

The Ky Fan Norms and Beyond: Dual Norms and Combinations for Matrix Optimization

Alexey Kravatskiy¹ Ivan Kozyrev¹ Nikolai Kozlov¹ Alexander Vinogradov¹ Daniil
Merkulov^{1,2,3,4} Ivan Oseledets^{5,2}

¹MIPT ²Skoltech ³HSE ⁴AI4Science ⁵AIRI

October 10, 2025

Moving beyond the spectral norm and Muon

Objective: $\min_{\mathbf{X} \in \mathbb{R}^{m \times n}} f(\mathbf{X})$

❓ What will happen if we change the spectral norm in the derivation of Muon?

F-Fanions: Muon, Neon, NSGD, Dion without EF, and so much more

From $\|\cdot\|_{\text{op}}$ and Muon:

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta UV^\top$$

To $\|\cdot\|_{\text{KF-k}}^\dagger$ and a general F-Fanion:

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \left(\alpha \sum_{i=1}^k u_i v_i^\top + (1 - \alpha) \frac{\mathbf{M}^t}{\|\mathbf{M}^t\|_F} \right), \alpha \in [0, 1]$$

Linear Minimization Oracle (LMO) and Trust Region

Let us equip $\mathbb{R}^{m \times n}$ with a norm $\|\cdot\|$. Its dual is $\|\mathbf{X}\|^\dagger = \sup_{\|\mathbf{X}'\| \leq 1} \langle \mathbf{X}, \mathbf{X}' \rangle$. $\langle \cdot, \cdot \rangle$ is a Frobenius product.

Both LMO and Trust Region lead to the update

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \arg \max_{\mathbf{X} \in \mathcal{B}_1} \langle \mathbf{M}^t, \mathbf{X} \rangle = \mathbf{X}^t - \eta \{\Delta \in \mathcal{B}_1 \mid \langle \mathbf{M}^t, \Delta \rangle = \|\mathbf{M}^t\|^\dagger\}$$

Recipe: It means we seek Δ from the 1-norm ball that delivers $\langle \mathbf{M}^t, \Delta \rangle = \|\mathbf{M}^t\|^\dagger$. We will often return to the SVD $\mathbf{M}^t = \mathbf{U}\Sigma\mathbf{V}^\top$.

Frobenius $\|\mathbf{M}^k\|_F$ and Normalized SGD

Deriving NSGD by Recipe

$\|\mathbf{M}^t\|_F^\dagger = \|\mathbf{M}^t\|_F$, and $\Delta = \frac{\mathbf{M}^t}{\|\mathbf{M}^t\|_F}$ with $\|\Delta\|_F = 1$ delivers it. Hence,

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \frac{\mathbf{M}^t}{\|\mathbf{M}^t\|_F}$$

Spectral $\|\mathbf{M}^k\|_{\text{op}}$ and Muon. Nuclear $\|\mathbf{M}^k\|_{\text{nuc}}$ and Neon.

Deriving Muon by Recipe

$\|\mathbf{M}^t\|_{\text{op}}^\dagger = \|\mathbf{M}^t\|_{\text{nuc}}$, and $\Delta = \mathbf{UV}^\top$ delivers it: $\|\Delta\|_{\text{op}} = 1$ and

$$\langle \mathbf{U}\Sigma\mathbf{V}^\top, \mathbf{UV}^\top \rangle = \text{tr}(\mathbf{V}\Sigma\mathbf{U}^\top\mathbf{UV}^\top) = \text{tr}\Sigma = \|\mathbf{M}^t\|_{\text{nuc}}$$

Hence,

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta\mathbf{UV}^\top$$

Deriving Neon by Recipe

$\|\mathbf{M}^t\|_{\text{nuc}}^\dagger = \|\mathbf{M}^t\|_{\text{op}}$, and $\Delta = u_1v_1^\top$ delivers it: $\|\Delta\|_{\text{nuc}} = 1$ and

$$\langle \mathbf{U}\Sigma\mathbf{V}^\top, u_1v_1^\top \rangle = \text{tr}(\mathbf{V}\Sigma\mathbf{U}^\top u_1v_1^\top) = \text{tr}\text{diag}(\sigma_1, 0, \dots, 0) = \sigma_1 = \|\mathbf{M}^t\|_{\text{op}}$$

Hence,

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta u_1v_1^\top$$

Of Matrix and Vector Algorithms

Table: lmo optimizers in Schatten S_p norms and in l_p norms.

