

Old Optimizer, New Norm: An Anthology

Jeremy Bernstein

Laker Newhouse

MIT CSAIL, United States

JBERNSTEIN@MIT.EDU

LAKERN@MIT.EDU

Abstract

Deep learning optimizers are often motivated through a mix of convex and approximate second-order theory. We select three such methods—Adam, Shampoo and Prodigy—and argue that each method can instead be understood as a squarely first-order method without convexity assumptions. In fact, after switching off exponential moving averages, each method is equivalent to *steepest descent* under a particular *norm*. By generalizing this observation, we chart a new design space for training algorithms. Different operator norms should be assigned to different tensors based on the role that the tensor plays within the network. For example, while linear and embedding layers may have the same weight space of $\mathbb{R}^{m \times n}$, these layers play different roles and should be assigned different norms. We hope that this idea of carefully metrizing the neural architecture might lead to more stable, scalable and indeed faster training.

Prologue

Deep learning optimizers are often motivated from the perspectives of convex and approximate second-order theory. These theoretical frameworks have been used to inspire algorithmic ideas, as well as providing means to analyse the convergence of various optimizers. However, we believe—and will attempt to demonstrate—that there is a wealth of untapped algorithmic opportunity in the simpler realm of exact first-order theory without convexity assumptions.

To make our case, we choose three optimizers that were originally analysed under convex or approximate second-order theory: Adam, Shampoo and Prodigy. After disabling their exponential moving averages (EMA), we show that each algorithm admits a parsimonious theoretical explanation as a variant of *steepest descent* under a certain norm. EMA can then be thought of as “smoothing out” the algorithm, or making it more robust to mini-batch noise, although nailing down the precise role of EMA is perhaps still an open problem.

By steepest descent, we mean the procedure of choosing a weight update $\Delta \mathbf{w}$ to minimise a local quadratic model of the loss function \mathcal{L} of the form $\mathcal{L}(\mathbf{w}) + \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w})^\top \Delta \mathbf{w} + \frac{\lambda}{2} \cdot \|\Delta \mathbf{w}\|^2$, visualized in Figure 1. Crucially, the *sharpness parameter* λ and *norm* $\|\cdot\|$ are chosen a priori, without touching an (approximate) Hessian during training. As such, we consider steepest descent to be a squarely first-order method and not an (approximate) second-order method.

Throughout the anthology, we rely on a dual description of steepest descent:

Proposition 1 (Steepest descent) For any $\mathbf{g} \in \mathbb{R}^n$ thought of as “the gradient” and any $\lambda \geq 0$ thought of as “the sharpness”, and for any norm $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ with dual norm $\|\cdot\|^\dagger$,

$$\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right] = -\frac{\|\mathbf{g}\|^\dagger}{\lambda} \cdot \arg \max_{\|\mathbf{t}\|=1} \mathbf{g}^\top \mathbf{t}. \quad (1)$$

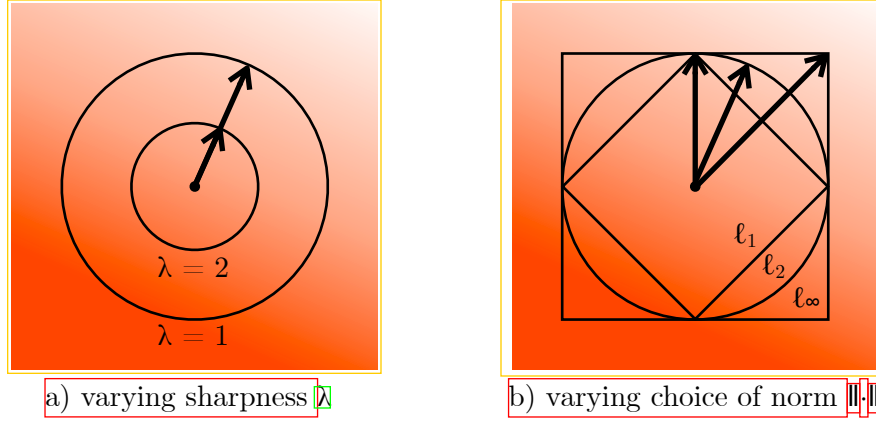


Figure 1: Steepest descent considers the problem of minimizing a linear functional under a quadratic penalty: $\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} [\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2]$ for $\mathbf{g} \in \mathbb{R}^n$. Here we show how the solution varies with the sharpness $\lambda > 0$ and the choice of norm $\|\cdot\|$. We overlay different norm balls on top of a linear color gradient, and use arrows to denote the solution, meaning the member of the norm ball that “minimizes the color”. a) Increasing the sharpness decreases the size of the solution vector. b) Changing the norm can change the direction of the solution vector. For different ℓ_p norms, the solution direction changes because the gradient is not axis-aligned. In practice, we should pick the sharpness and norm to fit the geometry of our loss.

Equation (1) separates the solution of the steepest descent problem into two pieces: first computing the *step size* as the dual norm of the gradient divided by the sharpness, and second solving for the *step direction* as the unit vector that maximizes the inner product with the gradient. The proof of this proposition is given in Appendix B.

Of course, the art of steepest descent lies in choosing a norm $\|\cdot\|$ and a sharpness λ suited to the optimization problem at hand. While it may be possible to turn this art into a science (Large et al., 2024), that ambition is beyond the scope of this anthology. Here we point out that past methods do implicitly make decisions about norms, and in a somewhat haphazard manner. In fact, they implicitly assign different *induced matrix norms* to the network layers:

Definition 1 (Induced operator norm) Given a matrix $\mathbf{M} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ and two normed vector spaces $(\mathbb{R}^{d_{\text{in}}}, \|\cdot\|_\alpha)$ and $(\mathbb{R}^{d_{\text{out}}}, \|\cdot\|_\beta)$, the “ α to β ” induced operator norm is given by:

$$\|\mathbf{M}\|_{\alpha \rightarrow \beta} = \max_{\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}} \frac{\|\mathbf{M}\mathbf{x}\|_\beta}{\|\mathbf{x}\|_\alpha}. \quad (2)$$

Definition 1 tells us that by varying the choice of vector norms $\|\cdot\|_\alpha$ and $\|\cdot\|_\beta$, we can induce a large family of matrix norms. In turn, this implies a correspondingly large family of steepest descent optimizers. By foregrounding this issue, we hope that algorithm designers may develop more suitable optimizers by becoming more intentional about their choice of norm.

Story I. Adam as Steepest Descent under the Max-of-Max Norm

ADAM is a widely used deep learning optimizer: the original paper of [Kingma and Ba \(2015\)](#) now has well over 100,000 citations. Adam has been motivated in various ways, including through convex analysis ([Kingma and Ba, 2015](#)) and as an approximate second-order method ([Sun and Spall, 2021](#)). However, there have been efforts to build a more direct understanding of Adam: for instance, with exponential moving averages (EMA) switched off, Adam is just *sign gradient descent* ([Balles and Hennig, 2018](#); [Bernstein et al., 2018](#)), which is equivalent to steepest descent under the infinity norm ([Carlson et al., 2015](#)). In this story, we connect Adam to a certain “max-of-max” norm, showing how Adam respects the tensor structure of a neural network in a very particular way.

To begin, we review how Adam connects to sign gradient descent. Ignoring bias corrections and numerical stabilizations, Adam is given by the following system of updates:

$$\mathbf{m}_t = \beta_1 \cdot \mathbf{m}_{t-1} + (1 - \beta_1) \cdot \mathbf{g}_t, \quad (3)$$

$$\mathbf{v}_t = \beta_2 \cdot \mathbf{v}_{t-1} + (1 - \beta_2) \cdot \mathbf{g}_t^2, \quad (4)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \cdot \mathbf{m}_t / \sqrt{\mathbf{v}_t}, \quad (5)$$

where t denotes the time step, $\mathbf{g}_t \in \mathbb{R}^n$ the gradient vector and $\eta > 0$ the step size. The EMA time scales of the first gradient moment \mathbf{m}_t and second moment \mathbf{v}_t are set by $0 \leq \beta_1, \beta_2 < 1$. All operations are conducted entry-wise. If we switch off EMA by setting $\beta_1 = \beta_2 = 0$, the Adam updates reduce to just sign gradient descent:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \cdot \mathbf{g}_t / \sqrt{\mathbf{g}_t^2} \quad (6)$$

$$= \mathbf{w}_t - \eta \cdot \text{sign}(\mathbf{g}_t). \quad (7)$$

This connection to sign descent should not be surprising since Adam, published in 2015, builds on the RMSprop optimizer that [Tieleman and Hinton \(2012\)](#) already called “the mini-batch version of just using the sign of the gradient”. And RMSprop itself built on the RPROP optimizer ([Riedmiller and Braun, 1993](#)), which also uses gradient signs.

