p(\theta) = -|\mathbb{D}| r(\theta)$ ,

$$\begin{aligned} p(\theta | \mathbb{D}) &= \frac{\exp [\log p(\mathbb{D} | \theta) + \log p(\theta)]}{p(\mathbb{D})} \\ &= \frac{\exp \{-|\mathbb{D}| [\mathcal{L}_{\mathbb{D}}(\theta) + r(\theta)]\}}{p(\mathbb{D})}. \end{aligned} \quad (2.19)$$

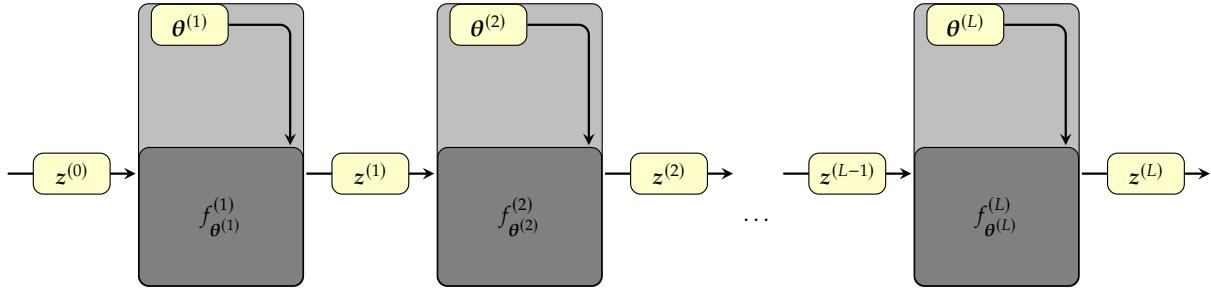
It underlines the aforementioned challenges to track the posterior. The exponent is non-linear in  $\theta$ , and  $p(\mathbb{D})$  requires computing an integral.

[Equation \(2.19\)](#) gives rise to posterior approximations that go beyond a delta distribution. The Laplace approximation [93] ([Section 3.3.1](#)) also

10: Mathematically, they translate into

$$\begin{aligned} p(\mathbb{D} | \theta) &= \prod_n p(x_n, y_n | \theta) \\ &= \prod_n p(y_n | x_n, \theta) p(x_n | \theta) \\ &= \prod_n p(y_n | x_n, \theta) p(x_n) \end{aligned}$$

with slightly different notation for  $\theta$  in comparison to the MLE discussion, as it is now treated probabilistically.



**Figure 2.4:** Forward pass of a sequential feedforward neural network (Equation (2.20a)). The computational graph indicates the data flow and dependencies of intermediate variables.

starts with the MAP estimate, but uses a quadratic Taylor expansion of the log-posterior around  $\theta_{\text{MAP}}$  to approximate the posterior by a Gaussian. This quadratic expansion requires higher-order information in form of second-order derivatives, presented in [Chapter 3](#).

## 2.2 Neural Networks

The previous section focused on structure in the empirical risk—like its sum structure, and interpretations of risk-based learning for specific loss functions—with assumptions about the model  $f_\theta$ . This section introduces structure in the model. While DNNs are generally highly over-parameterized, they usually rely on relatively simple components and construction principles.

This manuscript considers sequential feedforward neural networks that consist of layers. They comprise “classic” architectures like multi-layer perceptrons (MLPs), convolutional neural networks (CNNs), and established architectures like VGG [153], and ResNets [68]. Most of the discussion in this text also applies to other architectures, but would require heavier notation.

### 2.2.1 Layer-wise Notation

Sequential feedforward neural networks of depth  $L$  consist of *modules*, or *layers*,  $f_{\theta^{(l)}}$ ,  $l = 1, \dots, L$ , stacked on top of each other such that

$$f_\theta = f_{\theta^{(L)}}^{(L)} \circ f_{\theta^{(L-1)}}^{(L-1)} \circ \dots \circ f_{\theta^{(1)}}^{(1)} \quad (2.20a)$$

They map input features  $x =: z^{(0)}$  to predictions  $f_\theta(x) =: z^{(L)}$  via a sequence of intermediate hidden features  $z^{(1)}, \dots, z^{(L-1)}$ . In a forward pass, a module  $f_{\theta^{(l)}}$  receives the parental output  $z^{(l-1)} \in \mathbb{R}^{h^{(l-1)}}$  and applies an operation with (potentially empty) parameters  $\theta^{(l)} \in \mathbb{R}^{d^{(l)}}$ ,

$$z^{(l)} = f_{\theta^{(l)}}^{(l)}(z^{(l-1)}). \quad (2.20b)$$

The output features  $z^{(l)} \in \mathbb{R}^{h^{(l)}}$  serve as input to the next layer  $l + 1$ . This builds up dependencies in form of the computational graph shown in [Figure 2.4](#) that maps the leaf nodes  $z^{(0)}$  and  $\theta^{(1)}, \dots, \theta^{(L)}$  to the prediction

$z^{(L)}$ . The neural network parameters are often treated as a single vector  $\theta \in \mathbb{R}^D$  which results from layer-wise concatenation,

$$\theta = \begin{pmatrix} \theta^{(1)} \\ \theta^{(2)} \\ \vdots \\ \theta^{(L)} \end{pmatrix}. \quad (2.22)$$

To simplify the presentation, [Equation \(2.20b\)](#) assumes vector-shaped quantities. However, many neural networks process higher-dimensional data like images, represented by tensors. Sometimes, the tensor structure is convenient to work with. One can convert between the tensor and vector view without loss of generality by introducing conventions for tensor flattening ([Definition 2.2](#)) and reshaping ([Definition 2.3](#)); after all, multi-dimensional arrays are represented in a vector format in memory.

However, there exist different flattening conventions. Implementations often favor row-major ordering. This manuscript uses (the more common in literature) column-major order as it allows for elegant generalizations of derivative concepts for multi-variate functions, like the Jacobian ([Definition 2.5](#)), the Hessian ([Definition 3.2](#)), and their chain rules ([Theorems 2.2](#) and [3.1](#)). To translate analytical results into implementations, it is crucial to be aware of these differing conventions.

## 2.2.2 Modularity & Common Operations

An important strength of deep learning is its *modularity*. ML libraries provide a large number of operations, or modules, that can be combined in almost arbitrary ways through function composition, like in [Equation \(2.20a\)](#). Training the resulting models with first-order methods remains simple because their gradient can be automatically computed via AD ([Section 2.3](#)). New operations can easily be added because its implementation is decoupled to the modular level.

Modules are vaguely defined. Often, multiple operations that form a logical processing unit in a neural network are grouped into a single module, *e.g.* an MLP layer combines affine transformation and elementwise activation (see below). In extreme cases, even an entire neural network can be considered a single module that can be used in other neural networks; *e.g.* the neural network in [Figure 2.1](#) resembles a single layer in [Figure 2.4](#) and could act as one layer in a larger network.

For theoretical analyses, it is preferable to consider units with a small number of operations as modules. This, however, is inconvenient for constructing large architectures, where many operations are grouped into higher-level units. This text adapts a rather fine-grained view on modules that is close to their implementation in ML libraries like PyTorch.

A common categorization for modules distinguishes trainable functions with parameters, and parameter-free operations. [Table 2.1](#) lists the forward passes of common operations that will be illustrated in the following presentation of network architectures. To distinguish more clearly between input and output of an operation, the notation uses the symbols

### Definition 2.2 (Tensor flattening)

Let  $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2 \times \dots \times n_m}$  denote a tensor of rank  $m$  with dimensions  $n_1, n_2, \dots, n_m$ . The flattened tensor  $\text{vec}(\mathbf{A}) \in \mathbb{R}^{n_1 n_2 \cdots n_m}$  is a vector that concatenates  $\mathbf{A}$ 's elements in a first-index-varies-fastest fashion,

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} A_{1,1,1,\dots,1} \\ A_{2,1,1,\dots,1} \\ \vdots \\ A_{n_1,1,1,\dots,1} \\ A_{1,2,1,\dots,1} \\ \vdots \\ A_{n_1,2,1,\dots,1} \\ \vdots \\ A_{n_1,n_2,n_3,\dots,n_m} \end{pmatrix}. \quad (2.21)$$

The matrix case  $m = 2$  corresponds to column-stacking. Flattening a vector  $a$  leaves it unaffected, *i.e.*  $\text{vec}(a) = a$ .

### Definition 2.3 (Vector reshaping)

Let  $a \in \mathbb{R}^{n_1 n_2 \cdots n_m}$  be a vector. Reshaping that vector into a rank- $m$  tensor of shape  $S = (n_1, n_2, \dots, n_m)$  happens by filling  $a$ 's elements into the tensor in a first-index-varies-fastest fashion,

$$\mathbf{A} = \text{reshape}_S(a) \quad (2.23a)$$

with elements

$$\begin{aligned} A_{1,1,1,\dots,1} &= a_1, \\ A_{2,1,1,\dots,1} &= a_2, \\ &\vdots \\ A_{n_1,1,1,\dots,1} &= a_{n_1}, \\ A_{1,2,1,\dots,1} &= a_{n_1+1}, \\ &\vdots \\ A_{n_1,2,1,\dots,1} &= a_{2n_1}, \\ &\vdots \\ A_{n_1,n_2,n_3,\dots,n_m} &= a_{n_1 n_2 n_3 \cdots n_m}. \end{aligned} \quad (2.23b)$$

The matrix case  $m = 2$  corresponds to column-filling. With tensor flattening ([Definition 2.2](#)) this allows to define tensor reshaping: A tensor  $\mathbf{B}_1$  of shape  $S_1$  is rearranged into any tensor  $\mathbf{B}_2$  of compatible shape  $S_2$  by first flattening, then reshaping it, *i.e.*  $\mathbf{B}_2 := \text{reshape}_{S_2}(\mathbf{B}_1) := \text{reshape}_{S_2}(\text{vec } \mathbf{B}_1)$ .

**Table 2.1: Forward pass for common modules used in feedforward networks.**  
Input and output are denoted  $x, z$  rather than  $z^{(l)}, z^{(l+1)}$  to avoid clutter.  $I$  is the identity matrix. Bold upper-case symbols ( $W, X, Z, \dots$ ) denote matrices and bold upper-case sans serif symbols ( $\mathbf{W}, \mathbf{X}, \mathbf{Z}, \dots$ ) denote tensors. See Appendices A.2 to A.4 for details, and Tables 2.2 and 4.1 for extended versions of this table for the backward, and Hessian backward, pass.

OPERATION	FORWARD
Matrix-vector multiplication	$z(x, W) = Wx$
Matrix-matrix multiplication	$Z(X, W) = WX$
Addition	$z(x, b) = x + b$
Elementwise activation	$z(x) = \phi(x), \text{s.t. } z_i(x) = \phi(x_i)$
Skip-connection	$z(x, \theta) = x + s(x, \theta)$
Reshape/view	$Z(\mathbf{X}) = \text{reshape}(\mathbf{X})$
Index select/map $\pi$	$z(x) = \Pi x, \Pi_{j, \pi(j)} = 1,$
Convolution	$Z(\mathbf{X}, \mathbf{W}) = \mathbf{X} * \mathbf{W},$ $Z(W, \llbracket \mathbf{X} \rrbracket) = W \llbracket \mathbf{X} \rrbracket,$
Square loss	$\ell(f, y) = 1/C(y - f)^\top (y - f)$
Softmax cross-entropy	$\ell(f, y) = -\text{onehot}(y)^\top \log [p(f)]$

$x, z$  for module input and output instead of  $z^{(l-1)}, z^{(l)}$ , and neglects the layer superscript for the parameters, writing  $\theta$  instead of  $\theta^{(l)}$ .

### Deep Linear Networks & Multi-layer Perceptrons (MLPs)

*Linear layers* process inputs  $x$  by affine transformation, *i.e.* multiplication with a weight matrix  $W$ , followed by addition of a bias vector  $b$ ,

$$z = Wx + b \quad \text{where} \quad \theta = ((\text{vec } W)^\top \quad b^\top)^\top. \quad (2.24)$$

They are also referred to as *fully-connected* layers, because each output  $z_i$  depends on all inputs  $x$  through  $W_{i,:}$  and  $b_i$ .

*Deep linear networks* (Example 2.5) consist of only linear layers and are of interest for analytical studies [*e.g.* 16, 112, 143] as they are somewhat tractable. They describe a linear feature map, *i.e.* a linear function *w.r.t.* the input  $z^{(0)}$ , that is non-linear in the parameters. Therefore, such networks are as expressive as a single linear layer, but highly overparameterized.

A common technique to turn a deep linear network into a non-linear feature map is to interlace affine transformations with non-linear activations [136]. An *activation layer*  $\phi$  acts elementwise on its input, *i.e.* applies the same function  $\phi$  to each input element,

$$z = \phi(x) \quad \text{such that} \quad z_i = \phi(x_i).$$

There exist many activations (ReLU, sigmoid, tanh, *etc.* [60, Chapter 6]), and recent works proposing new choices (*e.g.* squared ReLU [156]).

*Multi-layer perceptrons* (MLPs, Example 2.6) combine affine transformation and activation in a layer. Activation functions  $\phi^{(l)}$  may vary between layers, but are often chosen identically, with no activation in the last layer. One way to interpret this design is that the mapping  $z^{(0)} \mapsto z^{(L-1)}$  acts as non-linear feature transformation, and the last layer  $z^{(L-1)} \mapsto z^{(L)}$  as linear classifier for the learned features.

### Convolutional Neural Networks (CNNs)

CNNs represent an important neural network architecture revolution and were the first class of deep neural network to beat “classical” methods on the ImageNet competition [41, 91, 138]. Broadly speaking, such networks contain convolutional layers with trainable parameters.

#### Example 2.5 (Deep linear network)

In notation of Equation (2.20), a deep linear network of depth  $L$  reads

$$z^{(l)} = W^{(l)} z^{(l-1)} + b^{(l)}$$

where

$$\theta^{(l)} = \left( (\text{vec } W^{(l)})^\top \quad b^{(l)^\top} \right)^\top, \quad l = 1, \dots, L.$$

#### Example 2.6 (Multi-layer perceptron (MLP))

In terms of Equation (2.20), an MLP of depth  $L$  reads

$$z^{(l)} = \phi^{(l)} \left( W^{(l)} z^{(l-1)} + b^{(l)} \right),$$

where

$$\theta^{(l)} = \left( (\text{vec } W^{(l)})^\top \quad b^{(l)^\top} \right)^\top, \quad l = 1, \dots, L,$$

$$\phi^{(L)} = \text{id}$$

and id denotes the identity.

**Convolutional layers:** Convolutions process multi-channel input features such as images and are parameterized by a kernel that can be imagined as a filter for patterns like edges, corners, *etc.* During the convolution operation, the kernel slides over the input features and produces an output element by contraction with the overlapping elements of the image. In most cases, each output channel is also shifted by a bias parameter which will be neglected in this presentation for simplicity (detailed discussion in [Appendix A.4](#), example in [Equation \(A.8\)](#)). Because the kernel moves over the image, it can detect patterns irrespective of their position. The process can be adjusted with various hyperparameters, such as stride, padding, groups, dilation (see [45] for a visual guide).

In contrast to the linear layer ([Equation \(2.24\)](#)) where each output is connected to all inputs via independent rows of the weight matrix, the parameters in the kernel are shared across all outputs. Therefore, convolutions usually require less parameters than fully-connected layers.

Nevertheless, both layers are related because convolution is a linear operation and can therefore be regarded as matrix multiplication. Due to the weight sharing, this matrix is structured by the kernel elements<sup>11</sup>. Alternatively, one can stack patches—input elements that overlap with the kernel at each stage—into columns of a matrix, which yields the unfolded input, denoted by  $\llbracket \mathbf{X} \rrbracket$  in [Table 2.1](#). Then, convolution is a matrix-matrix product between a matrix reshape of the kernel and the unfolded input [29] (see [Figure A.3c](#) for an illustration).

11: *E.g.*, in the one-dimensional case, the convolution of two vectors can be computed by expanding one into a Toeplitz matrix, and multiplying that onto the second vector [171].

**Padding & pooling layers:** Convolutions are often combined with other modules. *Padding layers* add pixels around the outer dimensions of an image, which helps to reduce information loss at the image boundaries during convolution. *Pooling layers* down-sample images by summarizing patches of pixels and reduce the number of hidden features. Similar to convolution, pooling considers patches of an input image. Two common summary operations are per-channel averaging and taking the per-channel maximum. They give rise to maximum and average pooling.

One can interpret padding and pooling as scatter operations, realized by multiplication with a sparse, binary matrix  $\Pi$  (compare [Table 2.1](#)):

- ▶ For padding,  $\Pi$  does not depend on the input, but only its shape and the hyperparameters. A row is empty if its index corresponds to the padded area, and otherwise contains a one at the element's index to be copied from the input.
- ▶ For maximum pooling,  $\Pi$  depends on the input. Each row contains a one at the index of the element with maximum value in the patch.
- ▶ For average pooling,  $\Pi$  does not depend on the input, but only its shape and the hyperparameters. Each row contains the inverse patch size at indices of the elements in the current patch.

### Residual Networks (ResNets)

The inclusion of skip (or residual) connections [68] represents another revolution in the design of CNNs, and enabled training of deeper architectures, with 100 or even 1000 layers. This lead to improved performance

of such CNNs on tasks like ImageNet [41, 138]. Skip connections branch off a hidden feature and feed it back after the residual block  $s(x, \theta)$ ,

$$z = x + s(x, \theta).$$

This can be seen as learning a small modification  $s(x, \theta)$ —a residual function—for  $x$ ; hence the name *residual connection*.

### Closing Remarks & Sources of Non-linearity

This short overview of common neural network layers and architectures is, of course, incomplete. Other famous layers include dropout [158], recursive [33, 47, 71], normalization [79, 174], attention layers [167], etc.

An interesting observation about the operations in Table 2.1 is that most of them are linear (linear, convolution, padding, and average pooling layers) or piece-wise linear (maximum pooling and ReLU activation layer) *w.r.t.* both their input and parameters. This implies that their second- and higher-order derivatives vanish. Non-linearity is often only introduced by activation layers (elementwise) and loss functions (after the model’s forward pass). While the properties of these components are not inherited by the entire neural network, they give rise to structure in a network’s higher-order derivatives, see *e.g.* Chapter 4.

## 2.3 Automatic Differentiation

Together with empirical risk minimization and neural networks, the last ingredient in the ML pipeline, Procedure 1.1, is computing the gradient. Contemporary methods to train neural networks (Section 3.1) rely on this quantity. ML libraries compute it via their built-in automatic differentiation (AD), built around the famous backpropagation algorithm [137], described in more detail here.

### 2.3.1 Foundations

#### An Example & Path Interpretation

Given a program that evaluates a function, AD produces a program to evaluate its derivative. The general idea is to consider the function as composition of atomic operations (*e.g.* addition, multiplication, . . . ) with known derivatives. To automatically compute the derivatives, one needs to track the dependencies between intermediate variables and combine their derivatives using the chain rule.

Figure 2.5 illustrates the basic principles of AD for an example from [119]. The starting point is a function defined by code (Figure 2.5a). Such a function transforms input to output variables through atomic operations and builds up intermediate variables along its execution. The relation between these operations are described by a directed graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with a set of nodes  $\mathcal{V} = \{z^{(1)}, z^{(2)}, \dots\}$  and a set of edges  $\mathcal{E} = \{(z^{(i_1)}, z^{(j_1)}), (z^{(i_2)}, z^{(j_2)}), \dots\}$  where  $(z^{(i)}, z^{(j)})$  denotes a directed edge from node  $z^{(i)}$  to  $z^{(j)}$  (Figure 2.5b).

(a) Function as Python program

```

1 from math import exp, sin
2
3 def z7(z1: float, z2: float):
4     """Example function."""
5
6     # intermediate variables
7     z3 = sin(z2)
8     z5 = z3 + z2
9     z4 = exp(z1)
10    z6 = z4 * z5
11
12    # output variable
13    z7 = z4 + z6
14
15    return z7

```

(b) Computation graph

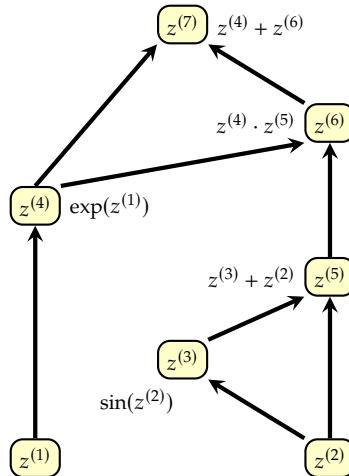


Figure 2.5: Basic AD principles [119].

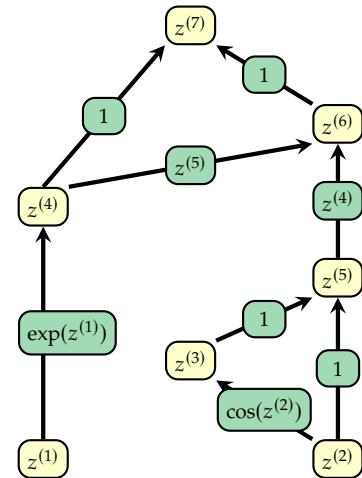
(a) Example input program represented by Python code. The atomic operations combined via composition are addition, multiplication, exponential function, and sine. The result  $z^{(7)}$  is computed from the inputs  $z^{(1)}, z^{(2)}$  through intermediates  $z^{(3)}, z^{(4)}, z^{(5)}, z^{(6)}$ ,

$$z^{(7)} = \exp(z^{(1)}) + \exp(z^{(1)}) (\sin(z^{(2)}) + z^{(2)}) .$$

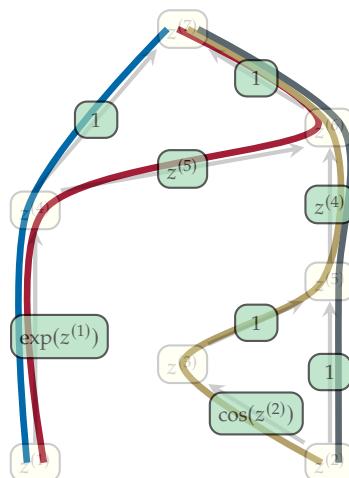
(b) Representation as computation graph to track dependencies between the intermediate variables on the level of atomic operations. (c) Computing derivatives relies on local derivatives  $\partial z^{(j)} / \partial z^{(i)}$  on edges  $(z^{(i)}, z^{(j)})$ , which need to be accumulated according to the chain rule. (d) Interpretation of the chain rule as sum over path products. Computing the derivatives of a node  $w.r.t.$  to another node in the graph requires summing the path product of local derivatives for all paths that connect them. In detail:

$$\begin{aligned} \frac{\partial z^{(7)}}{\partial z^{(1)}} &= \frac{\partial z^{(7)}}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial z^{(1)}} + \frac{\partial z^{(7)}}{\partial z^{(6)}} \frac{\partial z^{(6)}}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial z^{(1)}} \\ &= \exp(z^{(1)}) + z^{(5)} \exp(z^{(1)}) \\ &= \exp(z^{(1)}) \\ &\quad + (\sin(z^{(2)}) + z^{(2)}) \exp(z^{(1)}) , \\ \frac{\partial z^{(7)}}{\partial z^{(2)}} &= \frac{\partial z^{(7)}}{\partial z^{(6)}} \frac{\partial z^{(6)}}{\partial z^{(5)}} \frac{\partial z^{(5)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial z^{(2)}} \\ &\quad + \frac{\partial z^{(7)}}{\partial z^{(6)}} \frac{\partial z^{(6)}}{\partial z^{(5)}} \frac{\partial z^{(5)}}{\partial z^{(2)}} \\ &= z^{(4)} \cos(z^{(2)}) + z^{(4)} \\ &= \exp(z^{(1)}) \cos(z^{(2)}) + \exp(z^{(1)}) . \end{aligned}$$

(c) Local derivatives



(d) Bauer paths



To compute derivatives, the local derivatives on edges (Figure 2.5c) are combined according to the chain rule. For  $\partial z^{(j)} / \partial z^{(i)}$  between two variables in the graph, all paths connecting them need to be considered. A path between two nodes  $z^{(i)}$  and  $z^{(j)}$  is a sequence of edges that connect them: starting from  $z^{(i)}$ , following the edges in a path leads to  $z^{(j)}$ . Let  $[z^{(i)} \rightarrow z^{(j)}]$  denote the set of paths connecting  $z^{(i)}$  to  $z^{(j)}$ . Then the derivative is the sum of path products of local derivatives (Figure 2.5d),

$$\frac{\partial z^{(j)}}{\partial z^{(i)}} = \sum_{p \in [z^{(i)} \rightarrow z^{(j)}]} \prod_{(z^{(k)}, z^{(l)}) \in p} \frac{\partial z^{(l)}}{\partial z^{(k)}} . \quad (2.25)$$

The path formulation goes back to Bauer [12].

### The Jacobian Matrix & Its Chain Rule

For the computation graph of a neural network's empirical risk, tracking dependencies between variables at a scalar level would result in a considerable book-keeping overhead due to the large number of connections. This can be circumvented by using vector-valued (or tensor-valued) nodes  $z^{(1)}, z^{(2)}, \dots$  and tracking edges between vectors (or tensors), see

12: Definition 2.4 assumes vector-valued functions. With the flattening convention Definition 2.2, it generalizes to tensor-valued functions as follows:

**Definition 2.5 (Generalized Jacobian)** Let  $B : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}^{m \times p}$ ,  $A \mapsto B(A)$  be a differentiable matrix-to-matrix function. The Jacobian  $J_A B(A)$  of  $B$  w.r.t.  $A$  is an  $mp \times nq$  matrix

$$J_A B(A) = \frac{\partial \text{vec } B(A)}{\partial (\text{vec } A)^\top} \quad (2.27a)$$

with entries

$$[J_A B(A)]_{i,j} = \frac{\partial [\text{vec } B(A)]_i}{\partial [\text{vec } A]_j} \quad (2.27b)$$

and the flattening operation  $\text{vec}$  from Definition 2.2 [103, Chapter 9.4]. The analogous tensor case  $(A, B) \rightarrow (A, B)$  requires lengthy notation and is therefore omitted.

In the context of neural networks, the most common occurrences of Definition 2.5 involve vector-to-vector functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $x \mapsto f(x)$  with

$$J_x f(x) = \frac{\partial f(x)}{\partial x^\top}.$$

For instance,  $x$  can be considered the input or bias vector of a layer applying an affine transformation. Other cases involve matrix-to-vector mappings  $f : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}^m$ ,  $X \mapsto f(X)$  with

$$J_X f(X) = \frac{\partial f(X)}{\partial (\text{vec } X)^\top},$$

where  $X$  might correspond to the  $\mathbb{R}^{m \times q}$  weight matrix of a linear layer. See Table 2.2 for an overview.

13: Proper arrangement of partial derivatives leads to a generalized Jacobian chain rule for matrices/tensors:

**Theorem 2.2 (Generalized Jacobian chain rule)** Let  $B : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}^{m \times p}$  and  $C : \mathbb{R}^{m \times p} \rightarrow \mathbb{R}^{r \times s}$  be differentiable matrix-to-matrix functions. Let  $D = C \circ B : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}^{r \times s}$ ,  $A \mapsto D(A) = C(B(A))$  be their composition. Then,

$$J_A D(A) = [J_B C(B)] J_A B(A) \quad (2.30)$$

with the generalized Jacobian Definition 2.5 [103, Chapter 5.15]. The tensor case  $(D, C, B, A) \rightarrow (D, C, B, A)$  is analogous.

**Figure 2.6.** However, this requires a generalization of Equation (2.25) to multi-variate nodes. The accumulation is efficiently expressed as matrix multiplication by arranging partial derivatives into *Jacobians*.

**Definition 2.4 (Jacobian)** Let  $b : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $a \mapsto b(a)$  be a differentiable vector-to-vector function. The Jacobian  $J_a b(a)$  of  $b$  w.r.t.  $a$  is an  $m \times n$  matrix with partial derivatives,

$$J_a b(a) = \frac{\partial b(a)}{\partial a^\top}, \quad \text{with} \quad [J_a b(a)]_{i,j} = \frac{\partial b_i(a)}{\partial a_j}. \quad (2.26)$$

Matrix- and tensor-valued functions require flattening their arguments into vectors<sup>12</sup>. For a vector-to-scalar function  $b(a)$ , i.e.  $m = 1$ , the Jacobian has one row that contains the gradient,  $[J_a b(a)]^\top = \nabla_a b$ . The gradient will often be denoted by  $g$ . E.g.  $g_{p_{\text{data}}}(\theta) := \nabla_\theta \mathcal{L}_{p_{\text{data}}}(\theta)$  for the gradient of the population risk Equation (2.3a), and  $g_{\mathbb{D}}(\theta) := \nabla_\theta \mathcal{L}_{\mathbb{D}}(\theta)$  for the gradient of the empirical risk Equation (2.4b) on a dataset  $\mathbb{D}$  (with  $\mathbb{D} = \mathbb{D}_{\text{train}}$ ,  $\mathbb{B}$  for the train loss and the mini-batch gradient).

In the vector-valued case, one must accumulate Jacobians through matrix multiplies instead of scalar multiplications to compute derivatives,

$$J_{z^{(j)}} z^{(i)} = \sum_{p \in [z^{(i)} \rightarrow z^{(j)}]} \prod_{(z^{(k)}, z^{(l)}) \in p} J_{z^{(k)}} z^{(l)}(z^{(k)}). \quad (2.28)$$

The product term generalizes the chain rule to vector-valued functions.

**Theorem 2.1 (Jacobian chain rule)** Let  $b : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $a \mapsto b(a)$  and  $c : \mathbb{R}^m \rightarrow \mathbb{R}^r$ ,  $b \mapsto c(b)$  be differentiable vector-to-vector functions. Consider their composition  $d = c \circ b : \mathbb{R}^n \rightarrow \mathbb{R}^r$ ,  $a \mapsto d(a) = c(b(a))$ . The composition's Jacobian  $J_a d(a) \in \mathbb{R}^{r \times n}$  is related to the composite Jacobians via

$$J_a d(a) = [J_b c(b)] J_a b(a). \quad (2.29)$$

This can be generalized to tensor-valued functions<sup>13</sup>.

### Jacobian Accumulation (Automatic Differentiation Modes)

Given the computation graph  $\mathcal{G}$  of a function to be differentiated, Equation (2.28) describes the operations that need to be performed. But there are different schedules for carrying out these computations, with differing performance: e.g., it is possible to share accumulated derivative products between paths that share subpaths (like paths and in Figure 2.5d). And for a single path in the vector case, the optimal contraction order of the Jacobian matrix chain (one summand of Equation (2.28)) depends on the dimension of the nodes (Example 2.7). The following Jacobian accumulation schedules are of specific interest for AD:

- ▶ **Forward accumulation, or forward mode AD**, starts at the leafs, i.e. the nodes w.r.t. which the function is differentiated. Jacobians are accumulated in the same order as the function evaluation.

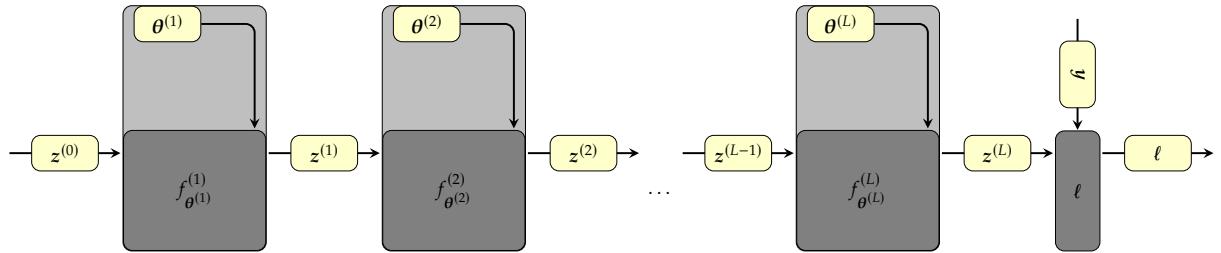


Figure 2.6: Computation graph of a sequential feedforward neural network’s loss for a single datum from Equation (2.31a).

- ▶ **Reverse accumulation, or reverse mode AD** [62, 99], starts at the root, *i.e.* the variable that is differentiated. Jacobians are accumulated from root to leaf nodes, traversing the graph backwards. This is often called a *backward pass*.
- ▶ **Optimal Jacobian accumulation** computes derivatives according to the optimal schedule which usually traverses the computation graph in a nontrivial fashion. For arbitrary computation graphs, finding this schedule is NP-hard [114].

Due to the specific structure of computation graphs in deep learning, reverse mode AD is often more practical than forward accumulation. This is outlined in the following section, that illustrates reverse mode for differentiation of a neural network’s loss *w.r.t.* its parameters, leading to the famous backpropagation algorithm [137].

### 2.3.2 Gradient Backpropagation

Gradient backpropagation [137] enables efficient differentiation of the training objective in deep learning. It is the central algorithm of popular ML libraries with built-in AD. This section presents backpropagation for chain-structured computation graphs (see [60, Chapter 6] for the general case) like the loss of a sequential feedforward neural network. Starting from the loss of a single datum, the goal is to show that ML libraries combine AD and batching to maximize efficiency. But this limits their functionality to computing the gradient, ignoring *e.g.* the per-sample structure in the loss. Alleviating this limitation to compute richer information (Chapter 3) using the existing implementation of gradient backpropagation is a main goal of Part II in this thesis.

#### Loss of a Single Datum

Consider the loss implied by a single datum  $(x, y)$ , a loss function  $\ell$ , and a neural network  $f_\theta$  depicted in Figure 2.6,

$$\ell(\theta) = \ell(f_\theta(x), y) \quad \text{with} \quad f_\theta = f_{\theta^{(L)}}^{(L)} \circ f_{\theta^{(L-1)}}^{(L-1)} \circ \dots \circ f_{\theta^{(1)}}^{(1)}. \quad (2.31a)$$

Its computation graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  has nodes

$$\mathcal{V} = \left\{\theta^{(l)}\right\}_{l=1}^L \cup \left\{z^{(l)}\right\}_{l=0}^L \cup \{y, \ell\} \quad (2.31b)$$

and edges

$$\mathcal{E} = \left\{ (\mathbf{z}^{(l-1)}, \mathbf{z}^{(l)}) \right\}_{l=1}^L \cup \left\{ (\boldsymbol{\theta}^{(l)}, \mathbf{z}^{(l)}) \right\}_{l=1}^L \cup \left\{ (\mathbf{z}^{(L)}, \ell), (\mathbf{y}, \ell) \right\}. \quad (2.31c)$$

Each edge implies a Jacobian, categorized as one of the following:

- ▶ The *input-output* Jacobian  $\mathbf{J}_{\mathbf{z}^{(l-1)}} \mathbf{z}^{(l)}(\mathbf{z}^{(l-1)})$  of a module  $l$ .
- ▶ The *parameter-output* Jacobian  $\mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{z}^{(l)}(\boldsymbol{\theta}^{(l)})$  of a module  $l$ .
- ▶ The *prediction-loss* Jacobian  $\mathbf{J}_{\mathbf{z}^{(L)}} \ell(\mathbf{z}^{(L)})$  has one column and will be written as gradient of the loss *w.r.t.* the model prediction,  $\mathbf{J}_{\mathbf{z}^{(L)}} \ell(\mathbf{z}^{(L)}) = [\nabla_{\mathbf{z}^{(L)}} \ell(\mathbf{z}^{(L)})]^\top$ .

The goal is to compute *parameter-loss* Jacobians, *i.e.* gradients of the loss *w.r.t.* parameters  $\nabla_{\boldsymbol{\theta}^{(l)}} \ell(\boldsymbol{\theta}^{(l)}) = [\nabla_{\boldsymbol{\theta}^{(l)}} \ell(\boldsymbol{\theta}^{(l)})]^\top$  for all layers  $l = 1, \dots, L$ .

First, consider only the gradient  $\nabla_{\boldsymbol{\theta}^{(l)}} \ell$  of one parameter  $\boldsymbol{\theta}^{(l)}$ . [Equation \(2.28\)](#) requires identifying all paths connecting  $\boldsymbol{\theta}^{(l)}$  to  $\ell$ . Due to the computation graph's chain structure, this is only a single path,

$$[\boldsymbol{\theta}^{(l)} \rightarrow \ell] = \{p\} \quad (2.32a)$$

with

$$p = \left( (\boldsymbol{\theta}^{(l)}, \mathbf{z}^{(l)}), (\mathbf{z}^{(l)}, \mathbf{z}^{(l+1)}), \dots, (\mathbf{z}^{(L-1)}, \mathbf{z}^{(L)}), (\mathbf{z}^{(L)}, \ell) \right). \quad (2.32b)$$

Plugging this into [Equation \(2.28\)](#) simplifies to

$$\mathbf{J}_{\boldsymbol{\theta}^{(l)}} \ell(\boldsymbol{\theta}) = \underbrace{[\mathbf{J}_{\mathbf{z}^{(L)}} \ell]}_{1 \times d^{(l)}} \underbrace{[\mathbf{J}_{\mathbf{z}^{(L-1)}} \mathbf{z}^{(L)}]}_{1 \times h^{(L)}} \cdots \underbrace{[\mathbf{J}_{\mathbf{z}^{(l)}} \mathbf{z}^{(l+1)}]}_{h^{(L)} \times h^{(L-1)}} \underbrace{[\mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{z}^{(l)}]}_{h^{(l+1)} \times h^{(l)}}, \quad (2.32c)$$

and in gradient notation

$$\nabla_{\boldsymbol{\theta}^{(l)}} \ell(\boldsymbol{\theta}) = \underbrace{[\mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{z}^{(l)}]^\top}_{d^{(l)}} \underbrace{[\mathbf{J}_{\mathbf{z}^{(l)}} \mathbf{z}^{(l+1)}]^\top}_{h^{(l)} \times h^{(l+1)}} \cdots \underbrace{[\mathbf{J}_{\mathbf{z}^{(L-1)}} \mathbf{z}^{(L)}]^\top}_{h^{(L-1)} \times h^{(L)}} \underbrace{\nabla_{\mathbf{z}^{(L)}} \ell}_{h^{(L)}}. \quad (2.32d)$$

In comparison to the general formulation [Equation \(2.28\)](#), the rather simple graphs of a neural network's loss yield much simpler expressions ([Equations \(2.32c\)](#) and [\(2.32d\)](#)) that are a result of the Jacobian chain rule ([Theorem 2.1](#)) applied to the loss [Equation \(2.31a\)](#). They also illustrate the impact of contraction order on performance due to the connection to matrix chains<sup>14</sup>, mentioned in [Section 2.3.1](#).

In forward mode, the matrix chain [Equation \(2.32d\)](#) would be evaluated from left to right, starting with a  $d^{(l)} \times h^{(l)}$  Jacobian that is transformed into  $d^{(l)} \times h^{(l')}$  matrices where  $l' > l$ . Since the parameter count  $d^{(l)}$  in a layer is large in DNNs, these intermediate matrices are costly to store.

In reverse mode, [Equation \(2.32d\)](#) is evaluated from right to left, starting with a  $h^{(L)}$ -dimensional vector that is transformed into vectors of dimension  $h^{(l')}$  with  $L > l' \geq l$ . These intermediate accumulations require less memory than forward mode.

Each approach has drawbacks, however. Reverse mode uses a more efficient matrix multiplication order, but the entire graph must have been evaluated and stored, or re-computed, to construct the Jacobians. While this is not needed for forward mode, forward accumulation starts with the Jacobian of an edge  $(\boldsymbol{\theta}^{(l)}, \mathbf{z}^{(l)})$  that is not shared with paths for other

14: Assuming no cost to compute a Jacobian, the optimal Jacobian contraction of [Equations \(2.32c\)](#) and [\(2.32d\)](#) are matrix chain problems that can be solved with dynamic programming: given  $n$  matrices  $A_1, A_2, \dots, A_n$ , the task is to find the optimal contraction schedule of  $A_1 A_2 \cdots A_n$ . This is crucial for performance, as this example from Cormen et al. [36, Chapter 15.2] illustrates:

**Example 2.7 (Matrix chain contraction)** Let  $A_1 \in \mathbb{R}^{10 \times 100}$ ,  $A_2 \in \mathbb{R}^{100 \times 5}$ ,  $A_3 \in \mathbb{R}^{5 \times 50}$ . There are two schedules to evaluate the chain  $A_1 A_2 A_3$  (cost for addition neglected for simplicity):

- ▶  $(A_1 A_2) A_3$ :  $B_1 = A_1 A_2 \in \mathbb{R}^{10 \times 5}$  costs 100 multiplications per element (5,000 in total).  $B_1 A_3 \in \mathbb{R}^{10 \times 50}$  costs 5 multiplications per element (2,500 in total).
- ▶  $A_1 (A_2 A_3)$ :  $B_2 = A_2 A_3 \in \mathbb{R}^{100 \times 50}$  costs 5 multiplications per element (25,000 in total).  $A_1 B_2 \in \mathbb{R}^{10 \times 50}$  costs 100 multiplications per element (50,000 in total).

The order  $(A_1 A_2) A_3$  uses 10x fewer operations (7,500 versus 75,000).

parameters. Therefore, intermediate accumulations can not be reused for other gradients. This is another crucial property of reverse mode, which does allow reuse of intermediate results: consider the accumulations to obtain the gradient  $\nabla_{\theta} \ell$  of all layers,

$$\begin{aligned}\nabla_{\theta^{(1)}} \ell &= \left[ J_{\theta^{(1)}} z^{(1)} \right]^T \left[ J_{z^{(1)}} z^{(2)} \right]^T \left[ J_{z^{(2)}} z^{(3)} \right]^T \cdots \left[ J_{z^{(L-1)}} z^{(L)} \right]^T \nabla_{z^{(L)}} \ell \\ \nabla_{\theta^{(2)}} \ell &= \left[ J_{\theta^{(2)}} z^{(2)} \right]^T \left[ J_{z^{(2)}} z^{(3)} \right]^T \cdots \left[ J_{z^{(L-1)}} z^{(L)} \right]^T \nabla_{z^{(L)}} \ell \\ \vdots &= \vdots \quad \ddots \quad \vdots \\ \nabla_{\theta^{(L-1)}} \ell &= \left[ J_{\theta^{(L-1)}} z^{(L-1)} \right]^T \left[ J_{z^{(L-1)}} z^{(L)} \right]^T \nabla_{z^{(L)}} \ell \\ \nabla_{\theta^{(L)}} \ell &= \left[ J_{\theta^{(L)}} z^{(L)} \right]^T \nabla_{z^{(L)}} \ell\end{aligned}$$

The paths for any two parameters  $\theta^{(l_1)}, \theta^{(l_2)}$  with  $l_2 > l_1$  share edges

$$(z^{(l_2)}, z^{(l_2+1)}), \dots, (z^{(L-1)}, z^{(L)}), (z^{(L)}, \ell).$$

Therefore, their matrix chains share the accumulated gradient

$$\nabla_{z^{(l_2)}} \ell = \left[ J_{z^{(l_2)}} z^{(l_2+1)} \right]^T \left[ J_{z^{(l_2+1)}} z^{(l_2+2)} \right]^T \cdots \left[ J_{z^{(L)}} z^{(L)} \right]^T \nabla_{z^{(L)}} \ell.$$

This gradient *w.r.t.* hidden features is passed backwards through the graph and used by a layer, before updating it and passing it to the next. The interpretation of the described accumulation scheme is therefore known as gradient backpropagation algorithm [137]:

**Definition 2.6 (Gradient backpropagation for sequential feedforward neural networks)** Given the computation graph of a loss  $\ell(f_{\theta}(x), y)$  implied by a datum  $(x, y)$ , a loss function  $\ell$ , and a sequential feedforward neural network  $f_{\theta}$  from Equation (2.31a), gradient backpropagation recovers the gradient vector  $\nabla_{\theta} \ell = ((\nabla_{\theta^{(1)}} \ell)^T, (\nabla_{\theta^{(2)}} \ell)^T, \dots, (\nabla_{\theta^{(L)}} \ell)^T)^T$  in stages by passing gradients backward through the graph (Figure 2.7a):

- ▶ Initialize the backpropagated vector with  $\nabla_{z^{(L)}} \ell$  at  $L$ .
- ▶ For layer  $l = L, \dots, 1$ 
  1. Receive the output gradient  $\nabla_{z^{(l)}} \ell$ .
  2. Recover the parameter gradient  $\nabla_{\theta^{(l)}} \ell = \left[ J_{\theta^{(l)}} z^{(l)} \right]^T \nabla_{z^{(l)}} \ell$ .
  3. Compute the input gradient  $\nabla_{z^{(l-1)}} \ell = \left[ J_{\theta^{(l)}} z^{(l)} \right]^T \nabla_{z^{(l)}} \ell$ .
  4. Free  $\nabla_{z^{(l)}} \ell$  and send  $\nabla_{z^{(l-1)}} \ell$  to layer  $l - 1$ .

Definition 2.6 is *modular*, as mentioned earlier in Section 2.2.2: it only relies on local derivatives. Supporting a new operation only requires specifying its forward pass and the vector-Jacobian products (VJPs<sup>15</sup>) with its input-output and parameter-output Jacobians. Given functionality to create and traverse computation graphs, implementations of backpropagation are very extensible due to its abstraction to the modular level.

Backpropagation itself performs a VJP with the network's parameter-output Jacobian  $J_{\theta} z^{(L)}$ . Choosing the vector to be  $\nabla_{z^{(L)}} \ell$  yields the gradient  $\nabla_{\theta} \ell = [J_{\theta} z^{(L)}]^T \nabla_{z^{(L)}} \ell$ . But one can also compute VJPs  $[J_{\theta} z^{(L)}]^T v$  with arbitrary vectors  $v$ . The model's parameter-output Jacobian is crucial for computing higher-order information (Chapter 3 and Part II).

Although backpropagation only requires VJPs, the Jacobian matrix is

15: During backpropagation, the transposed Jacobian is right-multiplied onto the backpropagated gradient vector (Equation (2.32d) from right to left). One can see this as left-multiplying the Jacobian to a column vector (Equation (2.32c) from left to right), *i.e.* computing a vector-Jacobian product. Forward accumulation instead left-multiplies the transposed Jacobian onto a matrix (Equation (2.32d) from left to right). One can view this as right-multiplying the Jacobian onto the transposed matrix (Equation (2.32c) from right to left). This requires multiple Jacobian-vector products (VJPs), or a Jacobian-matrix product (JMP).

**Table 2.2: Jacobians (Definition 2.5) for common modules in feedforward networks.** Input and output are denoted  $x, z$  rather than  $z^{(l)}, z^{(l+1)}$  to avoid clutter.  $I$  is the identity matrix. Matrices use bold upper-case symbols ( $\mathbf{W}, \mathbf{X}, \mathbf{Z}, \dots$ ), tensors use bold upper-case sans serif symbols ( $\mathbf{\mathbb{W}}, \mathbf{\mathbb{X}}, \mathbf{\mathbb{Z}}, \dots$ ). Most Jacobians can be elegantly derived with matrix differential calculus, see Appendix A.1 for details.

OPERATION	FORWARD	JACOBIAN (Definition 2.5)	DETAILS
Matrix-vector multiplication	$z(x, \mathbf{W}) = \mathbf{W}x$	$\mathbf{J}_x z = \mathbf{W},$ $\mathbf{J}_{\mathbf{W}} z = x^\top \otimes I$	Appendix A.2.1
Matrix-matrix multiplication	$\mathbf{Z}(\mathbf{X}, \mathbf{W}) = \mathbf{W}\mathbf{X}$	$\mathbf{J}_{\mathbf{X}} \mathbf{Z} = I \otimes \mathbf{W},$ $\mathbf{J}_{\mathbf{W}} \mathbf{Z} = \mathbf{X}^\top \otimes I$	Appendix A.2.1
Addition	$z(x, b) = x + b$	$\mathbf{J}_x z = \mathbf{J}_b z = I$	Appendix A.2.1
Elementwise activation	$z(x) = \phi(x)$ , s.t. $z_i(x) = \phi(x_i)$	$\mathbf{J}_x z = \text{diag}[\phi'(x)]$	Appendix A.2.2
Skip-connection	$z(x, \theta) = x + s(x, \theta)$	$\mathbf{J}_x z = I + \mathbf{J}_x s,$ $\mathbf{J}_{\theta} z = \mathbf{J}_{\theta} s$	Appendix A.2.3
Reshape/view	$\mathbf{Z}(\mathbf{X}) = \text{reshape}(\mathbf{X})$	$\mathbf{J}_{\mathbf{X}} \mathbf{Z} = I$	Appendix A.4.1
Index select/map $\pi$	$z(x) = \Pi x, \Pi_{j,\pi(j)} = 1$	$\mathbf{J}_x z = \Pi$	Appendix A.4.2
Convolution	$\mathbf{Z}(\mathbf{X}, \mathbf{W}) = \mathbf{X} \star \mathbf{W},$ $\mathbf{Z}(W, [\mathbf{X}]) = W [\mathbf{X}]$	$\mathbf{J}_{[\mathbf{X}]} \mathbf{Z} = I \otimes \mathbf{W},$ $\mathbf{J}_W \mathbf{Z} = [\mathbf{X}]^\top \otimes I$	Appendix A.4.3
Square loss	$\ell(f, y) = 1/c(y - f)^\top (y - f)$	$\mathbf{J}_f \ell = 2(f - y)^\top$	Appendix A.3.1
Softmax cross-entropy	$\ell(f, y) = -\text{onehot}(y)^\top \log[p(f)]$	$\mathbf{J}_f \ell = (y - p(f))^\top$	Appendix A.3.2

an interesting object for analytical studies into its structure, and for efficient implementation of functionality that goes beyond computing gradients (e.g. matrix-Jacobian products (MJP)). Table 2.2 contains the Jacobians of the common operations in neural networks from Section 2.2.2. These Jacobians are conveniently obtained using matrix differential calculus [103], presented in Appendix A.1.

### Backpropagation & Batching

To see the interplay between AD and batching (Section 2.1.2), consider differentiation of the mini-batch loss (Equation (2.5)). Popular ML libraries like PyTorch construct the computation graph on the level of tensor-valued variables with an additional batch axis.

**General case:** Starting from the previous differentiation of a single datum loss, the mini-batch scenario follows by the substitutions

$$\begin{aligned} x &=: z^{(0)} \leftrightarrow \mathbf{X} =: \mathbf{Z}^{(0)} \in \mathbb{R}^{|\mathbb{B}| \times h^{(0)}}, \\ y &\leftrightarrow \mathbf{Y} \in \mathbb{R}^{|\mathbb{B}| \times C}, \\ z^{(l)} &\leftrightarrow \mathbf{Z}^{(l)} \in \mathbb{R}^{|\mathbb{B}| \times h^{(l)}}, \\ \ell &\leftrightarrow \boldsymbol{\ell} \in \mathbb{R}^{|\mathbb{B}|}, \end{aligned} \quad (2.33a)$$

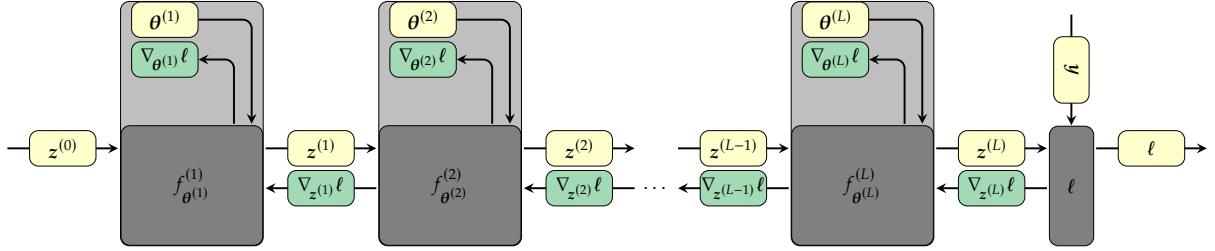
with stacked data  $\mathbf{X} = (x_1 \ x_2 \ \dots \ x_{|\mathbb{B}|})$  and  $\mathbf{Y} = (y_1 \ y_2 \ \dots \ y_{|\mathbb{B}|})$ , and assuming matrix-to-matrix layer functions. To account for the reduction of per-sample losses  $\ell$ , the graph is extended by

$$\mathcal{L}(\ell) = \text{mean}(\ell) = \frac{1}{|\mathbb{B}|} \sum_{n=1}^{|\mathbb{B}|} [\ell]_n. \quad (2.33b)$$

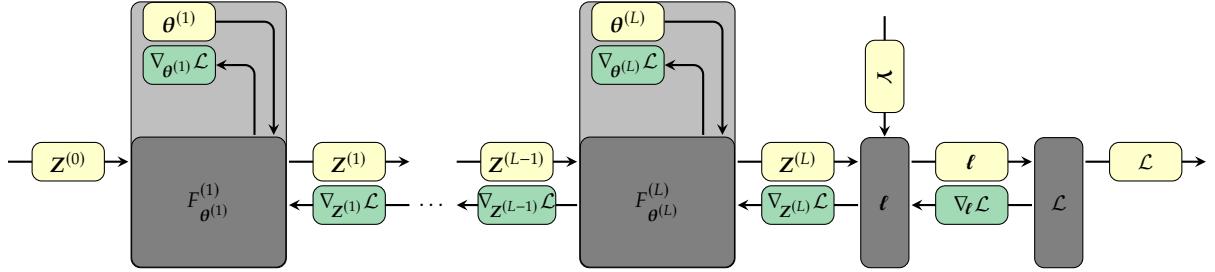
The computation graph, shown in Figure 2.7b, is  $\mathcal{G} = (\mathcal{E}, \mathcal{V})$  with nodes

$$\mathcal{V} = \left\{ \boldsymbol{\theta}^{(l)} \right\}_{l=1}^L \cup \left\{ \mathbf{Z}^{(l)} \right\}_{l=0}^L \cup \{ \mathbf{Y}, \boldsymbol{\ell}, \mathcal{L} \} \quad (2.33c)$$

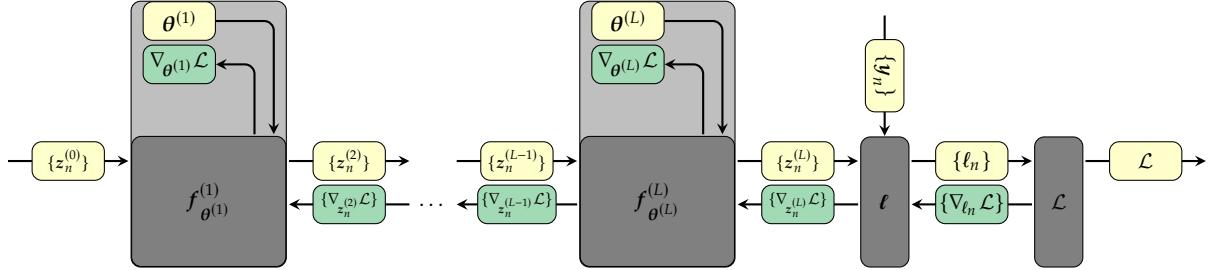
(a) One datum



(b) Batched data, arbitrary transformations



(c) Batched data, batched instructions



**Figure 2.7: Gradient backpropagation and (un)awareness of per-sample structure in many ML libraries.** (a) Computation graph of of a neural network’s loss on a single datum ( $x, y$ ). Gradients are backpropagated through the graph as described by Definition 2.6 to obtain  $\nabla_{\theta} \ell$ . (b) To exploit parallelism in the computations, multiple data are stacked into matrices ( $X, Y$ ) which are then processed by a sequence of matrix-to-matrix functions  $F_{\theta^{(l)}}$  into a batch of losses  $\ell$ , and reduced into a scalar  $\mathcal{L}$  via mean reduction. AD in popular ML libraries like PyTorch tracks variables on the level of batched tensors. Therefore, operations are allowed to build up dependencies between data—such that  $\mathcal{L} = 1/|\mathbb{B}| \sum_n [\ell(X, Y, \theta)]_n$  where each component of  $\ell$  may depend on *all* data (batch normalization [79] is such a case)—without breaking gradient backpropagation. ML libraries implement VJPS for the matrix-to-matrix functions  $F_{\theta^{(l)}}$ . This loses structure for operations that treat inputs independently along the batch axis. (c) The empirical risk on a mini-batch (Equation (2.5)) is such a case: all operations in the graph process inputs independently and with the same instructions along the batch axis. The following connections to the single datum case (a) hold:  $F_{\theta^{(l)}} \leftrightarrow f_{\theta^{(l)}} = \text{vmap}(f_{\theta^{(l)}})$ ,  $\ell = \text{vmap}(\ell)$  with vmap from Definition 2.1. Due to the more general support of AD in ML libraries for graphs of the form (b), their VJPs cannot be accessed per-sample.

and edges

$$\begin{aligned} \mathcal{E} = & \left\{ (\mathbf{Z}^{(l-1)}, \mathbf{Z}^{(l)}), \right\}_{l=1}^L \cup \left\{ (\boldsymbol{\theta}^{(l)}, \mathbf{Z}^{(l)}) \right\}_{l=1}^L \\ & \cup \left\{ (\mathbf{Z}^{(L)}, \ell), (\mathbf{Y}, \ell), (\ell, \mathcal{L}) \right\}. \end{aligned} \quad (2.33d)$$

Gradient backpropagation (Definition 2.6) carries over the mini-batch loss graph Equation (2.33) and efficiently recovers the gradient  $\nabla_{\theta} \mathcal{L}$ .

Popular libraries like PyTorch implement the required functionality, VJPs, for the matrix-to-matrix functions that process inputs with a batch axis (see Definition 2.5 for the Jacobian’s generaliation to matrix functions). Hence, operations  $\mathbf{Z}^{(l)} \mapsto \mathbf{Z}^{(l+1)}$  are allowed to create dependencies across the batch axis without breaking the gradient computation<sup>16</sup>.

16: E.g. batch normalization [79] introduces dependencies along the batch dimension by centering and re-scaling the input with statistics computed across the batch axis.

**Per-sample structure:** This work focuses on empirical risks that are averages over *per-sample* losses (recall [Equation \(2.4b\)](#)),

$$\mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta}) = \frac{1}{|\mathbb{B}|} \sum_{(x_n, y_n) \in \mathbb{B}} \ell(f_{\boldsymbol{\theta}}(x_n), y_n).$$

Hence, all operations act independently, and using the same instructions, along the batch dimension. All variables in the graph inherit this independence along their batch axis,

$$\mathbf{Z}^{(l)} = \begin{pmatrix} z_1^{(l)} & z_2^{(l)} & \dots & z_{|\mathbb{B}|}^{(l)} \end{pmatrix}, \quad l = 0, \dots, L, \quad (2.34a)$$

$$\boldsymbol{\ell} = \begin{pmatrix} \ell(f_{\boldsymbol{\theta}}(x_1), y_1) \\ \ell(f_{\boldsymbol{\theta}}(x_2), y_2) \\ \vdots \\ \ell(f_{\boldsymbol{\theta}}(x_{|\mathbb{B}|}), y_{|\mathbb{B}|}) \end{pmatrix}. \quad (2.34b)$$

For all layers  $l = 1, \dots, L$ , this independence across samples implies block-diagonal input-output Jacobians,

$$\mathbf{J}_{\mathbf{Z}^{(l-1)}} \mathbf{Z}^{(l)} = \begin{pmatrix} \mathbf{J}_{z_1^{(l-1)}} z_1^{(l)} & 0 & \dots & 0 \\ 0 & \mathbf{J}_{z_2^{(l-1)}} z_2^{(l)} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \mathbf{J}_{z_{|\mathbb{B}|}^{(l-1)}} z_{|\mathbb{B}|}^{(l)} \end{pmatrix} \quad (2.34c)$$

and per-sample block structure in the output-parameter Jacobian

$$\mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{Z}^{(l)} = \begin{pmatrix} \mathbf{J}_{\boldsymbol{\theta}^{(l)}} z_1^{(l)} & \mathbf{J}_{\boldsymbol{\theta}^{(l)}} z_2^{(l)} & \dots & \mathbf{J}_{\boldsymbol{\theta}^{(l)}} z_{|\mathbb{B}|}^{(l)} \end{pmatrix}. \quad (2.34d)$$

Many implementations of backpropagation make it difficult to access this per-sample structure. They only expose VJPs  $v \mapsto [\mathbf{J}_{\mathbf{Z}^{(l)}} \mathbf{Z}^{(l+1)}]^T v$  and  $v \mapsto [\mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{Z}^{(l+1)}]^T v$  for [Equations \(2.34c\)](#) and [\(2.34d\)](#). While this allows supporting AD of more general graphs than those of an empirical risk ([Equation \(2.4b\)](#)), it limits access to only the *average* gradient

$$\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta}) = \frac{1}{|\mathbb{B}|} \sum_{(x_n, y_n) \in \mathbb{B}} \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(x_n), y_n) \quad (2.35)$$

17: This notation is closer to recent developments in AD for ML libraries [23, 72] that separate batching and AD more clearly through vectorization via a vmap interface ([Definition 2.1](#)). A different way to arrive at the batched computation graph in [Figure 2.7c](#) is to start from the computation graph of a single datum's loss  $\ell(f_{\boldsymbol{\theta}}(x), y)$  and vectorize it to obtain the set of graphs  $\{\ell(f_{\boldsymbol{\theta}}(x_n), y_n)\}$ , which can be stacked into a single graph that produces  $\boldsymbol{\ell}$ . Appending the reduction node to this graph, one obtains the computation graph for the average loss.

when differentiating an empirical risk. Computing per-sample gradients is not more demanding than computing the average gradient—the only difference is taking the average—but their computation is not supported. More flexible access to the Jacobians [Equations \(2.34c\)](#) and [\(2.34d\)](#) through per-sample VJPs (or MJPs), enables the computation of various higher-order quantities. Their efficient realization will be the main focus of [Chapters 4, 5 and 7](#).

To highlight independence across the batch axis in a computation graph, a set notation will be preferred over the matrix notation ([Equations \(2.33a\)](#) and [\(2.34\)](#)) in the following. E.g. the text uses  $\{z_n^{(l)}\}_n$ , or just  $\{z_n^{(l)}\}$ , instead of  $\mathbf{Z}^{(l)}$ . [Figure 2.7c](#) illustrates this set notation<sup>17</sup>.

# 3.

## Higher-order Information

The previous chapter highlighted the important structure of deep neural networks and the empirical risk used to train them:

- ▶ **Sum structure & noise:** Given a model  $f_\theta$ , a convex loss function  $\ell$ , and a dataset  $\mathbb{D}$ , the empirical risk is an average over per-datum losses. Each of these losses is identically computed by feeding a datum through the model and the loss function. Batching these computations allows to efficiently evaluate them in parallel. Doing so on only a random subset of data—a mini-batch—reduces computational cost in exchange for noise.
- ▶ **Probabilistic interpretation:** For specific loss functions, which include regression and softmax classification, empirical risk minimization can be interpreted as maximum likelihood—or maximum a posteriori—estimation where the model parameterizes a likelihood distribution  $q(\mathbf{y} \mid f_\theta(\mathbf{x}))$ .
- ▶ **Layer-structure & modularity:** Neural networks consist of layers, or modules, which are glued together by function composition. Backpropagation abstracts the differentiation by the composition of the module-level derivatives, which simplifies supporting new operations and hides the complexity of AD from practitioners.

Popular ML libraries allow for efficient and automated gradient computation, but combine AD and batching in a way that allows to support more general computation graphs than the empirical risk. However, the added optimization complicates efficiently assessing the empirical risk’s per-sample structure, and higher-order information which requires slightly more flexible AD operations, like matrix-Jacobian products.

This chapter motivates why relying solely on the gradient has limitations, and—based on structure in the empirical risk—proposes quantities in the form of higher-order information that go beyond the gradient. It starts by presenting currently popular deep learning optimization methods ([Section 3.1](#)) which rely heavily on the average mini-batch gradient provided through AD in ML libraries. While training algorithms that use more information than the gradient have been historically mostly unexplored—because the quantities required are not as optimized and automated as gradient computations—they have exciting promises.

Second-order methods ([Section 3.2](#)) use more information than just the gradient. They are known to improve over first-order methods in “classic” optimization problems (convex optimization, general linear models). They incorporate curvature information in form of the Hessian ([Section 3.2.1](#)), or PSD approximations thereof, such as the generalized Gauss-Newton ([Section 3.2.3](#)) and Fisher information matrix ([Section 3.2.4](#)).

While optimization methods for neural networks are one main application in deep learning, there are other applications to quantities beyond the gradient. [Section 3.3](#) highlights additional use cases which underline the relevance of such information beyond optimization. [Chapter 6](#) presents another use case focused on improving the training of neural networks.

3.1 Popular Deep Learning Optimizers . . . . .	32
3.2 Second-order Optimization	35
3.3 Other Applications . . . . .	43

### 3.1 Popular Deep Learning Optimizers

[Section 2.1](#) formulates learning as minimizing the empirical risk, for which an optimization algorithm seeks to find a solution. Deep learning optimizers are iterative: they start from an initial parameter  $\theta_0$  and aim to improve the parameters in iterations, leading to a trajectory  $\theta_0, \theta_1, \dots$  in parameter space. An iteration is described by an update rule  $\mathcal{M} : \theta_t \mapsto \theta_{t+1}$  that relies on internal states of the optimizer from the observation history  $H_t$ , local observations at the current iterate  $\theta_t$ , and hyperparameters  $\phi_t$  [34]. The next iterate results from a simplified optimization problem where the actual objective is replaced by a computationally cheap approximation. Often, this direction is further adapted by outer-loop mechanisms, like line searches, or hyperparameters that require manual tuning to achieve good performance.

The update quality is influenced by the quality of the local approximation. Computing exact information about the objective  $\mathcal{L}_{\mathbb{D}}$  is expensive, but yields a more accurate local proxy. Due to the large-scale nature of deep learning, algorithms only rely on stochastic information from a mini-batch  $\mathbb{B}_t$  in form of the mini-batch loss  $\mathcal{L}_{\mathbb{B}_t}$  ([Section 2.1.2](#)). In addition to locality, noise further complicates deriving a precise local approximation.

#### Stochastic Gradient Descent

One of the most popular methods, according to [Figure 1.1](#), is stochastic gradient descent (SGD, [135]). Its update derives from

$$\mathcal{M} : \quad \theta_{t+1} = \arg \min_{\theta} m_{\theta_t}(\theta) = \theta_t - \eta g_{\mathbb{B}_t}(\theta_t) \quad (3.1a)$$

where  $m_{\theta_t}(\theta)$  is a first-order Taylor around the current iterate  $\theta_t$  based on the mini-batch gradient, regularized by a quadratic  $L_2$  penalty to discourage large steps with a learning rate  $\eta$ ,

$$m_{\theta_t}(\theta) = \mathcal{L}_{\mathbb{B}_t}(\theta_t) + g_{\mathbb{B}_t}(\theta_t)^T (\theta - \theta_t) + \frac{1}{2\eta} \|\theta - \theta_t\|_2^2 \quad (3.1b)$$

This leads to [Update Rule 3.1](#), which updates the parameters with the scaled negative mini-batch gradient. The negative gradient also corresponds to the direction of steepest descent, *i.e.* along which the loss decreases most rapidly around  $\theta_t$  (see [Section 3.2.4](#) and [Equation \(3.18\)](#)).

#### Momentum Methods (Incorporating Past Knowledge)

The descent direction in SGD becomes poorer as the mini-batch gradient is subject to more noise. Incorporating previous gradient observations can help reduce noise and find a better descent direction [60, Chapter 8.3]. From a physical interpretation, the optimizer builds up a velocity (referred to as *momentum*) that is adapted with new gradient observation. Two variations are the heavy ball [128] and Nesterov [117] momentum methods ([Update Rules 3.2](#) and [3.3](#)). Another intuition for momentum is that it suppresses SGD's oscillating behavior in the presence of two directions with differing curvature by averaging out gradients that point in opposite directions (see *e.g.* central panel of [Figure 6.1](#)).

**Update Rule 3.1 (SGD [135])**

$$\theta_{t+1} = \theta_t - \eta g_{\mathbb{B}_t}(\theta_t)$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

**Update Rule 3.2 (Momentum [128])**

$$\theta_{t+1} = \theta_t - v_t$$

where

$$v_t = \rho v_{t-1} + \eta g_{\mathbb{B}_t}(\theta_t)$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{momentum factor } \rho \in [0; 1)$$

$$\text{initial momentum } v_0 \in \mathbb{R}^D$$

**Update Rule 3.3 (NAG [117])**

$$\theta_{t+1} = \theta_t - v_t$$

where

$$v_t = \rho v_{t-1} + \eta g_{\mathbb{B}_t}(\theta_t - \rho v_{t-1})$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{momentum factor } \rho \in [0; 1)$$

$$\text{initial momentum } v_0 \in \mathbb{R}^D$$

**Update Rule 3.4 (AdaGrad [44])**

$$\theta_{t+1} = \theta_t - \eta g_{\mathbb{B}_t}(\theta_t) \oslash (s_t^{\circ 1/2} + \epsilon \mathbf{1})$$

where

$$s_t = s_{t-1} + (g_{\mathbb{B}_t}(\theta_t))^{\circ 2}$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{divide-by-zero safe guard } \epsilon \in \mathbb{R}^+$$

$$\text{initial } s_0 \in \mathbb{R}^D$$

**Update Rule 3.5 (RMSPProp [164])**

$$\theta_{t+1} = \theta_t - \eta g_{\mathbb{B}_t}(\theta_t) \oslash (s_t^{\circ 1/2} + \epsilon \mathbf{1})$$

where

$$s_t = \rho s_{t-1} + (1 - \rho)(g_{\mathbb{B}_t}(\theta_t))^{\circ 2}$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{decay rate } \rho \in [0; 1)$$

$$\text{divide-by-zero safe guard } \epsilon \in \mathbb{R}^+$$

$$\text{initial } s_0 \in \mathbb{R}^D$$

**Update Rule 3.6 (Adadelta [180])**

$$\theta_{t+1} = \theta_t - \eta \Delta_t$$

where

$$\Delta_t = (d_t + \epsilon \mathbf{1})^{\circ 1/2} \oslash (s_t + \epsilon \mathbf{1})^{\circ 1/2} \odot g_{\mathbb{B}_t}(\theta_t)$$

where

$$s_t = \rho s_{t-1} + (1 - \rho)(g_{\mathbb{B}_t}(\theta_t))^{\circ 2}$$

$$d_t = \rho d_{t-1} + (1 - \rho)\Delta_{t-1}^{\circ 2}$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{decay rate } \rho \in [0; 1)$$

$$\text{divide-by-zero safe guard } \epsilon \in \mathbb{R}^+$$

$$\text{initial } s_0, d_0, \Delta_0 \in \mathbb{R}^D$$

**Update Rule 3.7 (Adam [87])**

$$\theta_{t+1} = \theta_t - \eta \hat{m}_t \oslash (\hat{v}_t^{\circ 1/2} + \epsilon \mathbf{1})$$

where

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

where

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_{\mathbb{B}_t}(\theta_t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

where

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2)(g_{\mathbb{B}_t}(\theta_t))^{\circ 2}$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{decay rates } \beta_1, \beta_2 \in [0; 1)$$

$$\text{divide-by-zero safe guard } \epsilon \in \mathbb{R}^+$$

$$\text{initial } m_0, v_0 \in \mathbb{R}^D$$

**Update Rule 3.8 (AMSGrad [132])**

$$\theta_{t+1} = \theta_t - \eta \hat{m}_t \oslash (\hat{v}_t^{\circ 1/2} + \epsilon \mathbf{1})$$

where

$$\hat{m}_t = m_t / (1 - \beta_1^t)$$

where

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1)g_{\mathbb{B}_t}(\theta_t)$$

$$\hat{v}_t = v_t / (1 - \beta_2^t)$$

where

$$v_t = \max(\beta_2 v_{t-1} + (1 - \beta_2)(g_{\mathbb{B}_t}(\theta_t))^{\circ 2}, v_{t-1})$$

with

$$\text{learning rate } \eta \in \mathbb{R}^+$$

$$\text{decay rates } \beta_1, \beta_2 \in [0; 1)$$

$$\text{divide-by-zero safe guard } \epsilon \in \mathbb{R}^+$$

$$\text{initial } m_0, v_0 \in \mathbb{R}^D$$

**Figure 3.1:** Popular deep learning optimizers rely on the average mini-batch gradient. At iteration  $t$ , they incorporate information in form of the mini-batch average gradient  $g_{\mathbb{B}_t}$ . This is a representative subset of popular methods; see [145] for a more complete overview.

### Adaptive Methods (Per-parameter Learning Rate)

Adaptive methods like AdaGrad [44], RMSProp [164], Adadelta [180], Adam [87] & AMSGrad [132] keep a learning rate for each parameter that is adapted over time. This rescaling is driven by the elementwise square of the mini-batch gradient,  $\mathbf{g}_{\mathbb{B}_t}(\boldsymbol{\theta}_t)^{\odot 2}$  ([Update Rules 3.4 to 3.8](#)). These algorithms use this quantities in different ways, motivated by shortcomings of predecessors (e.g. too aggressive learning rate shrinking [164, 180] and convergence problems [132]).

### No Clear Winner & Mini-batch Gradient as Central Object

[Figure 3.1](#) summarizes the update rules of the above optimization methods. Their update rules to derive  $\boldsymbol{\theta}_{t+1}$  rely on three ingredients:

- ▶ The history of iterates  $\{\boldsymbol{\theta}_{t'}\}_{t'=0}^t$
- ▶ The history of gradients<sup>1</sup>  $\{\mathbf{g}_{\mathbb{B}_{t'}}(\boldsymbol{\theta}_{t'})\}_{t'=0}^t$
- ▶ The history of elementwise gradient squares  $\{\mathbf{g}_{\mathbb{B}_{t'}}(\boldsymbol{\theta}_{t'})^{\odot 2}\}_{t'=0}^t$

<sup>1</sup>: One mild exception is Nesterov momentum ([Update Rule 3.3](#)). It relies on the “lookahead” gradient history rather than the gradient history of iterates.

In summary, the main ingredient in all of these methods is the average mini-batch gradient (and its elementwise square). Their updates are cheap and only require computation of average mini-batch gradients, as efficiently provided by AD in popular ML libraries ([Section 2.3](#)):

$$\begin{aligned} \mathcal{M} &= \mathcal{M}(H_t, \phi_t) \\ \text{with } H_t &= \{\boldsymbol{\theta}_{t'}\}_{t'=0}^t, \cup \left\{ \mathbf{g}_{\mathbb{B}_{t'}}(\boldsymbol{\theta}_{t'}) \right\}_{t'=0}^t, \cup \left\{ \mathbf{g}_{\mathbb{B}_{t'}}(\boldsymbol{\theta}_{t'})^{\odot 2} \right\}_{t'=0}^t, \end{aligned} \quad (3.2)$$

While the mentioned methods are a representative subset of popular algorithms, there are more than one hundred algorithms with structurally similar update rules (see Table 2 in [145] and references therein). Some adaptive methods like Adam and RMSProp contain simpler methods like SGD and Momentum as special cases [34]. Therefore, they should, in principle, be able to perform better. However, recent work that compares deep learning optimizers through benchmarks to identify the best-performing method finds that the currently existing methods perform quite similarly [145].

One reason why newly developed optimizers of the structure [Equation \(3.2\)](#) seem to not be able to clearly improve over existing methods could be that their update rules are constrained to using the average gradient, as it is readily available in software. To overcome these limitations, it might be helpful to study the potential of methods that leverage information beyond the gradient. Second-order optimization incorporate curvature information and represent the state-of-the art in “classical” optimization problems (convex optimization, generalized linear models). Noise-reduction methods focus on improving the gradient estimator  $\mathbf{g}_{\mathbb{B}_t}$  by reducing its variance through per-sample information. The rest of this chapter provides definitions for the information used by such methods, and provides additional motivation for them by showcasing applications outside optimization.

## 3.2 Second-order Optimization

Second-order methods iteratively optimize an objective  $\mathcal{L}(\boldsymbol{\theta}) : \mathbb{R}^D \rightarrow \mathbb{R}$  using a local quadratic approximation  $\mathcal{L}(\boldsymbol{\theta}) \approx m_{\boldsymbol{\theta}_t}(\boldsymbol{\theta})$  around the current point  $\boldsymbol{\theta}_t$  after  $t$  iterations,

$$m_{\boldsymbol{\theta}_t}(\boldsymbol{\theta}) = a(\boldsymbol{\theta}_t) + \mathbf{b}(\boldsymbol{\theta}_t)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^\top \mathbf{C}(\boldsymbol{\theta}_t)(\boldsymbol{\theta} - \boldsymbol{\theta}_t), \quad (3.3a)$$

with an offset  $a(\boldsymbol{\theta}_t) \in \mathbb{R}$ , a slope vector  $\mathbf{b}(\boldsymbol{\theta}_t) \in \mathbb{R}^D$ , and a curvature matrix  $\mathbf{C}(\boldsymbol{\theta}_t) \in \mathbb{R}^{D \times D}$ . As described in [Section 3.1](#), the next iterate  $\boldsymbol{\theta}_{t+1}$  is obtained by minimizing the local approximation. For this proxy to possess a minimum, the curvature matrix  $\mathbf{C}(\boldsymbol{\theta}_t)$  must be PD. Then, [Equation \(3.3a\)](#) is minimized by

$$\boldsymbol{\theta}_{t+1} = \arg \min_{\boldsymbol{\theta}} m_{\boldsymbol{\theta}_t}(\boldsymbol{\theta}) = \boldsymbol{\theta}_t - \mathbf{C}(\boldsymbol{\theta}_t)^{-1} \mathbf{b}(\boldsymbol{\theta}_t). \quad (3.3b)$$

This update is computationally challenging, because the size of the curvature matrix is quadratic in  $D$  and generally infeasible to store. Additionally, the computational complexity of matrix inversion scales cubically in  $D$ . These problems can somewhat be addressed by approximately solving the linear system

$$\mathbf{C}(\boldsymbol{\theta}_t) \boldsymbol{\theta}_{t+1} = -\mathbf{b}(\boldsymbol{\theta}_t) \quad (3.3c)$$

for  $\boldsymbol{\theta}_{t+1}$ . This can be done with iterative solvers, such as CG, which only require matrix-vector products with  $\mathbf{C}(\boldsymbol{\theta}_t)$ , which can often be implemented without expanding the matrix representation in memory.

The first-order methods from [Section 3.1](#) circumvent these issues by using a diagonal—and often quite crude—curvature approximation that is cheap to store and invert (see [\[34\]](#) for an overview). *E.g.*, SGD's local approximation, [Equation \(3.1b\)](#), uses  $1/2\eta\mathbf{I}$  as curvature matrix. The following sections introduce common curvature matrices for local approximations of empirical risks.

### 3.2.1 Newton's Method & the Hessian Matrix

The Taylor series provides a meaningful local approximation of an analytic function. Its expansion up to second-order is

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}) &= \underbrace{\mathcal{L}(\boldsymbol{\theta}_t)}_{a(\boldsymbol{\theta}_t)} + \underbrace{\nabla_{\boldsymbol{\theta}_t} \mathcal{L}(\boldsymbol{\theta}_t)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_t)}_{\mathbf{b}_{\boldsymbol{\theta}_t}} \\ &\quad + \underbrace{\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^\top \nabla_{\boldsymbol{\theta}_t}^2 \mathcal{L}(\boldsymbol{\theta}_t)(\boldsymbol{\theta} - \boldsymbol{\theta}_t)}_{\mathbf{C}(\boldsymbol{\theta}_t)} + \mathcal{O}\left((\boldsymbol{\theta} - \boldsymbol{\theta}_t)^3\right) \end{aligned} \quad (3.4)$$

where  $\mathcal{O}((\boldsymbol{\theta} - \boldsymbol{\theta}_t)^3)$  denotes polynomial terms in the components of  $\boldsymbol{\theta} - \boldsymbol{\theta}_t$  of cubic order and above. This reveals the objective's Hessian  $\nabla_{\boldsymbol{\theta}_t}^2 \mathcal{L}(\boldsymbol{\theta}_t)$ , which collects its second-order partial derivatives in a matrix  $[\nabla_{\boldsymbol{\theta}_t}^2 \mathcal{L}(\boldsymbol{\theta}_t)]_{i,j} = \partial^2 \mathcal{L}(\boldsymbol{\theta}_t)/\partial [\boldsymbol{\theta}_t]_i \partial [\boldsymbol{\theta}_t]_j$ , as curvature matrix.

**Definition 3.1 (Hessian)** Let  $b : \mathbb{R}^D \rightarrow \mathbb{R}$ ;  $\mathbf{a} \mapsto b(\mathbf{a})$  be a differentiable vector-to-scalar function. The Hessian  $\nabla_{\mathbf{a}}^2 b \in \mathbb{R}^{D \times D}$  of  $b$  w.r.t.  $\mathbf{a}$  is a symmetric matrix containing the second-order partial derivatives

$$\nabla_{\mathbf{a}}^2 b = \frac{\partial^2 b}{\partial \mathbf{a} \partial \mathbf{a}^\top} \quad \text{with} \quad [\nabla_{\mathbf{a}}^2 b]_{i,j} = \frac{\partial^2 b}{\partial a_i \partial a_j} \quad (3.5)$$

**Definition 3.2 (Generalized Hessian)** Let  $B : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}^{m \times p}$  be a twice differentiable matrix function. The Hessian  $\nabla_A^2 B(A)$  is an  $(mnpq \times nq)$  matrix defined by

$$\begin{aligned} \nabla_A^2 B(A) &= J_A [J_A B(A)]^\top \\ &= \frac{\partial}{\partial (\text{vec } A)^\top} \text{vec} \left\{ \left[ \frac{\partial \text{vec } B(A)}{\partial (\text{vec } A)^\top} \right]^\top \right\} \end{aligned} \quad (3.6)$$

[103, Chapter 10.2] with flattening defined by [Definition 2.2](#). In element-wise notation, this is the matrix-stack of all output component Hessians

$$\nabla_A^2 B(A) = \begin{pmatrix} \nabla_{\text{vec } A}^2 [\text{vec } B]_1 \\ \nabla_{\text{vec } A}^2 [\text{vec } B]_2 \\ \vdots \\ \nabla_{\text{vec } A}^2 [\text{vec } B]_{np} \end{pmatrix},$$

with the Hessian from [Definition 3.1](#).

The tensor case is analogous but requires cluttered notation and is therefore omitted. Common forms for neural networks include vector-to-vector functions  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ ,  $\mathbf{x} \mapsto f(\mathbf{x})$  with

$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) = \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x}^\top \partial \mathbf{x}},$$

where  $\mathbf{x}$  can be considered the input or bias vector if a linear layer. In details,

$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) = \begin{pmatrix} \nabla_{\mathbf{x}}^2 f_1(\mathbf{x}) \\ \vdots \\ \nabla_{\mathbf{x}}^2 f_m(\mathbf{x}) \end{pmatrix}. \quad (3.7)$$

Others are matrix-to-vector mappings  $f : \mathbb{R}^{n \times q} \rightarrow \mathbb{R}^m$ ,  $\mathbf{X} \mapsto f(\mathbf{X})$  with

$$\nabla_{\mathbf{X}}^2 f(\mathbf{X}) = \frac{\partial}{\partial (\text{vec } \mathbf{X})^\top} \frac{\partial f(\mathbf{X})}{\partial \text{vec } \mathbf{X}},$$

e.g. with  $\mathbf{X}$  the weight of a linear layer.

