

THE NUCLEAR NORM AND BEYOND: DUAL NORMS AND COMBINATIONS FOR MATRIX OPTIMIZATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In this article, we explore the use of matrix norms for optimizing functions of weight matrices, a crucial problem in training large language models. Moving beyond the spectral norm that underlies the Muon update, we leverage the nuclear norm, its affine combinations with other norms, and their corresponding duals to develop a new family of Muon-like algorithms. We complement our theoretical analysis with an extensive empirical study of the algorithms across a wide range of tasks and settings.

1 INTRODUCTION

Minimizing loss functions in unprecedentedly high-dimensional spaces has recently become an integral and crucial part in training large language models. Hence, new scalable, time- and memory-efficient algorithms have been demanded. Besides well-known Adam and AdamW Kingma & Ba (2014), Loshchilov & Hutter (2017), recently proposed Muon has shown promising results on different benchmarks (CIFAR, NanoGPT, citations). Its key difference from earlier algorithms is that it has been constructed specifically for optimizing functions of weight matrices. Such functions are common in modern machine learning (citations), so it does not significantly restrict Muon’s applicability to modern problems.

That is what can be said from a practical point of view. From the perspective of theory, Muon’s main innovation was an intentional usage of matrix norms, i.e. the spectral norm, to derive the algorithm’s update (cite Deriving Muon). Several other attempts have been since made to construct new algorithms, mainly generalizations of Muon’s paradigm (Pethick et al. (2025) and Riabinin et al. (2025)).

Based upon recent theoretical advances that explain some theory behind Muon, Scion and Gluon (Bernstein (2025); Kovalev (2025); Pethick et al. (2025); Riabinin et al. (2025)), we explore application of other matrix norms to optimization of functions of matrices. As it has been done with Muon, we stipulate that our algorithms’ updates be fast to compute.

2 PROBLEM STATEMENT

2.1 FUNCTION OF A MATRIX AND ASSUMPTIONS

Warning: the text was copied, while formulas were adapted. We need to rephrase the subsection! We consider the problem of minimizing a differentiable matrix function $F(\cdot): \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

$$\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} f(\mathbf{X}). \quad (1)$$

To make the theoretical analysis possible in the future, we make three reasonable non-restrictive assumptions. Idea: move gradient assumptions and smoothness to the L-smooth part. Here we need only custom norm and their relation to the Frobenius norm.

Stochastic gradient estimator. We assume access to a stochastic estimator $g(\cdot; \xi): \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^{m \times n}$ of the gradient $\nabla f(\cdot)$, where $\xi \sim \mathcal{D}$ is a random variable sampled from a probability distribution \mathcal{D} . We assume that the stochastic gradient estimator $g(\cdot; \xi)$ is unbiased and has bounded

variance, that is, the following relations hold:

$$\mathbb{E}_{\xi \sim \mathcal{D}} [g(\mathbf{X}; \xi)] = \nabla f(\mathbf{X}) \quad \text{and} \quad \mathbb{E}_{\xi \sim \mathcal{D}} [\|g(\mathbf{X}; \xi) - \nabla f(\mathbf{X})\|_{\text{F}}^2] \leq \sigma^2 \quad \text{for all } \mathbf{X} \in \mathbb{R}^{m \times n}, \quad (\text{A1})$$

where $\sigma > 0$ is a positive variance parameter, and $\|\cdot\|_{\text{F}}$ is the standard Euclidean, i.e. Frobenius, norm induced by the inner product $\langle \cdot, \cdot \rangle$, i.e., $\|\mathbf{X}\|_{\text{F}} = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle} = \sqrt{\text{tr}(\mathbf{X}^\top \mathbf{X})}$. These assumptions have been widely adopted for the analysis of many stochastic gradient optimization algorithms (Ghadimi & Lan, 2013; Ghadimi et al., 2016; Cutkosky et al., 2020; Sun et al., 2023; Horváth et al., 2023; Gorbunov et al., 2020).