Still, *why should using the sign of the gradient be a good idea in deep learning?* In search of a motivation, it is interesting to consider that sign descent solves the problem of steepest descent under the vector ℓ_∞ norm, $\|\mathbf{v}\|_\infty := \max_i |\mathbf{v}_i|$ ([Carlson et al., 2015, 2016](#); [Fan, 2017](#)):

Proposition 2 (Sign descent as steepest descent under the infinity norm)

For any gradient vector $\mathbf{g} \in \mathbb{R}^n$ and sharpness $\lambda > 0$, it holds that:

$$\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|_\infty^2 \right] = -\frac{\|\mathbf{g}\|_1}{\lambda} \text{sign}(\mathbf{g}). \quad (8)$$

In words, the vector that minimizes a linear functional under an infinity norm penalty is a scalar multiple of a sign vector. The proof is given in [Appendix B](#).

While this connection between Adam, sign descent and steepest descent is perhaps cute, it does not answer a basic question: *Why does the vector ℓ_∞ norm have anything to do with neural network training?* In particular, taking the weight space to be \mathbb{R}^n equipped with

the simple infinity norm seems to “throw away” the fact that the weight space of a neural network is built in a structured way out of layers of matrices (and perhaps other tensors).

To resolve this conundrum, we suggest that in fact the vector ℓ_∞ norm on the flattened weight space doesn’t have anything to do with deep learning. Instead, there is a coincidence at play. The ℓ_∞ norm enjoys a special property summarized by the slogan “a max of a max is a max”. To see this, consider a neural network with a list of L weight matrices $\mathbf{W}_1, \dots, \mathbf{W}_L$. Let $\text{row}_r(\mathbf{W}_l)$ denote the r th row of the l th weight matrix, and let $\mathbf{w} = \text{flatten}(\mathbf{W}_1, \dots, \mathbf{W}_L) \in \mathbb{R}^n$ denote the full flattened weight vector. Then we have that:

$$\|\mathbf{w}\|_\infty = \max_l \max_r \|\text{row}_r(\mathbf{W}_l)\|_\infty = \max_l \|\mathbf{W}_l\|_{\ell_1 \rightarrow \ell_\infty}, \quad (9)$$

where the second equality follows via Proposition 8. In words, the infinity norm of the flattened weight vector coincides with the largest ℓ_1 to ℓ_∞ operator norm of the layers. So Equation (9) connects the *unstructured* space of the flattened weight vector to the *structured* space of the list of weight matrices. We refer to the object $\max_l \|\mathbf{W}_l\|_{\ell_1 \rightarrow \ell_\infty}$ as the “max-of-max norm”. And sign descent emerges as steepest descent under this norm:

Proposition 3 (Sign descent as steepest descent under the max-of-max norm)

For any list of gradient matrices $\mathbf{G}_1, \dots, \mathbf{G}_L$ and any sharpness $\lambda > 0$, consider the problem:

$$\arg \min_{\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L} \left[\sum_{l=1}^L \langle \mathbf{G}_l, \Delta \mathbf{W}_l \rangle + \frac{\lambda}{2} \max_{l=1}^L \|\Delta \mathbf{W}_l\|_{\ell_1 \rightarrow \ell_\infty}^2 \right], \quad (10)$$

where $\langle \cdot, \cdot \rangle$ denotes the Frobenius inner product, and $\Delta \mathbf{W}_l$ has the same shape as \mathbf{G}_l . For step size $\eta = \frac{1}{\lambda} \sum_{l=1}^L \|\mathbf{G}_l\|_{\ell_1 \rightarrow \ell_\infty}^\dagger$, where \dagger denotes the dual norm, Equation (10) is solved by:

$$\Delta \mathbf{W}_l = -\eta \cdot \text{sign}(\mathbf{G}_l) \quad \text{for each layer } l = 1, \dots, L. \quad (11)$$

In words, the matrix-aware steepest descent problem of Equation (10) is solved by layerwise sign descent as given in Equation (11). This observation—that sign descent updates are implicitly doing *per-matrix gradient normalization*—may be a major reason that Adam, sign descent and Lion (Chen et al., 2023) outperform vanilla gradient descent in large language model training (Zhao et al., 2024; Large et al., 2024). The proof is given in Appendix B.



All told, this story has shown that Adam without EMA is sign descent and that, coincidentally, sign descent solves two different steepest descent problems: one on the flattened weight space, and one that is aware of the matrix structure of neural architecture. But, at the end of this story, questions linger. *Why does the ℓ_1 to ℓ_∞ induced operator norm rear its head? What does it have to do with deep learning? Aren’t there other induced operator norms on matrices we could equally well consider?* For answers to these questions, dear reader, you’ll have to wait for our next story... a story about Shampoo!

Story II. Shampoo as Steepest Descent under the Spectral Norm

Now, dear reader, we turn our attention to Shampoo (Gupta et al., 2017, 2018). A variant of the Shampoo optimizer won the external tuning track of the 2024 *AlgoPerf: Training Algorithms* competition (Dahl et al., 2023). While the method was originally motivated as a generalization of the AdaGrad convex optimizer (Duchi et al., 2011) to tensor spaces, more recent work casts Shampoo as an approximate second-order method (Anil et al., 2020; Morwani et al., 2024). We will show that Shampoo—with accumulation disabled—is steepest descent under the max spectral norm over layers.

To begin, we show that Shampoo updates, without accumulation, are semi-orthogonal matrices. At time step t and for each layer, Shampoo collects the gradient matrix \mathbf{G}_t and makes the following update to the weight matrix \mathbf{W}_t :

$$\mathbf{L}_t = \mathbf{L}_{t-1} + \mathbf{G}_t \mathbf{G}_t^\top, \quad (12)$$

$$\mathbf{R}_t = \mathbf{R}_{t-1} + \mathbf{G}_t^\top \mathbf{G}_t, \quad (13)$$

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \cdot \mathbf{L}_t^{-1/4} \mathbf{G}_t \mathbf{R}_t^{-1/4}. \quad (14)$$

All operations, including the inverse fourth roots, are matrix operations. The accumulators \mathbf{L}_t and \mathbf{R}_t are referred to as the “left and right pre-conditioners”. Practitioners usually replace the simple sums in Equations (12) and (13) with EMAs (Shi et al., 2023). If we disable the accumulation, setting $\mathbf{L}_t = \mathbf{G}_t \mathbf{G}_t^\top$ and $\mathbf{R}_t = \mathbf{G}_t^\top \mathbf{G}_t$, Shampoo reduces to:

$$\mathbf{W}_{t+1} = \mathbf{W}_t - \eta \cdot (\mathbf{G}_t \mathbf{G}_t^\top)^{-1/4} \mathbf{G}_t (\mathbf{G}_t^\top \mathbf{G}_t)^{-1/4} \quad (15)$$

$$= \mathbf{W}_t - \eta \cdot \mathbf{U}_t \mathbf{V}_t^\top, \quad (16)$$

where Equation (16) is reached by substituting the reduced singular value decomposition (SVD) of the gradient $\mathbf{G}_t = \mathbf{U}_t \mathbf{\Sigma}_t \mathbf{V}_t^\top$ into Equation (15). Notice that there is a direct parallel between Equations (6) and (7) for Adam and Equations (15) and (16) for Shampoo. So, Shampoo without accumulation makes a *semi-orthogonal weight update*. In fact:

Proposition 4 (Projection to the closest semi-orthogonal matrix) *Consider the semi-orthogonal matrices $\mathcal{O}_{m \times n} := \{\mathbf{A} \in \mathbb{R}^{m \times n} : \mathbf{A} \mathbf{A}^\top = \mathbf{I}_m \text{ or } \mathbf{A}^\top \mathbf{A} = \mathbf{I}_n\}$ and let $\|\cdot\|_F$ denote the Frobenius norm. For any matrix $\mathbf{G} \in \mathbb{R}^{m \times n}$ with reduced SVD $\mathbf{G} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$:*

$$\arg \min_{\mathbf{A} \in \mathcal{O}_{m \times n}} \|\mathbf{A} - \mathbf{G}\|_F = \mathbf{U} \mathbf{V}^\top, \quad (17)$$

where the minimizer $\mathbf{U} \mathbf{V}^\top$ is unique if and only if the matrix \mathbf{G} has full rank.