Method	lmo constraint set \mathcal{D}	lmo	Reference
Normalized SGD	l_2 -ball, S_2 -ball	$-\eta \frac{g}{\ g\ _2} = -\eta \frac{g}{\ g\ _F}$	(Hazan et al., 2015)
Momentum Normalized SGD	Ball in l_2 , or Ball in S_2	$-\eta \frac{g}{\ g\ _2} = -\eta \frac{g}{\ g\ _F}$	(Cutkosky et al., 2020)
SignSGD	Ball in Max-norm l_∞	$-\eta \text{sign}(g)$	(Bernstein et al., 2018, Thm. 1)
Signum	Ball in Max-norm l_∞	$-\eta \text{sign}(g)$	(Bernstein et al., 2018, Thm. 3)
Muon	Ball in Spectral S_∞	$-\eta UV^\top$	(Jordan et al., 2024b)
Gauss-Southwell Coordinate Descent	Ball in l_1	$-\eta \{i : g_i \geq g_k \forall k\}$	(Shi et al., 2016, p.19)
Neon	Ball in Nuclear S_1	$-\eta u_1 v_1^\top$	This work

Understanding Dion by Thomas Pethick

Without momentum, Dion is simplified to

$$\begin{aligned}\Delta &\leftarrow g + e \\ e &\leftarrow \Delta - \sum_{i=1}^r \sigma_i u_i v_i^\top \\ x &\leftarrow x - \gamma \sum_{i=1}^r u_i v_i^\top\end{aligned}$$

where $\gamma > 0$ and $\sum_{i=1}^r \sigma_i u_i v_i^\top$ is the rank- r truncated SVD of $\Delta \in \mathbb{R}^{m \times n}$.

❓ What will we get if we try $\|M\|_{\text{KF-k}} := \sum_{i=1}^r \sigma_i$?

Ky Fan k -rank $\|\mathbf{M}^k\|_{\text{KF}-k}^\dagger$ and Fanions

Deriving Fanions by Recipe

$\|\mathbf{M}^t\|_{\text{KF}-k}^{\dagger\dagger} = \|\mathbf{M}^t\|_{\text{KF}-k}$, and $\Delta = \sum_{i=1}^k u_i v_i^\top$ with
 $\|\Delta\|_{\text{KF}-k}^\dagger = \max\{\frac{1}{k}\|\Delta\|_{\text{nuc}}, \|\Delta\|_{\text{op}}\} = \max\{\frac{1}{k}k, 1\} = 1$ delivers it:

$$\langle \mathbf{M}^t, \Delta \rangle = \langle \mathbf{U} \Sigma \mathbf{V}^\top, \sum_{i=1}^k u_i v_i^\top \rangle = \sum_{i,j=1}^{r,k} \langle u_i \sigma_i v_i^\top, u_j v_j^\top \rangle = \sum_{i=1}^k \sigma_i = \|\mathbf{M}^t\|_{\text{KF}-k}$$

Hence,

$$\mathbf{X}^{t+1} = \mathbf{X}^t - \eta \sum_{i=1}^k u_i v_i^\top$$

Computing the k -rank update: Lanczos method

Method	rtol	k	time (s)
Power Iterations	0.01	1	7.7
SVDS (TRLan)	0.01	1	0.18
PCA Low Rank (RSVD)	0.01	1	1.15
SVDS (TRLan)	0.01	10	0.47
PCA Low Rank (RSVD)	0.01	10	19.4
SVDS (TRLan)	0.01	100	1.96
PCA Low Rank (RSVD)	0.01	100	170

Comparison on a 5000×5000 matrix; rtol is the relative Frobenius error of $\sum_{i=1}^k u_i \sigma_i v_i^\top$

Frobeniusize the norms!

❗ $\|\cdot\|_F^\dagger = \|\cdot\|_F$ is not a Ky Fan norm, so NSGD is not a Fanion.

Let us consider the convex combination of the Ky Fan rank- k norm and the Frobenius norm, which we call F-KF- k -norm:

$$\|\cdot\|_{F-KF-k} = \alpha \|\cdot\|_{KF-k} + (1 - \alpha) \|\cdot\|_F, \alpha \in [0, 1]$$

Balls of Duals to Convex Combinations of Norms

Lemma

Let $\|\cdot\|_{(1)}$ and $\|\cdot\|_{(2)}$ be norms on a finite-dimensional Euclidean space, and let $\alpha, \beta \geq 0$. Define

$$\|x\| := \alpha\|x\|_{(1)} + \beta\|x\|_{(2)}.$$

Then the dual unit ball of $\|\cdot\|$ satisfies

$$B_{\|\cdot\|^\dagger} = \alpha B_{\|\cdot\|_{(1)}^\dagger} + \beta B_{\|\cdot\|_{(2)}^\dagger},$$

where $+$ denotes the Minkowski sum and $B_{\|\cdot\|_{(i)}^\dagger}$ is the unit ball of the dual norm $\|\cdot\|_{(i)}^\dagger$.

