## Connecting Parameter Magnitudes and Hessian Eigenspaces at Scale using Sketched Methods

Andres Fernandez, Frank Schneider, Maren Mahreseci, Philipp Hennig {a.fernandez, f.schneider, maren.mahsereci, philipp.hennig}@uni-tuebingen.de May 9, 2025









- ▶ **Dataset** of pairs (**x**<sub>i</sub>, **y**<sub>i</sub>)
- ▶ Neural Network  $\hat{y}_i = f(x_i, \theta)$  with parameters  $\theta \in \mathbb{R}^D$
- ▶ Loss  $\mathcal{L}(\boldsymbol{\theta}) = \sum_{i} \ell(\boldsymbol{y}_{i}, f(\boldsymbol{x}_{i}, \boldsymbol{\theta})) \in \mathbb{R}_{\geq 0}$  with gradient  $\boldsymbol{g} \in \mathbb{R}^{D}$  and Hessian  $\boldsymbol{H} \in \mathbb{R}^{D \times D}$

Early crystallization of parameters:

- > Parameters can be *pruned* (Blalock et al. 2020)
- ▶ Pruning masks  $\theta \odot m$  appear *early* in training (Frankle et al. 2019)
- Magnitude pruning masks don't change much during training! (You et al. 2020)

Early crystallization of loss landscape:

- ► **H** is rank-defficient, i.e.  $H \approx U_{top} \Lambda_{top} U_{top}^{\top}$  (e.g. Sagun et al. 2018)
- ▶ **g** resides mostly in **U**top (Gur-Ari et al. 2019)
- ▶ span(U<sub>top</sub>) doesn't change much during training! (Gur-Ari et al. 2019)



Are those connected?



### Research questions and contributions



#### Questions:

- ► Can this similarity be measured? If so, how?
- ▶ What similarity can be considered high? What are the implications?

#### **Contributions:**

- ▶ Methodology to compare arbitrary *k*-parameter masks to *top-k* Hessian eigenspaces
- $\blacktriangleright$  Algorithm and code to perform said measurements at scale  $\rightarrow$  Hessian eigendecompositions
- ▶ In DL, connection is orders of magnitude larger than random
- > Potential implications for pruning, optimization, UQ and loss landscape analysis

# Comparing parameters with Hessian subspaces

UNIVERSITAT

1

Top-*k* parameter pruning is a projection onto  $I_{D,k}$ :

Also recall the *top-k* eigenbasis  $U_{top}$ :

$$\mathbf{P}^{\top}(\mathbf{m}_{k} \odot \mathbf{\theta}) = \tilde{\mathbf{m}}_{k} \odot \tilde{\mathbf{\theta}} = \begin{pmatrix} \mathbf{I}_{k} & 0 \\ 0 & 0 \end{pmatrix} \tilde{\mathbf{\theta}} \eqqcolon \mathbf{I}_{D,k} \mathbf{I}_{D,k}^{\top} \tilde{\mathbf{\theta}} \qquad \mathbf{H} = \begin{pmatrix} \mathbf{U}_{top} & \mathbf{U}_{bulk} \\ \mathbf{U}_{bulk} & \mathbf{U}_{bulk} \end{pmatrix} \begin{pmatrix} \mathbf{D}_{top} & \mathbf{I}_{bulk} \\ - - - + \mathbf{I}_{bulk} & \mathbf{I}_{bulk} \\ - - - + \mathbf{I}_{bulk} & \mathbf{I}_{bulk} \end{pmatrix}$$

- We have same-shape, orthogonal matrices I<sub>D,k</sub> and U<sub>top</sub>
- **Grassmannian metrics** measure the distance between their **spaces**
- ▶ Theoretical and empirical analysis of several Grassmannian metrics
- ▶ The overlap metric is stable and has a random baseline value of  $\frac{k}{D}$ :

$$\frac{1}{k} \| \boldsymbol{J}_{D,k}^{\top} \boldsymbol{U}_{top} \|_{F}^{2} \in [0, 1] \quad \text{(higher} \iff \text{more similar)}$$