Newton's method uses the Hessian as curvature matrix.

**Update Rule 3.9 (Newton's method (simplified))** A Newton step is

$$\theta_{t+1} = \theta_t - H_{\mathbb{D}}(\theta_t)^{-1} g_{\mathbb{D}}(\theta_t) \quad (3.8)$$

with the gradient and Hessian from [Definitions 2.4](#) and [3.1](#) (practical implementations vary and often introduce additional hyperparameters such as a learning rate, damping term, mini-batch size, etc.).

Matrix-free multiplication with the Hessian [127] can be combined with CG to compute Newton steps via solving the linear system [Equation \(3.3c\)](#). This idea is known as Hessian-free optimization [106].

While the Taylor expansion motivates using the Hessian as curvature matrix, it leads to problems for non-convex functions like the empirical risk, as it is in general indefinite. Therefore [Equation \(3.3b\)](#) does not have a solution. In practice, PSD curvature matrices that are approximations of the Hessian are popular substitutes for the Hessian to avoid this issue ([Sections 3.2.2](#) to [3.2.4](#)).

### 3.2.2 The Gauss-Newton Method & Matrix

The Gauss-Newton (GN) method [22, chapter 6.3] tackles the nonlinear least squares regression task ([Example 2.1](#)) of fitting a function  $f_{\theta}$  to data in the form of vector-valued inputs  $\mathbf{x}_n \in \mathbb{R}^M$  and scalar-valued outputs  $y_n \in \mathbb{R}$  by minimizing the mean squared error,

$$\underset{\theta}{\text{minimize}} \mathcal{L}(\theta), \quad \text{where} \quad \mathcal{L}(\theta) = \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, y_n) \in \mathbb{D}} \underbrace{(f_{\theta}(\mathbf{x}_n) - y_n)^2}_{:= f_n}, \quad (3.9)$$

The objective's gradient is

$$\nabla_{\theta} \mathcal{L}(\theta) = \frac{2}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, y_n) \in \mathbb{D}} (\mathbf{J}_{\theta} f_n)^\top (f_n - y_n), \quad (3.10)$$

with the Jacobian  $\mathbf{J}_\theta f_n \in \mathbb{R}^{1 \times D}$  ([Definition 2.4](#)). The Hessian is

$$\nabla_\theta^2 \mathcal{L}(\theta) = \frac{2}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} (\mathbf{J}_\theta f_n)^\top \mathbf{J}_\theta f_n + \frac{2}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \nabla_\theta^2 f_n (f_n - y_n). \quad (3.11a)$$

The first term is the PSD *Gauss-Newton matrix* (up to scaling)

$$G_{\mathbb{D}}(\theta) = \frac{2}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} (\mathbf{J}_\theta f_n)^\top \mathbf{J}_\theta f_n, \quad (3.11b)$$

and the second term is the residual matrix

$$R_{\mathbb{D}}(\theta) = \frac{2}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \nabla_\theta^2 f_n (f_n - y_n). \quad (3.11c)$$

The *Gauss-Newton matrix* approximates the Hessian through first-order information of the model,  $(\mathbf{J}_\theta f_n)^\top = \nabla_\theta f_n$ , which is cheap to compute. For vanishing residual terms  $\lim_{(f_n - y_n) \rightarrow 0}$  (as the model predictions match the labels), or linear models ( $\nabla_\theta^2 f_\theta = \mathbf{0}$ ) it corresponds to the Hessian.

**Update Rule 3.10 (Gauss-Newton method (simplified))** For the nonlinear least squares task [Equation \(3.9\)](#), a Gauss-Newton step is<sup>2</sup>

$$\theta_{t+1} = \theta_t - G_{\mathbb{D}}(\theta_t)^{-1} \nabla_{\theta_t} \mathcal{L}(\theta_t), \quad (3.12)$$

with the gradient and Gauss-Newton matrix from [Definition 2.4](#) and [Equation \(3.11b\)](#) (practical implementations vary and often introduce additional hyperparameters such as a learning rate, damping term, mini-batch size, etc.).

2: Stacking Jacobians and residuals,

$$\begin{aligned} J(\theta) &= \begin{pmatrix} \mathbf{J}_\theta f_1 \\ \mathbf{J}_\theta f_2 \\ \vdots \\ \mathbf{J}_\theta f_{|\mathbb{D}|} \end{pmatrix} \in \mathbb{R}^{|\mathbb{D}| \times D}, \\ r(\theta) &= \begin{pmatrix} f_1 - y_1 \\ f_2 - y_2 \\ \vdots \\ f_{|\mathbb{D}|} - y_{|\mathbb{D}|} \end{pmatrix} \in \mathbb{R}^{|\mathbb{D}|}, \end{aligned}$$

absorbs the sums into matrix multiplies

$$\theta_{t+1} = \theta_t - (J(\theta_t)^\top J(\theta_t))^{-1} J(\theta_t)^\top r(\theta_t)$$

which can be solved through

$$J(\theta_t)^\top J(\theta_t)x = -J(\theta_t)^\top r(\theta_t)$$

via JVPs & VJPs in combination with CG.

### 3.2.3 The Generalized Gauss-Newton Matrix

The *generalized Gauss-Newton (GGN) matrix* is a PSD approximation to the Hessian that generalizes the GN through abstraction via empirical risk minimization (see [Section 2.1.1](#)). The GGN can be understood through different perspectives, and, using the probabilistic interpretation of the model, is related to the natural gradient method through its connections to the Fisher ([Section 3.2.4](#)).

#### From Gauss-Newton to Generalized Gauss-Newton

The GN matrix from [Section 3.2.3](#) stems from a nonlinear least squares problem that can be viewed as supervised regression ([Example 2.1](#)), i.e. mean squared error loss function, with scalar-valued labels ( $C = 1$ ). For the general case of empirical risk minimization ([Equation \(2.4b\)](#)), the Hessian decomposes due to the split between model and loss function,

$$\ell(\theta) = \ell(\cdot, y) \circ f_\theta(x) \quad (3.13a)$$

as a result of the Hessian chain rule [Theorem 3.1](#). For a single datum  $(x, y)$  with prediction  $f := f_\theta(x)$ , this yields

$$\nabla_\theta^2 \ell(\theta) = (\mathbf{J}_\theta f)^\top \left[ \nabla_f^2 \ell(f, y) \right] \mathbf{J}_\theta f + \sum_{c=1}^C (\nabla_\theta^2 [f]_c) \nabla_f \ell(f, y) \quad (3.13b)$$

The arrangement of partial derivatives in the generalizations of Jacobian [Definition 2.5](#) and Hessian [Definition 3.2](#) implies the following chain rule generalization for second-order derivatives:

**Theorem 3.1 (Chain rule for the generalized Hessian)** Let  $b : \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $c : \mathbb{R}^m \rightarrow \mathbb{R}^p$  be twice differentiable and  $d = c \circ b : \mathbb{R}^n \rightarrow \mathbb{R}^p$ ,  $a \mapsto d(a) = c(b(a))$ . The relation between the Hessian of  $d$  and the Jacobians and Hessians of the constituents  $c$  and  $b$  is given by

$$\begin{aligned} \nabla_a^2 d(a) &= [I_p \otimes \mathbf{J}_a b(a)]^\top [\nabla_b^2 c(b)] \mathbf{J}_a b(a) \\ &\quad + [\mathbf{J}_b c(b) \otimes I_n] \nabla_a^2 b(a) \end{aligned} \quad (3.14)$$

[restricted from [103](#), Chapter 6.10].

The matrix/vector case is analogous. Matrix differential calculus [\[103\]](#) is a useful tool to easily read off the Hessian from specific expressions, see [Appendix A.1](#). For applications of the Hessian chain rule, see [Chapter 4](#).

The first term carries curvature information of the loss function, while the second term contains curvature information of the model. Because  $\ell(f, y)$  is convex in  $f$ , the first term is PSD, whereas the second term is indefinite in general.

The GGN is the first term and neglects curvature information of the model. For the empirical risk [Equation \(2.4b\)](#), and using the shorthand  $f_n := f_\theta(x_n)$ , the Hessian is

$$H_{\mathbb{D}}(\theta) = \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \nabla_\theta^2 \ell(f_n, y_n) = G_{\mathbb{D}}(\theta) + R_{\mathbb{D}}(\theta) \quad (3.15a)$$

with the *generalized Gauss-Newton matrix*

$$G_{\mathbb{D}}(\theta) := \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} (\mathbf{J}_\theta f_n)^\top \left[ \nabla_{f_n}^2 \ell(f_n, y_n) \right] \mathbf{J}_\theta f_n \quad (3.15b)$$

and the *residual matrix*

$$R_{\mathbb{D}}(\theta) := \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \sum_{c=1}^C (\nabla_\theta^2 [f_n]_c) \nabla_{f_n} \ell(f_n, y_n). \quad (3.15c)$$

This decomposition reduces to [Equation \(3.11\)](#) for nonlinear least squares, where  $\nabla_f^2(x) \ell(f, y) = 2/c I$  ([Table 4.1](#)),  $\nabla_f \ell(f, y) = 2/c(f - y)$  ([Table 2.2](#)), and  $C = 1$ . Therefore, the GGN is a generalization of the GN via the chain rule applied to the model-loss function split.

### From Linearization to Generalized Gauss-Newton

Alternatively, one can replace the model in  $\ell(f_\theta(x), y)$  with a linear Taylor expansion around  $\theta$ ,

$$f_\theta(x) \leftrightarrow \hat{f}_{\theta'}(x) = f_\theta(x) + [\mathbf{J}_\theta f_\theta(x)] (\theta' - \theta).$$

This eliminates second-order terms in the model, i.e.  $\nabla_{\theta'}^2 \hat{f}_{\theta'}(x) = \mathbf{0}$ . Application of the Hessian chain rule to the loss with a linearized model,

$$\hat{\ell}(\theta') = \ell(\cdot, y) \circ \hat{f}_{\theta'}(x), \quad (3.16a)$$

and using the shorthand  $\hat{f}_{\theta'}(x) := \hat{f}$ , yields

$$\nabla_{\theta'}^2 \hat{\ell}(\theta') = \left( \mathbf{J}_{\theta'} \hat{f} \right)^\top \left[ \nabla_{\hat{f}}^2 \ell(\hat{f}, y) \right] \mathbf{J}_{\theta'} \hat{f}. \quad (3.16b)$$

At the expansion point where model predictions and Jacobians match, this yields the first term of the decomposition [Equation \(3.13b\)](#)

$$\left. \left( \nabla_{\theta'}^2 \ell(\hat{f}_{\theta'}(x), y) \right) \right|_{\theta'=\theta} = (\mathbf{J}_\theta f)^\top \left[ \nabla_f^2 \ell(f, y) \right] \mathbf{J}_\theta f. \quad (3.16c)$$

This carries over the empirical risk  $\hat{\mathcal{L}}(\theta') = 1/|\mathbb{D}| \sum_{(x_n, y_n) \in \mathbb{D}} \hat{\ell}(\theta')$  under a linearized model, whose Hessian is the GGN from [Equation \(3.15b\)](#),

$$\left. \left( \nabla_{\theta'}^2 \hat{\mathcal{L}}(\theta') \right) \right|_{\theta'=\theta} = G_{\mathbb{D}}(\theta) \quad (3.16d)$$

Hence, the GGN is a Hessian approximation that neglects curvature from the model, and becomes equivalent to the Hessian for linear models.

### 3.2.4 Natural Gradient Descent & the Fisher

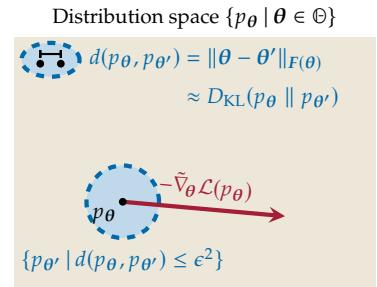
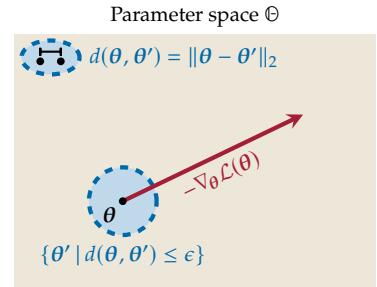
Natural gradient descent [5, NGD] uses the Fisher information matrix as curvature matrix. The natural gradient provides the direction of steepest in the space of probability distributions described by a model ([Figure 3.2](#), recall the probabilistic interpretation of certain empirical risks from [Section 2.1.3](#)). The Fisher is connected to the GGN from [Section 3.2.3](#) and for most models used in modern ML, is equivalent but provides a probabilistic interpretation.

#### Steepest Descent Direction & Notion of Distance

**Concept (steepest descent direction):** Consider an arbitrary objective function  $\mathcal{L} : \Theta \rightarrow \mathbb{R}$  with  $\Theta = \mathbb{R}^D$ . At a location  $\theta \in \Theta$ , the steepest descent direction is the direction in which the objective increases at the fastest rate, *i.e.* per infinitesimally small distance moved. Formally,

$$\Delta\theta^{(\text{steepest})} = \lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{\substack{\Delta\theta \\ d(\theta, \theta + \Delta\theta) \leq \epsilon}} \mathcal{L}(\theta + \Delta\theta) \quad (3.17)$$

and depends on the distance measure  $d(\cdot, \cdot)$  between elements in a small neighborhood around  $\theta$ .



**Gradient descent as steepest descent:** Using the Euclidean 2-norm to measure distances in  $\Theta$  via  $d(\theta_1, \theta_2) = \|\theta_1 - \theta_2\|_2$ , the steepest descent direction points along the negative gradient  $-\nabla_\theta \mathcal{L}(\theta)$  [[107](#), Chapter 6],

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{\substack{\Delta\theta \\ \|\Delta\theta\|_2 \leq \epsilon}} \mathcal{L}(\theta + \Delta\theta) = \frac{-\nabla_\theta \mathcal{L}(\theta)}{\|\nabla_\theta \mathcal{L}(\theta)\|_2}. \quad (3.18)$$

**Natural gradient as steepest descent:** [Section 2.1.3](#) showed that in many tasks, such as regression ([Example 2.1](#)) and softmax cross-entropy classification ([Example 2.2](#)), the parameters  $\theta$  model a probability distribution  $p_\theta(z)$  over a random variable  $z \in \Omega$ . One could therefore establish a different notion of distance by comparing probability distributions. The KL divergence is a similarity measure between densities, but it is not a proper metric; *e.g.* it is not symmetric in its arguments. The steepest descent direction only requires measuring distances within an infinitesimal ball though. For this purpose, a metric—described by the Fisher—can be established through Taylor expansion of the KL divergence.

**Figure 3.2:** Gradient descent and natural gradient descent via steepest descent. Gradient descent (*top*) follows the direction of steepest descent in Euclidean parameter space. NGD (*bottom*) considers the space of distributions, where local distances are measured via a quadratic expansion of the KL divergence, which gives rise to the Fisher. Details in the text. Figure inspired by [[107](#)].

Starting from the KL divergence between two infinitesimally close distributions  $p_{\theta}(z), p_{\theta+\Delta\theta}(z)$ ,

$$D_{\text{KL}}(p_{\theta+\Delta\theta} \parallel p_{\theta}) = \int_{\Omega} p_{\theta}(z) [\log p_{\theta}(z) - \log p_{\theta+\Delta\theta}(z)] dz \quad (3.19)$$

the first step is to Taylor-expand the logarithm around  $\theta$ ,

$$\begin{aligned} \log p_{\theta+\Delta\theta}(z) &= \log p_{\theta}(z) + (\Delta\theta)^{\top} \nabla_{\theta} \log p_{\theta}(z) \\ &\quad + \frac{1}{2} (\Delta\theta)^{\top} \nabla_{\theta}^2 \log p_{\theta}(z) (\Delta\theta) + \mathcal{O}((\Delta\theta)^3). \end{aligned}$$

Inserting this into Equation (3.19) results in

$$\begin{aligned} D_{\text{KL}}(p_{\theta}(z) \parallel p_{\theta+\Delta\theta}(z)) &= -\Delta\theta^{\top} \underbrace{\int_{\Omega} p_{\theta}(z) \nabla_{\theta} \log p_{\theta}(z) dz}_{=0, \text{ see Remark 3.1}} \\ &\quad - \frac{1}{2} \Delta\theta^{\top} \left( \int_{\Omega} p_{\theta}(z) \nabla_{\theta}^2 \log p_{\theta}(z) dz \right) \Delta\theta \\ &\quad + \mathcal{O}((\Delta\theta)^3) \end{aligned}$$

The Hessian in the second term is expressed as (Remark 3.2)

$$\nabla_{\theta}^2 \log p_{\theta}(z) = -\nabla_{\theta} \log p_{\theta}(z) (\nabla_{\theta} \log p_{\theta}(z))^{\top} + \frac{1}{p_{\theta}(z)} \nabla_{\theta}^2 p_{\theta}(z), \quad (3.20)$$

whose second term again vanishes in expectation (Remark 3.3). Hence,

$$D_{\text{KL}}(p_{\theta}(z) \parallel p_{\theta+\Delta\theta}(z)) = \frac{1}{2} \Delta\theta^{\top} F(\theta) \Delta\theta + \mathcal{O}((\Delta\theta)^3) \quad (3.21a)$$

with the two equivalent forms of the Fisher

$$F(\theta) = - \int_{\Omega} p_{\theta}(z) \nabla_{\theta}^2 \log p_{\theta}(z) dz \quad (3.21b)$$

$$= \int_{\Omega} p_{\theta}(z) (\nabla_{\theta} \log p_{\theta}(z)) (\nabla_{\theta} \log p_{\theta}(z))^{\top} dz \quad (3.21c)$$

Equation (3.21b) will be helpful to draw connections to the Hessian, and Equation (3.21c) provides links to the GGN.

**Remark 3.1 (The log-probability's  $\theta$ -gradient vanishes in expectation)**

$$\begin{aligned} - \int_{\Omega} p_{\theta}(z) \nabla_{\theta} \log p_{\theta}(z) dz \\ &= - \int_{\Omega} p_{\theta}(z) \frac{\nabla_{\theta} p_{\theta}(z)}{p_{\theta}(z)} dz \\ &= -\nabla_{\theta} \left( \int_{\Omega} p_{\theta}(z) dz \right) \\ &= -\nabla_{\theta} 1 = 0 \end{aligned}$$

**Remark 3.2 (Decomposition of the log-probability's  $\theta$ -Hessian)** Consider element  $(i, j)$  of the Hessian,

$$\begin{aligned} \frac{\partial^2 \log p_{\theta}(z)}{\partial \theta_i \partial \theta_j} \\ &= \frac{\partial}{\partial \theta_i} \left( \frac{1}{p_{\theta}(z)} \frac{\partial p_{\theta}(z)}{\partial \theta_j} \right) \\ &= -\frac{1}{p_{\theta}(z)^2} \frac{\partial p_{\theta}(z)}{\partial \theta_j} \frac{\partial p_{\theta}(z)}{\partial \theta_i} \\ &\quad + \frac{1}{p_{\theta}(z)} \frac{\partial^2 p_{\theta}(z)}{\partial \theta_i \partial \theta_j} \\ &= -\frac{\partial \log p_{\theta}(z)}{\partial \theta_j} \frac{\partial \log p_{\theta}(z)}{\partial \theta_i} \\ &\quad + \frac{1}{p_{\theta}(z)} \frac{\partial^2 p_{\theta}(z)}{\partial \theta_i \partial \theta_j} \end{aligned}$$

which, in vector notation, translates into Equation (3.20).

**Remark 3.3 (Hessian of the model distribution vanishes in expectation)**

$$\begin{aligned} \int_{\Omega} p_{\theta}(z) \frac{1}{p_{\theta}(z)} \nabla_{\theta}^2 p_{\theta}(z) dz \\ &= \nabla_{\theta}^2 \left( \int_{\Omega} p_{\theta}(z) dz \right) \\ &= \nabla_{\theta}^2 1 = 0 \end{aligned}$$

Locally, the KL divergence Equation (3.21a) gives rise to a metric induced by the Fisher-norm  $\|\cdot\|_{F(\theta)}$ ,

$$d(p_{\theta}, p_{\theta+\Delta\theta}) = \|\Delta\theta\|_{F(\theta)} := \sqrt{(\Delta\theta)^{\top} F(\theta) \Delta\theta} \quad (3.22)$$

Ollivier et al. [120] show that the steepest descent for a function  $\mathcal{L}(p_{\theta})$  points along the negative natural gradient  $-\tilde{\nabla}_{\theta} \mathcal{L}(p_{\theta}) := -F(\theta)^{-1} \nabla_{\theta} \mathcal{L}(p_{\theta})$ .

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} \arg \min_{\substack{\Delta\theta \\ \|\Delta\theta\|_{F(\theta)} \leq \epsilon^2/2}} \mathcal{L}(p_{\theta+\Delta\theta}) = -\frac{\tilde{\nabla}_{\theta} \mathcal{L}(p_{\theta})}{\|\tilde{\nabla}_{\theta} \mathcal{L}(p_{\theta})\|_{F(\theta)^{-1}}}. \quad (3.23)$$

### Natural Gradient Descent & Fisher in Empirical Risk Minimization

Section 2.1.3 presented connections of empirical risk minimization to learning  $p_{\text{data}}(\mathbf{x}, \mathbf{y})$  via  $p_{\theta}(\mathbf{x}, \mathbf{y}) = p_{\theta}(\mathbf{y} | \mathbf{x})p(\mathbf{x}) = q(\mathbf{y} | f)p(\mathbf{x})$  with a negative log-likelihood loss  $-\log q(\mathbf{y} | f) = \ell(f, \mathbf{y})$  and  $f_{\theta}(\mathbf{x}) := f$ . In these cases, the empirical risk Equation (2.4b) depends on a probability distribution, and the distribution space's geometry gives rise to NGD with  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$  in the Fisher Equation (3.21). After simplifying the  $\theta$ -derivatives and grouping dependencies of  $\mathbf{x}$  and  $\mathbf{y}$ , the Fisher reads

$$\begin{aligned} F(\theta) &= \int_{\mathbb{X}} p(\mathbf{x}) \left( - \int_{\mathbb{Y}} p_{\theta}(\mathbf{y} | \mathbf{x}) \nabla_{\theta}^2 \log p_{\theta}(\mathbf{y} | \mathbf{x}) d\mathbf{y} \right) d\mathbf{x} \\ &= \int_{\mathbb{X}} p(\mathbf{x}) \left( \int_{\mathbb{Y}} p_{\theta}(\mathbf{y} | \mathbf{x}) (\nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{x})) (\nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{x}))^\top d\mathbf{y} \right) d\mathbf{x}, \end{aligned}$$

or in short form

$$\begin{aligned} F(\theta) &= \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y} | \mathbf{x})} [-\nabla_{\theta}^2 \log p_{\theta}(\mathbf{y} | \mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y} | \mathbf{x})} [\nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{x}) (\nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{x}))^\top] \quad (3.24) \end{aligned}$$

Next, replace  $p(\mathbf{x}) \leftrightarrow p_{\mathbb{D}}(\mathbf{x})$  by empirical approximation through data,

$$\begin{aligned} F_{\mathbb{D}}(\theta) &= \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y} | \mathbf{x}_n)} [-\nabla_{\theta}^2 \log p_{\theta}(\mathbf{y} | \mathbf{x}_n)] \\ &= \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} \mathbb{E}_{\mathbf{y} \sim p_{\theta}(\mathbf{y} | \mathbf{x}_n)} [\nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{x}_n) (\nabla_{\theta} \log p_{\theta}(\mathbf{y} | \mathbf{x}_n))^\top] \quad (3.25) \end{aligned}$$

(note that the Fisher is independent of the labels  $\{\mathbf{y}_n\}$ !). Using the relation between log-likelihood and loss function leads to (with  $f_n := f_{\theta}(\mathbf{x}_n)$ )

$$\begin{aligned} F_{\mathbb{D}}(\theta) &= \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} | f_n)} [\nabla_{\theta}^2 \ell(f_n, \mathbf{y})] \\ &= \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} | f_n)} [\nabla_{\theta} \ell(f_n, \mathbf{y}) (\nabla_{\theta} \ell(f_n, \mathbf{y}))^\top]. \end{aligned}$$

The second equality views the Fisher as expected gradient outer product. With the Jacobian chain rule Theorem 2.2 applied to  $\ell \circ f_{\theta}$ , one obtains

$$\begin{aligned} F_{\mathbb{D}}(\theta) &= \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} | f_n)} [(\mathbb{J}_{\theta} f_n)^\top \nabla_{f_n} \ell(f_n, \mathbf{y}) (\mathbb{J}_{\theta} f_n) \nabla_{f_n} \ell(f_n, \mathbf{y})^\top] \\ &= \frac{1}{|\mathbb{D}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}} (\mathbb{J}_{\theta} f_n)^\top \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} | f_n)} [\nabla_{f_n} \ell(f_n, \mathbf{y}) (\nabla_{f_n} \ell(f_n, \mathbf{y}))^\top] \mathbb{J}_{\theta} f_n \quad (3.26) \end{aligned}$$

Equation (3.26) offers an interesting perspective to approximate the Fisher via MC sampling through computing gradients of the loss on targets drawn from  $q$  and will be used later, e.g. Chapters 5 and 7 and Appendix B.1.

3: It is important to stress the expectation over the model's likelihood  $\mathbf{y} \mid \mathbf{x}$ , and *not* over the empirical data. It may be tempting to replace the expectation over  $\mathbf{y}$  with the empirical distribution  $p_{\mathbb{D}}(\mathbf{y})$ . This leads to the empirical Fisher (Section 3.2.5), which is often used instead of the Fisher, but has limitations for applications like optimization [92].

The first equality in Equation (3.25) views the Fisher as an expected Hessian under the model's likelihood<sup>3</sup>. However, in general it does not coincide with the Hessian,  $F_{\mathbb{D}}(\boldsymbol{\theta}) \neq H_{\mathbb{D}}(\boldsymbol{\theta})$ , as  $\mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid f_n)}[\nabla_{\boldsymbol{\theta}}^2 \ell(f_n, \mathbf{y})] \neq \nabla_{\boldsymbol{\theta}}^2 \ell(f_n, \mathbf{y}_n)$ . After applying the Hessian chain rule Theorem 3.1 onto  $\ell \circ f_{\boldsymbol{\theta}}$ ,

$$\begin{aligned} F_{\mathbb{D}}(\boldsymbol{\theta}) &= \frac{1}{|\mathbb{D}|} \sum_{(x_n, \mathbf{y}_n) \in \mathbb{D}} (\mathbf{J}_{\boldsymbol{\theta}} f_n)^T \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid f_n)} [\nabla_{f_n}^2 \ell(f_n, \mathbf{y})] \mathbf{J}_{\boldsymbol{\theta}} f_n \\ &\quad + \frac{1}{|\mathbb{D}|} \sum_{(x_n, \mathbf{y}_n) \in \mathbb{D}} \sum_{c=1}^C \nabla_{\boldsymbol{\theta}}^2([f_n]_c) \underbrace{\mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid f_n)} [\nabla_{f_{\boldsymbol{\theta}}(x_n)} \ell(f_n, \mathbf{y})]}_{=0, \text{ same argument as Remark 3.1}} \\ &= \frac{1}{|\mathbb{D}|} \sum_{(x_n, \mathbf{y}_n) \in \mathbb{D}} (\mathbf{J}_{\boldsymbol{\theta}} f_n)^T \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid f_n)} [\nabla_{f_n}^2 \ell(f_n, \mathbf{y})] \mathbf{J}_{\boldsymbol{\theta}} f_n, \end{aligned} \quad (3.27)$$

one can see connections to the GGN (Equation (3.15b)).

### Connections Between Fisher & GGN

**Update Rule 3.11 (Natural gradient descent (simplified))** For empirical risks with loss functions that have a probabilistic interpretation, an NGD step is

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - F_{\mathbb{D}}(\boldsymbol{\theta}_t)^{-1} g_{\mathbb{D}}(\boldsymbol{\theta}_t)$$

with the gradient from Definition 2.4 and the Fisher from Equation (3.26) or Equation (3.27) (practical implementations vary and often introduce additional hyperparameters such as a learning rate, damping term, mini-batch size, etc.).

4: There are more scenarios in which Fisher and GGN coincide, for instance if  $q(\mathbf{y} \mid f)$  is an exponential family distribution with natural parameters  $f$ , see [107, Chapter 9] for a detailed presentation.

The Fisher (Equation (3.27)) and the GGN (Equation (3.15b)) are structurally similar. Both are identical if the expected Hessian of the loss *w.r.t.* the model's prediction under the model is identical to the empirical Hessian,

$$\nabla_{f_n}^2 \ell(f_n, \mathbf{y}_n) = \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid f_n)} [\nabla_{f_n}^2 \ell(f_n, \mathbf{y})] \implies G_{\mathbb{D}} = F_{\mathbb{D}}. \quad (3.28)$$

For both softmax cross-entropy and square loss, this  $f$ -Hessian is independent of  $\mathbf{y}$ , see Table 4.1. Therefore, the expectation has no effect, and the above equality is satisfied: for least squares regression (Example 2.1) and softmax cross-entropy classification (Example 2.1), the Fisher equals the GGN<sup>4</sup>. NGD (Update Rule 3.11) for those models is thus equivalent to the GGN and can be seen as an approximation of Newton's method (Update Rule 3.9) with the GGN instead of the Hessian.

### 3.2.5 The Gradient Covariance Matrix

The Fisher's form in Equation (3.26) reminds of an uncentered second moment of "would-be" gradients sampled from the likelihood implied by the model [124]. The uncentered gradient covariance matrix

$$\begin{aligned} K_{\mathbb{D}}(\boldsymbol{\theta}) &= \mathbb{E}_{(x, \mathbf{y}) \sim p_{\mathbb{D}}(x, \mathbf{y})} [\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(x), \mathbf{y}) (\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(x), \mathbf{y}))^T] \\ &= \frac{1}{|\mathbb{D}|} \sum_{(x_n, \mathbf{y}_n) \in \mathbb{D}} \nabla_{\boldsymbol{\theta}} \ell(f_n, \mathbf{y}_n) (\nabla_{\boldsymbol{\theta}} \ell(f_n, \mathbf{y}_n))^T \end{aligned} \quad (3.29)$$

if often referred to as empirical Fisher in some of the ML literature, because it follows from replacing the expectation over the model's distribution  $q(\mathbf{y} \mid f_n)$  by the expectation over the empirical data distribution  $p_{\mathbb{D}}(\mathbf{y})$  in the Fisher. However, empirical Fisher and Fisher are, almost always, *not* identical. The importance to distinguish them has been expressed in works like [92] and [162], and it is questionable whether the uncentered gradient covariance Equation (3.29) approximates curvature.

However, an alternative perspective to use  $K_{\mathbb{D}}$  for optimization is to make stochastic gradient-based optimization aware of the noise [9] (sketch in Figure 3.4). With the centered gradient covariance on data  $\mathbb{D}$

$$\begin{aligned}\Sigma_{\mathbb{D}}(\theta) &= \text{Var}_{(x,y) \sim p_{\mathbb{D}}(x,y)} [\nabla_{\theta} \ell(f_{\theta}(x), y)] \\ &= \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} (\nabla_{\theta} \ell(f_n, y_n) - g_{\mathbb{D}}(\theta)) (\nabla_{\theta} \ell(f_n, y_n) - g_{\mathbb{D}}(\theta))^{\top} \\ &= K_{\mathbb{D}}(\theta) - g_{\mathbb{D}}(\theta) g_{\mathbb{D}}(\theta)^{\top},\end{aligned}\quad (3.30)$$

an update step of the form (usually with  $\mathbb{D} = \mathbb{B}_t$ )

$$\theta_{t+1} = \theta_t - \Sigma_{\mathbb{D}}(\theta_t)^{-1} g_{\mathbb{D}}(\theta_t)$$

can be regarded as re-scaling the gradient according to fluctuations: directions that are subject to stronger noise will be shortened more strongly than directions of small noise.

The gradient covariance has also been proposed to adapt the batch size during training [6, 11, 19, 26] and to stop training before overfitting sets in [104]. Similar ideas of variance-adaptation can be found in currently popular deep learning optimizers such as Adam (Update Rule 3.7), that keep an exponential average of the squared mini-batch gradient. This shows up in the diagonal of the gradient covariance (second term)

$$\text{diag}(\Sigma_{\mathbb{D}}(\theta)) = \frac{1}{|\mathbb{D}|} \left( \sum_{(x_n, y_n) \in \mathbb{D}} (\nabla_{\theta} \ell(f_n, y_n))^{\odot 2} \right) - g_{\mathbb{D}}(\theta)^{\odot 2}. \quad (3.31)$$

See [10] for a precise analysis of the connections.

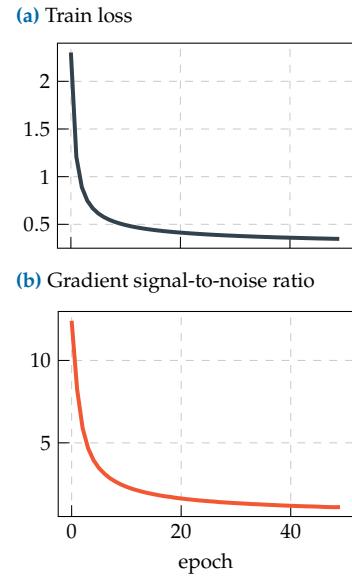
The empirical Fisher is formed by the gradients that make up the mini-batch average gradient  $\nabla_{\theta} \mathcal{L}(\theta)$ , and therefore basically free to compute on top of the gradient. But the required per-sample gradients are difficult to access in popular ML libraries because they are not agnostic to this per-sample structure (see Section 2.3.2). This complicates efficient usage of the gradient covariance matrix for applications like optimization, where performance is key.

### 3.3 Other Applications

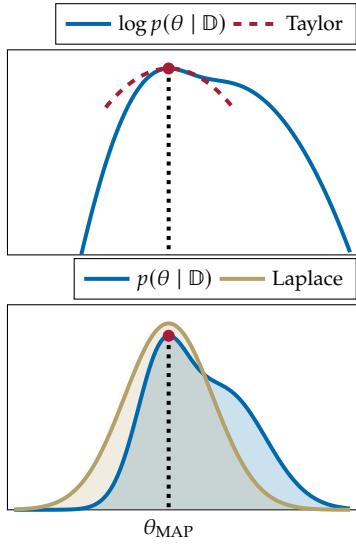
So far, this text showed use cases of higher-order information, in the form of curvature and gradient covariance, for optimization. While optimization is one main component of deep learning, there exist many other applications that require access to richer information. The following briefly describes some applications (mostly) outside the optimization field. Another argument for the utility of higher-order information to better understand the training of neural networks is made by Chapter 6.

#### 3.3.1 Bayesian Deep Learning & Laplace Approximations

Section 2.1.3 described the relation between regularized empirical risk minimization and maximum a posteriori (MAP) estimation. The MAP



**Figure 3.4: Illustration of gradient noise during training.** Evolution of (a) training loss and (b) gradient signal-to-noise ratio (computed with BackPACK, Chapter 5) during training (logistic regression on MNIST). As the loss decreases, the gradient noise increases.



**Figure 3.5:** Conceptual sketch of the Laplace approximation. The goal is to approximate  $p(\theta | \mathbb{D})$  with a Gaussian around a mode. (Top) First, one identifies a mode  $\theta_{\text{MAP}}$  of the log-posterior  $\log p(\theta | \mathbb{D})$ , then Taylor-expands the latter around it up to second order. (Bottom) The exponentiated Taylor expansion is proportional to a Gaussian. Normalizing it yields the desired Gaussian approximation.

estimate  $\theta_{\text{MAP}}$ , can be used to approximate the Bayesian posterior for the deep model’s parameters via  $p(\theta | \mathbb{D}) \approx \delta(\theta - \theta_{\text{MAP}})$ . The Laplace approximation [93] extends this posterior approximation to a multivariate Gaussian around  $\theta_{\text{MAP}}$  (see Figure 3.5 for an illustration).

Let  $\mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta) := \mathcal{L}_{\mathbb{D}}(\theta) + r(\theta)$  denote the regularized empirical risk from Equation (2.18) that gives rise to the MAP estimate  $\theta_{\text{MAP}}$ . Its quadratic Taylor expansion around  $\theta_{\text{MAP}}$  (assumed to be a stationary point) is

$$\begin{aligned} \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta) &= \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta_{\text{MAP}}) + \underbrace{\nabla_{\theta} \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta_{\text{MAP}})^{\top} (\theta - \theta_{\text{MAP}})}_{=0, \text{stationary point}} \\ &\quad + \frac{1}{2} (\theta - \theta_{\text{MAP}})^{\top} \nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta_{\text{MAP}}) (\theta - \theta_{\text{MAP}}) \\ &\quad + \mathcal{O}\left((\theta - \theta_{\text{MAP}})^3\right). \end{aligned}$$

Using up to quadratic term in the posterior Equation (2.19) yields

$$\begin{aligned} p(\theta | \mathbb{D}) &\approx \text{const}(\mathbb{D}, \theta_{\text{MAP}}) \\ &\quad \times \exp\left[-\frac{|\mathbb{D}|}{2} (\theta - \theta_{\text{MAP}})^{\top} \nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta_{\text{MAP}}) (\theta - \theta_{\text{MAP}})\right], \end{aligned}$$

which corresponds to a Gaussian

$$p(\theta | \mathbb{D}) \approx \mathcal{N}(\theta | \mu, \Sigma) = \frac{1}{Z} \exp\left[-\frac{1}{2} (\theta - \mu)^{\top} \Sigma^{-1} (\theta - \mu)\right] \quad (3.32a)$$

with mean, covariance, and normalization constant

$$\mu = \theta_{\text{MAP}}, \quad (3.32b)$$

$$\Sigma = \left(|\mathbb{D}| \nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta_{\text{MAP}})\right)^{-1}, \quad (3.32c)$$

$$Z = (2\pi)^{D/2} \det\left[\left(|\mathbb{D}| \nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}}^{\text{MAP}}(\theta_{\text{MAP}})\right)^{-1}\right]^{D/2}. \quad (3.32d)$$

The Gaussian posterior Equation (3.32a) is useful for Bayesian predictions with neural networks (Equation (2.14)) and various other tasks (e.g. model selection [75] and continual learning, see [40] for an overview).

The covariance matrix  $\Sigma$  in the Laplace approximation requires the Hessian  $\nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}}^{\text{MAP}} = \nabla_{\theta}^2 \mathcal{L}_{\mathbb{D}} + \nabla_{\theta}^2 r$  (specifically, its inverse). While the regularization term often has a simple Hessian, computing the empirical risk’s Hessian is expensive and suffers from the same issues of the Laplace approximation being undefined if the Hessian is not PSD. Applications often replace this Hessian with one of the PSD approximations presented in Section 3.2 because a covariance matrix must be PSD.

### 3.3.2 Model Compression

Neural networks are extremely overparameterized. Reducing their parameter count, also referred to as *pruning*, is important to deploy them on low-resource devices such as mobile phones and often decreases computational cost of inference. In general, the goal is to reduce parameters, while keeping the model’s performance (see [18] for a review).

One classic approach for post-training compression of models is the *Optimal Brain Damage/Surgeon (OBD/OBS)* framework of Hassibi and Stork [67] and LeCun et al. [96] that leverages approximate second-order information and is still used nowadays[43, 154, 161, 182]. The idea is to identify unimportant weights through a local quadratic Taylor approximation of the loss. Assuming the model is trained to a stationary point  $\theta_\star$ , the Taylor expansion's first-order term vanishes. A perturbation  $\delta := \theta - \theta_\star$  induces the change  $\Delta\mathcal{L}_D(\delta) := \mathcal{L}_D(\delta) - \mathcal{L}_D(\theta_\star)$  in the loss,

$$\Delta\mathcal{L}_D(\delta) = \underbrace{(\nabla_{\theta_\star}\mathcal{L}_D(\theta_\star))^\top \delta}_{=0, \text{stationary point}} + \frac{1}{2} \delta^\top H_D(\theta_\star) \delta + \mathcal{O}(\delta^3). \quad (3.33)$$

To simplify the discussion, consider eliminating only a single parameter. For a fixed direction  $i$ , the “best” perturbation  $\delta_\star(i)$  that eliminates parameter  $[\theta_\star]_i$  through satisfying the constraint  $\hat{e}_i^\top \delta_\star(i) = [\theta_\star]_i$  (with  $\hat{e}_i$  the unit vector in direction  $i$ ) introduces the smallest increase in the loss. This leads to the constrained optimization problem

$$\delta_\star(i) = \arg \min_{\delta} \frac{1}{2} \delta^\top H_D(\theta_\star) \delta \quad (3.34a)$$

$$\text{subject to } \hat{e}_i^\top \delta - [\theta_\star]_i = 0$$

As described in Remark 3.4, the solution to Equation (3.34a) is

$$\delta_\star(i) = -\frac{[\theta_\star]_i H_D(\theta_\star)^{-1} \hat{e}_i}{[H_D(\theta_\star)^{-1}]_{i,i}} \quad (3.34b)$$

which induces the change

$$\Delta\mathcal{L}_D(\delta_\star(i)) = \frac{[\theta_\star]_i^2}{2[H_D(\theta_\star)^{-1}]_{i,i}} \quad (3.34c)$$

Among these perturbations,  $\{\delta_\star(i)\}_{i=1}^D$ , the one introducing the smallest change in the loss identifies the parameter to eliminate<sup>5</sup>. This requires computing the pruning statistics Equation (3.34c) for all directions, and therefore evaluation of the inverse Hessian's diagonal. The perturbation Equation (3.34b) requires multiplication by the inverse Hessian.

To simplify these computations, the OBD framework [96, Equation (3)] approximates the Hessian by its diagonal for the induced change Equation (3.33), i.e.  $\Delta\mathcal{L}_D^{OBD}(\delta) := 1/2\delta^\top \text{diag}(H_D(\theta_\star))\delta$ . By substituting  $H_D(\theta_\star) \leftrightarrow \text{diag}(H_D(\theta_\star))$  into the above discussion, one obtains the approximate perturbations  $\delta_\star^{OBD}(i) = [\theta_\star]_i \hat{e}_i$  and pruning statistics  $\Delta\mathcal{L}_D^{OBD}(\delta_\star^{OBD}(i)) = [\theta_\star]_i^2 [H_D(\theta_\star)]_{i,i}/2$ . The diagonal approximation removes the need to access elements of the inverse Hessian.

### 3.3.3 Differential Privacy

The data used for training deep models sometimes contains sensitive information about individuals, like medical data, personal images *etc.*, which is incorporated into the model during training. Adversaries might be able to reconstruct such personal data given only black box access to the model [52]. Often though, such adversaries even have access to details about the training procedure, or explicit access to model parameters<sup>6</sup>. Such knowledge can be leveraged to make attacks more efficient [149].

**Remark 3.4 (Optimal perturbation to prune  $[\theta]_i$ )** The constraint in Equation (3.33) is incorporated via a Lagrange multiplier  $\lambda$ ,

$$\Delta\mathcal{L}_D(\delta, \lambda) = \frac{1}{2} \delta^\top H_D(\theta_\star) \delta - \lambda (\hat{e}_i^\top \delta - [\theta_\star]_i).$$

The optimal perturbation follows as

$$\nabla_\delta \Delta\mathcal{L}_D(\delta, \lambda) = H_D(\theta_\star) \delta + \lambda \hat{e}_i \stackrel{!}{=} 0$$

$$\implies \delta(\lambda) = -\lambda H_D(\theta_\star)^{-1} \hat{e}_i.$$

Substituting this into the induced change yields

$$\Delta\mathcal{L}_D(\lambda) = -\frac{1}{2} \lambda^2 [H_D(\theta_\star)^{-1}]_{i,i} + \lambda [\theta_\star]_i$$

and

$$\nabla_\lambda \Delta\mathcal{L}_D(\lambda) \stackrel{!}{=} 0$$

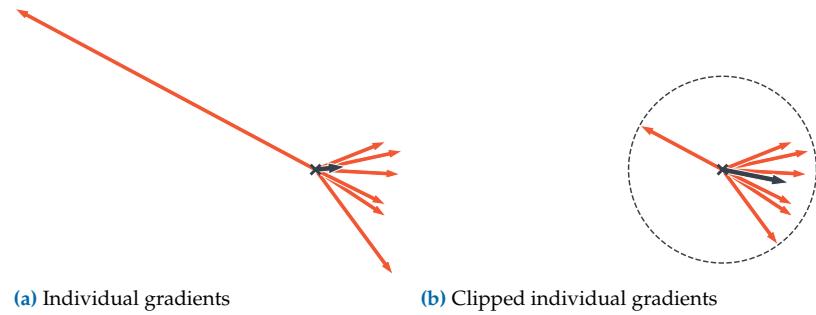
$$\implies \lambda = \frac{[\theta_\star]_i}{[H_D(\theta_\star)^{-1}]_{i,i}}.$$

Equations (3.34b) and (3.34c) follow by substitution of  $\delta_\star(i) := \delta(\lambda)$ .

5: Considering the removal of multiple parameters becomes prohibitive as the search space over indices increases exponentially (see [154] for details). It is therefore common to consider pruning candidates independently by sorting them according to their induced change.

6: E.g. distributed training on mobile devices that avoids sending data to a central server, but transfers the model to devices, which communicates updates to a central model.

**Figure 3.6:** (Sketch) DP-SGD requires access to per-sample gradients. (a) The average gradient (black arrow) might be dominated by the per-sample gradients (orange arrows) of a small number of data. Adversaries might use this to extract potentially sensitive information about that data. (b) Clipping per-sample gradients before averaging them limits the influence of individual data. The dashed circle's radius corresponds to the privacy threshold  $C$ .



Differential privacy [46, DP] aims at preserving privacy of an algorithm by injecting noise to limit the influence of individual data.

One prominent example for deep learning is differentially-private SGD (DP-SGD) [2, Algorithm 1], which uses a negative average of processed individual gradients over a mini-batch. To bound the influence of individual data, per-sample gradients whose norm exceeds a specific threshold  $C$  are clipped back to norm  $C$  (Figure 3.6),

$$\tilde{g}_n(\theta_t) = \frac{g_n(\theta_t)}{\max\left(1, \frac{\|g_n(\theta_t)\|_2}{C}\right)}.$$

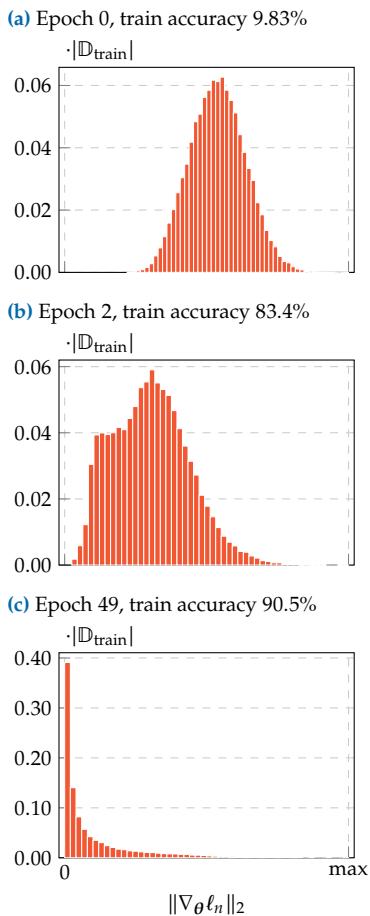
Gaussian noise of scale  $\sigma$  is then added to each gradient,

$$\hat{g}_n(\theta_t) = \tilde{g}_n(\theta_t) + \epsilon_n, \quad \epsilon_n \sim \mathcal{N}(\mathbf{0}, (\sigma C)^2 \mathbf{I}).$$

DP-SGD performs the update rule of SGD, with  $g_n$  replaced by  $\hat{g}_n$ ,

$$\theta_{t+1} = \theta_t - \eta_t \frac{1}{|\mathbb{B}|} \sum_{(x_n, y_n) \in \mathbb{B}} \hat{g}_n(\theta_t).$$

This requires access to per-sample gradient norms  $\{\|g_n(\theta_t)\|_2\}_n$ .



**Figure 3.8:** Importance distribution of samples during training. Importance is measured by individual gradient  $L_2$  norms, computed with BackPACK (Chapter 5). As training proceeds and more examples are correctly classified and become “unimportant”. Details: logistic regression on MNIST, trained with SGD,  $|\mathbb{B}| = 128$ , and learning rate 0.005.

### 3.3.4 Importance Sampling

A hypothesis about learning in ML is that the model first learns to correctly predict “easy” examples. Only in later phases are the “difficult” examples learned. Using these harder examples more frequently for training might help speed up the learning procedure. Put differently, some data matter more at certain stages of training than other, which is quantified through a measure of importance. Importance sampling realizes this idea of selecting important data more frequently. It does so by adapting the sampling procedure for mini-batches to reduce the gradient variance, which beneficially influences the convergence of stochastic first-order methods like SGD, and therefore speeds up training. A common strategy is to weight the importance of samples by the per-sample  $L_2$  norm (Equation (3.38) and Figure 3.8).

As a starting point to see how the sampling procedure affects optimization, consider the generalization of uniform sampling from Section 2.1.2 for  $|\mathbb{B}| = 1$ . At training iteration  $t$ , a sample  $n_t \sim p_t(n_t)$  is drawn from a current sampling distribution  $p_t(n_t)$  over  $\{1, \dots, |\mathbb{D}|\}$ . SGD with a

learning rate  $\eta$  uses the unbiased gradient estimator<sup>7</sup>

$$\hat{g}_{n_t}(\boldsymbol{\theta}_t) := \frac{1}{|\mathbb{D}| p_t(n_t)} \nabla_{\boldsymbol{\theta}_t} \ell_{n_t}(\boldsymbol{\theta}_t) \quad (3.35a)$$

with

$$\mathbb{E}_{n_t \sim p_t(n_t)} [\hat{g}_{n_t}(\boldsymbol{\theta}_t)] = \frac{1}{|\mathbb{D}|} \sum_{(x_n, y_n) \in \mathbb{D}} \nabla_{\boldsymbol{\theta}_t} \ell_{n_t}(\boldsymbol{\theta}_t) = g_{\mathbb{D}}(\boldsymbol{\theta}_t) \quad (3.35b)$$

to update the parameters

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \hat{g}_{n_t}(\boldsymbol{\theta}_t). \quad (3.35c)$$

Using a non-uniform distribution over stochastic gradients would introduce bias in the update, so in importance sampling the samples are weighted by their probability of being selected to undo the bias. The distribution can then be tuned to minimize the variance of the estimator and improve performance: one can assess convergence of Equation (3.35c) through the progress towards the minimizer  $\boldsymbol{\theta}_*$  in terms of squared  $L_2$  distance [85, 169],

$$\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_*\|_2^2 - \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_*\|_2^2. \quad (3.36a)$$

In expectation, this measure of convergence depends on the gradient estimator's variance (more specifically, its trace),

$$\begin{aligned} & \mathbb{E}_{n_t} [\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_*\|_2^2 - \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_*\|_2^2] \\ &= \mathbb{E}_{n_t} [\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_*\|_2^2 - \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_* - \eta \hat{g}_{n_t}\|_2^2] \\ &= \mathbb{E}_{n_t} [2\eta (\boldsymbol{\theta}_t - \boldsymbol{\theta}_*)^\top \hat{g}_{n_t} - \eta^2 \hat{g}_{n_t}^\top \hat{g}_{n_t}] \\ &= 2\eta (\boldsymbol{\theta}_t - \boldsymbol{\theta}_*)^\top g_{\mathbb{D}} - \eta^2 \mathbb{E}_{n_t} [\|\hat{g}_{n_t}\|_2^2] \\ &= 2\eta (\boldsymbol{\theta}_t - \boldsymbol{\theta}_*)^\top g_{\mathbb{D}} - \eta^2 \|g_{\mathbb{D}}\|_2^2 - \eta^2 \mathbb{E}_{n_t} [\|\hat{g}_{n_t} - g_{\mathbb{D}}\|_2^2] \\ &= 2\eta (\boldsymbol{\theta}_t - \boldsymbol{\theta}_*)^\top g_{\mathbb{D}} - \eta^2 \|g_{\mathbb{D}}\|_2^2 - \eta^2 \text{Tr} \{ \text{Var}_{n_t} [\hat{g}_{n_t}] \} \end{aligned} \quad (3.36b)$$

For what follows, the intuition is that sampling affects convergence through the variance term, and minimizing this term improves convergence. Hence, the goal is to identify the optimal sampling via

$$\underset{\substack{p_t(n_t) \\ \sum_{n=1}^{|\mathbb{D}|} p_t(n)=1}}{\text{minimize}} \text{Tr} \{ \text{Var}_{n_t} [\hat{g}_{n_t}] \} \Leftrightarrow \underset{\substack{p_t(n_t) \\ \sum_{n=1}^{|\mathbb{D}|} p_t(n)=1}}{\text{minimize}} \mathbb{E}_{n_t} [\|\hat{g}_{n_t}\|_2^2] \quad (3.37)$$

The solution, outlined in Remark 3.5, is

$$p_t(n_t) = \frac{\|\nabla_{\boldsymbol{\theta}_t} \ell_{n_t}(\boldsymbol{\theta}_t)\|_2}{\sum_{n=1}^{|\mathbb{D}|} \|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta}_t)\|_2}. \quad (3.38)$$

and depends on individual gradient  $L_2$  norms in the entire dataset. Samples with higher importance, *i.e.* gradient norm, are drawn more frequently. Because computing the sampling distribution requires a sweep over all data, practical versions further approximate Equation (3.38), *e.g.* by relaxing the optimization through bounds that are easier to compute, or updating the distribution only every few iterations [85].

7: For uniform sampling  $p_t(n_t) = 1/|\mathbb{D}|$ , which is the most common case in practice, the scaling factor cancels and yields the commonly used gradient estimator that would be used by "normal" SGD with  $|\mathbb{B}| = 1$ .

**Remark 3.5 (Optimal sampling distribution [116, 184])** (The iteration index  $t$  is suppressed for brevity) To derive Equation (3.38), one incorporates the constraint in Equation (3.37) via a Lagrange multiplier  $\mu \in \mathbb{R}$ , writing  $p(n)$  as  $|\mathbb{D}|$ -dimensional vector  $\mathbf{p}$ . Expanding  $\hat{g}$  with Equation (3.35a) leads to,

$$\begin{aligned} A(\mathbf{p}, \mu) &:= \frac{1}{|\mathbb{D}|^2} \sum_{n=1}^{|\mathbb{D}|} \frac{\|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2^2}{p_n} \\ &\quad + \mu \left( \sum_{n=1}^{|\mathbb{D}|} p_n - 1 \right). \end{aligned}$$

Setting the  $p_n$ -derivative of that to zero yields

$$\begin{aligned} \nabla_{p_n} A(\mathbf{p}, \mu) &= -\frac{\|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2^2}{|\mathbb{D}|^2 p_n^2} + \mu \stackrel{!}{=} 0 \\ \implies p_n &= \frac{\|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2}{\sqrt{\mu}}, \end{aligned}$$

using that all elements  $p_n \in (0, 1)$  of  $\mathbf{p}$ . Inserting into  $A$  gives

$$A(\mu) = \frac{2}{|\mathbb{D}|} \sqrt{\mu} \sum_{n=1}^{|\mathbb{D}|} \|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2 - \mu$$

Setting the  $\mu$ -derivative to zero yields

$$\begin{aligned} \nabla_\mu A(\mu) &= \frac{1}{|\mathbb{D}| \sqrt{\mu}} \sum_{n=1}^{|\mathbb{D}|} \|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2 - 1 \stackrel{!}{=} 0 \\ \implies \sqrt{\mu} &= \frac{1}{|\mathbb{D}|} \sum_{n=1}^{|\mathbb{D}|} \|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2 \end{aligned}$$

This then leads to

$$p_n = \frac{\|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2}{\sum_{n=1}^{|\mathbb{D}|} \|\nabla_{\boldsymbol{\theta}_t} \ell_n(\boldsymbol{\theta})\|_2}$$

which equals Equation (3.38) when switching back notation and introducing the iteration count.



## **Part II.**

# Backpropagation Beyond the Gradient



# Modular Block-diagonal Curvature Approximations for Feedforward Architectures

# 4.

## Abstract

We propose a modular extension of backpropagation for the computation of block-diagonal approximations to various curvature matrices of the training objective (in particular, the Hessian, generalized Gauss-Newton, and positive-curvature Hessian). The approach reduces the otherwise tedious manual derivation of these matrices into local modules, and is easy to integrate into existing machine learning libraries. Moreover, we develop a compact notation derived from matrix differential calculus. We outline different strategies applicable to our method. They subsume recently-proposed block-diagonal approximations as special cases, and are extended to convolutional neural networks in this work.

4.1 Introduction . . . . .	51
4.2 Notation . . . . .	53
4.3 Modular Hessian Backpropagation . . . . .	53
4.4 Experiments & Implementation . . . . .	59
4.5 Conclusion . . . . .	62

Code and experiments available at the Github repository [f-dangel/hbp](https://github.com/f-dangel/hbp)

## 4.1 Introduction

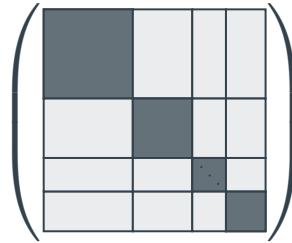
Gradient backpropagation is the central computational operation of contemporary deep learning. Its modular structure allows easy extension across network architectures, and thus automatic computation of gradients given the computational graph of the forward pass [for a review, see 13]. But optimization using only the first-order information of the objective’s gradient can be unstable and slow, due to “vanishing” or “exploding” behaviour of the gradient. Incorporating curvature, second-order methods can avoid such scaling issues and converge in fewer iterations. Such methods locally approximate the objective function  $\mathcal{L}$  by a quadratic  $\mathcal{L}(\theta) + \nabla_\theta \mathcal{L}(\theta)^\top (\theta_\star - \theta) + 1/2(\theta_\star - \theta)^\top C(\theta)(\theta_\star - \theta)$  around the current location  $\theta$ , using the gradient  $\nabla_\theta \mathcal{L}(\theta) = \partial \mathcal{L}(\theta) / \partial \theta$  and a positive semi-definite (PSD) curvature matrix  $C(\theta)$ —the Hessian of  $\mathcal{L}$  or approximations thereof. The quadratic is minimized by

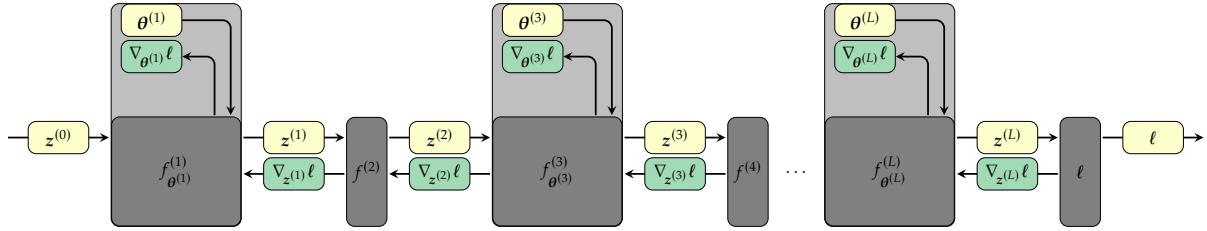
$$\theta_\star = \theta + \Delta\theta \quad \text{with} \quad \Delta\theta = -C(\theta)^{-1}\nabla_\theta \mathcal{L}. \quad (4.1)$$

Computing the update step requires solving the linear system  $C(\theta)\Delta\theta = -\nabla_\theta \mathcal{L}$ . To accomplish this task, providing a matrix-vector multiplication with the curvature matrix  $C(\theta)$  is sufficient.

### Approaches to Second-order Optimization

For some curvature matrices, exact multiplication can be performed at the cost of one backward pass by automatic differentiation [127, 148]. This *matrix-free* formulation can then be leveraged to solve Equation (4.1) using iterative solvers such as conjugate gradients (CG) [106]. However, since this linear solver can still require multiple iterations, the increased per-iteration progress of the resulting optimizer might be compensated by increased computational cost. Recently, a parallel version of Hessian-free optimization was proposed in [183], which only considers the content of Hessian sub-blocks along the diagonal. Reducing the Hessian to a block





**Figure 4.1: Standard sequential feedforward network architecture.** I.e. the repetition of affine transformations parameterized by  $\theta^{(l)} = ((\text{vec } W^{(l)})^\top, b^{(l)\top})^\top$  followed by element-wise activations. Arrows from left to right and vice versa indicate the data flow during forward pass and gradient backpropagation, respectively.

diagonal allows for parallelization, tends to lower the required number of CG iterations, and seems to improve the optimizer’s performance.

There have also been attempts to compute parts of the Hessian in an iterative fashion [110]. Storing these constituents efficiently often requires an involved manual analysis of the Hessian’s structure, leveraging its outer-product form in many scenarios [8, 115]. Recent works developed different block-diagonal approximations (BDA) of curvature matrices that provide fast multiplication [21, 31, 63, 109].

These works have repeatedly shown that, empirically, second-order information can improve the training of deep learning problems. Perhaps the most important practical hurdle to the adoption of second-order optimizers is that they tend to be tedious to integrate in existing ML frameworks because they require manual implementations. As efficient automated implementations have arguably been more important for the wide-spread use of deep learning than many conceptual advances, we aim to develop a framework that makes computation of Hessian approximations about as easy and automated as gradient backpropagation.

### Contribution

This chapter introduces a modular formalism for the computation of block-diagonal approximations of Hessian and curvature matrices, to various block resolutions, for feedforward neural networks. The framework unifies previous approaches in a form that, similar to gradient backpropagation, reduces implementation and analysis to local modules. Following the design pattern of gradient backprop also has the advantage that this formalism can readily be integrated into existing ML libraries, and flexibly modified for different block groupings and approximations.

The framework consists of three principal parts:

1. A modular formulation for *exact* computation of Hessian block diagonals of feedforward neural nets. We achieve a clear presentation by leveraging the notation of matrix differential calculus [103].
2. Projections onto the positive semi-definite cone by eliminating sources of concavity.
3. Backpropagation strategies to obtain (i) exact curvature matrix-vector products (with previously inaccessible BDAs of the Hessian) and (ii) further approximated multiplication routines that save computations by evaluating the matrix representations of intermediate quantities once, at the cost of additional memory consumption.

The first two contributions can be understood as an explicit formulation of well-known tricks for fast multiplication by curvature matrices using automatic differentiation [127, 148]. However, we also address a new class of curvature matrices, the positive-curvature Hessian (PCH) introduced in [31]. Our solutions to the latter two points are generalizations of previous works [21, 31] to the fully modular case, which become accessible due to the first contribution. They represent additional modifications to make the scheme computationally tractable and obtain curvature approximations with desirable properties for optimization.

## 4.2 Notation

We consider feedforward neural networks  $f_\theta(x)$  composed of  $L$  modules  $f_{\theta^{(l)}}, l = 1, \dots, L$ , which can be represented as a computational graph mapping the input  $z^{(0)} = x$  to the output  $z^{(L)}$  (Figure 4.1). A module  $f_{\theta^{(l)}}$  receives the parental output  $z^{(l-1)}$ , applies an operation involving the module parameters  $\theta^{(l)}$ , and sends the output  $z^{(l)}$  to its child. Thus,  $f_{\theta^{(l)}}$  is of the form  $z^{(l)} = f_{\theta^{(l)}}(z^{(l-1)})$ . Typical choices include element-wise nonlinear activation without any parameters and affine transformations  $z^{(l)} = W^{(l)} z^{(l-1)} + b^{(l)}$  with parameters given by the weights  $W^{(l)}$  and the bias  $b^{(l)}$ . Affine and activation modules are usually considered as a single conceptual unit, one *layer* of the network. However, for backpropagation of derivatives it is simpler to consider them separately as two *modules*.

Given the network output  $z^{(L)}(x, \theta^{(1, \dots, L)}) = f_\theta(x)$  of a datum  $x$  with label  $y$ , the goal is to minimize the expected risk of the loss function  $\ell(z^{(L)}, y)$ . Under the framework of empirical risk minimization, the parameters are tuned to optimize the loss on the training set  $\mathbb{D}_{\text{train}} = \{(x_n, y_n)\}_{n=1}^N$ ,

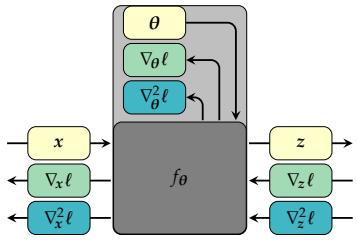
$$\min_{\theta^{(1, \dots, L)}} \mathcal{L}_{\mathbb{D}_{\text{train}}}(\theta^{(1, \dots, L)}) = \min_{\theta^{(1, \dots, L)}} \frac{1}{|\mathbb{D}_{\text{train}}|} \sum_{n=1}^N \ell(f_\theta(x_n), y_n). \quad (4.2)$$

In practice, the objective is typically further approximated stochastically by drawing a mini-batch  $\mathbb{B} \subset \mathbb{D}_{\text{train}}$  from the training set. We will treat both scenarios without further distinction, since the structure relevant to our purposes is that Equation (4.2) is an average of terms depending on individual data points. Quantities for optimization, be it gradients or second derivatives of the loss *w.r.t.* the network parameters, can be processed in parallel, then averaged.

## 4.3 Modular Hessian Backpropagation

First-order auto-differentiation for a custom module requires the definition of only two local operations, *forward* and *backward*, whose outputs are propagated along the computation graph. This modularity facilitates the extension of gradient backpropagation by new operations, which can then be used to build networks by composition. To illustrate the principle, we consider a single module from the net of Figure 4.1, depicted in Figure 4.2, in this section<sup>1</sup>. The forward pass  $f_\theta(x)$  maps the input  $x$  to the output  $z$  by means of the module parameters  $\theta$ . All quantities are

<sup>1</sup>: To simplify the presentation, we drop layer indices, choose distinct variable names for input and output ( $x \leftarrow z^{(l-1)}$ ,  $z \leftarrow z^{(l)}$ ,  $\theta \leftarrow \theta^{(l)}$ ,  $f_\theta \leftarrow f_{\theta^{(l)}}$ ) and focus on a single datum.



**Figure 4.2:** Forward pass, gradient backpropagation, and Hessian backpropagation for a single module. Arrows from left to right indicate the data flow in the forward pass  $z = f_\theta(x)$ , while the opposite orientation indicates the gradient backpropagation by Equation (4.4). We suggest to extend this by the backpropagation of the Hessian as indicated by Equation (4.7).

assumed to be vector-shaped (tensor-valued quantities can be flattened, see Definition 2.2). Optimization requires the gradient of the loss function w.r.t. the parameters,  $\partial \ell(\theta)/\partial \theta = V_\theta \ell$ . We will use the shorthand

$$\nabla_\cdot \ell = \frac{\partial \ell(\cdot)}{\partial \text{vec}(\cdot)}. \quad (4.3)$$

During gradient backpropagation the module receives the loss gradient with respect to its output,  $V_z \ell$ , from its child. The backward operation computes gradients w.r.t. the module parameters and input,  $V_\theta \ell$  and  $V_x \ell$  from  $V_z \ell$ . Backpropagation continues by sending the gradient w.r.t. the module's input to its parent, which proceeds in the same way (see Figure 4.1). By the chain rule, gradients w.r.t. an element  $x_i$  of the module's input can be computed as  $V_{x_i} \ell = \sum_j (\partial z_j / \partial x_i) V_{z_j} \ell$ . The vectorized version is compactly written in terms of the Jacobian matrix  $J_x z = \partial z / \partial x^\top$  (Definition 2.4), which contains all partial derivatives of  $z$  w.r.t.  $x$ . The arrangement of partial derivatives is  $[J_x z]_{j,i} = \partial z_j / \partial x_i$ , such that

$$\nabla_x \ell = (J_x z)^\top V_z \ell. \quad (4.4)$$

Analogously, the parameter gradients are given by  $V_{\theta_i} \ell = \sum_j (\partial z_j / \partial \theta_i) V_{z_j} \ell$ , or in vectorized form  $V_\theta \ell = (J_\theta z)^\top V_z \ell$ . This reflects the symmetry of both  $x$  and  $\theta$  acting as input to the module. Implementing gradient backpropagation thus requires multiplications by (transposed) Jacobians. We can apply the chain rule a second time to obtain expressions for second-order partial derivatives of the loss  $\ell$  w.r.t. elements of  $x$  or  $\theta$ ,

$$\begin{aligned} \frac{\partial^2 \ell}{\partial x_i \partial x_j} &= \frac{\partial}{\partial x_j} \left( \sum_k \frac{\partial z_k}{\partial x_i} V_{z_k} \ell \right) \\ &= \sum_{k,l} \frac{\partial z_k}{\partial x_i} \frac{\partial^2 \ell}{\partial z_k \partial z_l} \frac{\partial z_l}{\partial x_j} + \sum_k \frac{\partial^2 z_k}{\partial x_i \partial x_j} V_{z_k} \ell, \end{aligned} \quad (4.5)$$

by means of  $\partial / \partial x_j = \sum_l (\partial z_l / \partial x_j) \partial / \partial z_l$  and the product rule. The first term of Equation (4.5) propagates curvature information of the output further back, while the second term introduces second-order effects of the module itself. Using the Hessian matrix  $\nabla_x^2 \ell = \partial^2 \ell / (\partial x^\top \partial x)$  of a scalar function w.r.t. a vector-shaped quantity  $x$  (Definition 3.1),

$$\nabla_x^2 \ell(\cdot) = \frac{\partial^2 \ell(\cdot)}{\partial \text{vec}(\cdot)^\top \partial \text{vec}(\cdot)}, \quad (4.6)$$

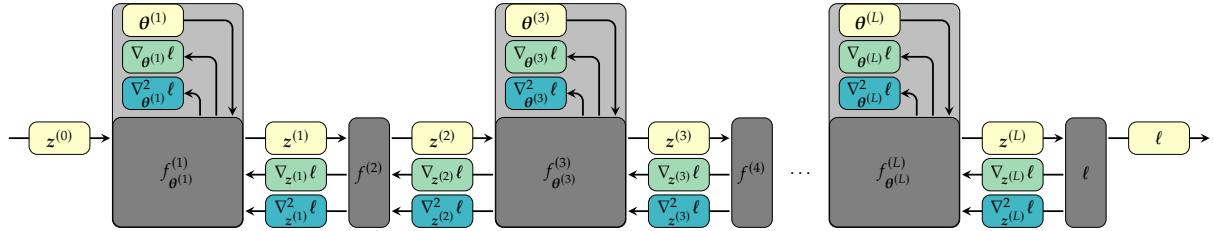
results in the matrix version of Equation (4.5),

$$\nabla_x^2 \ell = (J_x z)^\top V_z^2 \ell (J_x z) + \sum_k (\nabla_x^2 z_k) V_{z_k} \ell. \quad (4.7a)$$

Note that the second-order effect introduced by the module itself via  $\nabla_x^2 z_k$  vanishes if all components  $[f_\theta(x)]_k$  are linear in  $x$ . Because the layer parameters  $\theta$  can be regarded as inputs to the layer, they are treated in exactly the same way, replacing  $x$  by  $\theta$  in the above expression,

$$\nabla_\theta^2 \ell = (J_\theta z)^\top V_z^2 \ell (J_\theta z) + \sum_k (\nabla_\theta^2 z_k) V_{z_k} \ell. \quad (4.7b)$$

Equation (4.7) is the central functional expression herein, and will be referred to as the *Hessian backpropagation (HBP) equation*. Our suggested



**Figure 4.3:** Extension of backpropagation to Hessians. It yields diagonal blocks of the full parameter Hessian.

extension of gradient backpropagation is to also send the Hessian  $\nabla_z^2 \ell$  back through the graph. To do so, existing modules have to be extended by the HBP equation: *Given the Hessian  $\nabla_z^2 \ell$  of the loss w.r.t. all module outputs, an extended module has to extract the Hessians  $\nabla_\theta^2 \ell$ ,  $\nabla_x^2 \ell$  by means of Equation (4.7), and forward the Hessian w.r.t. its input  $\nabla_x^2 \ell$  to the parent module which proceeds likewise.* In this way, backprop of gradients can be extended to compute curvature information in modules. This corresponds to BDAs of the Hessian which ignore second-order derivatives w.r.t. parameters in different modules. Figure 4.3 shows the data flow. The computations required in Equation (4.7) depend only on *local quantities* that are, mostly, already being computed during gradient backpropagation<sup>2</sup>.

Before we proceed, we highlight the following aspects:

- ▶ The BDA of the Hessian need not be PSD. But our scheme can be modified to provide PSD curvature matrices by projection onto the positive semi-definite cone (see Section 4.3.1).
- ▶ Instead of evaluating all matrices during backpropagation, we can define matrix-vector products recursively. This yields exact curvature matrix multiplications by the block diagonals of the Hessian, the generalized Gauss-Newton (GGN) matrix and the PCH. Multiplications by the first two matrices can also be obtained by use of automatic differentiation [127, 148]. We also get access to the PCH which, in contrast to the GGN, considers curvature information introduced by the network (see Section 4.3.1).<sup>3</sup> For standard neural networks, only second derivatives of nonlinear activations have to be stored compared to gradient backpropagation.
- ▶ There are approaches [21, 31] that propagate matrix representations back through the graph to save repeated computations in the curvature matrix-vector product. The size of the matrices  $\nabla_{z^{(l)}}^2 \ell$  passed between layer  $l+1$  and  $l$  scales quadratically in the number of layer  $l$ 's output features. For convolutional layers, the dimension of these quantities exceeds computational budgets. And for batched computations, such a matrix has to be backpropagated for every datum in the mini-batch. Like previous schemes [21, 31], we introduce additional approximations for batch learning in Section 4.3.2. A connection to existing schemes is drawn in Appendix A.2.4.

HBP is easy to integrate into current ML libraries, so that curvature BDAs can be provided automatically for novel or existing second-order optimization methods<sup>4</sup>. Such methods have repeatedly been shown to be competitive with first-order methods [21, 31, 63, 109, 183].

2: By Faà di Bruno's formula [82] higher-order derivatives of function compositions are expressed recursively in terms of the composites' lower-order derivatives. Recycling these quantities can give significant speedup compared to repeatedly applying first-order auto-differentiation, which represents one key aspect of our work.

3: Implementations of HBP for exact matrix-vector products can reuse multiplication by the (transposed) Jacobian provided by many ML libraries. The second term Equation (4.7) needs special treatment though.

4: E.g. the matrix-vector products with block-diagonal curvature matrices described in this work have been integrated into the BackPACK library that operates on top of PyTorch and will be presented in Chapter 5 of this thesis.

**Table 4.1:** Hessian backpropagation for common modules used in feedforward networks.  $I$  denotes the identity matrix. We assign matrices to upper-case ( $\hat{W}, \hat{X}, \dots$ ) and tensors to upper-case sans serif symbols ( $W, X, \dots$ ).

OPERATION	FORWARD	HBP (Equation (4.7))	DETAILS
Matrix-vector multiplication	$z(x, W) = Wx$	$\nabla_x^2 \ell = W^\top (\nabla_z^2 \ell) W,$ $\nabla_W^2 \ell = x \otimes x^\top \otimes \nabla_z^2 \ell$	Appendix A.2.1
Matrix-matrix multiplication	$Z(X, W) = WX$	$\nabla_X^2 \ell = (I \otimes W)^\top \nabla_Z^2 \ell (I \otimes W),$ $\nabla_W^2 \ell = (X^\top \otimes I)^\top \nabla_Z^2 \ell (X^\top \otimes I)$	Appendix A.2.1
Addition	$z(x, b) = x + b$	$\nabla_x^2 \ell = \nabla_b^2 \ell = \nabla_z^2 \ell$	Appendix A.2.1
Elementwise activation	$z(x) = \phi(x), \text{s.t.,}$ $z_i(x) = \phi(x_i)$	$\nabla_x^2 \ell = \text{diag}[\phi'(x)] (\nabla_z^2 \ell) \text{diag}[\phi'(x)]$ + $\text{diag}[\phi''(x) \odot \nabla_z \ell]$	Appendix A.2.2
Skip-connection	$z(x, \theta) = x + s(x, \theta)$	$\nabla_x^2 \ell = (I + J_x s)^\top \nabla_z^2 \ell (I + J_x s)$ + $\sum_k (\nabla_x^2 s_k) \nabla_{z_k} \ell,$ $\nabla_\theta^2 \ell = (J_\theta s)^\top \nabla_z^2 \ell (J_\theta s)$ + $\sum_k (\nabla_\theta^2 s_k) \nabla_{z_k} \ell$	Appendix A.2.3
Reshape/view	$Z(X) = \text{reshape}(X)$	$\nabla_Z^2 \ell = \nabla_X^2 \ell$	Appendix A.4.1
Index select/map $\pi$	$z(x) = \Pi x, \Pi_{j,\pi(j)} = 1,$	$\nabla_x^2 \ell = \Pi^\top (\nabla_z^2 \ell) \Pi$	Appendix A.4.2
Convolution	$Z(X, W) = X \star W,$ $Z(W, \ X\ ) = W \ X\ ,$	$\nabla_{\ X\ }^2 \ell = (I \otimes W)^\top \nabla_Z^2 \ell (I \otimes W)$ $\nabla_W^2 \ell = (\ X\ ^\top \otimes I)^\top \nabla_Z^2 \ell (\ X\ ^\top \otimes I)$	Appendix A.4.3
Square loss	$\ell(f, y) = 1/C(y - f)^\top (y - f)$	$\nabla_f^2 \ell = 2/C I$	Appendix A.3.1
Softmax cross-entropy	$\ell(f, y) = -\text{onehot}(y)^\top \log [p(f)]$	$\nabla_f^2 \ell = \text{diag}[p(f)] - p(f)p(f)^\top$	Appendix A.3.2

### Relationship to Matrix Differential Calculus

To some extent, this paper is a re-formulation of earlier results [21, 31, 109] in the framework of matrix differential calculus [103], leveraged to achieve a new level of modularity. Matrix differential calculus is a set of notational rules that allow a concise construction of derivatives without the heavy use of indices. Equation (4.7) is a special case of the matrix chain rule of that framework (Theorem 2.1 and Equation (3.14)). A more detailed discussion of this connection can be found in Appendix A.1 of the Supplements, which also reviews definitions generalizing the concepts of Jacobian and Hessian in a way that preserves the chain rule. The elementary building block of our procedure is a *module* as shown in Figure 4.2. Like for gradient backprop, the operations required for HBP can be tabulated. Table 4.1 provides a selection of common modules. The derivations, which again leverage the matrix differential calculus framework, can be found in Appendices A.2 to A.4.

### 4.3.1 Obtaining Different Curvature Matrices

The HBP equation yields *exact* diagonal blocks  $\nabla_{\theta^{(1)}}^2 \ell, \dots, \nabla_{\theta^{(L)}}^2 \ell$  of the full parameter Hessian  $\nabla_\theta^2 \ell$ . They can be of interest in their own right for analysis of the loss function, but are not generally suitable for second-order optimization in the sense of equation 4.1, as they need neither be PSD nor invertible. For application in optimization, HBP can be modified to yield semi-definite BDAs of the Hessian. Equation (4.7) again provides the foundation for this adaptation, which is closely related to the concepts of KFRA [21], BDA-PCH [31], and, under certain conditions, KFAC [109]. We draw their connections by briefly reviewing them here.

### Generalized Gauss-Newton Matrix

The GGN emerges as the curvature matrix in the quadratic expansion of the loss function  $\ell(z^{(L)}, y)$  in the network output  $z^{(L)}$  (Section 3.2.3). It is also obtained by linearizing the network output  $z^{(L)}(\theta, x)$  in  $\theta$  before computing the loss Hessian [107], and reads (Equation (3.15b))

$$G(\theta) = \left( J_{\theta} z^{(L)} \right)^T \nabla_{z^{(L)}}^2 \ell \left( J_{\theta} z^{(L)} \right).$$

For diagonal blocks  $G(\theta^{(l)})$ , the Jacobian is unrolled using its chain rule (Theorem 2.2)  $J_{\theta^{(l)}} z^{(L)} = (J_{z^{(L-1)}} z^{(L)}) (J_{z^{(L-2)}} z^{(L-1)}) \cdots (J_{z^{(l)}} z^{(l+1)}) (J_{\theta^{(l)}} z^{(l)})$ . This shows that the Hessian  $\nabla_{z^{(L)}}^2 \ell$  of the loss function *w.r.t.* the network output is propagated back through a layer by multiplication from left and right with its Jacobian. This is accomplished in HBP by *ignoring second-order effects introduced by modules*, that is by setting the Hessian of the module function to zero, therefore neglecting the second term in Equation (4.7)<sup>5</sup>. In fact, if all activations in the network are piecewise linear (*e.g.* ReLUs), the GGN and Hessian blocks are equivalent. Moreover, diagonal blocks of the GGN are PSD if the loss function is convex (and thus  $\nabla_{z^{(L)}}^2 \ell$  is PSD). This is because blocks are recursively left- and right-multiplied with Jacobians, which does not alter the definiteness. Hessians of the loss functions listed in Table 4.1 are PSD. The resulting recursive scheme has been used by Botev et al. [21] under the acronym KFRA to optimize convex loss functions of fully-connected neural networks with piecewise linear activation functions.

5: Recall from Section 2.2 that many common neural network operations, such as affine transformations, convolutions, padding & pooling, have vanishing second-order derivatives *w.r.t.* both their input and parameters. Hence, the second term in Equation (4.7) vanishes exactly for these operations.

### Positive-curvature Hessian

Another concept of positive semi-definite BDAs of the Hessian (that additionally considers second-order module effects) was studied in Chen et al. [31] and named the PCH. It is obtained by modification of terms in the second summand of Equation (4.7) that can introduce concavity during HBP. This ensures positive semi-definiteness since the first summand is semi-definite, assuming the loss Hessian  $\nabla_{z^{(L)}}^2 \ell$  with respect to the network output to be positive semi-definite. Chen et al. [31] suggest to eliminate negative curvature of a matrix by computing the eigenvalue decomposition and either discard negative eigenvalues or cast them to their absolute value. This allows the construction of PSD curvature matrices even for non-convex loss functions. In the setting of Chen et al. [31], the PCH can empirically outperform optimization using the GGN. In usual feedforward neural networks, the concavity is introduced by nonlinear element-wise activations, and corresponds to a diagonal matrix (Table 4.1). Thus, convexity can be maintained during HBP by either clipping negative values to zero (PCH-clip), or taking their magnitude in the diagonal concave term (PCH-abs).