Non-Euclidean norm setting and Lipschitz continuous gradient. We assume that matrix space $\mathbb{R}^{m \times n}$ is equipped with a norm $\|\cdot\|: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_+$, which possibly does not coincide with the Frobenius, norm $\|\cdot\|_{\text{F}}$. In addition, we assume that the gradient $\nabla f(\cdot)$ is Lipschitz continuous with respect to the norm $\|\cdot\|$, that is, the following inequality holds:

$$\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{X}')\|^\dagger \leq L \|\mathbf{X} - \mathbf{X}'\| \quad \text{for all } \mathbf{X}, \mathbf{X}' \in \mathbb{R}^{m \times n}, \quad (\text{A2})$$

where $L > 0$ is the gradient Lipschitz constant, and $\|\cdot\|^\dagger: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}_+$ is the dual norm associated with $\|\cdot\|$, i.e., $\|\mathbf{X}\|^\dagger = \sup_{\|\mathbf{X}'\| \leq 1} \langle \mathbf{X}, \mathbf{X}' \rangle$ for all $\mathbf{X} \in \mathbb{R}^{m \times n}$. The assumption of gradient Lipschitz continuity is also widespread in the analysis of first-order optimization methods (Ghadimi & Lan, 2013; Gower et al., 2019; Cutkosky et al., 2020; Horváth et al., 2023; Gorbunov et al., 2020). It is important to highlight that while Assumption (A2) uses the dual norm $\|\cdot\|^\dagger$ to measure the difference between the gradients, the variance in Assumption (A1) is measured with respect to the Frobenius norm $\|\cdot\|_{\text{F}}^2$, which is necessary to properly utilize the unbiasedness property of the stochastic gradient estimator $g(\cdot; \xi)$. Therefore, we need to provide a connection between these norms using the following inequality:

$$\|\mathbf{X}\|^\dagger \leq \rho \cdot \|\mathbf{X}\|_{\text{F}} \quad \text{for all } \mathbf{X} \in \mathbb{R}^{m \times n}, \quad (\text{A3})$$

where $\rho > 0$ is a positive constant. Note that such a constant always exists due to the norm equivalence theorem, which always holds in the finite-dimensional space $\mathbb{R}^{m \times n}$. We recount ρ for different norms $\|\cdot\|$ in the appendix.

2.2 LINEAR MINIMIZATION ORACLE AND TRUST REGION

Let us look at the problem from the perspective of linear minimization oracle (lmo) and unconstrained stochastic conditional gradient descent (uSCG) (Pethick et al. (2025)). lmo is defined as:

$$\text{lmo}(\mathbf{S}) \in \arg \min_{\mathbf{X} \in \mathcal{S}} \langle \mathbf{S}, \mathbf{X} \rangle, \quad (2)$$

where \mathcal{S} is some set. We are interested in the case when \mathcal{S} is a ball in our $\|\cdot\|$ norm:

$$\mathcal{S} := \mathcal{B}_\eta := \{\mathbf{X} \mid \|\mathbf{X}\| \leq \eta\}. \quad (3)$$

uSCG update is defined as: $\mathbf{X}^{k+1} = \mathbf{X}^k + \gamma_k \text{lmo}(\mathbf{M}^k)$, where $\mathbf{M}^{k+1} = (1 - \alpha_{k+1})\mathbf{M}^k + \alpha_{k+1}g(\mathbf{X}^k, \xi_k)$ is a momentum.

It can be easily shown that the formula is equivalent to

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \gamma_k \eta \arg \max_{\mathbf{X} \in \mathcal{B}_1} \langle \mathbf{S}, \mathbf{X} \rangle = \mathbf{X}^k - \gamma_k \eta \{\Delta \in \mathcal{B}_1 \mid \langle \mathbf{M}^k, \Delta \rangle = \|\mathbf{M}^k\|^\dagger\}. \quad (4)$$

Let us set $\gamma_k = 1$. This transforms algorithm defined by eq. (4) into Algorithm 1 from Kovalev (2025). Therefore, we can view eq. (4) both as an lmo-based algorithm and as a trust-region algorithm.