So, Shampoo without accumulation projects the gradient matrix to the *closest semi-orthogonal matrix in Frobenius norm*. The proof is in Appendix B. Why might this be a good idea, you ask? Well, for one thing, it’s steepest descent—this time under the maximum spectral norm $\|\cdot\|_{\ell_2 \rightarrow \ell_2}$ (Definition 1) over all the matrices in the network:

Proposition 5 (Shampoo as steepest descent under the spectral norm) *For any list of gradient matrices $\mathbf{G}_1, \dots, \mathbf{G}_L$ and any sharpness $\lambda > 0$, consider the problem:*

$$\arg \min_{\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L} \left[\sum_{l=1}^L \langle \mathbf{G}_l, \Delta \mathbf{W}_l \rangle + \frac{\lambda}{2} \max_{l=1}^L \|\Delta \mathbf{W}_l\|_{\ell_2 \rightarrow \ell_2}^2 \right], \quad (18)$$

where $\langle \cdot, \cdot \rangle$ denotes the Frobenius inner product and $\Delta \mathbf{W}_l$ has the same shape as \mathbf{G}_l . Suppose that \mathbf{G}_l has reduced SVD given by $\mathbf{G}_l = \mathbf{U}_l \mathbf{\Sigma}_l \mathbf{V}_l^\top$ for each $l = 1, \dots, L$. Then Equation (18) is solved with a step size $\eta = \frac{1}{\lambda} \sum_{l=1}^L \text{tr } \mathbf{\Sigma}_l$ and an update:

$$\Delta \mathbf{W}_l = -\eta \cdot \mathbf{U}_l \mathbf{V}_l^\top \quad \text{for each } l = 1, \dots, L. \quad (19)$$

This solution for $\Delta \mathbf{W}_l$ is unique if and only if the matrix \mathbf{G}_l is of full rank.

The proof is given in Appendix B. A novelty of this proposition in contrast to prior work on stochastic spectral descent (Carlson et al., 2015, 2016) is our use of a max norm over layers to handle the multi-layer case. However, our main contribution here is to draw the connection between Proposition 5 and Shampoo as in Equations (15) and (16).

So, Shampoo without accumulation is steepest descent under the spectral norm. *Why might this be a good idea in deep learning?* The idea that we wish to advance is that one can derive upper bounds on the loss of machine learning models in terms of spectral norms. Here we present the simplest possible example: a linear model and the square loss.

Proposition 6 (Bounding the square loss of a linear predictor) Consider a matrix $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ that we shall think of as a linear predictor mapping an input $\mathbf{x} \in \mathbb{R}^{d_{\text{in}}}$ to an output $\mathbf{y} = \mathbf{W}\mathbf{x} \in \mathbb{R}^{d_{\text{out}}}$. Given a dataset of n samples $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$, where the i th input is normalized such that $\|\mathbf{x}_i\|_2 = \sqrt{d_{\text{in}}}$, we can construct the “square loss”:

$$\mathcal{L}(\mathbf{W}) := \frac{1}{2n} \sum_{i=1}^n \frac{1}{d_{\text{out}}} \|\mathbf{y}_i - \mathbf{W}\mathbf{x}_i\|_2^2. \quad (20)$$

Then, for any matrix $\Delta \mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ thought of as a weight update, it holds that:

$$\mathcal{L}(\mathbf{W} + \Delta \mathbf{W}) \leq \mathcal{L}(\mathbf{W}) + \langle \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}), \Delta \mathbf{W} \rangle + \frac{1}{2} \cdot \frac{d_{\text{in}}}{d_{\text{out}}} \cdot \|\Delta \mathbf{W}\|_{\ell_2 \rightarrow \ell_2}^2, \quad (21)$$

where $\langle \cdot, \cdot \rangle$ is the Frobenius inner product.

In words: the square loss of a linear predictor admits an upper bound that is quadratic in the spectral norm of the weight perturbation. Choosing the weight perturbation to minimize this upper bound is precisely steepest descent under the spectral norm! The proof is given in Appendix B. This optimizer design pattern, which starts by deriving an upper bound on the loss (as in Proposition 6) and then minimizes it (as in Proposition 5), is known generally as majorization-minimization (Lange, 2016). It is an exact and first-principles design pattern, without Hessian approximations or appeals to convex theory. This design pattern is used extensively by Carlson et al. (2015, 2016) to design optimizers for restricted Boltzmann machines and discrete graphical models. Generalizing the pattern to arbitrary network architectures and loss functions requires more advanced machinery (Bernstein et al., 2023; Streeter, 2023; Large et al., 2024).



And so, dear reader, we have reached the end of our second story. We have shown that Shampoo without accumulation corresponds to projecting the gradient matrix to the closest

Domain	Norm	Solution	Optimizer	Cousin
\mathbb{R}^n	Euclidean ℓ_2	$\Delta \mathbf{w} = -\frac{\ \mathbf{g}\ _2}{\lambda} \frac{\mathbf{g}}{\ \mathbf{g}\ _2}$	vanilla gradient descent	SGD
\mathbb{R}^n	infinity ℓ_∞	$\Delta \mathbf{w} = -\frac{\ \mathbf{g}\ _1}{\lambda} \text{sign}(\mathbf{g})$	sign descent	Adam
$\mathbb{R}^{m \times n}$	Frobenius S_2	$\Delta \mathbf{W} = -\frac{\ \mathbf{G}\ _F}{\lambda} \frac{\mathbf{G}}{\ \mathbf{G}\ _F}$	vanilla gradient descent	SGD
$\mathbb{R}^{m \times n}$	spectral S_∞	$\Delta \mathbf{W} = -\frac{\text{tr} \mathbf{\Sigma}}{\lambda} \mathbf{U} \mathbf{V}^\top$	spectral descent	Shampoo

Table 1: Popular optimizers are related to steepest descent under different norms. For vector-valued optimization problems, we consider the steepest descent problem $\arg \min_{\Delta \mathbf{w}} \mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \cdot \|\Delta \mathbf{w}\|_2^2$. For matrix-valued problems, we consider $\arg \min_{\Delta \mathbf{W}} \langle \mathbf{G}, \Delta \mathbf{W} \rangle + \frac{\lambda}{2} \cdot \|\Delta \mathbf{W}\|_2^2$, where $\langle \cdot, \cdot \rangle$ is the Frobenius inner product. We list the solution for different vector ℓ_p norms and Schatten S_p norms. The Schatten S_p norm of a matrix returns the ℓ_p norm of its vector of singular values. Finally, $\mathbf{G} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ is the reduced singular value decomposition of the gradient.

semi-orthogonal matrix, which solves the problem of steepest descent under the spectral norm. And we showed how steepest descent under the spectral norm emerges from upper bounding the square loss of a linear predictor. This perspective, of viewing Shampoo as a (smoothed out) projection to the space of semi-orthogonal matrices, grounds the algorithm in a prior literature on spectral descent (Carlson et al., 2015, 2016; Fan, 2017). And in Appendix A, we discuss how it might unlock new means for computing the Shampoo updates.

We summarize our first two stories in Table 1. And we still have one more left to tell...