- Computing overlap requires top-k Hessian eigendecomposition
- ▶ Intractable:  $\mathcal{O}(D^2)$  memory,  $\mathcal{O}(D^3)$  arithmetic (Golub et al. 2013)
- **Expensive measurements:** Each w = Hv costs 2 forward+backpropagations (Pearlmutter 1994)
- Sketched methods:  $\mathcal{O}(k)$  parallel measurements,  $\mathcal{O}(Dk)$  memory (Halko et al. 2011)
- PyTorch library: skerch





 $\frac{1}{k} \| \boldsymbol{I}_{Dk}^{\top} \boldsymbol{U}_{top} \|_{F}^{2}$ 

#### pip install skerch



#### Sketched SVD: Intuition

- Consider a LinOp, very large and matrix-free, with expensive measurements but some simpler sub-structure (low-rank in this case)
- We'll see how to sketch-and-solve:
  - Draw a few random measurements in a way that captures the sub-structure
  - Project measurements back into original space
- We will cover:
  - ▶ Sketching step: random measurements with guarantees
  - Solving step: cheap, traiditonal linear algebra
  - ▶ How this may help approximating full DL Hessians





# The Randomized Range Finder (RRF) (Halko et al. 2011)



Let's revisit our low-rank operator  $\mathbf{A} \in \mathbb{C}^{m \times n}$ :

$$\mathbf{A} := \left( \begin{array}{cc} \mathbf{U} & \bar{\mathbf{U}} \end{array} \right) \left( \begin{array}{c} \mathbf{D} \\ - & \bar{\mathbf{E}} \end{array} \right) \quad \left( \begin{array}{c} \mathbf{V}^* \\ \bar{\mathbf{V}}^* \end{array} \right) \quad \approx \quad \mathbf{U} \, \mathbf{D} \, \mathbf{V}^*$$

Consider  $\tilde{\boldsymbol{U}}, \tilde{\boldsymbol{V}}$  as *rotations* of  $\boldsymbol{U}, \boldsymbol{V}$ , such that  $\tilde{\boldsymbol{U}} := \boldsymbol{U}\boldsymbol{Z}$  for unitary  $\boldsymbol{Z} \in {}^{k \times k}$ . The key observation is that:

#### $\mathbf{A} pprox \mathbf{U} \mathbf{U}^* \mathbf{A} \mathbf{V} \mathbf{V}^{ op} = \widetilde{\mathbf{U}} \widetilde{\mathbf{U}}^* \mathbf{A} \widetilde{\mathbf{V}} \widetilde{\mathbf{V}}^*$

And the "magic" is that  $\tilde{U}$ ,  $\tilde{V}$  can be obtained from  $\mathcal{O}(k)$  random measurements! (Halko et al. 2011). Intuition:

- ▶ If we draw  $\{A\omega_1, A\omega_2, ...\}$  measurements, they are all likely to be **linearly independent**
- Furthermore, if they are random, they all likely to land on the top subspace
- ► The two conditions above mean that we are **fully covering the top subspace**
- $\blacktriangleright$  QR orthogonalization of our measurements yields  $ilde{m{U}}$
- Measurements from quasi-orthogonal iid noise are numerically stable and parallelizable



## Sketched SVD (Halko et al. 2011; Tropp et al. 2017)

Naively, all is left is to perform a  $k \times k$  SVD:

$$\mathbf{A} \approx \tilde{\mathbf{U}} \underbrace{\tilde{\mathbf{U}}_{\mathbf{C}=\mathbf{Z}_{1}\boldsymbol{\Sigma}\mathbf{Z}_{2}^{T}}^{*}}_{\mathbf{C}=\mathbf{Z}_{1}\boldsymbol{\Sigma}\mathbf{Z}_{2}^{T}} \tilde{\mathbf{V}}^{*}$$

Nice! But this requires us to do a **second, expensive pass** over  $A\tilde{V}$ , which is **unnecesary!** 

$$\begin{split} & \mathsf{Y} = \mathsf{A}\Omega = \mathsf{A}\tilde{\mathsf{V}}\tilde{\mathsf{V}}^*\Omega \\ \Longleftrightarrow \quad \underbrace{\mathsf{A}\tilde{\mathsf{V}}}_{\mathsf{U}\Sigma\mathsf{Z}^*} = \mathsf{Y}(\tilde{\mathsf{V}}^*\Omega)^\dagger \iff \mathsf{A} = \mathsf{A}\tilde{\mathsf{V}}\tilde{\mathsf{V}}^* = \mathsf{U}\Sigma\underbrace{(\mathsf{Z}^*\tilde{\mathsf{V}}^*)}_{\mathsf{V}^*} \end{split}$$

This is very good!