### Fisher Information Matrix

If the network output models a conditional probability density on the labels,  $q(y | z^{(L)})$ , maximum likelihood learning for the parameterized density  $p_{\theta}(y | x) = q(y | z^{(L)}(x, \theta))$  will correspond to choosing a negative log-likelihood loss function, *i.e.*  $\ell(z^{(L)}, y) =$

$-\log q(\mathbf{y} \mid \mathbf{z}^{(L)})$  (Section 2.1.3). Common loss functions like square and cross-entropy loss can be interpreted in this way (Examples 2.3 and 2.4), . Natural gradient descent [5] uses the Fisher information matrix  $\mathbf{F}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{y} \sim p_{\boldsymbol{\theta}}(\mathbf{y} \mid \mathbf{x})} \left[ (\partial \log p_{\boldsymbol{\theta}}(\mathbf{y} \mid \mathbf{x}) / \partial \boldsymbol{\theta}) (\partial \log p_{\boldsymbol{\theta}}(\mathbf{y} \mid \mathbf{x}) / \partial \boldsymbol{\theta})^T \right]$  as a PSD curvature matrix approximating the Hessian. It can be expressed as the log-predictive density's expected Hessian under  $p_{\boldsymbol{\theta}}$  itself:  $\mathbf{F}(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{y} \sim p_{\boldsymbol{\theta}}(\mathbf{y} \mid \mathbf{z}^{(L)})} [\nabla_{\boldsymbol{\theta}}^2 \log q(\mathbf{y} \mid \mathbf{z}^{(L)})]$ , see Equation (3.24). Assuming truly *i.i.d.* samples  $\mathbf{x}$ , the log-likelihood of multiple data decomposes and, after using the chain rule, results in the approximation

$$\mathbf{F}_{\mathbb{D}_{\text{train}}}(\boldsymbol{\theta}) \approx \frac{1}{|\mathbb{D}_{\text{train}}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbb{D}_{\text{train}}} \left( \mathbf{J}_{\boldsymbol{\theta}} \mathbf{z}^{(L)} \right)^T \mathbf{F}_q(\mathbf{z}^{(L)}) \left( \mathbf{J}_{\boldsymbol{\theta}} \mathbf{z}^{(L)} \right)$$

with  $\mathbf{F}_q(\mathbf{z}^{(L)}) = -\mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid \mathbf{z}^{(L)})} [\nabla_{\mathbf{z}^{(L)}}^2 \log q(\mathbf{y} \mid \mathbf{z}^{(L)})]$ , see Equation (3.27). In this form, the computational scheme for Fisher BDAs resembles the HBP of the GGN. However, instead of propagating back the loss Hessian *w.r.t.* the network, it uses the negative log-likelihood's expected Hessian the model's predictive distribution. Martens and Grosse [109] use Monte Carlo sampling to estimate this matrix in their KFAC optimizer. Relations between the Fisher and GGN are discussed in [107, 125]; for square and cross-entropy loss, they are equivalent (Section 3.2.4).

### 4.3.2 Batch Learning Approximations

In our HBP framework, exact multiplication by the curvature block of a module's parameter  $\boldsymbol{\theta}$  requires one backpropagation to this layer. The multiplication is recursively defined in terms of multiplication by the layer output Hessian  $\nabla_{\mathbf{z}}^2 \ell$ . If it were possible to have an explicit representation of this matrix in memory, the recursive computations hidden in  $\nabla_{\mathbf{z}}^2 \ell$  could be saved during the solution of the linear system implied by Equation (4.1). Unfortunately, the size of the backpropagated exact matrices scales quadratically in both the batch size<sup>6</sup> and the number of layer's output features. However, instead of propagating back the exact Hessian, a batch-averaged version can be used instead to circumvent the batch size scaling (originating from Botev et al. [21]). In combination with structural information about the parameter Hessian, this strategy is used in Botev et al. [21] and Chen et al. [31] to further approximate curvature multiplications, using quantities computed in a single backward pass and then kept in memory for application of the matrix-vector product. We can embed these explicit schemes into our modular approach. To do so, we denote averages over a batch  $\mathbb{B}$  by a bar, for instance  $1/|\mathbb{B}| \sum_{(\mathbf{x}, \mathbf{y}) \in \mathbb{B}} \nabla_{\boldsymbol{\theta}}^2 \ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y}) = \overline{\nabla_{\boldsymbol{\theta}}^2 \ell}$ . The modified backward pass of curvature information during HBP for a module receives a batch average of the Hessian *w.r.t.* the output,  $\overline{\nabla_{\mathbf{z}}^2 \ell}$ , which is used to formulate the matrix-vector product with the batch-averaged parameter Hessian  $\overline{\nabla_{\boldsymbol{\theta}}^2 \ell}$ . An average of the Hessian *w.r.t.* the module input,  $\overline{\nabla_{\mathbf{x}}^2 \ell}$ , is passed back. Existing work [21, 31] differs primarily in the specifics of how this batch average is computed. In HBP, these approximations can be formulated compactly within Equation (4.7). Relations to the cited works are discussed in more detail in Appendix A.2.4. The approximations

6: If samples are processed independently in every module, these matrices have block structure and scale linearly in batch size. Quadratic scaling is caused by transformations across different samples, like batch normalization.

amounting to relations used by Botev et al. [21] read

$$\overline{\nabla_x^2 \ell} \approx \overline{(\mathbf{J}_x z)^\top \nabla_z^2 \ell (\mathbf{J}_x z)} + \sum_k \overline{(\nabla_x^2 z_k) \nabla_{z_k} \ell}, \quad (4.8)$$

and likewise for  $\theta$ . In case of a linear layer  $z(x) = Wx + b$ , this approximation implies the relations  $\overline{\nabla_W^2 \ell} = \overline{x \otimes x^\top \otimes \nabla_z^2 \ell}$ ,  $\overline{\nabla_b^2 \ell} = \overline{\nabla_z^2 \ell}$ , and  $\overline{\nabla_x^2 \ell} = \overline{W^\top (\nabla_z^2 \ell) W}$ . Multiplication by this weight Hessian approximation with a vector  $v$  is achieved by storing  $\overline{x \otimes x^\top}$ ,  $\overline{\nabla_z^2 \ell}$  and performing the required contractions  $v \mapsto (\overline{x \otimes x^\top \otimes \nabla_z^2 \ell})v$ . Note that this approach is not restricted to curvature matrix-vector multiplication routines only. Kronecker structure in the approximation gives rise to optimization methods relying on direct inversion.

A cheaper approximation, used in Chen et al. [31],

$$\overline{\nabla_x^2 \ell} \approx \overline{(\mathbf{J}_x z)^\top \nabla_z^2 \ell (\mathbf{J}_x z)} + \sum_k \overline{(\nabla_x^2 z_k) \nabla_{z_k} \ell}, \quad (4.9)$$

leads to the modified relation  $\overline{\nabla_W^2 \ell} = \overline{x \otimes x^\top \otimes \nabla_z^2 \ell}$  for a linear layer. As this approximation is of the same rank as  $\overline{\nabla_z^2 \ell}$ , which is typically small, CG requires only a few iterations during optimization. It avoids large memory requirements for layers with numerous inputs, since it requires  $\overline{x}$  be stored instead of  $x \otimes x^\top$ .

Transformations that are linear in the module parameters (*e.g.* linear and convolutional layers), possess constant Jacobians *w.r.t.* the module input for each sample (see Table 4.1). Hence, in a network consisting of only these layers, both Equations (4.8) and (4.9) yield the same backpropagated Hessians  $\overline{\nabla_x^2 \ell}$ . This still leaves the degree of freedom for choosing the approximation scheme in the analogous equations for  $\theta$ .

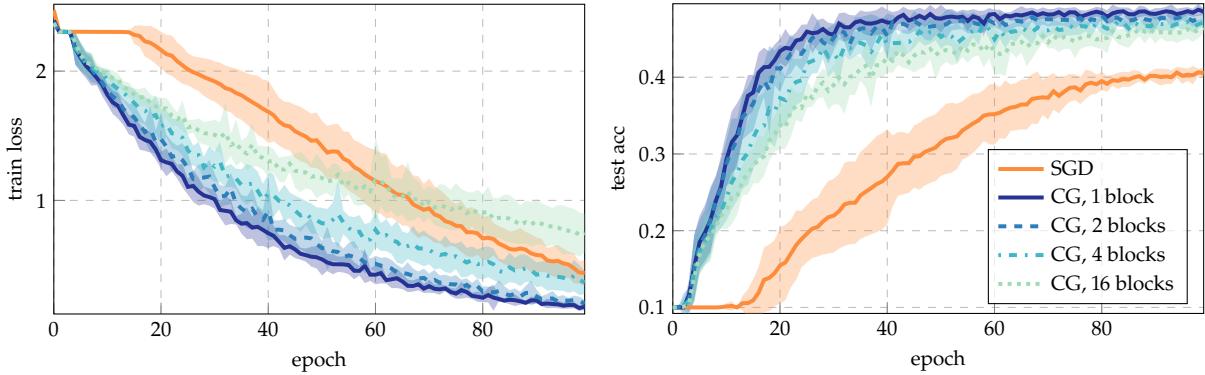
#### Remark

Both strategies for obtaining curvature matrix BDAs (implicit exact matrix-vector multiplications and explicit propagation of approximated curvature) are compatible. Regarding the connection to cited works, we note that the maximally modular structure of our framework changes the nature of these approximations and allows a more flexible formulation and implementation<sup>7</sup>.

<sup>7</sup>: The BackPACK library described in Chapter 5 uses the insights of this section to implement block-diagonal curvature approximations as extensions of gradient backpropagation on the modular level.

## 4.4 Experiments & Implementation

We illustrate the usefulness of incorporating curvature information with the two outlined strategies by experiments with a fully-connected and a convolutional neural network (CNN) on the CIFAR-10 dataset [90]. Following the guidelines of Schneider et al. [146], the training loss is estimated on a random subset of the training set of equal size as the test set. Each experiment is performed for 10 different random seeds and we show the mean values with shaded intervals of one standard deviation.



**Figure 4.4:** SGD and different Newton-style optimizers based on the PCH-abs with batch approximations. The same fully-connected neural network of [31] was used to generate the solid baseline results. Our modular approach allows further splitting the parameter blocks into sub-blocks that can independently be optimized in parallel (dashed lines).

For the loss function we use softmax cross-entropy. Details on the model architectures and hyperparameters are given in [Appendix A.5](#).

### Training Procedure & Update Rule

In comparison to a first-order optimization procedure, the training loop with HBP has to be extended by a single backward pass to backpropagate the batch-averaged or exact loss Hessian. This yields matrix-vector products with a curvature estimate  $C^{(l)}$  for each parameter block  $\theta^{(l)}$  of the network. Parameter updates  $\Delta\theta^{(l)}$  are obtained by applying CG to solve the linear system<sup>8</sup>

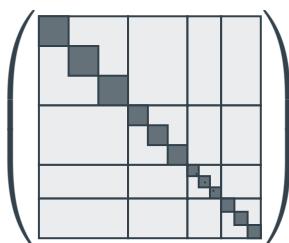
$$\left[ \alpha I + (1 - \alpha) C^{(l)} \right] \Delta\theta^{(l)} = -\nabla_{\theta^{(l)}} \mathcal{L}, \quad (4.10)$$

where  $\alpha$  acts as a step size limitation to improve robustness against noise. The CG routine terminates if the ratio of residual and gradient norm falls below a certain threshold or the maximum number of iterations has been reached. The solution returned by CG is scaled by a learning rate  $\eta$ , and parameters are updated by the relation  $\theta^{(l)} \leftarrow \theta^{(l)} + \eta \Delta\theta^{(l)}$ .

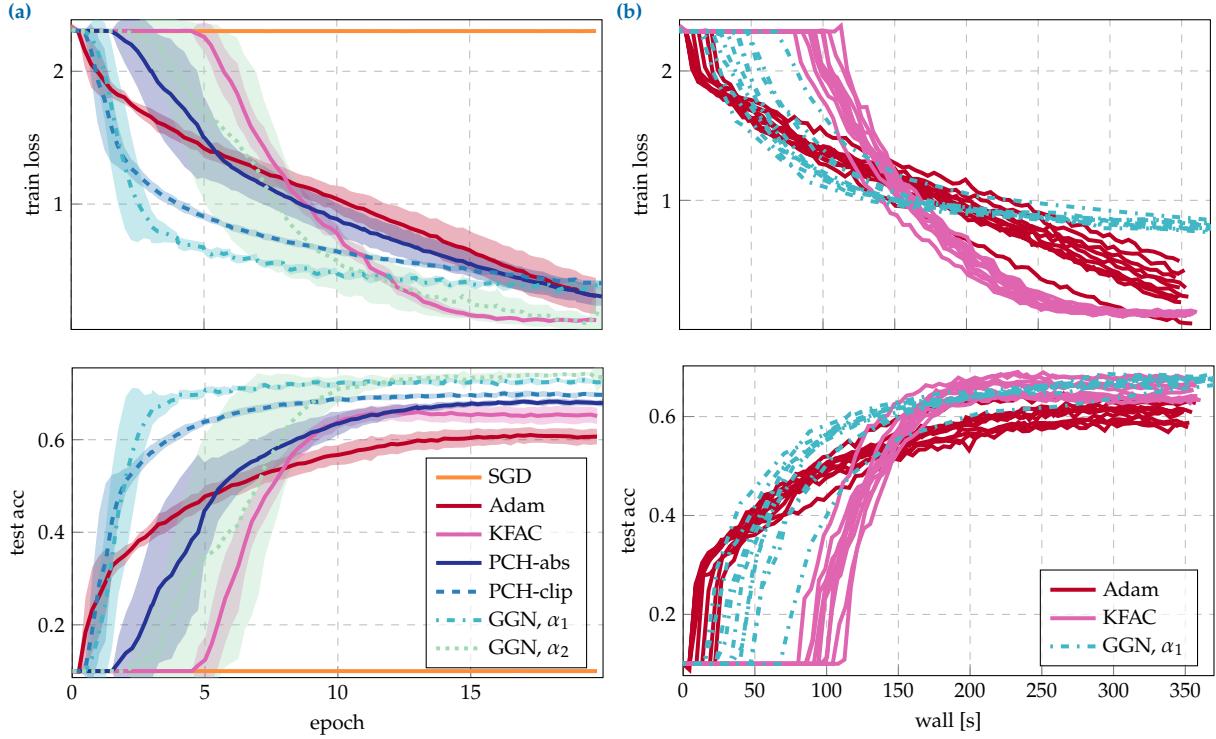
### FCNN, Batch Approximations & Sub-Blocking

The flexibility of HBP is illustrated by extending the results in Chen et al. [31]. Investigations are performed on a fully-connected network with sigmoid activations. Solid lines in [Figure 4.4](#) show the performance of the Newton-style optimizer and momentum SGD in terms of the training loss and test accuracy. The second-order method is capable to escape the initial plateau in fewer iterations.

The modularity of HBP allows for additional parallelism by splitting the linear system [Equation \(4.10\)](#) into smaller sub-blocks ([Figure 4.6](#), which then also need fewer iterations of CG). Doing so only requires a minor modification of the parameter Hessian computation by [Equation \(4.7\)](#). Consequently, we split weights and bias terms row-wise into a specified number of sub-blocks. [Figure 4.4](#) shows performance curves. In the initial phase, the BDA can be split into a larger number of sub-blocks without suffering from a loss in performance. The reduced curvature information



**Figure 4.6:** Illustration of sub-blocking. Here, each diagonal block is split into three blocks of equal size.



**Figure 4.5:** SGD, Adam, KFAC, and Newton-style methods with different exact curvature matrix-vector products (HBP). The architecture is a CNN with sigmoid activations (see Appendix A.5). (a) Comparison of train and test loss over iterations. SGD cannot train the net. (b) Wall-clock time comparison (on an RTX 2080 Ti GPU; same colors realize different random seeds).

is still sufficient to escape the initial plateau. However, larger blocks have to be considered in later stages to further reduce the loss efficiently.

The fact that this switch in modularity is necessary is an argument in favor of the HBP’s flexibility to efficiently realize such switches: for this experiment, the splitting for each block was artificially chosen to illustrate this flexibility. In principle, the splitting could be decided individually for each parameter block, and even changed at run time.

### CNN, Matrix-free Exact Curvature Multiplication

For convolutions, the large number of hidden features prohibits back-propagating a curvature matrix batch average. Instead, we use exact curvature matrix-vector products provided within HBP. The CNN possesses sigmoid activations and cannot be trained by SGD (Figure 4.5a). For comparison with another second-order method, we experiment with a public KFAC implementation [63, 109, see Appendix A.5 for details].

The matrix-free second-order methods progress fast in the initial stage of the optimization. However, progress in later phases stagnates. This may be caused by the limited sophistication of the update rule equation 4.10: if a small value for  $\alpha$  is chosen, the optimizer will perform well in the beginning (GGN,  $\alpha_1$ ). As the gradients become smaller, and hence more noisy, the step size limitation is too optimistic, which leads to a slow-down in optimization progress. A more conservative step size limitation improves the overall performance at the cost of fewer initial progress (GGN,  $\alpha_2$ ). In the training phase where damping is “effective”, our illustrative methods, and KFAC, exhibit better progress per iteration

on the objective than the first-order competitor Adam, underlining the usefulness of curvature even if only computed block-wise.

For an impression on performance in terms of run time, Figure 4.5b compares the wall-clock time of one matrix-free method and the baselines. The HBP-based optimizer can compete with existing methods and offers potential for further improvements, like sub-blocking and parallelized CG. Despite the more adaptive nature of second-order methods, their full power seems to still require adaptive damping, to account for the quality of the local quadratic approximation and restrict the update if necessary. The importance of these techniques to properly adapt the Newton direction has been emphasized in previous works [21, 106, 109] that aim to develop fully fletched second-order optimizers. Such adaptation, however, is beyond the scope of this text.

## 4.5 Conclusion

We have outlined a procedure to compute block-diagonal approximations of different curvature matrices for feedforward neural networks by a scheme that can be realized on top of gradient backpropagation. In contrast to other recently proposed methods, our implementation is aligned with the design of current machine learning frameworks and can flexibly compute Hessian sub-blocks to different levels of refinement. Its modular formulation facilitates closed-form analysis of Hessian diagonal blocks, and unifies previous approaches [21, 31].

Within our framework we presented two strategies: (i) obtaining exact curvature matrix-vector products that have not been accessible before by auto-differentiation (PCH), and (ii) backpropagation of further approximated matrix representations to save computations during training. As for gradient backpropagation, the Hessian backpropagation for different operations can be derived independently of the underlying graph. The extended modules can then be used as a drop-in replacement for existing modules to construct deep neural networks. Internally, backprop is extended by an additional Hessian backward pass through the graph to compute curvature information. It can be performed in parallel to, and reuse the quantities computed in, gradient backpropagation.

## Acknowledgments

The authors would like to thank Frederik Kunstner, Matthias Werner, Frank Schneider, and Agustinus Kristiadi for their constructive feedback on the manuscript, and gratefully acknowledge financial support by the European Research Council through ERC StG Action 757275 / PANAMA; the DFG Cluster of Excellence “Machine Learning - New Perspectives for Science”, EXC 2064/1, project number 390727645; the German Federal Ministry of Education and Research (BMBF) through the Tübingen AI Center (FKZ: 01IS18039A); and funds from the Ministry of Science, Research and Arts of the State of Baden-Württemberg. F. D. is grateful to the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for support.

# BackPACK: Packing More into Backprop

# 5.

## Abstract

Automatic differentiation frameworks are optimized for exactly one thing: computing the average mini-batch gradient. Yet, other quantities such as the variance of the mini-batch gradients or many approximations to the Hessian can, *in theory*, be computed efficiently, and at the same time as the gradient. While these quantities are of great interest to researchers and practitioners, current deep-learning software does not support their automatic calculation. Manually implementing them is burdensome, inefficient if done naïvely, and the resulting code is rarely shared. This hampers progress in deep learning, and unnecessarily narrows research to focus on gradient descent and its variants; it also complicates replication studies and comparisons between newly developed methods that require those quantities, to the point of impossibility. To address this problem, we introduce BackPACK, an efficient framework built on top of PyTorch, that extends the backpropagation algorithm to extract additional information from first- and second-order derivatives. Its capabilities are illustrated by benchmark reports for computing additional quantities on deep neural networks, and an example application by testing several recent curvature approximations for optimization.

## 5.1 Introduction

The success of deep learning and the applications it fuels can be traced to the popularization of automatic differentiation frameworks. Packages like TensorFlow [1], Chainer [165], MXNET [32], and PyTorch [126] provide efficient implementations of parallel, GPU-based gradient computations to a wide range of users, with elegant syntactic sugar.

However, this specialization also has its shortcomings: it assumes the user only wants to compute gradients or, more precisely, the average of gradients across a mini-batch of examples. Other quantities can also be computed with AD at a comparable cost or minimal overhead to the gradient backpropagation pass; for example, approximate second-order information or the variance of gradients within the batch. These quantities are valuable to understand the geometry of deep neural networks, for the identification of free parameters, and to push the development of more efficient optimization algorithms. But researchers who want to investigate their use face a chicken-and-egg problem: AD tools required to go beyond standard gradient methods are not available, but there is no incentive for their implementation in existing deep-learning software as long as no large portion of the users need it.

Second-order methods for deep learning have been continuously investigated for decades [*e.g.* 5, 14, 20, 109]. But still, the standard optimizers used in deep learning remain some variant of stochastic gradient descent (SGD); more complex methods have not found wide-spread, practical

5.1 Introduction . . . . .	63
5.2 Theory & Implementation . . . . .	65
5.3 Evaluation & Benchmarks . . . . .	70
5.4 Experiments . . . . .	70
5.5 Conclusion . . . . .	73

Code and experiments available at the  
Github repositories

[f-dangel/backpack](https://github.com/f-dangel/backpack),  
[f-dangel/backpack-experiments](https://github.com/f-dangel/backpack-experiments)



use. This is in stark contrast to domains like convex optimization and generalized linear models, where second-order methods are the default. There may of course be good scientific reasons for this difference; maybe second-order methods do not work well in the (non-convex, stochastic) setting of deep learning. And the computational cost associated with the high dimensionality of deep models may offset their benefits. Whether these are the case remains somewhat unclear though, because a much more direct road-block is that these methods are so complex to implement that few practitioners ever try them out.

Recent approximate second-order methods such as KFAC [109] show promising results, even on hard deep learning problems [166]. Their approach, based on the earlier work of Schraudolph [148], uses the structure of the network to compute approximate second-order information in a way that is similar to gradient backpropagation. This work sparked a new line of research to improve the second-order approximation [21, 54, 63, 108]. However, all of these methods require low-level applications of automatic differentiation to compute quantities other than the averaged gradient. It is a daunting task to implement them from scratch. Unless users spend significant time familiarizing themselves with the internals of their software tools, the resulting implementation is often inefficient, which also puts the original usability advantage of those packages into question. Even motivated researchers trying to develop new methods, who need not be expert software developers, face this problem. They often end up with methods that cannot compete in run time, not necessarily because the method is inherently bad, but because the implementation is not efficient. New methods are also frequently not compared to their predecessors and competitors because they are so hard to reproduce. Authors do not want to represent the competition in an unfair light caused by a bad implementation.

Another example is offered by a recent string of research to adapt to the *stochasticity* induced by mini-batch sampling. An empirical estimate of the (marginal) variance of the gradients within the batch has been found to be theoretically and practically useful for adapting hyperparameters like learning rates [105] and batch sizes [11], or regularize first-order optimization [10, 85, 94]. To get such a variance estimate, one simply has to square, then sum, the individual gradients after the backpropagation, but before they are aggregated to form the average gradient. Doing so should have negligible cost *in principle*, but is programmatically challenging in the standard packages.

1: See e.g. the Github issues [github.com/pytorch/pytorch/issues/1407](https://github.com/pytorch/pytorch/issues/1407), [7786](https://github.com/pytorch/pytorch/issues/7786), [8897](https://github.com/pytorch/pytorch/issues/8897) and forum discussions [discuss.pytorch.org/t/1433](https://discuss.pytorch.org/t/1433), [8405](https://discuss.pytorch.org/t/8405), [15270](https://discuss.pytorch.org/t/15270), [17204](https://discuss.pytorch.org/t/17204), [19350](https://discuss.pytorch.org/t/19350), [24955](https://discuss.pytorch.org/t/24955).

Members of the community have repeatedly asked for such features<sup>1</sup> but the established AD frameworks have yet to address such requests, as their focus has been—rightly—on improving their technical backbone. Features like those outlined above are not generally defined for arbitrary functions, but rather emerge from the specific structure of machine learning applications. General AD frameworks can not be expected to serve such specialist needs. This does not mean, however, that it is impossible to efficiently realize such features within these frameworks: in essence, backpropagation is a technique to compute multiplications with Jacobians. Methods to extract second-order information [110] or individual gradients from a mini-batch [59] have been known to a small group of specialists; they are just rarely discussed or implemented.

### Computing the gradient with PyTorch ...

```
X, y      = load_mnist_data()
model    = Linear(784, 10)
lossfunc = CrossEntropyLoss()

loss      = lossfunc(model(X), y)

loss.backward()

for param in model.parameters():
    print(param.grad)
```

### ... and the variance with BackPACK

```
X, y      = load_mnist_data()
model    = extend(Linear(784, 10))
lossfunc = extend(CrossEntropyLoss())

loss      = lossfunc(model(X), y)
with backpack(Variance()):
    loss.backward()

for param in model.parameters():
    print(param.grad)
    print(param.var)
```

#### 5.1.1 Our Contribution

To address this need for a specialized framework focused on machine learning, we propose a framework for the implementation of generalized backpropagation to compute additional quantities. The structure is based on the conceptual work of Dangel et al. [37] for modular backpropagation. This framework can be built on top of existing graph-based backpropagation modules; we provide an implementation on top of PyTorch, coined BackPACK, available at

<https://f-dangel.github.io/backpack/>.

The initial release supports efficient computation of individual gradients from a mini-batch, their  $L_2$  norm, an estimate of the variance, as well as diagonal and Kronecker factorizations of the generalized Gauss-Newton (GGN) matrix (see Table 5.1 for a feature overview). The library was designed to be minimally verbose to the user, easy to use (see Procedure 5.1), and to have low overhead (see Section 5.3). While other researchers are aiming to improve the flexibility of AD systems [23, 77, 78], our goal with this package is to provide access to quantities that are only byproducts of the backpropagation pass, rather than gradients themselves.

To illustrate the capabilities of BackPACK, we use it to implement preconditioned gradient descent optimizers with diagonal approximations of the GGN and recent Kronecker factorizations KFAC [109], KFLR, and KFRA [21]. Our results show that the curvature approximations based on Monte Carlo (MC) estimates of the GGN, the approach used by KFAC, give similar progress per iteration to their more accurate counterparts, but being much cheaper to compute. While the naïve update rule we implement does not surpass first-order baselines such as SGD with momentum and Adam [87], its implementation with various curvature approximations is made straightforward.

**Procedure 5.1:** BackPACK integrates with PyTorch to seamlessly extract more information from the backward pass. Instead of the variance (or alongside it, in the same pass), BackPACK can compute individual gradients in the mini-batch, their  $L_2$  norm and 2<sup>nd</sup> moment. It can also compute curvature approximations like diagonal or Kronecker factorizations of the GGN such as KFAC, KFLR & KFRA.

## 5.2 Theory & Implementation

We will distinguish between quantities that can be computed from information already present during a traditional backward pass (which we suggestively call *first-order extensions*), and quantities that need additional information (termed *second-order extensions*). The former group contains additional statistics such as the variance of the gradients within the mini-batch or the  $L_2$  norm of the gradient for each sample. Those can

be computed with minimal overhead during the backprop pass. The latter class contains approximations of second-order information, like the diagonal or Kronecker factorization of the generalized Gauss-Newton (GGN) matrix, which require the propagation of additional information through the graph. We will present those two classes separately:

#### First-order extensions

Extract more from the standard backward pass.

- ▶ Individual gradients from a mini-batch
- ▶  $L_2$  norm of the individual gradients
- ▶ Diagonal covariance and 2<sup>nd</sup> moment

#### Second-order extensions

Propagate new information along the graph.

- ▶ Diagonal of the GGN and the Hessian
- ▶ KFAC [109]
- ▶ KFRA and KFLR [21]

These quantities are only defined, or reasonable to compute, for a subset of models: the concept of individual gradients for each sample in a mini-batch or the estimate of the variance requires the loss for each sample to be independent. While such functions are common in machine learning, not all neural networks fit into this category. *E.g.* if the network uses batch normalization [79], the individual gradients in a mini-batch are correlated. Then, the variance is not meaningful anymore, and computing the individual contribution of a sample to the mini-batch gradient or the GGN becomes prohibitive. For those reasons, and to limit the scope of the project for version 1.0, BackPACK currently restricts the type of models it accepts. The supported models are traditional feedforward networks that can be expressed as a *sequence of modules*, for example a sequence of convolutional, pooling, linear and activation layers. Recurrent networks like LSTMs [71] or residual networks [68] are not yet supported, but the framework can be extended to cover them<sup>2</sup>.

2: BackPACK has been continuously developed since the initial release. Note-worthy added features include:

- ▶ New extensions: per-sample Hessian/GGN diagonal (version 1.3), matrix-free multiplication with block-diagonal curvature matrices from [Chapter 4](#) (version 1.2), and GGN low-rank factors (versions 1.4, 1.5), see [Chapter 7](#).
- ▶ Broader support of modules and hyperparameters, especially basic support for residual and recurrent networks (version 1.4).

We assume a sequential model  $f_{\theta} : \times \mathbb{X} \rightarrow \mathbb{F}$  and a dataset of  $N$  samples  $(x_n, y_n) \in \mathbb{X} \times \mathbb{Y}$  with  $n = 1, \dots, N$ . The model maps each sample  $x_n$  to a prediction  $f_{\theta}(x_n)$  using some parameters  $\theta \in \Theta$ . The predictions are evaluated with a loss function  $\ell : \mathbb{F} \times \mathbb{Y} \rightarrow \mathbb{R}$ , for example the softmax cross-entropy ([Equation \(2.4\)](#)), which compares them to the ground truth  $y_n$ . This leads to the objective function  $\mathcal{L} : \Theta \rightarrow \mathbb{R}$ ,

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(f_{\theta}(x_n), y_n). \quad (5.1)$$

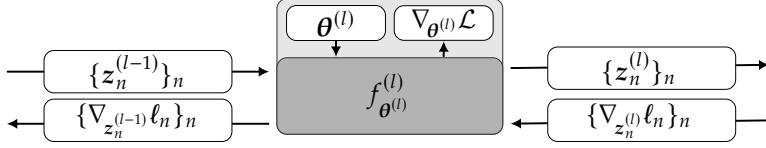
As a shorthand, we will use  $\ell_n(\theta) = \ell(f_{\theta}(x_n), y_n)$  for the loss and  $f_n(\theta) = f_{\theta}(x_n)$  for the model output of individual samples. Our goal is to provide more information about the derivatives of  $\{\ell_n(\theta)\}_{n=1}^N$  w.r.t. the parameters  $\theta$  of the model  $f_{\theta}$ .

### 5.2.1 Primer on Backpropagation

ML libraries with integrated automatic differentiation use the modular structure of  $f_n(\theta)$  to compute derivatives (see [13] for an overview). If  $f_{\theta}$  is a sequence of  $L$  transformations, it can be expressed as

$$f_n(\theta) = \left( f_{\theta^{(L)}}^{(L)} \circ \dots \circ f_{\theta^{(1)}}^{(1)} \right) (x_n), \quad (5.2)$$

where  $f_{\theta^{(l)}}^{(l)}$  is the  $l$ th transformation with parameters  $\theta^{(l)}$ , such that  $\theta = (\theta^{(1)\top}, \dots, \theta^{(L)\top})^\top$ . The loss function can also be seen as another



**Figure 5.1:** Schematic representation of the standard backpropagation pass for module  $l$  with  $N$  samples.

transformation, appended to the network. Let  $z_n^{(l-1)}, z_n^{(l)}$  denote the input and output of the operation  $f_{\theta^{(l)}}$  for sample  $n$ , such that  $z_n^{(0)}$  is the original data and  $z_n^{(1)}, \dots, z_n^{(L)}$  represent the transformed output of each layer, leading to the computation graph

$$z_n^{(0)} \xrightarrow{f_{\theta^{(1)}}(z_n^{(0)})} z_n^{(1)} \xrightarrow{f_{\theta^{(2)}}(z_n^{(1)})} \dots \xrightarrow{f_{\theta^{(L)}}(z_n^{(L-1)})} z_n^{(L)} \xrightarrow{\ell(z_n^{(L)}, y_n)} \ell_n(\boldsymbol{\theta}).$$

To compute the gradient of  $\ell_n$  w.r.t.  $\boldsymbol{\theta}^{(l)}$ , one unrolls the chain rule,

$$\begin{aligned} \nabla_{\boldsymbol{\theta}^{(l)}} \ell_n(\boldsymbol{\theta}) &= \left( J_{\boldsymbol{\theta}^{(l)}} z_n^{(l)} \right)^T \left( J_{z_n^{(l)}} z_n^{(l+1)} \right)^T \cdots \left( J_{z_n^{(L-1)}} z_n^{(L)} \right)^T \nabla_{z_n^{(L)}} \ell_n(\boldsymbol{\theta}) \\ &= \left( J_{\boldsymbol{\theta}^{(l)}} z_n^{(l)} \right)^T \nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta}), \end{aligned} \quad (5.3)$$

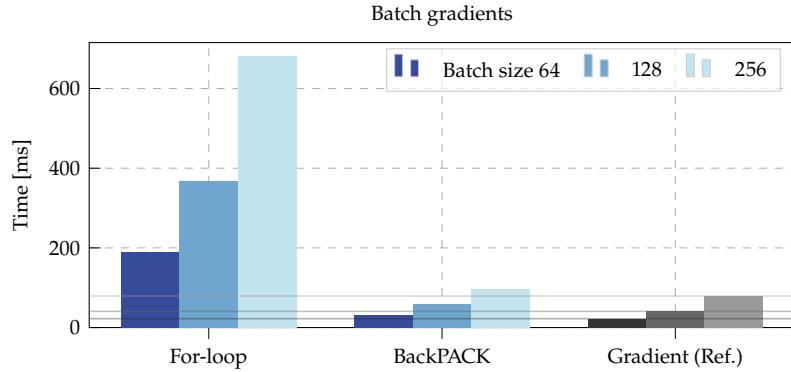
where  $J_a b$  is the Jacobian of  $b$  w.r.t.  $a$ ,  $[J_a b]_{i,j} = \partial[b]_i / \partial[a]_j$  (see [Definition 2.4](#)). A similar expression exists for the module inputs  $z_n^{(l-1)}$ :  $\nabla_{z_n^{(l-1)}} \ell_n(\boldsymbol{\theta}) = (J_{z_n^{(l-1)}} z_n^{(l)})^T \nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta})$ . This recursive structure makes it possible to extract the gradient by propagating the gradient of the loss. In the backpropagation algorithm, a module  $l$  receives the loss gradient w.r.t. its output,  $\nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta})$ . It then extracts the gradient with respect to its parameters and inputs,  $\nabla_{\boldsymbol{\theta}^{(l)}} \ell_n(\boldsymbol{\theta})$  and  $\nabla_{z_n^{(l-1)}} \ell_n(\boldsymbol{\theta})$ , according to [Equation \(5.3\)](#). The gradient w.r.t. its input is sent further down the graph. This process, illustrated in [Figure 5.1](#), is repeated for each transformation until all gradients are computed. To implement backpropagation, each module only needs to know how to multiply with its Jacobians.

For second-order quantities, we rely on the work of Mizutani and Dreyfus [110] and Dangel et al. [37], who showed that a scheme similar to [Equation \(5.3\)](#) exists for the *block diagonal* of the Hessian. A block w.r.t. the parameters of a module,  $\nabla_{\boldsymbol{\theta}^{(l)}}^2 \ell_n(\boldsymbol{\theta})$ , can be obtained by the recursion

$$\begin{aligned} \nabla_{\boldsymbol{\theta}^{(l)}}^2 \ell_n(\boldsymbol{\theta}) &= \left( J_{\boldsymbol{\theta}^{(l)}} z_n^{(l)} \right)^T \nabla_{z_n^{(l)}}^2 \ell_n(\boldsymbol{\theta}) \left( J_{\boldsymbol{\theta}^{(l)}} z_n^{(l)} \right) \\ &\quad + \sum_j \left( \nabla_{\boldsymbol{\theta}^{(l)}}^2 \left[ z_n^{(l)} \right]_j \right) \left[ \nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta}) \right]_j. \end{aligned} \quad (5.4)$$

A similar relation holds for the module's output Hessian  $\nabla_{z_n^{(l)}}^2 \ell_n(\boldsymbol{\theta})$ .

Both backpropagations of [Equations \(5.3\)](#) and [\(5.4\)](#) hinge on the multiplication by Jacobians to both vectors and matrices. However, the design of AD limits the application of Jacobians to vectors only. This prohibits the exploitation of vectorization in the matrix case, which is needed for second-order information. The lacking flexibility of Jacobians is one motivation for our work. Since all quantities needed to compute statistics of the derivatives are already computed during the backward pass, another motivation is to provide access to them at minor overhead.



**Figure 5.2:** Computing individual gradients in a batch using a for-loop (*i.e.* one individual forward and backward pass per sample) or using vectorized operations with BackPACK. The plot shows computation time, comparing to a traditional gradient computation, on the 3c3d network (see Section 5.4) for the CIFAR-10 dataset [146].

### 5.2.2 First-order Extensions

As the principal first-order extension, consider computing the *per-sample* gradients in a batch of size  $N$ . These individual gradients are implicitly computed during a traditional backward pass because the batch gradient is their sum, but they are not directly accessible. The naïve way to compute  $N$  individual gradients is to do  $N$  separate forward and backward passes. This (inefficiently) replaces every matrix-matrix multiplication by  $N$  matrix-vector multiplications. BackPACK batches computations to obtain large efficiency gains, as illustrated by Figure 5.2<sup>3</sup>.

As the quantities necessary to compute the individual gradients are already propagated through the computation graph, we can reuse them by inserting code in the standard backward pass. With access to this information, before it is cleared for memory efficiency, BackPACK computes the Jacobian-multiplications for each sample

$$\{\nabla_{\theta^{(l)}} \ell_n(\theta)\}_{n=1}^N = \left\{ \left( J_{\theta^{(l)}} z_n^{(l)} \right)^T \nabla_{z_n^{(l)}} \ell_n(\theta) \right\}_{n=1}^N, \quad (5.5)$$

without summing the result—see Figure 5.3 for a schematic representation. This duplicates some of the computation performed by the backpropagation, as the Jacobian is applied twice (once by PyTorch and BackPACK with and without summation over the samples, respectively). However, the associated overhead is small compared to the for-loop approach: the major computational cost arises from the propagation of information required for each layer, rather than the formation of the gradient *within* each layer.

This scheme for individual gradient computation is the basis for all first-order extensions. In this direct form, however, it is expensive in memory: if the model is  $D$ -dimensional, storing  $\mathcal{O}(ND)$  elements is prohibitive for large batches. For the variance, 2<sup>nd</sup> moment and  $L_2$  norm, BackPACK takes advantage of the Jacobian’s structure to directly compute them without forming the individual gradient, reducing memory overhead. See Appendix B.1.1 for details.

### 5.2.3 Second-order Extensions

Second-order extensions require propagation of more information through the graph. As an example, we will focus on the GGN matrix

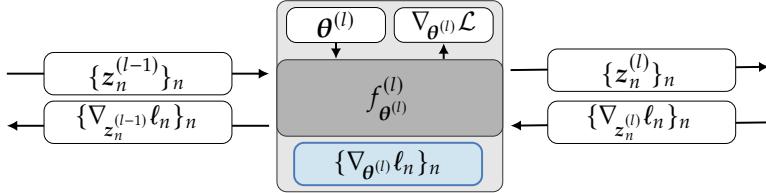
**Example 5.1 (Symmetric decomposition of the softmax cross-entropy loss Hessian [122])** Consider the softmax cross-entropy loss (Equation (2.2)) Hessian from Table Table 4.1. Its symmetric decomposition  $S \in \mathbb{R}^{C \times C}$  is

$$S = (\mathbf{I} - p \mathbf{1}^\top) \text{diag}(\sqrt{p}) = \text{diag}(\sqrt{p}) - p \sqrt{p} \mathbf{p}^\top, \quad (5.6)$$

where  $\sqrt{p(f)} = p(f)^{\odot 1/2}$ . It satisfies  $\nabla_f^2 \ell(f, y) = SS^\top$ ,

$$\begin{aligned} SS^\top &= [\text{diag}(\sqrt{p}) - p \sqrt{p} \mathbf{p}^\top] \\ &\quad [\text{diag}(\sqrt{p}) - \sqrt{p} \mathbf{p}^\top] \\ &= \text{diag}(p) - 2p \mathbf{p}^\top + p \sqrt{p} \mathbf{p}^\top \sqrt{p} \mathbf{p}^\top \\ &= \text{diag}(p) - p \mathbf{p}^\top, \end{aligned}$$

using that  $p$ ’s elements sum to one, *i.e.*  $\sqrt{p}^\top \sqrt{p} = 1$ . This is the expression from Table 4.1.



**Figure 5.3:** Schematic representation of the individual gradients' extraction in addition to the standard backward pass at the  $l$ th module for  $N$  samples.

[148]. It is guaranteed to be PSD and is a reasonable approximation of the Hessian near the minimum, which motivates its use in approximate second-order methods. For popular loss functions, it coincides with the Fisher information matrix used in natural gradient methods [5]; for a more in depth discussion of the equivalence, see Section 3.2.4 and the reviews of Martens [107] and Kunstner et al. [92]. For an objective function that can be written as the composition of a loss function  $\ell$  and a model  $f_\theta$ , such as Equation (5.1), the GGN of  $1/N \sum_n \ell(f_\theta(x_n), y_n)$  is

$$G(\theta) = \frac{1}{N} \sum_n (\mathbf{J}_\theta f_n)^\top \nabla_{f_n}^2 \ell(f_n, y_n) (\mathbf{J}_\theta f_n). \quad (5.7)$$

The full matrix is too large to compute and store. Current approaches focus on its diagonal blocks, where each block corresponds to a layer in the network. Every block itself is further approximated, for example using a Kronecker factorization. The approach used by BackPACK for their computation is a refinement of the *Hessian Backpropagation equations* of Dangel et al. [37]. It relies on two insights: firstly, the computational bottleneck in the GGN's computation is the multiplication with the Jacobian of the network,  $\mathbf{J}_\theta f_n$ , while the network output Hessian is easy to compute for most popular loss functions. Secondly, it is not necessary to compute and store each of the  $N D \times D$  matrices for a network with  $D$  parameters, as Equation (5.7) is a quadratic expression. Given a symmetric factorization  $S_n$  of the Hessian,  $S_n S_n^\top = \nabla_{f_n}^2 \ell(f_n, y_n)$  (e.g. Example 5.1), it is sufficient to compute  $(\mathbf{J}_\theta f_n)^\top S_n$  and square the result. A network output is typically small compared to its inner layers; networks on CIFAR-100 need  $C = 100$  class outputs but could use convolutional layers with more than 100,000 parameters.

The factorization leads to a  $D \times C$  matrix, which makes it possible to efficiently compute GGN block diagonals. Also, the computation is very similar to that of a gradient, which computes  $(\mathbf{J}_\theta f_n)^\top \nabla_{f_n} \ell_n$ . A module  $f_{\theta^{(l)}}$  receives the symmetric factorization of the GGN w.r.t. its output,  $z_n^{(l)}$ , and multiplies it with the Jacobians w.r.t. the parameters  $\theta^{(l)}$  and inputs  $z_n^{(l-1)}$  to produce a symmetric factorization of the GGN w.r.t. the parameters and inputs, as shown in Figure 5.4.

This propagation serves as the basis of the second-order extensions. If the full symmetric factorization is not wanted, for memory reasons, it is possible to extract more specific information such as the diagonal. If  $B$  is the symmetric factorization for a GGN block, the diagonal can be computed as  $[\text{diag}(BB^\top)]_i = [BB^\top]_{i,i} = \sum_j [B]_{i,j}^2$ .

This framework can be used to extract the main Kronecker factorizations of the GGN, KFAC and KFLR, which we extend to convolution using

**Example 5.2 (MC approximation of the softmax-cross entropy loss Hessian)** Consider the softmax cross-entropy loss (Equation (2.2)) Hessian from Table 4.1. An MC approximation of the symmetric decomposition (Equation (5.6)) is constructed by the vectors

$$\tilde{s} = \tilde{y}(c) - p(f) \quad (5.8)$$

with  $\tilde{y}(c) = \text{onehot}(c)$  and  $c$  drawn from a categorical distribution implied by the softmax-probabilities,  $c \sim \text{Cat}(c; p(f))$ . The random vector  $\tilde{y}$  satisfies  $\mathbb{E}_c [\tilde{y}] = p(f)$  and  $\mathbb{E}_c [\tilde{y}\tilde{y}^\top] = \text{diag}[p(f)]$ . With these properties, we can show  $\mathbb{E}_{\tilde{s}} [\tilde{s}\tilde{s}^\top] = \nabla_f^2 \ell$ , i.e. the expected outer product of Equation (5.8) is the Hessian,

$$\begin{aligned} \mathbb{E}_{\tilde{s}} [\tilde{s}\tilde{s}^\top] &= \mathbb{E}_c [\tilde{y}\tilde{y}^\top] - \mathbb{E}_c [\tilde{y}] p^\top \\ &\quad - p \mathbb{E}_c [\tilde{y}^\top] + pp^\top \\ &= \text{diag}(p) - pp^\top. \end{aligned}$$

Instead of using the  $C \times C$  matrix square root  $S$ , we can draw  $M < C$  samples  $c_1, \dots, c_M$  and stack their  $\tilde{s}$ -vectors into a smaller  $C \times M$  matrix

$$\tilde{S} = \frac{1}{\sqrt{M}} (\tilde{s}(c_1) \dots \tilde{s}(c_M)) \quad (5.9)$$

that approximates Equation (5.6),

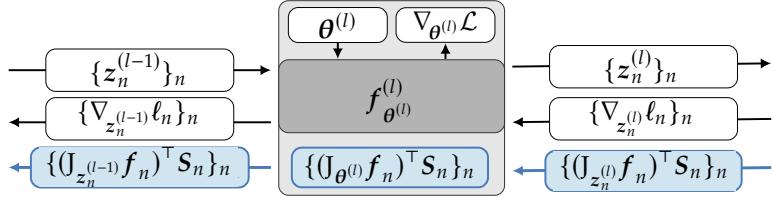
$$\begin{aligned} \tilde{S}\tilde{S}^\top &= \frac{1}{M} \sum_{i=1}^M \tilde{s}(c_i)\tilde{s}(c_i)^\top \\ &\approx \mathbb{E}_{\tilde{s}} [\tilde{s}\tilde{s}^\top] = SS^\top, \end{aligned}$$

using less memory. From the connection between Fisher and Hessian, Equation (5.8) are ‘would-be gradients’ under targets sampled from the model’s likelihood, i.e.  $\tilde{s}^\top = \mathbf{J}_f \ell(f, \tilde{y})$  with  $\tilde{y} \sim q(\cdot | f)$  and the Jacobian from Table 2.2.

**Figure 5.4:** Schematic of the additional backward pass to compute a symmetric factorization of the GGN.

$$G(\theta) = \frac{1}{N} \sum_n (\mathbf{J}_\theta f_n)^\top S_n S_n^\top (\mathbf{J}_\theta f_n)$$

alongside the gradient at the  $l$ th module, for  $N$  samples.



the approach of Grosse and Martens [63]. The important difference between the two methods is the initial matrix factorization  $S_n$ . Using a full symmetric factorization of the initial Hessian,  $S_n S_n^\top = \nabla_{f_n}^2 \ell_n$ , yields the KFLR approximation. KFAC uses an MC approximation by sampling a vector  $s_n$  such that  $\mathbb{E}_{s_n} [s_n s_n^\top] = \nabla_{f_n}^2 \ell_n$  (see Example 5.2). KFLR is therefore more precise but more expensive than KFAC, especially for networks with high-dimensional outputs, which is reflected in our benchmark on CIFAR-100 in Section 5.3. The technical details on how Kronecker factors are extracted and information is propagated for second-order BackPACK extensions are documented in Appendix B.1.2.

## 5.3 Evaluation & Benchmarks

4: 3c3d is a sequence of three convolutions and three dense linear layers with 895,210 parameters. All-CNN-C is a sequence of nine convolutions with 1,387,108 parameters.

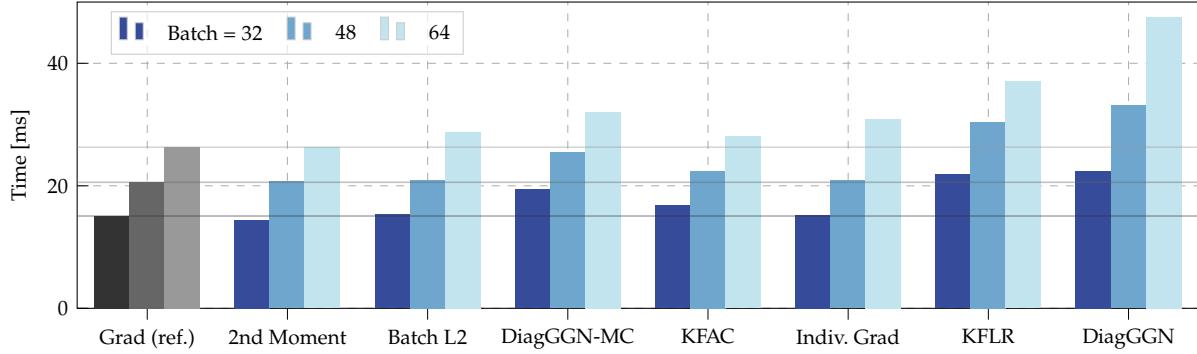
We benchmark the overhead of BackPACK on CIFAR-10 and CIFAR-100, using the 3c3d network and the All-CNN-C network of Springenberg et al. [157] provided by DeepOBS [146]<sup>4</sup>. Figure 5.5 shows the results.

For first-order extensions, the computation of individual gradients from a mini-batch adds noticeable overhead due to the additional memory requirements to store them. But more specific quantities such as the  $L_2$  norm, 2<sup>nd</sup> moment and variance can be extracted efficiently. Regarding second-order extensions, the GGN computation can be expensive for nets with large outputs like CIFAR-100, regardless of the approximation being diagonal or Kronecker-factored. Thankfully, the MC approximation used by KFAC, which we also implement for a diagonal approximation, can be computed at minimal overhead—much less than two backward passes. This last point is encouraging, as our optimization experiment in Section 5.4 suggest that this approximation is reasonably accurate.

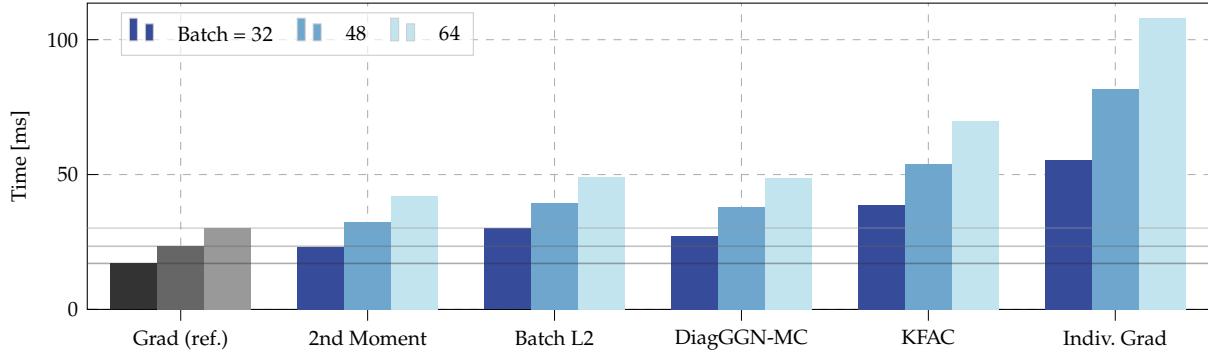
## 5.4 Experiments

To illustrate the utility of BackPACK, we implement preconditioned gradient descent optimizers using diagonal and Kronecker approximations of the GGN. To our knowledge, and despite their apparent simplicity, results using diagonal approximations or the naïve damping update rule we chose have not been reported in publications so far. However, this section is not meant to introduce a bona-fide new optimizer. Our goal is to show that BackPACK can enable research of this kind. The update rule we implement uses a curvature matrix  $C(\theta_t^{(l)})$ , which could be a

(a) 3c3d, CIFAR-10



(b) All-CNN-C, CIFAR-100



**Figure 5.5:** Overhead benchmark for computing the gradient and first- or second-order extensions on real networks, compared to just the gradient. Most quantities add little overhead. KFLR and DiagGGN propagate 100× more information than KFAC and DiagGGN-MC on CIFAR-100 and are two orders of magnitude slower. We report benchmarks on those, and the Hessian’s diagonal, in Appendix B.2.

diagonal or Kronecker factorization of the GGN blocks, and a damping parameter  $\lambda$  to precondition the gradient:

$$\boldsymbol{\theta}_{t+1}^{(l)} = \boldsymbol{\theta}_t^{(l)} - \eta \left( \mathbf{C}(\boldsymbol{\theta}_t^{(l)}) + \lambda \mathbf{I} \right)^{-1} \nabla_{\boldsymbol{\theta}_t^{(l)}} \mathcal{L}(\boldsymbol{\theta}_t^{(l)}), \quad l = 1, \dots, L. \quad (5.10)$$

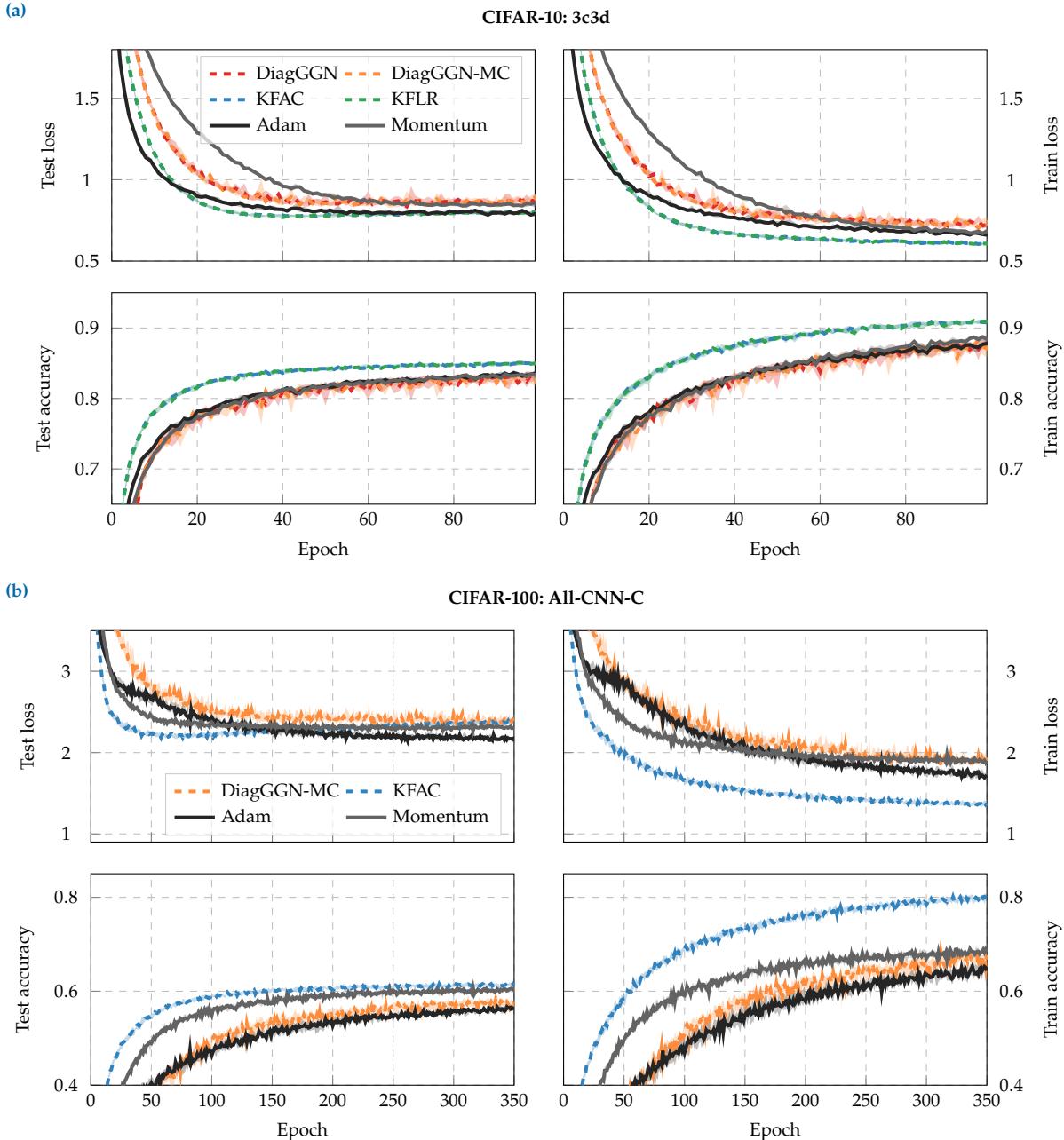
We run the update rule with the following approximations of the generalized Gauss-Newton: the exact diagonal (DiagGGN) and an MC estimate (DiagGGN-MC), and the Kronecker factorizations KFAC [109], KFLR and KFRA<sup>5</sup> [21]. The inversion required by the update rule is straightforward for the diagonal curvature. For the Kronecker-factored quantities, we use the approximation introduced by [109] (see Appendix B.3.3).

These curvature estimates are tested for the training of deep neural networks by running the corresponding optimizers on the main test problems of the benchmarking suite DeepOBS<sup>6</sup> [146]. We use the setup (batch size, number of training epochs) of DeepOBS’ baselines, and tune the learning rate  $\eta$  and damping parameter  $\lambda$  with a grid search for each optimizer (details in Appendix B.3.2). The best hyperparameter settings is chosen according to the final accuracy on a validation set. We report the median and quartiles of the performance for ten random seeds.

Figure 5.6a shows the results for the 3c3d network trained on CIFAR-10. The optimizers that leverage Kronecker-factored curvature approximations beat the baseline performance in terms of per-iteration progress on the training loss, training and test accuracy. Using the same hyperparam-

5: KFRA was not originally designed for convolutions; we extend it using the Kronecker factorization of Grosse and Martens [63]. While it can be computed for small networks on MNIST, which we report in Appendix B.3.4, the approximate backward pass of KFRA does not seem to scale to large convolution layers.

6: [deepobs.github.io](https://deepobs.github.io). We cannot run BackPACK on all test problems in this benchmark due to the limitations outlined in Section 5.2. Despite this limitation, we still run on models spanning a representative range of image classification problems.



**Figure 5.6:** Median performance with shaded quartiles of the DeepOBS benchmark for (a) 3c3d (895,210 parameters) on CIFAR-10 and (a) All-CNN-C (1,387,108 parameters) on CIFAR-100. Solid lines show DeepOBS' baselines of momentum SGD and Adam.

eters, there is little difference between KFAC and KFLR, or DiagGGN and DiagGGN-MC. Given that the quantities based on MC-sampling are considerably cheaper, this experiment suggests it being an important technique to reduce the computational burden of curvature proxies.

Figure 5.6b shows benchmarks for the All-CNN-C network trained on CIFAR-100. Due to the high-dimensional output, the curvatures using a full matrix propagation rather than an MC sample cannot be run on this problem due to memory issues. Both DiagGGN-MC and KFAC can compete with the baselines in terms of progress per iteration. As the update rule we implemented is simplistic on purpose, this is promising for future applications of second-order methods that can more efficiently use the additional information given by curvature approximations.

Feature	Details
Individual gradients	$\frac{1}{N} \nabla_{\theta^{(l)}} \ell_n(\theta), \quad n = 1, \dots, N$
Batch variance	$\frac{1}{N} \sum_{n=1}^N [\nabla_{\theta^{(l)}} \ell_n(\theta)]_j^2 - [\nabla_{\theta^{(l)}} \mathcal{L}(\theta)]_j^2$
2 <sup>nd</sup> moment	$\frac{1}{N} \sum_{n=1}^N [\nabla_{\theta^{(l)}} \ell_n(\theta)]_j^2, \quad j = 1, \dots, d^{(l)}$ .
Indiv. gradient $L_2$ norm	$\left\  \frac{1}{N} \nabla_{\theta^{(l)}} \ell_n(\theta) \right\ _2^2, \quad n = 1, \dots, N$
DiagGGN	$\text{diag}(\mathbf{G}(\theta^{(l)}))$
DiagGGN-MC	$\text{diag}(\tilde{\mathbf{G}}(\theta^{(l)}))$
Hessian diagonal	$\text{diag}(\nabla_{\theta}^2 \mathcal{L}(\theta))$
KFAC	$\tilde{\mathbf{G}}^{(l)}(\theta^{(l)}) \approx \mathbf{A}^{(l)} \otimes \mathbf{B}_{\text{KFAC}}^{(l)}$
KFLR	$\mathbf{G}^{(l)}(\theta^{(l)}) \approx \mathbf{A}^{(l)} \otimes \mathbf{B}_{\text{KFLR}}^{(l)}$
KFRA	$\mathbf{G}^{(l)}(\theta^{(l)}) \approx \mathbf{A}^{(l)} \otimes \mathbf{B}_{\text{KFRA}}^{(l)}$

**Table 5.1:** Overview of the features supported in the first release of BackPACK.

## 5.5 Conclusion

Machine learning’s coming-of-age has been accompanied, and in part driven, by a maturing of the software ecosystem. This has drastically simplified the lives of developers and researchers alike, but has also crystallized parts of the algorithmic landscape. This has dampened research in cutting-edge areas that are far from mature, like second-order optimization for deep neural networks. To ensure that good ideas can bear fruit, researchers must be able to compute new quantities without an overwhelming software development burden. To support research and development in optimization for deep learning, we have introduced BackPACK, an efficient implementation in PyTorch of recent conceptual advances and extensions to backpropagation (Table 5.1 lists all features). BackPACK enriches the syntax of AD packages to offer additional observables to optimizers beyond the batch-averaged gradient. Our experiments demonstrate that BackPACK’s implementation offers drastic efficiency gains over the kind of naïve implementation within reach of the typical researcher. As a demonstrative example, we “invented” a few optimization routines that, without BackPACK, would require demanding implementation work and can now be tested with ease. We hope that studies like this allow BackPACK to help mature the ML software ecosystem further.

## Acknowledgments

We would like to thank Aaron Bahde, Ludwig Bald, and Frank Schneider for their help with DeepOBS and Lukas Balles, Simon Bartels, Filip de Roos, Tim Fischer, Nicolas Krämer, Agustinus Kristiadi, Frank Schneider, Jonathan Wenger, and Matthias Werner for constructive feedback.

The authors gratefully acknowledge financial support by the European Research Council through ERC StG Action 757275 / PANAMA; the DFG Cluster of Excellence “Machine Learning - New Perspectives for Science”, EXC 2064/1, project number 390727645; the German Federal Ministry of Education and Research (BMBF) through the Tübingen AI Center (FKZ: 01IS18039A); and funds from the Ministry of Science, Research and Arts of the State of Baden-Württemberg. F. D. is grateful to the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for support.



# Cockpit: A Practical Debugging Tool for the Training of Deep Neural Networks

# 6.

## Abstract

When engineers train deep learning models, they are very much “flying blind”. Commonly used methods for real-time training diagnostics, such as monitoring the train/test loss, are limited. Assessing a net’s training process solely through these performance indicators is akin to debugging software without access to internal states through a debugger. To address this, we present Cockpit, a collection of instruments that enable a closer look into the inner workings of a learning machine, and a more informative and meaningful status report for practitioners. It facilitates the identification of learning phases and failure modes, like ill-chosen hyperparameters. The instruments leverage novel higher-order information about the gradient distribution and curvature, which has only recently become efficiently accessible. We believe that such a debugging tool, which we open-source for PyTorch, is valuable in troubleshooting the training process. By revealing new insights, it also more generally contributes to explainability and interpretability of deep nets.

6.1 Introduction & Motivation	75
6.2 Cockpit’s Instruments . . . . .	77
6.3 Experiments . . . . .	81
6.4 Showcase . . . . .	84
6.5 Benchmark . . . . .	84
6.6 Conclusion . . . . .	86

Code and experiments available at the Github repositories [f-dangel/cockpit](https://github.com/f-dangel/cockpit), [f-dangel/cockpit-experiments](https://github.com/f-dangel/cockpit-experiments)

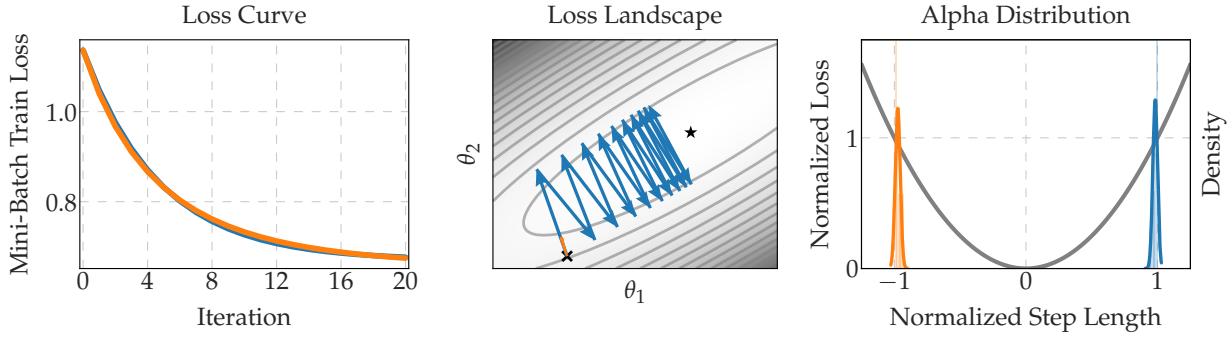
## 6.1 Introduction & Motivation

Deep learning represents a new programming paradigm: instead of deterministic programs, users design models and “simply” train them with data. In this metaphor, deep learning is a meta-programming form, where *coding* is replaced by *training*. Here, we ponder the question how we can provide more insight into this process by building a *debugger* specifically designed for deep learning.

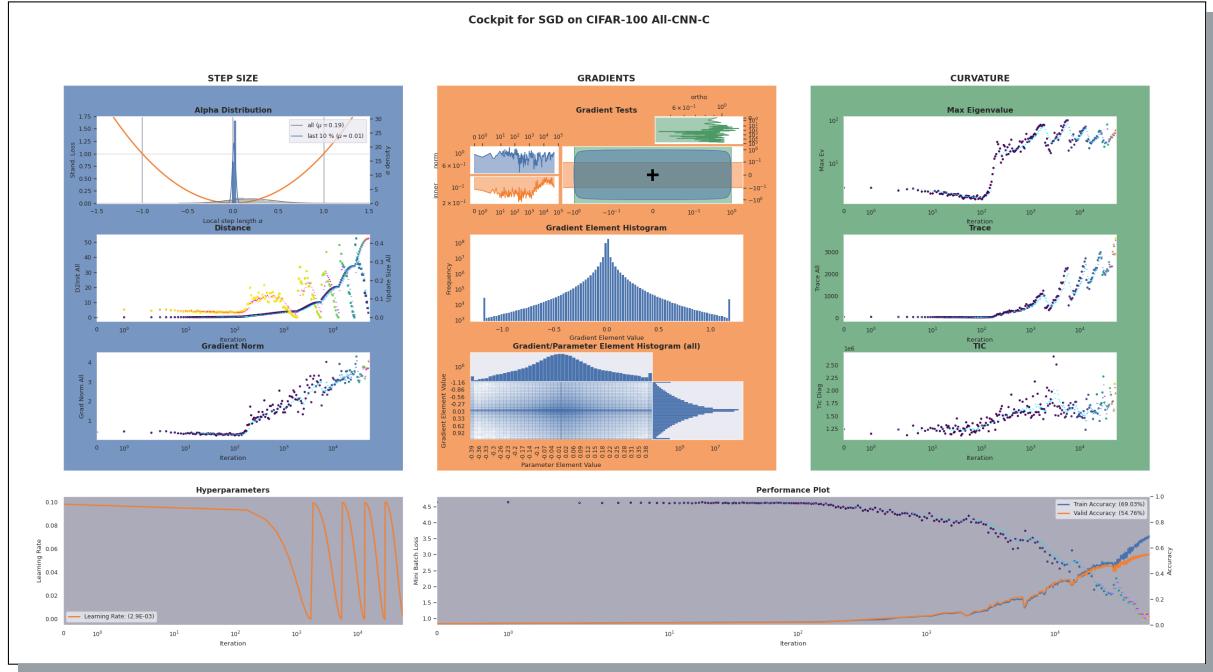
Debuggers are crucial for traditional software development. When things fail, they provide access to the internal workings of the code, allowing a look “into the box”. This is much more efficient than re-running the program with different inputs. And yet, deep learning is arguably closer to the latter. If the attempt to train a deep net fails, a machine learning engineer faces various options: should they change the training hyperparameters (how?); the optimizer (to which one?); the model (how?); or just re-run with a different seed? Machine learning toolboxes provide scant help to guide these decisions.

Of course, traditional debuggers can be applied to deep learning. They will give access to every single weight of a neural net, or the individual pixels of its training data. But this rarely yields insights towards successful training. Extracting meaningful information requires a statistical approach and distillation of the bewildering complexity into a manageable summary. Tools like TensorBoard [1] or Weights & Biases [17] were built in part to streamline this visualization. Yet, the quantities that are widely monitored (mainly train/test loss & accuracy), provide only scant explanation for relative differences between multiple training runs, because *they do not show the network’s internal state*. Figure 6.1 illustrates





**Figure 6.1:** Illustrative example: learning curves do not tell the whole story. Two different optimization runs (—/—) can lead to virtually the same loss curve (left). However, the actual optimization trajectories (middle), exhibit vastly different behaviors. In practice, the trajectories are intractably large and cannot be visualized directly. Recommendable actions for both scenarios (increase/decrease the learning rate) cannot be inferred from the loss curve. The  $\alpha$ -distribution, one Cockpit instrument (right), not only clearly distinguishes the two scenarios, but also allows for taking decisions how the learning rate should be adapted. See Section 6.3.3 for further details.



**Figure 6.2:** Screenshot of Cockpit’s full view while training the All-CNN-C [157] on CIFAR-100 with SGD using a cyclical learning rate schedule. This figure and its labels are not meant to be legible, but rather give an impression of Cockpit’s user experience. Gray panels (bottom row) show the information currently tracked by most practitioners. The individual instruments are discussed in Section 6.2, and observations are described in Section 6.4. An animated version can be found in the accompanying Github repository.

how such established learning curves can describe the *current* state of the model – whether it is performing well or not – while failing to inform about training state and dynamics. They tell the user *that* things are going well or badly, but not *why*. The situation is similar to flying a plane by sight, without instruments to provide feedback. It is not surprising, then, that achieving state-of-the-art performance in deep learning requires expert intuition, or plain trial & error.

We aim to enrich the deep learning pipeline with a visual and statistical debugging tool that uses newly proposed observables as well as several established ones (Section 6.2). We leverage and augment recent extensions to AD (*i.e.* BackPACK [38] for PyTorch [126]) to efficiently access second-order statistical (*e.g.* gradient variances) and geometric (*e.g.* Hessian) information. We show how these quantities can aid the deep learning

engineer in tasks, like learning rate selection, as well as detecting common bugs with data processing or model architectures (Section 6.3).

Concretely, we introduce Cockpit, a flexible and efficient framework for online-monitoring these observables during training in carefully designed plots we call “instruments” (Figure 6.2). To practical, such visualization must have a manageable computational overhead. We show that Cockpit scales well to real-world deep learning problems (see Figure 6.2 and Section 6.4). We also provide three different configurations of varying computational complexity and demonstrate that their instruments keep the computational cost *well below* a factor of 2 in run time (Section 6.5). It is available as open-source code, extendable, and seamlessly integrates into existing PyTorch training loops (Appendix C.1).

## 6.2 Cockpit’s Instruments

### Setting

We consider supervised regression/classification with labeled data  $(\mathbf{x}, \mathbf{y}) \in \mathbb{X} \times \mathbb{Y}$  generated by a distribution  $p_{\text{data}}(\mathbf{x}, \mathbf{y})$ . The training set  $\mathbb{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$  consists of  $N$  *i.i.d.* samples from  $p_{\text{data}}$  and the deep model  $f_{\theta} : \mathbb{X} \rightarrow \mathbb{F}$  maps inputs  $\mathbf{x}_n$  to predictions  $f_{\theta}(\mathbf{x}_n)$  by parameters  $\theta \in \Theta := \mathbb{R}^D$ . This prediction is evaluated by a loss function  $\ell : \mathbb{F} \times \mathbb{Y} \rightarrow \mathbb{R}$  which compares to the label  $\mathbf{y}_n$ . The goal is minimizing an inaccessible expected risk  $\mathcal{L}_{p_{\text{data}}}(\theta) = \int \ell(f_{\theta}(\mathbf{x}), \mathbf{y}) d p_{\text{data}}(\mathbf{x}, \mathbf{y})$  by empirical approximation through  $\mathcal{L}_{\mathbb{D}}(\theta) = 1/N \sum_{n=1}^N \ell(f_{\theta}(\mathbf{x}_n), \mathbf{y}_n) := 1/N \sum_{n=1}^N \ell_n(\theta)$ , which in practice is stochastically sub-sampled on mini-batches  $\mathbb{B} \subset \mathbb{D}$ ,

$$\mathcal{L}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{B}} \ell_n(\theta). \quad (6.1)$$

As is standard practice, we use first- and second-order information of the mini-batch loss, described by its gradient  $\mathbf{g}_{\mathbb{B}}(\theta)$  and Hessian  $\mathbf{H}_{\mathbb{B}}(\theta)$ ,

$$\mathbf{g}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{B}} \nabla_{\theta} \ell_n(\theta), \quad \mathbf{H}_{\mathbb{B}}(\theta) = \frac{1}{|\mathbb{B}|} \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{B}} \nabla_{\theta}^2 \ell_n(\theta). \quad (6.2)$$

### Design Choices

To minimize computational and design overhead, we restrict the metrics to quantities that require no additional model evaluations. This means that, at training step  $t \rightarrow t+1$  with mini-batches  $\mathbb{B}_t, \mathbb{B}_{t+1}$  and parameters  $\theta_t, \theta_{t+1}$ , we access information about the mini-batch losses  $\mathcal{L}_{\mathbb{B}_t}(\theta_t)$  and  $\mathcal{L}_{\mathbb{B}_{t+1}}(\theta_{t+1})$ , but no cross-terms that require additional forward passes.

### Key Point

$\mathcal{L}_{\mathbb{B}}(\theta), \mathbf{g}_{\mathbb{B}}(\theta)$ , and  $\mathbf{H}_{\mathbb{B}}(\theta)$  are just expected values of a *distribution* over the batch. Only recently, this distribution has begun to attract attention [48] as its computation has become more accessible [23, 38]. Contemporary optimizers leverage only the *mean* gradient and neglect higher moments.

**Table 6.1:** Overview of Cockpit quantities. They range from cheap byproducts, to nonlinear transformations of first-order information and Hessian-based measures. Some quantities have already been proposed, others are first to be considered in this work. They are categorized into configurations *economy*  $\subseteq$  *business*  $\subseteq$  *full* based on their run time overhead (see Section 6.5 for a detailed evaluation).

Name	Short Description	Config	Pos. in Fig. 6.2
Alpha	Normalized step size on a noisy quadratic interpolation between $\theta_t, \theta_{t+1}$	economy	top left
Distance	Distance from initialization $\ \theta_t - \theta_0\ _2$	economy	middle left
UpdateSize	Update size $\ \theta_{t+1} - \theta_t\ _2$	economy	middle left
GradNorm	Mini-batch gradient norm $\ g_{\mathbb{B}}(\theta)\ _2$	economy	bottom left
NormTest	Normalized fluctuations of residual norms $\ g_{\mathbb{B}} - g_n\ _2$ , proposed in [26]	economy	top center
InnerTest	Normalized fluctuations of $g_n$ 's parallel to $g_{\mathbb{B}}$ , proposed in [19]	economy	top center
OrthoTest	Like InnerTest but using the orthogonal components, proposed in [19]	economy	top center
GradHist1d	Histogram of individual gradient elements, $\{(g_n)_j\}_{(x_n, y_n) \in \mathbb{B}}^{j=1, \dots, D}$	economy	middle center
TICDiag	Relation of (diagonal) curvature and gradient noise, inspired by [162]	business	bottom right
HessTrace	Exact or approximate Hessian trace, $\text{Tr}(H_{\mathbb{B}}(\theta))$ , inspired by [177]	business	middle right
HessMaxEV	Maximum Hessian eigenvalue, $\lambda_{\max}(H_{\mathbb{B}}(\theta))$ , inspired by [177]	full	top right
GradHist2d	Histogram of weights & per-sample gradients, $\{(\theta_j, [g_n]_j)\}_{(x_n, y_n) \in \mathbb{B}}^{j=1, \dots, D}$	full	bottom center

One core point of our work is making extensive use of these distribution properties, trying to visualize them in various ways. This distinguishes Cockpit from being “just a collection of plots” that could be built in tools like TensorBoard. Leveraging these distributional quantities, we create instruments and show how they can help adapt hyperparameters (Section 6.2.1), analyze the loss landscape (Section 6.2.2), and track network dynamics (Section 6.2.3). Instruments can sometimes be built from already-computed information or are efficient variants of previously proposed observables. To keep the presentation concise, we highlight the instruments shown in Figure 6.2 and listed in Table 6.1. Appendix C.3 defines them formally and contains more extensions, such as the mean GSNR [100], the early stopping [104] and CABS [11] criterion, which can all be used in Cockpit.

### 6.2.1 Adapting Hyperparameters

One big challenge in deep learning is setting the hyperparameters correctly, which is currently mostly done by trial & error through parameter searches. We aim to augment this process with instruments that inform the user about the effect that the chosen parameters have on the current training process.

#### Alpha : Are We Crossing the Valley?

Using individual loss and gradient observations at the start and end point of each iteration, we build a noise-informed uni-variate quadratic approximation along the step direction (*i.e.* the loss as a function of the step size), and assess to which point on this parabola our optimizer moves. We standardize this value  $\alpha$  such that stepping to the valley-floor is assigned  $\alpha = 0$ , the starting point is at  $\alpha = -1$  and updates to the point exactly opposite of the starting point have  $\alpha = 1$  (see Appendix C.3.2 for a more detailed visual and mathematical description of  $\alpha$ ). Figure 6.1 illustrates the scenarios  $\alpha = \pm 1$  and how monitoring the

$\alpha$ -distribution (right panel) can help distinguish between two training runs with similar performance but distinct failure sources. By default, this Cockpit instrument shows the  $\alpha$ -distribution for the last 10 % of training and the entire training process (top left plot in Figure 6.2). In Section 6.3.3 we show empirically that, *counter-intuitively*, it is generally *not* a good idea to choose the step size such that  $\alpha$  is close to zero.

### Distances: Are We Making Progress?

Another way to discern the trajectories in Figure 6.1 is by measuring the  $L_2$  distance from initialization [113] and the update size [4, 51] in parameter space. Both are shown together in one Cockpit instrument (see also center-left plot in Figure 6.2) and are far larger for the blue line in Figure 6.1. These distance metrics are also able to disentangle phases for the blue path. Using the same step size, it will continue to “jump back and forth” between the loss valley’s walls but at some point cease to make progress. During this “surfing of the walls”, the *distance from initialization* increases, ultimately though, it will stagnate, with the *update size* remaining non-zero, indicating diffusion. While the initial “surfing the wall”-phase benefits training (see Section 6.3.3), achieving stationarity may require adaptation once the optimizer reaches that diffusion.

### Gradient Norm: How Steep Is the Wall?

The *update size* will show that the orange trajectory is stuck. But why? Such slow-down can result from both a bad learning rate and from loss landscape plateaus. The *gradient norm* (bottom left panel in Figure 6.2) distinguishes these two causes.

### Gradient Tests: How Noisy Is the Batch?

The batch size trades off gradient accuracy versus computational cost. Recently, adaptive sampling strategies based on testing geometric constraints between mean and individual gradients have been proposed [19, 26]. The *norm*, *inner product*, and *orthogonality tests* use a standardized radius and two band widths (parallel and orthogonal to the gradient mean) that indicate how strongly individual gradients scatter around the mean. The original works use these values to adapt batch sizes. Instead, Cockpit combines all three tests into a single gauge (top middle plot of Figure 6.2) using the standardized noise radius and band widths for visualization. These noise signals can be used to guide batch size adaptation on- and offline, or to probe the influence of gradient alignment on training speed [142] and generalization [27, 28, 100].

## 6.2.2 Hessian Properties for Local Loss Geometry

An intuition for the local loss landscape helps in many ways. It can help diagnose whether training is stuck, to adapt the step size, and explain stability or regularization [56, 81]. The key challenge is the large number of weights: low-dimensional projections of surfaces can behave

unintuitively [112], but tracking the extreme or average behaviors may help in debugging, especially if first-order metrics fail.

### Hessian Eigenvalues: A Gorge or a Lake?

In convex optimization, the maximum Hessian eigenvalue crucially determines the appropriate step size [144]. Many works have studied the Hessian spectrum in machine learning [e.g. 55, 56, 112, 139, 140, 177]. In short: curvature matters. Established [127] and recent autodiff frameworks [38] can compute Hessian properties without requiring the full matrix. Cockpit leverages this to provide the *Hessian's largest eigenvalue* and *trace* (right top and middle plots in Figure 6.2). The former resembles the loss surface's sharpest valley and can thus hint at training instabilities [81]. The *trace* describes a notion of "average curvature", since the eigenvalues  $\lambda_i$  relate to it by  $\sum_i \lambda_i = \text{Tr}(H_{\mathbb{B}}(\theta))$ , which might correlate with generalization [80].

### TIC: How Do Curvature & Gradient Noise Interact?

There is an ongoing debate about curvature's link to generalization [e.g. 42, 70, 86]. The *Takeuchi Information Criterion* (TIC) [160, 162] estimates the generalization gap by a ratio between Hessian and non-central second gradient moment. It also provides intuition for changes in the objective function implied by gradient noise. Inspired by [162], Cockpit provides mini-batch TIC estimates (bottom right plot of Figure 6.2).

### 6.2.3 Visualizing Internal Network Dynamics

Histograms are a natural visual compression of the high-dimensional  $|\mathbb{B}| \times D$  individual gradient values. They give insights into the gradient *distribution* and hence offer a more detailed view of the learning signal. Together with the parameter associated to each individual gradient, the entire model status and dynamics can be visualized in a single plot and be monitored during training. This provides a more fine-grained view of training compared to tracking parameters and gradient norms [51].

### Gradient & Parameter Histograms: What Is Happening in Our Net?

Cockpit offers a uni-variate *histogram of the gradient elements*, i.e. the numbers  $\{[g_n(\theta)]_j\}_{(x_n, y_n) \in \mathbb{B}}^{j=1, \dots, D}$ . Additionally, a combined *histogram of parameter-gradient pairs*  $\{([\theta]_j, [g_n(\theta)]_j)\}_{(x_n, y_n) \in \mathbb{B}}^{j=1, \dots, D}$  provides a two-dimensional look into the network's gradient and parameter values in a mini-batch. Section 6.3.1 shows an example use-case of the gradient histogram; Section 6.3.2 makes the case for the layer-wise variants of the instruments.