Story III. Prodigy: Automatically Computing the Escape Velocity

FOR our final story, we speak of Prodigy (Mishchenko and Defazio, 2023). The Prodigy optimizer falls amid a series of recent works (Defazio and Mishchenko, 2023; Khaled et al., 2023; Ivgi et al., 2023) that attempt to apply convex theory to design and analyse deep learning optimizers that do not require tuning. In contrast, we argue that Prodigy (without EMA) is but another example of steepest descent, where instead of using the step size $\eta = \|\mathbf{g}\|^\dagger / \lambda$ from Proposition 1, Prodigy uses a heuristic to automatically warm up to a good step size. This demonstrates the value of Proposition 1 for disentangling the optimizer design problem. If one knows a good norm $\|\cdot\|$ but is ignorant of the sharpness parameter λ , then one may obtain the step direction by solving $\arg \max_{\|\mathbf{t}\|=1} \mathbf{g}^\top \mathbf{t}$ from Proposition 1, while using another means to find a good step size.

Then let us make our case. We focus on Algorithm 3 in the Prodigy paper, since this is the version used in their experiments. We first show that with EMA switched off, Prodigy implements sign gradient descent with a step size that warms up automatically. Ignoring the numerical stabilization and learning rate schedule, Prodigy is given by:

$$\mathbf{m}_t = \beta_1 \cdot \mathbf{m}_{t-1} + (1 - \beta_1) \cdot \eta_t \mathbf{g}_t, \quad (22)$$

$$\mathbf{v}_t = \beta_2 \cdot \mathbf{v}_{t-1} + (1 - \beta_2) \cdot \eta_t^2 \mathbf{g}_t^2, \quad (23)$$

$$\mathbf{r}_t = \sqrt{\beta_2} \cdot \mathbf{r}_{t-1} + (1 - \sqrt{\beta_2}) \cdot \eta_t^2 \mathbf{g}_t^\top (\mathbf{w}_0 - \mathbf{w}_t), \quad (24)$$

$$\mathbf{s}_t = \sqrt{\beta_2} \cdot \mathbf{s}_{t-1} + (1 - \sqrt{\beta_2}) \cdot \eta_t^2 \mathbf{g}_t, \quad (25)$$

$$\eta_{t+1} = \max \left(\eta_t, \frac{r_t}{\|\mathbf{s}_t\|_1} \right), \quad (26)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \cdot \mathbf{m}_t / \sqrt{\mathbf{v}_t}, \quad (27)$$

where t denotes the time step and $\mathbf{g}_t \in \mathbb{R}^n$ the gradient vector. While this system of updates may seem intimidating, if we switch off EMA by setting $\beta_1 = \beta_2 = 0$, the Prodigy updates simplify dramatically to just sign gradient descent with a dynamical step size as follows:

$$\eta_{t+1} = \max \left(\eta_t, \frac{\mathbf{g}_t^\top (\mathbf{w}_0 - \mathbf{w}_t)}{\|\mathbf{g}_t\|_1} \right), \quad (28)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \cdot \text{sign}(\mathbf{g}_t). \quad (29)$$

But Proposition 2 showed that sign descent is steepest descent under the infinity norm. Therefore Equations (28) and (29) prove our claim that Prodigy without EMA is steepest descent, although with a dynamically chosen step size denoted η_t .

All that remains is to understand the dynamical rule, given by Equation (28), for choosing the step size η_t . We shall argue that this dynamical rule can be understood to approximate a heuristic algorithm for achieving, but not exceeding, what we shall call *escape velocity*:

- Choose a very small initial step size η_0 —small enough to be *a priori* sure that $\eta_0 \ll \eta_*$, where η_* denotes *escape velocity*: the unknown but optimal initial step size;
- At each step, check if the weights \mathbf{w}_t have escaped the linearization of the loss around the initial weights \mathbf{w}_0 —if not, double the step size according to $\eta_{t+1} = 2 \times \eta_t$;
- Once the weights \mathbf{w}_t have escaped the initial linearization, stop increasing the step size. We say that the step size η_t has reached escape velocity η_* .

The rationale behind this procedure is that if we knew the optimal initial step size η_* , then the weights should escape the initial linearization of the loss in a single step. Formally, the directional derivative $(\mathbf{w}_1 - \mathbf{w}_0)^\top \mathbf{g}_1$ must vanish if the step size is chosen optimally (Cauchy, 1847). If the directional derivative in the direction of the first weight update is still negative $(\mathbf{w}_1 - \mathbf{w}_0)^\top \mathbf{g}_1 < 0$, then we could have taken a larger step. Said another way, we can use the angle that the gradient \mathbf{g}_1 makes with the change in weights $\mathbf{w}_1 - \mathbf{w}_0$ to tell us whether or not we should increase the step size. Notice that procedure has no reliance on convexity. With this in mind, let us massage Prodigy’s step size update (Equation (28)) as follows:

$$\eta_{t+1} = \max \left(\eta_t, \frac{\mathbf{g}_t^\top (\mathbf{w}_0 - \mathbf{w}_t)}{\|\mathbf{g}_t\|_1} \right) = \max \left(\eta_t, \frac{\|\mathbf{g}_t\|_2}{\|\mathbf{g}_t\|_1} \times \|\mathbf{w}_t - \mathbf{w}_0\|_2 \times \cos \theta \right), \quad (30)$$

where θ denotes the angle between the gradient \mathbf{g}_t and the difference in weights $\mathbf{w}_0 - \mathbf{w}_t$. To help make sense of this expression, we make two assumptions:

1. The gradient is a “dense” vector in \mathbb{R}^n , meaning that $\|\mathbf{g}_t\|_2 / \|\mathbf{g}_t\|_1 \approx 1/\sqrt{n}$;
2. \mathbf{w}_t is still close enough to the initialization \mathbf{w}_0 that $\cos \theta \approx 1$.

Under these assumptions, Equation (30) becomes just $\eta_{t+1} \approx \max(\eta_t, \|\mathbf{w}_t - \mathbf{w}_0\|_{\text{RMS}})$, where the root mean square (RMS) norm is defined via $\|\cdot\|_{\text{RMS}} := \frac{1}{\sqrt{n}} \|\cdot\|_2$. Combined with Equation (29), this allows us to estimate the size of the weight change at step $t+1$:

$$\|\mathbf{w}_{t+2} - \mathbf{w}_{t+1}\|_{\text{RMS}} = \eta_{t+1} \cdot \|\text{sign}(\mathbf{g}_t)\|_{\text{RMS}} \approx \max(\eta_t, \|\mathbf{w}_t - \mathbf{w}_0\|_{\text{RMS}}) \geq \|\mathbf{w}_t - \mathbf{w}_0\|_{\text{RMS}},$$

where we have used the fact that a sign vector has unit RMS norm. In words, while assumptions (1) and (2) hold, the step size at time $t+1$ is equivalent to the whole progress up to step t . This suggests exponential growth in the step size that continues until assumption (2) breaks, which we think of as the step size reaching the escape velocity η_* .

Now we wish to point out that this procedure is just one amongst a space of line search methods that one might consider (Armijo, 1966; Riedmiller and Braun, 1993; Kenneweg et al., 2024). For instance, Prodigy’s decision to only let η_t increase and never decrease could be sub-optimal. And the decision to measure the angle between the gradient and the weight difference $\mathbf{w}_t - \mathbf{w}_0$ has alternatives. One could instead use the most recent weight difference $\mathbf{w}_t - \mathbf{w}_{t-1}$. Lastly, in place of relying on the norm ratio $\|\mathbf{g}_t\|_2 / \|\mathbf{g}_t\|_1$ to implicitly convert the ℓ_2 norm $\|\mathbf{w}_t - \mathbf{w}_0\|_2$ into the RMS norm $\|\mathbf{w}_t - \mathbf{w}_0\|_{\text{RMS}}$, one could consider a more explicit method. For instance, we found a rule akin to $\eta_{t+1} = \eta_t \times (1 + \cos \theta)$ to work well in some preliminary experiments.



Our time grows short, dear reader, and our third story draws to an end. We have argued that Prodigy without EMA is sign descent—an example of steepest descent—with a particular mechanism for warming up the step size. Starting with a tiny initial step size, Prodigy multiplicatively increases the step size until the weights escape the initial locally linear region of the loss. Prodigy’s step size adjustment is based on the angle between the gradient and the total weight change. This is a form of online line search. This highlights that once one has chosen a norm, the steepest descent framework allows freedom to estimate the step size in various different ways.