3 DIFFERENT NORMS $\|\cdot\|$ IMPLY DIFFERENT UPDATES

Based on different norms $\|\cdot\|$, we simplify the update defined by the aforementioned equation:

$$\mathbf{X}^{k+1} = \mathbf{X}^k - \eta \{\Delta \in \mathcal{B}_1 \mid \langle \mathbf{M}^k, \Delta \rangle = \|\mathbf{M}^k\|^\dagger\} \quad (5)$$

In all the work, we define U, Σ, V^\top as components of the singular value decomposition of M^k : $M^k = U\Sigma V^\top$. We use common notations: $U = [u_1, u_2, \dots, u_r]$, $\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_r)$, and $V = [v_1, v_2, \dots, v_r]$.

$\|M^k\|_F$ and NSGD

Lemma 1. When $\|\cdot\| = \|\cdot\|_F$, eq. (5) turns into:

$$X^{k+1} = X^k - \eta \frac{M^k}{\|M^k\|_F} \quad (6)$$

It is an interesting observation, because in other works (Pethick et al. (2025)), $\|\cdot\|_F$ was used to recover SGD. The difference is in how one states the problem.

$\|M^k\|_{\text{op}}$ and Muon

Lemma 2. When $\|\cdot\| = \|\cdot\|_{\text{op}}$, eq. (5) turns into:

$$X^{k+1} = X^k - \eta UV^\top \quad (7)$$

$\|M^k\|_{\text{nuc}}$ and Neon

Lemma 3. When $\|\cdot\| = \|\cdot\|_{\text{nuc}}$, eq. (5) turns into:

$$X^{k+1} = X^k - \eta u_1 v_1^\top \quad (8)$$

We name the derived algorithm *Neon*. In the section Matrix side of updates, we will discuss how to compute an update efficiently.

$\|M^k\|_{F*}^\dagger$ and F-Muon We define $\|\cdot\|_{F*}$ as a convex combination of $\|\cdot\|_{\text{nuc}}$ and $\|\cdot\|_F$:

$$\|X\|_{F*} = \alpha \|X\|_{\text{nuc}} + (1 - \alpha) \|X\|_F, \quad (9)$$

where $\alpha \in [0, 1]$ defines a specific norm of F^* -family.

Lemma 4. When $\|\cdot\| = \|\cdot\|_{F*}^\dagger$, eq. (5) turns into:

$$X^{k+1} = X^k - \eta (\alpha UV^\top + (1 - \alpha) \frac{M^k}{\|M^k\|_F}) \quad (10)$$

We name the derived algorithm *F-Muon*. It turns out that F-Muon is a convex combination of Normalized SGD and Muon, which is curious. The implications are significant and discussed in the following sections.

$\|M^k\|_{F2}^\dagger$ and F-Neon We define $\|\cdot\|_{F2}$ as a convex combination of $\|\cdot\|_{\text{op}}$ and $\|\cdot\|_F$:

$$\|X\|_{F2} = \alpha \|X\|_{\text{op}} + (1 - \alpha) \|X\|_F, \quad (11)$$

where $\alpha \in [0, 1]$ defines a specific norm of F2-family.

Lemma 5. When $\|\cdot\| = \|\cdot\|_{F2}^\dagger$, eq. (5) turns into:

$$X^{k+1} = X^k - \eta (\alpha u_1 v_1^\top + (1 - \alpha) \frac{M^k}{\|M^k\|_F}) \quad (12)$$

We name the derived algorithm *F-Neon*. It turns out that F-Neon is a convex combination of Normalized SGD and Neon, which is curious. The implications are significant and discussed in the following sections.

3.1 ALGORITHMS FOR MATRICES \leftrightarrow ALGORITHMS FOR VECTORS

Table that compares Muons and Neons to vector algorithms (NSGD).