## 6.3 Experiments

The diverse information provided by Cockpit can help users and researchers in many ways, some of which, just like for a traditional debugger, only become apparent in practical use. In this section, we present a few motivating examples, selecting specific instruments and scenarios in which they are practically useful. Specifically, we show that Cockpit can help the user discern between, and thus fix, common training bugs (Sections 6.3.1 and 6.3.2) that are otherwise hard to distinguish as they lead to the same failure: bad training. We demonstrate that Cockpit can guide practitioners to choose efficient hyperparameters *within a single training run* (Sections 6.3.2 and 6.3.3). Finally, we highlight that Cockpit’s instruments can provide research insights about the optimization process (Section 6.3.3). Our empirical findings are demonstrated on problems from the DeepOBS [146] benchmark collection.

### 6.3.1 Incorrectly Scaled Data

One prominent source of bugs is the data pipeline. To pick a relatively simple example: for standard optimizers to work at their usual learning rates, network inputs must be standardized (*i.e.* between zero and one, or have zero mean and unit variance [*e.g.* 15]). If the user forgets to do this, optimizer performance is likely to degrade. It can be difficult to identify the source of this problem as it does not cause obvious failures, `NaN` or `Inf` gradients, *etc.*. We now construct a semi-realistic example, to show how using Cockpit can help diagnose this problem upon observing slow training performance.

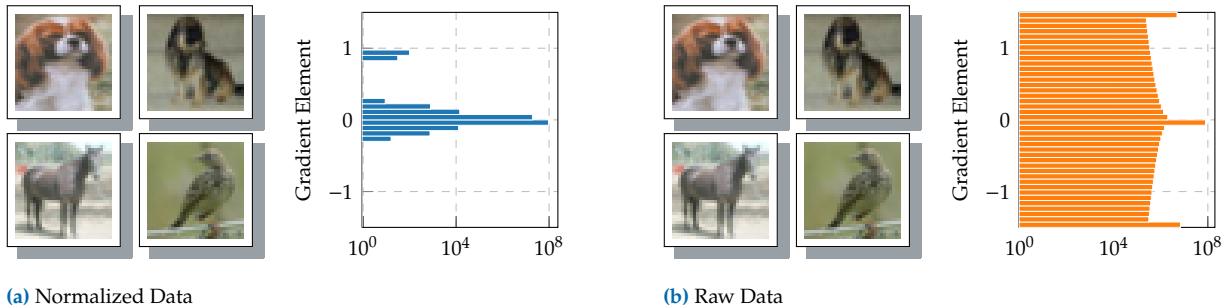
By default<sup>1</sup>, the popular image datasets CIFAR-10/100 [90] are provided as NumPy [66] arrays that consist of integers in the interval  $[0, 255]$ . This *raw* data, instead of the widely used version with floats in  $[0, 1]$ , changes the data scale and thus the gradients by a factor of 255. Therefore, the optimizer’s optimal learning rate is scaled as well. In other words, the default parameters of popular optimization methods may not work well anymore, or good hyperparameters may take extreme values. Even if the user directly inspects the training images, this may not be apparent (Figure 6.3). But the gradient histogram instrument of Cockpit, which has a deliberate default plotting range around  $[-1, 1]$  to highlight such problems, immediately and prominently shows that there is an issue.

1: See the documentation, available at [cs.toronto.edu/~kriz/cifar.html](http://cs.toronto.edu/~kriz/cifar.html)

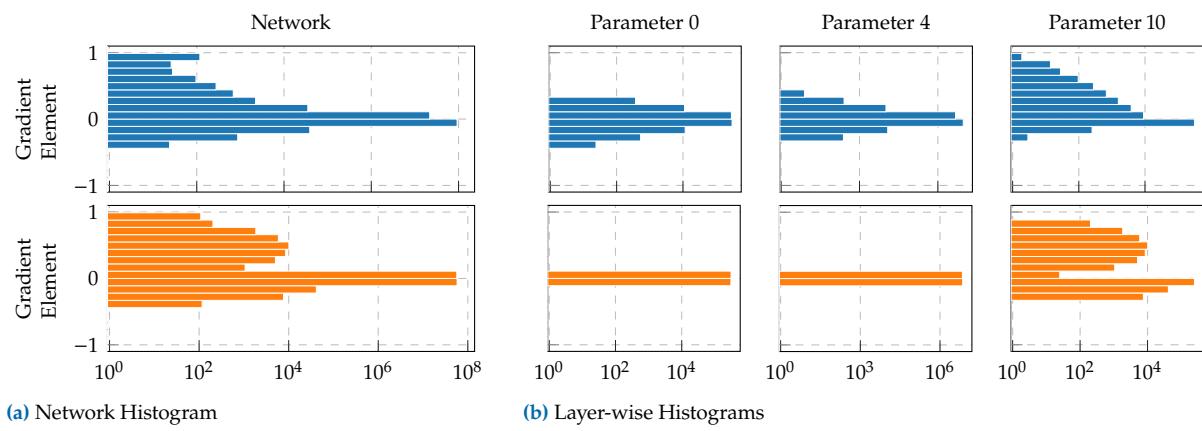
Of course, this particular data is only a placeholder for real practical data sets. While this problem may not frequently arise in the highly pre-processed, packaged CIFAR-10, it is not a rare problem for practitioners who work with their personal datasets. This is particularly likely in domains outside standard computer vision, *e.g.* when working with mixed-type data without obvious natural scales.

### 6.3.2 Vanishing Gradients

The model itself can be a source of training bugs. As before, such problems mostly arise with novel datasets, where well-working architectures are



**Figure 6.3: Same inputs, different gradients; Catching data bugs with Cockpit.** (a) *normalized* ( $[0, 1]$ ) and (b) *raw* ( $[0, 255]$ ) images look identical in auto-scaled front-ends like matplotlib’s `imshow`. The gradients on 3c3d, however, are crucially affected by this scaling.

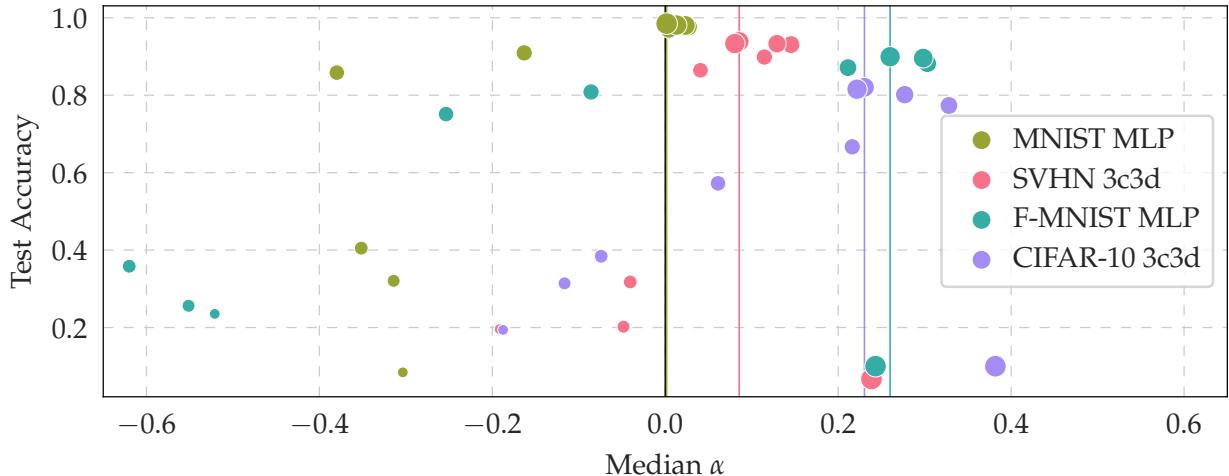


**Figure 6.4: Gradient distributions of two similar architectures on the same problem.** (a) Distribution of individual gradient elements summarized over the entire network. Both seem similar. (b) Layer-wise histograms for a subset of layers. Parameter 0 is the layer closest to the network’s input, parameter 10 closest to its output. Only the layer-wise view reveals that there are several degenerated gradient distributions for the orange network making training unnecessarily hard.

unknown. The following example shows how even small (in terms of code) model modifications may severely harm the training.

Figure 6.4a shows the distribution of gradient values of two different network architectures in blue and orange. Although the blue model trains considerably better than the orange one, their gradient distributions look quite similar. The difference becomes evident when inspecting the histogram *layer-wise*. We can see that multiple layers have a degenerated gradient distribution with many elements being practically zero (see Figure 6.4b, bottom row). Since the fully connected layers close to the output have far more parameters (a typical pattern of convolutional networks), they dominate the network-wide histogram. This obscures that a major part of the model is effectively unable to train.

Both the blue and orange networks follow DeepOBS’s 3c3d architecture. The only difference is the non-linearity: the blue net uses standard ReLU activations, while the orange one has sigmoid activations. Here, the layer-wise histogram instrument of Cockpit highlights which part of the architecture makes training unnecessarily hard. Accessing information layer-wise is also essential due to the strong overparameterization in deep models where training can happen in small subspaces [65]. Again, this is hard to do with common monitoring tools, such as the loss curve.



**Figure 6.5:** Test accuracy as a function of standardized step size  $\alpha$ . For four DeepOBS problems (see Appendix C.5), final test accuracy is shown versus the median  $\alpha$ -value over the entire training. Marker size indicates the magnitude of the raw learning rate, marker color identifies tasks (see legend). For each problem, the best-performing setting is highlighted by a vertical colored line.

### 6.3.3 Tuning Learning Rates

Once the architecture is defined, the optimizer’s learning rate is the most important hyperparameter to tune. Getting it right requires extensive hyperparameter searches at high resource costs. Cockpit’s instruments can provide intuition and information to streamline this process: in contrast to the raw learning rate, the curvature-standardized step size  $\alpha$ -quantity (see Section 6.2.1) has a natural scale.

Across multiple optimization problems, we observe, perhaps surprisingly, that the best runs and indeed all good runs have a median  $\alpha > 0$  (Figure 6.5). This illustrates a fundamental difference between stochastic optimization, as is typical for machine learning, and classic deterministic optimization. Instead of locally stepping “to the valley floor” (optimal in the deterministic case), stochastic optimizers should *overshoot* the valley somewhat. This need to “surf the walls” has been hypothesized before [e.g. 173, 176] as a property of neural network training. Frequently, learning rates are adapted during training, which fits with our observation about positive  $\alpha$ -values: “overshooting” allows fast early progression towards areas of lower loss, but it does not yield convergence in the end. Real-time visualizations of the training state, as offered by Cockpit, can augment these fine-tuning processes.

Figure 6.5 also indicates a major challenge preventing simple automated tuning solutions: the optimal  $\alpha$ -value is problem-dependent, and simpler problems, such as a multi-layer perceptron (MLP) on MNIST [95], behave much more similar to classic optimization problems. Algorithmic research on small problems can thus produce misleading conclusions. The figure also shows that the  $\alpha$ -gauge is not sufficient by itself: extreme overshooting with a too-large learning rate leads to poor performance, which however can be prevented by taking additional instruments into account. This makes the case for the cockpit metaphor of increasing interpretability from several instruments in conjunction. By combining the  $\alpha$ -instrument with other gauges that capture the local geometry or network dynamics, the user can better identify good choices of the learning rate and other hyperparameters.

## 6.4 Showcase

Having introduced the tool, we can now return to [Figure 6.2](#) for a closer look. The figure shows a snapshot from training the All-CNN-C [157] on CIFAR-100 using SGD with a cyclic learning rate schedule (bottom left panel). Diagonal curvature instruments are configured to use an MC approximation to save run time (here,  $C = 100$ , compare [Section 6.5](#)).

A glance at all panels shows that the learning rate schedule is reflected in the metrics. However, the instruments also provide insights into the early phase of training (first  $\sim 100$  iterations), where the learning rate is still unaffected by the schedule: there, the loss plateaus and the optimizer takes relatively small steps (compared to later, as can be seen in the small gradient norms, and small distance from initialization). Based on these low-cost instruments, one may thus at first suspect that training was poorly initialized; but training indeed succeeds after iteration 100! Viewing Cockpit entirely though, it becomes clear that optimization in these first steps is not stuck at all: while loss, gradient norms, and distance in parameter space remain almost constant, curvature changes, which expresses itself in a clear downward trend of the maximum Hessian eigenvalue (top right panel).

The importance of early training phases has recently been hypothesized [51], suggesting a logarithmic timeline. Not only does our showcase support this hypothesis, but it also provides an explanation from the curvature-based metrics, which in this particular case are the only meaningful feedback in the first few training steps. It also suggests monitoring training at log-spaced intervals. Cockpit provides the flexibility to do so, indeed, [Figure 6.2](#) has been created with log-scheduled tracking events.

As a final note, we recognize that the approach taken here promotes an amount of *manual work* (monitoring metrics, deliberately intervening, *etc.*) that may seem ironic and at odds with the paradigm of automation that is at the heart of machine learning. However, we argue that this might be what is needed at this point in the evolution of the field. Deep learning has been driven notably by scaling compute resources [163], and fully automated one-shot training may still be some way out. To develop better training methods, researchers, not just users, need *algorithmic* interpretability and explainability: direct insights and intuition about the processes taking place “inside” neural nets. To highlight how Cockpit might provide this, we contrast in [Appendix C.6](#) the view of two convex DeepOBS problems (noisy quadratic & logistic regression on MNIST). In both cases, the instruments behave differently compared to the deep learning problem in [Figure 6.2](#). In particular, the gradient norm increases (left column, bottom panel) during training, and individual gradients become less scattered (center column, top panel). This is diametrically opposed to the convex problems and shows that deep learning differs even qualitatively from well-understood optimization problems.

## 6.5 Benchmark

[Section 6.3](#) made a case for Cockpit as an effective debugging and tuning tool. To make the library useful in practice, it must also have limited

computational cost. We now show that it is possible to compute all quantities at reasonable overhead. The user can control the absolute cost along two dimensions, by reducing the number of instruments, or by reducing their update frequency.

All benchmark results show SGD without momentum. Cockpit's quantities, however, work for generic optimizers and can mostly be used identically without increased costs. One current exception is `Alpha` which can be computed more efficiently given the update rule.<sup>2</sup>

2: This is currently implemented for vanilla SGD. Otherwise, Cockpit falls back to a less efficient scheme.

### Complexity Analysis

Computing more information adds computational overhead, of course. However, recent work [38] has shown that first-order information, like distributional statistics on the batch gradients, can be computed on top of the mean gradient at little extra cost. Similar savings apply for most quantities in Table 6.1, as they are (non-)linear transformations of individual gradients. A subset of Cockpit's quantities also uses second-order information from the Hessian diagonal. For ReLU networks on a classification task with  $C$  classes, the additional work is proportional to  $C$  gradient backpropagations (*i.e.*  $C = 10$  for CIFAR-10,  $C = 100$  for CIFAR-100). Parallel processing can, to some extent, process these extra backpropagations in parallel without significant overhead. If this is no longer possible, we can fall back to a Monte Carlo (MC) sampling approximation, which reduces the number of extra backprop passes to the number of samples (1 by default).<sup>3</sup>

While parallelization is possible for the gradient instruments, computing the maximum Hessian eigenvalue is inherently sequential. Similar to Yao et al. [177], we use matrix-free Hessian-vector products by automatic differentiation [127], where each product's costs are proportional to one gradient computation. Regardless of the underlying iterative eigensolver, multiple such products must be queried to compute the spectral norm (the number depends on the spectral gap to the second-largest eigenvalue).

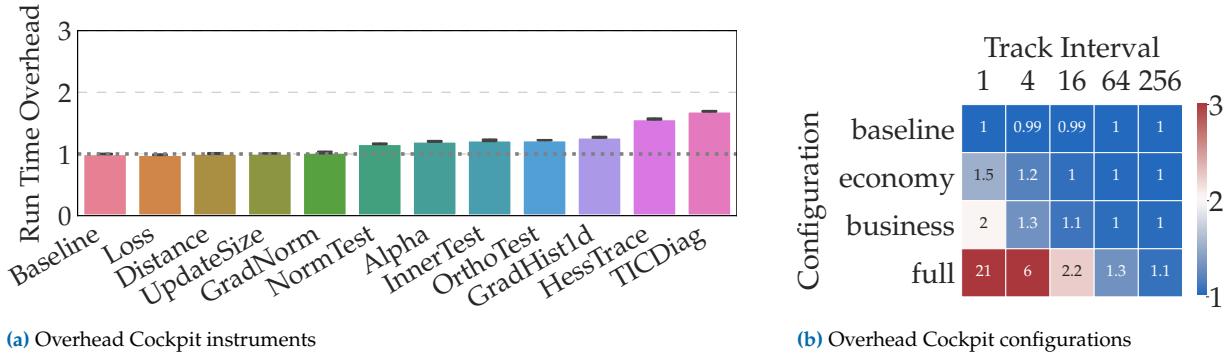
3: An MC-sampled approximation of the Hessian/generalized Gauss-Newton has been used in Figure 6.2 to reduce the prohibitively large number of extra backprops on CIFAR-100 ( $C = 100$ ).

### Run Time Benchmark

Figure 6.6a shows the wall-clock computational overhead for individual instruments (details in Appendix C.5).<sup>4</sup> As expected, byproducts are virtually free, and quantities that rely solely on first-order information add little overhead (at most roughly 25 % on this problem). Thanks to parallelization, the ten extra backward passes required for Hessian quantities reduce to less than 100 % overhead. Individual overheads also do not simply add up when multiple quantities are tracked, because quantities relying on the same information share computations.

To allow a rough cost control, Cockpit currently offers three configurations, called `"economy"`, `"business"`, and `"full"`, in increasing order of cost (Table 6.1). As a basic guideline, we consider a factor of two to be an acceptable limit for the increase in training time and benchmark the configurations' run times for different tracking intervals. Figure 6.6b shows a run time matrix for the CIFAR-10 3c3d problem, where settings that meet this limit are set in blue (more problems including ImageNet

4: To improve readability, we exclude `HessMaxEV` here, because its overhead is large compared to other quantities. Surprisingly, we also observed significant cost for the 2D histogram on GPU. It is caused by an implementation bottleneck for histogram shapes observed in deep models. We thus also omit `GradHist2d` here, as we expect it to be eliminated with future implementations (see Appendix C.5.2 for a detailed analysis and further benchmarks). Both quantities, however, are part of the benchmark shown in Figure 6.6b.



**Figure 6.6:** Run time overhead for individual Cockpit instruments and configurations as shown on CIFAR-10 3c3d on a GPU. **(a)** The run time overheads for individual instruments are shown as multiples of the *baseline* (no tracking). Most instruments add little overhead. This plot shows the overhead in one iteration, determined by averaging over multiple iterations and random seeds. **(b)** Overhead for different Cockpit configurations. Adjusting the tracking interval and re-using the computation shared by multiple instruments can make the overhead orders of magnitude smaller. Blue fields mark settings that allow tracking without doubling the training time.

are shown in Appendix C.5). Speedups due to shared computations are easy to read off: summing all the individual overheads shown in Figure 6.6a would result in a total overhead larger than 200 %, while the joint overhead (*business*) reduces to 140 %. The *economy* configuration can easily be tracked at every step of this problem and stay well below our threshold of doubling the execution time. Cockpit’s full view, shown in Figure 6.2, can be updated every 64-th iteration without a major increase in training time (this corresponds to about five updates per epoch). Finally, tracking any configuration about once per epoch—which is common in practice—adds overhead close to zero (rightmost column).

This good performance is largely due to the efficiency of the BackPACK package [38], which we leverage with custom and optimized modification, that compacts information layer-wise and then discards unneeded buffers. Using layer-wise information (Section 6.3.2) scales better to large networks, where storing the entire model’s individual gradients all at once becomes increasingly expensive (see Appendix C.5). To the best of our knowledge, many of the quantities in Table 6.1, especially those relying on individual gradients, have only been explored on rather small problems. With Cockpit they can now be accessed at a reasonable rate for deep learning models outside the toy problem category.

## 6.6 Conclusion

Contemporary machine learning, in particular deep learning, remains a craft and an art. High dimensionality, stochasticity, and non-convexity require constant tracking and tuning, often resulting in a painful process of trial and error. When things fail, popular performance measures, like the training loss, do not provide enough information by themselves. These metrics only tell *whether* the model is learning, but not *why*. Alternatively, traditional debugging tools can provide access to individual weights and data. However, in models whose power only arises from possessing myriad weights, this approach is hopeless, like looking for the proverbial needle in a haystack.

To mitigate this, we proposed Cockpit, a practical visual debugging tool

for deep learning. It offers instruments to monitor the network's internal dynamics during training, in real-time. In its presentation, we focused on two crucial factors affecting user experience: Firstly, such a debugger must provide meaningful insights. To demonstrate Cockpit's utility, we showed how it can identify bugs where traditional tools fail. Secondly, it must come at a feasible computational cost. Although Cockpit uses rich second-order information, efficient computation keeps the necessary run time overhead cheap. The open-source PyTorch package can be added to many existing training loops.

Obviously, such a tool is never complete. Just like there is no perfect universal debugger, the list of current instruments is naturally incomplete. Further practical experience with the tool, for example in the form of a future larger user study, could provide additional evidence for its utility. However, our analysis shows that Cockpit provides useful tools and extracts valuable information presently not accessible to the user. We believe that this improves algorithmic interpretability – helping practitioners understand how to make their models work – but may also inspire new research. The code is designed flexibly, deliberately separating the computation and visualization. New instruments can be added easily and also be shown by the user's preferred visualization tool, *e.g.* TensorBoard. Of course, instead of just showing the data, the same information can be used by novel algorithms directly, side-stepping the human in the loop.

## Acknowledgments

The authors gratefully acknowledge financial support by the European Research Council through ERC StG Action 757275 / PANAMA; the DFG Cluster of Excellence "Machine Learning - New Perspectives for Science", EXC 2064/1, project number 390727645; the German Federal Ministry of Education and Research (BMBF) through the Tübingen AI Center (FKZ: 01IS18039A); and funds from the Cyber Valley Initiative of the Ministry for Science, Research and Arts of the State of Baden-Württemberg. Moreover, the authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Felix Dangel and Frank Schneider. Further, we are grateful to Agustinus Kristiadi, Alexandra Gessner, Christian Fröhlich, Filip de Roos, Jonathan Wenger, Julia Grosse, Lukas Tatzel, Marius Hobbahn, and Nicholas Krämer for providing feedback to the manuscript.



# ViViT: Curvature Access Through the Generalized Gauss-Newton's Low-rank Structure

7.

## Abstract

Curvature in form of the Hessian or its generalized Gauss-Newton (GGN) approximation is valuable for algorithms that rely on a local model for the loss to train, compress, or explain deep networks. Existing methods based on implicit multiplication via automatic differentiation or Kronecker-factored block diagonal approximations do not consider noise in the mini-batch. We present ViViT, a curvature model that leverages the GGN's low-rank structure without further approximations. It allows for efficient computation of eigenvalues, eigenvectors, as well as per-sample first- and second-order directional derivatives. The representation is computed in parallel with gradients in one backward pass and offers a fine-grained cost-accuracy trade-off, which allows it to scale. We demonstrate this by conducting performance benchmarks and substantiate ViViT's usefulness by studying the impact of noise on the GGN's structural properties during neural network training.

7.1 Introduction & Motivation	89
7.2 Notation & Method	91
7.3 Experiments	95
7.4 Related Work	101
7.5 Use Cases	101
7.6 Conclusion	102

Code and experiments available at the Github repositories [f-dangel/vivit](https://github.com/f-dangel/vivit), [f-dangel/vivit-experiments](https://github.com/f-dangel/vivit-experiments)

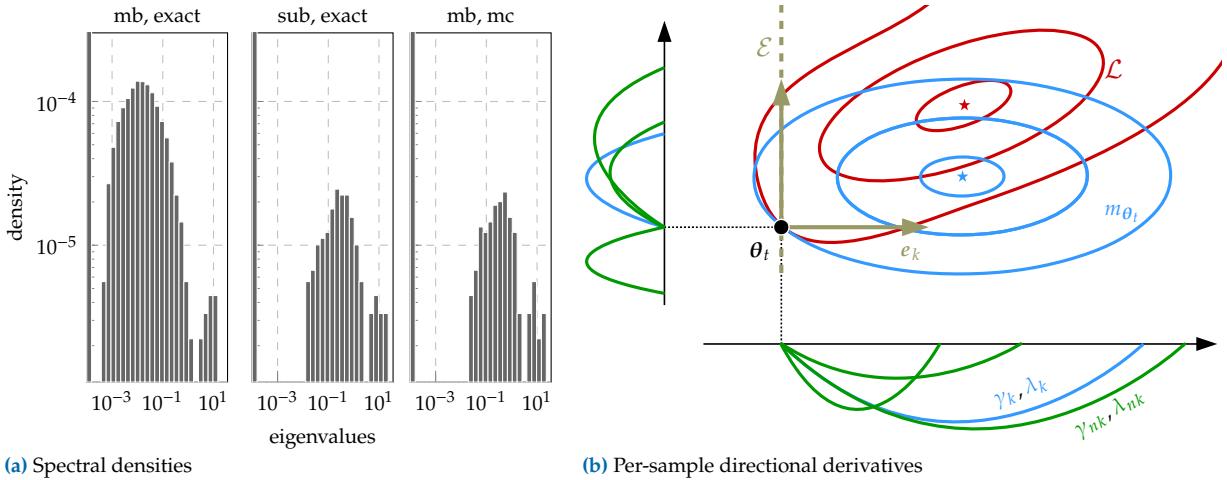
## 7.1 Introduction & Motivation

The large number of trainable parameters in deep neural networks imposes computational constraints on the information that can be made available to optimization algorithms. Standard machine learning libraries [1, 126] mainly provide access to first-order information in the form of *average* mini-batch gradients. This is a limitation that complicates the development of novel methods that may outperform the state-of-the-art: They must use the same objects to remain easy to implement and use, and to rely on the highly optimized code of those libraries. There is evidence that this has led to stagnation in the performance of first-order optimizers [145]. Here, we thus study how to provide efficient access to richer information, namely higher-order derivatives and their distribution across the mini-batch.



Recent advances in AD [23, 38] have made such information more readily accessible through leveraging algebraic structure in the differentiated loss. We use and extend this functionality to efficiently access curvature in form of the Hessian's generalized Gauss-Newton (GGN) approximation. It offers practical advantages over the Hessian and is established for training [106, 109], compressing [154], or adding uncertainty to [89, 133, 134] neural networks. It is also linked theoretically to the natural gradient method [5] via the Fisher information matrix [107, Section 9.2].

Traditional ways to access curvature fall into two categories. Firstly, repeated automatic differentiation allows for matrix-free exact multiplication with the Hessian [127] and GGN [148]. Iterative linear and eigensolvers can leverage such functionality to compute Newton steps [53, 106, 183] and spectral properties [3, 55, 61, 123, 139, 140, 177] on



**Figure 7.1:** Overview of ViViT’s quantities: **(a)** GGN eigenvalue distribution of DeepOBS’ 3c3d architecture on CIFAR-10 [146] for settings with different costs on a mini-batch of size  $N = 128$ . From left to right: exact GGN, exact GGN on a mini-batch fraction, MC approximation of the GGN. **(b)** Pictorial illustration: loss function  $\mathcal{L}$  from Equation (7.1), quadratic model  $m_{\theta_t}$  around  $\theta_t \in \mathbb{R}^2$  from Equation (7.6) (both represented by their contour lines). The low-rank structure provides efficient access to the GGN’s eigenvectors  $\{e_k\}$ , along which  $m_{\theta_t}$  decouples into one-dimensional parabolas characterized by the directional derivatives  $\gamma_k, \lambda_k$  and per-sample contributions  $\gamma_{n,k}, \lambda_{n,k}$  (Equation (7.8)).  $E$  is the GGN’s top-1 eigenspace.

arbitrary architectures thanks to the generality of AD. However, repeated matrix-vector products are potentially detrimental to performance.

Secondly, K-FAC (Kronecker-factored approximate curvature) [21, 63, 108, 109] constructs an explicit light-weight representation of the GGN based on its algebraic Kronecker structure. The computations are streamlined via gradient backpropagation and the resulting matrices are cheap to store and invert. This allows K-FAC to scale: It has been used successfully with large mini-batches [121]. One reason for this efficiency is that K-FAC only approximates the GGN’s block diagonal, neglecting interactions across layers. Such terms could be useful, however, for applications like uncertainty quantification with Laplace approximations [40, 89, 133, 134] that currently rely on K-FAC. Moreover, due to its specific design for optimization, the Kronecker representation does not become more accurate with more data. It remains a simplification, exact only under assumptions unlikely to be met in practice [109]. This might be a downside for applications that depend on a precise curvature proxy.

Here, we propose ViViT (inspired by  $VV^\top$  in Equation (7.3)), a curvature model that leverages the GGN’s low-rank structure. Like K-FAC, its representation is computed in parallel with gradients. But it allows a cost-accuracy trade-off, ranging from the *exact* GGN to an approximation that costs a single gradient computation. Our contributions are:

- ▶ We highlight the GGN’s low-rank structure, and the structural limit for the inherent curvature information contained in a mini-batch.
- ▶ We present how to compute various GGN properties efficiently by exploiting this structure (Figure 7.1): The exact eigenvalues, eigenvectors, and per-sample directional derivatives. In contrast to other methods, these quantities allow modeling curvature noise.
- ▶ We introduce approximations that allow a flexible trade-off between computational cost and accuracy. We also provide a fully-featured efficient implementation in PyTorch [126] on top of BackPACK [38].

- We empirically demonstrate scalability and efficiency of leveraging the GGN’s low-rank structure through benchmarks on different deep neural network architectures. Finally, we use ViViT’s quantities to study the GGN, and how it is affected by noise, during training.

The main focus is demonstrating that many interesting curvature properties, including uncertainty, can be computed efficiently. Practical applications of this curvature uncertainty are discussed in [Section 7.5](#).

## 7.2 Notation & Method

Consider a model  $f_\theta : \mathbb{X} \rightarrow \mathbb{Y}$  and a dataset  $\{(x_n, y_n) \in \mathbb{X} \times \mathbb{Y}\}_{n=1}^N$ . For simplicity we use  $N$  for both the mini-batch and training set size. The network, parameterized by  $\theta \in \Theta$ , maps a sample  $x_n$  to a prediction  $f_\theta(x_n) \in \mathbb{F}$ . Predictions are scored by a convex loss function  $\ell : \mathbb{F} \times \mathbb{Y} \rightarrow \mathbb{R}$  (*e.g.* cross-entropy or square loss), which compares to the ground truth  $y_n$ . The training objective  $\mathcal{L} : \Theta \rightarrow \mathbb{R}$  is the empirical risk

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(x_n), y_n). \quad (7.1)$$

We use  $\ell_n(\theta) = \ell(f_\theta(x_n), y_n)$  and  $f_n(\theta) = f_\theta(x_n)$  for per-sample losses and predictions. For gradients, we write  $g_n(\theta) = \nabla_\theta \ell_n(\theta)$  and  $g(\theta) = \nabla_\theta \mathcal{L}(\theta)$ , suppressing  $\theta$  if unambiguous. We also set  $\Theta = \mathbb{R}^D$  and  $\mathbb{F} = \mathbb{R}^C$  with  $D, C$  the model parameter and prediction space dimension, respectively. For classification,  $C$  is the number of classes.

### Hessian & GGN

Two-fold chain rule application to the split  $\ell \circ f$  decomposes the Hessian of [Equation \(7.1\)](#) into two parts  $\nabla_\theta^2 \mathcal{L}(\theta) = G(\theta) + R(\theta) \in \mathbb{R}^{D \times D}$ ; the positive semi-definite GGN

$$G(\theta) = \frac{1}{N} \sum_{n=1}^N (\mathbf{J}_\theta f_n)^\top \nabla_{f_n}^2 \ell_n (\mathbf{J}_\theta f_n) = \frac{1}{N} \sum_{n=1}^N G_n(\theta) \quad (7.2)$$

and a residual  $R = 1/N \sum_{n=1}^N \sum_{c=1}^C (\nabla_\theta^2 [f_n]_c)[\nabla_{f_n} \ell_n]_c$ . Here, we use the Jacobian  $\mathbf{J}_\theta \mathbf{b}$  that contains partial derivatives of  $\mathbf{b}$  w.r.t.  $a$ ,  $[\mathbf{J}_\theta \mathbf{b}]_{i,j} = \partial b_i / \partial a_j$  ([Definition 2.4](#)). As the residual may alter the Hessian’s definiteness—undesirable in many applications—we focus on the GGN. [Section 7.3.2](#) provides empirical evidence that the curvature’s top eigenspace is largely unaffected by this simplification.

### Low-rank Structure

By basic inequalities, [Equation \(7.2\)](#) has  $\text{rank}(G) \leq NC$ .<sup>1</sup> To make this explicit, we factorize the positive semi-definite Hessian  $\nabla_{f_n}^2 \ell_n = \sum_{c=1}^C \mathbf{s}_{n,c} \mathbf{s}_{n,c}^\top$ , where  $\mathbf{s}_{n,c} \in \mathbb{R}^C$  and denote its backpropagated version by  $\mathbf{v}_{n,c} = (\mathbf{J}_\theta f_n)^\top \mathbf{s}_{n,c} \in \mathbb{R}^D$ . Absorbing sums into matrix multiplications,

1: We assume the overparameterized deep learning setting ( $NC < D$ ) and suppress the trivial rank bound  $D$ .

we arrive at the GGN’s outer product representation that lies at the heart of the ViViT concept,

$$\mathbf{G} = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \mathbf{v}_{n,c} \mathbf{v}_{n,c}^\top = \mathbf{V} \mathbf{V}^\top \quad (7.3)$$

with  $\mathbf{V} = 1/\sqrt{N} (\mathbf{v}_{1,1} \ \mathbf{v}_{1,2} \ \dots \ \mathbf{v}_{N,C}) \in \mathbb{R}^{D \times NC}$ .  $\mathbf{V}$  allows for *exact* computations with the explicit GGN matrix, at linear rather than quadratic memory cost in  $D$ . We first formulate the extraction of relevant GGN properties from this factorization, before addressing how to further approximate  $\mathbf{V}$  to reduce memory and computation costs.

### 7.2.1 Computing the Full GGN Eigenspectrum

Each GGN eigenvalue  $\lambda \in \mathbb{R}$  is a root of the characteristic polynomial  $\det(\mathbf{G} - \lambda \mathbf{I}_D)$  with identity matrix  $\mathbf{I}_D \in \mathbb{R}^{D \times D}$ . Leveraging the factorization of [Equation \(7.3\)](#) and the matrix determinant lemma, the  $D$ -dimensional eigenproblem reduces to that of the much smaller Gram matrix  $\tilde{\mathbf{G}} = \mathbf{V}^\top \mathbf{V} \in \mathbb{R}^{NC \times NC}$  which contains pairwise scalar products of  $\mathbf{v}_{n,c}$  (see [Appendix D.1.1](#)),

$$\det(\mathbf{G} - \lambda \mathbf{I}_D) = 0 \iff \det(\tilde{\mathbf{G}} - \lambda \mathbf{I}_{NC}) = 0. \quad (7.4)$$

With at least  $D - NC$  trivial solutions, the GGN curvature is zero along most directions in parameter space. Nontrivial solutions that give rise to curved directions are fully-contained in the Gram matrix, and hence *much* cheaper to compute.

Despite various Hessian spectral studies which rely on iterative eigensolvers and implicit matrix multiplication [3, 55, 61, 123, 139, 140, 177], we are not aware of works that efficiently extract the *exact* GGN spectrum from its Gram matrix. In contrast to those techniques, this matrix can be computed in parallel with gradients in a single backward pass, which results in less sequential overhead. We demonstrate in [Section 7.3.1](#) that exploiting the low-rank structure for computing the leading eigenpairs is superior to a power iteration based on matrix-free multiplication in terms of run time.

Eigenvalues themselves can help identify reasonable hyperparameters, like learning rates [97]. But we can also reconstruct the associated eigenvectors. These are directions along which curvature information is contained in the mini-batch. Let  $\tilde{\mathbb{S}}_+ = \{(\lambda_k, \tilde{\mathbf{e}}_k) \mid \lambda_k \neq 0, \tilde{\mathbf{G}}\tilde{\mathbf{e}}_k = \lambda_k \tilde{\mathbf{e}}_k\}_{k=1}^K$  denote the nontrivial Gram spectrum<sup>2</sup> with orthonormal eigenvectors  $\tilde{\mathbf{e}}_j^\top \tilde{\mathbf{e}}_k = \delta_{j,k}$  ( $\delta$  represents the Kronecker delta and  $K = \text{rank}(G)$ ). Then, the transformed vectors  $\mathbf{e}_k = 1/\sqrt{\lambda_k} \mathbf{V}\tilde{\mathbf{e}}_k$  ( $k = 1, \dots, K$ ) are orthonormal eigenvectors of  $\mathbf{G}$  associated to eigenvalues  $\lambda_k$  (see [Appendix D.1.2](#)), *i.e.* for all  $(\lambda_k, \tilde{\mathbf{e}}_k) \in \tilde{\mathbb{S}}_+$

$$\tilde{\mathbf{G}}\tilde{\mathbf{e}}_k = \lambda_k \tilde{\mathbf{e}}_k \implies \mathbf{G}\mathbf{e}_k = \lambda_k \mathbf{e}_k. \quad (7.5)$$

2: In the following, we assume ordered eigenvalues, *i.e.*  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_K$ , for convenience.

3: [Appendix D.3.2](#) describes implicit multiplication with  $\mathbf{G}^{-1}$ .

The eigenspectrum also provides access to the GGN’s pseudo-inverse based on  $\mathbf{V}$  and  $\tilde{\mathbb{S}}_+$ , required by *e.g.* second-order methods.<sup>3</sup>

## 7.2.2 Computing Directional Derivatives

Various algorithms rely on a local quadratic approximation of the loss landscape. For instance, optimization methods adapt their parameters by stepping into the minimum of the local proxy. Curvature, in the form of the Hessian or GGN, allows to build a quadratic model given by the Taylor expansion. Let  $m_{\theta_t}$  denote the quadratic model for the loss around position  $\theta_t \in \Theta$  that uses curvature represented by the GGN,

$$m_{\theta_t}(\theta) = \text{const} + (\theta - \theta_t)^\top g(\theta_t) + \frac{1}{2}(\theta - \theta_t)^\top G(\theta_t)(\theta - \theta_t). \quad (7.6)$$

At its base point  $\theta_t$ , the shape of  $m_{\theta_t}$  along an arbitrary normalized direction  $e \in \Theta$  (*i.e.*  $\|e\|_2 = 1$ ) is determined by the local gradient and curvature. Specifically, the projection of [Equation \(7.6\)](#) onto  $e$  gives rise to the (scalar) first-and second-order directional derivatives

$$\gamma_e = e^\top \nabla_\theta m_{\theta_t}(\theta_t) = e^\top g(\theta_t) \in \mathbb{R}, \quad (7.7a)$$

$$\lambda_e = e^\top \nabla_\theta^2 m_{\theta_t}(\theta_t) e = e^\top G(\theta_t) e \in \mathbb{R}. \quad (7.7b)$$

As  $G$ 's characteristic directions are its eigenvectors, they form a natural basis for the quadratic model. Denoting  $\gamma_k = \gamma_{e_k}$  and  $\lambda_k = \lambda_{e_k}$  the directional gradient and curvature along eigenvector  $e_k$ , we see from [Equation \(7.7b\)](#) that the directional curvature indeed coincides with the GGN's eigenvalue.

Analogous to the gradient and GGN, the directional derivatives  $\gamma_k$  and  $\lambda_k$  inherit the sum structure of the loss function from [Equation \(7.1\)](#), *i.e.* they decompose into contributions from individual samples. Let  $\gamma_{n,k}$  and  $\lambda_{n,k}$  denote these first- and second-order derivatives contributions of sample  $x_n$  in direction  $k$ , *i.e.*

$$\gamma_{n,k} = e_k^\top g_n = \frac{\tilde{e}_k^\top V^\top g_n}{\sqrt{\lambda_k}}, \quad (7.8a)$$

$$\lambda_{n,k} = e_k^\top G_n e_k = \frac{\|V_n^\top V \tilde{e}_k\|_2^2}{\lambda_k}, \quad (7.8b)$$

where  $V_n \in \mathbb{R}^{D \times C}$  is a scaled sub-matrix of  $V$  with fixed sample index. Note that directional derivatives can be evaluated efficiently with the Gram matrix eigenvectors without explicit access to the associated directions in parameter space.

In [Equation \(7.7\)](#), gradient  $g$  and curvature  $G$  are sums over  $g_n$  and  $G_n$ , respectively, from which follows the relationship between directional derivatives and per-sample contributions  $\gamma_k = 1/N \sum_{n=1}^N \gamma_{n,k}$  and  $\lambda_k = 1/N \sum_{n=1}^N \lambda_{n,k}$ . [Figure 7.1b](#) shows a pictorial view of the quantities provided by ViViT.

Access to per-sample directional gradients  $\gamma_{n,k}$  and curvatures  $\lambda_{n,k}$  along  $G$ 's natural directions is a distinct feature of ViViT. They provide geometric information about the local loss landscape *as well as* about the model's directional curvature stochasticity over the mini-batch.

### 7.2.3 Computational Complexity

So far, we have formulated the computation of the GGN’s eigenvalues (Equation (7.4)), eigenvectors (Equation (7.5)), and per-sample directional derivatives (Equation (7.8)). Now, we analyze their computational complexity in more detail to identify critical performance factors. Those limitations can effectively be addressed with approximations that allow the costs to be decreased in a fine-grained fashion. We substantiate our theoretical analysis with empirical performance measurements in Section 7.3.1.

#### Relation to Gradient Computation

Machine learning libraries are optimized to backpropagate signals  $1/N \nabla_{f_n} \ell_n$  and accumulate the result into the mini-batch gradient  $\mathbf{g} = 1/N \sum_{n=1}^N [\mathbf{J}_{\theta} f_n]^T \nabla_{f_n} \ell_n$ . Each column  $\mathbf{v}_{n,c}$  of  $\mathbf{V}$  also involves applying the Jacobian, but to a different vector  $\mathbf{s}_{n,c}$  from the loss Hessian’s symmetric factorization. For popular loss functions, like square and cross-entropy loss, this factorization is analytically known and available at negligible overhead. Hence, computing  $\mathbf{V}$  basically costs  $C$  gradient computations as it involves  $NC$  backpropagations, while the gradient requires  $N$ . However, the practical overhead is expected to be smaller: computations can re-use information from BackPACK’s vectorized Jacobians and enjoy additional speedup on parallel processors like GPUs.

#### Stage-wise Discarding $\mathbf{V}$

$\mathbf{V}$ ’s columns correspond to backpropagated vectors. During backpropagation, sub-matrices of  $\mathbf{V}$ , associated to parameters in the current layer, become available once at a time and can be discarded immediately after their use. This allows for memory savings without any approximations.

One example is the Gram matrix  $\tilde{\mathbf{G}}$  formed by pairwise scalar products of  $\{\mathbf{v}_{n,c}\}_{n=1,c=1}^{N,C}$  in  $\mathcal{O}((NC)^2D)$  operations. The spectral decomposition  $\tilde{\mathbb{S}}_+$  has additional cost of  $\mathcal{O}((NC)^3)$ . Similarly, the terms for the directional derivatives in Equation (7.8) can be built up stage-wise: first-order derivatives  $\{\gamma_{n,k}\}_{n=1,k=1}^{N,K}$  require the vectors  $\{\mathbf{V}^T \mathbf{g}_n \in \mathbb{R}^{NC}\}_{n=1}^N$  that cost  $\mathcal{O}(N^2CD)$  operations. Second-order derivatives are basically for free, as  $\{\mathbf{V}_n^T \mathbf{V} \in \mathbb{R}^{C \times NC}\}_{n=1}^N$  is available from  $\tilde{\mathbf{G}}$ .

#### GGN Eigenvectors

Transforming an eigenvector  $\tilde{\mathbf{e}}_k$  of the Gram matrix to the GGN eigenvector  $\mathbf{e}_k$  through application of  $\mathbf{V}$  (Equation (7.5)) costs  $\mathcal{O}(NCD)$  operations. However, repeated application of  $\mathbf{V}$  can be avoided for sums of the form  $\sum_k (c_k / \sqrt{\lambda_k}) \mathbf{e}_k$  with arbitrary weights  $c_k \in \mathbb{R}$ . The summation can be performed in the Gram space at negligible overhead, and only the resulting vector  $\sum_k c_k \tilde{\mathbf{e}}_k$  needs to be transformed. For a practical example – computing damped Newton steps – see Appendix D.2.1.

## 7.2.4 Approximations & Implementation

Although the GGN’s representation by  $V$  has linear memory cost in  $D$ , it requires memory equivalent to  $NC$  model copies.<sup>4</sup> Of course, this is infeasible for many networks and datasets, *e.g.* ImageNet ( $C = 1000$ ). So far, our formulation was concerned with *exact* computations. We now present approximations that allow  $N$ ,  $C$  and  $D$  in the above cost analysis to be replaced by smaller numbers, enabling ViViT to trade-off accuracy and performance.

<sup>4</sup>: Our implementation uses a more memory-efficient approach that avoids expanding  $V$  for linear layers by leveraging structure in their Jacobian (see Appendix D.3.1).

### MC approximation & Curvature Sub-sampling

To reduce the scaling in  $C$ , we can approximate the factorization  $\nabla_{f_n}^2 \ell_n(\theta) = \sum_{c=1}^C s_{n,c} s_{n,c}^\top$  by a smaller set of vectors. One principled approach is to draw MC samples  $\{\tilde{s}_{n,m}\}$  with  $\mathbb{E}_m[\tilde{s}_{n,m} \tilde{s}_{n,m}^\top] = \nabla_{f_n}^2 \ell_n(\theta)$  as in [38] or Example 5.2. This reduces the scaling of backpropagated vectors from  $C$  to the number of MC samples  $M$  ( $= 1$  in the following if not specified). A common independent approximation to reduce the scaling in  $N$  is computing curvature on a mini-batch subset [25, 183].

### Parameter Groups (Block-diagonal Approximation)

Some applications, *e.g.* computing Newton steps, require  $V$  to be kept in memory for performing the transformation from Gram space into the parameter space. Still, we can reduce costs by using the GGN’s diagonal blocks  $\{G^{(l)}\}_{l=1}^L$  of each layer, rather than the full matrix  $G$ . Such blocks are available during backpropagation and can thus be used and discarded step by step. In addition to the previously described approximations for reducing the costs in  $N$  and  $C$ , this technique tackles scaling in  $D$ .

### Implementation Details

BackPACK’s functionality allows us to efficiently compute individual gradients and  $V$  in a single backward pass, using either an exact or MC-factorization of the loss Hessian. To reduce memory consumption, we extend its implementation with a protocol to support mini-batch sub-sampling and parameter groups. By hooks into the package’s extensions, we can discard buffers as soon as possible during backpropagation, effectively implementing all discussed approximations and optimizations.

In Section 7.3, we specifically address how the above approximations affect run time and memory requirements, and study their impact on structural properties of the GGN.

## 7.3 Experiments

For the practical use of the ViViT concept, it is essential that (i) the computations are efficient and (ii) that we gain an understanding of how sub-sampling noise and the approximations introduced in Section 7.2.4 alter the structural properties of the GGN. In the following, we therefore

empirically investigate ViViT’s scalability and approximation properties in the context of deep learning. The insights from this analysis substantiate ViViT’s value as a monitoring tool for deep learning optimization.

### Experimental Setting

Architectures include three deep CNNs from DeepOBS [146] (2c2d on Fashion-MNIST, 3c3d on CIFAR-10 and All-CNN-C on CIFAR-100), as well as ResNets from He et al. [68] on CIFAR-10 based on Idelbayev [74]—all architectures use cross-entropy loss. Based on the approximations presented in [Section 7.2.4](#), we distinguish the following cases:

- ▶ **mb, exact:** Exact GGN with all mini-batch samples. Backpropagates  $N C$  vectors.
- ▶ **mb, mc:** MC-approximated GGN with all mini-batch samples. Backpropagates  $N M$  vectors with  $M$  the number of MC-samples.
- ▶ **sub, exact:** Exact GGN on a subset of mini-batch samples ( $\lfloor N/8 \rfloor$  as in [183]). Backpropagates  $\lfloor N/8 \rfloor C$  vectors.
- ▶ **sub, mc:** MC-approximated GGN on a subset of mini-batch samples. Backpropagates  $\lfloor N/8 \rfloor M$  vectors with  $M$  the number of MC-samples.

#### 7.3.1 Scalability

We now complement the theoretical computational complexity analysis from [Section 7.2.3](#) with empirical studies. Results were generated on a workstation with an Intel Core i7-8700K CPU (32 GB) and one NVIDIA GeForce RTX 2080 Ti GPU (11 GB). We use  $M = 1$  in the following.

### Memory Performance

We consider two tasks:

1. **Computing eigenvalues:** The nontrivial eigenvalues  $\{\lambda_k \mid (\lambda_k, \tilde{e}_k) \in \tilde{\mathbb{S}}_+\}$  are obtained by forming and eigen-decomposing the Gram matrix  $\tilde{G}$ , allowing stage-wise discarding of  $V$  (see [Sections 7.2.1](#) and [7.2.3](#)).
2. **Computing the top eigenpair:** For  $(\lambda_1, e_1)$ , we compute the Gram matrix spectrum  $\tilde{\mathbb{S}}_+$ , extract its top eigenpair  $(\lambda_1, \tilde{e}_1)$ , and transform it into parameter space by [Equation \(7.5\)](#), *i.e.*  $(\lambda_1, e_1 = 1/\sqrt{\lambda_1} V \tilde{e}_1)$ . This requires more memory than task 1 as  $V$  must be stored.

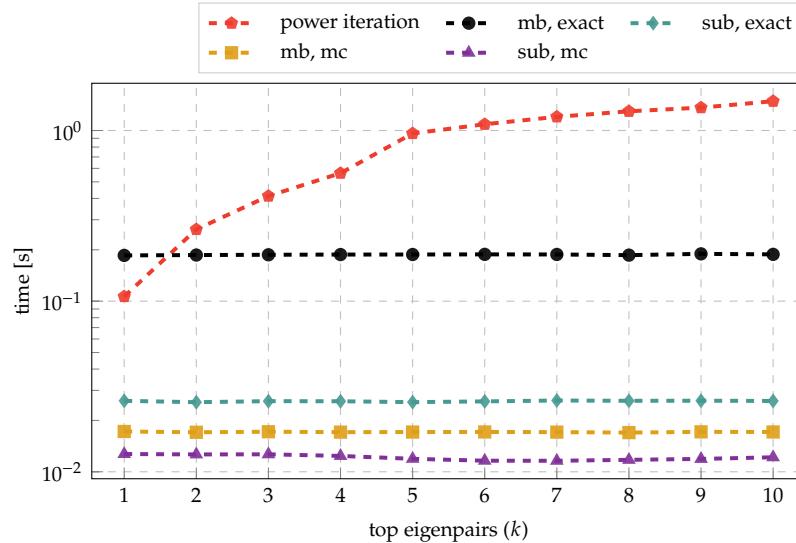
As a comprehensive memory performance measure, we use the largest batch size before our system runs out of memory—we call this the *critical batch size*  $N_{\text{crit}}$ .

[Figure 7.2a](#) tabularizes the critical batch sizes on GPU for the 3c3d architecture on CIFAR-10. As expected, computing eigenpairs requires more memory and leads to consistently smaller critical batch sizes in comparison to computing only eigenvalues. Yet, they all exceed the traditional batch size used for training ( $N = 128$ , see [146]), even when using the exact GGN. With ViViT’s approximations, the memory overhead can be reduced to significantly increase the applicable batch size.

(a) Memory performance

$N_{\text{crit}}$ (eigenvalues)			$N_{\text{crit}}$ (top eigenpair)		
GGN <sup>Data</sup>	mb	sub	GGN <sup>Data</sup>	mb	sub
exact	909	4375	exact	677	3184
mc	3840	6626	mc	3060	6029

(b) Run time performance



**Figure 7.2:** GPU memory and run time performance: Performance measurements for the 3c3d architecture ( $D = 895,210$ ) on CIFAR-10 ( $C = 10$ ). (a) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (b) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a batch of size  $N = 128$ .

We report similar results for more architectures, a block-diagonal approximation (as in Zhang et al. [183]), and on CPU in Appendix D.2.1, where we also benchmark a third task—computing damped Newton steps.

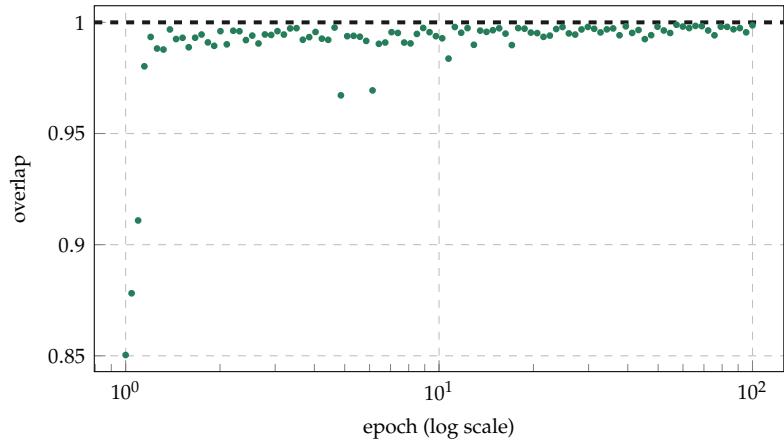
### Run Time Performance

Next, we consider computing the  $k$  leading eigenvectors and eigenvalues of a matrix. A power iteration that computes eigenpairs iteratively via matrix-vector products serves as a reference. For a fixed value of  $k$ , we repeat both approaches 20 times and report the shortest time.

For the power iteration, we adapt the implementation from the PyHessian library [177] and replace its Hessian-vector product by a matrix-free GGN-vector product [148] through PyTorch’s AD. We use the same default hyperparameters for the termination criterion. Similar to task 1, our method obtains the top- $k$  eigenpairs<sup>5</sup> by computing  $\tilde{\mathbb{S}}_+$ , extracting its leading eigenpairs and transforming the eigenvectors  $\tilde{e}_1, \tilde{e}_2, \dots, \tilde{e}_k$  into parameter space by application of  $V$  (see Equation (7.5)).

Figure 7.2b shows the GPU run time for the 3c3d architecture on CIFAR-10, using a mini-batch of size  $N = 128$ . Without any approximations to the GGN, our method already outperforms the power iteration for  $k > 1$  and increases *much* slower in run time as more leading eigenpairs are requested. This means that, relative to the transformation of each eigenvector from the Gram space into the parameter space through  $V$ , the run time mainly results from computing  $V$ ,  $\tilde{G}$ , and eigendecomposing the latter. This is consistent with the computational complexity of those operations in  $NC$  (compare Section 7.2.3) and allows for efficient extraction of a large number of eigenpairs. The run time curves of the

5: In contrast to the power iteration that is restricted to dominating eigenpairs, our approach allows choosing arbitrary eigenpairs.



**Figure 7.3: Full-batch GGN versus full-batch Hessian:** Overlap between the top- $C$  eigenspaces of the full-batch GGN and full-batch Hessian during training of the 3c3d network on CIFAR-10 with SGD.

approximations confirm this behavior by featuring the same flat profile. Additionally, they require significantly less time than the exact mini-batch computation. Results for more network architectures, a block-diagonal approximation and on CPU are reported in [Appendix D.2.1](#).

### 7.3.2 Approximation Quality

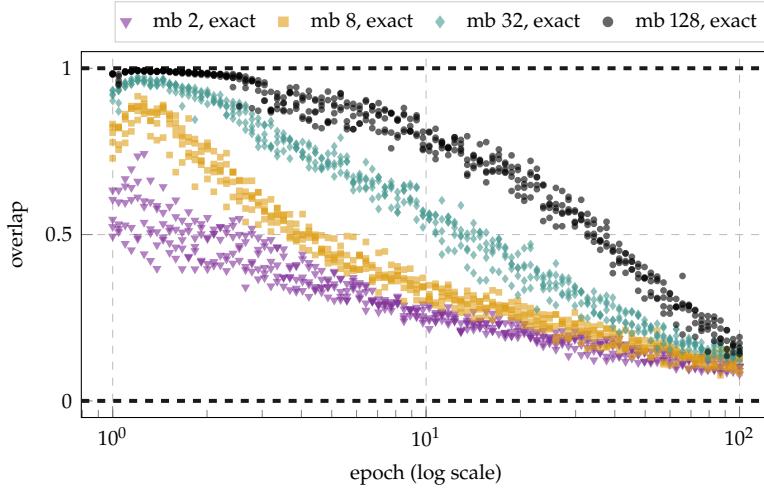
ViViT is based on the Hessian’s generalized Gauss-Newton approximation (see [Equation \(7.2\)](#)). In practice, the GGN is only computed on a mini-batch which yields a statistical estimator for the *full-batch GGN* (*i.e.* the GGN evaluated on the entire training set). Additionally, we introduce curvature sub-sampling and an MC approximation (see [Section 7.2.4](#)), *i.e.* further approximations that alter the curvature’s structural properties. In this section, we compare quantities at different stages within this hierarchy of approximations. We use the test problems from above and train the networks with both SGD and Adam (details in [Appendix D.2.2](#)).

#### GGN Versus Hessian

First, we empirically study the relationship between the GGN and the Hessian in the deep learning context. To capture *solely* the effect of neglecting the residual  $R$  (see [Equation \(7.2\)](#)), we consider the noise-free case and compute  $H$  and  $G$  on the entire training set.

We characterize both curvature matrices by their top- $C$  eigenspace: the space spanned by the eigenvectors to the  $C$  largest eigenvalues. This is a  $C$ -dimensional subspace of the parameter space  $\Theta$ , on which the loss function is subject to particularly strong curvature. The *overlap* between these spaces serves as the comparison metric. Let  $\{e_c^U\}_{c=1}^C$  the set of orthonormal eigenvectors to the  $C$  largest eigenvalues of some symmetric matrix  $U$  and  $\mathcal{E}^U = \text{span}(e_1^U, \dots, e_C^U)$ . The projection onto this subspace  $\mathcal{E}^U$  is given by the projection matrix  $P^U = (e_1^U, \dots, e_C^U)(e_1^U, \dots, e_C^U)^\top$ . As in Gur-Ari et al. [65], we define the overlap between two top- $C$  eigenspaces  $\mathcal{E}^U$  and  $\mathcal{E}^V$  of the matrices  $U$  and  $V$  by

$$\text{overlap}(\mathcal{E}^U, \mathcal{E}^V) = \frac{\text{Tr}(P^U P^V)}{\sqrt{\text{Tr}(P^U) \text{Tr}(P^V)}} \in [0, 1]. \quad (7.9)$$



**Figure 7.4: Mini-batch GGN versus full-batch GGN:** Overlap between the top- $C$  eigenspaces of the mini-batch GGN and full-batch GGN during training of the 3c3d network on CIFAR-10 with SGD. For each mini-batch size, 5 different mini-batches are drawn.

If  $\text{overlap}(\mathcal{E}^U, \mathcal{E}^V) = 0$ , then  $\mathcal{E}^U$  and  $\mathcal{E}^V$  are orthogonal to each other; if the overlap is 1, the subspaces are identical.

Figure 7.3 shows the overlap between the full-batch GGN and Hessian during training of the 3c3d network on CIFAR-10 with SGD. Except for a short phase at the beginning of the optimization procedure (note the log scale for the epoch-axis), a strong agreement (overlap  $\geq 0.85$ ) between the top- $C$  eigenspaces is observed. We make similar observations with the other test problems (see Appendix D.2.3), yet to a slightly lesser extent for CIFAR-100. Consequently, we identify the GGN as an interesting object, since it consistently shares relevant structure with the Hessian matrix.

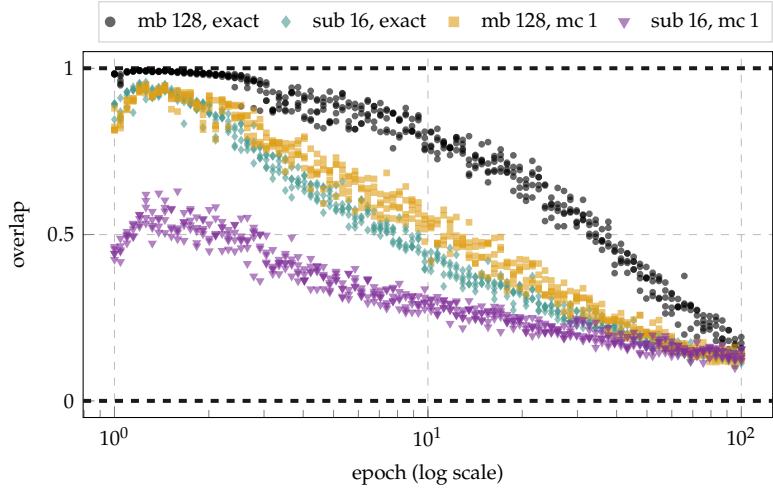
### Eigenspace Under Noise & Approximations

ViViT uses mini-batching to compute a statistical estimator of the full-batch GGN. This approximation alters the top- $C$  eigenspace, as shown in Figure 7.4: with decreasing mini-batch size, the approximation carries less and less structure of its full-batch counterpart, as indicated by dropping overlaps. In addition, at constant batch size, a decrease in approximation quality can be observed over the course of training. This might be a valuable insight for the design of second-order optimization methods, where this structural decay could lead to performance degradation over the course of the optimization, which has to be compensated for by a growing batch-size (e.g. Martens [106] reports that the optimal batch size grows during training).

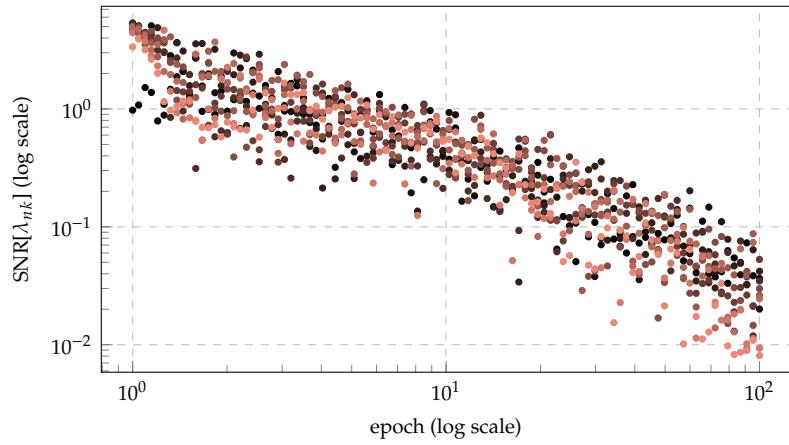
To allow for a fine-grained cost-accuracy trade-off, ViViT introduces further approximations to the mini-batch GGN (see Section 7.2.4). Figure 7.5 shows the overlap between these GGN approximations and the full-batch GGN<sup>6</sup>. The order of the approximations is as expected: with increasing computational effort, the approximations improve and, despite the greatly reduced computational effort compared to the exact mini-batch GGN, significant structure of the top- $C$  eigenspace is preserved. Details and results for the other test problems are reported in Appendix D.2.4.

So far, our analysis is based on the top- $C$  eigenspace of the curvature matrices. We extend it by studying the effect of noise and approximations on the curvature *magnitude* along the top- $C$  directions in Appendix D.2.5.

6: A comparison with the mini-batch GGN as ground truth can be found in Appendix D.2.4



**Figure 7.5:** Approximations versus full-batch GGN: Overlap between the top-C eigenspaces of the mini-batch GGN, ViViT’s approximations and the full-batch GGN during training of the 3c3d network on CIFAR-10 with SGD. Each approximation is evaluated on 5 mini-batches.



**Figure 7.6:** Directional curvature SNRs: Curvature SNRs along each of the mini-batch GGN’s top-C eigenvectors during training of the 3c3d network on CIFAR-10 with SGD. At fixed epoch, the SNR for the most curved direction is shown in ● and the SNR for the direction with the smallest curvature is shown in ○.

### 7.3.3 Per-sample Directional Derivatives

A unique feature of ViViT’s quantities is that they provide a notion of *curvature uncertainty* through *per-sample* first- and second-order directional derivatives (Equation (7.8)). To quantify noise in these derivatives, we compute their signal-to-noise ratios (SNRs). For each direction  $e_k$ , the SNR is given by the squared empirical mean divided by the empirical variance of the  $N$  mini-batch samples  $\{\gamma_{n,k}\}_{n=1}^N$  and  $\{\lambda_{n,k}\}_{n=1}^N$ , respectively.

Figure 7.6 shows curvature SNRs during training the 3c3d network on CIFAR-10 with SGD. The curvature signal along the top-C eigenvectors decreases from  $\text{SNR} > 1$  by two orders of magnitude. In comparison, the directional gradients do not exhibit such a pattern (see Appendix D.2.6). Results for the other test cases can be found in Appendix D.2.6.

In this section, we have given a glimpse of the *very rich* quantities that can be efficiently computed under ViViT’s concept. In Section 7.5, we discuss their practical use—curvature uncertainty in particular.

## 7.4 Related Work

### GGN Spectrum & Low-rank Structure

Other works point out the GGN’s low-rank structure. Botev et al. [21] present the rank bound ( $NC$ ) and propose an alternative to K-FAC based on backpropagating a decomposition of the loss Hessian. Papyan [122] presents the factorization in [Equation \(7.3\)](#) and studies the eigenvalue spectrum’s hierarchy for cross-entropy loss. In this setting, the GGN further decomposes into summands, some of which are then analyzed through similar Gram matrices. These can be obtained as contractions of  $\tilde{G}$ , but our approach goes beyond them as it does not neglect terms. We are not aware of works that obtain the exact spectrum *and* leverage a highly-efficient fully-parallel implementation. This may be because, until recently [23, 38], vectorized Jacobians required to perform those operations efficiently were not available.

### Efficient Operations with Low-rank Matrices in Deep Learning

Chen et al. [30] use [Equation \(7.3\)](#) for element-wise evaluation of the GGN in FCNNs. They also present a variant based on MC sampling. This element-wise evaluation is then used to construct hierarchical matrix approximations of the GGN. ViViT instead leverages the global low-rank structure that also enables efficient eigendecomposition.

Another prominent low-rank matrix in deep learning is the un-centered gradient covariance (sometimes called empirical Fisher). Singh and Alistarh [154] describe implicit multiplication with its inverse and apply it for neural network compression, assuming the empirical Fisher as a Hessian proxy. However, this assumption has limitations, specifically for optimization [92]. In principle though, the low-rank structure also permits the application of our methods from [Section 7.2](#).

## 7.5 Use Cases

Aiming to provide a well-founded, theoretical and empirical evaluation, we have consciously focused on studying the approximation quality of ViViT’s quantities, as well as on demonstrating the efficiency of their computation. We believe it is interesting in itself that the low-rank structure provides access to quantities that would otherwise be costly. Still, we want to briefly address possible use cases—their full development and assessment, however, will amount to separate paper(s):

- ▶ **Monitoring tool:** Our computationally efficient curvature model provides geometric *and* stochastic information about the local loss landscape and can be used by tools like Cockpit [147] to debug optimizers or to gain insights into the optimization problem itself (as in [Sections 7.3.2](#) and [7.3.3](#)).
- ▶ **Second-order optimization:** The quantities provided by ViViT, in particular the first- and second-order directional derivatives, can be used to build a stochastic quadratic model of the loss function and perform Newton-like parameter updates. In contrast to existing

second-order methods, *per-sample* quantities contain information about the reliability of that quadratic model. This offers a new dimension for improving second-order methods through statistics on the mini-batch *distribution* of the directional derivatives (*e.g.* for variance-adapted step sizes), potentially increasing the method’s performance and stability.

## 7.6 Conclusion

We have presented ViViT, a curvature model based on the low-rank structure of the Hessian’s generalized Gauss-Newton (GGN) approximation. This structure allows for efficient extraction of *exact* curvature properties, such as the GGN’s full eigenvalue spectrum and directional gradients and curvatures along the associated eigenvectors. ViViT’s quantities scale by approximations that allow for a fine-grained cost-accuracy trade-off. In contrast to alternatives, these quantities offer a notion of curvature uncertainty across the mini-batch in the form of directional derivatives.

We empirically demonstrated the efficiency of leveraging the GGN’s low-rank structure and substantiated its usefulness by studying characteristics of curvature noise on various deep learning architectures.

The low-rank representation is efficiently computed in parallel with gradients during a single backward pass. As it mainly relies on vectorized Jacobians, it is general enough to be integrated into existing machine learning libraries in the future. For now, we provide an efficient open-source implementation in PyTorch [126] by extending the existing BackPACK [38] library.

## **Part III.**

### **Conclusion & Future Directions**



# 8.

## Conclusion & Future Directions

Contemporary deep learning is powered by methods that solely rely on the gradient. This is reflected in popular machine learning libraries which prioritize its computation. However, it narrows research to focus on gradient-based algorithms that are not agnostic to the empirical risk's stochasticity and geometry beyond first order. To advance the field, we need to explore the potential of higher-order information beyond the gradient. One main hindrance to further explore its utility has been that it is complicated to implement, which makes it difficult for practitioners to try it out. Therefore, one major goal of this work was to ease experimentation with higher-order information by making it as conveniently accessible as the gradient. This thesis demonstrates that rich information beyond the gradient is *affordable*, can be made *readily available* in existing machine learning libraries, and is *useful* to enable novel approaches for advancing deep learning:

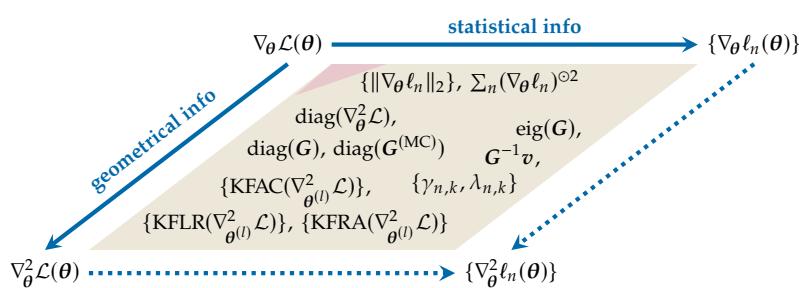
### (Q1) Which information beyond the gradient is efficiently accessible?

Figure 8.1 provides an overview of the higher-order information made accessible in this work: per-sample gradients—whose empirical mean is the mini-batch gradient—can be explicitly computed, or reduced into higher-order statistics like the gradient variance. Light-weight structural approximations through diagonal and Kronecker matrices enable the computation of per-layer second-order derivatives in the Hessian or generalized Gauss-Newton. The generalized Gauss-Newton's outer product structure allows to go beyond per-layer terms, even to compute with the full matrix, and to access curvature noise.

### (Q2) How to compute this information—conveniently, automatically, and efficiently—re-using the existing backpropagation implementation of ML frameworks?

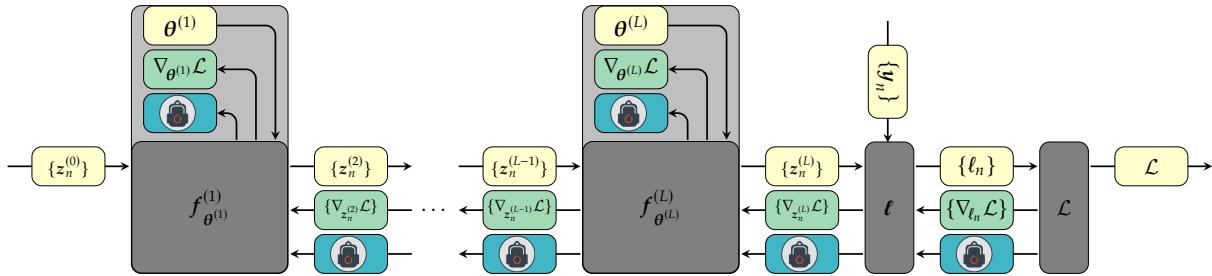
All quantities presented in this thesis are phrased as extensions of a standard backward pass.

Therefore, they can be computed *efficiently* and at the same time as the gradient (Figure 8.2). Gradient statistics share most computations with the gradient, and can recycle information from the standard backward



8.1 Summary & Impact . . . . .	107
8.2 Future Work . . . . .	108

**Figure 8.1:** Higher-order information made available by this work (extends Figure 1.2). Quantities are roughly mapped onto the landscape spanned by the stochasticity and geometry of the loss. They can be categorized as (i) gradient statistics, (ii) diagonal curvature approximations, (iii) Kronecker-factored curvature approximations, and (iv) noise-aware curvature from the generalized Gauss-Newton's low-rank structure.



**Figure 8.2:** All information in Figure 8.1 is an extension of the standard backward pass for the gradient. Forward (○) and backward (●) pass are implemented by machine learning libraries. Minimal invasion by adding an extended backward pass (●) allows to compute various other quantities. The specifics (◎) on what and how to backpropagate depend on the quantity of interest.

pass. Approximate second-order derivatives require sending additional information through the computation graph. In runtime, the additional work is often significantly reduced through hardware parallelism.

Their formulations vary in what, and how, objects are being backpropagated through layers, and how the target quantities are extracted for each parameter. As these rules are defined on a per-module basis (Figure 8.3), the approach is extensible and guarantees *automatic* computation.

From an implementation perspective, this is achieved through lightweight extension of an existing gradient backpropagation implementation without implementing a new framework. This enables easy integration into existing machine learning libraries, and makes it *convenient* for practitioners to extend their code at minimal overhead.

**Figure 8.3:** Modularity is crucial for automation and extensibility. Given a functioning backpropagation procedure, the extended backward pass in Figure 8.2 reduces to specifying a set of rules for each module and each quantity. This simplifies adding new operations, while preserving automatic computation.

In addition to second-order optimization, various deep learning methods require higher-order information: Laplace approximations [e.g. 40], model compression (pruning) [e.g. 154], differential privacy [e.g. 2], importance sampling [e.g. 85], variance adaptation [e.g. 9], parameter initialization [e.g. 155], batch size adaptation [e.g. 11], etc. Some of them were briefly outlined in this thesis to motivate the demand that this information be more readily available. Easier access to this information enables more efficient and creative research in these areas and helps establish the resulting methods through user-friendly implementations.

As a use case, this thesis focused on enabling a deeper look into the inner workings of neural networks during training through the lens of higher-order information. This allows to identify common failure modes, which makes training less painful, and thereby highlights the utility of higher-order information for deep learning.

Going further, this work also underlines the unexplored nature of higher-order information: the developed extended automatic differentiation functionality enables novel efficient computational schemes that address shortcomings in existing approaches. This allows to build more powerful local approximations of the loss landscape which are agnostic to noise, curvature, and even noise in the curvature. Such works help to identify important techniques to reduce run time, make algorithms that rely on such information more competitive with gradient-based methods, and shape the development of future machine learning frameworks.

## 8.1 Summary & Impact

### Extending Backpropagation to the Hessian

Gradient backpropagation is efficient, automated, and extensible. Chapter 4 presented how—for sequential feedforward architectures—this carries over to second-order derivatives: just like backpropagation recovers the layer-wise gradients, Hessian backpropagation recovers the per-layer Hessians. It fully aligns the computation of local Hessians with gradients and unifies the view on block-diagonal curvature approximations like the block-diagonal generalized Gauss–Newton [148], Fisher [5], and its Kronecker-factorized approximations [21, 31, 63, 109].

### Packing More into Backprop

The BackPACK library, presented in Chapter 5, provides efficient access to various deep learning quantities through implementing the insights on extended backpropagation for the Hessian on top of PyTorch. During a standard backward pass that computes the average gradient, it extracts (i) per-sample gradients and gradient statistics, and (ii) approximate second-order derivatives in the form of diagonal and Kronecker-factorized curvature. This often adds only little overhead. BackPACK easily integrates into existing code and simplifies experimentation with the above quantities: it powers other libraries for Bayesian applications with Laplace approximations [40, 76], out-of-distribution generalization [64, 130], and differential privacy [179], as well as the follow-up works in this thesis (Cockpit [147], ViViT [39]). More than two years after its release, the library is still actively used, with multiple hundred downloads per week at the time of writing (July 2022).

### Enabling a Closer Look Into Neural Nets

Higher-order information as provided by BackPACK is valuable to guide neural network training. Common methods for real-time training diagnostics, such as monitoring the loss, are limited because they only indicate whether a model is training, but not why. The Cockpit library, presented in Chapter 6, enables a closer look into neural networks during training. The live-monitoring tool visualizes established, recently proposed [11, 19, 26, 100, 104, 162, 177], and novel summary statistics that are efficiently computed by BackPACK. It allows to identify common bugs in the machine learning pipeline, such as improper data pre-processing or vanishing gradients, but also to guide learning rate selection, and to study implicit regularization [56, 112]. This showcases the potential of higher-order information to assist practitioners.

### Enabling Novel Ways to Compute with Curvature

BackPACK’s extended automatic differentiation functionality enables algorithmic advances to tackle limitations of existing curvature proxies: diagonal or Kronecker-factorized curvatures are (i) not agnostic to noise in the mini-batch, (ii) strict approximations that do not become exact in

any limit, and (iii) restricted to the block diagonal. ViViT’s quantities, presented in [Chapter 7](#), address this through the generalized Gauss–Newton’s low-rank structure, which allows for exact computation with the full—rather than block-diagonal—matrix, and principled approximations to reduce cost in exchange for less accuracy. ViViT enables efficient computation of spectral properties, as well as directional gradients and curvatures on a per-sample basis that quantify noise. Monitoring this noise through signal-to-noise ratios helps understand its characteristics in deep learning [48] and to identify challenges for optimization and generalization from the interplay between noise and curvature [162].

## 8.2 Future Work

Deep learning needs more than just the gradient. To leverage the full potential of higher-order information, we need to (i) build more—and refine existing—tools to (ii) study and better understand algorithmic challenges in deep learning, and (iii) amplify the practicality of such next-generation algorithms through user-friendly and highly-efficient implementations like gradient-based methods.

### Extending Cockpit

[Chapter 6](#) demonstrated Cockpit’s utility for identifying common bugs in the machine learning pipeline. These failures were deliberately designed on well-known, standardized machine learning problems, to illustrate Cockpit’s purpose. Since they were implicitly tuned over many years to work well with currently popular methods, they rarely exhibit failure modes. Cockpit will be even more useful for debugging unknown, non-standardized problems that have *not* undergone such tuning, and are therefore extremely likely to exhibit failures. We want to establish Cockpit as a tool for practitioners facing such problems and are looking forward to its first success stories “in the wild”.

Cockpit also provides functionality for scientific analyses of neural networks. Training such models is often wasteful in computation, *e.g.* using large grid searches whose computations are effectively discarded after identifying one well-performing set of hyperparameters. Cockpit’s summary statistics are condensed and could be stored for a large number of training trajectories corresponding to different hyperparameter settings. These trajectories could be collected into a dataset that could serve for the analysis of optimization algorithms to understand their implicit bias, to study properties of neural networks that generalize well, or for meta-learning optimization strategies.

To further improve the meaningfulness and interpretability of Cockpit’s instruments, its control over parts of the network could be made more customizable: currently, most quantities can be computed either on all parameters, or per parameter. For very large networks, it will be more practical to group parameters, and to compute and visualize Cockpit’s instruments per parameter group.

### Noise-aware Second-order Methods

Newton steps are powerful, but their stability is strongly affected by noise: one corrupted step might undo all previous progress. Improving their stability is thus one key challenge to make them work in the mini-batch setting. To do that, we need to quantify noise in the mini-batch. But popular curvature proxies used in second-order methods are not noise-agnostic. Therefore, we need curvature approximations that give access to noise, *e.g.* through per-sample information as provided through ViViT. Such information could be used to develop noise-aware stabilizing mechanisms for Newton steps, like damping.

### Optimizing Run Time & Advancing Automatic Differentiation

In contrast to gradient-based methods, the run time performance of higher-order methods can still be significantly improved: *e.g.* BackPACK outperforms naive implementations (like for-loops) and achieves practical overheads. But it relies on PyTorch’s Python API, which is sometimes not flexible enough, and therefore realizes some functionality through less efficient workarounds. Recent advances in automatic differentiation, like JAX [23] and functorch [72], rely on function transformation to achieve a clean separation of automatic differentiation and batching, and allow for more efficient implementations through just-in-time compilation.

Often, performance improvements are achieved through leveraging linear algebra, like properties of the Kronecker product [101] and matrix decompositions [*e.g.* 38, 39]. Recent work suggests that there is potential for further improvements, as a number of these optimizations are not yet realized [141]. Therefore, one direction would be to improve the automated optimization of operations in second-order methods in these libraries. This would further reduce the overhead of second-order methods stemming from poor implementation, and make these performance gains widely available to the machine learning community.



## **Part IV.**

# Appendix



# Additional Material for Chapter 4

# A.

Here, we provide additional details and derivations for Hessian back-propagation (HBP).

Appendix A.1 relates the HBP Equation (4.7) to the chain rule for matrix derivatives [103]. It relies on the clean definitions of generalized Jacobian and Hessian matrices for multi-variate functions (Equations (2.27) and (3.6)), and their chain rules (Theorems 2.2 and 3.1).

With matrix derivatives, the HBP equation for a variety of module functions can be derived elegantly. Appendices A.2 to A.4 contain the HBP derivations for all operations in Table 4.1. We split the considered operations into different categories to achieve a cleaner structure. Appendix A.2 contains details on operations used for the construction of fully-connected neural networks (FCNNs) and skip-connections. Appendix A.2.4 illustrates the analytic composition of multiple modules by combining the backward passes of a nonlinear elementwise activation function and an affine transformation. This yields the recursive schemes of Botev et al. [21] and Chen et al. [31], the latter of which has been used in the experiment of Section 4.4. The analysis of the Hessian for common loss functions is provided in Appendix A.3. Operations occurring in convolutional neural networks (CNNs) are subject of Appendix A.4.

Appendix A.5 provides details on model architectures, training procedures used in the experiments of Section 4.4, and an additional experiment on a modified test problem of the DeepOBS benchmark library [146].

## A.1 Matrix Derivatives

Index notation for higher-order derivatives of multi-variate matrix functions can become heavy [8, 31, 110, 115]. We tackle this by embedding our approach in the notation of matrix differential calculus, which

1. yields notation consistent with established literature on matrix derivatives [103] and clarifies the origin of the symbols  $J$  and  $\nabla^2$ , used extensively in the main text (Equations (2.27) and (3.6))<sup>1</sup>.
2. allows for using a multi-dimensional generalization of the chain rule (Theorems 2.2 and 3.1).
3. lets us extract first- and second-order derivatives from differentials using the identification rules of Magnus and Neudecker [103] without bothering to deal with index notation.

A.1 Matrix Derivatives . . . . .	113
A.2 HBP for FCNNs . . . . .	115
A.3 HBP for Loss Functions .	118
A.4 HBP for CNNs . . . . .	119
A.5 Experimental Details . . .	124

<sup>1</sup>: Magnus and Neudecker [103] use  $DB(A)$ ,  $HB(A)$  for the Jacobian  $J_A B$  and Hessian  $\nabla_A^2 B$  of matrix variables  $A, B$ .

With these techniques it is easy to see how structure, like Kronecker products, appears in the derivatives.

## Preliminaries & Notation

Equations (2.27) and (3.6) and theorems 2.2 and 3.1 represent a collection of results from the book of Magnus and Neudecker [103]. They generalize the concept of first- and second-order derivatives to multi-variate matrix functions in terms of the Jacobian and Hessian matrix. While there exist multiple ways to arrange the partial derivatives, the presented definitions allows for a multi-variate generalization of the chain rule.