Epilogue

This anthology has presented new ways of understanding old optimizers. [Proposition 1](#) decouples the optimizer design problem into two pieces: first choosing a norm and second finding a step size. This design space is already broad. We have argued that Adam chooses the infinity norm ([Proposition 2](#)) or equivalently the max-of-max norm ([Proposition 3](#)), which respects a layered matrix structure. Shampoo chooses the spectral norm ([Proposition 5](#)). Prodigy chooses the same norm as Adam, and then uses a heuristic to automatically warm up to a good step size, as in [Equation \(28\)](#), which we term *reaching escape velocity*.

Through the lens of steepest descent, the decisions that Adam, Shampoo and Prodigy make may seem arbitrary. In fact, we think that they *are* somewhat arbitrary. And there may be more principled ways to make these decisions. To demonstrate this point, we now introduce a tool called the modular norm ([Large et al., 2024](#)) and its corresponding steepest descent algorithm. The modular norm generalizes the norms that appeared in [Proposition 3](#) for Adam and [Proposition 5](#) for Shampoo. Formally:

Proposition 7 (Steepest descent under the modular norm) *Given scalar coefficients $s_1, \dots, s_L > 0$ and norms $\|\cdot\|_1, \dots, \|\cdot\|_L$, we define the modular norm as the mapping:*

$$\mathbf{W}_1, \dots, \mathbf{W}_L \mapsto \max \{s_1 \|\mathbf{W}_1\|_1, \dots, s_L \|\mathbf{W}_L\|_L\}. \quad (31)$$

The corresponding steepest descent problem is given by:

$$\arg \min_{\Delta \mathbf{W}_1, \dots, \Delta \mathbf{W}_L} \left[\sum_{l=1}^L \langle \mathbf{G}_l, \Delta \mathbf{W}_l \rangle + \frac{\lambda}{2} \max_{l=1}^L s_l^2 \|\Delta \mathbf{W}_l\|_l^2 \right], \quad (32)$$

where $\langle \cdot, \cdot \rangle$ denotes the Frobenius inner product, and for each $l = 1, \dots, L$ the two matrices $\Delta \mathbf{W}_l$ and \mathbf{G}_l are of the same shape. If we define the global step size $\eta = \frac{1}{\lambda} \sum_{k=1}^L \frac{1}{s_k} \|\mathbf{G}_k\|_k^2$, then the solution to [Equation \(32\)](#) is given by:

$$\Delta \mathbf{W}_l = -\frac{\eta}{s_l} \cdot \arg \max_{\|\mathbf{T}_l\|_l=1} \langle \mathbf{G}_l, \mathbf{T}_l \rangle \quad \text{for each layer } l = 1, \dots, L. \quad (33)$$

In words, steepest descent under the modular norm updates each layer in a direction informed by that layer's norm and with a global step size computed as a weighted sum of the dual norms of the gradients over layers. The proof of this proposition is given in [Appendix B](#).

When confronted with the modular norm, it's natural to ask how one should assign norms to layers. And there are so many norms to choose from! Beyond the familiar $\ell_2 \rightarrow \ell_2$ spectral norm, many other induced operator norms are computationally tractable:

Proposition 8 ($\ell_1 \rightarrow \ell_p$ and $\ell_p \rightarrow \ell_\infty$ induced operator norms are tractable) *For a matrix $\mathbf{M} \in \mathbb{R}^{m \times n}$ with m rows $\{\text{row}_i(\mathbf{M})\}_{i=1}^m$ and n columns $\{\text{col}_j(\mathbf{M})\}_{j=1}^n$, and $1 \leq p \leq \infty$:*

$$\|\mathbf{M}\|_{\ell_1 \rightarrow \ell_p} = \max_j \|\text{col}_j(\mathbf{M})\|_p; \quad \|\mathbf{M}\|_{\ell_p \rightarrow \ell_\infty} = \max_i \|\text{row}_i(\mathbf{M})\|_{\frac{p}{p-1}}. \quad (34)$$

In words, the $\ell_1 \rightarrow \ell_p$ operator norm is the largest ℓ_p norm of the columns; the $\ell_p \rightarrow \ell_\infty$ operator norm is the largest dual ℓ_p norm over the rows. The proof is given in [Appendix B](#).

To assign a norm to a layer, we believe that one should consider the role that layer plays in the neural network. For instance, since linear layers are typically used to map to and from vectors with roughly unit RMS norm, it is appropriate to equip linear layers with the induced RMS to RMS operator norm (Yang et al., 2023), which resolves to a rescaled spectral norm. And since embedding layers map from one-hot vectors to vectors with roughly unit RMS norm, it is appropriate to equip embedding layers with the ℓ_1 to RMS operator norm, which resolves to a rescaled ℓ_1 to ℓ_2 operator norm. So embedding layers and linear layers should be equipped with different norms despite the weight space being a matrix space in both cases. In short, the algorithm designer has freedom to choose input and output norms for layers that capture differences in how the layers are used; inducing the corresponding operator norm on the layer’s weights provides control over how the optimizer learns representations.

We believe that picking the right norms could improve the speed and scalability of neural network training. We are seeing evidence that equipping neural network layers with better norms can lead to learning rate transfer across scale (Yang et al., 2023; Large et al., 2024). And since Shampoo won the external tuning track of the 2024 AlgoPerf competition (Dahl et al., 2023), it is garnering interest as a fast training method. The second story in our anthology shows that Shampoo is closely connected to the spectral norm.

In conclusion, this work highlights a perspective on optimizer design as choosing two things: a *norm* and a *step size*. We have shown that three popular methods—Adam, Shampoo and Prodigy—fit within this perspective. We hope that researchers can design improved training algorithms by choosing norms and step sizes more intentionally.

“Though this be madness, yet there is method in’t.”

Hamlet

Acknowledgements

We are grateful to Tim Large and Phillip Isola for invaluable discussions on the stories in this anthology. We also thank Jack Gallagher, Keller Jordan and Victor Butoi for very helpful conversations.

References

- Rohan Anil, Vineet Gupta, Tomer Koren, Kevin Regan, and Yoram Singer. Scalable second order optimization for deep learning. *arXiv:2002.09018*, 2020. Cited on pages 5 and 15.
- Larry Armijo. Minimization of functions having Lipschitz continuous first partial derivatives. *Pacific Journal of Mathematics*, 1966. Cited on page 9.
- Lukas Balles and Philipp Hennig. Dissecting Adam: The sign, magnitude and variance of stochastic gradients. In *International Conference on Machine Learning*, 2018. Cited on page 3.
- Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signSGD: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, 2018. Cited on page 3.
- Jeremy Bernstein, Chris Mingard, Kevin Huang, Navid Azizan, and Yisong Yue. Automatic Gradient Descent: Deep Learning without Hyperparameters. *arXiv:2304.05187*, 2023. Cited on page 6.
- David Carlson, Volkan Cevher, and Lawrence Carin. Stochastic spectral descent for Restricted Boltzmann Machines. In *International Conference on Artificial Intelligence and Statistics*, 2015. Cited on pages 3, 6, and 7.
- David Carlson, Ya-Ping Hsieh, Edo Collins, Lawrence Carin, and Volkan Cevher. Stochastic spectral descent for discrete graphical models. *Selected Topics in Signal Processing*, 2016. Cited on pages 3, 6, and 7.
- Augustin-Louis Cauchy. Méthode générale pour la résolution des systèmes d’équations simultanées. *Comptes Rendus Hebdomadaires des Séances de l’Académie des Sciences*, 1847. Cited on page 9.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V Le. Symbolic discovery of optimization algorithms. In *Neural Information Processing Systems*, 2023. Cited on page 4.
- George E. Dahl, Frank Schneider, Zachary Nado, Naman Agarwal, Chandramouli Shama Sastry, Philipp Hennig, Sourabh Medapati, Runa Eschenhagen, Priya Kasimbeg, Daniel Suo, Juhan Bae, Justin Gilmer, Abel L. Peirson, Bilal Khan, Rohan Anil, Mike Rabbat, Shankar Krishnan, Daniel Snider, Ehsan Amid, Kongtao Chen, Chris J. Maddison, Rakshith Vasudev, Michal Badura, Ankush Garg, and Peter Mattson. Benchmarking neural network training algorithms. *arXiv:2306.07179*, 2023. Cited on pages 5 and 11.
- Aaron Defazio and Konstantin Mishchenko. Learning-rate-free learning by D-adaptation. In *International Conference on Machine Learning*, 2023. Cited on page 8.
- John C. Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal Machine Learning Research*, 2011. Cited on page 5.