$\|M^k\|_{F-KF-k}^\dagger$ and F-Fanions

Deriving F-Fanions by Recipe

$\|M^t\|_{F-KF-k}^{\dagger\dagger} = \|M^t\|_{F-KF-k}$, and $\Delta = \alpha \sum_{i=1}^k u_i v_i^\top + (1 - \alpha) \frac{M^k}{\|M^k\|_F}$ delivers it:

- ① $\|\Delta\|_{F-KF-k}^\dagger \leq 1$ because $\sum_{i=1}^k u_i v_i^\top$ lies in $B_{\|\cdot\|_{KF-k}}$ and $\frac{M^k}{\|M^k\|_F}$ lies in $B_{\|\cdot\|_F}$, so Δ lies in the Minkowski sum
- ② $\langle M^t, \Delta \rangle = \alpha \langle U \Sigma V^\top, \sum_{i=1}^k u_i v_i^\top \rangle + (1 - \alpha) \langle M^t, \frac{M^t}{\|M^t\|_F} \rangle = \alpha \sum_{i=1}^k \sigma_i + (1 - \alpha) \|M^t\|_F = \|M^t\|_{F-KF-k}$

Hence,

$$X^{t+1} = X^t - \eta \left(\alpha \sum_{i=1}^k u_i v_i^\top + (1 - \alpha) \frac{M^t}{\|M^t\|_F} \right)$$

F-Muon and F-Neon

Convex combinations with Frobenius norm:

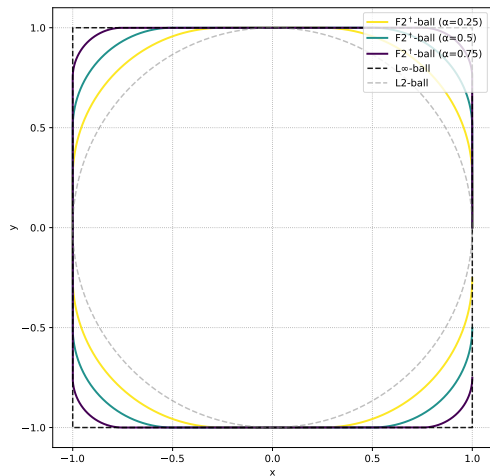
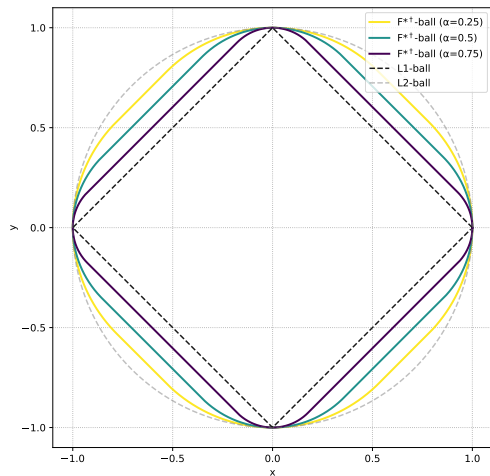
$$\|\mathbf{X}\|_{F*} = \alpha\|\mathbf{X}\|_* + (1 - \alpha)\|\mathbf{X}\|_F, \quad \|\mathbf{X}\|_{F2} = \alpha\|\mathbf{X}\|_{op} + (1 - \alpha)\|\mathbf{X}\|_F.$$

Dual-induced updates:

$$\text{F-Muon: } \mathbf{X}^{k+1} = \mathbf{X}^k - \eta \left(\alpha \mathbf{U} \mathbf{V}^\top + (1 - \alpha) \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_F} \right),$$

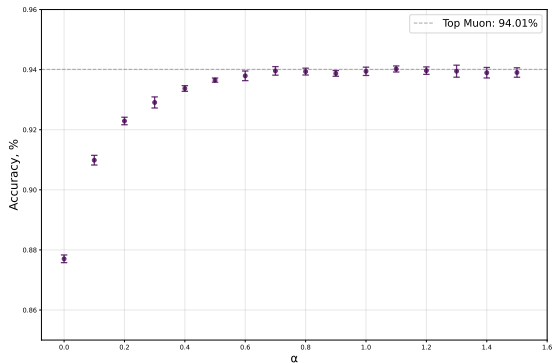
$$\text{F-Neon: } \mathbf{X}^{k+1} = \mathbf{X}^k - \eta \left(\alpha u_1 v_1^\top + (1 - \alpha) \frac{\mathbf{M}^k}{\|\mathbf{M}^k\|_F} \right).$$