4 MATRIX SIDE OF UPDATES

4.1 THEORY

Not only formulas, but also citations! Solutions to efficiently find some parts of SVD:

- Basic SVD (complexity asymptotic!)
- Newton-Schulz (complexity asymptotic!)
- Our job: Power iterations, Randomized SVD and Lanczos

4.2 EXPERIMENTS

Nikolay's task: SVD, RSVD, Lanczos, Newton-Schulz. We use torch and cupy to test it. Options: all possible on torch, or additionally to compare all on cupy. The goal: to find how much do we lose due to cupy

5 TRUST REGION BOUNDS FOR L-SMOOTH FUNCTIONS

First, we analyze the problem in the unstochastic case. From Corollary 1 of Kovalev (2025), we directly get the following result that matches lower bounds, as was noted in Kovalev (2025).

Lemma 6. *To reach the precision $\min_{k=1\dots K} \|\nabla f(\mathbf{X}_k)\|^\dagger \leq \varepsilon$ by the iterations equation 5 under the conditions of Assumption (A2), it is sufficient to choose the stepsize η and the number of iterations K as follows:*

$$\eta = \mathcal{O}\left(\frac{\varepsilon}{L}\right), \quad K = \mathcal{O}\left(\frac{L\Delta_0}{\varepsilon^2}\right). \quad (13)$$

In the stochastic case, from Corollary of 2 of Kovalev (2025), we directly get the following result that matches lower bounds, as was noted in Kovalev (2025):

Lemma 7. *To reach the precision $\mathbb{E} \min_{k=1\dots K} \|\nabla f(\mathbf{X}_k)\|^\dagger \leq \varepsilon$ by equation 5 under the assumptions Assumption (A1), Assumption (A2), Assumption (A3), it is sufficient to choose the parameters as follows:*

$$\eta = \mathcal{O}\left(\min\left\{\frac{\varepsilon}{L}, \frac{\varepsilon^3}{\rho^2\sigma^2L}\right\}\right), \quad \alpha = \mathcal{O}\left(\min\left\{1, \frac{\varepsilon^2}{\rho^2\sigma^2}\right\}\right), \quad (14)$$

$$K = \mathcal{O}\left(\max\left\{\frac{\rho\sigma}{\varepsilon}, \frac{\rho^3\sigma^3}{\varepsilon^3}, \frac{L\Delta_0}{\varepsilon^2}, \frac{L\Delta_0\rho^2\sigma^2}{\varepsilon^4}\right\}\right). \quad (15)$$

As the norms $\|\cdot\|_F$, $\|\cdot\|_{\text{nuc}}$, $\|\cdot\|_F$ are almost proportional to each other when $m, n \rightarrow \infty$ (with high probability for random matrices), the expected convergence guarantees in terms of $\|\cdot\|_F$ are the same (it can be easily shown by noting that $\|\mathbf{X}\| \sim \alpha \|\mathbf{X}\|_F$, $\|\nabla f(\mathbf{X})\|^\dagger \sim \frac{1}{\alpha} \|\nabla f(\mathbf{X})\|_F$, and expressing L -constant via L_F -constant for the Frobenius norm).

From the theory of random martices and the Marchenko-Pastur law, we get that random $\mathbf{M} \in \mathbb{R}^{m \times n}$: $\mathbf{M} \sim \mathcal{N}(0, \sigma^2 \mathbb{R}^{m \times n})$ has the following asymptotics of its norms:

Nuclear: $\sigma n \sqrt{m}$

Frobenius: $\sigma \sqrt{mn}$

Spectral: $\sigma(\sqrt{m} + \sqrt{n})$

This means that for square random matrices $n \times n$ the following asymptotics take place: $\|\cdot\|_F \sim \frac{\sqrt{n}}{2} \|\cdot\|_{\text{op}}$ and $\|\cdot\|_{\text{nuc}} \sim \frac{n}{2} \|\cdot\|_{\text{op}}$.

6 EXPERIMENTS

6.1 RANDOMIZED LINEAR LEAST SQUARED

Since the provided by other authors Kovalev (2025); Riabinin et al. (2025) theoretical guarantees are almost norm-independent, we have to test them in practice.