- Kai Fan. Unifying the stochastic spectral descent for Restricted Boltzmann Machines with Bernoulli or Gaussian inputs. *arXiv:1703.09766*, 2017. Cited on pages 3 and 7.
- Vladimir Feinberg, Xinyi Chen, Y. Jennifer Sun, Rohan Anil, and Elad Hazan. Sketchy: Memory-efficient adaptive regularization with frequent directions. In *Neural Information Processing Systems*, 2023. Cited on page 15.
- Vineet Gupta, Tomer Koren, and Yoram Singer. A unified approach to adaptive regularization in online and stochastic optimization. Technical report, Google Brain, 2017. Cited on page 5.
- Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *International Conference on Machine Learning*, 2018. Cited on page 5.
- Nicholas J. Higham. *Functions of Matrices*. Society for Industrial and Applied Mathematics, 2008. Cited on page 15.
- Maor Ivgi, Oliver Hinder, and Yair Carmon. DoG is SGD’s best friend: A parameter-free dynamic step size schedule. In *International Conference on Machine Learning*, 2023. Cited on page 8.
- Philip Kenneweg, Tristan Kenneweg, and Barbara Hammer. Improving line search methods for large scale neural network training. In *International Conference on Artificial Intelligence, Computer, Data Sciences and Applications*, 2024. Cited on page 9.
- Ahmed Khaled, Konstantin Mishchenko, and Chi Jin. DoWG unleashed: An efficient universal parameter-free gradient descent method. In *Neural Information Processing Systems*, 2023. Cited on page 8.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. Cited on page 3.
- Slobodan Lakić. On the computation of the matrix k-th root. *Journal of Applied Mathematics and Mechanics*, 1998. Cited on page 15.
- Kenneth Lange. *MM Optimization Algorithms*. Society for Industrial and Applied Mathematics, 2016. Cited on page 6.
- Tim Large, Yang Liu, Minyoung Huh, Hyojin Bahng, Phillip Isola, and Jeremy Bernstein. Scalable optimization in the modular norm. *arXiv:2405.14813*, 2024. Cited on pages 2, 4, 6, 10, and 11.
- Per-Gunnar Martinsson and Joel A. Tropp. Randomized numerical linear algebra: Foundations and algorithms. *Acta Numerica*, 2020. Cited on page 15.
- Konstantin Mishchenko and Aaron Defazio. Prodigy: An expeditiously adaptive parameter-free learner. *arXiv:2306.06101*, 2023. Cited on page 8.
- Depen Morwani, Itai Shapira, Nikhil Vyas, Eran Malach, Sham Kakade, and Lucas Janson. A new perspective on Shampoo’s preconditioner. *arXiv:2406.17748*, 2024. Cited on page 5.

Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning: The RPROP algorithm. In *International Conference on Neural Networks*, 1993. Cited on pages 3 and 9.

Hao-Jun Michael Shi, Tsung-Hsien Lee, Shintaro Iwasaki, Jose Gallego-Posada, Zhijing Li, Kaushik Rangadurai, Dheevatsa Mudigere, and Michael Rabbat. A distributed data-parallel PyTorch implementation of the distributed Shampoo optimizer for training neural networks at-scale. *arXiv:2309.06497*, 2023. Cited on page 5.

Matthew Streeter. Universal majorization-minimization algorithms. *arXiv:2308.00190*, 2023. Cited on page 6.

Shiqing Sun and James C. Spall. Connection of diagonal Hessian estimates to natural gradients in stochastic optimization. In *Information Sciences and Systems*, 2021. Cited on page 3.

Tijmen Tieleman and Geoffrey Hinton. RMSprop. *Coursera: Neural Networks for Machine Learning*, Lecture 6.5, 2012. Cited on page 3.

Greg Yang, James B. Simon, and Jeremy Bernstein. A spectral condition for feature learning. *arXiv:2310.17813*, 2023. Cited on page 11.

Rosie Zhao, Depen Morwani, David Brandfonbrener, Nikhil Vyas, and Sham Kakade. Deconstructing what makes a good optimizer for language models. *arXiv:2407.07972*, 2024. Cited on page 4.

Appendix A. Computational Strategies for Shampoo

Let $\mathbf{G} \in \mathbb{R}^{m \times n}$ be a gradient matrix with reduced SVD $\mathbf{G} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$. By Equations (15) and (16), the corresponding Shampoo update (with EMA disabled) is given by:

$$\Delta \mathbf{W} = -\eta \cdot (\mathbf{G} \mathbf{G}^\top)^{-1/4} \mathbf{G} (\mathbf{G}^\top \mathbf{G})^{-1/4} = -\eta \cdot \mathbf{U} \mathbf{V}^\top. \quad (35)$$

Here we list every means we know of computing or approximating this equation. First, we mention that $(\mathbf{G} \mathbf{G}^\top)^{-1/4} \mathbf{G} (\mathbf{G}^\top \mathbf{G})^{-1/4} = (\mathbf{G} \mathbf{G}^\top)^{-1/2} \mathbf{G} = \mathbf{G} (\mathbf{G}^\top \mathbf{G})^{-1/2}$, so if one is willing to compute inverse matrix roots, one need only compute either $(\mathbf{G} \mathbf{G}^\top)^{-1/2}$ or $(\mathbf{G}^\top \mathbf{G})^{-1/2}$, whichever has smaller dimension. With that said, to compute Equation (35), one may:

1. **Do the SVD.** Apply an SVD routine to compute \mathbf{U} , $\mathbf{\Sigma}$ and \mathbf{V}^\top and just discard $\mathbf{\Sigma}$.
2. **Do sketching.** Sketching is a randomized method (Martinsson and Tropp, 2020) that can be used to approximate the SVD. See, for instance, Sketchy (Feinberg et al., 2023).
3. **Do Newton iteration for inverse p th roots.** Inverse matrix roots such as $(\mathbf{G} \mathbf{G}^\top)^{-1/2}$ can be computed via Newton iteration (Lakić, 1998). This is discussed in Chapter 7 of Higham (2008)’s book. And see Anil et al. (2020)’s paper.
4. **Do Newton-Schulz iteration.** We developed a “Newton-Schulz iteration” for computing $\mathbf{U} \mathbf{V}^\top$, adapted from Equation 5.22 in Higham (2008)’s book. In short, if we set $\mathbf{X}_0 = \mathbf{G} / \|\mathbf{G}\|_{\ell_2 \rightarrow \ell_2}$ (or alternatively $\mathbf{X}_0 = \mathbf{G} / \|\mathbf{G}\|_F$) and iterate:

$$\mathbf{X}_{t+1} = \frac{3}{2} \cdot \mathbf{X}_t - \frac{1}{2} \cdot \mathbf{X}_t \mathbf{X}_t^\top \mathbf{X}_t, \quad (36)$$

then as $t \rightarrow \infty$, the sequence $\mathbf{X}_t \rightarrow \mathbf{U} \mathbf{V}^\top$. To see this, one should plot the univariate cubic function $f(x) := \frac{3}{2} \cdot x - \frac{1}{2} \cdot x^3$ and see that, for $0 < x < \sqrt{3}$, iterating this cubic will push x closer and closer to 1 . The final step is to realize that the effect of the iteration in Equation (36) is to apply this cubic $f(x)$ to each singular value of \mathbf{X}_t . This also shows that the spectral normalization $\mathbf{X}_0 = \mathbf{G} / \|\mathbf{G}\|_{\ell_2 \rightarrow \ell_2}$ is stronger than what is required: we need only ensure that \mathbf{X}_0 has all singular values greater than zero and less than $\sqrt{3}$ in order for the iteration to converge.