We denote matrix, vector, and scalar functions by  $F$ ,  $f$ , and  $\phi$ , respectively. Matrix (vector) arguments are written as  $X(x)$ . Vectorization ( $\text{vec}$ , Definition 2.2) applies column-stacking, such that for matrices  $A, B, C$ ,

$$\text{vec}(ABC) = (C^\top \otimes A) \text{vec}(B). \quad (\text{A.1})$$

We assign vectors to bold lower-case ( $x, \theta, \dots$ ), matrices to bold upper-case ( $W, X, \dots$ ), and tensors to bold upper-case sans serif symbols ( $\mathbf{W}, \mathbf{X}, \dots$ ).  $\odot$  means elementwise multiplication (Hadamard product).

### Remark on Vectorization

The generalized Jacobian and Hessian from [103] rely on vectorization of matrices. Convolutional neural networks usually act on tensors and we incorporate these by assuming them to be flattened such that the first index varies fastest. For a matrix (tensor of order two), this is consistent with column-stacking. E.g., the vectorized version of the tensor  $\mathbf{A} \in \mathbb{R}^{n_1 \times n_2 \times n_3}$  with  $n_1, n_2, n_3 \in \mathbb{N}$  is  $\text{vec} \mathbf{A} = (A_{1,1,1}, A_{2,1,1}, \dots, A_{n_1,1,1}, A_{1,2,1}, \dots, A_{n_1,n_2,n_3})^\top$ . To formulate the generalized Jacobian or Hessian for tensor operations, its action on a vector or matrix view of the original tensor is considered. Consequently, all operations can be reduced to vector-valued functions, which we consider in the following.

The vectorization scheme is not unique. Most of the linear algebra literature assumes column-stacking. However, when it comes to implementations, a lot of programming languages store tensors in row-major order, corresponding to row-stacking vectorization (last index varies fastest). Thus, special attention has to be paid in implementations.

### A.1.1 Relation to the Modular Approach

Theorem 3.1 can directly be applied to the graph  $\ell \circ f_{\theta^{(L)}}^{(L)} \circ f_{\theta^{(L-1)}}^{(L-1)} \circ \dots \circ f_{\theta^{(1)}}^{(1)}$  of the sequential feedforward net under investigation. For any module function  $f_{\theta^{(l)}}^{(l)}$ , the loss can be expressed as a composition of two functions by squashing preceding modules in the graph into a single function  $f_{\theta^{(l-1)}}^{(l-1)} \circ \dots \circ f_{\theta^{(1)}}^{(1)}$ , and likewise composing the module itself and all subsequent functions, i.e.  $\ell \circ f_{\theta^{(L)}}^{(L)} \circ \dots \circ f_{\theta^{(1)}}^{(1)}$ .

The analysis can therefore be reduced to the module shown in Figure 4.2 receiving an input  $x \in \mathbb{R}^n$  that is used to compute the output  $z \in \mathbb{R}^m$ .

The scalar loss is then expressed as a mapping  $\ell(z(x), y) : \mathbb{R}^n \rightarrow \mathbb{R}^p$  with  $p = 1$ . Suppressing the label  $y$ , Equation (3.14) implies

$$\begin{aligned}\nabla_x^2 \ell &= (I_p \otimes J_x z)^\top (\nabla_z^2 \ell) J_x z + (J_z \ell \otimes I_n) \nabla_x^2 z \\ &= (J_x z)^\top (\nabla_z^2 \ell) J_x z + (J_z \ell \otimes I_n) \nabla_x^2 z.\end{aligned}\quad (\text{A.2})$$

The HBP Equation (4.7), which contains Hessians of elements of  $z$ , is obtained by substituting the form Equation (3.7) into Equation (A.2).

## A.2 HBP for FCNNs

### A.2.1 Linear Layer (Matrix-vector Multiplication, Matrix-Matrix Multiplication, Addition)

Consider the function  $f$  of a module applying an affine transformation to a vector. Apart from the input  $x$ , additional parameters of the module are given by the weight matrix  $W$  and the bias term  $b$ ,

$$\begin{aligned}f : \quad \mathbb{R}^n \times \mathbb{R}^{m \times n} \times \mathbb{R}^m &\rightarrow \mathbb{R}^m \\ (x, W, b) &\mapsto z = Wx + b.\end{aligned}$$

To compute the Jacobians *w.r.t.* each variable, we use the differentials

$$\begin{aligned}dz(x) &= Wdx, \\ dz(b) &= db, \\ dz(W) = (dW)x &\implies d\text{vec } z(W) = (x^\top \otimes I_m) \text{vec}(dW),\end{aligned}$$

using Equation (A.1) to establish the implication in the last line. With the *first identification tables* provided in Magnus and Neudecker [103, Chapter 9.6], the Jacobians can be read off from the differentials as  $J_x z = W$ ,  $J_b z = I_m$ ,  $J_W z = x^\top \otimes I_m$ . All second module derivatives  $\nabla_x^2 z$ ,  $\nabla_W^2 z$ , and  $\nabla_b^2 z$  vanish since  $f$  is linear in all inputs. Inserting the Jacobians into Equation (4.7) results in

$$\nabla_x^2 \ell = W^\top (\nabla_z^2 \ell) W, \quad (\text{A.3a})$$

$$\nabla_b^2 \ell = \nabla_z^2 \ell, \quad (\text{A.3b})$$

$$\begin{aligned}\nabla_W^2 \ell &= (x^\top \otimes I_m)^\top (\nabla_z^2 \ell) (x^\top \otimes I_m) \\ &= xx^\top \otimes \nabla_z^2 \ell = x \otimes x^\top \otimes \nabla_z^2 \ell.\end{aligned}\quad (\text{A.3c})$$

The HBP relations for matrix-vector multiplication and addition listed in Table 4.1 are special cases of Equation (A.3). HBP for matrix-matrix multiplication is derived in a completely analogous fashion.

### A.2.2 Elementwise Activation

Next, consider the elementwise application of a nonlinear function  $\phi$ ,

$$\begin{aligned}\phi : \quad \mathbb{R}^m &\rightarrow \mathbb{R}^m \\ x \mapsto z = \phi(x) &\text{ such that } [\phi(x)]_k = \phi([x]_k),\end{aligned}$$

For the matrix differential *w.r.t.*  $\mathbf{x}$ , this implies

$$d\phi(\mathbf{x}) = \phi'(\mathbf{x}) \odot d\mathbf{x} = \text{diag}[\phi'(\mathbf{x})] d\mathbf{x},$$

where  $\phi'$  means elementwise application of  $\phi'$ , and consequently, the Jacobian is a diagonal matrix  $J_x \phi(\mathbf{x}) = \text{diag}[\phi'(\mathbf{x})]$ . For the Hessian, note that the function value  $\phi_k(\mathbf{x}) := [\phi(\mathbf{x})]_k$  only depends on  $x_k := [\mathbf{x}]_k$  and thus  $\nabla_x^2 \phi_k = \phi''(x_k) \hat{\mathbf{e}}_k \hat{\mathbf{e}}_k^\top$ , with the one-hot unit vector  $\hat{\mathbf{e}}_k \in \mathbb{R}^m$  in coordinate direction  $k$ . Inserting all quantities into Equation (4.7) yields

$$\begin{aligned} \nabla_x^2 \ell &= \text{diag}[\phi'(\mathbf{x})] (\nabla_z^2 \ell) \text{diag}[\phi'(\mathbf{x})] + \sum_k \phi''(x_k) \hat{\mathbf{e}}_k \hat{\mathbf{e}}_k^\top (\nabla_{z_k} \ell) \\ &= \text{diag}[\phi'(\mathbf{x})] (\nabla_z^2 \ell) \text{diag}[\phi'(\mathbf{x})] + \text{diag}[\phi''(\mathbf{x}) \odot \nabla_z \ell], \end{aligned} \quad (\text{A.4})$$

where  $\phi''$  means elementwise application of  $\phi''$ .

### A.2.3 Skip-connection

Residual learning [68] uses skip-connection units to facilitate the training of DNNs. In its simplest form, the mapping  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$  reads

$$z(\mathbf{x}, \theta) = \mathbf{x} + s(\mathbf{x}, \theta),$$

with a potentially nonlinear operation  $(\mathbf{x}, \theta) \mapsto s$ . The input and parameter Jacobians are given by  $J_z z(\mathbf{x}) = I_m + J_x s(\mathbf{x})$  and  $J_\theta z(\theta) = J_\theta s(\theta)$ . Using Equation (4.7), one finds

$$\begin{aligned} \nabla_x^2 \ell &= [I_m + J_x s(\mathbf{x})]^\top (\nabla_z^2 \ell) [I_m + J_x s(\mathbf{x})] + \sum_k [\nabla_x^2 s_k(\mathbf{x})] (\nabla_{z_k} \ell), \\ \nabla_\theta^2 \ell &= [J_\theta s(\theta)]^\top (\nabla_z^2 \ell) [J_\theta s(\theta)] + \sum_k [\nabla_\theta^2 s_k(\theta)] (\nabla_{z_k} \ell). \end{aligned}$$

### A.2.4 Relation to Recursive Schemes in Previous Work

The modular decomposition of curvature backpropagation facilitates the analysis of modules composed of multiple operations. Now, we analyze the composition of two modules. This yields the recursive schemes presented by Botev et al. [21] (KFRA) and Chen et al. [31] (BDA-PCH).

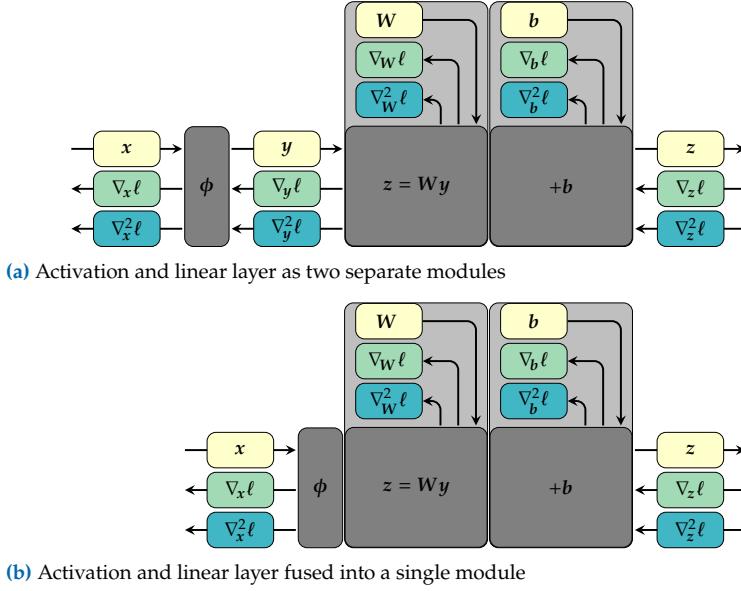
#### Analytic Composition of Multiple Modules

Consider the module  $g = f \circ \phi$ ,  $x \mapsto y := \phi(x) \mapsto z = f(y(x))$ . Assume  $\phi$  to act elementwise on the input, followed by a linear layer  $f : z(y) = W y + b$  (Figure A.1a). Analytic elimination of the intermediate backward pass recovers the backward pass of the fused module that consists of both operations (Figure A.1b). The first Hessian backward pass through the linear module  $f$  (Equation (A.3)) implies

$$\nabla_y^2 \ell = W^\top (\nabla_z^2 \ell) W, \quad (\text{A.5a})$$

$$\nabla_W^2 \ell = y \otimes y^\top \otimes \nabla_z^2 \ell = \phi(x) \otimes \phi(x)^\top \otimes \nabla_z^2 \ell, \quad (\text{A.5a})$$

$$\nabla_b^2 \ell = \nabla_z^2 \ell. \quad (\text{A.5b})$$



**Figure A.1:** The sequence of element-wise activation  $\phi$  and linear layer can be interpreted as two modules, or a single one. (a) Both operations are analyzed separately to derive the HBP. (b) Back-propagation of  $\nabla_z^2 \ell$  is expressed in terms of  $\nabla_x^2 \ell$  without intermediate message.

Further backpropagation through  $\phi$  with Equation (A.4) results in

$$\begin{aligned}
 \nabla_x^2 \ell &= \text{diag}[\phi'(x)] \left( \nabla_y^2 \ell \right) \text{diag}[\phi'(x)] + \text{diag}[\phi''(x) \odot (\nabla_y \ell)] \\
 &= \text{diag}[\phi'(x)] [W^\top (\nabla_z^2 \ell) W] \text{diag}[\phi'(x)] \\
 &\quad + \text{diag}[\phi''(x) \odot W^\top (\nabla_z \ell)] \\
 &= \{W \text{diag}[\phi'(x)]\}^\top (\nabla_z^2 \ell) \{W \text{diag}[\phi'(x)]\} \\
 &\quad + \text{diag}[\phi''(x) \odot W^\top (\nabla_z \ell)].
 \end{aligned} \tag{A.5c}$$

We use invariance of a diagonal matrix under transposition and  $\nabla_y \ell = W^\top (\nabla_z \ell)$  for the backpropagated gradient for the last equality. The Jacobian  $J_x g(x)$  of the module shown in Figure A.1b is  $J_x g(x) = W \text{diag}[\phi(x)] = [W^\top \odot \phi'(x)]^\top$  (broadcasting  $\phi'(x)$ ). In summary, HBP for the composite layer  $z(x) = W\phi(x) + b$  is given by Equation (A.5).

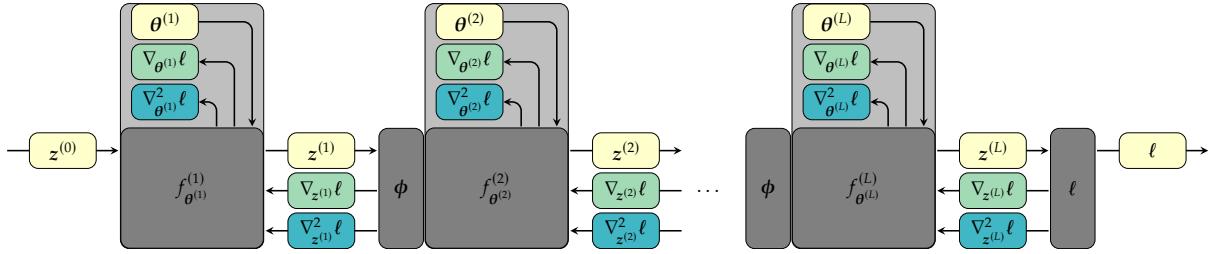
### Obtaining the Relations of KFRA and BDA-PCH

The derivations for the composite module given above are closely related to the recursive schemes of Botev et al. [21] and Chen et al. [31]. Their relations are obtained from a straightforward conversion of the HBP rules equation A.5. Consider a sequence of a linear layer  $f_{\theta^{(1)}}^{(1)}$  and multiple composite modules  $f_{\theta^{(2)}}, \dots, f_{\theta^{(L)}}$  as shown in Figure A.2.

According to Equation (A.5b) both the linear layer and the composite  $f_{\theta^{(l)}}$  identify the gradient (Hessian) w.r.t. their outputs,  $\nabla_{z^{(l)}} \ell$  ( $\nabla_{z^{(l)}}^2 \ell$ ), as the gradient (Hessian) w.r.t. their bias term,  $\nabla_{b^{(l)}} \ell$  ( $\nabla_{b^{(l)}}^2 \ell$ ). Introducing layer indices for all quantities, one finds the recursion

$$\nabla_{b^{(l)}}^2 \ell = \nabla_{z^{(l)}}^2 \ell, \tag{A.6a}$$

$$\nabla_{W^{(l)}}^2 \ell = \phi(z^{(l-1)}) \otimes \phi(z^{(l-1)})^\top \otimes \nabla_{b^{(l)}}^2 \ell, \tag{A.6b}$$



**Figure A.2:** Grouping scheme for the recursive Hessian computation proposed by KFRA and BDA-PCH. The backward messages between the linear layer and the preceding nonlinear activation are analytically fused.

for \$l = L - 1, \dots, 1\$, and

$$\begin{aligned} \nabla_{z^{(l-1)}}^2 \ell &= \left\{ \mathbf{W}^{(l)} \operatorname{diag} [\phi'(z^{(l-1)})] \right\}^\top \nabla_{b^{(l)}}^2 \ell \left\{ \mathbf{W}^{(l)} \operatorname{diag} [\phi'(z^{(l-1)})] \right\} \\ &\quad + \operatorname{diag} [\phi''(z^{(l-1)}) \odot \mathbf{W}^{(l)\top} \nabla_{b^{(l)}} \ell] \\ &= \left\{ \mathbf{W}^{(l)\top} \odot \phi'(z^{(l-1)}) \right\} \nabla_{b^{(l)}}^2 \ell \left\{ \mathbf{W}^{(l)\top} \odot \phi'(z^{(l-1)}) \right\}^\top \\ &\quad + \operatorname{diag} [\phi''(z^{(l-1)}) \odot \mathbf{W}^{(l)\top} \nabla_{b^{(l)}} \ell] \end{aligned} \quad (\text{A.6c})$$

for \$l = L - 1, \dots, 2\$. It is initialized with the loss function's gradient (Hessian) \$\nabla\_{z^{(L)}} \ell\$ (\$\nabla\_{z^{(L)}}^2 \ell\$).

Equation (A.6) are equivalent to the expressions provided in [21, 31]. Their emergence from compositions of HBP equations of simple operations represents one key insight of Chapter 4. Both works use the batch average strategy presented in Section 4.3.2 to obtain curvature estimates.

## A.3 HBP for Loss Functions

### A.3.1 Square Loss

Square loss of the model prediction \$f \in \mathbb{R}^C\$ and the true label \$y \in \mathbb{R}^C\$ is computed by (Equation (2.1))

$$\ell(f, y) = \frac{1}{C} (y - f)^\top (y - f).$$

Differentiating

$$d\ell(f) = \frac{1}{C} [-(df)^\top (y - f) - (y - f)^\top df] = -\frac{2}{C} (y - f)^\top df$$

once more yields

$$d^2\ell(f) = \frac{2}{C} (df)^\top df = \frac{2}{C} (df)^\top I_C (df).$$

The Hessian is extracted with the *second identification tables* from Magnus and Neudecker [103, Chapter 10.4] and reproduces the expected result

$$\nabla_f^2 \ell(f) = \frac{2}{C} I_C. \quad (\text{A.7})$$

### A.3.2 Softmax Cross-entropy Loss

The computation of cross-entropy from logits (Equation (2.2)) is composed of two operations. First, the neural network outputs are transformed into log-probabilities by the softmax function. Then, the cross-entropy with the distribution implied by the label is computed.

#### Log-softmax

The output's elements  $f \in \mathbb{R}^C$  of a neural network are assigned to log-probabilities  $z(f) = \log p(f) \in \mathbb{R}^C$  by means of the softmax function  $p(f) = \exp(f)/\sum_i \exp(f_i)$ . Consequently,

$$z(f) = f - \log \left[ \sum_i \exp(f_i) \right],$$

and the Jacobian reads  $J_f z(f) = I_C - \mathbf{1}_C p(f)^\top$  with  $\mathbf{1}_C \in \mathbb{R}^C$  denoting a vector of ones. The log-probability Hessians w.r.t.  $f$  are given by

$$\nabla_f^2[z(f)]_k = -\text{diag}[p(f)] + p(f)p(f)^\top.$$

#### Cross-entropy

The negative log-probabilities are used for the cross-entropy with the probability distribution of the label  $y \in \{1, \dots, c\}$ ,

$$\ell(z, y) = -\text{onehot}(y)^\top z.$$

Since  $\ell$  is linear in the log-probabilities  $z$ , i.e.  $\nabla_z^2 \ell(z) = \mathbf{0}$ , the HBP is

$$\begin{aligned} \nabla_f^2 \ell(f) &= [J_f z(f)]^\top \nabla_z^2 \ell(z) [J_f z(f)] + \sum_k \nabla_f^2[z(f)]_k \frac{\partial \ell(z)}{\partial [z]_k} \\ &= \{-\text{diag}[p(f)] + p(f)p(f)^\top\} \sum_k [-\text{onehot}(y)]_k \\ &= \text{diag}[p(f)] - p(f)p(f)^\top. \end{aligned}$$

## A.4 HBP for CNNs

The recursive approaches in [21, 31] tackle the computation of curvature blocks for FCNNs. To the best of our knowledge, an extension to CNNs has not been achieved so far. One reason might be that convolutions come with heavy notation that is difficult to deal with in index notation.

Martens and Grosse [109] provide a procedure for computing a Kronecker-factored approximation of the Fisher for convolutions (KFC). This scheme relies on the property of the Fisher to describe the covariance of the log-likelihood's gradients under the model's distribution. Curvature information is thus encoded in the expectation value, and not by back-propagation of second-order derivatives.

To derive the HBP for convolutions, we use that an efficient implementation of the forward pass is decomposed into multiple operations (see

[Figure A.4](#)), which can be considered independently by means of our modular approach (see [Figure A.5](#) for details). Our analysis starts by considering the backpropagation of curvature through operations that occur frequently in CNNs. This includes the reshape ([Appendix A.4.1](#)) and extraction operation ([Appendix A.4.2](#)). In [Appendix A.4.3](#) we outline the modular decomposition and curvature backpropagation of convolution in two dimensions. The approach carries over to other dimensions.

All operations under consideration in this section are linear. Hence the second terms in [Equation \(4.7\)](#) vanish. Again, we use the framework of matrix differential calculus [103] to avoid index notation.

### A.4.1 Reshape/View

The reshape operation reinterprets a tensor  $\mathbf{X} \in \mathbb{R}^{n_1 \times \dots \times n_x}$  as another tensor  $\mathbf{Z} \in \mathbb{Z}^{m_1 \times \dots \times m_z}$  ([Definition 2.3](#))

$$\mathbf{Z}(\mathbf{X}) = \text{reshape}(\mathbf{X}),$$

which possesses the same number of elements, *i.e.*  $\prod_i n_i = \prod_i m_i$ . One example is given by the vec operation from [Definition 2.2](#). It corresponds to a reshape into a tensor of order one. As the arrangement of elements remains unaffected,  $\text{vec } \mathbf{Z} = \text{vec } \mathbf{X}$ , and reshaping corresponds to the identity map on the vectorized input. Consequently, one finds (remember that  $\nabla_{\mathbf{X}}^2 \ell$  is a shorthand notation for  $\nabla_{\text{vec } \mathbf{X}}^2 \ell$ )

$$\nabla_{\mathbf{X}}^2 \ell = \nabla_{\mathbf{Z}}^2 \ell.$$

### A.4.2 Index Select/Map

Selecting elements of a tensor can be phrased as matrix-vector multiplication of a binary matrix  $\Pi$  and the vectorized tensor. The mapping is described by an index map  $\pi$ . Element  $j$  of the output  $\mathbf{z} \in \mathbb{R}^m$  is selected as element  $\pi(j)$  from the input  $\mathbf{x} \in \mathbb{R}^n$ . Only elements  $[\Pi]_{j,\pi(j)}$  in the selection matrix  $\Pi \in \mathbb{R}^{m \times n}$  are one, while all other entries vanish. Consequently, index selection can be expressed as

$$[\mathbf{z}]_j = [\mathbf{x}]_{\pi(j)} \Leftrightarrow \mathbf{z}(\mathbf{x}) = \Pi \mathbf{x} \quad \text{with} \quad [\Pi]_{j,\pi(j)} = 1.$$

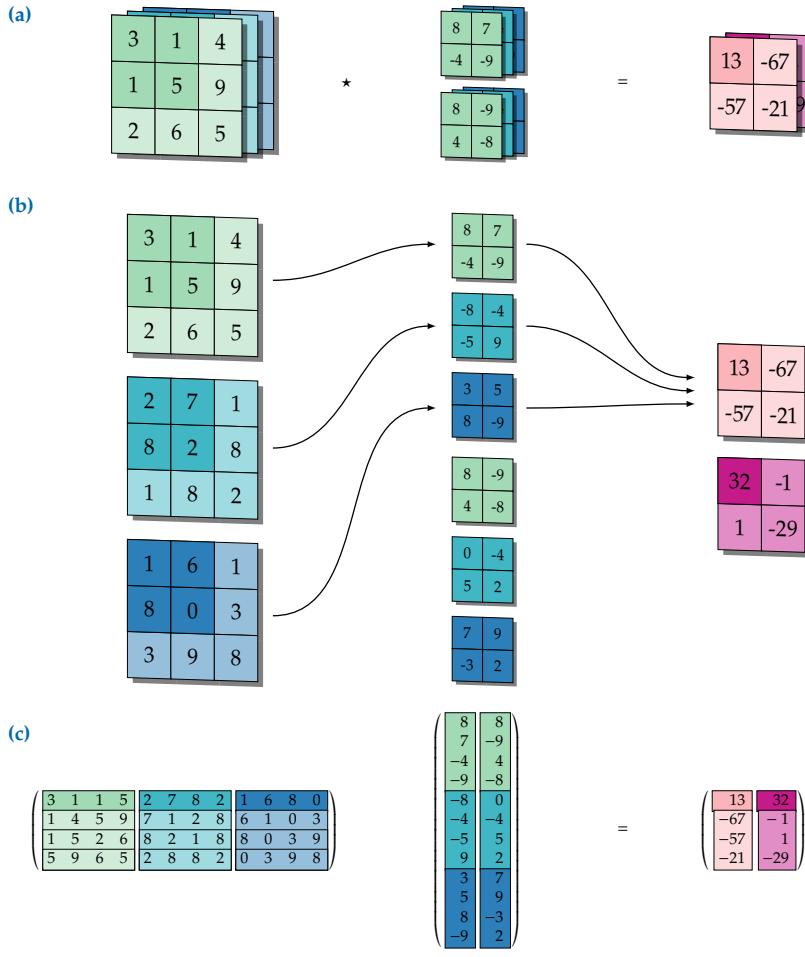
The HBP is equivalent to the linear layer discussed in [Appendix A.2.1](#),

$$\nabla_{\mathbf{z}}^2 \ell = \Pi^T (\nabla_{\mathbf{x}}^2 \ell) \Pi.$$

Applications include max-pooling and the `im2col` / `unfold` operation (see [Appendix A.4.3](#)). Average-pooling represents a weighted sum of index selections and can be treated analogously.

### A.4.3 Convolution

The convolution operation acts on local patches of a multi-channel input of sequences, images, or volumes. In the following, we restrict the discussion to two-dimensional convolution. [Figure A.3a](#) illustrates the



**Figure A.3:** Two-dimensional convolution  $\mathbf{Y} = \mathbf{X} * \mathbf{W}$  without bias term. (a) The input  $\mathbf{X}$  consists of  $C_{\text{in}} = 3$  channels (different colors) of  $(3 \times 3)$  images. Filter maps of size  $(2 \times 2)$  are provided by the kernel  $\mathbf{W}$  for the generation of  $C_{\text{out}} = 2$  output channels. Patch and kernel volumes that are contracted in the first step of the convolution are highlighted. Assuming no padding and a stride of one results in four patches. New features  $\mathbf{Y}$  consist of  $C_{\text{out}} = 2$  channels of  $(2 \times 2)$  images. (b) Detailed view. All tensors are unrolled along the first axis. (c) Convolution as matrix multiplication. From left to right, the matrices  $[\mathbf{X}]^T$ ,  $\mathbf{W}^T$ , and  $\mathbf{Y}^T$  are shown.

setting. A collection of filter maps, the *kernel*  $\mathbf{W}$ , is slid over the spatial coordinates of the input tensor  $\mathbf{X}$ . In each step, the kernel is contracted with the current area of overlap (the *patch*).

Both the sliding process as well as the structure of the patch area can be controlled by hyperparameters of the operation (kernel size, stride, dilation). Moreover, it is common practice to extend the input tensor, for instance by zero-padding [for an introduction to the arithmetics of convolutions, see 45]. The approach presented here is not limited to a certain choice of convolution hyperparameters.

### Forward Pass & Notation

We now introduce the quantities involved in the process along with their dimensions. For a summary, see Table A.1. A forward pass of convolution proceeds as follows (see Figure A.3b for an example):

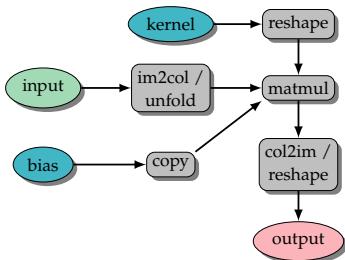
- ▶ The input  $\mathbf{X}$ , a tensor of order three, stores a collection of two-dimensional data. Its components  $\mathbf{X}_{c_{\text{in}}, x_1, x_2}$  are referenced by indices for the channel  $c_{\text{in}}$  and the spatial location  $(x_1, x_2)$ .  $C_{\text{in}}$ ,  $X_1$ ,  $X_2$  denote input channels, width & height of the image, respectively.
- ▶ The kernel  $\mathbf{W}$  is a tensor of order four with shape  $(C_{\text{out}}, C_{\text{in}}, K_1, K_2)$ . Kernel width  $K_1$  and height  $K_2$  determine the patch size  $P = K_1 K_2$

**Table A.1:** Important quantities for the convolution operation. The number of patches equals the product of the outputs' spatial dimensions, *i.e.*  $P = Y_1 Y_2$ .

Tensor	Dimensionality	Description
$\mathbf{X}$	$(C_{in}, X_1, X_2)$	Input
$\mathbf{W}$	$(C_{out}, C_{in}, K_1, K_2)$	Kernel
$\mathbf{Y}$	$(C_{out}, Y_1, Y_2)$	Output
$\llbracket \mathbf{X} \rrbracket$	$(C_{in} K_1 K_2, P)$	Expanded input
$\mathbf{W}$	$(C_{out}, C_{in} K_1 K_2)$	Matricized kernel
$\mathbf{Y}$	$(C_{out}, P)$	Matricized output
$\mathbf{b}$	$C_{out}$	Bias vector
$\mathbf{B}$	$(C_{out}, P)$	Bias matrix

for each channel. New features are obtained by contracting the patch and kernel. This is repeated for a collection of  $C_{out}$  output channels stored along the first axis of  $\mathbf{W}$ .

- ▶ Each output channel  $c_{out}$  is shifted by a bias  $b_{c_{out}}$ , stored in the  $C_{out}$ -dimensional vector  $\mathbf{b}$ .
- ▶ The output  $\mathbf{Y} = \mathbf{X} \star \mathbf{W}$  with components  $[\mathbf{Y}]_{c_{out}, y_1, y_2}$  is of the same structure as the input. We denote the spatial dimensions of  $\mathbf{Y}$  by  $Y_1, Y_2$ , respectively. Hence  $\mathbf{Y}$  is of dimension  $(C_{out}, Y_1, Y_2)$ .



**Figure A.4:** Decomposition of the convolution operation's forward pass.

**Example (index notation):** A special case where input and output have the same spatial dimensions [63] uses a stride of one, kernel widths  $K_1 = K_2 = 2K+1$ , ( $K \in \mathbb{N}$ ), and padding  $K$ . Elements of the filter  $[\mathbf{W}]_{c_{out}, c_{in}, :, :}$  are addressed with the index set  $\{-K, \dots, 0, \dots, K\} \times \{-K, \dots, 0, \dots, K\}$ :

$$[\mathbf{Y}]_{c_{out}, y_1, y_2} = \sum_{k_1=-K}^K \sum_{k_2=-K}^K [\mathbf{X}]_{c_{in}, x_1+k_1, x_2+k_2} [\mathbf{W}]_{c_{out}, c_{in}, k_1, k_2} + [b]_{c_{out}}. \quad (\text{A.8})$$

Elements of  $\mathbf{X}$  addressed out of bounds evaluate to zero. Arbitrary convolutions come with even heavier notation.

### Convolution as Matrix Multiplication

Evaluating convolutions by sums of the form [Equation \(A.8\)](#) leads to poor memory locality [63]. For improved performance, the computation is mapped to a matrix multiplication [29]. To do so, patches of the input  $\mathbf{X}$  are extracted and flattened into columns of a matrix. The patch extraction is indicated by the operator  $\llbracket \cdot \rrbracket$  and the resulting matrix  $\llbracket \mathbf{X} \rrbracket$  is of dimension  $(C_{in} K_1 K_2 \times P)$  (see left part of [Figure A.3c](#) showing  $\llbracket \mathbf{X} \rrbracket^\top$ ). In other words, elements contracted with the kernel are stored along the first axis of  $\llbracket \mathbf{X} \rrbracket$ .  $\llbracket \cdot \rrbracket$  is also referred to as `im2col` or `unfold` operation<sup>2</sup>, and accounts for padding.

The kernel tensor  $\mathbf{W}$  is reshaped into a  $(C_{out} \times C_{in} K_1 K_2)$  matrix  $\mathbf{W}$ , and elements of the bias vector  $vb$  are copied columnwise into a  $(C_{out} \times P)$  matrix  $\mathbf{B} = b\mathbf{1}_P^\top$ , where  $\mathbf{1}_P$  is a  $P$ -dimensional vector of ones. Patchwise contractions are carried out by matrix multiplication and yield a matrix  $\mathbf{Y}$  of shape  $(C_{out}, P)$  with  $P = Y_1 Y_2$ ,

$$\mathbf{Y} = \mathbf{W} \llbracket \mathbf{X} \rrbracket + \mathbf{B} \quad (\text{A.9})$$

([Figure A.3c](#) shows  $\mathbf{W}^\top$ ,  $\llbracket \mathbf{X} \rrbracket^\top$ , and  $\mathbf{Y}$  from left to right). Reshaping  $\mathbf{Y}$  into a  $(C_{out}, Y_1, Y_2)$  tensor, referred to as `col2im`, yields  $\mathbf{Y}$ . [Figure A.4](#) summarizes the outlined decomposition of the forward pass.

2: Our definition of the `unfold` operator slightly differs from [63], where flattened patches are stacked rowwise. This lets us achieve an analogous form to a linear layer. Conversion is achieved by transposition.

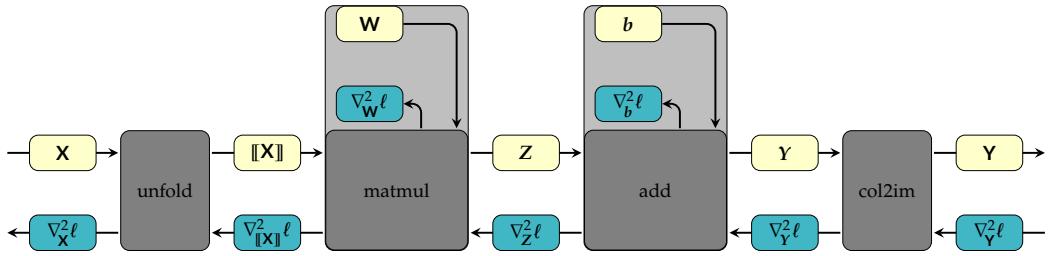


Figure A.5: Decomposition of convolution with notation for the study of curvature backpropagation.

#### A.4.4 HBP for Convolution

We now compose the HBP for convolution, proceeding from right to left with the operations depicted in Figure A.5, by analyzing the backpropagation of curvature for each module, adopting the figure's notation.

Reshape/col2im

The col2im operation takes a matrix  $\mathbf{Y} \in \mathbb{R}^{C_{\text{out}} \times Y_1 Y_2}$  and reshapes it into the tensor  $\mathbf{Y} \in \mathbb{R}^{C_{\text{out}} \times Y_1 \times Y_2}$ . According to Appendix A.4.1,  $\nabla_{\mathbf{Y}}^2 \ell = \nabla_{\mathbf{Y}}^2 \ell$ .

##### Bias Hessian

Forward pass  $\mathbf{Y} = \mathbf{Z} + \mathbf{B}$  and Equation (A.3b) imply  $\nabla_{\mathbf{Y}}^2 \ell = \nabla_{\mathbf{B}}^2 \ell = \nabla_{\mathbf{Z}}^2 \ell$ . To obtain the Hessian w.r.t.  $\mathbf{b}$  from  $\nabla_{\mathbf{B}}^2 \ell$ , consider the columnwise copy operation  $\mathbf{B}(\mathbf{b}) = \mathbf{b} \mathbf{1}_P^\top$ , whose matrix differential is  $d\mathbf{B}(\mathbf{b}) = (\mathbf{d}\mathbf{b}) \mathbf{1}_P^\top$ . Vectorization yields  $d\text{vec } \mathbf{B}(\mathbf{b}) = (\mathbf{1}_P \otimes \mathbf{I}_{C_{\text{out}}}) \mathbf{d}\mathbf{b}$ . Hence, the Jacobian is  $J_{\mathbf{b}} \mathbf{B}(\mathbf{b}) = \mathbf{1}_P \otimes \mathbf{I}_{C_{\text{out}}}$ , and Equation (4.7) yields

$$\nabla_{\mathbf{b}}^2 \ell = (\mathbf{1}_P \otimes \mathbf{I}_{C_{\text{out}}})^\top \nabla_{\mathbf{B}}^2 \ell (\mathbf{1}_P \otimes \mathbf{I}_{C_{\text{out}}}).$$

This performs a linewise and columnwise summation over  $\nabla_{\mathbf{B}}^2 \ell$ , summing entities that correspond to copies of the same entry of  $\mathbf{b}$  in the matrix  $\mathbf{B}$ . It can also be regarded as a reshape of  $\nabla_{\mathbf{B}}^2 \ell$  into a  $(C_{\text{out}}, P, C_{\text{out}}, P)$  tensor, which is then contracted over the second and fourth axis.

##### Weight Hessian

HBP for the matrix-matrix multiplication  $\mathbf{Z}(\mathbf{W}, \|\mathbf{X}\|) = \mathbf{W} \|\mathbf{X}\|$  was discussed in Appendix A.2.1. The Jacobians are given by  $J_{\|\mathbf{X}\|} \mathbf{Z}(\|\mathbf{X}\|) = \mathbf{I}_P \otimes \mathbf{W}$  and  $J_{\mathbf{W}} \mathbf{Z}(\mathbf{W}) = \|\mathbf{X}\|^\top \otimes \mathbf{I}_S$  with the patch size  $S = C_{\text{in}} K_1 K_2$ . Hence, HBP for the weight matrix & unfolded input are

$$\begin{aligned} \nabla_{\|\mathbf{X}\|}^2 \ell &= (\mathbf{I}_P \otimes \mathbf{W})^\top \nabla_{\mathbf{Z}}^2 \ell (\mathbf{I}_P \otimes \mathbf{W}), \\ \nabla_{\mathbf{W}}^2 \ell &= (\|\mathbf{X}\|^\top \otimes \mathbf{I}_S)^\top \nabla_{\mathbf{Z}}^2 \ell (\|\mathbf{X}\|^\top \otimes \mathbf{I}_S). \end{aligned}$$

From what has been said about the reshape operation in Appendix A.4.1, it follows that  $\nabla_{\mathbf{W}}^2 \ell = \nabla_{\mathbf{W}}^2 \ell$ .

### Im2col/Unfold

The patch extraction operation  $\llbracket \cdot \rrbracket$  copies all patch elements into the column of a matrix. It thus represents a selection of elements by an index map which is hard to express in notation. Numerically, it is obtained by calling `im2col` on a  $(C_{\text{in}}, X_1, X_2)$  index tensor whose entries correspond to the indices. The resulting tensor contains all information about the index map. HBP follows the relation of [Appendix A.4.2](#).

### Discussion

Although convolution can be understood as a matrix multiplication, the parameter Hessian is not identical to that of the linear layer discussed in [Appendix A.2.1](#). The difference is due to the parameter sharing of the convolution. In the case of a linear layer  $z = Wx + b$ , the Hessian of the weight matrix for a single sample possesses Kronecker structure [8, 21, 31], *i.e.*  $\nabla_W^2 \ell = x \otimes x^\top \otimes \nabla_z^2 \ell$ . For convolutional layers, however, it has been argued by [8] that block diagonals of the Hessian do not inherit the same structure. Rephrasing the forward pass equation [A.9](#) in terms of vectorized quantities, we find

$$\text{vec } Y = (I_P \otimes W) \text{vec}[\llbracket X \rrbracket] + \text{vec } B.$$

In this perspective, convolution corresponds to a fully-connected linear layer, with the additional constraints that the weight and bias matrix be circular. Defining  $\hat{W} := I_P \otimes W$ , one then finds the Hessian  $\nabla_{\hat{W}}^2 \ell$  to possess Kronecker structure. Parameterization with a kernel tensor encodes the circularity constraint in weight sharing.

For the kernel Hessian  $\nabla_W^2 \ell$  to possess Kronecker structure, the output Hessian  $\nabla_Z^2 \ell$  must be assumed to factorize into a Kronecker product of  $S \times S$  and  $C_{\text{out}} \times C_{\text{out}}$  matrices. These assumptions are somewhat in parallel with the approximations introduced in [63] to obtain KFC.

## A.5 Experimental Details

### Fully-connected Neural Network

The same model as in [31] (see [Table A.2a](#)) is used to extend the experiment performed therein. The weights of each linear layer are initialized with the Xavier method of Glorot and Bengio [57]. Bias terms are initialized to zero. Backpropagation of the Hessian uses approximation [Equation \(4.9\)](#) of [Equation \(A.6\)](#) to compute the curvature blocks  $\overline{\nabla_{W^{(l)}}^2 \ell}$  and  $\overline{\nabla_{b^{(l)}}^2 \ell}$ .

Hyperparameters are chosen as follows to obtain consistent results with the original work: all runs shown in [Figure 4.4](#) use a batch size of  $|\mathcal{B}| = 500$ . For SGD, the learning rate is assigned to  $\eta = 0.1$  with momentum  $\rho = 0.9$ . Block-splitting experiments with the second-order method use the PCH-abs. All runs were performed with a learning rate  $\eta = 0.1$  and a regularization strength of  $\alpha = 0.02$ . For the convergence criterion of CG, the maximum number of iterations is restricted to  $n_{\text{CG}} = 50$ ; convergence is reached at a relative tolerance  $\epsilon_{\text{CG}} = 0.1$ .

**Table A.2: Model architectures under consideration.** We use `Conv2d(in_channels, out_channels, kernel_size, padding)`, `ZeroPad2d(padding_left, padding_right, padding_top, padding_bottom)`, `Linear(in_features, out_features)`, and `MaxPool2d(kernel_size, stride)` as patterns to describe module hyperparameters. Convolution strides are always one. (a) FCNN used to extend the experiment in Chen et al. [31] (3 846 810 parameters). (b) CNN architecture (1 099 226 parameters). (c) DeepOBS 3c3d test problem with three convolutional and three dense layers (895 210 parameters). ReLU activation functions are replaced by sigmoids.

(a) FCNN (Figure 4.4)

#	Module
1	<code>Flatten()</code>
2	<code>Linear(3072, 1024)</code>
3	<code>Sigmoid()</code>
4	<code>Linear(1024, 512)</code>
5	<code>Sigmoid()</code>
6	<code>Linear(512, 256)</code>
7	<code>Sigmoid()</code>
8	<code>Linear(256, 128)</code>
9	<code>Sigmoid()</code>
10	<code>Linear(128, 64)</code>
11	<code>Sigmoid()</code>
12	<code>Linear(64, 32)</code>
13	<code>Sigmoid()</code>
14	<code>Linear(32, 16)</code>
15	<code>Sigmoid()</code>
16	<code>Linear(16, 10)</code>

(b) CNN (Figure 4.5)

#	Module
1	<code>Conv2d(3, 16, 3, 1)</code>
2	<code>Sigmoid()</code>
3	<code>Conv2d(16, 16, 3, 1)</code>
4	<code>Sigmoid()</code>
5	<code>MaxPool2d(2, 2)</code>
6	<code>Conv2d(16, 32, 3, 1)</code>
7	<code>Sigmoid()</code>
8	<code>Conv2d(32, 32, 3, 1)</code>
9	<code>Sigmoid()</code>
10	<code>MaxPool2d(2, 2)</code>
11	<code>Flatten()</code>
12	<code>Linear(2048, 512)</code>
13	<code>Sigmoid()</code>
14	<code>Linear(512, 64)</code>
15	<code>Sigmoid()</code>
16	<code>Linear(64, 10)</code>

(c) DeepOBS 3c3d (Figure A.6)

#	Module
1	<code>Conv2d(3, 64, 5, 0)</code>
2	<code>Sigmoid()</code>
3	<code>ZeroPad2d(0, 1, 0, 1)</code>
4	<code>MaxPool2d(3, 2)</code>
5	<code>Conv2d(64, 96, 3, 0)</code>
6	<code>Sigmoid()</code>
7	<code>ZeroPad2d(0, 1, 0, 1)</code>
8	<code>MaxPool2d(3, 2)</code>
9	<code>ZeroPad2d(1, 1, 1, 1)</code>
10	<code>Conv2d(96, 128, 3, 0)</code>
11	<code>Sigmoid()</code>
12	<code>ZeroPad2d(0, 1, 0, 1)</code>
13	<code>MaxPool2d(3, 2)</code>
14	<code>Flatten()</code>
15	<code>Linear(1152, 512)</code>
16	<code>Sigmoid()</code>
17	<code>Linear(512, 256)</code>
18	<code>Sigmoid()</code>
19	<code>Linear(256, 10)</code>

## Convolutional Neural Net

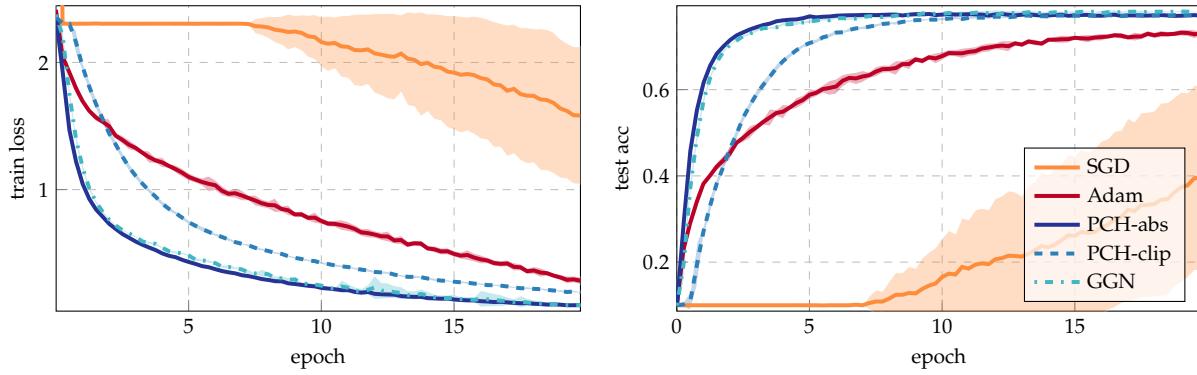
The CNN architecture shown in Table A.2b is trained on a hyperparameter grid. Runs with smallest final training loss are selected to rerun on different random seeds. The curves in Figure 4.5b represent mean values and standard deviations for ten different realizations over the random seed. All layer parameters were initialized with PyTorch’s default.

For the first-order methods (SGD, Adam), we considered batch sizes  $|\mathbb{B}| \in \{100, 200, 500\}$ . For SGD, momentum  $\rho$  was tuned over the set  $\{0, 0.45, 0.9\}$ . Although we varied the learning rate over a large range of values  $\eta \in \{10^{-3}, 10^{-2}, 0.1, 1, 10\}$ , losses kept plateauing and did not decrease. In particular, the loss even increased for the large learning rates. For Adam, we only vary the learning rate  $\eta \in \{10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1, 10\}$ .

As second-order methods scale better to large batch sizes, we used  $|\mathbb{B}| \in \{200, 500, 1000\}$  for them. The CG convergence parameters were set to  $n_{CG} = 200$  and  $\epsilon_{CG} = 0.1$ . For all curvature matrices, we varied the learning rates over  $\eta \in \{0.05, 0.1, 0.2\}$  and  $\alpha \in \{10^{-4}, 10^{-3}, 10^{-2}\}$ .

To compare with another second-order method, we experimented with a public PyTorch implementation of the KFAC optimizer [63, 109] ([github.com/alecwangcq/KFAC-Pytorch](https://github.com/alecwangcq/KFAC-Pytorch)). All hyperparameters were kept at their default setting. The learning rate was varied over  $\eta \in \{10^{-4}, 10^{-3}, 10^{-2}, 0.1, 1, 10\}$ .

The hyperparameters of results shown in Figure 4.5 read as follows:



**Figure A.6:** Comparison of SGD, Adam, and Newton-style methods with different exact curvature matrix-vector products (HBP). The architecture is the DeepOBS 3c3d network [146] with sigmoid activations (Table A.2c).

- ▶ SGD ( $|\mathbb{B}| = 100, \eta = 10^{-3}, \rho = 0.9$ ). The particular choice of these hyperparameters is artificial. This run is representative for SGD, which does not achieve any progress at all.
- ▶ Adam ( $|\mathbb{B}| = 100, \eta = 10^{-3}$ )
- ▶ KFAC ( $|\mathbb{B}| = 500, \eta = 0.1$ )
- ▶ PCH-abs ( $|\mathbb{B}| = 1000, \eta = 0.2, \alpha = 10^{-3}$ ),  
PCH-clip ( $|\mathbb{B}| = 1000, \eta = 0.1, \alpha = 10^{-4}$ )
- ▶ GGN,  $\alpha_1$  ( $|\mathbb{B}| = 1000, \eta = 0.1, \alpha = 10^{-4}$ ). This run does not yield the minimum training loss on the grid; it is shown to illustrate that the second-order method can escape the flat regions in early stages.
- ▶ GGN,  $\alpha_2$  ( $|\mathbb{B}| = 1000, \eta = 0.1, \alpha = 10^{-3}$ ). Compared to  $\alpha_1$ , the second-order method requires more iterations to escape the initial plateau, caused by the larger regularization strength. However, this leads to improved robustness against noise in later training stages.

**Additional experiment:** Another experiment conducted with HBP considers the 3c3d architecture (Table A.2c) of DeepOBS [146] on CIFAR-10. ReLU activations are replaced by sigmoids to make the problem more challenging. The hyperparameter grid is chosen identically to the CNN experiment above, and results are summarized in Figure A.6. In particular, the hyperparameter settings for each competitor are:

- ▶ SGD ( $|\mathbb{B}| = 100, \eta = 1, \rho = 0$ )
- ▶ Adam ( $|\mathbb{B}| = 100, \eta = 10^{-3}$ )
- ▶ PCH-abs ( $|\mathbb{B}| = 500, \eta = 0.1, \alpha = 10^{-3}$ ),  
PCH-clip ( $|\mathbb{B}| = 500, \eta = 0.1, \alpha = 10^{-2}$ )
- ▶ GGN ( $|\mathbb{B}| = 500, \eta = 0.05, \alpha = 10^{-3}$ )

# Additional Material for Chapter 5

# B.

## B.1 BackPACK Extensions

This section provides technical details on BackPACK's quantities.

B.1 BackPACK Extensions . . .	127
B.2 Benchmark Details . . .	134
B.3 Experimental Details . . .	135
B.4 BackPACK Cheat Sheet . .	139

### Notation

Consider an arbitrary module  $f_{\theta^{(l)}}^{(l)}$  of a network  $l = 1, \dots, L$ , parameterized by  $\theta^{(l)}$ . It transforms the output of its parent layer for sample  $n$ ,  $z_n^{(l-1)}$ , to its output  $z_n^{(l)}$ , i.e.

$$z_n^{(l)} = f_{\theta^{(l)}}^{(l)}(z_n^{(l-1)}) , \quad n = 1, \dots, N , \quad (\text{B.1})$$

where  $N$  is the number of samples. In particular,  $z_n^{(0)} = x_n$  and  $z_n^{(L)}(\theta) = f_\theta$ , where  $f_\theta$  is the transformation of the whole network. The dimension of the hidden layer  $l$ 's output  $z_n^{(l)}$  is written  $h^{(l)}$  and  $\theta^{(l)}$  is of dimension  $d^{(l)}$ . The dimension of the network output, the prediction  $z^{(L)}$ , is  $h^{(L)} = C$ . For classification,  $C$  corresponds to the number of classes.

All quantities are assumed to be vector-shaped. For image-processing transformations that usually act on tensor-shaped inputs, we can reduce to the vector scenario by vectorizing all quantities ([Definition 2.2](#)); this discussion does not rely on a specific flattening scheme. However, for an efficient implementation, vectorization should match the layout of the memory of the underlying arrays.

### Jacobian

The Jacobian ([Definition 2.4](#))  $J_a b$  of a vector  $b \in \mathbb{R}^B$  w.r.t. another vector  $a \in \mathbb{R}^A$  is a  $[B \times A]$  matrix of partial derivatives,  $[J_a b]_{i,j} = \partial [b]_i / \partial [a]_j$ .

#### B.1.1 First-order Quantities

The basis for the extraction of additional information about first-order derivatives is [Equation \(5.3\)](#), which we state again for multiple samples,

$$\nabla_{\theta^{(l)}} \mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \nabla_{\theta^{(l)}} \ell_n(\theta) = \frac{1}{N} \sum_{n=1}^N \left( J_{\theta^{(l)}} z_n^{(l)} \right)^\top \left( \nabla_{z_n^{(l)}} \ell_n(\theta) \right) .$$

During the backpropagation step of module  $l$ , we have access to  $\nabla_{z_n^{(l)}} \ell(\theta)$ ,  $n = 1, \dots, N$ . To extract more quantities involving the gradient, we use additional information about the transformation  $f_{\theta^{(l)}}^{(l)}$  in our custom implementation of Jacobian  $J_{\theta^{(l)}} z_n^{(l)}$  and transposed Jacobian  $(J_{\theta^{(l)}} z_n^{(l)})^\top$ .

### Individual Gradients

The contribution of each sample to the overall gradient,  $1/N \nabla_{\theta^{(l)}} \ell_n(\boldsymbol{\theta})$ , is computed by application of the transposed Jacobian,

$$\frac{1}{N} \nabla_{\theta^{(l)}} \ell_n(\boldsymbol{\theta}) = \frac{1}{N} \left( J_{\theta^{(l)}} z_n^{(l)} \right)^\top \left( \nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta}) \right), \quad n = 1, \dots, N. \quad (\text{B.2})$$

For each parameter  $\theta^{(l)}$  the individual gradients are of size  $[N \times d^{(l)}]$ .

### Individual Gradient $L_2$ Norm

The quantity  $\|1/N \nabla_{\theta^{(l)}} \ell_n(\boldsymbol{\theta})\|_2^2$ , for  $n = 1, \dots, N$ , could be extracted from the individual gradients (Equation (B.2)) as

$$\begin{aligned} & \left\| \frac{1}{N} \nabla_{\theta^{(l)}} \ell_n(\boldsymbol{\theta}) \right\|_2^2 \\ &= \left[ \frac{1}{N} \left( J_{\theta^{(l)}} z_n^{(l)} \right)^\top \left( \nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta}) \right) \right]^\top \left[ \frac{1}{N} \left( J_{\theta^{(l)}} z_n^{(l)} \right)^\top \left( \nabla_{z_n^{(l)}} \ell_n(\boldsymbol{\theta}) \right) \right], \end{aligned}$$

which is an  $N$ -dimensional object per parameter  $\theta^{(l)}$ . However, this is memory inefficient as the individual gradients are  $[N \times d^{(l)}]$ . BackPACK circumvents this by using the Jacobian's structure whenever possible.

For a specific example, take a linear layer with parameters  $\boldsymbol{\theta}$  as an  $[A \times B]$  matrix. The layer transforms the inputs  $\{z_n^{(l-1)}\}_{n=1}^N$ , an  $[N \times A]$  matrix which we will now refer to as  $A$ . During the backward pass, it receives the gradient of the individual losses *w.r.t.* its output,  $\{1/N \nabla_{z_n^{(l)}} \ell_n\}_{n=1}^N$ , as an  $[N \times B]$  matrix which we will refer to as  $B$ . The overall gradient, an  $[A \times B]$  matrix, can be computed as  $A^\top B$ , and the individual gradients are a set of  $N$   $[A \times B]$  matrices,  $\{[A]_{n,:}[B]_{n,:}^\top\}_{n=1}^N$ . We want to avoid storing that information. To reduce the memory requirement, note that the individual gradient norm can be written as

$$\left\| \frac{1}{N} \nabla_{\theta} \ell_n \right\|^2 = \sum_i \sum_j ([A]_{n,i} [B]_{n,j})^2,$$

and that the summation can be done independently for each matrix, as  $\sum_i \sum_j ([A]_{n,i} [B]_{n,j})^2 = (\sum_i [A]_{n,i})^2 (\sum_j [B]_{n,j}^2)$ . Therefore, we can square each matrix (element-wise) and sum over non-batch dimensions. This yields vectors  $\mathbf{a}, \mathbf{b}$  of  $N$  elements, where  $[\mathbf{a}]_n = \sum_i [A]_{n,i}^2$ . The individual gradients'  $L_2$  norm is then given by  $\mathbf{a} \odot \mathbf{b}$ .

### Second moment

The gradient second moment (or more specifically, the diagonal of the second moment) is the sum of the squared elements of the individual gradients in a mini-batch, *i.e.*

$$\frac{1}{N} \sum_{n=1}^N [\nabla_{\theta^{(l)}} \ell_n(\boldsymbol{\theta})]_j^2, \quad j = 1, \dots, d^{(l)}. \quad (\text{B.3})$$

It can be used to evaluate the variance of individual elements of the gradient (see below). The second moment is of dimension  $d^{(l)}$ , like the layer parameter  $\theta^{(l)}$ . Similarly to the  $L_2$  norm, it can be computed from individual gradients, but is more efficiently computed implicitly.

Revisiting the linear layer example from the individual  $L_2$  norm computation, the second moment of the parameters  $\theta[i, j]$  is given by  $\sum_n ([A]_{n,i} [B]_{n,j})^2$ , which can be directly computed by taking the element-wise square of  $A$  and  $B$ ,  $A^{\odot 2}$ ,  $B^{\odot 2}$ , and computing  $(A^{\odot 2})^\top B^{\odot 2}$ .

## Variance

Gradient variances over a mini-batch (or more precisely, the covariance diagonal) are computed from the second moment and the gradient,

$$\frac{1}{N} \sum_{n=1}^N [\nabla_{\theta^{(l)}} \ell_n(\theta)]_j^2 - [\nabla_{\theta^{(l)}} \mathcal{L}(\theta)]_j^2, \quad j = 1, \dots, d^{(l)}. \quad (\text{B.4})$$

The element-wise gradient variance is of same dimension as the layer parameter  $\theta^{(l)}$ , *i.e.*  $d^{(l)}$ .

### B.1.2 Second-order Quantities Based on the GGN

#### Backpropagation for the GGN's Block Diagonal

The computation of quantities that originate from approximations of the Hessian require an additional backward pass (see [37]). Most curvature approximations supported by BackPACK rely on the generalized Gauss-Newton (GGN) matrix [148] from Equation (3.15b)

$$G(\theta) = \frac{1}{N} \sum_{n=1}^N (\mathbf{J}_\theta f_n)^\top \nabla_{f_n}^2 \ell(f_n, y_n) (\mathbf{J}_\theta f_n). \quad (\text{B.5})$$

One interpretation of the GGN is that it corresponds to the empirical risk Hessian when the model is replaced with its first-order Taylor expansion, *i.e.* by linearizing the network and ignoring second-order effects. Hence, the effect of module curvature in the recursive scheme of Equation (5.4) can be ignored to obtain the simpler expression

$$\begin{aligned} G(\theta^{(l)}) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{J}_{\theta^{(l)}} f_n)^\top \nabla_{f_n}^2 \ell(f_n, y_n) (\mathbf{J}_{\theta^{(l)}} f_n) \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{J}_{\theta^{(l)}} z_n^{(l)})^\top G(z_n^{(l)}) (\mathbf{J}_{\theta^{(l)}} z_n^{(l)}) \end{aligned} \quad (\text{B.6})$$

for the exact block diagonal of the full GGN. In analogy to  $G(\theta^{(l)})$  we have introduced the  $[d^{(l)} \times d^{(l)}]$ -dimensional quantity

$$G(z_n^{(l)}) = (\mathbf{J}_{z_n^{(l)}} f_n)^\top \nabla_{f_n}^2 \ell(f_n, y_n) (\mathbf{J}_{z_n^{(l)}} f_n)$$

that needs to be backpropagated following [Equation \(5.4\)](#) as

$$\mathbf{G}(\mathbf{z}_n^{(l-1)}) = \left( \mathbf{J}_{\mathbf{z}_n^{(l-1)}} \mathbf{z}_n^{(l)} \right)^\top \mathbf{G}(\mathbf{z}_n^{(l)}) \left( \mathbf{J}_{\mathbf{z}_n^{(l-1)}} \mathbf{z}_n^{(l)} \right), \quad l = 1, \dots, L, \quad (\text{B.7a})$$

initialized with the loss Hessian *w.r.t.* to the network prediction, *i.e.*

$$\mathbf{G}(\mathbf{z}_n^{(L)}) = \nabla_{f_n}^2 \ell(f_n, \mathbf{y}_n). \quad (\text{B.7b})$$

Although this scheme is exact, it is computationally infeasible as it requires backpropagation of  $N [h^{(l)} \times h^{(l)}]$  matrices between layer  $l+1$  and  $l$ . Even for small  $N$ , this is impossible for nets with large convolutions.

As an example, the first layer of the All-CNN-C network outputs  $29 \times 29$  images with 96 channels, which already gives  $h^{(l)} = 80,736$ , which leads to half a Gigabyte per sample. Moreover, storing all the  $[d^{(l)} \times d^{(l)}]$ -dimensional blocks  $\mathbf{G}(\theta^{(l)})$  is not possible. BackPACK implements different approximation strategies, developed by Martens and Grosse [\[109\]](#) and Botev et al. [\[21\]](#) that address both of these complexity issues from different perspectives.

### Symmetric Factorization Scheme

One way to improve the memory footprint of the backpropagated matrices in the case where the model prediction's dimension  $C$  (the number of classes in an image classification task) is small compared to all hidden features  $h^{(l)}$  is to propagate a symmetric factorization of the GGN instead. It relies on the observation that if the loss function itself is convex, even though its composition with the network might not be, its Hessian *w.r.t.* the network output can be decomposed as (*e.g.* [Example 5.1](#))

$$\nabla_{f_n}^2 \ell(f_n, \mathbf{y}_n) = \mathbf{S}(\mathbf{z}_n^{(L)}) \mathbf{S}(\mathbf{z}_n^{(L)})^\top \quad (\text{B.8})$$

with the  $[C \times C]$ -dimensional matrix factorization of the loss Hessian,  $\mathbf{S}(\mathbf{z}_n^{(L)})$ , for sample  $n$ . Consequently, the GGN in [Equation \(B.5\)](#) reduces to an outer product,

$$\mathbf{G}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{n=1}^N \left[ (\mathbf{J}_{\boldsymbol{\theta}} f_n)^\top \mathbf{S}(\mathbf{z}_n^{(L)}) \right] \left[ (\mathbf{J}_{\boldsymbol{\theta}} f_n)^\top \mathbf{S}(\mathbf{z}_n^{(L)}) \right]^\top. \quad (\text{B.9})$$

The analogue for diagonal blocks follows from [Equation \(B.6\)](#) and reads

$$\mathbf{G}(\boldsymbol{\theta}^{(l)}) = \frac{1}{N} \sum_{n=1}^N \left[ \left( \mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{z}_n^{(l)} \right)^\top \mathbf{S}(\mathbf{z}_n^{(l)}) \right] \left[ \left( \mathbf{J}_{\boldsymbol{\theta}^{(l)}} \mathbf{z}_n^{(l)} \right)^\top \mathbf{S}(\mathbf{z}_n^{(l)}) \right]^\top, \quad (\text{B.10})$$

where we defined the  $[h^{(l)} \times C]$ -dimensional matrix square root  $\mathbf{S}(\mathbf{z}_n^{(l)}) := (\mathbf{J}_{\mathbf{z}_n^{(l)}} f_n)^\top \mathbf{S}(\mathbf{z}_n^{(L)})$ . Instead of having layer  $l$  backpropagate  $N$  objects of shape  $[h^{(l)} \times h^{(l)}]$  according to [Equation \(B.7\)](#), we instead backpropagate the matrix square root via

$$\mathbf{S}(\mathbf{z}_n^{(l-1)}) = \left( \mathbf{J}_{\mathbf{z}_n^{(l-1)}} \mathbf{z}_n^{(l)} \right)^\top \mathbf{S}(\mathbf{z}_n^{(l)}), \quad l = 1, \dots, L, \quad (\text{B.11})$$

starting with [Equation \(B.8\)](#). This reduces the backpropagated matrix of layer  $l$  to  $[h^{(l)} \times C]$  for each sample.

### Diagonal Curvature Approximations

**Diagonal of the GGN (DiagGGN):** The factorization trick for the loss Hessian reduces the size of the backpropagated quantities, but does not address the intractable size of the GGN diagonal blocks  $G(\theta^{(l)})$ . In BackPACK, we can extract  $\text{diag}(G(\theta^{(l)}))$  given the backpropagated quantities  $S(z_n^{(l)})$ ,  $l = 1, \dots, N$ , without building up the matrix representation of Equation (B.10). In particular, we compute

$$\begin{aligned} & \text{diag}\left[G(\theta^{(l)})\right] \\ &= \frac{1}{N} \sum_{n=1}^N \text{diag}\left\{\left[\left(J_{\theta^{(l)}} z_n^{(l)}\right)^T S(z_n^{(l)})\right] \left[\left(J_{\theta^{(l)}} z_n^{(l)}\right)^T S(z_n^{(l)})\right]^T\right\}. \end{aligned} \quad (\text{B.12})$$

**Diagonal of the GGN with MC-sampled loss Hessian (DiagGGN-MC):** We use the same backpropagation strategy of Equation (B.11), replacing the symmetric factorization of Equation (B.8) with an approximation by a smaller matrix  $\tilde{S}(z_n^{(L)})$  of size  $[C \times \tilde{C}]$  and  $\tilde{C} < C$  (e.g. Example 5.2),

$$\nabla_{f_n}^2 \ell(f_n, y_n) \approx \tilde{S}(z_n^{(L)}) \left(\tilde{S}(z_n^{(L)})\right)^T. \quad (\text{B.13})$$

This further reduces the size of backpropagated curvature quantities. Martens and Grosse [109] introduced such a sampling scheme with KFAC based on the connection between the GGN and the Fisher. Most loss functions used in machine learning have a probabilistic interpretation as negative log-likelihood of a probabilistic model (see Section 2.1.3). The squared error of regression is equivalent to a Gaussian noise assumption and the cross-entropy is linked to the categorical distribution. In this case, the loss Hessian *w.r.t.* the network output is equal, in expectation, to the outer products of gradients *if the target is sampled according to a particular distribution*,  $q(y | f)$ , defined by the network output. Sampling targets  $\hat{y} \sim q(y | f)$  for a datum  $(x, y)$  with  $f := f_\theta(x)$ , we have

$$\mathbb{E}_{\hat{y} \sim q(\cdot | f)} \left[ (\nabla_f \ell(f, \hat{y})) (\nabla_f \ell(f, \hat{y}))^T \right] = \nabla_f^2 \ell(f, y). \quad (\text{B.14})$$

Sampling one such gradient leads to a rank-1 MC approximation of the loss Hessian. With the substitution  $S \leftrightarrow \tilde{S}$ , we compute an MC approximation of the GGN diagonal in BackPACK as

$$\begin{aligned} & \text{diag}\left[G(\theta^{(l)})\right] \\ & \approx \frac{1}{N} \sum_{n=1}^N \text{diag}\left\{\left[\left(J_{\theta^{(l)}} z_n^{(l)}\right)^T \tilde{S}(z_n^{(l)})\right] \left[\left(J_{\theta^{(l)}} z_n^{(l)}\right)^T \tilde{S}(z_n^{(l)})\right]^T\right\}. \end{aligned} \quad (\text{B.15})$$

### Kronecker-factored Curvature Approximations

A different approach to reduce memory complexity of the GGN blocks  $G(\theta^{(l)})$ , apart from diagonal curvature approximations, is representing them as Kronecker products (KFAC for linear [109] and convolutional layers [63] KFLR and KFRA for linear layers by [21]),

$$G(\theta^{(l)}) \approx A^{(l)} \otimes B^{(l)}. \quad (\text{B.16})$$

1: We keep the PyTorch convention that weights and bias are treated as separate parameters. For the bias terms, we can store the full matrix representation of the GGN. This factor reappears in the Kronecker factorization of the GGN *w.r.t.* the weights.

2: In the case of convolutions, one has to sum over the spatial indices of a single channel of  $z_n^{(l)}$  as the bias is added to an entire channel, see [63] for details.

For both linear and convolutional layers, the first Kronecker factor  $A^{(l)}$  is obtained from the inputs  $z_n^{(l-1)}$  to layer  $l$ . Instead of repeating the technical details of the aforementioned references, we will focus on how they differ in (i) the backpropagated quantities and (ii) the backpropagation strategy. As a result, we will be able to extend KFLR and KFRA to CNNs<sup>1</sup>.

**KFAC & KFLR:** KFAC uses an MC-sampled estimate of the loss Hessian with a square root factorization  $\tilde{S}(z_n^{(L)})$  like in [Equation \(B.13\)](#). The backpropagation is equivalent to the computation of the GGN diagonal. For the weights of a linear layer  $l$ , the second Kronecker term is

$$\mathbf{B}_{\text{KFAC}}^{(l)} = \frac{1}{N} \sum_{n=1}^N \tilde{S}(z_n^{(l)}) \left( \tilde{S}(z_n^{(l)}) \right)^{\top},$$

which at the same time corresponds to the GGN of the layer's bias<sup>2</sup>.

In contrast to KFAC, KFLR backpropagates the exact square root factorization  $S(z_n^{(L)})$ , *i.e.* for the weights of a linear layer<sup>2</sup> (details in [21])

$$\mathbf{B}_{\text{KFLR}}^{(l)} = \frac{1}{N} \sum_{n=1}^N S(z_n^{(l)}) \left( S(z_n^{(l)}) \right)^{\top}.$$

**KFRA:** The backpropagation strategy for KFRA eliminates the scaling of the backpropagated curvature quantities with the batch size  $N$  in [Equation \(B.7\)](#). Instead of layer  $l$  receiving the  $N$  exact  $[h^{(l)} \times h^{(l)}]$  matrices  $G(z_n^{(l)})$ ,  $n = 1, \dots, N$ , only a single averaged object, denoted  $\bar{G}^{(l)}$ , is used as an approximation. In particular, the recursion changes to

$$\bar{G}^{(l-1)} = \frac{1}{N} \sum_{n=1}^N \left( J_{z_n^{(l-1)}} z_n^{(l)} \right)^{\top} \bar{G}^{(l)} \left( J_{z_n^{(l-1)}} z_n^{(l)} \right), \quad l = 1, \dots, L, \quad (\text{B.17a})$$

and is initialized with the batch-averaged loss Hessian

$$\bar{G}^{(L)} = \frac{1}{N} \sum_{n=1}^N \nabla_{f_n}^2 \ell(f_n, y_n). \quad (\text{B.17b})$$

For a linear layer, KFRA uses<sup>2</sup> (see [21] for more details)

$$\mathbf{B}_{\text{KFRA}}^{(l)} = \bar{G}^{(l)}.$$

### B.1.3 The Exact Hessian Diagonal

For neural networks consisting only of piecewise linear activation functions, computing the diagonal of the Hessian is equivalent to computing the GGN diagonal. This is because for these activations the second term in the Hessian backpropagation recursion ([Equation \(5.4\)](#)) vanishes.

However, for activation functions with non-vanishing second derivative, these residual terms have to be accounted for in the backpropagation.

The Hessian backpropagation for module  $l$  reads

$$\nabla_{\theta^{(l)}}^2 \ell(\theta) = \left( J_{\theta^{(l)}} z_n^{(l)} \right)^T \nabla_{z_n^{(l)}}^2 \ell(\theta) \left( J_{\theta^{(l)}} z_n^{(l)} \right) + R_n^{(l)}(\theta^{(l)}), \quad (\text{B.18a})$$

$$\nabla_{z_n^{(l-1)}}^2 \ell(\theta) = \left( J_{z_n^{(l-1)}} z_n^{(l)} \right)^T \nabla_{z_n^{(l)}}^2 \ell(\theta) \left( J_{z_n^{(l-1)}} z_n^{(l)} \right) + R_n^{(l)}(z_n^{(l-1)}), \quad (\text{B.18b})$$

for  $n = 1, \dots, N$ . Those  $[h^{(l)} \times h^{(l)}]$  residuals are

$$\begin{aligned} R_n^{(l)}(\theta^{(l)}) &= \sum_j \left( \nabla_{\theta^{(l)}}^2 [z_n^{(l)}]_j \right) \left[ \nabla_{z_n^{(l)}} \ell(\theta) \right]_j, \\ R_n^{(l)}(z_n^{(l-1)}) &= \sum_j \left( \nabla_{z_n^{(l-1)}}^2 [z_n^{(l)}]_j \right) \left[ \nabla_{z_n^{(l)}} \ell(\theta) \right]_j. \end{aligned}$$

Common parameterized layers (linear and convolution) have  $R_n^{(l)}(\theta^{(l)}) = \mathbf{0}$ . For elementwise activations,  $R_n^{(l)}(z_n^{(l-1)})$  are diagonal matrices.

Storing these quantities becomes very memory-intensive for high-dimensional nonlinear activation layers. In BackPACK, this complexity is reduced by application of the aforementioned matrix square root factorization trick. To do so, we divide the symmetric factorization of  $R_n^{(l)}(z_n^{(l-1)})$  into the positive- and negative-definite terms

$$\begin{aligned} R_n^{(l)}(z_n^{(l-1)}) \\ = P_n^{(l)}(z_n^{(l-1)}) \left( P_n^{(l)}(z_n^{(l-1)}) \right)^T - N_n^{(l)}(z_n^{(l-1)}) \left( N_n^{(l)}(z_n^{(l-1)}) \right)^T. \end{aligned} \quad (\text{B.19})$$

$P_n^{(l)}(z_n^{(l-1)})$ ,  $N_n^{(l)}(z_n^{(l-1)})$  represent the matrix square root of  $R_n^{(l)}(z_n^{(l-1)})$  projected on its positive and negative eigenspace, respectively.

This composition allows for the extension of the GGN backpropagation: in addition to  $S(z_n^{(l)})$ , the residual decompositions  $P_n^{(l)}(z_n^{(l-1)})$ ,  $N_n^{(l)}(z_n^{(l-1)})$  also have to be backpropagated according to [Equation \(B.11\)](#). All diagonals are extracted from the backpropagated matrix square roots (see [Equation \(B.12\)](#)). Diagonals from decompositions in the negative residual eigenspace have to be weighted by a factor of  $-1$  before summation.

In terms of complexity, one backpropagation for  $R_n^{(l)}(z^{(l-1)})$  changes the dimensionality as follows

$$R_n^{(l)}(z^{(l-1)}) : [h^{(l)} \times h^{(l)}] \rightarrow [h^{(l-1)} \times h^{(l-1)}] \rightarrow [h^{(l-2)} \times h^{(l-2)}] \rightarrow \dots$$

With the square root factorization, one instead obtains

$$P_n^{(l)}(z_n^{(l-1)}) : [h^{(l)} \times h^{(l)}] \rightarrow [h^{(l-1)} \times h^{(l)}] \rightarrow [h^{(l-2)} \times h^{(l)}] \rightarrow \dots,$$

$$N_n^{(l)}(z_n^{(l-1)}) : [h^{(l)} \times h^{(l)}] \rightarrow [h^{(l-1)} \times h^{(l)}] \rightarrow [h^{(l-2)} \times h^{(l)}] \rightarrow \dots.$$

Roughly speaking, this is more efficient if the hidden dimension of a nonlinear activation layer exceeds the net's largest hidden dimension.

## Algorithm

Consider one backpropagation step of module  $l$ . Assume  $R_n^{(l)}(\theta^{(l)}) = \mathbf{0}$ , i.e. a linear, convolution, or non-parameterized layer. Then the following

computations are performed in the protocol for the diagonal Hessian:

- Receive the following from the child module  $l+1$  (for  $n = 1, \dots, N$ ):

$$\begin{aligned} \Phi = & \left\{ S(z_n^{(l)}) , \right. \\ & P_n^{(l+1)}(z_n^{(l)}) , \\ & N_n^{(l+1)}(z_n^{(l)}) , \\ & (J_{z_n^{(l)}} z_n^{(l+1)})^\top P_n^{(l+2)}(z_n^{(l+1)}) , \\ & (J_{z_n^{(l)}} z_n^{(l+1)})^\top N_n^{(l+2)}(z_n^{(l+1)}) , \\ & \dots \\ & (J_{z_n^{(l)}} z_n^{(l+1)})^\top (J_{z_n^{(l+1)}} z_n^{(l+2)})^\top \dots (J_{z_n^{(L-3)}} z_n^{(L-2)})^\top P_n^{(L-1)}(z_n^{(L-2)}) , \\ & \left. (J_{z_n^{(l)}} z_n^{(l+1)})^\top (J_{z_n^{(l+1)}} z_n^{(l+2)})^\top \dots (J_{z_n^{(L-3)}} z_n^{(L-2)})^\top N_n^{(L-1)}(z_n^{(L-2)}) \right\} \end{aligned}$$

- Extract the parameter Hessian diagonal  $\text{diag}(\nabla_{\theta^{(l)}}^2 \mathcal{L}(\theta))$
- For each quantity  $A \in \Phi$  extract the diagonal from the square root factorization and sum over the samples, *i.e.* compute

$$\frac{1}{N} \sum_{n=1}^N \text{diag} \left\{ \left[ \left( J_{\theta^{(l)}} z_n^{(l)} \right)^\top A_n \right] \left[ \left( J_{\theta^{(l)}} z_n^{(l)} \right)^\top A_n \right]^\top \right\} .$$

Multiply the expression by  $-1$  if  $A$  stems from backpropagation of a residual's negative eigenspace's factorization.

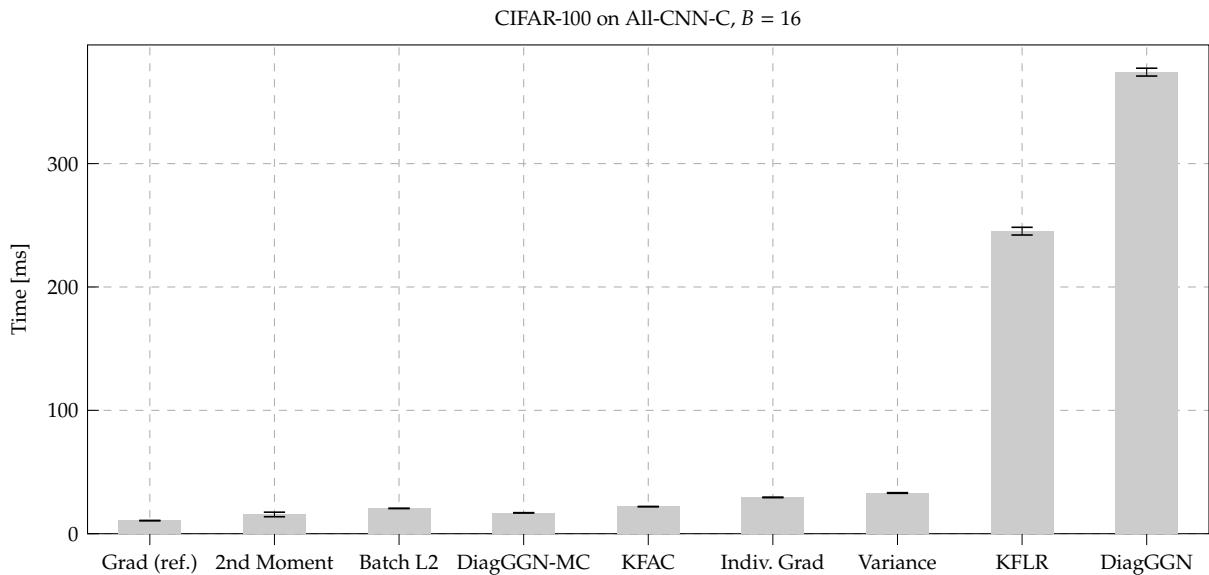
- Sum all expressions to obtain the block Hessian's diagonal  $\text{diag}(\nabla_{\theta^{(l)}}^2 \mathcal{L}(\theta))$

- Backpropagate the received quantities to the parent module  $l-1$
- For each quantity  $A_n \in \Phi$ , apply  $(J_{z_n^{(l-1)}} z_n^{(l)})^\top A_n$
- Append  $P_n^{(l+1)}(z_n^{(l)})$  and  $N_n^{(l+1)}(z_n^{(l)})$  to  $\Phi$

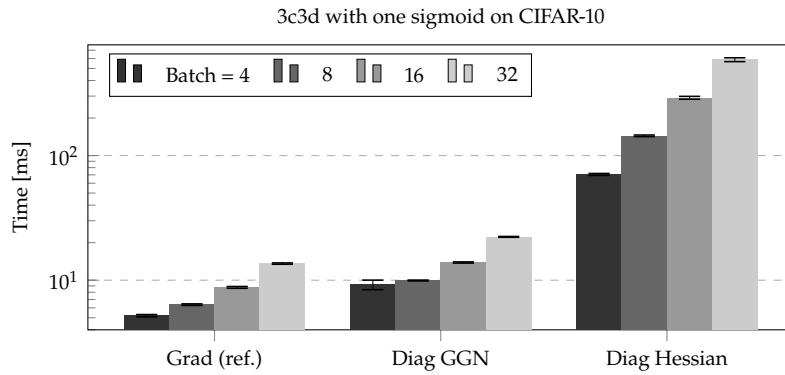
## B.2 Benchmark Details

**KFAC versus KFLR:** As the KFLR of Botev et al. [21] is orders of magnitude more expensive to compute than the KFAC of Martens and Grosse [109] on CIFAR-100, it was not included in the main plot. This is not an implementation error; it follows from the definition of those methods. To approximate the GGN,  $G(\theta) = 1/N \sum_n [J_\theta f_n]^\top \nabla_{f_n}^2 \ell_n [J_\theta f_n]$ , KFAC uses a rank-1 approximation for each of the inner Hessian  $\nabla_{f_n}^2 \ell_n \approx s_n s_n^\top$ , and needs to propagate a *vector* through the computation graph for each sample. KFLR uses the complete inner Hessian instead. For CIFAR-100, the network has 100 output nodes—one for each class—and the inner Hessians are  $[100 \times 100]$  matrices. KFLR needs to propagate a *matrix* through the computation graph for each sample, which is 100× more expensive as shown in Figure B.1.

**GGN diagonal versus Hessian diagonal:** Most nets used in deep learning use ReLU activations. ReLU functions have no *curvature* as they are piecewise linear. Because of this, the diagonal of the GGN is



**Figure B.1:** KFLR and DiagGGN are more expensive to run on large networks. The gradient takes less than 20 ms to compute, but KFLR and DiagGGN are approximately 100× more expensive.



**Figure B.2:** Diagonal of the Hessian versus the GGN. If the network contains a single sigmoid activation function, the diagonal of the Hessian is an order of magnitude more computationally intensive than the diagonal of the GGN.

equivalent to the diagonal of the Hessian [107]. However, for networks that use non piecewise linear activation functions like sigmoids or tanh, computing the Hessian diagonal can be much more expensive than the GGN diagonal. To illustrate this, we modify the smaller net used in our benchmarks to include a single sigmoid activation function before the last classification layer. The results in Figure B.2 show that the Hessian diagonal computation is already an order of magnitude more expensive than for the GGN.