To test the bounds from Kovalev (2025) in practice, we construct the following L-smooth problem:

$$F(\mathbf{X}) = \frac{1}{2} \langle (\mathbf{X} - \mathbf{S}), \mathbf{M}(\mathbf{X} - \mathbf{S})\mathbf{N} \rangle \quad (16)$$

where $\mathbf{X} \in \mathbb{R}^{m \times n}$, $m = 10$, $n = 10$, $\mathbf{S} \in \mathbb{R}^{m \times n}$, $\mathbf{M} \in \mathbb{S}_+^m$ and $\mathbf{N} \in \mathbb{S}_+^n$ are positive-semidefinite matrices. Spectra of \mathbf{M} and \mathbf{N} are uniformly distributed in $(0, 1)$ interval.

It is easy to derive the gradient

$$\nabla F(\mathbf{X}) = \mathbf{M}(\mathbf{X} - \mathbf{S})\mathbf{N}, \quad (17)$$

Let us define γ as $\|\cdot\| \sim \gamma \|\cdot\|_F$, which is asymptotics from the previous section. Then $\|\cdot\|^\dagger \sim \frac{1}{\gamma} \|\cdot\|_F$, as $\|\cdot\|_F^\dagger = \|\cdot\|_F$. Hence, $\|\cdot\| \sim \gamma^2 \|\cdot\|^\dagger$

Then $\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{Y})\|^\dagger = \|\mathbf{M}(\mathbf{X} - \mathbf{Y})\mathbf{N}\|^\dagger \leq \|\mathbf{M}\|^\dagger \|\mathbf{N}\|^\dagger \|\mathbf{X} - \mathbf{Y}\|^\dagger \sim \|\mathbf{M}\|^\dagger \|\mathbf{N}\|^\dagger \gamma^2 \|\mathbf{X} - \mathbf{Y}\|$, and $L \sim \gamma^2 \|\mathbf{M}\|^\dagger \|\mathbf{N}\|^\dagger$.

Now for the known norms:

Frobenius (NSGD): $L = \|\mathbf{M}\|_F \|\mathbf{N}\|_F$

Spectral (Muon): $\gamma \sim \frac{2}{\sqrt{n}} \implies L \sim \frac{2}{n} \|\mathbf{M}\|_{\text{nuc}} \|\mathbf{N}\|_{\text{nuc}}$

Nuclear (Neon): $\gamma \sim \sqrt{n} \implies L \sim n \|\mathbf{M}\|_{\text{op}} \|\mathbf{N}\|_{\text{op}}$

We take learning rate η and iteration number K from eq. (13).

We are interested to get $\varepsilon_2 = 1\text{e-}1$ precision in $\|\cdot\|_{\text{op}}$. Hence, we target $\varepsilon = \varepsilon_2$ for the nuclear norm (Neon), $\varepsilon = \frac{n}{2}$ for the spectral norm (Muon), and $\varepsilon = \frac{\sqrt{n}}{2}$ for the Frobenius norm (NSGD).

An independent experiment To make the problem more real-world, we do not use known information about smoothness. We run NormalizedSGD, Muon, F-Muon with $\alpha = 1/2$, Neon, and F-Neon with $\alpha = 1/2$ for 100 000 iterations with learning rate = 1e-3.

6.2 LOGISTIC REGRESSION

6.3 BENCHMARKS

CNN benchmark NanoGPT benchmark

7 CONCLUSION

- Future work: (L0, L1)-smoothness, probably extrapolation (to make dependency on eps milder, but we don't to spend memory here so never mind), hyperparameter-free (Bernstein2023)

8 APPENDIX

8.1 NORM DERIVATIONS

8.2 UPDATES DERIVATIONS

Derivation of eq. (7) follows from eq. (10) with $\alpha = 1$. Indeed, $\|\cdot\|_{\text{op}}^\dagger = 1 \cdot \|\cdot\|_{\text{nuc}} + 0 \cdot \|\cdot\|_F$.