- First we obtain thin outer matrices  $\tilde{\boldsymbol{U}} \in \mathcal{N}^{m \times \alpha}$ ,  $\tilde{\boldsymbol{V}} \in \mathcal{N}^{m \times \alpha}$  from random measurements and QR
- For the the term of term
- ► For **single-pass**, solve a well-conditioned least squares problem







- **Scalable**: Rank-1500 eigendecompositions on 12M-parameter networks
- Orders of magnitude higher for all observed splits, steps, rank sizes and problems
- $\blacktriangleright$  Parameter inspection cheaply informs about curvature  $\rightarrow$  training, pruning, UQ, analysis
- Still, spaces are far from identical  $(\frac{k}{D}$  is small), so no direct mapping

## skerch: Sketched matrix decompositions for PyTorch

"I love it, 5/5!" — Diederik P. Kingma, probably talking about the Adam optimizer



# https://github.com/ andres-fr/skerch

Fernandez, Schneider, Mahsereci, Hennig – Connecting Parameter Magnitudes and Hessian Eigenspaces at Scale using Scheduler and Scheduler and Scale using Scheduler and Scheduler and Scale using Scheduler and Scheduler and Scale using Scheduler and Scale using Scheduler and Scale using

## Thank you!



#### **Conclusions:**

- ► Grassmannian metrics to compare arbitrary parameters and Hessian eigenspaces
- $\blacktriangleright$  Sketched eigendecompositions to measure overlap at scale  $\rightarrow \texttt{skerch}$
- > DL overlap orders-of-magnitude larger than baseline (albeit far from identical)
- > Connecting expensive Hessian quantities with cheap parameter observations

#### Future work:

- ► Scalability: We also explore faster alternatives like perturbation-based and GGN
- ► Explaining why do we observe high overlap
- ► Leveraging this effect in downstream applications

#### How can this be useful to you?

- ► Insights on comparing spaces of seemingly unrelated quantities
- ► Scalable linear algebra
- ▶ DL Hessians: Optimization, pruning, UQ, learning theory

Fernandez, Schneider, Mahsereci, Hennig – Connecting Parameter Magnitudes and Hessian Eigenspaces at Scale using Sketched Methods (TMLR 2025)



Blalock, Davis et al. (2020), "What is the State of Neural Network Pruning?" In: Proceedings of Machine Learning and Systems (MLSys). Candès, Emmanuel J. (2008). "The restricted isometry property and its implications for compressed sensing". In: Comptes Rendus Mathematique 346.9, pp. 589–592. ISSN: 1631-073X. DOI: https://doi.org/10.1016/j.crma.2008.03.014.URL https://www.sciencedirect.com/science/article/pii/S1631073X08000964. Frankle, Jonathan and Michael Carbin (2019). "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks". In: International Conference on Learning Representations. URL: https://openreview.net/forum?id=r.Ul-b3RcF7 Golub, Gene H. and Charles F. Van Loan (2013), Matrix Computations, Fourth, The Johns Hopkins University Press, Gur-Ari, Guy, Daniel A. Roberts, and Ethan Dver (2019). Gradient Descent Happens in a Tiny Subspace, URL: https://openreview.net/forum?id=BveTHsAgtX. Halko, N., P. G. Martinsson, and J. A. Troop (2011). "Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions". In Hardt, Moritz, Beniamin Recht, and Yoram Singer (2016). "Train Faster, Generalize Better: Stability of Stochastic Gradient Descent". In: Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48. ICML'16. New York, NY, USA: JMLR.org, pp. 1225-1234. Pearlmutter, Barak A. (1994), "Fast Exact Multiplication by the Hessian", In. Sagun, Levent et al. (2018), Empirical Analysis of the Hessian of Over-Parametrized Neural Networks, ICLR 2018 Workshop, URL: https://openreview.net/forum?id=rJrTwxbCb. Tropp, Joel A. et al. (2017). "Practical Sketching Algorithms for Low-Rank Matrix Approximation". In: SIAM Journal on Matrix Analysis and Applications. (2019) "Streaming Low-Rank Matrix Approximation with an Application to Scientific Simulation". In: SIAM Journal on Scientific Computing You, Haoran et al. (2020), "Drawing Early-Bird Tickets: Toward More Efficient Training of Deep Networks", In: International Conference on Learning Representations, URL: https://openreview.net/forum?id=B.lxsrgStvr