## B.3 Experimental Details

### B.3.1 Protocol

The optimizer experiments are performed according to the protocol suggested by DeepOBS:

- ▶ Train the neural network with the investigated optimizer and vary its hyperparameters on a specified grid. This training is performed for a single random seed only.

- ▶ DeepOBS evaluates metrics during the training procedure. From all runs of the grid search, it selects the best run automatically. The results shown in this work were obtained with the default strategy, favoring highest final accuracy on the validation set.
- ▶ For a better understanding of the optimizer performance with respect to randomized routines in the training process, DeepOBS reruns the best hyperparameter setting for ten different random seeds. The results show mean values over these repeated runs, with standard deviations as uncertainty indicators.
- ▶ Along with the benchmarked optimizers, we show the DeepOBS base line performances for Adam and momentum SGD (Momentum). They are provided by DeepOBS.

The optimizers built upon BackPACK's curvature estimates were benchmarked on DeepOBS's image classification problems from [Table B.1](#).

**Table B.1:** Image classification test problems considered from the DeepOBS library [146].

Codename	Description	Dataset	# Parameters
LogReg	Linear model	MNIST	7,850
2c2d	2 convolutional and 2 dense linear layers	Fashion-MNIST	3,274,634
3c3d	3 convolutional and 3 dense linear layers	CIFAR-10	895,210
All-CNN-C	9 convolutional layers [157]	CIFAR-100	1,387,108

### B.3.2 Grid Search & Best Hyperparameter Setting

Both the learning rate  $\eta$  and damping  $\lambda$  are tuned over the grid

$$\eta \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\} , \quad \lambda \in \{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10\} .$$

We use the same batch size ( $N = 128$  for all problems, except  $N = 256$  for All-CNN-C on CIFAR-100) as the base lines and the optimizers run for the identical number of epochs.

The best hyperparameter settings are summarized in [Table B.2](#).

### B.3.3 Update rule

We use a simple update rule with a constant damping parameter  $\lambda$ . Consider the parameters  $\theta$  of a single module in a neural network with

**Table B.2:** Best hyperparameter settings for optimizers and baselines shown in this work. In the Momentum baselines, the momentum was fixed to 0.9. Parameters for computation of the running averages in Adam use the default values  $(\beta_1, \beta_2) = (0.9, 0.999)$ . The symbols ✓ and ✗ denote whether the hyperparameter setting is an interior point of the grid or not, respectively.

Curvature	mnist_logreg			fmnist_2c2d			cifar10_3c3d			cifar100_allcnn		
	$\eta$	$\lambda$	int									
DiagGGN	$10^{-3}$	$10^{-3}$	✓	$10^{-4}$	$10^{-4}$	✗	$10^{-3}$	$10^{-2}$	✓	-	-	-
DiagGGN-MC	$10^{-3}$	$10^{-3}$	✓	$10^{-4}$	$10^{-4}$	✗	$10^{-3}$	$10^{-2}$	✓	$10^{-3}$	$10^{-3}$	✓
KFAC	$10^{-2}$	$10^{-2}$	✓	$10^{-3}$	$10^{-3}$	✓	1	10	✗	1	1	✓
KFLR	$10^{-2}$	$10^{-2}$	✓	$10^{-2}$	$10^{-3}$	✓	1	10	✗	-	-	-
KFRA	$10^{-2}$	$10^{-2}$	✓	-	-	-	-	-	-	-	-	-
Baseline	$\eta$			$\eta$			$\eta$			$\eta$		
Momentum	$\approx 2.07 \cdot 10^{-2}$			$\approx 2.07 \cdot 10^{-2}$			$\approx 3.79 \cdot 10^{-3}$			$\approx 4.83 \cdot 10^{-1}$		
Adam	$\approx 2.98 \cdot 10^{-4}$			$\approx 1.27 \cdot 10^{-4}$			$\approx 2.98 \cdot 10^{-4}$			$\approx 6.95 \cdot 10^{-4}$		

$L_2$ -regularization of strength  $\delta$ . Let  $C(\theta_t)$  denote the curvature matrix and  $\nabla_{\theta_t} \mathcal{L}(\theta_t)$  the gradient at step  $t$ . One iteration applies

$$\theta_{t+1} \leftarrow \theta_t + [C(\theta_t) + (\lambda + \delta)I]^{-1} [\nabla_{\theta_t} \mathcal{L}(\theta_t) + \delta \theta_t]. \quad (\text{B.20})$$

The inverse cannot be computed exactly (in reasonable time) for the Kronecker-factored curvatures KFAC, KFLR, and KFRA. We use the scheme of [109] to approximately invert  $C(\theta_t) + (\lambda + \delta)I$  if  $C(\theta_t)$  is Kronecker-factored;  $C(\theta_t) = A(\theta_t) \otimes B(\theta_t)$ . It replaces the expression  $(\lambda + \delta)I$  by diagonal terms added to each Kronecker factor. In summary, this replaces

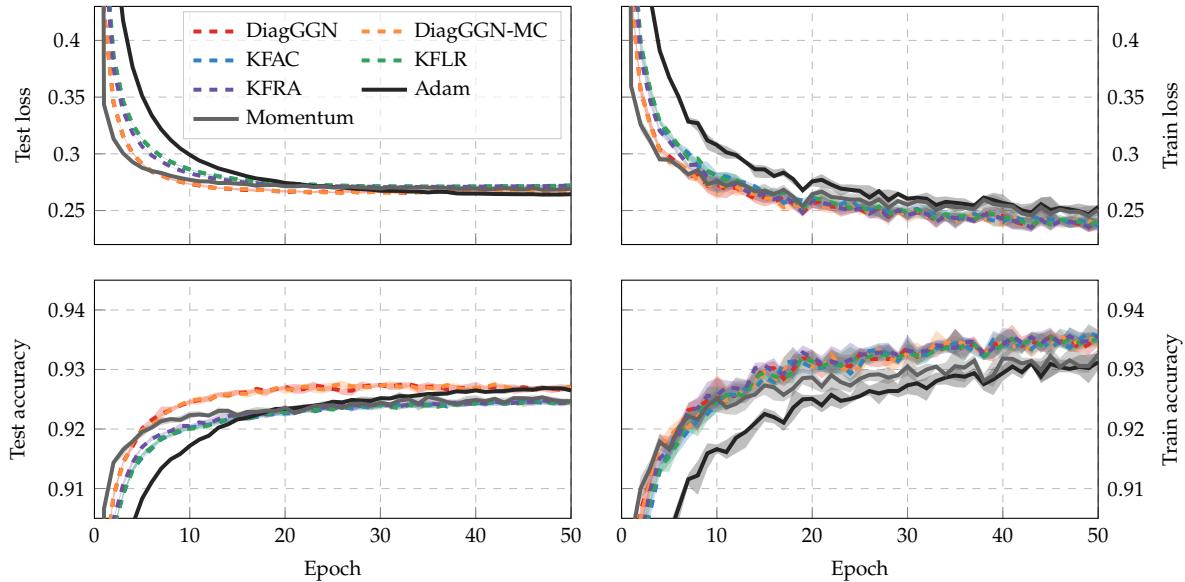
$$\begin{aligned} & [A(\theta_t) \otimes B(\theta_t) + (\lambda + \delta)I]^{-1} \\ \text{by } & \left[ A(\theta_t) + \pi_t \sqrt{\lambda + \delta} I \right]^{-1} \otimes \left[ B(\theta_t) + \frac{1}{\pi_t} \sqrt{\lambda + \delta} I \right]^{-1}. \end{aligned} \quad (\text{B.21})$$

A principled choice for the parameter  $\pi_t$  is  $\pi_t = \sqrt{\|A(\theta_t) \otimes I_B\| / \|I_A \otimes B(\theta_t)\|}$  for any matrix norm  $\|\cdot\|$ . We follow [109] and choose the trace norm,

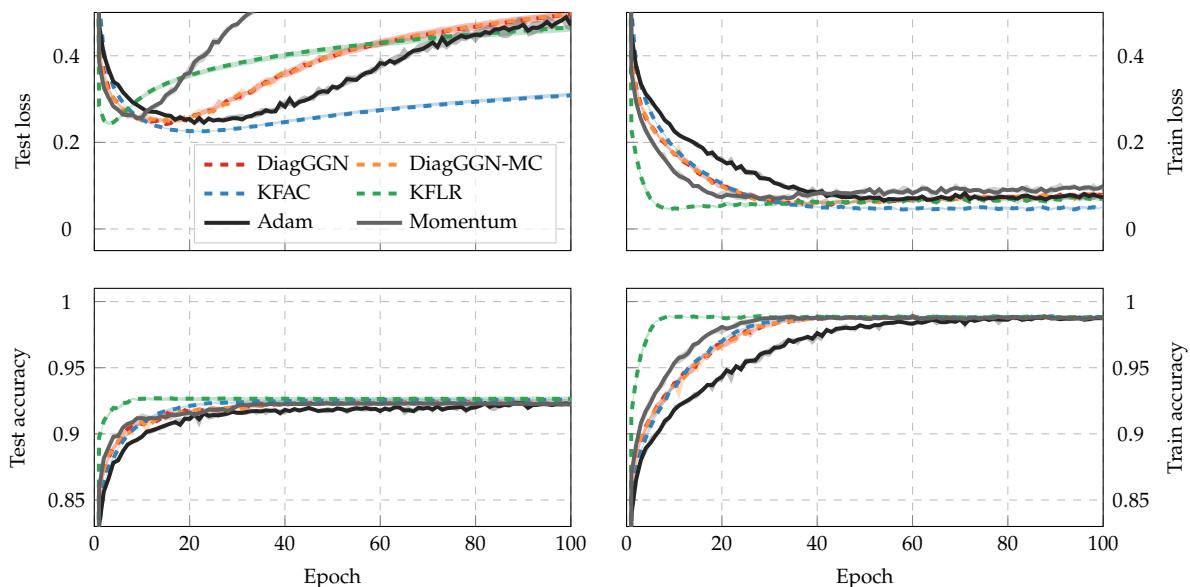
$$\pi_t = \sqrt{\frac{\text{Tr}(A(\theta_t)) \dim(B)}{\dim(A) \otimes \text{Tr}(B(\theta_t))}}. \quad (\text{B.22})$$

### B.3.4 Additional results

This section presents the results for MNIST using a logistic regression in [Figure B.3a](#) and Fashion-MNIST using the 2c2d network, composed of two convolutional and two linear layers, in [Figure B.3b](#).



(a) LogReg (7,850 parameters) on MNIST



(b) 2c2d (3,274,634 parameters) on Fashion-MNIST

**Figure B.3: Additional results.** Median performance with shaded quartiles of the best hyperparameter settings chosen by DeepOBS. Solid lines show well-tuned baselines of momentum SGD and Adam that are provided by DeepOBS.

## B.4 BackPACK Cheat Sheet

### ► Assumptions

- Sequential feedforward network

$$z_n^{(0)} \xrightarrow{f_{\theta^{(1)}}^{(1)}(z_n^{(0)})} z_n^{(1)} \xrightarrow{f_{\theta^{(2)}}^{(2)}(z_n^{(1)})} \dots \xrightarrow{f_{\theta^{(L)}}^{(L)}(z_n^{(L-1)})} z_n^{(L)} \xrightarrow{\ell(z_n^{(L)}, y)} \ell(\theta)$$

- $d^{(l)}$  : Dimension of parameter  $\theta^{(l)}$

- Empirical risk

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{n=1}^N \ell(f_\theta(x_n), y_n)$$

### ► Shorthands

$$\begin{aligned} \ell_n(\theta) &:= \ell(f_\theta(x_n), y_n), \quad n = 1, \dots, N, \\ f_n &:= f_\theta(x_n) = z_n^{(L)}(\theta), \quad n = 1, \dots, N \end{aligned}$$

### ► Generalized Gauss-Newton matrix

$$G(\theta) = \frac{1}{N} \sum_{n=1}^N (\mathbf{J}_\theta f_n)^\top \nabla_{f_n}^2 \ell_n(\theta) (\mathbf{J}_\theta f_n)$$

### ► Approximative GGN via MC sampling

$$\tilde{G}(\theta) = \frac{1}{N} \sum_{n=1}^N (\mathbf{J}_\theta f_n)^\top \left[ (\nabla_{f_n} \ell(f_n, \hat{y})) (\nabla_{f_n} \ell(f_n, \hat{y}_n))^\top \right]_{\hat{y}_n \sim q(\cdot | f_n)} (\mathbf{J}_\theta f_n)$$

**Table B.3:** Overview of the features supported in the first release of BackPACK. The quantities are computed separately for all module parameters, *i.e.*  $l = 1, \dots, L$ .

Feature	Details
Individual gradients	$\frac{1}{N} \nabla_{\theta^{(l)}} \ell_n(\theta), \quad n = 1, \dots, N$
Batch variance	$\frac{1}{N} \sum_{n=1}^N [\nabla_{\theta^{(l)}} \ell_n(\theta)]^2 - [\nabla_{\theta^{(l)}} \mathcal{L}(\theta)]_j^2, \quad j = 1, \dots, d^{(l)}$
2 <sup>nd</sup> moment	$\frac{1}{N} \sum_{n=1}^N [\nabla_{\theta^{(l)}} \ell_n(\theta)]_j^2, \quad j = 1, \dots, d^{(l)}$ .
Individual gradient $L_2$ norms	$\left\  \frac{1}{N} \nabla_{\theta^{(l)}} \ell_n(\theta) \right\ _2^2, \quad n = 1, \dots, N$
DiagGGN	$\text{diag}(\mathbf{G}(\theta^{(l)}))$
DiagGGN-MC	$\text{diag}(\tilde{\mathbf{G}}(\theta^{(l)}))$
Hessian diagonal	$\text{diag}(\nabla_{\theta^{(l)}}^2 \mathcal{L}(\theta))$
KFAC	$\tilde{\mathbf{G}}(\theta^{(l)}) \approx \mathbf{A}^{(l)} \otimes \mathbf{B}_{\text{KFAC}}^{(l)}$
KFLR	$\mathbf{G}(\theta^{(l)}) \approx \mathbf{A}^{(l)} \otimes \mathbf{B}_{\text{KFLR}}^{(l)}$
KFRA	$\mathbf{G}(\theta^{(l)}) \approx \mathbf{A}^{(l)} \otimes \mathbf{B}_{\text{KFRA}}^{(l)}$



# Additional Material for Chapter 6

## C.1 Code Example

One design principle of Cockpit is its easy integration with conventional PyTorch training loops. [Procedure C.1](#) shows a working example of a standard training loop with Cockpit. More examples and tutorials are described in the documentation. Cockpit's syntax is inspired by BackPACK: it can be used interchangeably with the library responsible for most back-end computations. Changes to the code are straightforward:

```

1  """Example: Training Loop using Cockpit."""
2
3  import torch
4  from _utils_examples import cnn, fmnist_data, get_logpath
5  from backpack import extend
6  from cockpit import Cockpit, CockpitPlotter
7  from cockpit.utils.configuration import configuration as config
8
9  fmnist_data = fmnist_data()
10 model = extend(cnn())
11 loss_fn = extend(torch.nn.CrossEntropyLoss(reduction="mean"))
12 losses_fn = torch.nn.CrossEntropyLoss(reduction="none")
13 opt = torch.optim.SGD(model.parameters(), lr=1e-2)
14
15 cockpit = Cockpit(model.parameters(), quantities=config("full"))
16 plotter = CockpitPlotter()
17
18 max_steps, global_step = 50, 0
19 for inputs, labels in iter(fmnist_data):
20     opt.zero_grad()
21
22     outputs = model(inputs)
23     loss = loss_fn(outputs, labels)
24     losses = losses_fn(outputs, labels)
25
26     with cockpit(
27         global_step,
28         info={
29             "batch_size": inputs.shape[0],
30             "individual_losses": losses,
31             "loss": loss,
32             "optimizer": opt,
33         },
34     ):
35         loss.backward(
36             create_graph=cockpit.create_graph(global_step),
37         )
38
39     opt.step()
40     global_step += 1
41
42     if global_step >= max_steps:
43         break
44
45 cockpit.write(get_logpath())
46 plotter.plot(get_logpath())

```

C.1 Code Example . . . . .	141
C.2 Instrument Overview . .	142
C.3 Mathematical Details .	143
C.4 Additional Experiments .	155
C.5 Implementation Details & Additional Benchmarks .	157
C.6 Cockpit View of Convex Stochastic Problems . . .	166

**Procedure C.1:** Complete training loop with Cockpit in PyTorch. Line changes are highlighted in blue (█).

- **Importing** (Lines 5, 6 and 7): besides importing Cockpit we also

need to import BackPACK which is required for extending (parts of) the model (see next step).

- ▶ **Extending** (*Lines 10 and 11*): when defining the model and the loss function, we need to *extend* both of them using BackPACK. This is as trivial as wrapping them in the `extend()` function provided by BackPACK and lets BackPACK know that additional quantities (such as the individual gradients) should be computed for them. Note, that while applying BackPACK is easy, it currently does not support all possible model architectures and layer types. Specifically, *batch norm* layers are not supported since using them results in ill-defined individual gradients.
- ▶ **Individual losses** (*Line 12*): for the `Alpha` quantity, Cockpit also requires the individual loss values (to estimate the variance of the loss estimate). This can be computed cheaply but is not usually part of a conventional training loop. Creating this loss is done analogously to creating any other loss, with the only exception of setting `reduction="none"`. Since we don't differentiate this loss, we don't need to extend it.
- ▶ **Cockpit configuration** (*Line 15 and 16*): Initializing the Cockpit requires passing them (extended) model parameters as well as a list of quantities that should be tracked. [Table 6.1](#) provides an overview of all possible quantities. In this example, we use one of the pre-defined configurations offered by Cockpit. Separately, we initialize the plotting part of Cockpit. We deliberately detached the visualization from the tracking to allow greater flexibility.
- ▶ **Quantity computation** (*Line 26 to 37*): Performing the training is very similar to a regular training loop, with the only difference being that the backward pass should be surrounded by the Cockpit context (`with cockpit():`). Additionally to the `global_step` we also pass a few additional information to the Cockpit that are computed anyway and can be re-used, such as the batch size, the individual losses, or the optimizer itself. After the backward pass (when the context is left) all Cockpit quantities are automatically computed.
- ▶ **Logging & visualizing** (*Line 45 and 46*): at any point during training—here at the end—we can write all quantities to a log file. We can use this log file, or alternatively the Cockpit directly, to visualize all quantities in a status screen similar to [Figure 6.2](#).

## C.2 Instrument Overview

[Table C.1](#) lists all quantities available in the first public release of Cockpit. If necessary, we provide references to their mathematical definition. This table contains additional quantities, compared to [Table 6.1](#) in the main text. To improve the presentation of this work, we decided to not describe every quantity available in Cockpit in the main part and instead focus on the investigated metrics. Custom quantities can be added easily without having to understand the inner-workings.

**Table C.1:** Overview of all Cockpit quantities with a short description and, if necessary, a reference to mathematical definition.

Name	Description	Math
Loss	Mini-batch training loss at current iteration, $\mathcal{L}_{\mathbb{B}}(\theta)$	(6.1)
Parameters	Parameter values $\theta_t$ at the current iteration	-
Distance	$L_2$ distance from initialization $\ \theta_t - \theta_0\ _2$	-
UpdateSize	Update size of the current iteration $\ \theta_{t+1} - \theta_t\ _2$	-
GradNorm	Mini-batch gradient norm $\ g_{\mathbb{B}}(\theta)\ _2$	-
Time	Time of the current iteration (e.g. used in benchmark of Appendix C.5)	-
Alpha	Normalized step on a noisy quadratic interpolation between two iterates $\theta_t, \theta_{t+1}$	(C.7)
CABS	Adaptive batch size for SGD, optimizes expected objective gain per cost, from [11]	(C.9)
EarlyStopping	Evidence-based early stopping criterion for SGD, proposed in [104]	(C.11d)
GradHist1d	Histogram of individual gradient elements, $\{[g_n(\theta)]_j\}_{n \in \mathbb{B}}^{j=1, \dots, D}$	(C.12)
GradHist2d	Histogram of weights and individual gradient elements, $\{([\theta]_j, [g_n(\theta)]_j)\}_{n \in \mathbb{B}}^{j=1, \dots, D}$	(C.13)
NormTest	Normalized fluctuations of the residual norms $\ g_{\mathbb{B}} - g_n\ $ , proposed in [26]	(C.16c)
InnerTest	Normalized fluctuations of $g_n$ 's parallel components along $g_{\mathbb{B}}$ , proposed in [19]	(C.19c)
OrthoTest	Normalized fluctuations of $g_n$ 's orthogonal components along $g_{\mathbb{B}}$ , proposed in [19]	(C.22b)
HessMaxEV	Maximum Hessian eigenvalue, $\lambda_{\max}(H_{\mathbb{B}}(\theta))$ , inspired by [177]	(C.23)
HessTrace	Exact or approximate Hessian trace, $\text{Tr}(H_{\mathbb{B}}(\theta))$ , inspired by [177]	-
TICDiag	Relation between (diagonal) curvature and gradient noise, inspired by [162]	(C.26)
TICTrace	Relation between curvature and gradient noise trace, inspired by [162]	(C.25)
MeanGSNR	Average gradient signal-to-noise ratio (GSNR), inspired by [100]	(C.28b)

## C.3 Mathematical Details

Here, we provide the mathematical background for each instrument in Table C.1. This complements the more informal description presented in Section 6.2, which focused more on the expressiveness of the individual quantities. We start by setting up the necessary notation in addition to the one introduced in Section 6.2. See Sections 2.1 and 3.2 for more details.

### C.3.1 Additional Notation

#### Population Properties

The population risk  $\mathcal{L}_{p_{\text{data}}}(\theta) \in \mathbb{R}$  and its variance  $\Lambda(\theta) \in \mathbb{R}$  are given by

$$\begin{aligned} \mathcal{L}_{p_{\text{data}}}(\theta) &= \mathbb{E}_{(x,y) \sim p_{\text{data}}} [\ell(f_{\theta}(x), y)] \\ &= \int \ell(f_{\theta}(x), y) p_{\text{data}}(x, y) dx dy, \end{aligned} \quad (\text{C.1a})$$

$$\begin{aligned} \Lambda_{p_{\text{data}}}(\theta) &= \text{Var}_{(x,y) \sim p_{\text{data}}} [\ell(f_{\theta}(x), y)] \\ &= \int (\ell(f_{\theta}(x), y) - \mathcal{L}_{p_{\text{data}}}(\theta))^2 p_{\text{data}}(x, y) dx dy, \end{aligned} \quad (\text{C.1b})$$

and its gradient  $\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta}) \in \mathbb{R}^D$  and variance  $\boldsymbol{\Sigma}_{p_{\text{data}}}(\boldsymbol{\theta}) \in \mathbb{R}^{D \times D}$  are

$$\begin{aligned}\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta}) &= \mathbb{E}_{(x,y) \sim p_{\text{data}}} [\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(x), y)] \\ &= \int \nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}), \mathbf{y}) p_{\text{data}}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y},\end{aligned}\tag{C.2a}$$

$$\begin{aligned}\boldsymbol{\Sigma}_{p_{\text{data}}}(\boldsymbol{\theta}) &= \text{Var}_{(x,y) \sim p_{\text{data}}} [\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(x), y)] \\ &= \int (\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}), \mathbf{y}) - \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})) \\ &\quad (\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{y}), \mathbf{y}) - \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta}))^\top p_{\text{data}}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}.\end{aligned}\tag{C.2b}$$

### Empirical Approximations

Let  $\mathbb{D}$  denote a set of samples drawn *i.i.d.* from  $p_{\text{data}}$ , *i.e.*  $\mathbb{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^{|\mathbb{D}|}$ . With a slight abuse of notation ( $n \in \mathbb{D}$  instead of  $(\mathbf{x}_n, \mathbf{y}_n) \in \mathbb{D}$ ) the empirical risk approximated with  $\mathbb{D}$  is

$$\mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}) = \frac{1}{|\mathbb{D}|} \sum_{n \in \mathbb{D}} \ell_n(\boldsymbol{\theta})\tag{C.3a}$$

(later,  $\mathbb{D}$  will represent either a mini-batch  $\mathcal{B}$ , or the train set  $\mathbb{D}_{\text{train}}$ ). The empirical risk gradient  $\mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta}) \in \mathbb{R}^D$  on  $\mathbb{D}$  is

$$\mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}) = \frac{1}{|\mathbb{D}|} \sum_{n \in \mathbb{D}} \nabla_{\boldsymbol{\theta}} \ell_n(\boldsymbol{\theta}) = \frac{1}{|\mathbb{D}|} \sum_{n \in \mathbb{D}} \mathbf{g}_n(\boldsymbol{\theta}),\tag{C.3b}$$

with individual gradients  $\mathbf{g}_n(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \ell_n(\boldsymbol{\theta}) \in \mathbb{R}^D$  from a sample  $n$ . Population risk and gradient variances  $\Lambda_{p_{\text{data}}}(\boldsymbol{\theta})$ ,  $\boldsymbol{\Sigma}_{p_{\text{data}}}(\boldsymbol{\theta})$  can be empirically estimated on  $\mathbb{D}$  with the sample variances  $\hat{\Lambda}_{\mathbb{D}}(\boldsymbol{\theta}) \in \mathbb{R}$ ,  $\hat{\boldsymbol{\Sigma}}_{\mathbb{D}}(\boldsymbol{\theta}) \in \mathbb{R}^{D \times D}$ ,

$$\Lambda_{p_{\text{data}}}(\boldsymbol{\theta}) \approx \frac{1}{|S|-1} \sum_{n \in \mathbb{D}} (\ell_n(\boldsymbol{\theta}) - \mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta}))^2 := \hat{\Lambda}_{\mathbb{D}}(\boldsymbol{\theta}),\tag{C.4a}$$

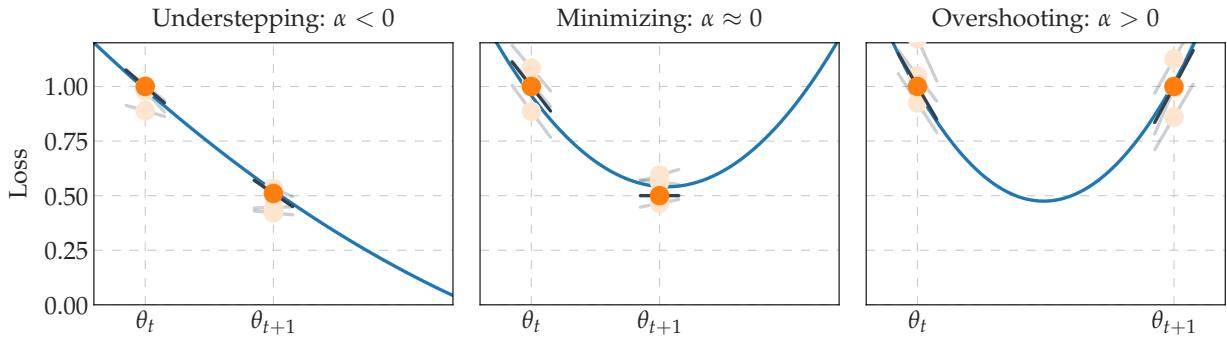
$$\begin{aligned}\boldsymbol{\Sigma}_{p_{\text{data}}}(\boldsymbol{\theta}) &\approx \frac{1}{|\mathbb{D}|-1} \sum_{n \in \mathbb{D}} (\mathbf{g}_n(\boldsymbol{\theta}) - \mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta})) (\mathbf{g}_n(\boldsymbol{\theta}) - \mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta}))^\top := \hat{\boldsymbol{\Sigma}}_{\mathbb{D}}(\boldsymbol{\theta}) \\ &\approx \frac{1}{|\mathbb{D}|-1} \left[ \left( \sum_{n \in \mathbb{D}} \mathbf{g}_n(\boldsymbol{\theta}) \mathbf{g}_n(\boldsymbol{\theta})^\top \right) - |\mathbb{D}| \mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta}) \mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta})^\top \right].\end{aligned}\tag{C.4b}$$

Often, gradient elements are assumed independent and hence their variance is diagonal ( $\odot^2$  denotes element-wise square),

$$\text{diag}[\boldsymbol{\Sigma}_{\mathbb{D}}(\boldsymbol{\theta})] \approx \frac{1}{|S|-1} \sum_{n \in \mathbb{D}} (\mathbf{g}_n(\boldsymbol{\theta}) - \mathbf{g}_{\mathbb{D}}(\boldsymbol{\theta}))^{\odot 2} = \text{diag}[\hat{\boldsymbol{\Sigma}}_{\mathbb{D}}(\boldsymbol{\theta})] \in \mathbb{R}^D.\tag{C.5}$$

### Slicing

To avoid confusion between  $\boldsymbol{\theta}_t$  (parameter at iteration  $t$ ) and  $\theta_j$  ( $j$ -th parameter entry), we denote the latter as  $[\boldsymbol{\theta}]_j$ .



**Figure C.1: Motivational sketch for the  $\alpha$  quantity.** In each iteration of the optimizer we observe the loss function at two positions  $\theta_t$  and  $\theta_{t+1}$  (shown in ●). The black lines (—) show the observed slope at this position, which we can get from projecting the gradients onto the current step direction  $\theta_{t+1} - \theta_t$ . Note, that all four observations (two loss and two slope values) are noisy, due to being computed on a mini-batch. With access to the individual losses and gradients (some samples shown in ●/—), we can estimate their noise level and build a noise-informed quadratic fit (—). Using this fit, we determine whether the optimizer minimizes the local uni-variate loss (middle plot), or whether we understep (left plot) or overshoot (right plot) the minimum.

### C.3.2 Normalized Step Length ([Alpha](#))

#### Motivation

The goal of the  $\alpha$ -quantity is to estimate and quantify the effect that a selected learning rate has on the optimizer's steps. Consider the optimizer's step at training iteration  $t$ . This parameter update from  $\theta_t$  to  $\theta_{t+1}$  happens in a one-dimensional space, defined by the update direction  $\theta_{t+1} - \theta_t = s_t$ . The update direction depends on the update rule of the optimizer, e.g. for SGD with learning rate  $\eta$  it is simply  $s_t = -\eta g_{\mathbb{B}_t}(\theta_t)$ .

We build a noise-informed uni-variate quadratic approximation along this update step ( $\theta_t \rightarrow \theta_{t+1}$ ) based on the two noisy loss function observations at  $\theta_t$  and  $\theta_{t+1}$  and the two noisy slope observation at these two points. Examining this quadratic fit, we are able to determine where on this parabola our optimizer steps. Standardizing this, we express a step to the minimum of the loss in the update direction as  $\alpha = 0$ . Analogously, steps that end short of this minimum result in  $\alpha < 0$ , and a step over the minimum in  $\alpha > 0$ . These three different scenarios are illustrated in Figure C.1 also showing the underlying observations that would lead to them. Figure 6.1 shows the distribution of  $\alpha$ -values for two very different optimization trajectories.

#### Noisy Observations

In order to build an approximation for the loss function in the update direction, we leverage the four observations of the function (and its derivative) that are available in each iteration. Due to the stochasticity of deep learning optimization, we also take into account the noise-level of all observations by estimating them. The first two observations are the mini-batch training losses  $\mathcal{L}_{\mathbb{B}_t}(\theta_t)$ ,  $\mathcal{L}_{\mathbb{B}_{t+1}}(\theta_{t+1})$  at point  $\theta_t$  and  $\theta_{t+1}$ , which are computed in every standard training loop. The mini-batch

losses are averages over individual losses,

$$\begin{aligned}\mathcal{L}_{\mathbb{B}_t}(\boldsymbol{\theta}_t) &= \mathbb{E}_{\mathbb{B}_t} [\ell(\boldsymbol{\theta}_t)] = \frac{1}{|\mathbb{B}_t|} \sum_{n \in \mathbb{B}_t} \ell_n(\boldsymbol{\theta}_t), \\ \mathcal{L}_{\mathbb{B}_{t+1}}(\boldsymbol{\theta}_{t+1}) &= \mathbb{E}_{\mathbb{B}_{t+1}} [\ell(\boldsymbol{\theta}_{t+1})] = \frac{1}{|\mathbb{B}_{t+1}|} \sum_{n \in \mathbb{B}_{t+1}} \ell_n(\boldsymbol{\theta}_{t+1}),\end{aligned}$$

and using these individual losses, we can also compute the variances to estimate the noise-level of our loss observation,

$$\begin{aligned}\text{Var}_{\mathbb{B}_t} [\ell(\boldsymbol{\theta}_t)] &= \left( \frac{1}{|\mathbb{B}_t|} \sum_{n \in \mathbb{B}_t} \ell_n(\boldsymbol{\theta}_t)^2 \right) - \left( \frac{1}{|\mathbb{B}_t|} \sum_{n \in \mathbb{B}_t} \ell_n(\boldsymbol{\theta}_t) \right)^2, \\ \text{Var}_{\mathbb{B}_{t+1}} [\ell(\boldsymbol{\theta}_{t+1})] &= \left( \frac{1}{|\mathbb{B}_{t+1}|} \sum_{n \in \mathbb{B}_{t+1}} \ell_n(\boldsymbol{\theta}_{t+1})^2 \right) - \left( \frac{1}{|\mathbb{B}_{t+1}|} \sum_{n \in \mathbb{B}_{t+1}} \ell_n(\boldsymbol{\theta}_{t+1}) \right)^2.\end{aligned}$$

Similarly, we proceed with the slope in the update direction. To compute the slope of the loss function in the direction of the optimizer's update  $\mathbf{s}_t$ , we project the current gradient along this update direction

$$\begin{aligned}\mathbb{E}_{\mathbb{B}_t} \left[ \frac{\mathbf{s}_t^\top \mathbf{g}(\boldsymbol{\theta}_t)}{\|\mathbf{s}_t\|^2} \right] &= \frac{1}{|\mathbb{B}_t|} \sum_{n \in \mathbb{B}_t} \frac{\mathbf{s}_t^\top \mathbf{g}_n(\boldsymbol{\theta}_t)}{\|\mathbf{s}_t\|^2}, \\ \mathbb{E}_{\mathbb{B}_{t+1}} \left[ \frac{\mathbf{s}_t^\top \mathbf{g}(\boldsymbol{\theta}_{t+1})}{\|\mathbf{s}_t\|^2} \right] &= \frac{1}{|\mathbb{B}_{t+1}|} \sum_{n \in \mathbb{B}_{t+1}} \frac{\mathbf{s}_t^\top \mathbf{g}_n(\boldsymbol{\theta}_{t+1})}{\|\mathbf{s}_t\|^2}.\end{aligned}$$

Just like before, we can also compute the variance of this slope, by leveraging individual gradients,

$$\begin{aligned}\text{Var}_{\mathbb{B}_t} \left[ \frac{\mathbf{s}_t^\top \mathbf{g}(\boldsymbol{\theta}_t)}{\|\mathbf{s}_t\|^2} \right] &= \frac{1}{|\mathbb{B}_t|} \sum_{n \in \mathbb{B}_t} \left( \frac{\mathbf{s}_t^\top \mathbf{g}_n(\boldsymbol{\theta}_t)}{\|\mathbf{s}_t\|^2} \right)^2 - \left( \frac{1}{|\mathbb{B}_t|} \sum_{n \in \mathbb{B}_t} \frac{\mathbf{s}_t^\top \mathbf{g}_n(\boldsymbol{\theta}_t)}{\|\mathbf{s}_t\|^2} \right)^2, \\ \text{Var}_{\mathbb{B}_{t+1}} \left[ \frac{\mathbf{s}_t^\top \mathbf{g}(\boldsymbol{\theta}_{t+1})}{\|\mathbf{s}_t\|^2} \right] &= \frac{1}{|\mathbb{B}_{t+1}|} \sum_{n \in \mathbb{B}_{t+1}} \left( \frac{\mathbf{s}_t^\top \mathbf{g}_n(\boldsymbol{\theta}_{t+1})}{\|\mathbf{s}_t\|^2} \right)^2 - \left( \frac{1}{|\mathbb{B}_{t+1}|} \sum_{n \in \mathbb{B}_{t+1}} \frac{\mathbf{s}_t^\top \mathbf{g}_n(\boldsymbol{\theta}_{t+1})}{\|\mathbf{s}_t\|^2} \right)^2.\end{aligned}$$

## Quadratic Fit & Normalization

Using our (noisy) observations, we are now ready to build an approximation for the loss as a function of the step size, which we will denote as  $f(\tau)$ . We assume a quadratic function for  $f$ , which follows recent reports for the loss landscape of neural networks [176], *i.e.* a function  $f(\tau) = w_0 + w_1 \tau + w_2 \tau^2$  parameterized by  $\mathbf{w} \in \mathbb{R}^3$ . We further assume a Gaussian likelihood of the form

$$p(\tilde{f} | \mathbf{w}, \Phi) = \mathcal{N}(\tilde{f} | \Phi^\top \mathbf{w}, \Lambda) \tag{C.6}$$

for observations  $\tilde{f}$  of the loss and its slope. The observation matrix  $\Phi$  and the noise matrix of the observations  $\Lambda$  are

$$\Phi = \begin{pmatrix} 1 & 1 & 0 & 0 \\ \tau_1 & \tau_2 & 1 & 1 \\ \tau_1^2 & \tau_2^2 & 2\tau_1 & 2\tau_2 \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \sigma_{\tilde{f}_1} & 0 & 0 & 0 \\ 0 & \sigma_{\tilde{f}_2} & 0 & 0 \\ 0 & 0 & \sigma_{\tilde{f}'_1} & 0 \\ 0 & 0 & 0 & \sigma_{\tilde{f}'_2} \end{pmatrix},$$

where  $\tau$  denotes the position and  $\sigma$  denotes the noise-level estimate of the observation. The maximum likelihood solution of [Equation \(C.6\)](#) for the parameters of our quadratic fit is given by

$$w = (\Phi \Lambda^{-1} \Phi^\top)^{-1} \Phi \Lambda^{-1} \tilde{f}. \quad (\text{C.7})$$

Once we have the quadratic fit of the uni-variate loss along the update direction, we normalize the scales such that the resulting  $\alpha$  expresses the effective step taken by the optimizer sketched in [Figure C.1](#).

## Usage

The  $\alpha$ -quantity is related to recent line search approaches [105, 168]. However, instead of searching for an acceptable step by repeated attempts, we instead report the effect of the current step size selection. This could, for example, be used to disentangle the two optimization runs in [Figure 6.1](#). Additionally, this information could also be used to automatically adapt the learning rate during the training process. But, as discussed in [Section 6.3.3](#), it isn't trivial what the "correct" decision is, as it might depend on the optimization problem, the training phase, and other factors. Having this  $\alpha$ -quantity can, however, provide more insight into what kind of steps are used in well-tuned runs with traditional optimizers such as SGD.

### C.3.3 CABS Criterion: Coupling Adaptive Batch Sizes with Learning Rates ([CABS](#))

The CABS criterion, proposed by Balles et al. [11], can be used to adapt the mini-batch size during training with SGD. It relies on the gradient noise and approximately optimizes the objective's expected gain per cost. The adaptation rule is (with learning rate  $\eta$ )

$$|\mathbb{B}| \leftarrow \eta \frac{\text{Tr}(\Sigma_{p_{\text{data}}}(\theta))}{\mathcal{L}_{p_{\text{data}}}(\theta)}, \quad (\text{C.8})$$

and the practical implementation approximates  $\mathcal{L}_{p_{\text{data}}}(\theta) \approx \mathcal{L}_{\mathbb{B}}(\theta)$  and  $\text{Tr}(\Sigma_{p_{\text{data}}}(\theta)) \approx (|\mathbb{B}| - 1)/|\mathbb{B}| \text{Tr}(\hat{\Sigma}_{\mathbb{B}}(\theta))$  (compare equations (10, 22) and first paragraph of Section 4 in [11]). This yields the quantity computed in Cockpit's [CABS](#) instrument,

$$|\mathbb{B}| \leftarrow \eta \frac{\frac{1}{|\mathbb{B}|} \sum_{j=1}^D \sum_{n \in \mathbb{B}} [\mathbf{g}_n(\theta) - \mathbf{g}_{\mathbb{B}}(\theta)]_j^2}{\mathcal{L}_{\mathbb{B}}(\theta)}. \quad (\text{C.9})$$

### Usage

The CABS criterion suggests a batch size which is optimal under certain assumptions. This suggestion can support practitioners in the batch size selection for their deep learning task.

#### C.3.4 Early-stopping Criterion for SGD (`EarlyStopping`)

The empirical risk  $\mathcal{L}_{\mathbb{D}}(\boldsymbol{\theta})$ , and the mini-batch loss  $\mathcal{L}_{\mathbb{B}}(\boldsymbol{\theta})$  are only estimators of the target objective  $\mathcal{L}_{p_{\text{data}}}(\boldsymbol{\theta})$ . Mahsereci et al. [104] motivate  $p(g_{\mathbb{B}, \mathbb{D}}(\boldsymbol{\theta}) \mid g_{p_{\text{data}}}(\boldsymbol{\theta}) = 0)$  as a measure for detecting noise in the finite datasets  $\mathbb{B}, \mathbb{D}$  due to sampling from  $p_{\text{data}}$ . They propose an evidence-based (EB) criterion for early stopping the training procedure based on mini-batch statistics, and model  $p(g_{\mathbb{B}}(\boldsymbol{\theta}))$  with a sampled diagonal variance approximation (compare [Equation \(C.5\)](#)),

$$p(g_{\mathbb{B}}(\boldsymbol{\theta})) \approx \prod_{j=1}^D \mathcal{N}\left(g_{\mathbb{B}}(\boldsymbol{\theta}) \mid \left[g_{p_{\text{data}}}(\boldsymbol{\theta})\right]_j, \frac{[\hat{\Sigma}_{\mathbb{B}}(\boldsymbol{\theta})]_{j,j}}{|\mathbb{B}|}\right). \quad (\text{C.10})$$

Their SGD stopping criterion is

$$\frac{2}{D} [\log p(g_{\mathbb{B}}(\boldsymbol{\theta})) - \mathbb{E}_{g_{\mathbb{B}}(\boldsymbol{\theta}) \sim p(g_{\mathbb{B}}(\boldsymbol{\theta}))} [\log p(g_{\mathbb{B}}(\boldsymbol{\theta}))]] > 0, \quad (\text{C.11a})$$

and translates into

$$1 - \frac{|\mathbb{B}|}{D} \sum_{d=1}^D \frac{[g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2}{[\hat{\Sigma}_{\mathbb{B}}(\boldsymbol{\theta})]_{d,d}} > 0, \quad (\text{C.11b})$$

$$1 - \frac{|\mathbb{B}|}{D} \sum_{d=1}^D \frac{[g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2}{\frac{1}{|\mathbb{B}|-1} \sum_{n \in \mathbb{B}} [g_n(\boldsymbol{\theta}) - g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2} > 0, \quad (\text{C.11c})$$

$$1 - \frac{|\mathbb{B}|(|\mathbb{B}| - 1)}{D} \sum_{d=1}^D \frac{[g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2}{\left(\sum_{n \in \mathbb{B}} [g_n(\boldsymbol{\theta})]_d^2\right) - |\mathbb{B}| [g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2} > 0. \quad (\text{C.11d})$$

Cockpit's `EarlyStopping` quantity computes the left side of [Equation \(C.11d\)](#).

### Usage

Cockpit's `EarlyStopping` quantity can inform practitioners that training is about to be completed and the model might be at risk of overfitting.

#### C.3.5 Individual Gradient Element Histograms (`GradHist1d`, `GradHist2d`)

For the  $|\mathbb{B}| \times D$  individual gradient elements, Cockpit's `GradHist1d` instrument displays a histogram of

$$\{[g_n(\boldsymbol{\theta})]_d\}_{n \in \mathbb{B}, d=1, \dots, D}. \quad (\text{C.12})$$

Cockpit's `GradHist2d` instrument displays a two-dimensional histogram of the  $|\mathbb{B}| \times D$  tuples

$$\{([\theta]_d, [g_n(\theta)]_d)\}_{n \in \mathbb{B}, d=1,\dots,D} \quad (\text{C.13})$$

and the marginalized one-dimensional histograms over the parameter and gradient axes.

## Usage

Sections 6.3.1 and 6.3.2 provide use cases (identifying data pre-processing issues and vanishing gradients) for both the gradient histogram as well as its layer-wise extension.

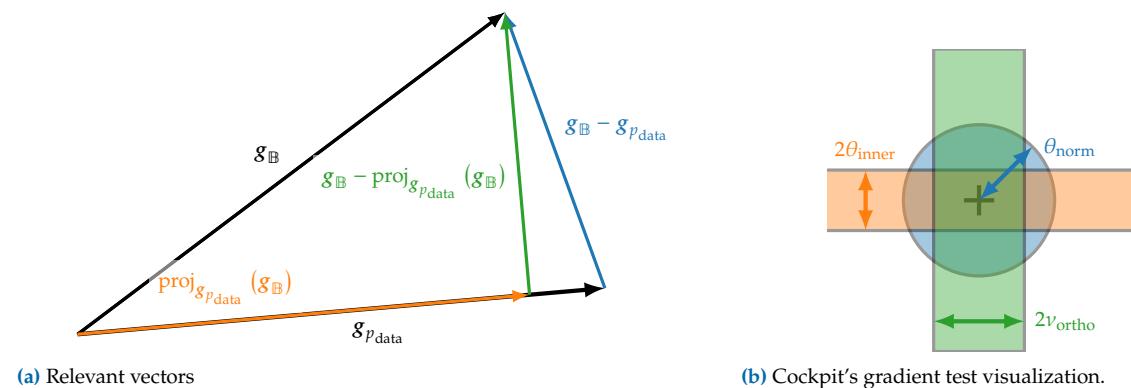
### C.3.6 Gradient Tests (`NormTest`, `InnerTest`, `OrthoTest`)

Bollapragada et al. [19] and Byrd et al. [26] propose batch size adaptation schemes based on the gradient noise. They formulate geometric constraints between population and mini-batch gradient and accessible approximations that can be probed to decide whether the mini-batch size should be increased. Because mini-batches are *i.i.d.* from  $p_{\text{data}}$ , it holds that

$$\mathbb{E}[g_{\mathbb{B}}(\theta)] = g_{p_{\text{data}}}(\theta), \quad (\text{C.14a})$$

$$\mathbb{E}[g_{\mathbb{B}}(\theta)^T g_{p_{\text{data}}}(\theta)] = \|g_{p_{\text{data}}}(\theta)\|^2. \quad (\text{C.14b})$$

The above works propose enforcing other weaker similarity in expectation during optimization. These geometric constraints reduce to basic vector geometry (see Figure C.2a for an overview of the relevant vectors). We recall their formulation here for consistency and derive the practical versions, which can be computed from training observables and are used in Cockpit (consult Figure C.2b for the visualization).



**Figure C.2:** Conceptual sketch for gradient tests. (a) Relevant vectors to formulate the geometric constraints between population and mini-batch gradient probed by the gradient tests. (b) Gradient test visualization in Cockpit.

## Usage

All three gradient tests describe the noise level of the gradients. Bolapragada et al. [19] and Byrd et al. [26] adapt the batch size so that the proposed geometric constraints are fulfilled. Practitioners can use the combined gradient test plot, *i.e.* top center plot in [Figure 6.2](#), to monitor gradient noise during training and adjust hyperparameters such as the batch size.

### Norm Test ([NormTest](#))

The norm test [26] constrains the residual norm  $\|g_{\mathbb{B}}(\theta) - g_{p_{\text{data}}}(\theta)\|$ , rescaled by  $\|g_{p_{\text{data}}}(\theta)\|$ . This gives rise to a standardized ball of radius  $\theta_{\text{norm}} \in (0, \infty)$  around the population gradient, where the mini-batch gradient should reside. Byrd et al. [26] set  $\theta_{\text{norm}} = 0.9$  in their experiments and increase the batch size if (in the practical version, see below) the following constraint is not fulfilled

$$\mathbb{E} \left[ \frac{\|g_{\mathbb{B}}(\theta) - g_{p_{\text{data}}}(\theta)\|^2}{\|g_{p_{\text{data}}}(\theta)\|^2} \right] \leq \theta_{\text{norm}}^2. \quad (\text{C.15a})$$

Instead of taking the expectation over mini-batches, Byrd et al. [26] note that the above will be satisfied if

$$\frac{1}{|\mathbb{B}|} \mathbb{E} \left[ \frac{\|g_n(\theta) - g_{p_{\text{data}}}(\theta)\|^2}{\|g_{p_{\text{data}}}(\theta)\|^2} \right] \leq \theta_{\text{norm}}^2. \quad (\text{C.15b})$$

They propose a practical form of this test,

$$\frac{1}{|\mathbb{B}|(|\mathbb{B}| - 1)} \frac{\sum_{n \in \mathbb{B}} \|g_n(\theta) - g_{\mathbb{B}}(\theta)\|^2}{\|g_{\mathbb{B}}(\theta)\|^2} \leq \theta_{\text{norm}}^2, \quad (\text{C.16a})$$

which can be computed from mini-batch statistics. Rearranging

$$\sum_{n \in \mathbb{B}} \|g_n(\theta) - g_{\mathbb{B}}(\theta)\|^2 = \left( \sum_{n \in \mathbb{B}} \|g_n(\theta)\|^2 \right) - |\mathbb{B}| \|g_{\mathbb{B}}(\theta)\|^2, \quad (\text{C.16b})$$

we arrive at

$$\frac{1}{|\mathbb{B}|(|\mathbb{B}| - 1)} \left[ \frac{\sum_{n \in \mathbb{B}} \|g_n(\theta)\|^2}{\|g_{\mathbb{B}}(\theta)\|^2} - |\mathbb{B}| \right] \leq \theta_{\text{norm}}^2 \quad (\text{C.16c})$$

that leverages the norm of both the mini-batch and the individual gradients, which can be aggregated over parameters during a backward pass. Cockpit's [NormTest](#) corresponds to the maximum radius  $\theta_{\text{norm}}$  for which the above inequality holds.

### Inner Product Test (`InnerTest`)

The inner product test [19] constrains the projection of  $\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})$  onto  $\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})$  (compare Figure C.2a),

$$\text{proj}_{\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}(\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})) := \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta}), \quad (\text{C.17})$$

rescaled by  $\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|$ . This restricts the mini-batch gradient to reside in a standardized band of relative width  $\theta_{\text{inner}} \in (0, \infty)$  around the population risk gradient. Bollapragada et al. [19] use  $\theta_{\text{inner}} = 0.9$  (in the practical version, see below) to adapt the batch size if the parallel component's variance does not satisfy the condition

$$\text{Var}\left(\frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2}\right) = \mathbb{E}\left[\left(\frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} - 1\right)^2\right] \leq \theta_{\text{inner}}^2 \quad (\text{C.18a})$$

(note that by Equation (C.14) we have  $\mathbb{E}[\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})/\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2] = 1$ ). Bollapragada et al. [19] bound Equation (C.18a) by the individual gradient variance,

$$\begin{aligned} & \frac{1}{|\mathbb{B}|} \text{Var}\left(\frac{\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2}\right) \\ &= \frac{1}{|\mathbb{B}|} \mathbb{E}\left[\left(\frac{\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} - 1\right)^2\right] \leq \theta_{\text{inner}}^2. \end{aligned} \quad (\text{C.18b})$$

They then propose a practical form of Equation (C.18b), which uses the mini-batch sample variance,

$$\begin{aligned} & \frac{1}{|\mathbb{B}|} \text{Var}\left(\frac{\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^2}\right) \\ &= \frac{1}{|\mathbb{B}|(|\mathbb{B}| - 1)} \left[ \sum_{n \in \mathbb{B}} \left( \frac{\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^2} - 1 \right)^2 \right] \leq \theta_{\text{inner}}^2. \end{aligned} \quad (\text{C.19a})$$

Expanding

$$\sum_{n \in \mathbb{B}} \left( \frac{\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^2} - 1 \right)^2 = \frac{\sum_{n \in \mathbb{B}} (\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}))^2}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^4} - |\mathbb{B}| \quad (\text{C.19b})$$

and inserting Equation (C.19b) into Equation (C.19a) yields

$$\frac{1}{|\mathbb{B}|(|\mathbb{B}| - 1)} \left[ \frac{\sum_{n \in \mathbb{B}} (\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}))^2}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^4} - |\mathbb{B}| \right] \leq \theta_{\text{inner}}^2. \quad (\text{C.19c})$$

It relies on pairwise scalar products between individual gradients, which can be aggregated over layers during backpropagation. Cock-

Cockpit's `InnerTest` quantity computes the maximum band width  $\theta_{\text{inner}}$  that satisfies [Equation \(C.19c\)](#).

### Orthogonality Test (`OrthoTest`)

In contrast to the inner product test ([Appendix C.3.6](#)) which constrains the projection ([Equation \(C.17\)](#)), the orthogonality test [19] constrains the orthogonal part (see [Figure C.2](#) (a))

$$\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}) - \text{proj}_{\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}(\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})) , \quad (\text{C.20})$$

rescaled by  $\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|$ . This restricts the mini-batch gradient to a standardized band of relative width  $\nu_{\text{ortho}} \in (0, \infty)$  parallel to the population gradient. Bollapragada et al. [19] use  $\nu = \tan(80^\circ) \approx 5.84$  (in the practical version, see below) to adapt the batch size if the following condition is violated,

$$\mathbb{E} \left[ \left\| \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}) - \text{proj}_{\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}(\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}))}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|} \right\|^2 \right] \leq \nu_{\text{ortho}}^2 . \quad (\text{C.21a})$$

Expanding the norm, and inserting [Equation \(C.17\)](#), this simplifies to

$$\begin{aligned} \mathbb{E} \left[ \left\| \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|} - \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} \frac{\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|} \right\|^2 \right] &\leq \nu_{\text{ortho}}^2 , \\ \mathbb{E} \left[ \frac{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^2}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} - \frac{(\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta}))^2}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^4} \right] &\leq \nu_{\text{ortho}}^2 . \end{aligned} \quad (\text{C.21b})$$

Bollapragada et al. [19] bound this inequality with individual gradients,

$$\frac{1}{|\mathbb{B}|} \mathbb{E} \left[ \left\| \frac{\mathbf{g}_n(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} - \frac{\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|^2} \frac{\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})}{\|\mathbf{g}_{p_{\text{data}}}(\boldsymbol{\theta})\|} \right\|^2 \right] \leq \nu_{\text{ortho}}^2 . \quad (\text{C.21c})$$

They propose the practical form

$$\frac{1}{|\mathbb{B}|(|\mathbb{B}| - 1)} \mathbb{E} \left[ \left\| \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|} - \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^2} \frac{\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|} \right\|^2 \right] \leq \nu_{\text{ortho}}^2 , \quad (\text{C.22a})$$

which simplifies to

$$\frac{1}{|\mathbb{B}|(|\mathbb{B}| - 1)} \sum_{n \in \mathbb{B}} \left( \frac{\|\mathbf{g}_n(\boldsymbol{\theta})\|^2}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^2} - \frac{(\mathbf{g}_n(\boldsymbol{\theta})^\top \mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta}))^2}{\|\mathbf{g}_{\mathbb{B}}(\boldsymbol{\theta})\|^4} \right) \leq \nu_{\text{ortho}}^2 . \quad (\text{C.22b})$$

It relies on pairwise scalar products between individual gradients which can be aggregated over layers during a backward pass. Cockpit's `OrthoTest` quantity computes the maximum band width  $\nu_{\text{ortho}}$  which satisfies [Equation \(C.22b\)](#).

### Relation to Acute Angle Test

Recently, a novel “acute angle test” was proposed by Bahamou and Goldfarb [6]. While the theoretical constraint between  $\mathbf{g}_{\mathbb{B}}(\theta)$  and  $\mathbf{g}_{p_{\text{data}}}(\theta)$  differs from the orthogonality test, the practical versions coincide. Hence, we do not incorporate the acute angle here.

### C.3.7 Hessian Maximum Eigenvalue (`HessMaxEV`)

The Hessian’s maximum eigenvalue  $\lambda_{\max}(\mathbf{H}_{\mathbb{B}}(\theta))$  is computed with an iterative eigensolver from Hessian-vector products through PyTorch’s automatic differentiation [127]. Like Yao et al. [177], we employ power iterations with similar **default stopping parameters** (stop after at most 100 iterations, or if the iterate does converge with a relative and absolute tolerance of  $10^{-3}, 10^{-6}$ , respectively) to compute  $\lambda_{\max}(\mathbf{H}_{\mathbb{B}}(\theta))$  through the `HessMaxEV` quantity in Cockpit.

In principle, more sophisticated eigensolvers (for example Arnoldi’s method) could be applied to converge in fewer iterations or compute eigenvalues other than the leading ones. Warsa et al. [170] empirically demonstrate that the FLOP ratio between power iteration and implicitly restarted Arnoldi method can reach values larger than 100. While we can use such a beneficial method on a CPU through `scipy.sparse.linalg.eigsh` we are restricted to the GPU-compatible power iteration for GPU training. We expect that extending the support of popular machine learning libraries like PyTorch for such iterative eigensolvers on GPUs can help to save computation time.

$$\lambda_{\max}(\mathbf{H}_{\mathbb{B}}(\theta)) = \max_{\|v\|_2=1} \|\mathbf{H}_{\mathbb{B}}(\theta)v\|_2 = \max_{v \in \mathbb{R}^D} \frac{v^\top \mathbf{H}_{\mathbb{B}}(\theta)v}{v^\top v}. \quad (\text{C.23})$$

### Usage

The Hessian’s maximum eigenvalue describes the loss surface’s sharpest direction and thus provides an understanding of the current loss landscape. Additionally, in convex optimization, the largest Hessian eigenvalue crucially determines the appropriate step size [144]. In Section 6.4, we can observe that although training seems stuck in the very first few iterations progress is visible when looking at the maximum Hessian eigenvalue.

### C.3.8 Hessian Trace (`HessTrace`)

In comparison to Yao et al. [177], who leverage Hessian-vector products [127] to estimate the Hessian trace, we compute the exact value  $\text{Tr}(\mathbf{H}_{\mathbb{B}}(\theta))$  with the `HessTrace` quantity in Cockpit by aggregating the output of BackPACK’s `DiagHessian` extension, which computes the diagonal entries of  $\mathbf{H}_{\mathbb{B}}(\theta)$ . Alternatively, the trace can also be estimated with the GGN matrix, or an MC-sampled approximation thereof.

## Usage

The Hessian trace equals the sum of the eigenvalues and thus provides a notion of “average curvature” of the current loss landscape. It has long been theorized and discussed that curvature and generalization performance may be linked [70, e.g.].

### C.3.9 Takeuchi Information Criterion (`TICDiag`, `TICTrace`)

Recent work by Thomas et al. [162] suggests that optimizer convergence speed and generalization is mainly influenced by curvature and gradient noise; and hence their interaction is crucial to understand the generalization and optimization behavior of deep neural networks. They reinvestigate the Takeuchi Information criterion [160], an estimator for the generalization gap in overparameterized maximum likelihood estimation. At a local minimum  $\theta_\star$ , the generalization gap is estimated by the TIC

$$\frac{1}{|\mathbb{D}|} \text{Tr}(\mathbf{H}_{p_{\text{data}}}(\theta_\star)^{-1} \mathbf{K}_{p_{\text{data}}}(\theta_\star)) , \quad (\text{C.24})$$

where  $\mathbf{H}_{p_{\text{data}}}(\theta_\star)$  is the population Hessian and  $\mathbf{K}_{p_{\text{data}}}(\theta_\star)$  is the gradient’s uncentered second moment,

$$\mathbf{K}_{p_{\text{data}}}(\theta_\star) = \int \nabla_{\theta_\star} \ell(f_{\theta_\star}(x), y) (\nabla_{\theta_\star} \ell(f_{\theta_\star}(x), y))^\top p_{\text{data}}(x, y) dx dy.$$

Both matrices are inaccessible in practice. In their experiments, Thomas et al. [162] propose the approximation  $\text{Tr}(\mathbf{K})/\text{Tr}(\mathbf{H})$  for  $\text{Tr}(\mathbf{H}^{-1}\mathbf{K})$ . They also replace the Hessian by the Fisher as it is easier to compute. With these practical simplifications, they investigate the TIC of trained neural networks where the curvature and noise matrix are evaluated on a large dataset.

The TIC provided in Cockpit differs from this setting, since by design we want to observe quantities during training, while avoiding additional model predictions. Also, BackPACK provides access to the Hessian; hence we don’t need to use the Fisher. We propose the following two approximations of the TIC from a mini-batch:

- ▶ `TICTrace`: uses the approximation of Thomas et al. [162] which replaces the matrix-product trace by the product of traces,

$$\frac{\text{Tr}(\mathbf{K}_\mathbb{B}(\theta))}{\text{Tr}(\mathbf{H}_\mathbb{B}(\theta))} = \frac{\frac{1}{|\mathbb{B}|} \sum_{n \in \mathbb{B}} \|g_n(\theta)\|^2}{\text{Tr}(\mathbf{H}_\mathbb{B}(\theta))} . \quad (\text{C.25})$$

- ▶ `TICDiag`: uses a diagonal approximation of the Hessian, which is cheap to invert,

$$\begin{aligned} & \text{Tr} \left( \text{diag}(\mathbf{H}_\mathbb{B}(\theta))^{-1} \mathbf{K}_\mathbb{B}(\theta) \right) \\ &= \frac{1}{|\mathbb{B}|} \sum_{d=1}^D [\mathbf{H}_\mathbb{B}(\theta)]_{d,d}^{-1} \left[ \sum_{n \in \mathbb{B}} g_n(\theta)^{\odot 2} \right]_d . \end{aligned} \quad (\text{C.26})$$

## Usage

The TIC is a proxy for the generalization gap, see Thomas et al. [162].

### C.3.10 Gradient Signal-to-noise Ratio (`MeanGSNR`)

The gradient signal-to-noise ratio  $\text{GSNR}([\boldsymbol{\theta}]_d) \in \mathbb{R}$  for a single parameter  $[\boldsymbol{\theta}]_d$  is defined as

$$\text{GSNR}([\boldsymbol{\theta}]_d) = \frac{\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P} \left[ [\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y})]_d \right]^2}{\text{Var}_{(\mathbf{x}, \mathbf{y}) \sim P} \left[ [\nabla_{\boldsymbol{\theta}} \ell(f_{\boldsymbol{\theta}}(\mathbf{x}), \mathbf{y})]_d \right]} = \frac{[g_p(\boldsymbol{\theta})]_d^2}{[\Sigma_p(\boldsymbol{\theta})]_{d,d}}. \quad (\text{C.27})$$

Liu et al. [100] use it to explain generalization properties of models in the early training phase. We apply their estimation to mini-batches,

$$\text{GSNR}([\boldsymbol{\theta}]_d) \approx \frac{[g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2}{\frac{|\mathbb{B}| - 1}{|\mathbb{B}|} [\hat{\Sigma}_{\mathbb{B}}(\boldsymbol{\theta})]_{d,d}} = \frac{[g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2}{\frac{1}{|\mathbb{B}|} \left( \sum_{n \in \mathbb{B}} [g_n(\boldsymbol{\theta})]_d^2 \right) - [g_{\mathbb{B}}(\boldsymbol{\theta})]_d^2}. \quad (\text{C.28a})$$

Inspired by Liu et al. [100], Cockpit's `MeanGSNR` computes the average GSNR over all parameters,

$$\frac{1}{D} \sum_{j=1}^D \text{GSNR}([\boldsymbol{\theta}]_j). \quad (\text{C.28b})$$

## Usage

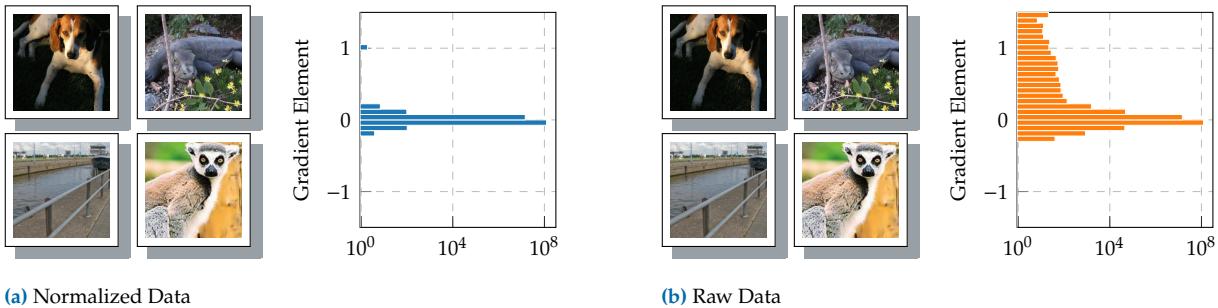
The GSNR describes the gradient noise level which is influenced, among other things, by the batch size. Using the GSNR, perhaps in combination with the gradient tests or the CABS criterion could provide practitioners a clearer picture of suitable batch sizes for their particular problem. As shown by Liu et al. [100], the GSNR is also linked to generalization of neural networks.

## C.4 Additional Experiments

In this section, we present additional experiments and use cases that showcase Cockpit's utility. Appendix C.4.1 shows that Cockpit is able to scale to larger datasets by running the experiment with incorrectly scaled data (see Section 6.3.1) on ImageNet instead of CIFAR-10. Appendix C.4.2 provides another concrete use case similar to Figure 6.1: detecting regularization during training.

### C.4.1 Incorrectly Scaled Data for ImageNet

We repeat the experiment of Section 6.3.1 on the ImageNet [41] dataset instead of CIFAR-10. We also use a larger neural network model, switching from 3c3d to VGG16 [153]. This demonstrates that Cockpit is able to scale



(a) Normalized Data

(b) Raw Data

**Figure C.3: Same inputs, different gradients on ImageNet.** This is structurally the same plot as Figure 6.3, but using ImageNet and VGG16. (a) normalized ([0, 1]) and (b) raw ([0, 255]) images look identical in auto-scaled front-ends like matplotlib's `imshow`. The gradient distribution on the VGG16 model, however, is affected by this scaling.

to both larger models and datasets. The input size of the images is almost fifty times larger ( $224 \times 224$  instead of  $32 \times 32$ ). The model size increased by roughly a factor of 150 (VGG16 for ImageNet has roughly 138 million parameters, 3c3d has less than a million).

Similar to the example in the main text, the gradients are affected by the scaling introduced via the input images, albeit less drastically (see Figure C.3). Due to the gradient scaling, default optimization hyperparameters might not work well anymore for the model using the raw data.

#### C.4.2 Detecting Implicit Regularization of The Optimizer

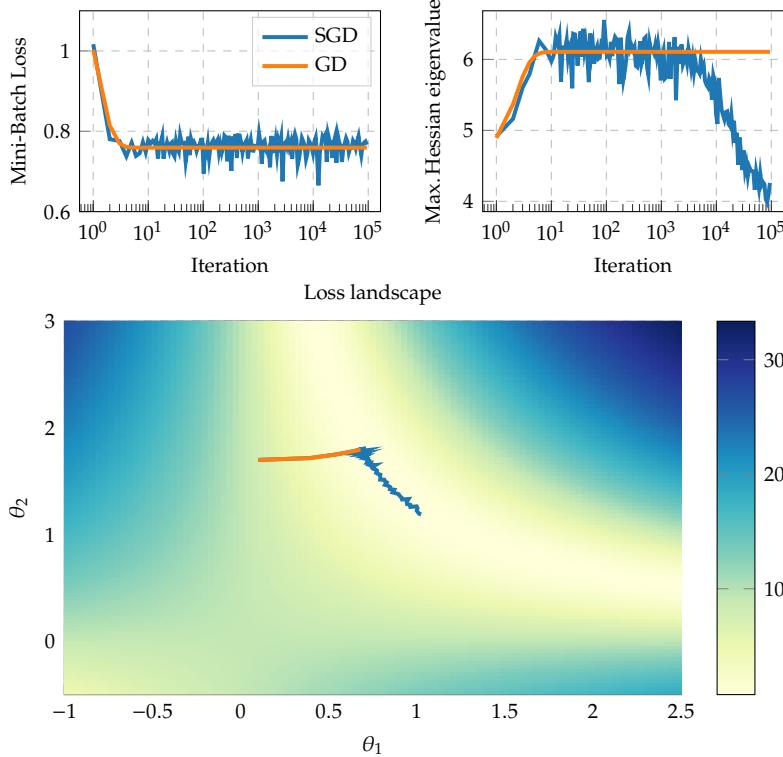
In non-convex optimization, optimizers can converge to local minima with different properties. Here, we illustrate this by investigating the effect of sub-sampling noise on a simple task from [56, 112].

We generate synthetic data  $\mathbb{D} = \{(x_n, y_n) \in \mathbb{R} \times \mathbb{R}\}_{n=1}^{N=100}$  for a regression task with  $x \sim \mathcal{N}(x \mid 0, 1)$  with noisy observations  $y = 1.4x + \epsilon$  where  $\epsilon \sim \mathcal{N}(\epsilon \mid 0, 1)$ . The model is a scalar net with parameters  $\theta = (w_1 \ w_2)^\top \in \mathbb{R}^2$ , initialized at  $\theta_0 = (0.1 \ 1.7)^\top$ , that produces predictions  $f_\theta(x) = w_2 w_1 x$ . We seek to minimize the mean squared error

$$\mathcal{L}_{\mathbb{D}}(\theta) = \frac{1}{N} \sum_{n=1}^N (f_\theta(x_n) - y_n)^2$$

and compare SGD ( $|\mathbb{B}| = 95$ ) with GD ( $|\mathbb{B}| = N = 100$ ) at a learning rate of 0.1 (see Figure C.4).

We observe that the loss of both SGD and GD is almost identical. Using a noisy gradient regularizes the Hessian's maximum eigenvalue though. It decreases in later stages where the loss curve suggests that training has converged. This regularization effect constitutes an important phenomenon that cannot be observed by monitoring only the loss.



**Figure C.4:** Observing implicit regularization of the optimizer with Cockpit through a comparison of SGD and GD on a synthetic problem inspired by [56, 112] (details in the text). *Top left:* The mini-batch loss of both optimizers looks similar. *Top right:* Noise due to mini-batching regularizes the Hessian’s maximum eigenvalue in stages where the loss suggests that training has converged. *Bottom:* Optimization trajectories in parameter space. SGD is attracted to the flattest minimum.

## C.5 Implementation Details & Additional Benchmarks

In this section, we provide more details about our implementation (Appendix C.5.1) to access the desired quantities with as little overhead as possible. Additionally, we present more benchmarks for individual instruments (Appendix C.5.2) and Cockpit configurations (Appendix C.5.2). These are similar but extended versions of the ones presented in Figures 6.6a and 6.6b in the main text. Lastly, we benchmark different implementations of computing the two-dimensional gradient histogram (Appendix C.5.3), identifying a computational bottleneck for its current GPU implementation.

### Hardware Details

We conducted benchmarks on the following setup:

- ▶ **CPU:** Intel Core i7-8700K CPU @ 3.70 GHz × 12 (32 GB)
- ▶ **GPU:** NVIDIA GeForce RTX 2080 Ti (11 GB)

### Test Problem Details

The experiments in this paper rely mostly on optimization problems provided by the DeepOBS benchmark suite [146]. If not stated otherwise, we use the default training details suggested by DeepOBS, that are summarized below. For more details see the original paper.

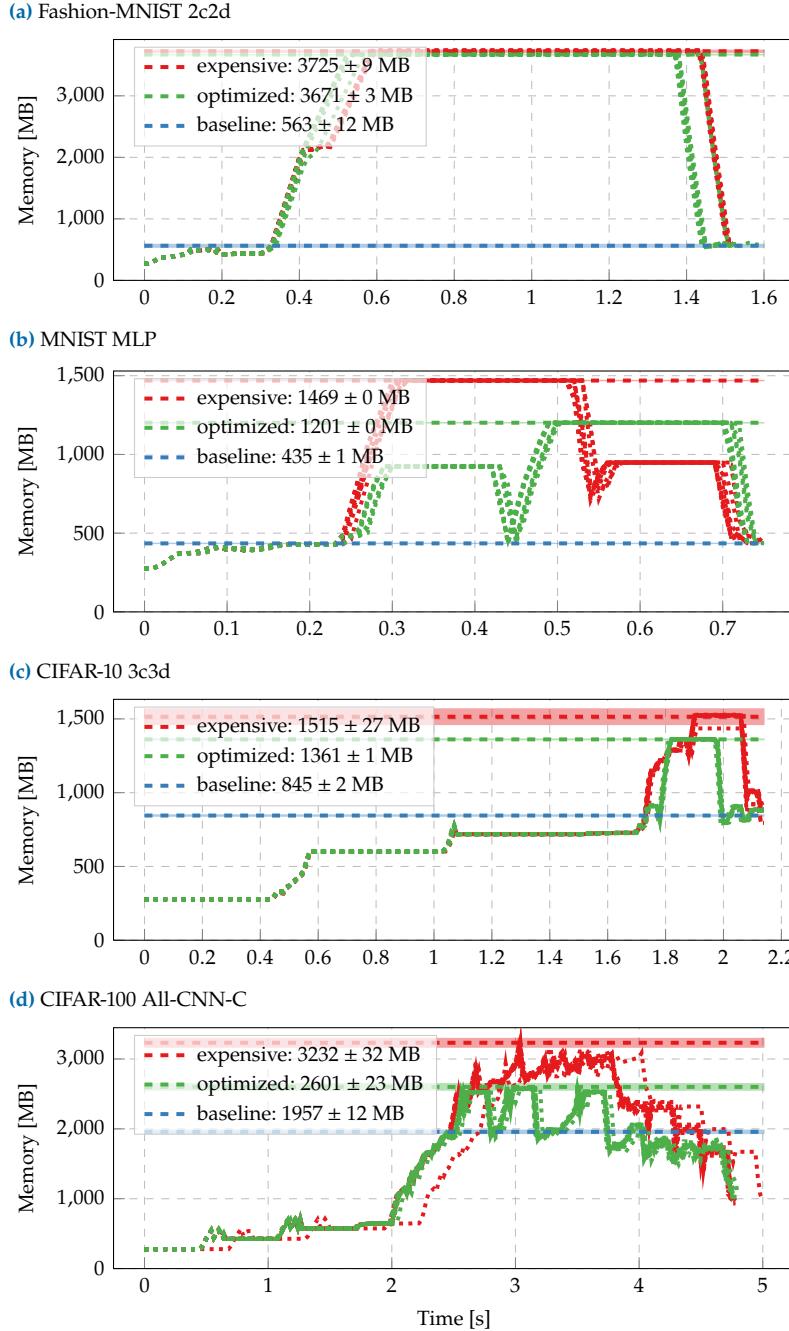
- ▶ **Quadratic Deep:** A stochastic quadratic problem with an eigen-spectrum similar to what has been reported for neural nets. Default batch size 128, default number of epochs 100.
- ▶ **MNIST Log. Reg.:** Multinomial logistic regression on MNIST [95]. Default batch size 128, default number of epochs 50.
- ▶ **MNIST MLP:** Multi-layer perceptron on MNIST. Default batch size 128, default number of epochs 100.
- ▶ **Fashion-MNIST MLP:** Multi-layer perceptron on Fashion-MNIST [175]. Default batch size 128, default number of epochs 100.
- ▶ **Fashion-MNIST 2c2d:** A two convolutional and two dense layered neural network on Fashion-MNIST. Default batch size 128, default number of epochs 100.
- ▶ **CIFAR-10 3c3d:** A three convolutional and three dense layered neural network on CIFAR-10 [90]. Default batch size 128, default number of epochs 100.
- ▶ **CIFAR-100 All-CNN-C:** All Convolutional Neural Network C (All-CNN-C [157]) on CIFAR-100 [90]. Default batch size 256, default number of epochs 350.
- ▶ **SVHN 3c3d:** A three convolutional and three dense layered neural network on SVHN [118]. Default batch size 128, default number of epochs 100.

### C.5.1 Hooks & Memory Benchmarks

To improve memory consumption, we compact information during the backward pass by adding hooks to the neural network's layers. These are executed after BackPACK extensions and have access to the quantities computed therein. They compress information to what is requested by a quantity and free the memory occupied by BackPACK buffers. Such savings primarily depend on the parameter distribution over layers, and are bigger for more balanced architectures (compare Figure C.5).

#### Example

Say, we want to compute a histogram over the  $|\mathbb{B}| \times D$  individual gradient elements of a network. Suppose that  $|\mathbb{B}| = 128$  and the model is DeepOBS's CIFAR-10 3c3d test problem with 895,210 parameters. Given that every parameter is stored in single precision, the model requires  $895,210 \times 4 \text{ Bytes} \approx 3.41 \text{ MB}$ . Storing the individual gradients will require  $128 \times 895,210 \times 4 \text{ Bytes} \approx 437 \text{ MB}$  (for larger networks this quickly exceeds the available memory as the individual gradients occupy  $|\mathbb{B}|$  times the model size). If instead, the layer-wise individual gradients are condensed into histograms of negligible size and immediately freed afterwards during backpropagation, the maximum memory overhead reduces to storing the individual gradients of the largest layer. For our example, the largest layer has 589,824 parameters, and the associated individual gradients will require  $128 \times 589,824 \times 4 \text{ Bytes} \approx 288 \text{ MB}$ , saving roughly 149 MB of RAM. In practice, we observe these expected savings, see Figure C.5c.



**Figure C.5:** Memory consumption and savings with hooks during one forward-backward step on a CPU for different DeepOBS problems. We compare three settings; i) without Cockpit (baseline); ii) Cockpit with `GradHist1d` with BackPACK (expensive); iii) Cockpit with `GradHist1d` with BackPACK and additional hooks (optimized). Peak memory consumptions are highlighted by horizontal dashed bars and shown in the legend. Shaded areas, if visible, fill two standard deviations above and below the mean value, all of them result from ten independent runs. Dotted lines indicate individual runs. Our optimized approach allows to free obsolete tensors during backpropagation and thereby reduces memory consumption. From top to bottom: the effect is less pronounced for architectures that concentrate the majority of parameters in a single layer ((a) 3,274,634 total, 3,211,264 largest layer) and increases for more balanced networks (b) 1,336,610 total, 784,000 largest layer, (c): 895,210 total, 589,824 largest layer).

## C.5.2 Additional Run Time Benchmarks

### Individual Instrument Overhead

To estimate the computational overhead for individual instruments, we run Cockpit with that instrument for 32 iterations, tracking at every step. Training proceeds with the default batch size specified by the DeepOBS problem and uses SGD with learning rate  $10^{-3}$ . We measure the time between iterations 1 and 32, and average for the overhead per step. Every such estimate is repeated over 10 random seeds to obtain mean and error bars as reported in Figure 6.6a.

Note that this protocol does *not* include initial overhead for setting up

data loading and also does *not* include the time for evaluating train/test loss on a larger dataset, which is usually done by practitioners. Hence, we even expect the shown overheads to be smaller in a conventional training loop which includes the above steps.

### Individual Overhead on GPU Versus CPU

[Figure C.6](#) and [Figure C.7](#) show the individual overhead for four different DeepOBS problems on GPU and CPU, respectively. The left part of [Figure C.6](#) (c) corresponds to [Figure 6.6a](#). Right panels show the expensive quantities, which we omitted in the main text as they were expected to be expensive due to their computational work (`HessMaxEV`) or bottlenecks in the implementation (`GradHist2d`, see [Appendix C.5.3](#) for details). We see that they are in many cases equally or more expensive than computing all other instruments. Another expected feature of the GPU-to-CPU comparison is that parallelism on the CPU is significantly less pronounced. Hence, we observe an increased overhead for all quantities that contain non-linear transformations and contractions of the high-dimensional individual gradients, or require additional backpropagations (curvature).

### Configuration Overhead

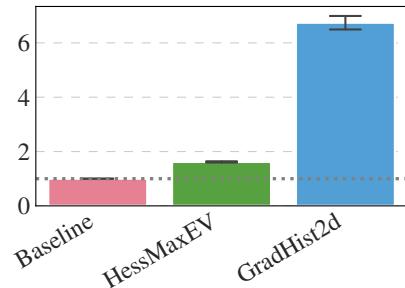
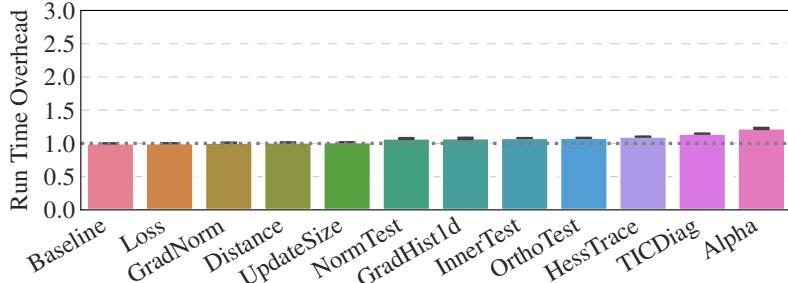
For the estimation of different Cockpit configuration overheads, we use almost the same setting as described above, training for 512 iterations and tracking only every specified interval.

### Configuration Overhead on GPU Versus CPU

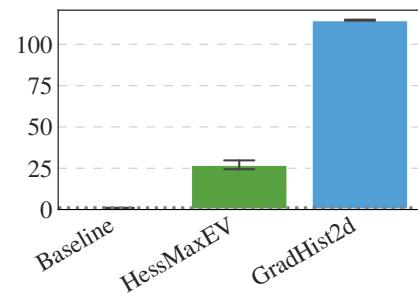
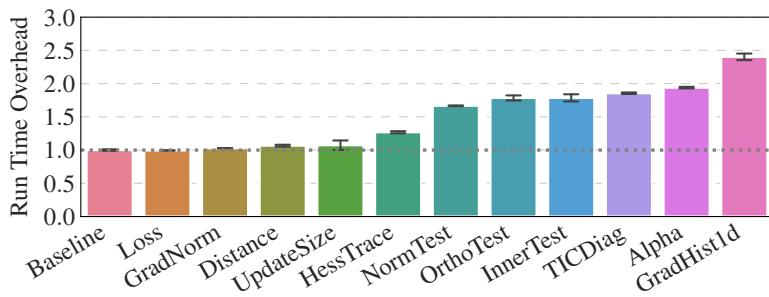
[Figure C.8](#) and [Figure C.9](#) show the configuration overhead for four different DeepOBS problems. The bottom left part of [Figure C.8](#) corresponds to [Figure 6.6b](#). In general, increased parallelism can be exploited on a GPU, leading to smaller overheads in comparison to a CPU.

Cockpit can even scale to significantly larger problems, such as a ResNet-50 on ImageNet-like data. [Figure C.10](#) shows the computational overhead for different tracking intervals on such a large-scale problem. Using the *economy* configuration, we can achieve our self-imposed goal of at most doubling the run time even when tracking every fourth step. More extensive configurations (such as the *full* set) would indeed have almost prohibitively large costs associated. However, these costs could be dramatically reduced when one decides to only inspect a part of the network using Cockpit. Note, individual gradients are not properly defined when using batch norm, therefore, we replaced these batch norm layers with identity layers when using the ResNet-50.