Geometric intuition: Dual balls

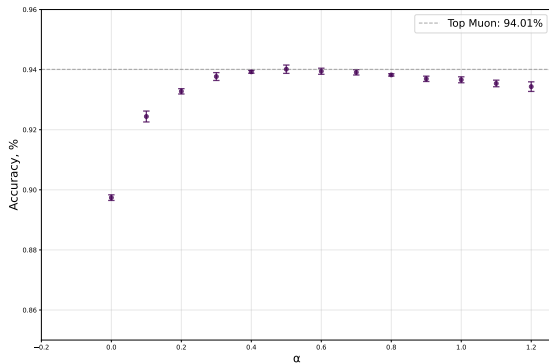


Dual balls for F-Muon and F-Neon across α in a 2D singular-value space

F-Muon on CIFAR-10 airbench



With params tuned for Muon, as in (Keller, 2023)

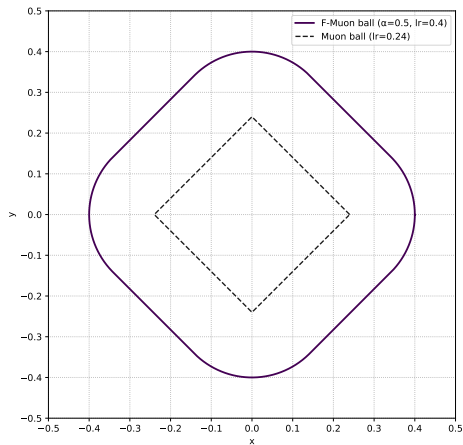


With params tuned for F-Muon

F-Muon with $\alpha = 0.5$ matches Muon after tuning

❓ How should we interpret cases with $\alpha > 1$?

LMO balls on CIFAR-10



LMO balls for learning rates from the experiment

Setup from Jordan et al. (2024a): 1750 iterations and cross-entropy loss lower than 3.28

- Muon: $\text{lr}=0.05$, $\text{momentum}=0.95$, final loss = 3.279
- F-Muon with $\alpha = 0.5$: $\text{lr}=0.07$, $\text{momentum}=0.95$, final loss = 3.281
- NSGD: $\text{lr}=0.07$, $\text{momentum}=0.96$, final loss = 3.4651
- F-Muon with $\alpha = 0.5$: $\text{lr}=0.07$, $\text{momentum}=0.96$, final loss = 3.2824!

Formal assumptions

L-smoothness:

$$\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{X}')\|^\dagger \leq L\|\mathbf{X} - \mathbf{X}'\| \quad \text{for all } \mathbf{X}, \mathbf{X}' \in \mathbb{R}^{m \times n}, \quad (\text{A1})$$

Bounded variance:

$$\mathbb{E}_{\xi \sim \mathcal{D}} [g(\mathbf{X}; \xi)] = \nabla f(\mathbf{X}) \quad \text{and} \quad \mathbb{E}_{\xi \sim \mathcal{D}} [\|g(\mathbf{X}; \xi) - \nabla f(\mathbf{X})\|_F^2] \leq \sigma^2 \quad \text{for all } \mathbf{X} \in \mathbb{R}^{m \times n}, \quad (\text{A2})$$

Bounded norm:

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

From Marchenko-Pastur law, $\|\cdot\|_F \sim \frac{\sqrt{n}}{2} \|\cdot\|_{\text{op}}$ and $\|\cdot\|_{\text{nuc}} \sim \frac{n}{2} \|\cdot\|_{\text{op}}$ for large n .

Bounds from (Kovalev, 2025)