Finally, there are in fact a family of degree $2n+1$ polynomial iterations of the form

$$\mathbf{X}_{t+1} = a \cdot \mathbf{X}_t + b \cdot \mathbf{X}_t \mathbf{X}_t^\top \mathbf{X}_t + c \cdot (\mathbf{X}_t \mathbf{X}_t^\top)^2 \mathbf{X}_t + \dots + z \cdot (\mathbf{X}_t \mathbf{X}_t^\top)^n \mathbf{X}_t \quad (37)$$

for suitable a, b, c, \dots, z that could be used instead of Equation (36). One should choose coefficients a, b, c, \dots, z so that the univariate polynomial $g(x) = a \cdot x + b \cdot x^3 + c \cdot x^5 + \dots + z \cdot x^{2n+1}$ is a suitable approximation to $\text{sign}(x)$.

Which of these methods is most useful in practice may depend on factors such as the condition number of the matrix \mathbf{G} or the nature of the available computational resources.

Appendix B. Proofs

Proposition 1: Steepest descent

First, let's study the minimization under the change of variables $\Delta \mathbf{w} = c \cdot \mathbf{t}$, where $c \geq 0$ encodes the “magnitude” and \mathbf{t} is a unit vector ($\|\mathbf{t}\| = 1$) encoding the “direction”:

$$\min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right] = \min_{c \geq 0} \min_{\mathbf{t} \in \mathbb{R}^n: \|\mathbf{t}\|=1} \left[c \cdot \mathbf{g}^\top \mathbf{t} + \frac{\lambda}{2} c^2 \|\mathbf{t}\|^2 \right] \quad (38)$$

$$= \min_{c \geq 0} \left[c \cdot \min_{\mathbf{t} \in \mathbb{R}^n: \|\mathbf{t}\|=1} [\mathbf{g}^\top \mathbf{t}] + \frac{\lambda}{2} c^2 \right] \quad (39)$$

$$= \min_{c \geq 0} \left[-c \cdot \|\mathbf{g}\|^\dagger + \frac{\lambda}{2} c^2 \right], \quad (40)$$

Inspecting Equation (39), we see that the minimizer for the direction \mathbf{t} is given by:

$$\mathbf{t} = \arg \min_{\mathbf{t} \in \mathbb{R}^n: \|\mathbf{t}\|=1} [\mathbf{g}^\top \mathbf{t}] = - \arg \max_{\mathbf{t} \in \mathbb{R}^n: \|\mathbf{t}\|=1} [\mathbf{g}^\top \mathbf{t}] \quad (41)$$

And similarly, by inspecting Equation (40), the minimizer for the magnitude c is given by:

$$c = \arg \min_{c \geq 0} \left[-c \cdot \|\mathbf{g}\|^\dagger + \frac{\lambda}{2} c^2 \right] = \frac{\|\mathbf{g}\|^\dagger}{\lambda}. \quad (42)$$

Multiplying these expressions, we obtain the minimizer for $\Delta \mathbf{w}$, yielding the result. \blacksquare

Proposition 2: Sign descent as steepest descent under the infinity norm

The result follows by applying Proposition 1. We just need that $\arg \max_{\|\mathbf{t}\|_\infty=1} \mathbf{g}^\top \mathbf{t} = \text{sign}(\mathbf{g})$, and also that the dual norm $\|\mathbf{g}\|_\infty^\dagger := \max_{\|\mathbf{t}\|_\infty=1} \mathbf{g}^\top \mathbf{t} = \mathbf{g}^\top \text{sign}(\mathbf{g}) = \|\mathbf{g}\|_1$. \blacksquare

Proposition 3: Sign descent as steepest descent under the max-of-max norm

The result follows from Proposition 7 by setting all the scalars s_1, \dots, s_L to one and all the norms $\|\cdot\|_1, \dots, \|\cdot\|_L$ to the ℓ_1 to ℓ_∞ operator norm. All we need is to show that the argmax at each matrix space $l = 1, \dots, L$ satisfies:

$$\arg \max_{\|\mathbf{T}_l\|_{\ell_1 \rightarrow \ell_\infty} = 1} \text{tr}(\mathbf{G}_l^\top \mathbf{T}_l) = \text{sign}(\mathbf{G}_l). \quad (43)$$

But this holds because, by Proposition 8, $\|\mathbf{T}\|_{\ell_1 \rightarrow \ell_\infty} = \max_i \|\text{col}_i(\mathbf{T})\|_\infty = \max_{ij} |\mathbf{T}_{ij}|$, and therefore all components in the argmax must be of unit size and gradient aligned. \blacksquare

Proposition 4: Projection to the closest semi-orthogonal matrix

To begin, we observe that the minimizer over semi-orthogonal matrices of the “distance” $\|\mathbf{A} - \mathbf{G}\|_F$ is the same as the maximizer over semi-orthogonal matrices of the “alignment” $\langle \mathbf{A}, \mathbf{G} \rangle$, where $\langle \cdot, \cdot \rangle$ denotes the Frobenius inner product. This is because:

$$\|\mathbf{A} - \mathbf{G}\|_F^2 = \|\mathbf{A}\|_F^2 - 2 \cdot \langle \mathbf{A}, \mathbf{G} \rangle + \|\mathbf{G}\|_F^2, \quad (44)$$

and the term $\|\mathbf{A}\|_F^2$ is fixed at $\|\mathbf{A}\|_F^2 = \min(m, n)$ for a semi-orthogonal matrix $\mathbf{A} \in \mathcal{O}_{m \times n}$. Now, let $\mathbf{G} = \sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ denote the SVD of \mathbf{G} . Then the alignment satisfies:

$$\langle \mathbf{A}, \mathbf{G} \rangle = \text{tr} \sum_i \sigma_i \mathbf{v}_i \mathbf{u}_i^\top \mathbf{A} = \sum_i \sigma_i \mathbf{u}_i^\top \mathbf{A} \mathbf{v}_i \leq \sum_i \sigma_i, \quad (45)$$

where the second equality follows by the cyclic property of the trace, and the inequality is since \mathbf{A} being semi-orthogonal means that $\mathbf{u}^\top \mathbf{A} \mathbf{v} \leq 1$ for any two unit vectors \mathbf{u} and \mathbf{v} .

Next, observe that for the semi-orthogonal matrix $\mathbf{A}_* = \sum_i \mathbf{u}_i \mathbf{v}_i^\top$, we have that:

$$\langle \mathbf{A}_*, \mathbf{G} \rangle = \sum_i \sigma_i \sum_j \mathbf{u}_i^\top \mathbf{u}_j \mathbf{v}_j^\top \mathbf{v}_i = \sum_i \sigma_i, \quad (46)$$

since the $\{\mathbf{u}_i\}$ and $\{\mathbf{v}_i\}$ are orthonormal. Comparing against Equation (45), we see that \mathbf{A}_* indeed maximizes the alignment, since it achieves the upper bound of $\sum_i \sigma_i$. And \mathbf{A}_* therefore also minimizes the distance $\|\mathbf{A} - \mathbf{G}\|_F$ amongst semi-orthogonal matrices \mathbf{A} . Note that if \mathbf{U} is the matrix that has the $\{\mathbf{u}_i\}$ as columns, and likewise for \mathbf{V} and the $\{\mathbf{v}_i\}$, then this solution may equivalently be expressed as $\mathbf{A}_* = \mathbf{U} \mathbf{V}^\top$.

All that remains is to explore the uniqueness of this solution:

- If \mathbf{G} is full rank, the solution \mathbf{A}_* is unique. \mathbf{G} being full rank means that all the singular values σ_i are positive. In this case, we see from Equation (45) that to maximize the alignment the semi-orthogonal matrix \mathbf{A} must satisfy $\mathbf{u}_i^\top \mathbf{A} \mathbf{v}_i = 1$ for all i . Since \mathbf{A} has spectral norm one, in turn this requires that $\mathbf{A} \mathbf{v}_i = \mathbf{u}_i$ and $\mathbf{A}^\top \mathbf{u}_i = \mathbf{v}_i$ for all i . These conditions uniquely pick out $\mathbf{A} = \sum_i \mathbf{u}_i \mathbf{v}_i^\top$.
- If \mathbf{G} is not full rank then the solution \mathbf{A}_* is not unique. This solution is just as good:

$$\mathbf{A}_\dagger = \sum_{i: \sigma_i > 0} \mathbf{u}_i \mathbf{v}_i^\top + \sum_{i: \sigma_i = 0} \mathbf{u}_i (-\mathbf{v}_i)^\top. \quad (47)$$

This completes the proof. ■

Proposition 5: Shampoo as steepest descent under the spectral norm

First, we apply Proposition 7 with scalars s_1, \dots, s_L set to one and all norms set to $\|\cdot\|_{\ell_2 \rightarrow \ell_2}$. This tells us that the solution is given by $\Delta \mathbf{W}_l = -\eta \cdot \arg \max_{\|\mathbf{T}_l\|=1} \text{tr}(\mathbf{G}_l^\top \mathbf{T}_l)$ for each $l = 1, \dots, L$ and with $\eta = \frac{1}{\lambda} \sum_{k=1}^L \|\mathbf{G}_k\|_{\ell_2 \rightarrow \ell_2}^\dagger$. We just need to resolve the dual norm and evaluate the argmax.