Derivation of eq. (10): Since $\|\cdot\|^\dagger = \|\cdot\|_{F*}^{\dagger\dagger} = \|\cdot\|_{F*}$, the goal is to reach $\|M^k\|_{F*} = \alpha \text{tr } \Sigma + (1 - \alpha)\|M^k\|_F$.

Let us note that $\Delta = \alpha(UV^\top) + (1 - \alpha)\frac{M^k}{\|M^k\|_F}$ delivers this value. Indeed, by the trace property, $\langle M^k, \Delta \rangle = \langle U\Sigma V^\top, \alpha UV^\top + (1 - \alpha)\frac{U\Sigma V^\top}{\|M^k\|_F} \rangle = \alpha \text{tr } \Sigma + (1 - \alpha)\|M^k\|_F = \|M^k\|_{F*}$, which completes the proof.

Derivation of eq. (8) follows from eq. (12) with $\alpha = 1$. Indeed, $\|\cdot\|_{\text{nuc}}^\dagger = 1 \cdot \|\cdot\|_{\text{op}} + 0 \cdot \|\cdot\|_F$.

Derivation of eq. (12): Since $\|\cdot\|^\dagger = \|\cdot\|_{F2}^{\dagger\dagger} = \|\cdot\|_{F2}$, the goal is to reach $\|M^k\|_{F2} = \alpha\sigma_1 + (1 - \alpha)\|M^k\|_F$.

Let us note that $\Delta = \alpha(u_1 v_1^\top) + (1 - \alpha)\frac{M^k}{\|M^k\|_F}$ delivers this value. Indeed, by the trace property and singular vectors orthogonality, $\langle M^k, \Delta \rangle = \langle U\Sigma V^\top, \alpha u_1 v_1^\top + (1 - \alpha)\frac{U\Sigma V^\top}{\|M^k\|_F} \rangle = \alpha \text{tr } \text{diag}(\sigma_1, 0, \dots, 0) + (1 - \alpha)\|M^k\|_F = \|M^k\|_{F2}$, which completes the proof.

8.3

Note: nuSAM and their update. They had no Lanczos Technical details on experiments

REFERENCES

- Jeremy Bernstein. Deriving muon, 2025. URL <https://jeremybernste.in/writing/deriving-muon>.
- Ashok Cutkosky, Harsh Mehta, and Francesco Orabona. Momentum-based variance reduction in nonconvex sgd. *Advances in Neural Information Processing Systems*, 2020.
- Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 2013.
- Saeed Ghadimi, Guanghui Lan, and Hongchao Zhang. Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 2016.
- Eduard Gorbunov, Dmitry Kovalev, and Peter Richtárik. Linearly converging error compensated sgd. *Advances in Neural Information Processing Systems*, 2020.
- Robert M Gower et al. Sgd: General analysis and improved rates. *International Conference on Machine Learning*, 2019.
- Samuel Horváth, Dmitry Kovalev, and Peter Richtárik. Stochastic recursive momentum for nonconvex optimization. *arXiv preprint arXiv:2302.07731*, 2023.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Dmitry Kovalev. Understanding gradient orthogonalization for deep learning via non-euclidean trust-region optimization. *arXiv preprint arXiv:2503.12645*, 2025.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Thomas Pethick, Wanyun Xie, Kimon Antonakopoulos, Zhenyu Zhu, Antonio Silveti-Falls, and Volkan Cevher. Training deep learning models with norm-constrained lmos. *arXiv preprint arXiv:2502.07529*, 2025.
- Artem Riabinin, Egor Shulgin, Kaja Grutkowska, and Peter Richtárik. Gluon: Making muon & scion great again!(bridging theory and practice of lmo-based optimizers for llms). *arXiv preprint arXiv:2505.13416*, 2025.
- Ruoyu Sun et al. Momentum methods for stochastic optimization: A survey and new results. *arXiv preprint arXiv:2302.06675*, 2023.