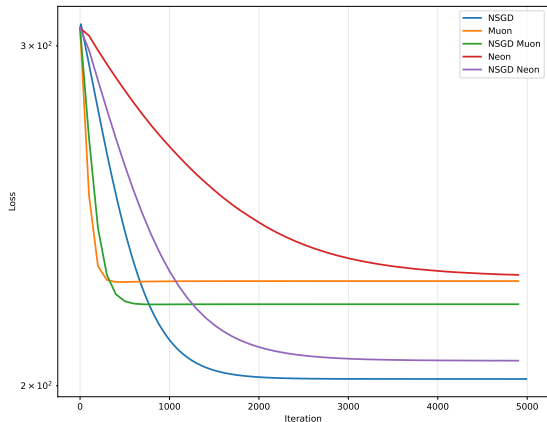
Lemma

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

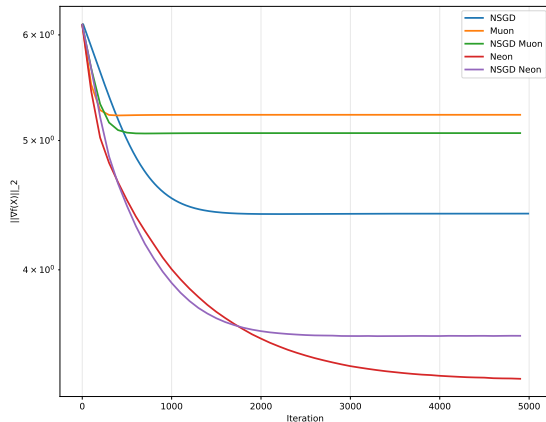
$$\eta = \mathcal{O} \left(\min \left\{ \frac{\varepsilon}{L}, \frac{\varepsilon^3}{\rho^2 \sigma^2 L} \right\} \right), \quad \alpha = \mathcal{O} \left(\min \left\{ 1, \frac{\varepsilon^2}{\rho^2 \sigma^2} \right\} \right), \quad (1)$$

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

Random linear least squares: Loss and the Spectral Norm of the gradient



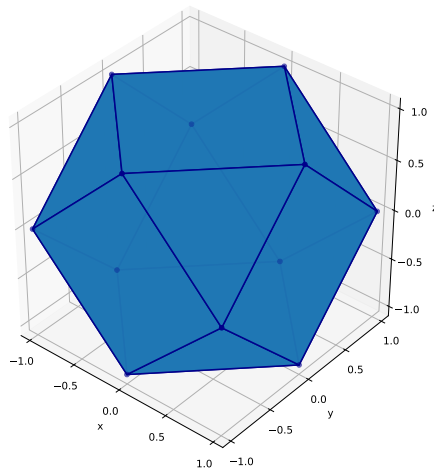
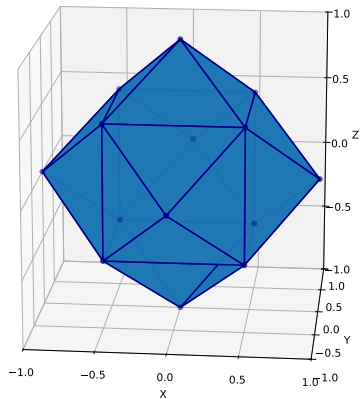
Loss vs iteration



Spectral norm of gradient vs iteration

Linear least squares problem for a 50×50 matrix. Common $\text{lr}=0.01$.

Geometric intuition: Ky Fan norms









Ky Fan rank-2 ball and its dual in a 3D singular-value space

Conclusions

- Muon and NSGD are special cases of F-Fanions
- Mixture of norm-based updates is a norm-based update!
- No existing bounds describe the superiority of the spectral norm

For future experiments: Exploration of matrix and vector algorithms correspondence

References I

-  Bernstein, Jeremy et al. (2018). “signSGD: Compressed optimisation for non-convex problems”. In: *International conference on machine learning*. PMLR, pp. 560–569.
-  Cutkosky, Ashok, Harsh Mehta, and Francesco Orabona (2020). “Momentum-based variance reduction in nonconvex SGD”. In: *Advances in Neural Information Processing Systems*.
-  Hazan, Elad, Kfir Levy, and Shai Shalev-Shwartz (2015). “Beyond convexity: Stochastic quasi-convex optimization”. In: *Advances in neural information processing systems* 28.
-  Jordan, Keller et al. (2024a). *modded-nanogpt: Speedrunning the NanoGPT baseline*. URL: <https://github.com/KellerJordan/modded-nanogpt>.
-  Jordan, Keller et al. (2024b). *Muon: An optimizer for hidden layers in neural networks*. URL: <https://kellerjordan.github.io/posts/muon/>.
-  Keller, Jordan (2023). *cifar10-airbench*. URL: <https://github.com/KellerJordan/cifar10-airbench>.

References II



Kovalev, Dmitry (2025). “Understanding gradient orthogonalization for deep learning via non-euclidean trust-region optimization”. In: *arXiv preprint arXiv:2503.12645*.



Shi, Hao-Jun Michael et al. (2016). “A primer on coordinate descent algorithms”. In: *arXiv preprint arXiv:1610.00040*.