## Backup

A tale as old as Geoffrey Hinton



- ► Dataset  $\mathcal{D} := \{(x_d, y_d)\}_{n=1}^{D} \in (\mathcal{X} \times \mathcal{Y})^{D}$ , model  $f_{\theta} : \mathcal{X} \ni x_n \mapsto \hat{y}_n \in \mathcal{Y}$  and loss function  $l(y_d, \hat{y}_d) : (\mathcal{Y} \times \mathcal{Y}) \mapsto \mathbb{R}_{\geq 0}$
- ► Local loss landscape:  $\mathcal{L}(\theta+\delta) \approx \hat{\mathcal{L}}(\theta+\delta) = \frac{1}{2}\delta^{\top}H_{\theta}\delta + \nabla_{\theta}^{\top}\delta + \mathcal{L}(\theta)$
- Empirical risk minimization on training data leads to update  $\theta_{t+1} = \theta_t H_{\theta_t}^{-1} \nabla_{\theta_t}$

Often, *D* and  $dim(\theta) = N$  are very large  $\rightarrow$  exact updates unfeasible!

Speedup techniques also generalize well (Hardt et al. 2016), but optimizers often present **slow convergence** and **hyperparameter instability**.





Recall we want to do measurements  $A\Upsilon$  with random  $\Upsilon$ . If we e.g. draw Gaussian noise, *this is prohibitively slow and inefficient*! And what we *really* want is  $\Upsilon$  to be *uncorrelated* with A, but ideally also *orthogonal* to ensure numerical stability (Halko et al. 2011, p. 6.2). This is well characterized via the *Restricted Isometry Property* (Candès 2008).

Luckily, there is a fast, competitive, stable and matrix-free way of doing random measurements! The Scrambled Subsampled Randomized Fourier Transform (SSRFT) (Tropp et al. 2019, p. 3.2), related to the *Fast Johnson-Lindenstrauss Transform (FJLT)*, achieves exactly this:

 $\mathsf{SSRFT:} \ \mathcal{RFII}\mathcal{FII}'$ 

Here,  $\mathcal{F}$  is the Fourier transform,  $\Pi$  are signed random permutations and  $\mathcal{R}$  is a random index picker. Furthermore:

- ▶ It is matrix-free and leverages the FFT, requiring only O(n) memory and  $O(n \log n)$  time
- ▶ It is composed by unitary projections, hence it is unitary
- As a bonus, the adjoint operation is also trivial
- Empirically, it has been shown to outperform Gaussian noise

Fernandez, Schneider, Mahsereci, Hennig – Connecting Parameter Magnitudes and Hessian Eigenspaces at Scale using Sketched Methods (TMLR 2025)

The forced ammendment

- Double-edged sword: Arbitrarily large, but become implicit
- Satisfy linearity via mat-vec products:  $A(c_1v_1 + c_2v_2) = c_1Av_1 + c_2Av_2$
- Examples: differentiation, Laplace, Fourier, convolution...
- This restricts the type of algorithms and analysis available

# both support matrix-vector multiplications
w1 = mat @ v1
w2 = H @ v2

```
# and both are linear
w12 = H @ (a*v1 + b*v2)
assert w12 == a*(H @ v1) + b*(H @ v2)
```

```
# But H is not explicit in memory!
x1 = mat[3, 4:] # good
x2 = H[3, 4:] # not supported!
```

```
# instead, we have to run an HVP:
e3 = one_hot_vector(idx=3)
x2 = (e3 @ H)[4:]
```