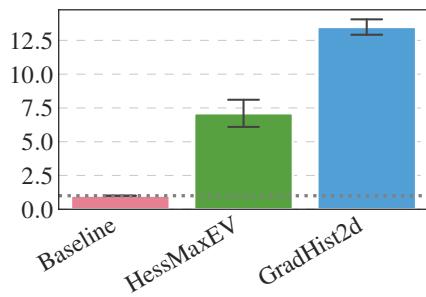
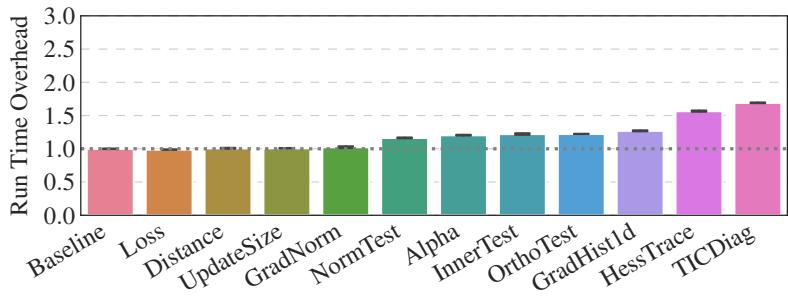
(a) Computational overhead for MNIST LogReg (GPU)



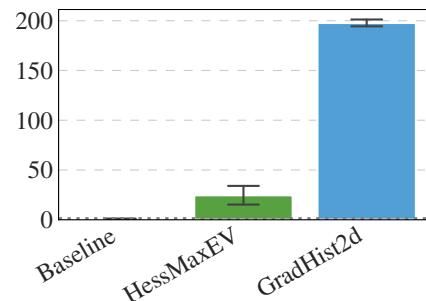
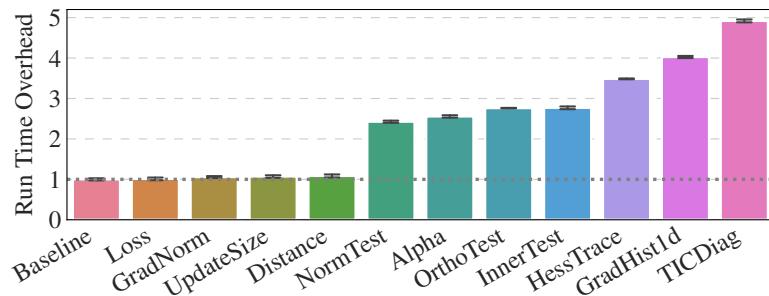
(b) Computational overhead for MNIST MLP (GPU)



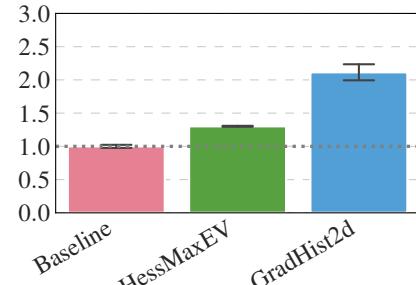
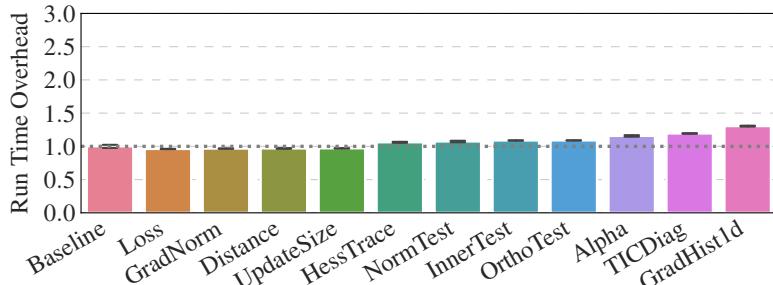
(c) Computational overhead for CIFAR-10 3c3d (GPU)



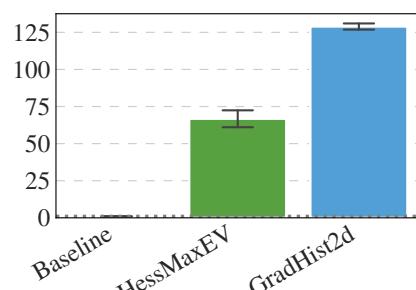
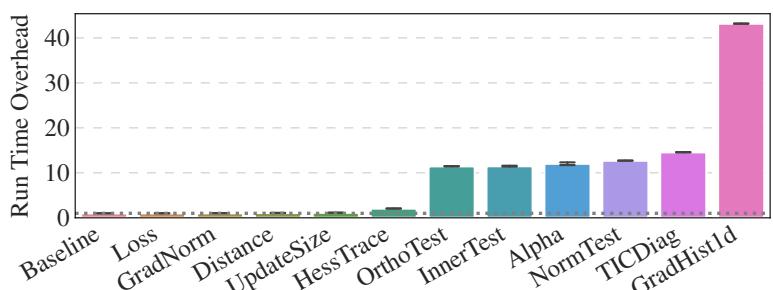
(d) Computational overhead for Fashion-MNIST 2c2d (GPU)

**Figure C.6:** Individual overhead of Cockpit's instruments on GPU for four different problems. All run times are shown as multiples of the baseline without tracking. Expensive quantities are displayed in separate panels on the right. Experimental details in the text.

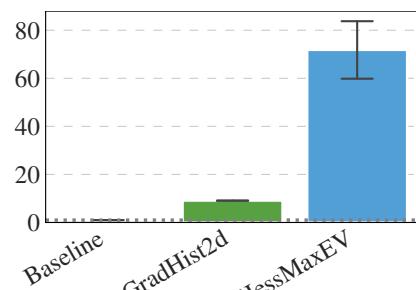
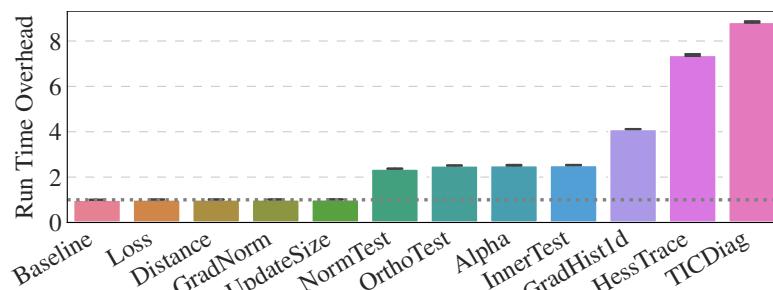
(a) Computational overhead for MNIST LogReg (CPU)



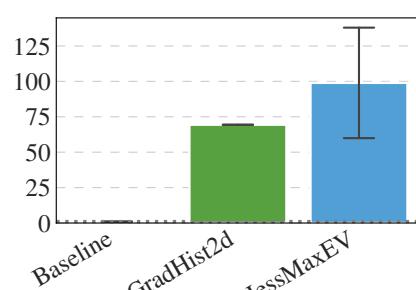
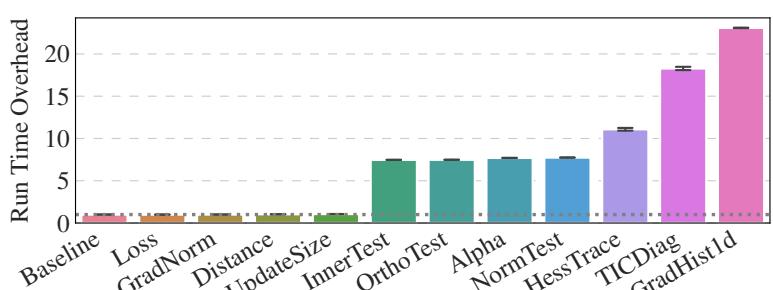
(b) Computational overhead for MNIST MLP (CPU)



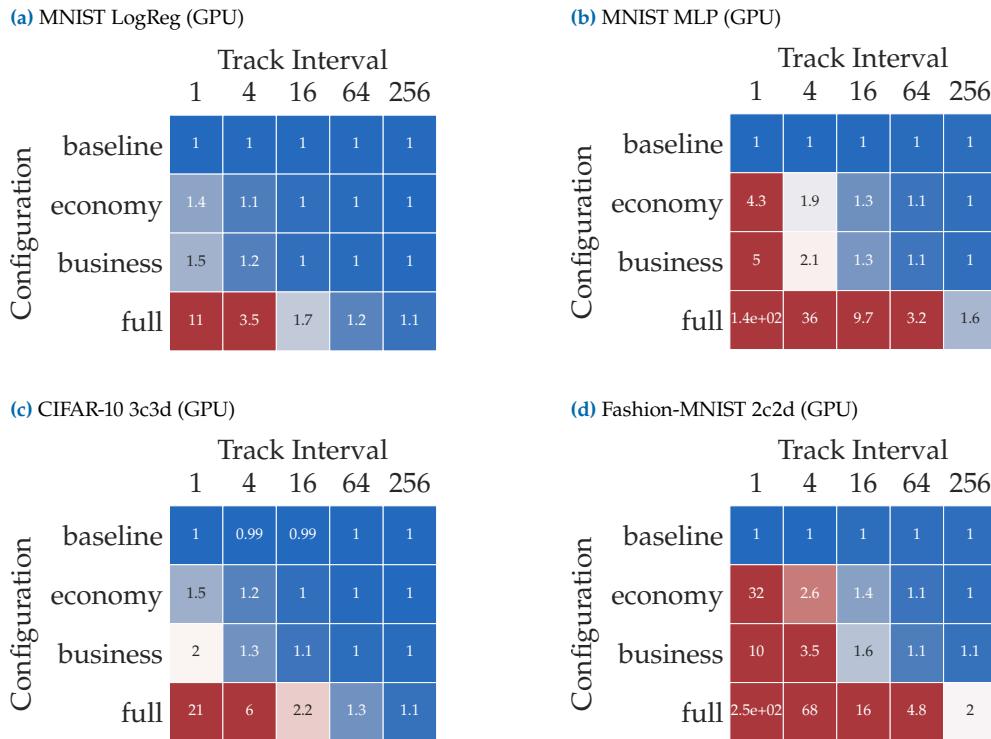
(c) Computational overhead for CIFAR-10 3c3d (CPU)



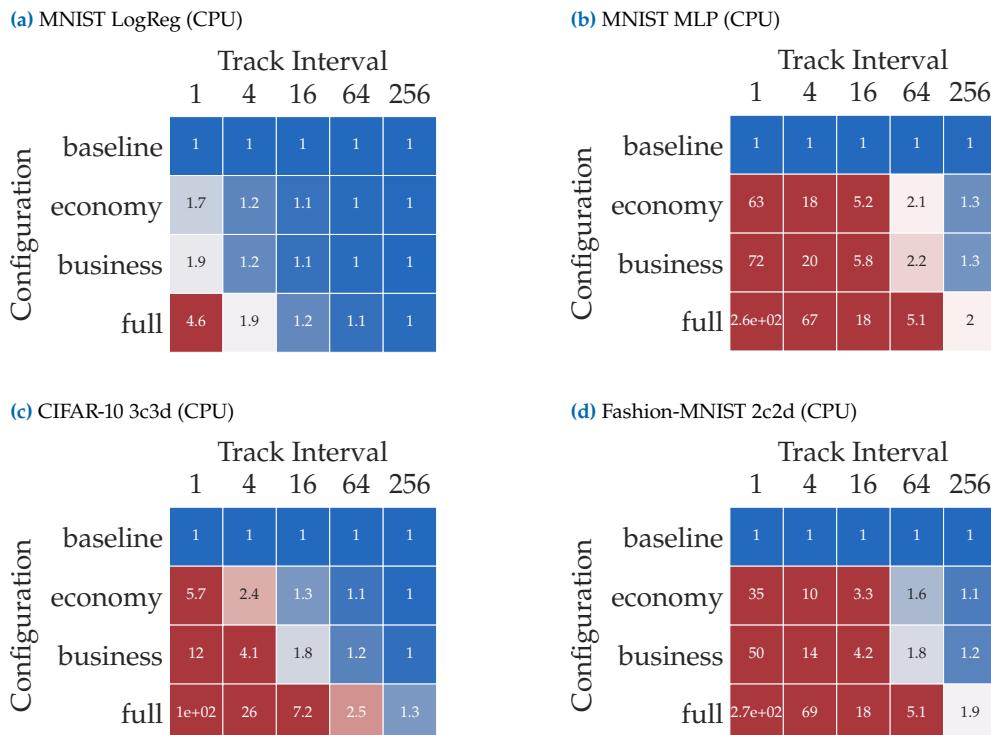
(d) Computational overhead for Fashion-MNIST 2c2d (CPU)



**Figure C.7:** Individual overhead of Cockpit's instruments on CPU for four different problems. All run times are shown as multiples of the *baseline* without tracking. Expensive quantities are displayed in separate panels on the right. Experimental details in the text.



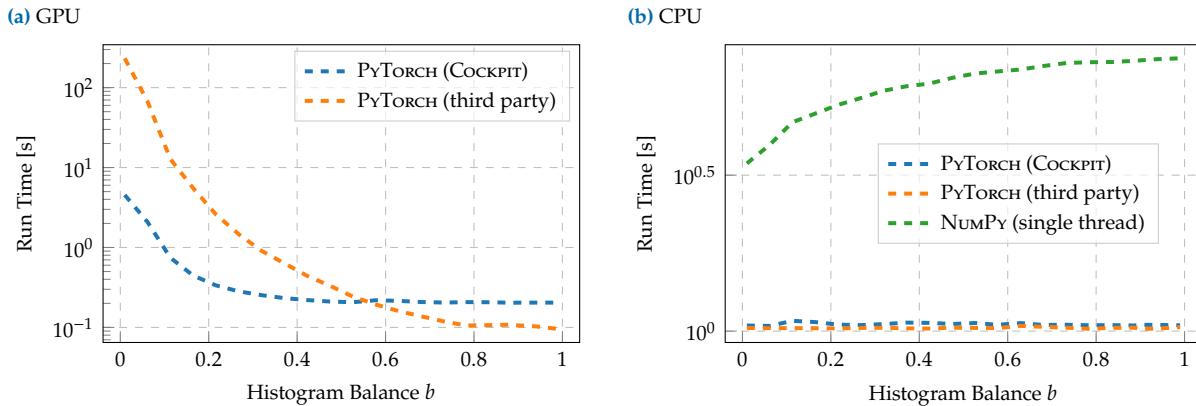
**Figure C.8:** Overhead of Cockpit configurations on GPU for four different problems with varying tracking interval. Color bar is the same as in Figure 6.6.



**Figure C.9:** Overhead of Cockpit configurations on CPU for four different problems with varying tracking interval. Color bar is the same as in Figure 6.6.

**Figure C.10:** Overhead of Cockpit configurations on GPU for ResNet-50 on ImageNet. Cockpit’s instruments scale efficiently even to very large problems (here: 1000 classes, (3, 224, 224)-sized inputs, and a batch size of 64. For individual gradients to be defined, we replaced the batch norm layers of the ResNet-50 model with identities.) Color bar is the same as in Figure 6.6.

	Track Interval				
	1	4	16	64	256
baseline	1	1	1	1	1
economy	3.7	1.9	1.2	1.1	1



**Figure C.11:** Performance of two-dimensional histogram GPU implementations depends on the data. (a) Run time for two different GPU implementations with histograms of different imbalance. Cockpit’s implementation outperforms the third party solution by more than one order of magnitude in the deep learning regime ( $b \ll 1$ ). (b) On CPU, performance is robust to histogram balance. The run time difference between NumPy and PyTorch is due to multi-threading. Data has the same size as DeepOBS’s CIFAR-10 3c3d problem ( $D = 895, 210$ ,  $|\mathbb{B}| = 128$ ). Curves represent averages over 10 independent runs. Error bars are omitted to improve legibility.

### C.5.3 Performance of Two-dimensional Histograms

Both one- and two-dimensional histograms require  $|\mathbb{B}| \times D$  elements be accessed, and hence perform similarly. However, we observed different behavior on GPU and decided to omit the two-dimensional histogram’s run time in the main text. As explained here, this performance lack is not fundamental, but a shortcoming of the GPU implementation. PyTorch provides built-in functionality for computing one-dimensional histograms at the time of writing, but is not yet featuring multi-dimensional histograms. We experimented with three implementations:

- ▶ **PyTorch (third party):** A third party implementation\* under review for being integrated into PyTorch†. It relies on `torch.bincount`, which uses `atomicAdd`s that represent a bottleneck for histograms where most counts are contained in one bin.‡ This occurs often for over-parameterized deep models, as most of the gradient elements are zero.
- ▶ **PyTorch (Cockpit):** Our implementation uses a workaround, computes bin indices and scatters the counts into their associated bins with `torch.Tensor.put_`. This circumvents `atomicAdd`s, but has poor memory locality.
- ▶ **NumPy:** The single-threaded `numpy.histogram2d` serves as baseline, but does not run on GPUs.

\* Permission granted by the authors of [github.com/miranov25/histogramdd\\_pytorch.py](https://github.com/miranov25/histogramdd_pytorch.py).

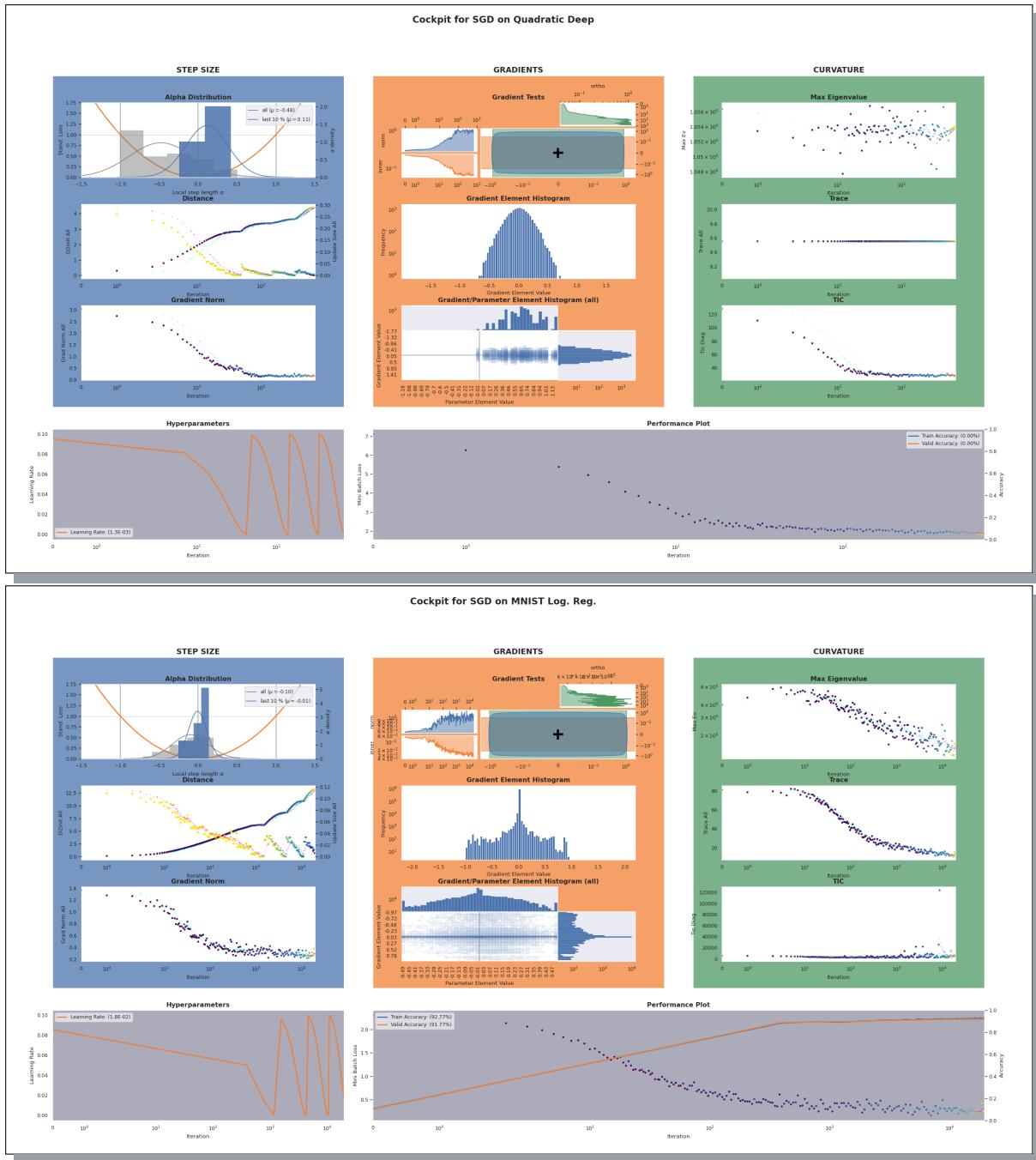
† See <https://github.com/pytorch/pytorch/pull/44485>.

‡ See <https://discuss.pytorch.org/t/torch-bincount-1000x-slower-on-cuda/42654>

To demonstrate the strong performance dependence on the data, we generate data from a uniform distribution over  $[0, b] \times [0, b]$ , where  $b \in (0, 1)$  parametrizes the histogram's balance, and compute two-dimensional histograms on  $[0, 1] \times [0, 1]$ . Figure C.11a shows a clear increase in run time of both GPU implementations for more imbalanced histograms. Note that even though our implementation outperforms the third party by more than one order of magnitude in the deep neural network regime ( $b \ll 1$ ), it is still considerably slower than a one-dimensional histogram (see Figure C.6 (c)), and even slower on GPU than on CPU (Figure C.11 (b)). As expected, the CPU implementations do not significantly depend on the data (Figure C.11b). The performance difference between PyTorch and NumPy is likely due to multi-threading versus single-threading.

Although a carefully engineered histogram GPU implementation is currently not available, we think it will reduce the computational overhead to that of a one-dimensional histogram in future releases.

## C.6 Cockpit View of Convex Stochastic Problems



**Figure C.12:** Screenshot of Cockpit's full view for convex DeepOBS problems. Top Cockpit shows training on a noisy quadratic loss function. Bottom shows training on logistic regression on MNIST. Figure and labels are not meant to be legible. It is evident, that there is a fundamental difference in the optimization process, compared to training deep networks, i.e. Figure 6.2. This is, for example, visible when comparing the gradient norms, which converge to zero for convex problems but not for deep learning.

# Additional Material for Chapter 7

# D.

## D.1 Mathematical Details

D.1 Mathematical Details . . .	167
D.2 Experimental Details . . .	168
D.3 Implementation Details . . .	193

### D.1.1 The GNN's Eigenvalues & the Gram Matrix

For [Equation \(7.4\)](#), consider the left hand side of the GGN's characteristic polynomial  $\det(G - \lambda I_D) = 0$ . Inserting the ViViT factorization ([Equation \(7.3\)](#)) and using the matrix determinant lemma yields

$$\begin{aligned}
 & \det(-\lambda I_D + G) \\
 &= \det(-\lambda I_D + VV^\top) && \text{(Low-rank structure (7.3))} \\
 &= \det(I_{NC} + V^\top(-\lambda I_D)^{-1}V) \det(-\lambda I_D) && \text{(Matrix determinant lemma)} \\
 &= \det\left(I_{NC} - \frac{1}{\lambda}V^\top V\right)(-\lambda)^D \\
 &= \left(-\frac{1}{\lambda}\right)^{NC} \det(V^\top V - \lambda I_{NC}) (-\lambda)^D \\
 &= (-\lambda)^{D-NC} \det(\tilde{G} - \lambda I_{NC}) . && \text{(Gram matrix)}
 \end{aligned}$$

Setting the above expression to zero reveals that the GGN's spectrum decomposes into  $D-NC$  zero eigenvalues and the Gram matrix spectrum obtained from  $\det(\tilde{G} - \lambda I_{NC}) = 0$ .

### D.1.2 Relation Between GGN & Gram Matrix Eigenvectors

Assume the nontrivial Gram matrix spectrum  $\tilde{\mathbb{S}}_+ := \{(\lambda_k, \tilde{e}_k) \mid \lambda_k \neq 0, \tilde{G}\tilde{e}_k = \lambda_k\tilde{e}_k\}_{k=1}^K$  with orthonormal eigenvectors  $\tilde{e}_j^\top \tilde{e}_k = \delta_{j,k}$  ( $\delta$  is the Kronecker delta) and  $K = \text{rank}(G)$ . We now show that  $e_k = 1/\sqrt{\lambda_k}V\tilde{e}_k$  are normalized eigenvectors of  $G$  and inherit orthogonality from  $\tilde{e}_k$ .

To see the first, consider right-multiplication of the GGN with  $e_k$ , then expand the low-rank structure,

$$\begin{aligned}
 Ge_k &= \frac{1}{\sqrt{\lambda_k}}VV^\top V\tilde{e}_k && \text{(Equation (7.3) and definition of } e_k) \\
 &= \frac{1}{\sqrt{\lambda_k}}V\tilde{G}\tilde{e}_k && \text{(Gram matrix)} \\
 &= \lambda_k \frac{1}{\sqrt{\lambda_k}}V\tilde{e}_k && \text{(Eigenvector property of } \tilde{e}_k) \\
 &= \lambda_k e_k .
 \end{aligned}$$

Orthonormality of the  $e_k$  results from the Gram matrix eigenvector

orthonormality,

$$\begin{aligned}
 e_j^\top e_k &= \left( \frac{1}{\sqrt{\lambda_j}} \tilde{e}_j^\top V^\top \right) \left( \frac{1}{\sqrt{\lambda_k}} V \tilde{e}_k \right) && \text{(Definition of } e_j, e_k\text{)} \\
 &= \frac{1}{\sqrt{\lambda_j \lambda_k}} \tilde{e}_j^\top \tilde{G} \tilde{e}_k && \text{(Gram matrix)} \\
 &= \frac{\lambda_k}{\sqrt{\lambda_j \lambda_k}} \tilde{e}_j^\top \tilde{e}_k && \text{(Eigenvector property of } \tilde{e}_k\text{)} \\
 &= \delta_{j,k}. && \text{(Orthonormality)}
 \end{aligned}$$

## D.2 Experimental Details

This section uses the notation from [Section 7.3](#) (see [Table D.1](#)).

**Table D.1: Notation for curvature approximations.** The notation is introduced in [Section 7.3](#). This table recapitulates the abbreviations (referring to the approximations introduced in [Section 7.2.4](#)) and provides corresponding explanations.

Abbreviation	Explanation
<b>mb, exact</b>	Exact GGN with all mini-batch samples. Backpropagates $NC$ vectors.
<b>mb, mc</b>	MC-approximated GGN with all mini-batch samples. Backpropagates $NM$ vectors with $M$ the number of MC-samples.
<b>sub, exact</b>	Exact GGN on a subset of mini-batch samples ( $\lfloor N/8 \rfloor$ as in <a href="#">[183]</a> ). Backpropagates $\lfloor N/8 \rfloor C$ vectors.
<b>sub, mc</b>	MC-approximated GGN on a subset of mini-batch samples. Backpropagates $\lfloor N/8 \rfloor M$ vectors with $M$ the number of MC-samples.

### GGN Spectra ([Figure 7.1a](#))

To obtain the spectra of [Figure 7.1a](#) we initialize the respective architecture, then draw a mini-batch and evaluate the GGN eigenvalues under the described approximations, clipping the Gram matrix eigenvalues at  $10^{-4}$ . [Figure D.1](#) provides the spectra for all used architectures with both the full GGN and a per-layer block-diagonal approximation.

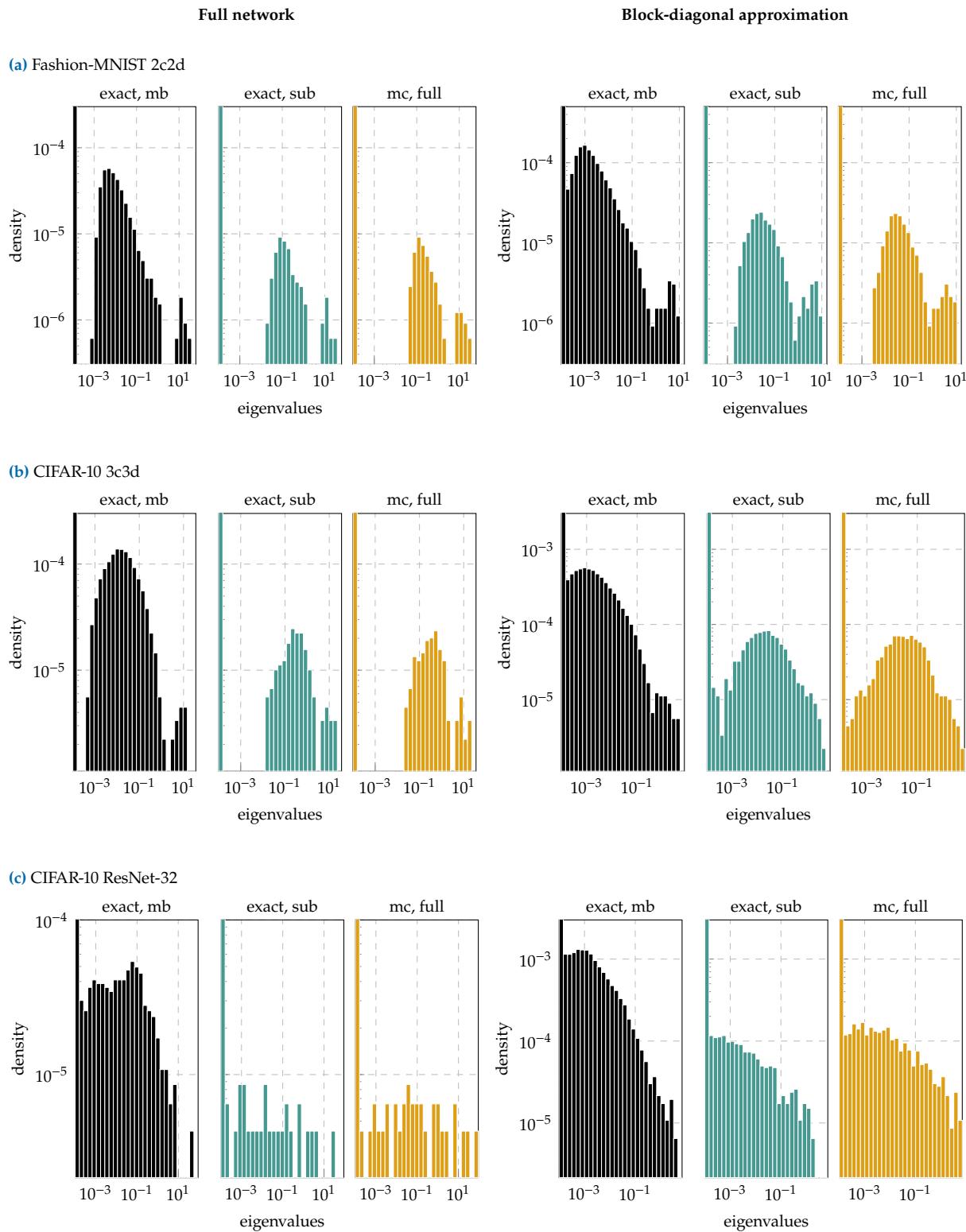
#### D.2.1 Performance Evaluation

##### Hardware Details

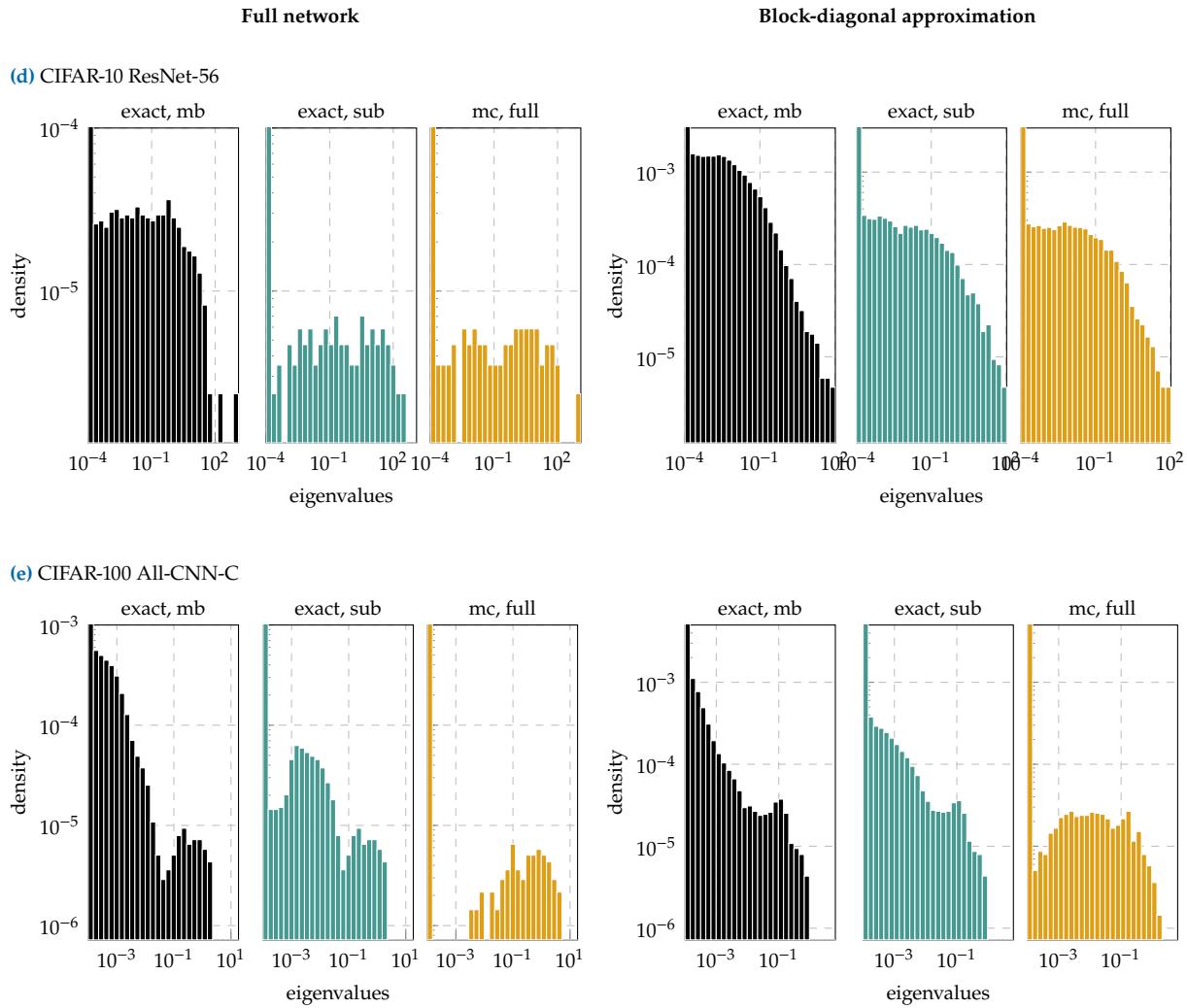
Results were generated on a workstation with an Intel Core i7-8700K CPU (32 GB) and one NVIDIA GeForce RTX 2080 Ti GPU (11 GB).

##### Note

ViViT's quantities are implemented through BackPACK, which is triggered by PyTorch's gradient computation. Consequently, they can only be computed together with PyTorch's mini-batch gradient.



**Figure D.1: GGN spectra of different architectures under ViViT's approximations:** Left and right columns contain results with the full network's GGN and a per-layer block-diagonal approximation, respectively.



**Figure D.1: GGN spectra of different architectures under ViViT’s approximations:** Left and right columns contain results with the full network’s GGN and a per-layer block-diagonal approximation, respectively.

### Architectures

We use untrained deep convolutional and residual networks from DeepOBS [146] and [74]. If a net has batch normalization layers, we set them to evaluation mode. Otherwise, the loss would not obey the sum structure of Equation (7.1). The batch normalization layers’ internal moving averages, required for evaluation mode, are initialized by performing five forward passes with the current mini-batch in training mode before.

In experiments with fixed mini-batches the batch sizes correspond to DeepOBS’ default value for training where possible (CIFAR-10:  $N = 128$ , Fashion-MNIST:  $N = 128$ ). The ResNets use a batch size of  $N = 128$ . On CIFAR-100 (trained with  $N = 256$ ), we reduce the batch size to  $N = 64$  to fit the exact computation on the full mini-batch, used as baseline, into memory. If the GGN approximation is evaluated on a subset of the mini-batch (**sub**),  $\lfloor N/8 \rfloor$  of the samples are used (as in [183]). The MC approximation is always evaluated with a single sample ( $M = 1$ ).

### Memory Performance (Critical Batch Sizes)

Two tasks are considered (see [Section 7.3.1](#)):

1. **Computing eigenvalues:** Compute the nontrivial eigenvalues  $\{\lambda_k \mid (\lambda_k, \tilde{e}_k) \in \mathbb{S}_+\}$ .
2. **Computing the top eigenpair:** Compute the top eigenpair  $(\lambda_1, e_1)$ .

We repeat the tasks above and vary the mini-batch size until the device runs out of memory. The largest mini-batch size that can be handled by our system is denoted as  $N_{\text{crit}}$ , the critical batch size. We determine this number by bisection on the interval  $[1; 32768]$ .

Subfigures (a) and (b) of [Figures D.2](#) to [D.11](#) present the results. As described in [Section 7.2.3](#), computing eigenvalues is more memory-efficient than computing eigenvectors and exhibits larger critical batch sizes. In line with the description in [Section 7.2.4](#), a block-diagonal approximation is usually more memory-efficient and results in a larger critical batch size. Curvature sub-sampling and MC approximation further increase the applicable batch sizes.

In summary, there always exists a combination of approximations which allows for critical batch sizes larger than the traditional size used for training (some architectures even permit exact computation). Different accuracy-cost trade-offs may be preferred, depending on the application and the computational budget. By the presented approximations, ViViT's representation is capable to adapt over a wide range.

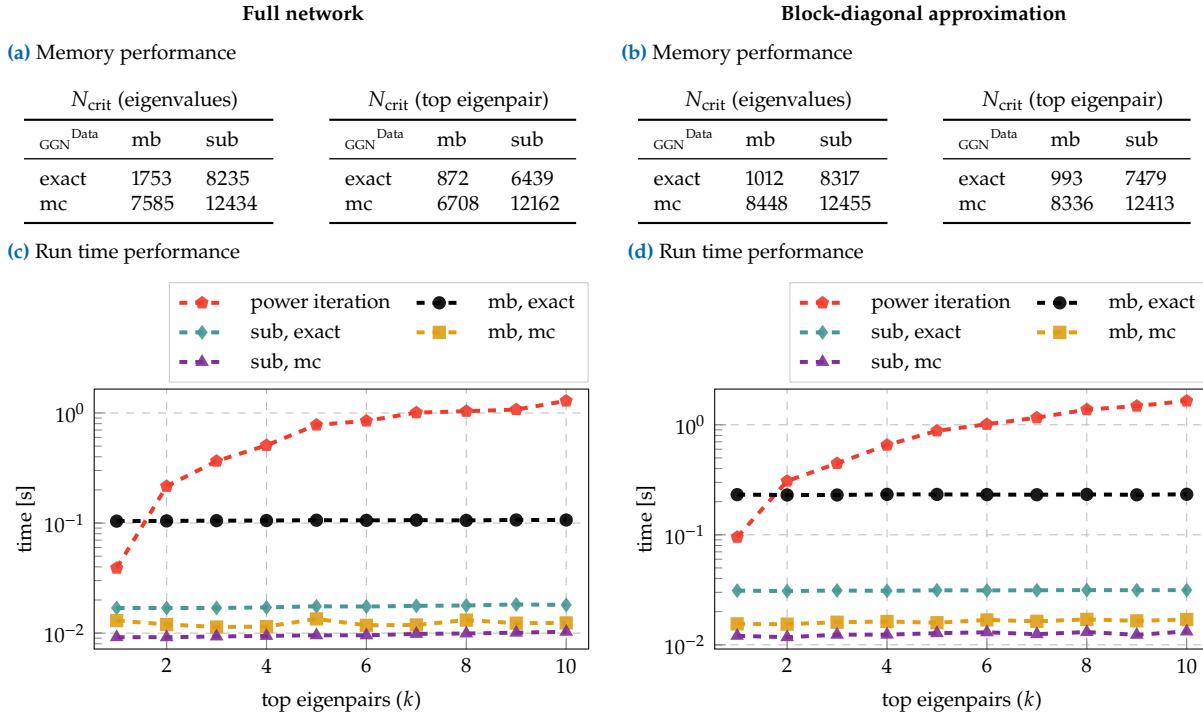
### Run Time Performance

Here, we consider the task of computing the  $k$  leading eigenvectors and eigenvalues of a matrix. ViViT's eigenpair computation is compared with a power iteration that computes eigenpairs iteratively via matrix-vector products. The power iteration baseline is based on the PyHessian library [\[177\]](#) and uses the same termination criterion (at most 100 matrix-vector products per eigenvalue; stop if the eigenvalue estimate's relative change is less than  $10^{-3}$ ). In contrast to PyHessian, we use a different data format and stack the computed eigenvectors. This reduces the number of `for`-loops in the orthonormalization step. We repeat each run time measurement 20 times and report the shortest execution time as result.

Subfigures (c) and (d) of [Figures D.2](#) to [D.11](#) show the results. For most architectures, our exact method outperforms the power iteration for  $k > 1$  and increases only marginally in run time as the number of requested eigenvectors grows. The proposed approximations share this property, and further reduce run time.

### Note On CIFAR-100 (Large $C$ )

For datasets with a large number of classes, like CIFAR-100 ( $C = 100$ ), computations with the exact GGN are costly. In particular, constructing the Gram matrix  $\tilde{G}$  has quadratic memory cost in  $C$ , and its eigendecomposition has cubic cost in time with  $C$  (see [Section 7.2.3](#)).



**Figure D.2:** GPU memory and run time performance for the 2c2d architecture on Fashion-MNIST. Left and right columns show results with the full network's GGN ( $D = 3,274,634$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .

As a result, the exact computation only works with batch sizes smaller than DeepOBS' default ( $N = 256$  for CIFAR-100, see subfigures (a) and (b) of Figures D.10 and D.11). For the GGN block-diagonal approximation, which fits into CPU memory for  $N = 64$ , the exact computation of top eigenpairs is slower than a power iteration and only becomes comparable if a large number of eigenpairs is requested, see Figure D.11d.

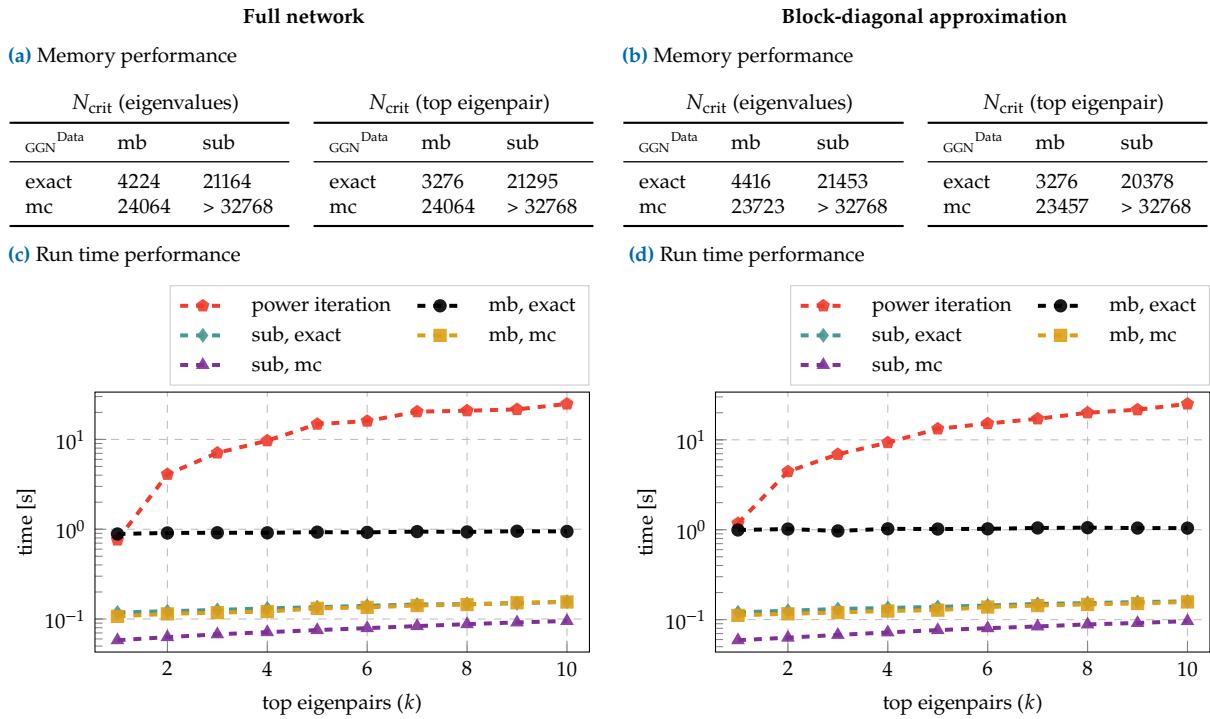
For such datasets, the approximations proposed in Section 7.2.4 are essential to reduce costs. The most effective approximation to eliminate the scaling with  $C$  is using an MC approximation. Figures D.10 and D.11 confirm that the approximate computations scale to batch sizes used for training and that computing eigenpairs takes less time than a power iteration.

### Computing Damped Newton Steps

A Newton step  $-(G + \delta I)^{-1}g$  with damping  $\delta > 0$  can be decomposed into updates along the eigenvectors of the GGN  $G$ ,

$$-(G + \delta I)^{-1}g = \sum_{k=1}^K \frac{-\gamma_k}{\lambda_k + \delta} e_k + \sum_{k=K+1}^D \frac{-\gamma_k}{\delta} e_k. \quad (\text{D.1})$$

It corresponds to a Newton update along nontrivial eigendirections that uses the first- and second-order directional derivatives described in Section 7.2.2 and a gradient descent step with learning rate  $1/\delta$  along trivial directions (with  $\lambda_k = 0$ ). In the following, we refer to the first summand of Equation (D.1) as Newton step. As described in Section 7.2.3,



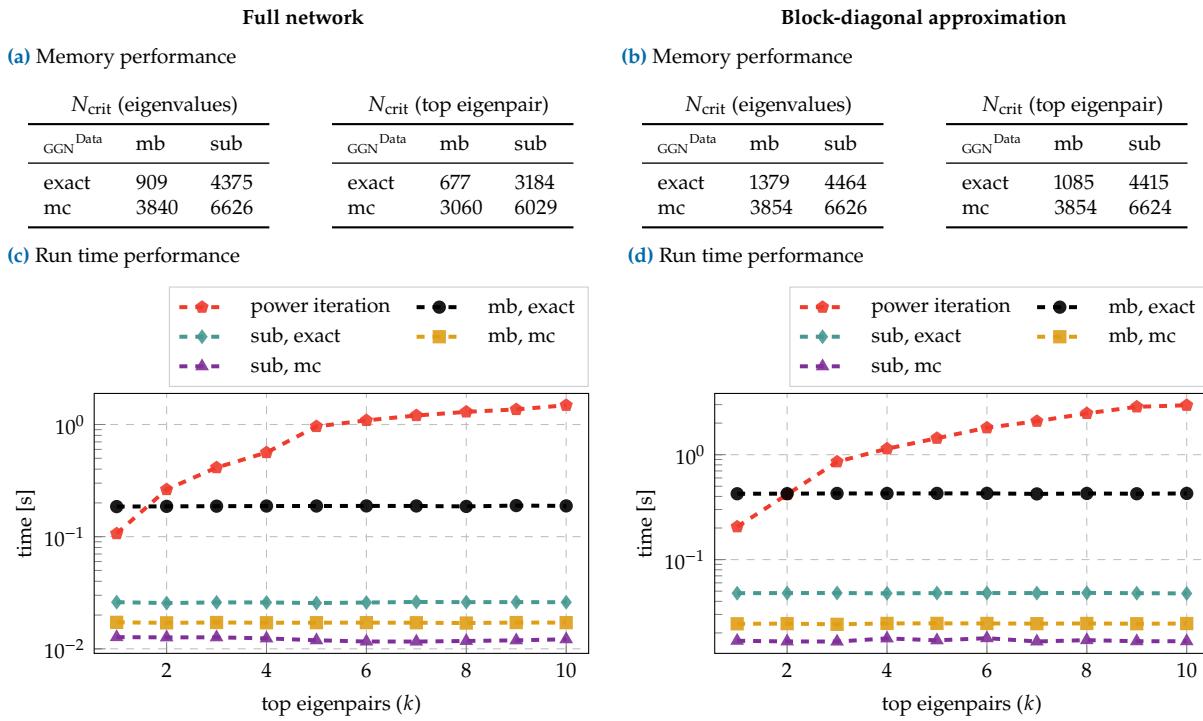
**Figure D.3:** CPU memory and run time performance for the 2c2d architecture on Fashion-MNIST. Left and right columns show results with the full network’s GGN ( $D = 3,274,634$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .

we can perform the weighted sum in the Gram matrix space, rather than the parameter space, by computing

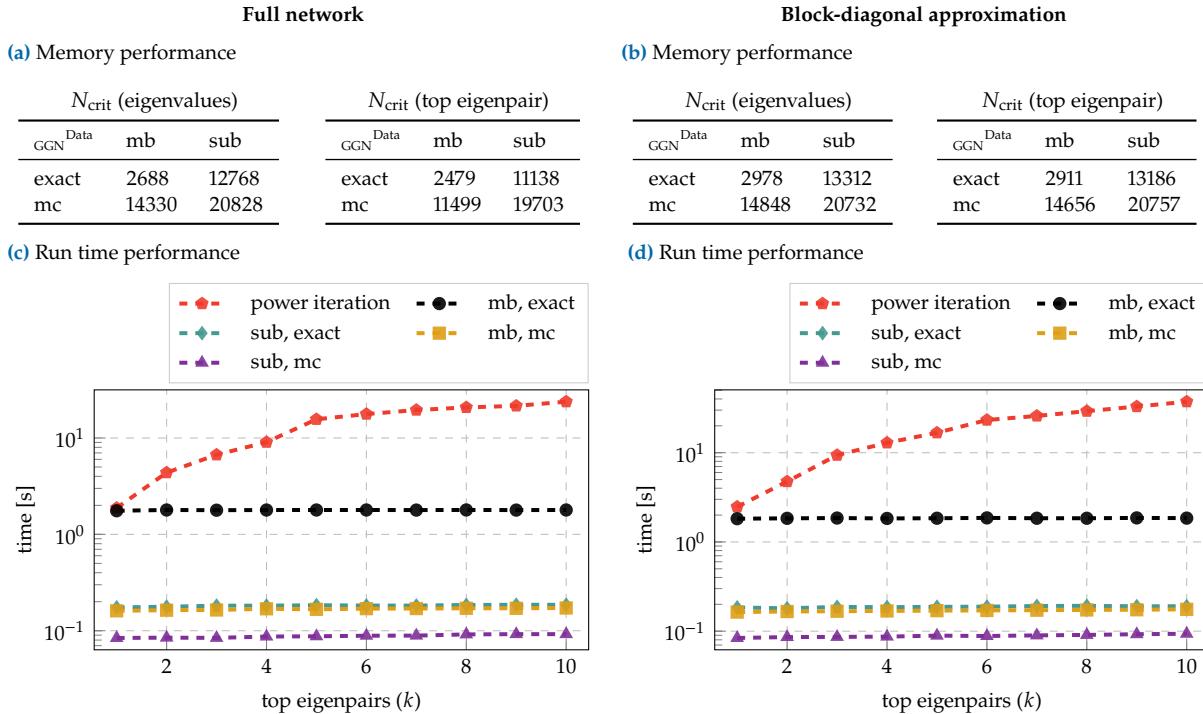
$$\sum_{k=1}^K \frac{-\gamma_k}{\lambda_k + \delta} e_k = \sum_{k=1}^K \frac{-\gamma_k}{\lambda_k + \delta} \frac{1}{\sqrt{\lambda_k}} V \tilde{e}_k = V \left( \sum_{k=1}^K \frac{-\gamma_k}{(\lambda_k + \delta)\sqrt{\lambda_k}} \tilde{e}_k \right).$$

This way, only a single vector needs to be transformed from Gram space into parameter space.

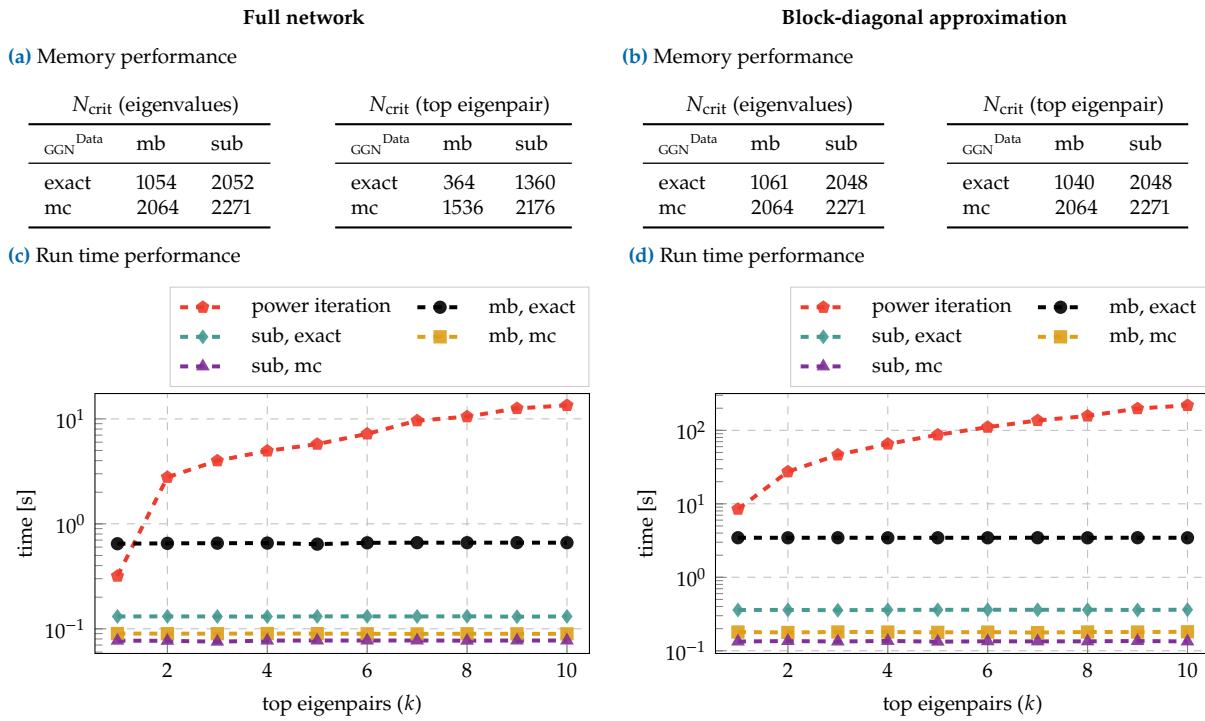
Table D.2 shows critical batch sizes for the Newton step computation (first term on the right side of Equation (D.1)), using Gram matrix eigenvalues larger than  $10^{-4}$  and constant damping  $\delta = 1$ . Second-order directional derivatives  $\lambda_k$  are evaluated on the same samples as the GGN eigenvectors, but we *always* use all mini-batch samples to compute the directional gradients  $\gamma_n$ . Using our approximations, the Newton step computation scales to batch sizes beyond the traditional sizes used for training.



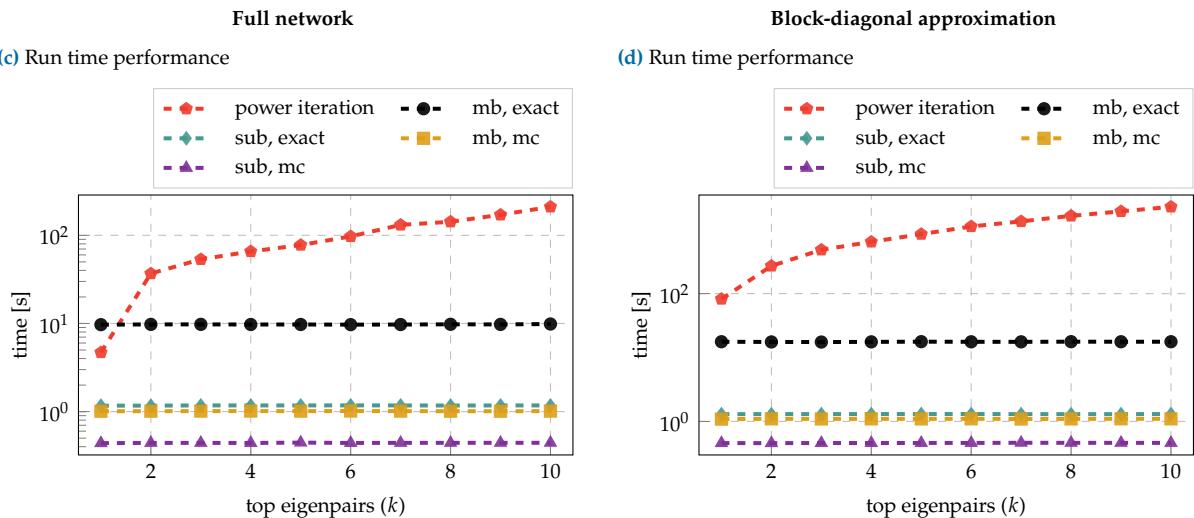
**Figure D.4:** GPU memory and run time performance for the 3c3d architecture on CIFAR-10. Left and right columns show results with the full network’s GGN ( $D = 895,210$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .



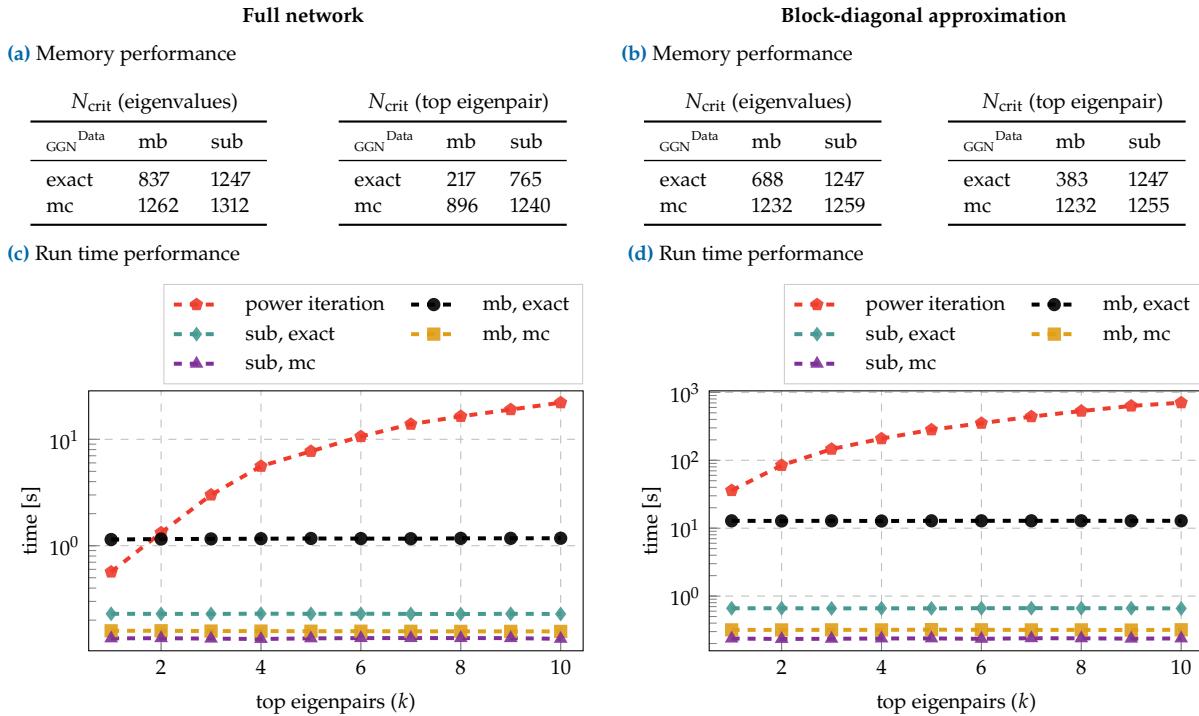
**Figure D.5:** CPU memory and run time performance for the 3c3d architecture on CIFAR-10. Left and right columns show results with the full network’s GGN ( $D = 895,210$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .



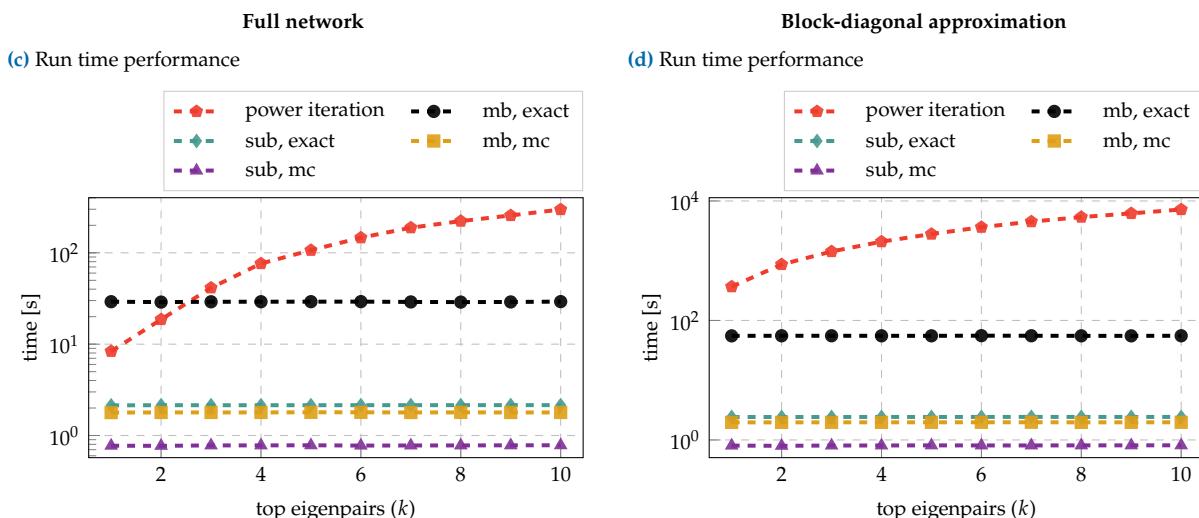
**Figure D.6:** GPU memory and run time performance for the ResNet-32 architecture on CIFAR-10. Left and right columns show results with the full network's GGN ( $D = 464,154$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .



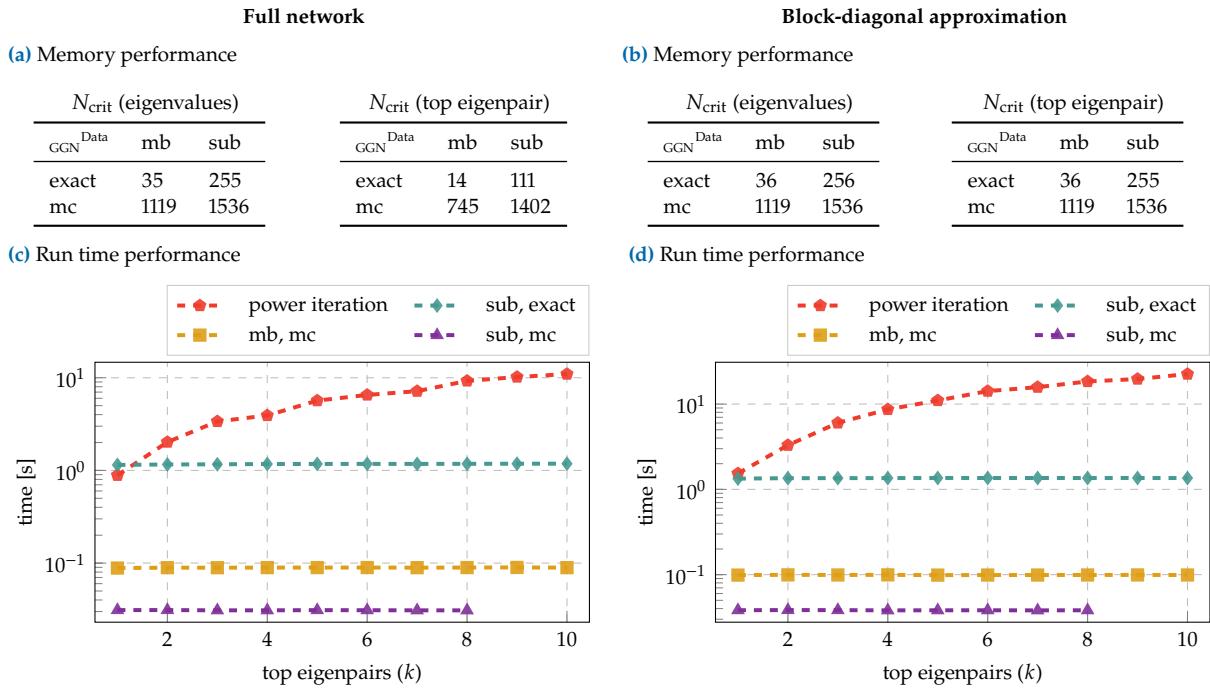
**Figure D.7:** CPU memory and run time performance for the ResNet-32 architecture on CIFAR-10. Left and right columns show results with the full network's GGN ( $D = 464,154$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .



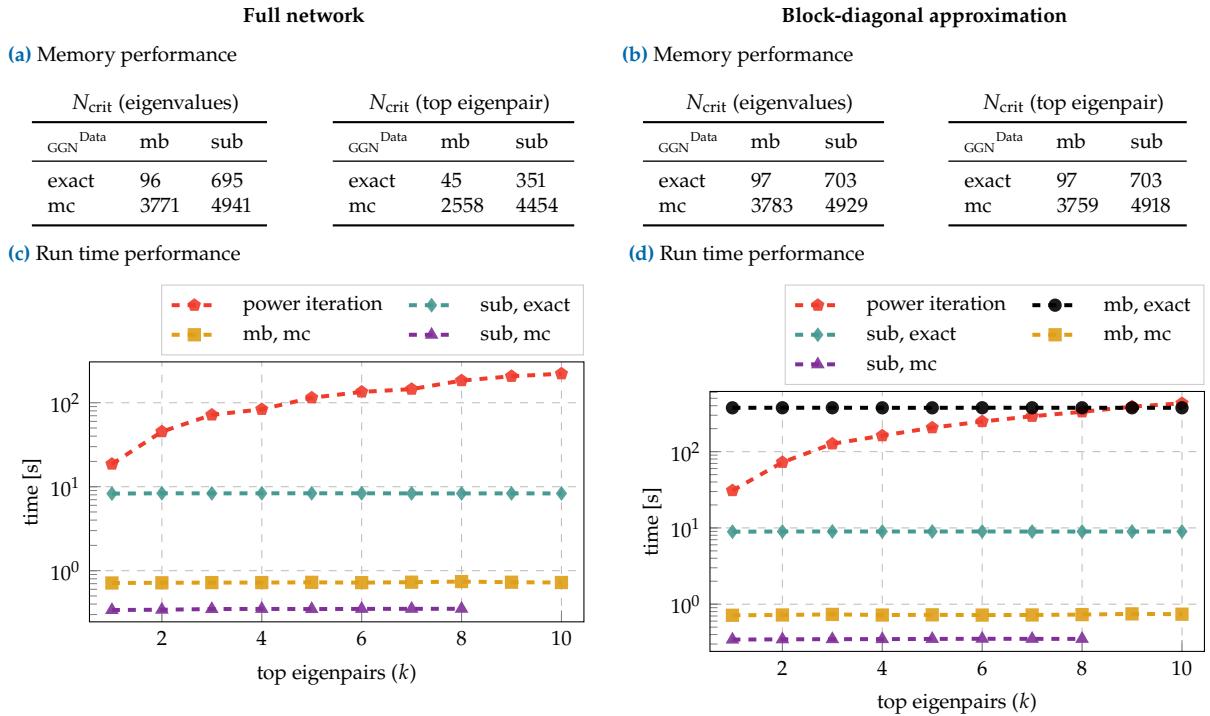
**Figure D.8:** GPU memory and run time performance for the ResNet-56 architecture on CIFAR-10. Left and right columns show results with the full network’s GGN ( $D = 853,018$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .



**Figure D.9:** CPU memory and run time performance for the ResNet-56 architecture on CIFAR-10. Left and right columns show results with the full network’s GGN ( $D = 853,018$ ,  $C = 10$ ) and a per-layer block-diagonal approximation, respectively. (a, b) (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 128$ .



**Figure D.10:** GPU memory and run time performance for the All-CNN-C architecture on CIFAR-100. Left and right columns show results with the full network’s GGN ( $D = 1,387,108$ ,  $C = 100$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 64$ .



**Figure D.11:** CPU memory and run time performance for the All-CNN-C architecture on CIFAR-100. Left and right columns show results with the full network’s GGN ( $D = 1,387,108$ ,  $C = 100$ ) and a per-layer block-diagonal approximation, respectively. (a),(b) Critical batch sizes  $N_{\text{crit}}$  for computing eigenvalues and the top eigenpair. (c),(d) Run time comparison with a power iteration for extracting the  $k$  leading eigenpairs using a mini-batch of size  $N = 64$ .

**Table D.2: Memory performance for computing damped Newton steps:** Left and right columns show the critical batch sizes with the full network's GGN and a per-layer block-diagonal approximation, respectively.

Fashion-MNIST 2c2d													
Full network			Block-diagonal approximation										
$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)			$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)				
GGN	Data	mb	sub	GGN	Data	mb	exact	202	487	GGN	Data	mb	sub
exact	66	159		exact	1107	1639	exact	68	159	exact	210	487	
mc	362	528		mc			mc	368	528	mc	1137	1643	

CIFAR-10 3c3d													
Full network			Block-diagonal approximation										
$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)			$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)				
GGN	Data	mb	sub	GGN	Data	mb	exact	667	2215	GGN	Data	mb	sub
exact	208	727		exact	3473	5632	exact	349	795	exact	1046	2423	
mc	1055	1816		mc			mc	1659	2112	mc	4997	6838	

CIFAR-10 ResNet-32													
Full network			Block-diagonal approximation										
$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)			$N_{\text{crit}}$ (GPU)							
GGN	Data	mb	sub	GGN	Data	mb	exact	1051	1851	GGN	Data	mb	sub
exact	344	1119		exact	2048	2208	exact	1205	1535	mc			

CIFAR-10 ResNet-56													
Full network			Block-diagonal approximation										
$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)			$N_{\text{crit}}$ (GPU)							
GGN	Data	mb	sub	GGN	Data	mb	exact	767	1165	GGN	Data	mb	sub
exact	209	640		exact	1232	1255	mc	687	890	mc			

CIFAR-100 All-CNN-C													
Full network			Block-diagonal approximation										
$N_{\text{crit}}$ (GPU)			$N_{\text{crit}}$ (CPU)			$N_{\text{crit}}$ (GPU)							
GGN	Data	mb	sub	GGN	Data	mb	exact	43	309	GGN	Data	mb	sub
exact	13	87		exact	2015	2865	exact	35	135	exact	95	504	
mc	640	959		mc			mc	1079	1536	mc	3360	3920	

## D.2.2 Training of Neural Networks

### Procedure

We train the following DeepOBS [146] architectures with SGD and Adam: 3c3d on CIFAR-10, 2c2d on Fashion-MNIST and All-CNN-C on CIFAR-100; all are equipped with cross-entropy loss. To ensure successful training, we use the hyperparameters from [38] (see Table D.3).

We also train a residual network ResNet-32 [68] with cross-entropy loss on CIFAR-10 with both SGD and Adam. For this, we use a batch size of 128 and train for 180 epochs. Momentum for SGD was fixed to 0.9, and Adam uses the default parameters ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ ). For both optimizers, the learning rate was determined via grid search. Following [146], we use a log-equidistant grid from  $10^{-5}$  to  $10^2$  and 36 grid points. As performance metric, the best test accuracy during training (evaluated once every epoch) is used.

### Results

The results for the hyperparameter grid search are reported in Table D.3. The training metrics training/test loss/accuracy for all eight test problems are shown in Figures D.12 and D.13.

Problem	SGD	Adam	Batch size	Epochs
Fashion-MNIST 2c2d	$\eta \approx 2.07 \cdot 10^{-2}$	$\eta \approx 1.27 \cdot 10^{-4}$	$N = 128$	100
CIFAR-10 3c3d	$\eta \approx 3.79 \cdot 10^{-3}$	$\eta \approx 2.98 \cdot 10^{-4}$	$N = 128$	100
CIFAR-10 ResNet-32	$\eta \approx 6.31 \cdot 10^{-2}$	$\eta \approx 2.51 \cdot 10^{-3}$	$N = 128$	180
CIFAR-100 All-CNN-C	$\eta \approx 4.83 \cdot 10^{-1}$	$\eta \approx 6.95 \cdot 10^{-4}$	$N = 256$	350

**Table D.3: Hyperparameter settings for training runs.** For both SGD and Adam, we report their learning rates  $\eta$  (taken from the baselines in [38] or, for ResNet-32, determined via grid search). Momentum for SGD is fixed to 0.9. Adam uses the default parameters  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ . We also report the batch size used for training and the number of training epochs.

## D.2.3 GGN versus Hessian

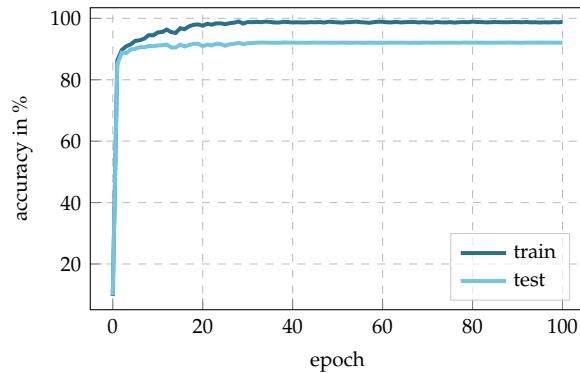
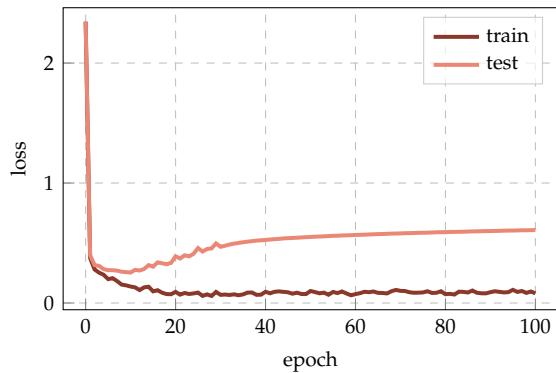
### Checkpoints

During training of the neural networks (see Appendix D.2.2), we store a copy of the model (*i.e.* the network's current parameters) at specific checkpoints. This grid defines the temporal resolution for all subsequent computations. Since training progresses much faster in the early training stages, we use a log-grid with 100 grid points between 1 and the number of training epochs and shift this grid by  $-1$ .

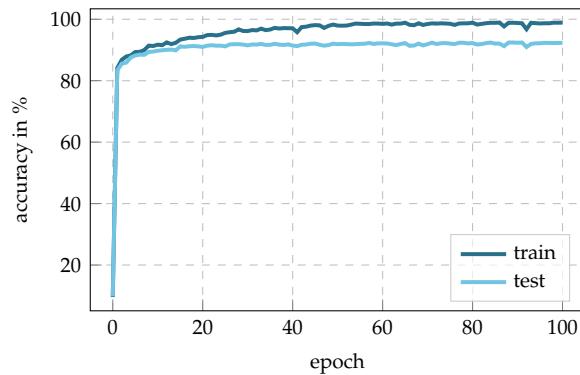
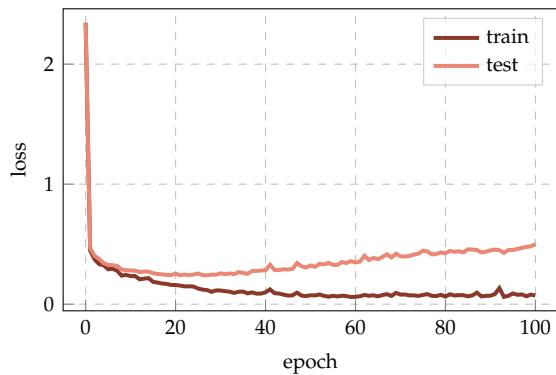
### Overlap

Recall from Section 7.3.2: for the set of orthonormal eigenvectors  $\{e_c^U\}_{c=1}^C$  to the  $C$  largest eigenvalues of some symmetric matrix  $U$ , let  $P^U = (e_1^U, \dots, e_C^U)(e_1^U, \dots, e_C^U)^\top$ . As in [65], the overlap between two subspaces

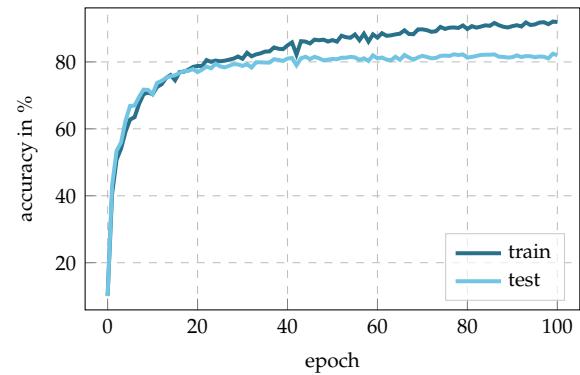
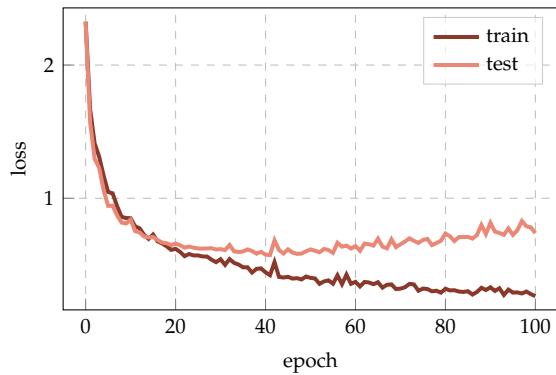
(a) Fashion-MNIST 2c2d SGD



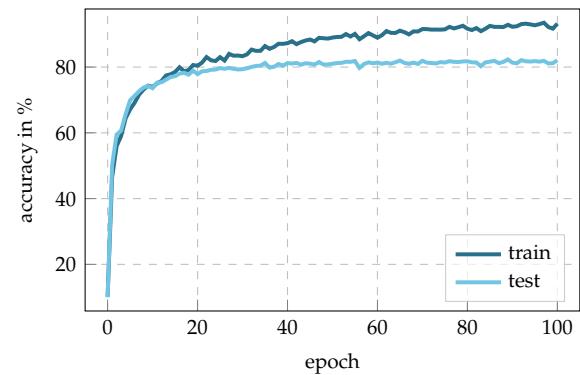
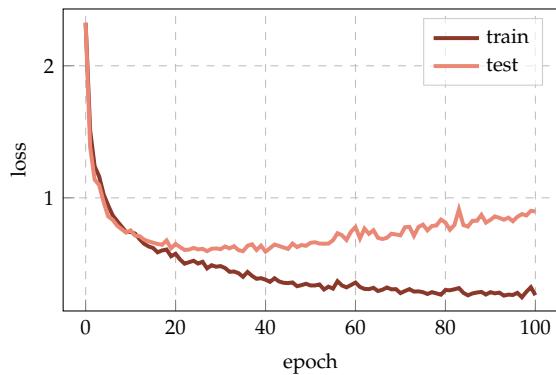
(b) Fashion-MNIST 2c2d Adam



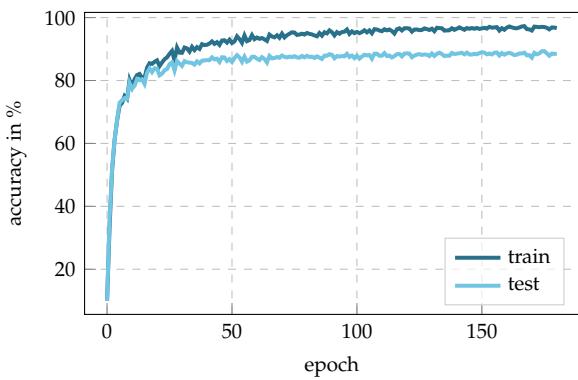
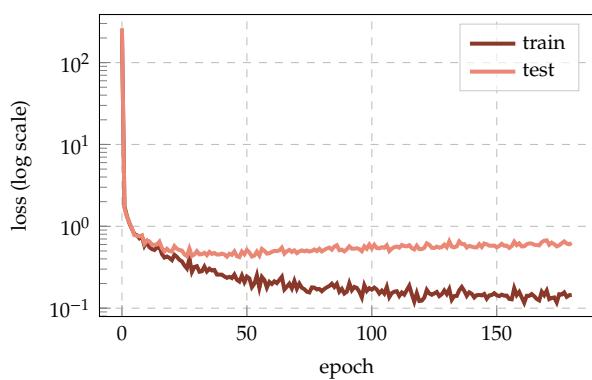
(c) CIFAR-10 3c3d SGD



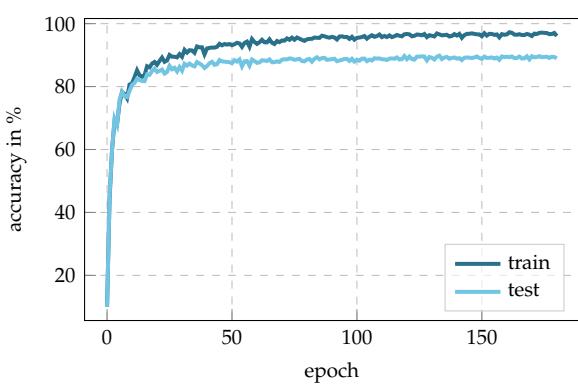
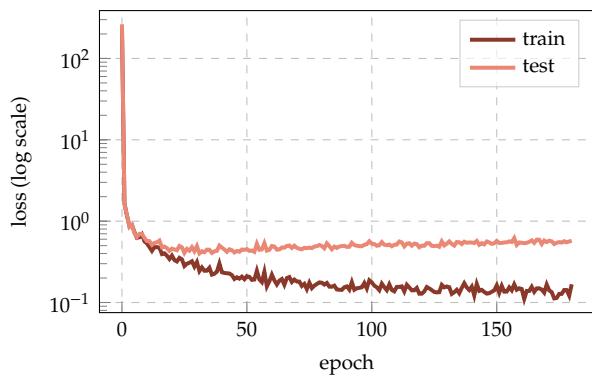
(d) CIFAR-10 3c3d Adam

**Figure D.12:** Training metrics (1). Training/test loss/accuracy for all test problems.

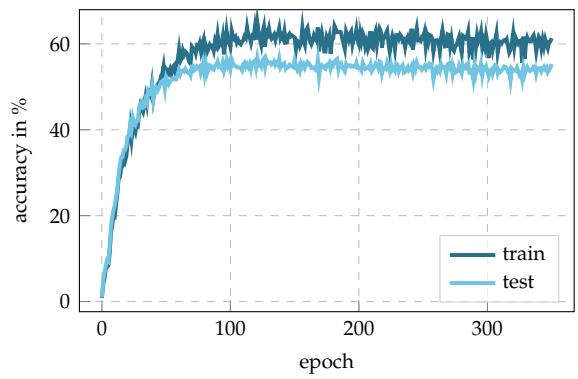
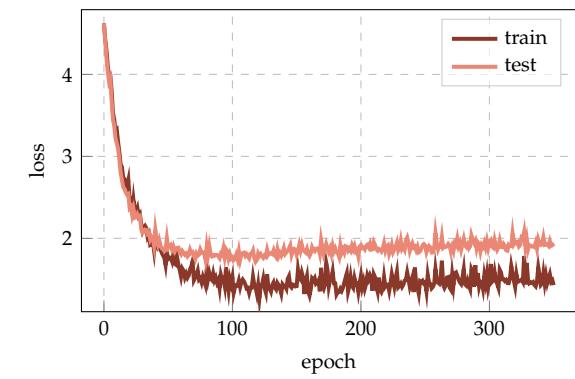
(a) CIFAR-10 ResNet-32 SGD



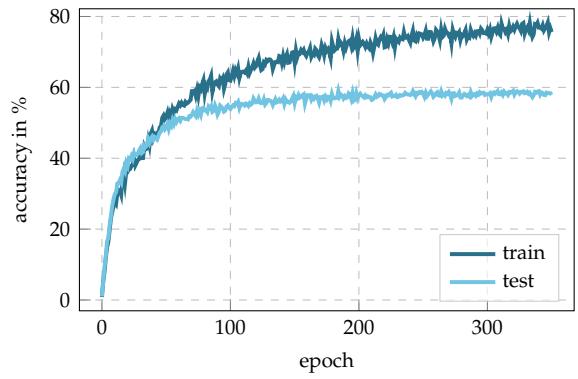
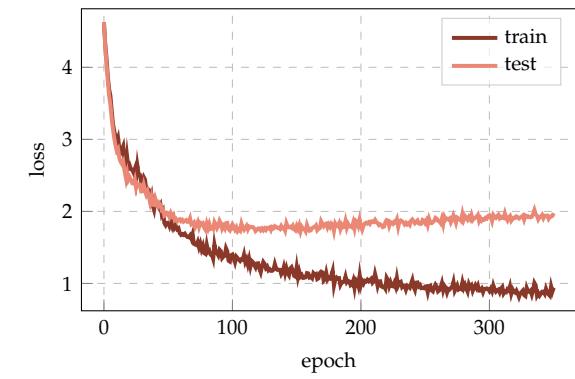
(b) CIFAR-10 ResNet-32 Adam



(c) CIFAR-100 All-CNN-C SGD



(d) CIFAR-100 All-CNN-C Adam

**Figure D.13: Training metrics (2).** Training/test loss/accuracy for all test problems.

$\mathcal{E}^U = \text{span}(e_1^U, \dots, e_C^U)$  and  $\mathcal{E}^V = \text{span}(e_1^V, \dots, e_C^V)$  of the matrices  $U$  and  $V$  is defined by

$$\text{overlap}(\mathcal{E}^U, \mathcal{E}^V) = \frac{\text{Tr}(P^U P^V)}{\sqrt{\text{Tr}(P^U) \text{Tr}(P^V)}} \in [0, 1].$$

The overlap can be computed efficiently by using the trace's cyclic property: it holds  $\text{Tr}(P^U P^V) = \text{Tr}(W^\top W)$  with  $W = (e_1^U, \dots, e_C^U)^\top (e_1^V, \dots, e_C^V) \in \mathbb{R}^{C \times C}$ . Note that this is a small  $C \times C$  matrix, whereas  $P^U, P^V \in \mathbb{R}^{D \times D}$ . Since

$$\begin{aligned} \text{Tr}(P^U) &= \text{Tr}((e_1^U, \dots, e_C^U)(e_1^U, \dots, e_C^U)^\top) \\ &= \text{Tr}((e_1^U, \dots, e_C^U)^\top (e_1^U, \dots, e_C^U)) \quad (\text{Cyclic property of trace}) \\ &= \text{Tr}(I_C) \quad (\text{Orthonormality of the eigenvectors}) \\ &= C \end{aligned}$$

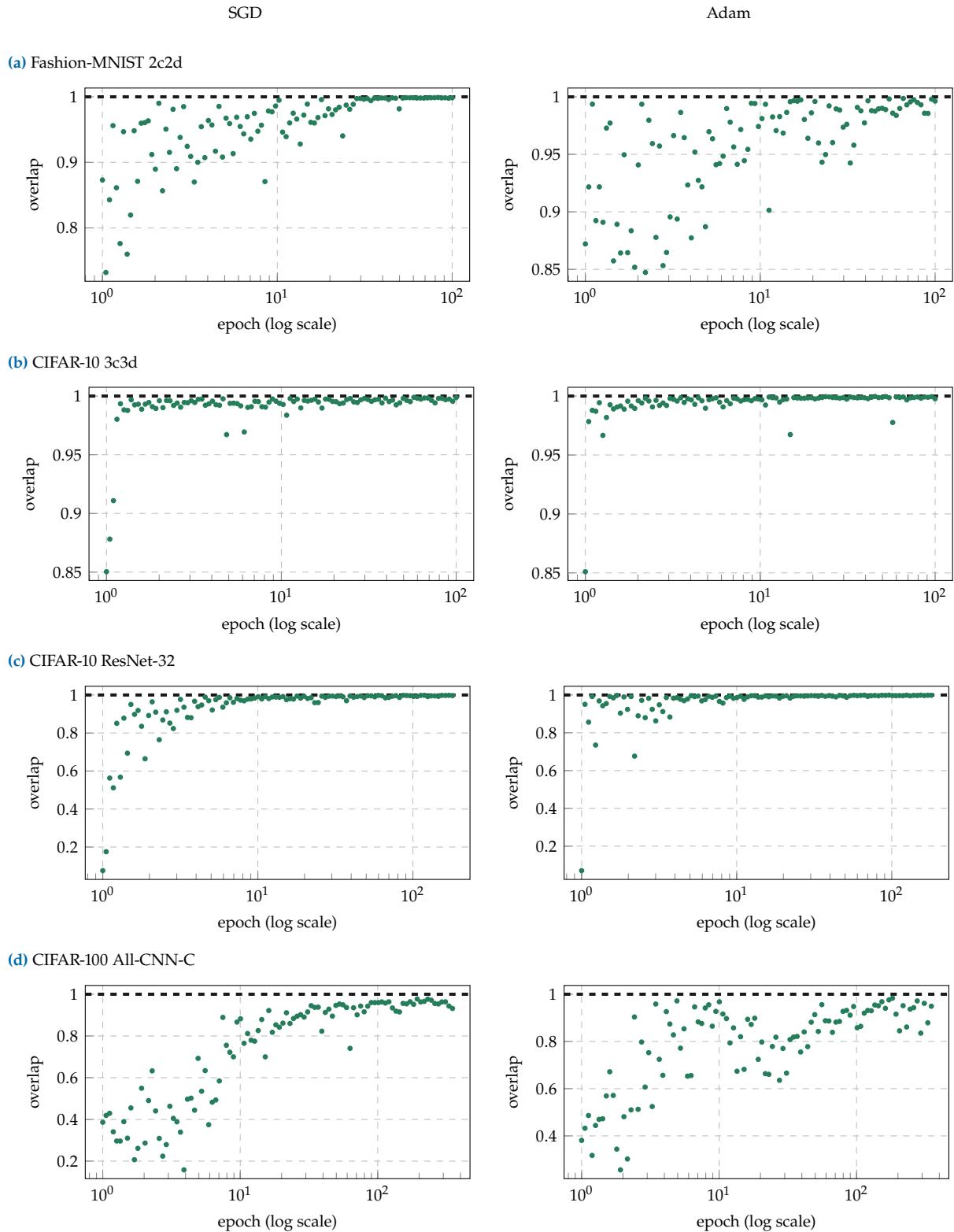
(and analogous  $\text{Tr}(P^V) = C$ ), the denominator simplifies to  $C$ .