Let's start with the dual norm. For a matrix \mathbf{G} with SVD $\sum_i \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ we have:

$$\|G\|_{\ell_2 \rightarrow \ell_2}^\dagger := \max_{\|T\|_{\ell_2 \rightarrow \ell_2}=1} \text{tr } G^\top T = \text{tr} \sum_i \sigma_i v_i u_i^\top T = \sum_i \sigma_i u_i^\top T v_i \leq \sum_i \sigma_i = \text{tr } \Sigma, \quad (48)$$

where the upper bound follows from the spectral norm constraint on T . But this upper bound is attained by setting $T = UV^\top$ (also resolving the argmax) and so $\|G\|_{\ell_2 \rightarrow \ell_2}^\dagger = \text{tr } \Sigma$. The uniqueness claim follows by the same argument as for Proposition 4. ■

Proposition 6: Bounding the square loss of a linear predictor

First observe that the square loss is quadratic in W so there are no cubic terms or higher. The bound must agree to first-order with the first-order Taylor expansion of $\mathcal{L}(W + \Delta W)$, which is precisely $\mathcal{L}(W) + \langle \nabla_W \mathcal{L}(W), \Delta W \rangle$, since otherwise the bound would be violated for sufficiently small ΔW . To obtain the second-order piece of the bound, it's easiest just to multiply out $\mathcal{L}(W + \Delta W)$ and see that the second-order piece of $\mathcal{L}(W + \Delta W)$ satisfies:

$$\frac{1}{2n} \sum_{i=1}^n \frac{1}{d_{\text{out}}} \|\Delta W x^{(i)}\|_2^2 \leq \frac{1}{2n} \sum_{i=1}^n \frac{1}{d_{\text{out}}} \|\Delta W\|_{\ell_2 \rightarrow \ell_2}^2 \cdot \|x^{(i)}\|_2^2 = \frac{1}{2} \frac{d_{\text{in}}}{d_{\text{out}}} \|\Delta W\|_{\ell_2 \rightarrow \ell_2}^2, \quad (49)$$

where the last equality uses the input normalization $\|x^{(i)}\|_2 = \sqrt{d_{\text{in}}}$. We are done. ■

Proposition 7: Steepest descent under the modular norm

For each layer $l = 1, \dots, L$, we decompose ΔW_l into its magnitude and direction: $\Delta W_l = c_l \cdot T_l$ for $c_l \geq 0$ and $\|T_l\|_2 = 1$. Under this change of variables, the minimization becomes:

$$\min_{\Delta W_1, \dots, \Delta W_L} \left[\sum_{l=1}^L \langle G_l, \Delta W_l \rangle + \frac{\lambda}{2} \sum_{l=1}^L s_l^2 \|\Delta W_l\|_2^2 \right] \quad (50)$$

$$= \min_{c_1, \dots, c_L \geq 0} \left[\sum_{l=1}^L c_l \min_{\|T_l\|_2=1} \langle G_l, T_l \rangle + \frac{\lambda}{2} \sum_{l=1}^L s_l^2 c_l^2 \right] \quad (51)$$

$$= \min_{c_1, \dots, c_L \geq 0} \left[- \sum_{l=1}^L c_l \|G_l\|_l^\dagger + \frac{\lambda}{2} \sum_{l=1}^L s_l^2 c_l^2 \right] \quad (52)$$

$$= \min_{\eta \geq 0} \left[- \sum_{l=1}^L \frac{\eta}{s_l} \|G_l\|_l^\dagger + \frac{\lambda}{2} \eta^2 \right], \quad (53)$$

where Equation (53) follows by observing that at the minimum we must have $s_1 c_1, \dots, s_L c_L$ all taking the same value of $\eta \geq 0$ (still to be determined), since otherwise we could increase the sum $\sum_l c_l \|G_l\|_l^\dagger$ by increasing any of the slack c_l without paying a penalty in terms of

the max. We can now read off the minimizers from [Equations \(51\)](#) to [\(53\)](#):

$$\mathbf{T}_l = \arg \min_{\|\mathbf{T}_l\|=1} \langle \mathbf{G}_l, \mathbf{T}_l \rangle = - \arg \max_{\|\mathbf{T}_l\|=1} \langle \mathbf{G}_l, \mathbf{T}_l \rangle; \quad (54)$$

$$c_l = \frac{\eta}{s_l}; \quad (55)$$

$$\eta = \frac{1}{\lambda} \sum_{k=1}^L \frac{1}{s_k} \|\mathbf{G}_k\|_k^\dagger. \quad (56)$$

Combining, we obtain the overall minimizer for each $l = 1, \dots, L$ via $\Delta \mathbf{W}_l = c_l \cdot \mathbf{T}_l = -\frac{\eta}{s_l} \arg \max \langle \mathbf{G}_l, \mathbf{T}_l \rangle$, where η is given by [Equation \(56\)](#), proving the result. ■

Proposition 8: $\ell_1 \rightarrow \ell_p$ and $\ell_p \rightarrow \ell_\infty$ induced operator norms are tractable

Let's start with the $\ell_1 \rightarrow \ell_p$ operator norm. Here we observe that, in matrix-vector multiplication, each component of an input vector selects and scales a column of the matrix:

$$\|\mathbf{M}\|_{\ell_1 \rightarrow \ell_p} = \max_{\|\mathbf{x}\|_1=1} \|\mathbf{M}\mathbf{x}\|_p = \max_{\|\mathbf{x}\|_1=1} \left\| \sum_j \text{col}_j(\mathbf{M}) \mathbf{x}_j \right\|_p \leq \max_{\|\mathbf{x}\|_1=1} \sum_j |\mathbf{x}_j| \cdot \|\text{col}_j(\mathbf{M})\|_p \quad (57)$$

$$\leq \max_{\|\mathbf{x}\|_1=1} \|\mathbf{x}\|_1 \cdot \max_j \|\text{col}_j(\mathbf{M})\|_p \quad (58)$$

$$= \max_j \|\text{col}_j(\mathbf{M})\|_p, \quad (59)$$

by the triangle inequality and Hölder's inequality. But the upper bound in [Equation \(59\)](#) is attained by selecting the column index $j_\star = \arg \max_j \|\text{col}_j(\mathbf{M})\|_p$ with the largest norm, then setting $\mathbf{x}_{j_\star} = 1$ and the other input components to zero. So $\|\mathbf{M}\|_{\ell_1 \rightarrow \ell_p} = \max_j \|\text{col}_j(\mathbf{M})\|_p$.

Next, let's deal with the $\ell_p \rightarrow \ell_\infty$ operator norm. Here we break up a matrix-vector product in terms of the dot product between the vector and the matrix rows:

$$\|\mathbf{M}\|_{\ell_p \rightarrow \ell_\infty} = \max_{\|\mathbf{x}\|_p=1} \|\mathbf{M}\mathbf{x}\|_\infty = \max_{\|\mathbf{x}\|_p=1} \max_i |\mathbf{x}^\top \text{row}_i(\mathbf{M})| \quad (60)$$

$$= \max_i \max_{\|\mathbf{x}\|_p=1} |\mathbf{x}^\top \text{row}_i(\mathbf{M})| \quad (61)$$

$$= \max_i \|\text{row}_i(\mathbf{M})\|_p^\dagger. \quad (62)$$

The proof is completed by recalling that the vector ℓ_p norm is dual to the vector ℓ_q norm for $1/p + 1/q = 1$. In other words, $\|\cdot\|_p^\dagger = \|\cdot\|_{\frac{p}{p-1}}$. ■