### Procedure

For each checkpoint, we compute the top- $C$  eigenvalues and associated eigenvectors of the full-batch GGN and Hessian (*i.e.* GGN and Hessian are both evaluated on the entire training set) using an iterative matrix-free approach. We then compute the overlap between the top- $C$  eigenspaces as described above. The eigenspaces (*i.e.* the top- $C$  eigenvalues and associated eigenvectors) are stored on disk such that they can be used as a reference by subsequent experiments.

### Results

The results for all test problems are presented in [Figure D.14](#). Except for a short phase at the beginning of the optimization procedure (note the log scale for the epoch-axis), a strong agreement (note the different limits for the overlap-axis) between the top- $C$  eigenspaces is observed. We make similar observations for all test problems, yet to a slightly lesser extent for CIFAR-100. A possible explanation for this would be that the 100-dimensional eigenspaces differ in the eigenvectors associated with relatively small curvature. The corresponding eigenvalues already transition into the bulk of the spectrum, where the "sharpness of separation" decreases. However, since all directions are equally weighted in the overlap, overall slightly lower values are obtained.



**Figure D.14: Full-batch GGN versus full-batch Hessian.** Overlap between the top-C eigenspaces of the full-batch GGN and full-batch Hessian during training for all test problems.

### D.2.4 Eigenspace Under Noise & Approximations

#### Procedure (1)

We use the checkpoints and the definition of overlaps between eigenspaces from [Appendix D.2.3](#). For the approximation of the GGN, we consider the cases listed in [Table D.4](#).

**Table D.4:** Considered cases for approximation of the eigenspace: We use a different set of cases for the approximation of the GGN's full-batch eigenspace depending on the test problem. For the test problems with  $C = 10$ , we use  $M = 1$  MC-sample, for the CIFAR-100 All-CNN-C test problem ( $C = 100$ ), we use  $M = 10$  MC-samples in order to reduce the computational costs by the same factor.

Problem	Cases
Fashion-MNIST 2c2d CIFAR-10 3c3d and CIFAR-10 ResNet-32	<b>mb, exact</b> with mini-batch sizes $N \in \{2, 8, 32, 128\}$ <b>mb, mc</b> with $N = 128$ and $M = 1$ MC-sample <b>sub, exact</b> using 16 samples from the mini-batch <b>sub, mc</b> using 16 samples from the mini-batch and $M = 1$ MC-sample
CIFAR-100 All-CNN-C	<b>mb, exact</b> with mini-batch sizes $N \in \{2, 8, 32, 128\}$ <b>mb, mc</b> with $N = 128$ and $M = 10$ MC-samples <b>sub, exact</b> using 16 samples from the mini-batch <b>sub, mc</b> using 16 samples from the mini-batch and $M = 10$ MC-samples

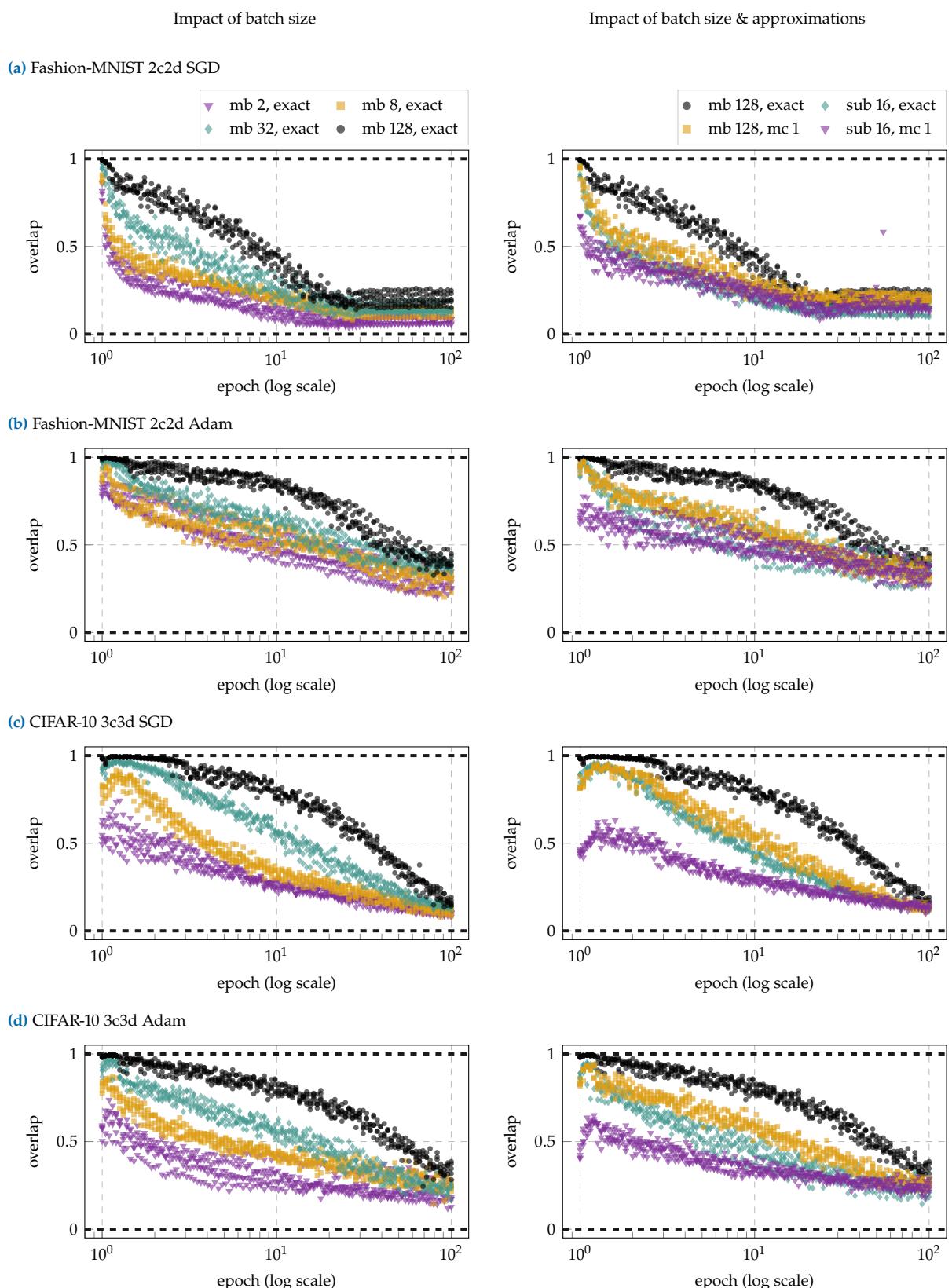
For every checkpoint and case, we compute the top- $C$  eigenvectors of the respective approximation to the GGN. The eigenvectors are either computed directly using ViViT (by transforming eigenvectors of the Gram matrix into parameter space, see [Section 7.2.1](#)) or, if not applicable (because memory requirements exceed  $N_{\text{crit}}$ , see [Section 7.3.1](#)), using an iterative matrix-free approach. The overlap is computed in reference to the GGN's full-batch top- $C$  eigenspace (see [Appendix D.2.3](#)). We extract 5 mini-batches from the training data and repeat the above procedure for each mini-batch (*i.e.* we obtain 5 overlap measurements for every checkpoint and case). The same 5 mini-batches are used over all checkpoints and cases.

#### Results (1)

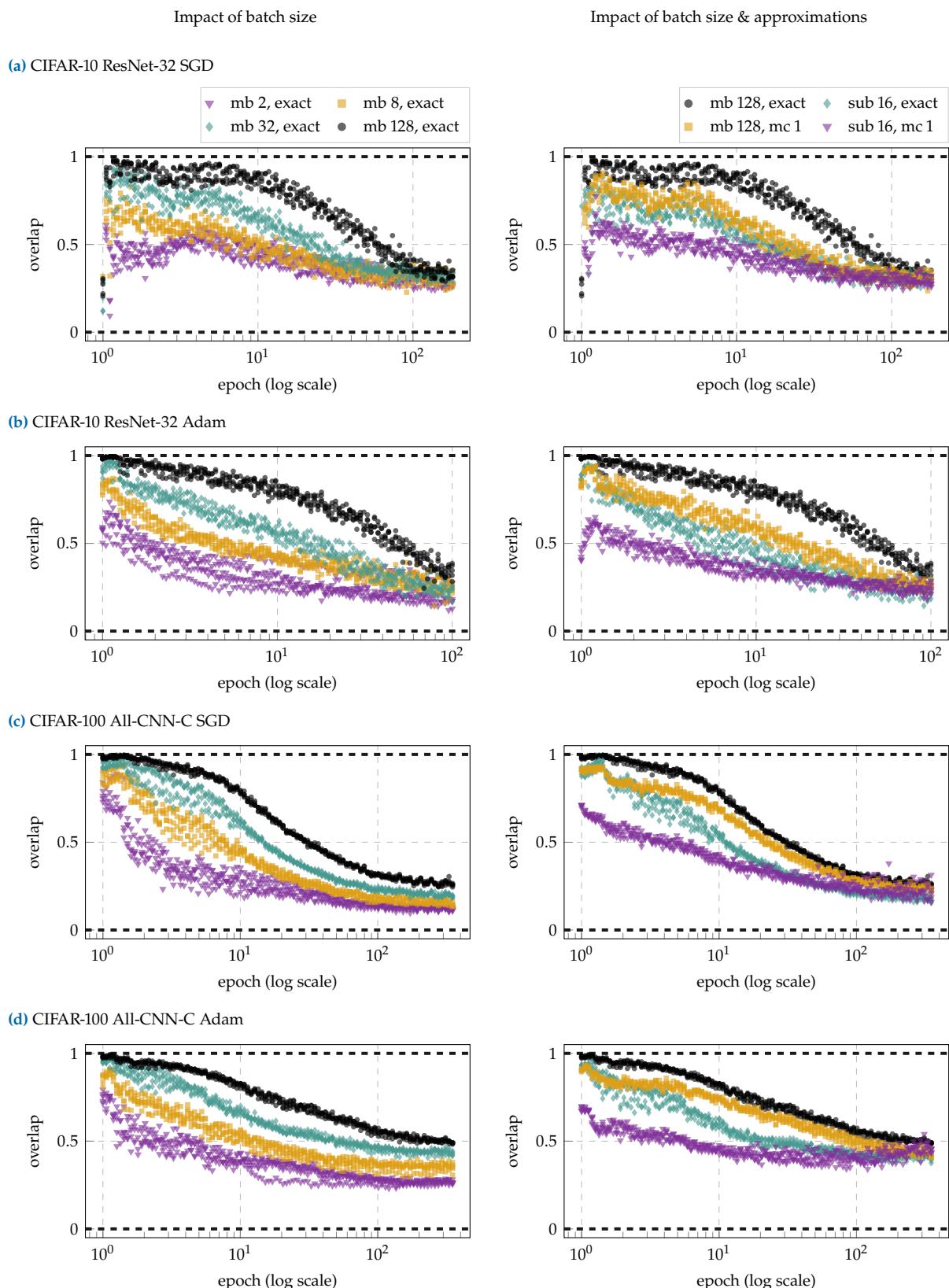
The results can be found in [Figure D.15](#) and [D.16](#). All test problems show the same characteristics: with decreasing computational effort, the approximation carries less and less structure of its full-batch counterpart, as indicated by dropping overlaps. In addition, for a fixed approximation method, a decrease in approximation quality can be observed over the course of training.

#### Procedure (2)

Since ViViT's GGN approximations using curvature sub-sampling and/or the MC approximation (the cases **mb**, **mc** as well as **sub, exact** and **sub, mc** in [Table D.4](#)) are based on the *mini*-batch GGN, we cannot expect them to perform better than this baseline. We thus repeat the analysis from above but use the *mini*-batch GGN with batch-size  $N = 128$  as ground truth instead of the full-batch GGN. Of course, the *mini*-batch reference top- $C$  eigenspace is always evaluated on the same *mini*-batch as the approximation.



**Figure D.15: ViViT versus full-batch GGN (1).** Overlap between the top-C eigenspaces of different GGN approximations and the full-batch GGN during training for all test problems. Each approximation is evaluated on 5 different mini-batches.



**Figure D.16:** ViViT versus full-batch GGN (2). Overlap between the top- $C$  eigenspaces of different GGN approximations and the full-batch GGN during training for all test problems. Each approximation is evaluated on 5 different mini-batches.

## Results (2)

[Figure D.17](#) shows the results. Over large parts of the optimization (note the log scale for the epoch-axis), the MC approximation seems to be better suited than curvature sub-sampling (which has comparable computational cost). For the CIFAR-100 All-CNN-C problem, the MC approximation stands out particularly early from the other approximations and consistently yields higher overlaps with the mini-batch GGN.

### D.2.5 Curvature Under Noise & Approximations

GGN and Hessian are predominantly used to locally approximate the loss by a quadratic model  $q$  (see [Equation \(7.6\)](#)). Even if the curvature's eigenspace is completely preserved in spite of the approximations, they can still alter the curvature *magnitude* along the eigenvectors.

#### Procedure

[Table D.5](#) shows the cases considered in this experiment.

**Table D.5: Considered cases for approximation of curvature:** We use a different set of cases for the approximation of the GGN depending on the test problem. For the test problems with  $C = 10$ , we use  $M = 1$  MC-sample, for the CIFAR-100 All-CNN-C test problem ( $C = 100$ ), we use  $M = 10$  MC-samples in order to reduce the computational costs by the same factor.

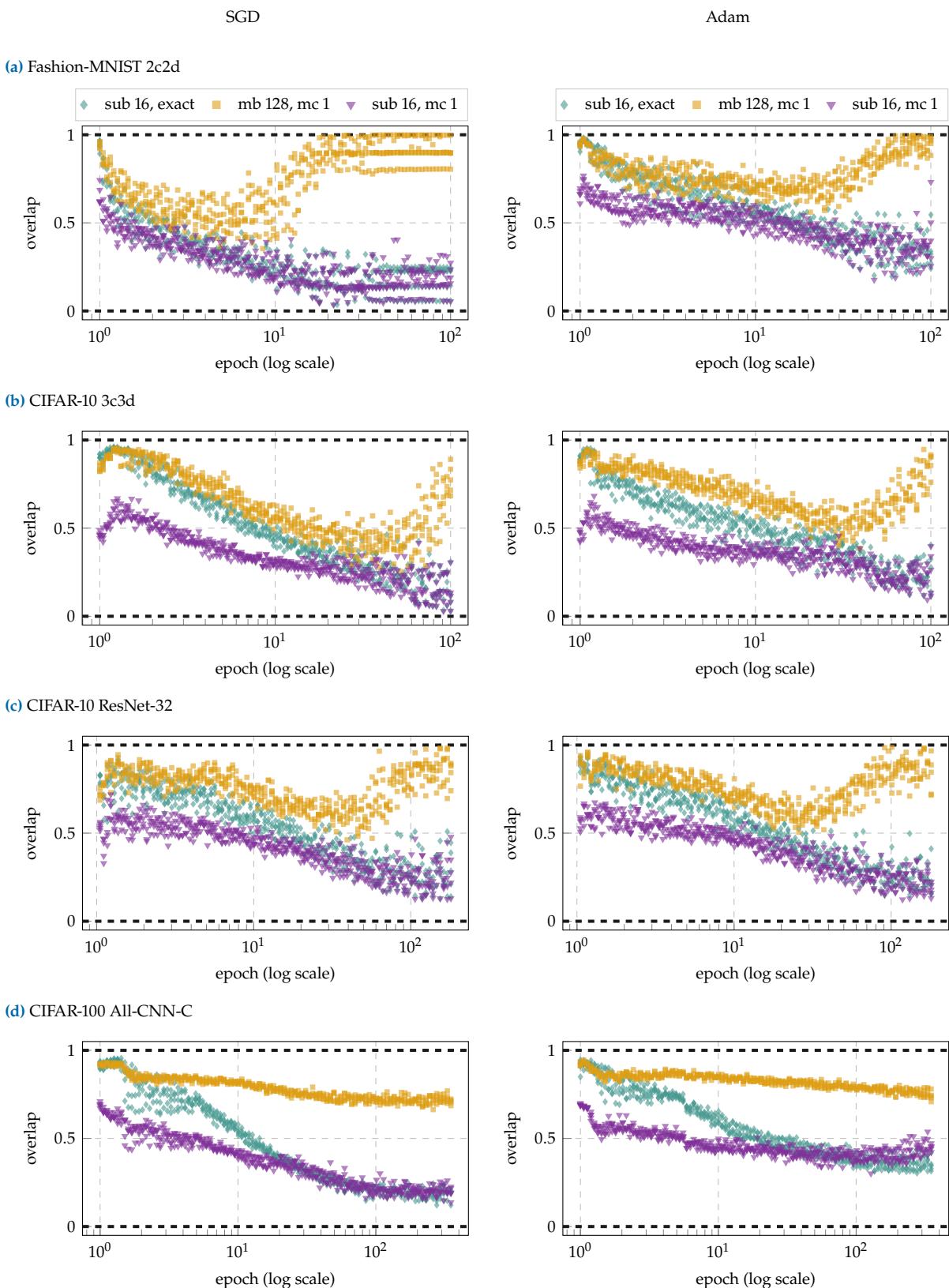
Problem	Cases
Fashion-MNIST 2c2d CIFAR-10 3c3d and CIFAR-10 ResNet-32	<b>mb, exact</b> with mini-batch size $N = 128$ <b>mb, mc</b> with $N = 128$ and $M = 1$ MC-sample <b>sub, exact</b> using 16 samples from the mini-batch <b>sub, mc</b> using 16 samples from the mini-batch and $M = 1$ MC-sample
CIFAR-100 All-CNN-C	<b>mb, exact</b> with mini-batch size $N = 128$ <b>mb, mc</b> with $N = 128$ and $M = 10$ MC-samples <b>sub, exact</b> using 16 samples from the mini-batch <b>sub, mc</b> using 16 samples from the mini-batch and $M = 10$ MC-samples

Due to the large computational effort for evaluating the full-batch directional derivatives, a subset of the checkpoints from [Appendix D.2.3](#) is used for two problems: we use every second checkpoint for CIFAR-10 ResNet-32 and every forth checkpoint for CIFAR-100 All-CNN-C.

For each checkpoint and case, we compute the top- $C$  eigenvectors  $\{\mathbf{e}_k\}_{k=1}^C$  of the GGN approximation  $\mathbf{G}^{(\text{ap})}$  either with ViViT or using an iterative matrix-free approach (as in [Appendix D.2.4](#)). The second-order directional derivative of the corresponding quadratic model along direction  $\mathbf{e}_k$  is then given by  $\lambda_k^{(\text{ap})} = \mathbf{e}_k^\top \mathbf{G}^{(\text{ap})} \mathbf{e}_k$  (see [Equation \(7.7\)](#)). As a reference, we compute the full-batch GGN  $\mathbf{G}^{(\text{fb})}$  and the resulting directional derivatives along the same eigenvectors  $\{\mathbf{e}_k\}_{k=1}^C$ , i.e.  $\lambda_k^{(\text{fb})} = \mathbf{e}_k^\top \mathbf{G}^{(\text{fb})} \mathbf{e}_k$ . The average (over all  $C$  directions) relative error is given by

$$\epsilon = \frac{1}{C} \sum_{k=1}^C \frac{|\lambda_k^{(\text{ap})} - \lambda_k^{(\text{fb})}|}{\lambda_k^{(\text{fb})}}.$$

The procedure above is repeated on 3 mini-batches from the training data (i.e. we obtain 3 average relative errors for every checkpoint and case) –



**Figure D.17: ViViT versus mini-batch GGN.** Overlap between the top- $C$  eigenspaces of different GGN approximations and the mini-batch GGN during training for all test problems. Each approximation is evaluated on 5 different mini-batches.

except for the CIFAR-100 All-CNN-C test problem, where we perform only a single run to keep the computational effort manageable.

## Results

Figure D.18 shows the results. We observe similar results as in Appendix D.2.4: with increasing computational effort, the approximated directional derivatives become more precise and the overall approximation quality decreases over the course of the optimization. For the ResNet-32 architecture, the average errors are particularly large.

### D.2.6 Directional Derivatives

#### Procedure

We use the checkpoints from Appendix D.2.3. For every checkpoint, we compute the top- $C$  eigenvectors of the mini-batch GGN ( $N = 128$ ) using an iterative matrix-free method. We also compute the mini-batch gradient. The first- and second-order directional derivatives of the resulting quadratic model (Equation (7.6)) are given by Equation (7.8).

We use these directional derivatives  $\{\gamma_{n,k}\}_{n=1,k=1}^{N,C}$ ,  $\{\lambda_{n,k}\}_{n=1,k=1}^{N,C}$  to compute signal-to-noise ratios (SNRs) along the top- $C$  eigenvectors. The curvature SNR along direction  $e_k$  is given by the squared sample mean divided by the empirical variance of the samples  $\{\lambda_{n,k}\}_{n=1}^N$ , i.e.

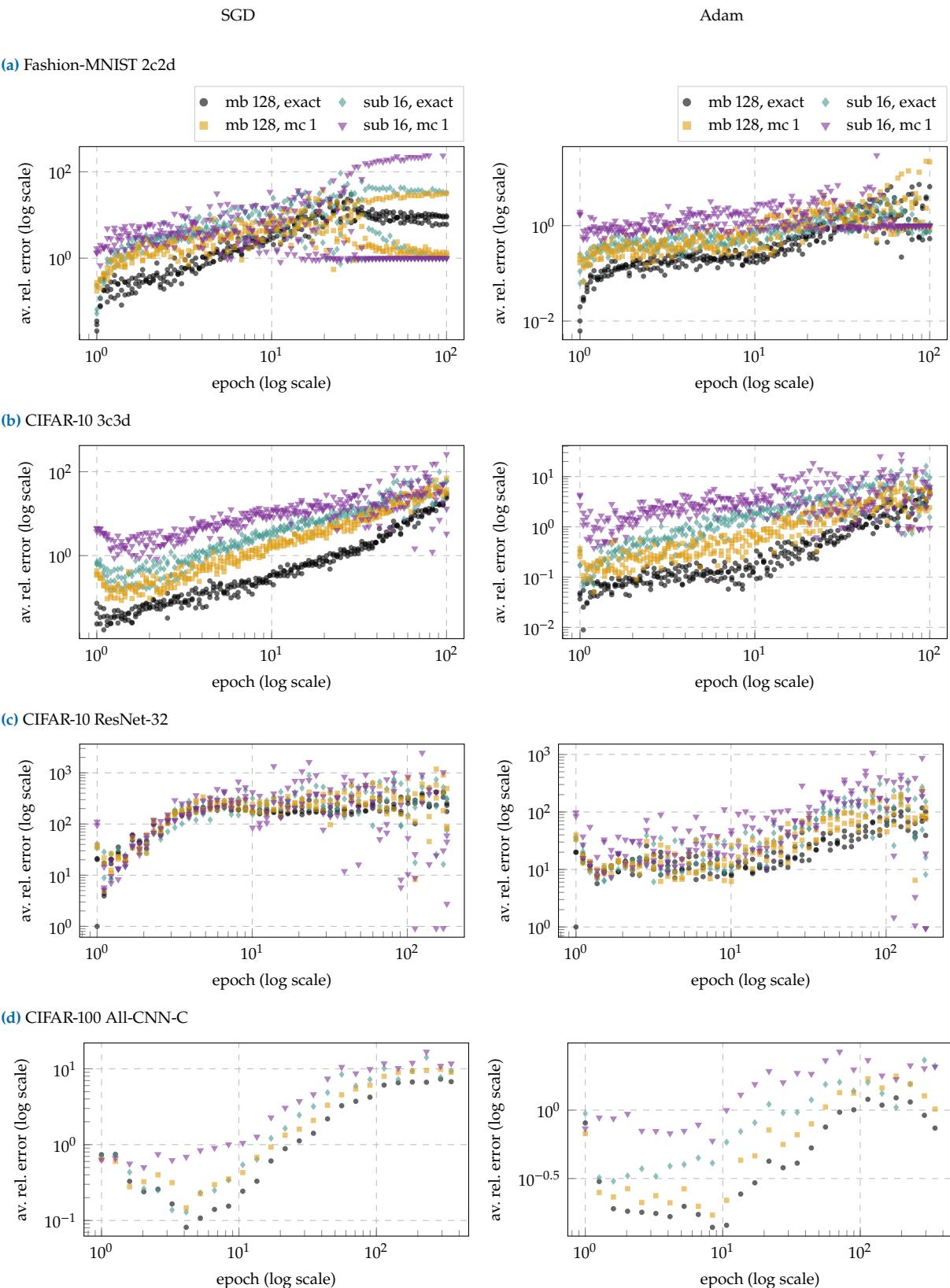
$$\text{SNR} = \frac{\lambda_k^2}{\sum_{n=1}^N (\lambda_{n,k} - \lambda_k)^2} \quad \text{where} \quad \lambda_k = \frac{1}{N} \sum_{n=1}^N \lambda_{n,k}.$$

(and similarly for  $\{\gamma_{n,k}\}_{n=1}^N$ ).

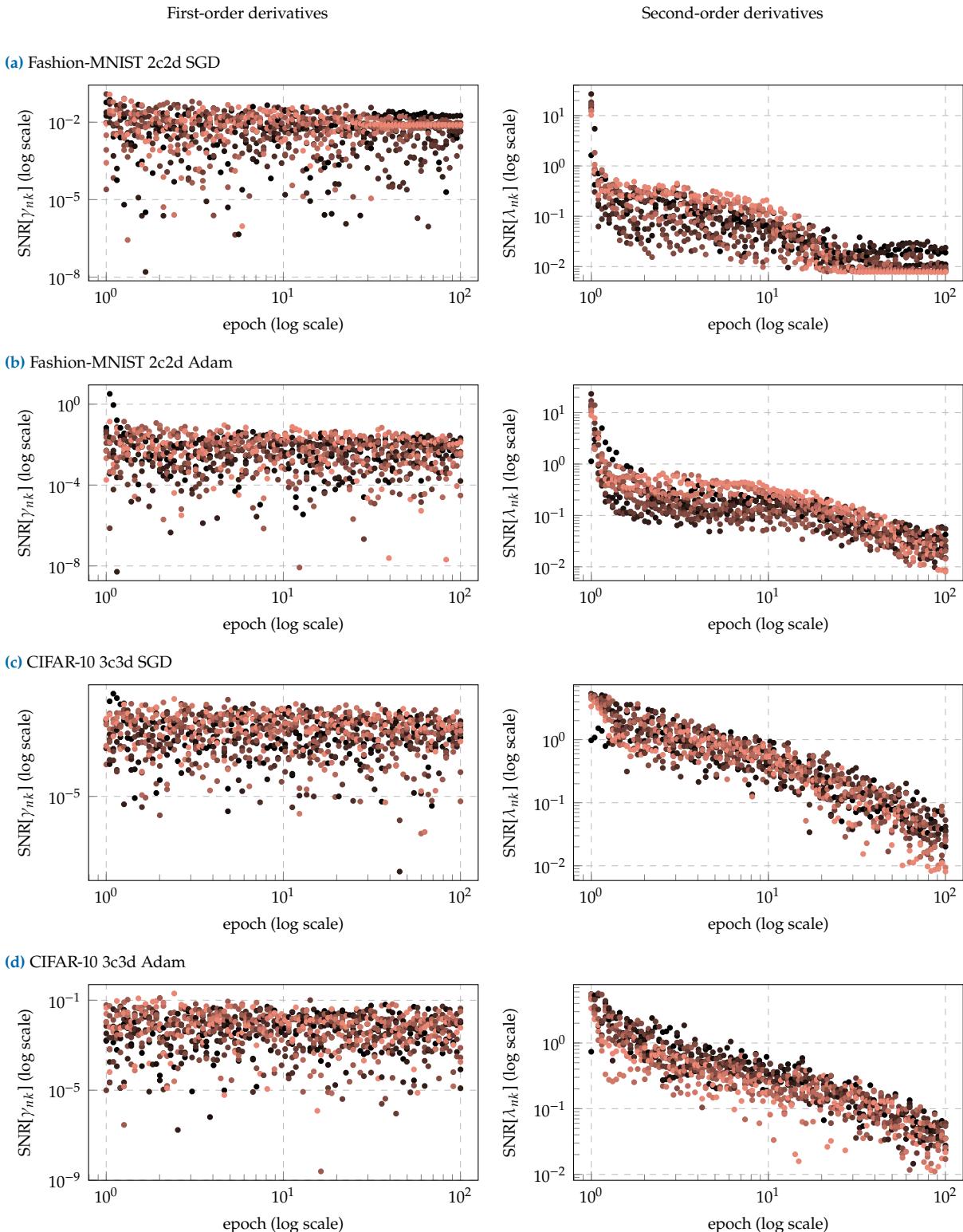
#### Results

Figures D.19 and D.20 show the results. These plots show the SNRs in  $C$  distinct colors that generated from linear interpolation in the RGB color space from black (●) to light red (○). At each checkpoint, the colors are assigned based on the *order* of the respective directional curvature  $\lambda_k$ : the SNR that belongs to the direction with the smallest curvature is shown in black and the SNR that belongs to the direction with the largest curvature is shown in light red. The color thus encodes only the order of the top- $C$  directional curvatures – *not* their magnitude. We use this color encoding to reveal potential correlations between SNR and curvature.

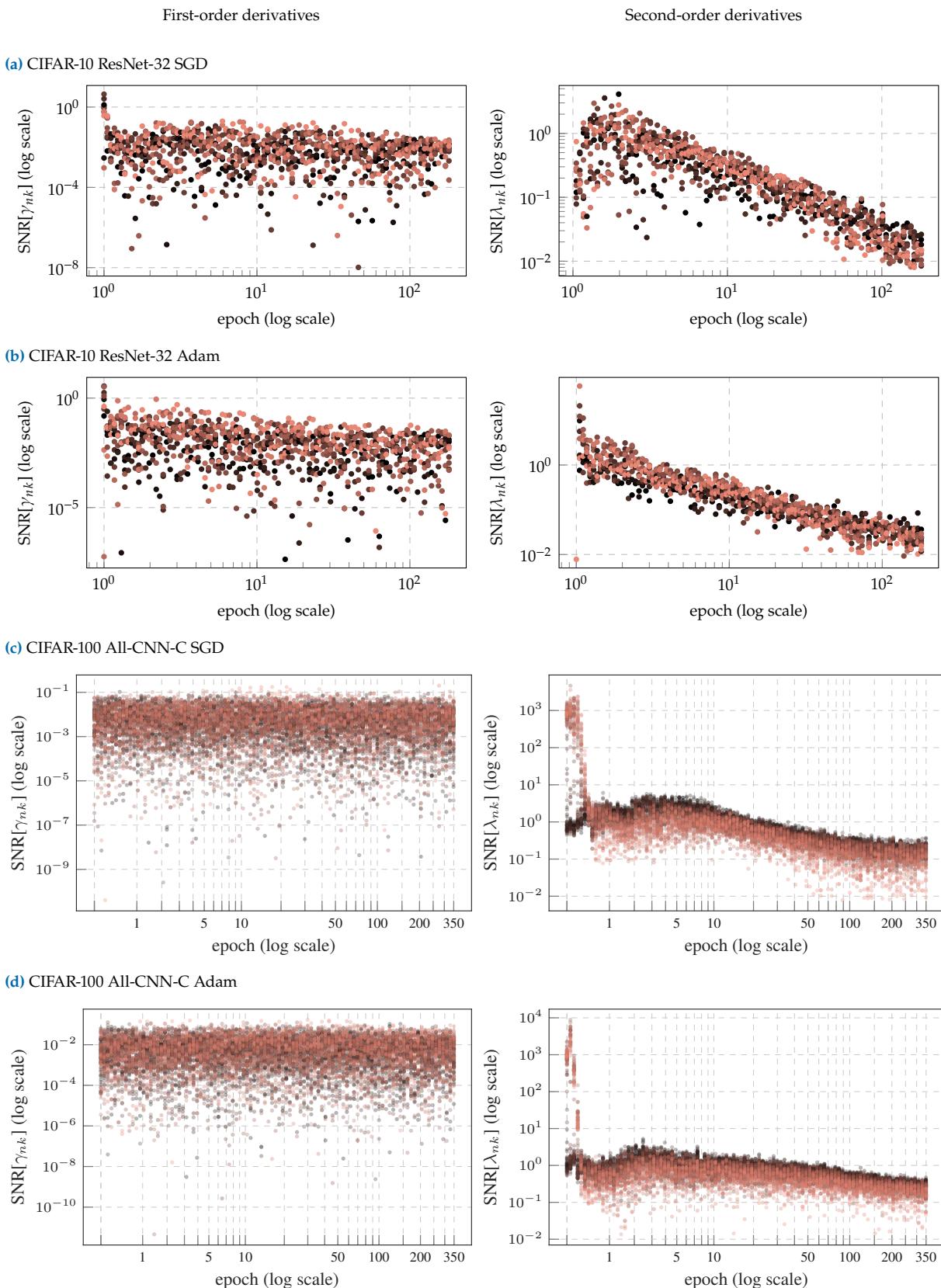
We find that the gradient SNR along the top- $C$  eigenvectors is consistently small (in comparison to the curvature SNR) and remains roughly on the same level during the optimization. The curvature signal decreases as training proceeds. The SNRs along the top- $C$  eigendirections do not appear to show any significant correlation with the corresponding curvatures. Only for the CIFAR-100 test problems we can suspect a correlation between strong curvature and small curvature SNR.



**Figure D.18: ViViT's versus full-batch quadratic model.** Comparison between approximations to the quadratic model and the full-batch model in terms of the average relative error for the directional curvature during training for all test problems.



**Figure D.19: Directional SNRs (1).** SNR along each of the mini-batch GGN's top-C eigenvectors during training for all test problems. At fixed epoch, the SNR for the most curved direction is shown in ● and for the least curved direction in ○.



**Figure D.20: Directional SNRs (2).** SNR along each of the mini-batch GGN's top-C eigenvectors during training for all test problems. At fixed epoch, the SNR for the most curved direction is shown in ● and for the least curved direction in ○.

## D.3 Implementation Details

### Layer View of Backpropagation

Consider a single layer  $f_{\theta^{(l)}}$  that transforms inputs  $\mathbf{z}_n^{(l-1)} \in \mathbb{R}^{h^{(l-1)}}$  into outputs  $\mathbf{z}_n^{(l)} \in \mathbb{R}^{h^{(l)}}$  by means of a parameter  $\theta^{(l)} \in \mathbb{R}^{d^{(l)}}$ . During backpropagation for  $V$ , the layer receives vectors  $\mathbf{s}_{n,c}^{(l)} = (\mathbf{J}_{\mathbf{z}_n^{(l)}} f_n)^\top \mathbf{s}_{n,c}$  from the previous stage (recall  $\nabla_{f_n}^2 \ell_n = \sum_{c=1}^C \mathbf{s}_{n,c} \mathbf{s}_{n,c}^\top$ ). Parameter contributions  $\mathbf{v}_{nc}^{(l)}$  to  $V$  are obtained by application of its Jacobian,

$$\begin{aligned} \mathbf{v}_{nc}^{(l)} &= (\mathbf{J}_{\theta^{(l)}} f_n)^\top \mathbf{s}_{n,c} \\ &= (\mathbf{J}_{\theta^{(l)}} \mathbf{z}_n^{(l)})^\top (\mathbf{J}_{\mathbf{z}_n^{(l)}} f_n)^\top \mathbf{s}_{n,c} \quad (\text{Chain rule}) \\ &= (\mathbf{J}_{\theta^{(l)}} \mathbf{z}_n^{(l)})^\top \mathbf{s}_{n,c}^{(l)}. \quad (\text{Definition of } \mathbf{s}_{n,c}^{(l)}) \end{aligned} \quad (\text{D.2})$$

Consequently,  $\theta^{(l)}$ 's contribution to  $V$ , denoted by  $\mathbf{V}^{(l)} \in \mathbb{R}^{d^{(l)} \times NC}$ , is

$$\begin{aligned} \mathbf{V}^{(l)} &= \frac{1}{\sqrt{N}} \begin{pmatrix} \mathbf{v}_{1,1}^{(l)} & \mathbf{v}_{1,2}^{(l)} & \dots & \mathbf{v}_{N,C}^{(l)} \end{pmatrix} \\ \text{with } \mathbf{v}_{nc}^{(l)} &= (\mathbf{J}_{\theta^{(l)}} f_n)^\top \mathbf{s}_{n,c}. \end{aligned} \quad (\text{D.3})$$

### D.3.1 Optimized Gram Matrix for Linear Layers

Our goal is to efficiently extract  $\theta^{(l)}$ 's contribution to the Gram matrix,

$$\tilde{\mathbf{G}}^{(l)} = \mathbf{V}^{(l)\top} \mathbf{V}^{(l)} \in \mathbb{R}^{NC \times NC}. \quad (\text{D.4})$$

#### Gram Matrix via Expanding $\mathbf{V}^{(l)}$

One way to construct  $\mathbf{G}^{(l)}$  is to first expand  $\mathbf{V}^{(l)}$  (Equation (D.3)) via the Jacobian  $\mathbf{J}_{\theta^{(l)}} \mathbf{z}_n^{(l)}$ , then contract it (Equation (D.4)). This can be a memory bottleneck for large linear layers which are common in many architectures close to the network output. However if only the Gram matrix rather than  $V$  is required, structure in the Jacobian can be used to construct  $\tilde{\mathbf{G}}^{(l)}$  without expanding  $\mathbf{V}^{(l)}$  and thus reduce this overhead.

#### Optimization for Linear Layers

Now, let  $f_{\theta^{(l)}}$  be a linear layer with weights  $\mathbf{W}^{(l)} \in \mathbb{R}^{h^{(l)} \times h^{(l-1)}}$ , i.e.  $\theta^{(l)} = \text{vec}(\mathbf{W}^{(l)}) \in \mathbb{R}^{d^{(l)} = h^{(l)} h^{(l-1)}}$  with column stacking as flattening,

$$f_{\theta^{(l)}}^{(l)} : \quad \mathbf{z}_n^{(l)} = \mathbf{W}^{(l)} \mathbf{z}_n^{(l-1)}.$$

The Jacobian is

$$\mathbf{J}_{\theta^{(l)}} \mathbf{z}_n^{(l)} = \mathbf{z}_n^{(l-1)\top} \otimes \mathbf{I}_{h^{(l)}}, \quad (\text{D.5})$$

and can be used to compute Gram matrix entries without expanding  $\mathbf{V}^{(l)}$ ,

$$\begin{aligned} & \left[ \tilde{\mathbf{G}}^{(l)} \right]_{(n,c),(n',c')} \\ &= \mathbf{v}_{n,c}^{(l)\top} \mathbf{v}_{n',c'}^{(l)} \end{aligned} \quad (\text{Equation (D.4)})$$

$$= \mathbf{s}_{n,c}^{(l)\top} \left( \mathbf{J}_{\theta^{(l)}} \mathbf{z}_n^{(l)} \right) \left( \mathbf{J}_{\theta^{(l)}} \mathbf{z}_{n'}^{(l)} \right)^{\top} \mathbf{s}_{n',c'}^{(l)} \quad (\text{Equation (D.5)})$$

$$= \mathbf{s}_{n,c}^{(l)\top} \left( \mathbf{z}_n^{(l-1)\top} \otimes \mathbf{I}_{h^{(l)}} \right) \left( \mathbf{z}_{n'}^{(l-1)\top} \otimes \mathbf{I}_{h^{(l)}} \right)^{\top} \mathbf{s}_{n',c'}^{(l)} \quad (\text{Equation (D.2)})$$

$$= \mathbf{s}_{n,c}^{(l)\top} \left( \mathbf{z}_n^{(l-1)\top} \mathbf{z}_{n'}^{(l-1)} \otimes \mathbf{I}_{h^{(l)}} \right) \mathbf{s}_{n',c'}^{(l)} \quad (\text{Equation (D.2)})$$

$$= \mathbf{z}_n^{(l-1)\top} \mathbf{z}_{n'}^{(l-1)} \mathbf{s}_{n,c}^{(l)\top} \mathbf{I}_{h^{(l)}} \mathbf{s}_{n',c'}^{(l)} \quad (\mathbf{z}_n^{(l-1)\top} \mathbf{z}_{n'}^{(l-1)} \in \mathbb{R})$$

$$= \left( \mathbf{z}_n^{(l-1)\top} \mathbf{z}_{n'}^{(l-1)} \right) \left( \mathbf{s}_{n,c}^{(l)\top} \mathbf{s}_{n',c'}^{(l)} \right).$$

We see that the Gram matrix is built from two Gram matrices based on  $\{\mathbf{z}_n^{(l-1)}\}_{n=1}^N$  and  $\{\mathbf{s}_{n,c}^{(l)}\}_{n=1,c=1}^{N,C}$ , that require  $\mathcal{O}(N^2)$  and  $\mathcal{O}((NC)^2)$  memory, respectively. In comparison, the naïve approach via  $\mathbf{V}^{(l)} \in \mathbb{R}^{d^{(l)} \times NC}$  scales with the number of weights, which is often comparable to  $D$ . For instance, the 3c3d architecture on CIFAR-10 has  $D = 895,210$  and the largest weight matrix has  $d^{(l)} = 589,824$ , whereas  $NC = 1,280$  during training [146].

### D.3.2 Implicit Multiplication by the Inverse (Block-Diagonal) GGN

#### Inverse GGN-Vector Products

1: Instead of  $\delta \mathbf{I}_D$ , any other easy-to-invert matrix can be used.

A damped Newton step<sup>1</sup> requires multiplication by  $(\mathbf{G} + \delta \mathbf{I}_D)^{-1}$ . By means of Equation (7.3) and the matrix inversion lemma,

$$\begin{aligned} & (\delta \mathbf{I}_D + \mathbf{G})^{-1} \\ &= (\delta \mathbf{I}_D + \mathbf{V} \mathbf{V}^{\top})^{-1} \end{aligned} \quad (\text{Equation (7.3)}) \quad (\text{D.6})$$

$$= \frac{1}{\delta} \left( \mathbf{I}_D + \frac{1}{\delta} \mathbf{V} \mathbf{V}^{\top} \right)^{-1}$$

$$= \frac{1}{\delta} \left[ \mathbf{I}_D - \frac{1}{\delta} \mathbf{V} \left( \mathbf{I}_{NC} + \mathbf{V}^{\top} \frac{1}{\delta} \mathbf{V} \right)^{-1} \mathbf{V}^{\top} \right] \quad (\text{Matrix inversion lemma})$$

$$= \frac{1}{\delta} \left[ \mathbf{I}_D - \mathbf{V} \left( \delta \mathbf{I}_{NC} + \mathbf{V}^{\top} \mathbf{V} \right)^{-1} \mathbf{V}^{\top} \right] \quad (\text{Gram matrix})$$

$$= \frac{1}{\delta} \left[ \mathbf{I}_D - \mathbf{V} \left( \delta \mathbf{I}_{NC} + \tilde{\mathbf{G}} \right)^{-1} \mathbf{V}^{\top} \right]. \quad (\text{D.7})$$

Inverse GGN-vector products require inversion of the damped Gram matrix as well as applications of  $\mathbf{V}, \mathbf{V}^{\top}$  for the transformations between Gram and parameter space.

### Inverse Block-Diagonal GGN-Vector Products

Next, we replace the full GGN by its block diagonal approximation  $\mathbf{G} \approx \mathbf{G}_{\text{BDA}} = \text{diag}(\mathbf{G}^{(1)}, \mathbf{G}^{(2)}, \dots)$  with

$$\mathbf{G}^{(l)} = \mathbf{V}^{(l)} \mathbf{V}^{(l)\top} \in \mathbb{R}^{d^{(l)} \times d^{(l)}}$$

and  $\mathbf{V}^{(l)}$  as in [Equation \(D.3\)](#). Then, inverse multiplication reduces to each block,

$$\mathbf{G}_{\text{BDA}}^{-1} = \text{diag}\left(\mathbf{G}^{(1)^{-1}}, \mathbf{G}^{(2)^{-1}}, \dots\right).$$

If again a damped Newton step is considered, we can reuse [Equation \(D.7\)](#) with the substitutions

$$(\mathbf{G}, D, \mathbf{V}, \mathbf{V}^\top, \tilde{\mathbf{G}}) \leftrightarrow \left(\mathbf{G}^{(l)}, d^{(l)}, \mathbf{V}^{(l)}, \mathbf{V}^{(l)\top}, \tilde{\mathbf{G}}^{(l)}\right)$$

to apply the inverse and immediately discard the ViViT factors: At backpropagation of layer  $T_{\theta^{(l)}}^{(l)}$

1. Compute  $\mathbf{V}^{(l)}$  using [Equation \(D.3\)](#).
2. Compute  $\tilde{\mathbf{G}}^{(l)}$  using [Equation \(D.4\)](#).
3. Compute  $(\delta \mathbf{I}_{NC} + \tilde{\mathbf{G}}^{(l)})^{-1}$ .
4. Apply the inverse in [Equation \(D.7\)](#) with the above substitutions to the target vector.
5. Discard  $\mathbf{V}^{(l)}, \mathbf{V}^{(l)\top}, \tilde{\mathbf{G}}^{(l)}$ , and  $(\delta \mathbf{I}_{NC} + \tilde{\mathbf{G}}^{(l)})^{-1}$ . Proceed to layer  $i-1$ .

Note that the above scheme should only be used for parameters that satisfy  $d^{(l)} > NC$ , i.e.  $\dim(\mathbf{G}^{(l)}) > \dim(\tilde{\mathbf{G}}^{(l)})$ . Low-dimensional parameters can be grouped with others to increase their joint dimension, and to control the block structure of  $\mathbf{G}_{\text{BDA}}$ .



# Bibliography

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng. "TensorFlow: A System for Large-Scale Machine Learning". *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*. 2016.
- [2] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. "Deep Learning with Differential Privacy". *ACM SIGSAC Conference on Computer and Communications Security*. 2016.
- [3] R. P. Adams, J. Pennington, M. J. Johnson, J. Smith, Y. Ovadia, B. Patton, and J. Saunderson. "Estimating the Spectral Density of Large Implicit Matrices". 2018.
- [4] A. M. Agrawal, A. Tindle, H. Sikka, S. Singh, and A. Kayid. "Investigating Learning in Deep Neural Networks using Layer-Wise Weight Change". 2020.
- [5] S.-I. Amari. "Natural Gradient works Efficiently in Learning". *Neural computation* (1998).
- [6] A. Bahamou and D. Goldfarb. "A Dynamic Sampling Adaptive-SGD Method for Machine Learning". 2019.
- [7] D. Bahdanau, K. Cho, and Y. Bengio. "Neural machine translation by jointly learning to align and translate" (2014).
- [8] C. Bakker, M. J. Henry, and N. O. Hodas. "The Outer Product Structure of Neural Network Derivatives" (2018).
- [9] L. Balles. "Noise-Aware Stochastic Optimization" (2022). PhD thesis.
- [10] L. Balles and P. Hennig. "Dissecting Adam: The Sign, Magnitude and Variance of Stochastic Gradients". *International Conference on Machine Learning (ICML)*. 2018.
- [11] L. Balles, J. Romero, and P. Hennig. "Coupling Adaptive Batch Sizes with Learning Rates". *Conference on Uncertainty in Artificial Intelligence (UAI)*. 2017.
- [12] F. L. Bauer. "Computational Graphs and Rounding Error". *SIAM Journal on Numerical Analysis* (1974).
- [13] A. G. Baydin, B. A. Pearlmutter, A. A. Radul, and J. M. Siskind. "Automatic Differentiation in Machine Learning: A Survey". *Journal of Machine Learning Research (JMLR)* (2018).
- [14] S. Becker and Y. Le Cun. "Improving the Convergence of Back-Propagation Learning with Second Order Methods". *Connectionist Models Summer School*. 1989.
- [15] Y. Bengio. "Practical recommendations for gradient-based training of deep architectures". *Neural Networks* (2012).
- [16] A. Bernacchia, M. Lengyel, and G. Hennequin. "Exact natural gradient in deep linear networks and its application to the nonlinear case". *Advances in Neural Information Processing Systems (NeurIPS*. 2018.
- [17] L. Biewald. "Experiment Tracking with Weights and Biases". 2020.
- [18] D. Blalock, J. J. Gonzalez Ortiz, J. Frankle, and J. Guttag. "What is the state of neural network pruning?" *Proceedings of machine learning and systems* (2020).
- [19] R. Bollapragada, R. H. Byrd, and J. Nocedal. "Adaptive Sampling Strategies for Stochastic Optimization". *SIAM Journal on Optimization* (2017).
- [20] A. Bordes, L. Bottou, and P. Gallinari. "SGD-QN: Careful Quasi-Newton Stochastic Gradient Descent". *Journal of Machine Learning Research (JMLR)* (2009).
- [21] A. Botev, H. Ritter, and D. Barber. "Practical Gauss-Newton Optimisation for Deep Learning". *International Conference on Machine Learning (ICML)*. 2017.
- [22] L. Bottou, F. E. Curtis, and J. Nocedal. "Optimization Methods for Large-Scale Machine Learning" (2016).

- [23] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, and S. Wanderman-Milne. “JAX: composable transformations of Python + NumPy programs” (2018).
- [24] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. “Language Models are Few-Shot Learners” (2020).
- [25] R. H. Byrd, G. M. Chin, W. Neveitt, and J. Nocedal. “On the Use of Stochastic Hessian Information in Optimization Methods for Machine Learning”. *SIAM Journal on Optimization* (2011).
- [26] R. H. Byrd, G. M. Chin, J. Nocedal, and Y. Wu. “Sample Size Selection in Optimization Methods for Machine Learning”. *Mathematical Programming* (2012).
- [27] S. Chatterjee. “Coherent Gradients: An Approach to Understanding Generalization in Gradient Descent-based Optimization”. *International Conference on Learning Representations (ICLR)*. 2020.
- [28] S. Chatterjee and P. Zielinski. “Making Coherence Out of Nothing At All: Measuring the Evolution of Gradient Alignment”. 2020.
- [29] K. Chellapilla, S. Puri, and P. Simard. “High Performance Convolutional Neural Networks for Document Processing”. *International Workshop on Frontiers in Handwriting Recognition*. 2006.
- [30] C. Chen, S. Reiz, C. D. Yu, H.-J. Bungartz, and G. Biros. “Fast Approximation of the Gauss–Newton Hessian Matrix for the Multilayer Perceptron”. *SIAM Journal on Matrix Analysis and Applications* (2021).
- [31] S.-W. Chen, C.-N. Chou, and E. Chang. “BDA-PCH: Block-Diagonal Approximation of Positive-Curvature Hessian for Training Neural Networks” (2018).
- [32] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang. “MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems”. *Neural Information Processing Systems (NIPS), Workshop on Machine Learning Systems*. 2015.
- [33] K. Cho, B. Merriënboer, D. Bahdanau, and Y. Bengio. “On the Properties of Neural Machine Translation: Encoder-Decoder Approaches” (2014).
- [34] D. Choi, C. J. Shallue, Z. Nado, J. Lee, C. J. Maddison, and G. E. Dahl. “On Empirical Comparisons of Optimizers for Deep Learning”. 2020.
- [35] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha. “StarGAN v2: Diverse Image Synthesis for Multiple Domains”. *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2020.
- [36] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. “Introduction to Algorithms”. 2001.
- [37] F. Dangel, S. Harmeling, and P. Hennig. “Modular Block-diagonal Curvature Approximations for Feedforward Architectures”. *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2020.
- [38] F. Dangel, F. Kunstner, and P. Hennig. “BackPACK: Packing more into Backprop”. *International Conference on Learning Representations (ICLR)*. 2020.
- [39] F. Dangel, L. Tatze, and P. Hennig. “ViViT: Curvature Access Through The Generalized Gauss–Newton’s Low-Rank Structure”. *Transactions on Machine Learning Research (TMLR)* (2022).
- [40] E. Daxberger, A. Kristiadi, A. Immer, R. Eschenhagen, M. Bauer, and P. Hennig. “Laplace Redux–Effortless Bayesian Deep Learning”. *Advances in Neural Information Processing Systems (NeurIPS)*. 2021.
- [41] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. “ImageNet: A Large-Scale Hierarchical Image Database”. *IEEE conference on computer vision and pattern recognition (CVPR)*. 2009.
- [42] L. Dinh, R. Pascanu, S. Bengio, and Y. Bengio. “Sharp Minima Can Generalize For Deep Nets”. *International Conference on Machine Learning (ICML)*. 2017.
- [43] X. Dong, S. Chen, and S. Pan. “Learning to Prune Deep Neural Networks via Layer-wise Optimal Brain Surgeon”. *Advances in Neural Information Processing Systems (NIPS)*. 2017.

- [44] J. Duchi, E. Hazan, and Y. Singer. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization". *Journal of Machine Learning Research (JMLR)* (2011).
- [45] V. Dumoulin and F. Visin. "A guide to convolution arithmetic for deep learning" (2016).
- [46] C. Dwork, A. Roth, et al. "The algorithmic foundations of differential privacy". *Foundations and Trends in Theoretical Computer Science* (2014).
- [47] J. L. Elman. "Finding structure in time". *Cognitive Science* (1990).
- [48] F. Faghri, D. Duvenaud, D. J. Fleet, and J. Ba. "A Study of Gradient Variance in Deep Learning". 2020.
- [49] W. Fedus, B. Zoph, and N. Shazeer. "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity". *Journal of Machine Learning Research (JMLR)* (2022).
- [50] H. Fischer. "A History of the Central Limit Theorem: From Classical to Modern Probability Theory". 2010.
- [51] J. Frankle, D. J. Schwab, and A. S. Morcos. "The Early Phase of Neural Network Training". *International Conference on Learning Representations (ICLR)*. 2020.
- [52] M. Fredrikson, S. Jha, and T. Ristenpart. "Model Inversion Attacks That Exploit Confidence Information and Basic Countermeasures". *ACM Conference on Computer and Communications Security (CCS)*. 2015.
- [53] M. Gargiani, A. Zanelli, M. Diehl, and F. Hutter. "On the Promise of the Stochastic Generalized Gauss-Newton Method for Training DNNs". 2020.
- [54] T. George, C. Laurent, X. Bouthillier, N. Ballas, and P. Vincent. "Fast Approximate Natural Gradient Descent in a Kronecker-factored Eigenbasis" (2018).
- [55] B. Ghorbani, S. Krishnan, and Y. Xiao. "An Investigation into Neural Net Optimization via Hessian Eigenvalue Density". *International Conference on Machine Learning (ICML)*. 2019.
- [56] B. Ginsburg. "On regularization of gradient descent, layer imbalance and flat minima". 2020.
- [57] X. Glorot and Y. Bengio. "Understanding the Difficulty of Training Deep Feedforward Neural Networks". *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2010.
- [58] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. "Generative Adversarial Nets". *Neural Information Processing Systems (NIPS)*. 2014.
- [59] I. J. Goodfellow. "Efficient Per-Example Gradient Computations" (2015).
- [60] I. J. Goodfellow, Y. Bengio, and A. Courville. "Deep Learning". 2016.
- [61] D. Granziol, X. Wan, and T. Garipov. "Deep Curvature Suite". 2021.
- [62] A. Griewank. "Who invented the reverse mode of differentiation?" *Documenta Mathematica* (2012).
- [63] R. Grosse and J. Martens. "A Kronecker-factored approximate Fisher matrix for convolution layers". *International Conference on Machine Learning (ICML)*. 2016.
- [64] I. Gulrajani and D. Lopez-Paz. "In Search of Lost Domain Generalization". *International Conference on Learning Representations (ICLR)*. 2021.
- [65] G. Gur-Ari, D. A. Roberts, and E. Dyer. "Gradient Descent Happens in a Tiny Subspace". 2018.
- [66] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. "Array programming with NumPy". *Nature* (2020).
- [67] B. Hassibi and D. Stork. "Second order derivatives for network pruning: Optimal Brain Surgeon". *Advances in Neural Information Processing Systems (NIPS)*. 1992.
- [68] K. He, X. Zhang, S. Ren, and J. Sun. "Deep Residual Learning for Image Recognition". *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016.
- [69] J. F. Henriques, S. Ehrhardt, S. Albanie, and A. Vedaldi. "Small Steps and Giant Leaps: Minimal Newton Solvers for Deep Learning". *International Conference on Computer Vision (ICCV)*. 2019.
- [70] S. Hochreiter and J. Schmidhuber. "Flat Minima". *Neural Computation* (1997).

- [71] S. Hochreiter and J. Schmidhuber. "Long Short-Term Memory". *Neural Computation* (1997).
- [72] R. Z. Horace He. "functorch: JAX-like composable function transforms for PyTorch". 2021.
- [73] J. Hughes. "Why Functional Programming Matters". *Computer Journal* (1989).
- [74] Y. Idelbayev. "Proper ResNet Implementation for CIFAR10/CIFAR100 in PyTorch". 2018.
- [75] A. Immer, M. Bauer, V. Fortuin, G. Rätsch, and K. M. Emtiyaz. "Scalable Marginal Likelihood Estimation for Model Selection in Deep Learning". *International Conference on Machine Learning (ICML)*. 2021.
- [76] A. Immer, M. Korzepa, and M. Bauer. "Improving predictions of Bayesian neural nets via local linearization". *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2021.
- [77] M. Innes. "Don't Unroll Adjoint: Differentiating SSA-Form Programs" (2018).
- [78] M. Innes. "Flux: Elegant machine learning with Julia". *Journal of Open Source Software* (2018).
- [79] S. Ioffe and C. Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". *International Conference on Machine Learning (ICML)*. 2015.
- [80] S. Jastrzebski, D. Arpit, O. Astrand, G. B. Kerg, H. Wang, C. Xiong, R. Socher, K. Cho, and K. J. Geras. "Catastrophic Fisher Explosion: Early Phase Fisher Matrix Impacts Generalization". *International Conference on Machine Learning (ICML)*. 2021.
- [81] S. Jastrzebski, M. Szymczak, S. Fort, D. Arpit, J. Tabor, K. Cho, and K. Geras. "The Break-Even Point on Optimization Trajectories of Deep Neural Networks". *International Conference on Learning Representations (ICLR)*. 2020.
- [82] W. P. Johnson. "The curious history of Faà di Bruno's formula". *The American mathematical monthly* (2002).
- [83] N. P. Jouppi, C. Young, N. Patil, D. Patterson, G. Agrawal, R. Bajwa, S. Bates, S. Bhatia, N. Boden, A. Borchers, R. Boyle, P.-I. Cantin, C. Chao, C. Clark, J. Coriell, M. Daley, M. Dau, J. Dean, B. Gelb, T. V. Ghaemmaghami, R. Gottipati, W. Gulland, R. Hagmann, C. R. Ho, D. Hogberg, J. Hu, R. Hundt, D. Hurt, J. Ibarz, A. Jaffey, A. Jaworski, A. Kaplan, H. Khaitan, D. Killebrew, A. Koch, N. Kumar, S. Lacy, J. Laudon, J. Law, D. Le, C. Leary, Z. Liu, K. Lucke, A. Lundin, G. MacKean, A. Maggiore, M. Mahony, K. Miller, R. Nagarajan, R. Narayanaswami, R. Ni, K. Nix, T. Norrie, M. Omernick, N. Penukonda, A. Phelps, J. Ross, M. Ross, A. Salek, E. Samadiani, C. Severn, G. Sizikov, M. Snelham, J. Souter, D. Steinberg, A. Swing, M. Tan, G. Thorson, B. Tian, H. Toma, E. Tuttle, V. Vasudevan, R. Walter, W. Wang, E. Wilcox, and D. H. Yoon. "In-Datacenter Performance Analysis of a Tensor Processing Unit". *Annual International Symposium on Computer Architecture*.
- [84] T. Karras, T. Aila, S. Laine, and J. Lehtinen. "Progressive Growing of GANs for Improved Quality, Stability, and Variation". *International Conference on Learning Representations (ICLR)*. 2018.
- [85] A. Katharopoulos and F. Fleuret. "Not All Samples Are Created Equal: Deep Learning with Importance Sampling". *International Conference on Machine Learning (ICML)*. 2018.
- [86] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang. "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima". *International Conference on Learning Representations (ICLR)*. 2017.
- [87] D. P. Kingma and J. Ba. "Adam: A Method for Stochastic Optimization". *International Conference on Learning Representations (ICLR)*. 2015.
- [88] D. Kirk. "NVIDIA CUDA software and GPU parallel computing architecture". 2007.
- [89] A. Kristiadi, M. Hein, and P. Hennig. "Being Bayesian, even just a bit, fixes overconfidence in ReLU networks". *International Conference on Machine Learning (ICML)*. 2020.
- [90] A. Krizhevsky. *Learning multiple layers of features from tiny images*. Technical report. 2009.
- [91] A. Krizhevsky, I. Sutskever, and G. E. Hinton. "ImageNet Classification with Deep Convolutional Neural Networks". *Neural Information Processing Systems (NIPS)*. 2012.
- [92] F. Kunstner, L. Balles, and P. Hennig. "Limitations of the Empirical Fisher Approximation". *Advances in Neural Information Processing Systems (NeurIPS)*. 2019.

- [93] P. S. Laplace. "Mémoire sur la probabilité de causes par les événements". *Mémoire de l'académie royale des sciences* (1774).
- [94] N. Le Roux, P. Manzagol, and Y. Bengio. "Topmoumoute Online Natural Gradient Algorithm". *Advances in Neural Information Processing Systems* 20. 2007.
- [95] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition". *Proceedings of the IEEE*. 1998.
- [96] Y. LeCun, J. Denker, and S. Solla. "Optimal Brain Damage". *Neural Information Processing Systems (NIPS)*. 1989.
- [97] Y. LeCun, P. Simard, and B. Pearlmutter. "Automatic Learning Rate Maximization by On-Line Estimation of the Hessian's Eigenvectors". *Neural Information Processing Systems (NIPS)*. 1993.
- [98] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. "Microsoft COCO: Common Objects in Context". *European Conference on Computer Vision (ECCV)*. 2014.
- [99] S. Linnainmaa. "Taylor expansion of the accumulated rounding error". *BIT Numerical Mathematics* (1976).
- [100] J. Liu, Y. Bai, G. Jiang, T. Chen, and H. Wang. "Understanding Why Neural Networks Generalize Well Through GSNR of Parameters". *International Conference on Learning Representations (ICLR)*. 2020.
- [101] C. F. Loan. "The ubiquitous Kronecker product". *Journal of Computational and Applied Mathematics* (2000).
- [102] M.-T. Luong, H. Pham, and C. D. Manning. "Effective approaches to attention-based neural machine translation" (2015).
- [103] J. R. Magnus and H. Neudecker. "Matrix Differential Calculus with Applications in Statistics and Econometrics". 1999.
- [104] M. Mahsereci, L. Balles, C. Lassner, and P. Hennig. "Early Stopping without a Validation Set". 2017.
- [105] M. Mahsereci and P. Hennig. "Probabilistic Line Searches for Stochastic Optimization". *Journal of Machine Learning Research (JMLR)* (2017).
- [106] J. Martens. "Deep Learning via Hessian-Free Optimization." *International Conference on Machine Learning (ICML)*. 2010.
- [107] J. Martens. "New perspectives on the natural gradient method" (2014).
- [108] J. Martens, J. Ba, and M. Johnson. "Kronecker-factored Curvature Approximations for Recurrent Neural Networks". *International Conference on Learning Representations (ICLR)*. 2018.
- [109] J. Martens and R. Grosse. "Optimizing Neural Networks with Kronecker-Factored Approximate Curvature". *International Conference on Machine Learning (ICML)*. 2015.
- [110] E. Mizutani and S. E. Dreyfus. "Second-Order Stagewise Backpropagation for Hessian-Matrix Analyses and Investigation of Negative Curvature". *Neural Networks* (2008).
- [111] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. "Playing Atari with Deep Reinforcement Learning" (2013).
- [112] R. Mulayoff and T. Michaeli. "Unique Properties of Flat Minima in Deep Networks". *International Conference on Machine Learning (ICML)*. 2020.
- [113] V. Nagarajan and J. Z. Kolter. "Generalization in Deep Networks: The Role of Distance from Initialization". 2019.
- [114] U. Naumann. "Optimal Jacobian Accumulation is NP-Complete". *Mathematical Programming* (2008).
- [115] M. Naumov. "Feedforward and Recurrent Neural Networks Backward Propagation and Hessian in Matrix Form" (2017).
- [116] D. Needell, R. Ward, and N. Srebro. "Stochastic Gradient Descent, Weighted Sampling, and the Randomized Kaczmarz algorithm". *Advances in Neural Information Processing Systems (NIPS)*. 2014.
- [117] Y. Nesterov. "A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ ". *Proceedings of the USSR Academy of Sciences* (1983).

- [118] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. "Reading Digits in Natural Images with Unsupervised Feature Learning". 2011.
- [119] D. Oktay, N. McGreivy, J. Aduol, A. Beatson, and R. P. Adams. "Randomized Automatic Differentiation". *International Conference on Learning Representations (ICLR)*. 2021.
- [120] Y. Ollivier, L. Arnold, A. Auger, and N. Hansen. "Information-Geometric Optimization Algorithms: A Unifying Picture via Invariance Principles" (2011).
- [121] K. Osawa, Y. Tsuji, Y. Ueno, A. Naruse, R. Yokota, and S. Matsuoka. "Large-Scale Distributed Second-Order Optimization Using Kronecker-Factored Approximate Curvature for Deep Convolutional Neural Networks". *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2019.
- [122] V. Papyan. "Measurements of Three-Level Hierarchical Structure in the Outliers in the Spectrum of Deepnet Hessians". *International Conference on Machine Learning (ICML)*. 2019.
- [123] V. Papyan. "The Full Spectrum of Deepnet Hessians at Scale: Dynamics with SGD Training and Sample Size". 2019.
- [124] V. Papyan. "Traces of Class/Cross-Class Structure Pervade Deep Learning Spectra". *Journal of Machine Learning Research (JMLR)* (2020).
- [125] R. Pascanu and Y. Bengio. "Revisiting natural gradient for deep networks" (2013).
- [126] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. "PyTorch: An Imperative Style, High-Performance Deep Learning Library". *Advances in Neural Information Processing Systems (NeurIPS)*. 2019.
- [127] B. A. Pearlmutter. "Fast Exact Multiplication by the Hessian". *Neural Computation* (1994).
- [128] B. Polyak. "Some methods of speeding up the convergence of iteration methods". *USSR Computational Mathematics and Mathematical Physics* (1964).
- [129] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. "Language Models are Unsupervised Multitask Learners" (2019).
- [130] A. Rame, C. Dancette, and M. Cord. "Fishr: Invariant Gradient Variances for Out-of-distribution Generalization". *International Conference on Machine Learning (ICML)*. 2022.
- [131] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. "Zero-Shot Text-to-Image Generation". *International Conference on Machine Learning (ICML)*. 2021.
- [132] S. J. Reddi, S. Kale, and S. Kumar. "On the Convergence of Adam and Beyond". *International Conference on Learning Representations (ICLR)*. 2018.
- [133] H. Ritter, A. Botev, and D. Barber. "A Scalable Laplace Approximation for Neural Networks". *International Conference on Learning Representations (ICLR)*. 2018.
- [134] H. Ritter, A. Botev, and D. Barber. "Online Structured Laplace Approximations for Overcoming Catastrophic Forgetting". *Advances in Neural Information Processing Systems (NeurIPS)*. 2018.
- [135] H. Robbins and S. Monro. "A Stochastic Approximation Method". *The Annals of Mathematical Statistics* (1951).
- [136] F. Rosenblatt. "The perceptron: a probabilistic model for information storage and organization in the brain." *Psychological Review* (1958).
- [137] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. "Learning Internal Representations by Error Propagation". *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Volume 1: Foundations*. 1986.
- [138] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. "Imagenet large scale visual recognition challenge". *International journal of computer vision* (2015).
- [139] L. Sagun, L. Bottou, and Y. LeCun. "Eigenvalues of the Hessian in Deep Learning: Singularity and Beyond". 2017.

- [140] L. Sagun, U. Evci, V. U. Guney, Y. Dauphin, and L. Bottou. "Empirical Analysis of the Hessian of Over-Parametrized Neural Networks". 2018.
- [141] A. Sankaran, N. A. Alashti, C. Psarras, and P. Bientinesi. "Benchmarking the Linear Algebra Awareness of TensorFlow and PyTorch". 2022.
- [142] K. A. Sankararaman, S. De, Z. Xu, W. R. Huang, and T. Goldstein. "The Impact of Neural Network Over-parameterization on Gradient Confusion and Stochastic Gradient Descent". *International Conference on Machine Learning (ICML)*. 2020.
- [143] A. M. Saxe, J. L. McClelland, and S. Ganguli. "Exact solutions to the nonlinear dynamics of learning in deep linear neural networks". *International Conference on Learning Representations (ICLR)*. 2014.
- [144] M. Schmidt. "Convergence rate of stochastic gradient with constant step size" (2014).
- [145] R. M. Schmidt, F. Schneider, and P. Hennig. "Descending through a Crowded Valley - Benchmarking Deep Learning Optimizers". *International Conference on Machine Learning (ICML)*. 2021.
- [146] F. Schneider, L. Balles, and P. Hennig. "DeepOBS: A Deep Learning Optimizer Benchmark Suite". *International Conference on Learning Representations (ICLR)*. 2019.
- [147] F. Schneider, F. Dangel, and P. Hennig. "Cockpit: A Practical Debugging Tool for Training Deep Neural Networks". *Advances in Neural Information Processing Systems (NeurIPS)*. 2021.
- [148] N. N. Schraudolph. "Fast Curvature Matrix-Vector Products for Second-Order Gradient Descent". *Neural Computation* (2002).
- [149] R. Shokri and V. Shmatikov. "Privacy-Preserving Deep Learning". *ACM SIGSAC Conference on Computer and Communications Security*. 2015.
- [150] C. Shorten and T. M. Khoshgoftaar. "A survey on image data augmentation for deep learning". *Journal of big data* (2019).
- [151] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. "Mastering the Game of Go with Deep Neural Networks and Tree Search". *Nature* (2016).
- [152] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, and D. Hassabis. "A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play". *Science* (2018).
- [153] K. Simonyan and A. Zisserman. "Very Deep Convolutional Networks for Large-Scale Image Recognition". *International Conference on Learning Representations (ICLR)*. 2015.
- [154] S. P. Singh and D. Alistarh. "WoodFisher: Efficient Second-Order Approximation for Neural Network Compression". *Advances in Neural Information Processing Systems (NeurIPS)*. 2020.
- [155] M. Skorski, A. Temperoni, and M. Theobald. "Revisiting Weight Initialization of Deep Neural Networks". *Asian Conference on Machine Learning (ACML)*. 2021.
- [156] D. So, W. Mańke, H. Liu, Z. Dai, N. Shazeer, and Q. V. Le. "Searching for Efficient Transformers for Language Modeling". *Advances in Neural Information Processing Systems (NeurIPS)*. 2021.
- [157] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller. "Striving for Simplicity: The All Convolutional Net". *International Conference on Learning Representations (ICLR)*. 2015.
- [158] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. "Dropout: A Simple Way to Prevent Neural Networks from Overfitting". *Journal of Machine Learning Research (JMLR)* (2014).
- [159] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. "Going deeper with convolutions". *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2015.
- [160] K. Takeuchi. "The distribution of information statistics and the criterion of goodness of fit of models". *Mathematical Science* (1976).
- [161] L. Theis, I. Korshunova, A. Tejani, and F. Huszár. "Faster gaze prediction with dense networks and Fisher pruning" (2018).

- [162] V. Thomas, F. Pedregosa, B. van Merriënboer, P.-A. Manzagol, Y. Bengio, and N. L. Roux. “On the interplay between noise and curvature and its effect on optimization and generalization”. *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2020.
- [163] N. C. Thompson, K. Greenewald, K. Lee, and G. F. Manso. “The Computational Limits of Deep Learning”. 2020.
- [164] T. Tieleman, G. Hinton, et al. “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude”. COURSERA: *Neural networks for machine learning* (2012).
- [165] S. Tokui, K. Oono, S. Hido, and J. Clayton. “Chainer: A next-generation open source framework for deep learning”. *Neural Information Processing Systems (NIPS), Workshop on Machine Learning Systems*. 2015.
- [166] Y. Tsuji, K. Osawa, Y. Ueno, A. Naruse, R. Yokota, and S. Matsuoka. “Performance Optimizations and Analysis of Distributed Deep Learning with Approximated Second-Order Optimization Method”. *International Conference on Parallel Processing, Workshop Proceedings*. 2019.
- [167] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin. “Attention is All you Need”. *Advances in Neural Information Processing Systems (NeurIPS)*. 2017.
- [168] S. Vaswani, A. Mishkin, I. Laradji, M. Schmidt, G. Gidel, and S. Lacoste-Julien. “Painless Stochastic Gradient: Interpolation, Line-Search, and Convergence Rates”. *Neural Information Processing Systems (NeurIPS)*. 2019.
- [169] L. Wang, Y. Yang, M. R. Min, and S. T. Chakradhar. “Accelerating Deep Neural Network Training with Inconsistent Stochastic Gradient Descent”. *Neural networks* (2017).
- [170] J. Warsa, T. Wareing, J. Morel, J. Mcghee, and R. Lehoucq. “Krylov Subspace Iterations for Deterministic k-Eigenvalue Calculations”. *Nuclear Science and Engineering* (2004).
- [171] Wikipedia. “Toeplitz matrix”. [https://en.wikipedia.org/w/index.php?title=Toeplitz\\_matrix&oldid=1094551484](https://en.wikipedia.org/w/index.php?title=Toeplitz_matrix&oldid=1094551484). (online) accessed August 11. 2022.
- [172] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al. “Google’s neural machine translation system: Bridging the gap between human and machine translation” (2016).
- [173] Y. Wu, M. Ren, R. Liao, and R. B. Grosse. “Understanding short-horizon bias in stochastic meta-optimization”. *International Conference on Learning Representations (ICLR)* (2018).
- [174] Y. Wu and K. He. “Group Normalization”. *International Journal of Computer Vision* (2019).
- [175] H. Xiao, K. Rasul, and R. Vollgraf. “Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms”. 2017.
- [176] C. Xing, D. Arpit, C. Tsirigotis, and Y. Bengio. “A Walk with SGD”. 2018.
- [177] Z. Yao, A. Gholami, K. Keutzer, and M. W. Mahoney. “PyHessian: Neural Networks Through the Lens of the Hessian”. *IEEE International Conference on Big Data*. 2020.
- [178] Z. Yao, A. Gholami, S. Shen, M. Mustafa, K. Keutzer, and M. W. Mahoney. “ADAHESSIAN: An Adaptive Second Order Optimizer for Machine Learning”. 2021.
- [179] A. Yousefpour, I. Shilov, A. Sablayrolles, D. Testuggine, K. Prasad, M. Malek, J. Nguyen, S. Ghosh, A. Bharadwaj, J. Zhao, G. Cormode, and I. Mironov. “Opacus: User-Friendly Differential Privacy Library in PyTorch”. *Advances in Neural Information Processing Systems (NeurIPS), Workshop Privacy in Machine Learning*. 2021.
- [180] M. Zeiler. “ADADELTA: An adaptive learning rate method” (2012).
- [181] M. D. Zeiler and R. Fergus. “Visualizing and Understanding Convolutional Networks”. *European Conference on Computer Vision (ECCV)*. 2014.
- [182] W. Zeng and R. Urtasun. “MLPrune: Multi-Layer Pruning for Automated Neural Network Compression”. 2018.
- [183] H. Zhang, C. Xiong, J. Bradbury, and R. Socher. “Block-diagonal Hessian-free Optimization for Training Neural Networks” (2017).

- [184] P. Zhao and T. Zhang. "Stochastic Optimization with Importance Sampling for Regularized Loss Minimization". *International Conference on Machine Learning (ICML)*. 2015.