
Modular Optimization

The Great Mind of Jeremy Bernstein

Donghu Kim

NanoGPT Speedrun

NanoGPT (124M) by Andrej Karpathy

The screenshot shows the GitHub repository for NanoGPT. The repository is public and has 402 watches, 6.9k forks, and 41.6k stars. It is managed by karpathy. The repository contains 6 branches and 0 tags. The file list includes assets, config, data, .gitattributes, .gitignore, LICENSE, README.md, bench.py, configurator.py, model.py, sample.py, scaling_laws.ipynb, train.py, and transformer_sizing.ipynb. The right sidebar shows the repository's description, about section, releases, packages, contributors, and languages.

karpathy / nanoGPT

Code Issues 223 Pull requests 68 Actions Projects Security Insights

nanoGPT Public

Watch 402 Fork 6.9k Star 41.6k

master 6 Branches 0 Tags

Go to file Add file Code

karpathy Merge pull request #578 from devin-open-source/devin/1733728337-fix-w... 93a43d9 · 6 months ago 209 Commits

assets	adjust teaser figure with a more tuned result	2 years ago
config	Fix for gradient_accumulation_steps training slow	2 years ago
data	Merge pull request #420 from vinjn/fix-371-enc-is-not-defin...	last year
.gitattributes	keep only what's needed	2 years ago
.gitignore	feature: .gitignore - added venv folders	last year
LICENSE	Add MIT LICENSE file	3 years ago
README.md	Merge branch 'master' into test1	last year
bench.py	Fix AssertionError on macOS - need to check CUDA availabili...	2 years ago
configurator.py	shuttling the poor mans configurator aside into its own file a...	3 years ago
model.py	Merge pull request #274 from apivovarov/gelu	2 years ago
sample.py	Fix AssertionError on macOS - need to check CUDA availabili...	2 years ago
scaling_laws.ipynb	fix typo (params -> tokens)	2 years ago
train.py	fix: ensure non-zero learning rate during warmup at iteratio...	6 months ago
transformer_sizing.ipynb	oops forgot to subtract embedding params, which don't ent...	2 years ago

README MIT license

About

The simplest, fastest repository for training/finetuning medium-sized GPTs.

Readme MIT license Activity 41.6k stars 402 watching 6.9k forks Report repository

Releases

No releases published

Packages

No packages published

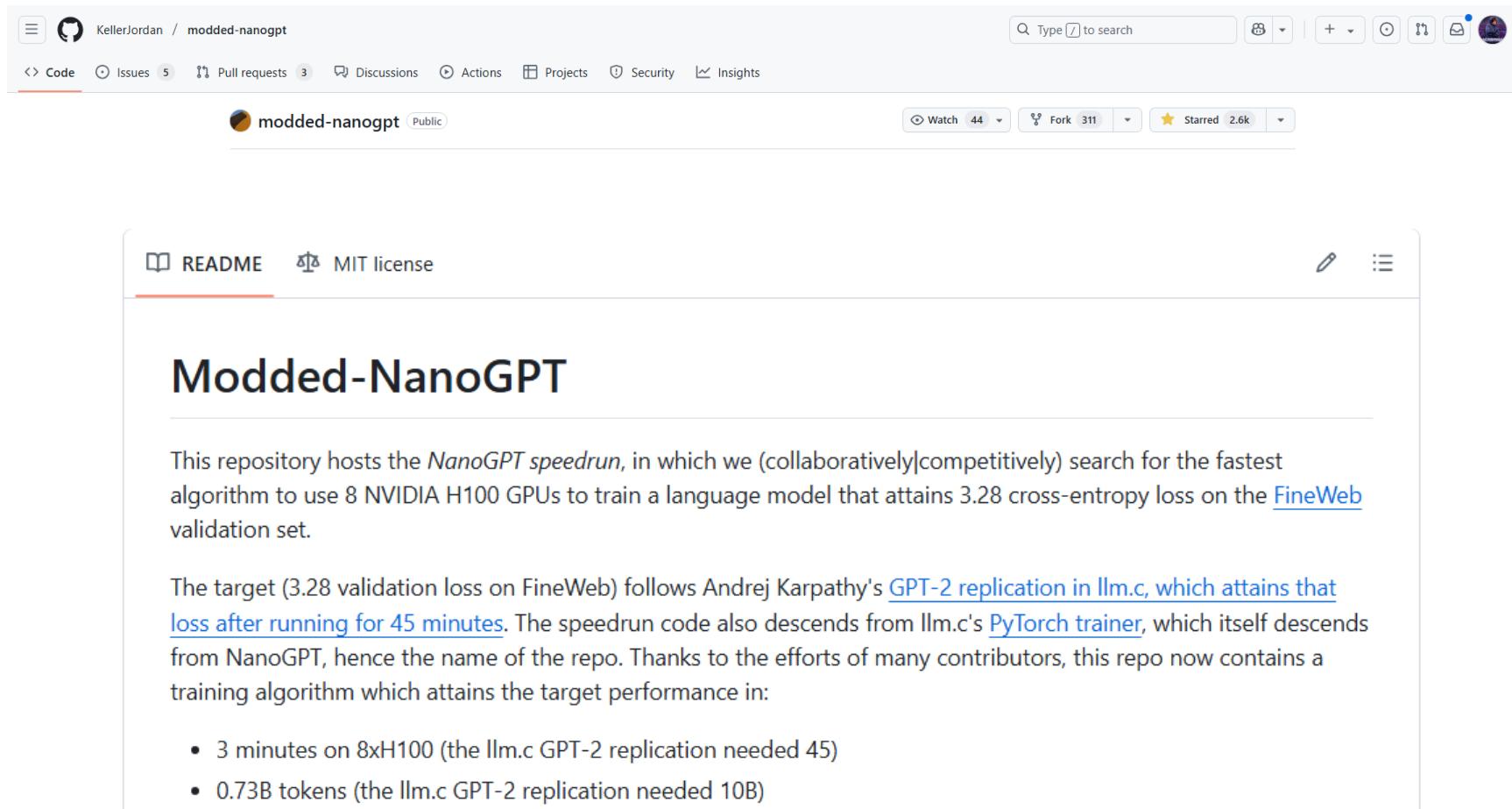
Contributors 37

+ 23 contributors

Languages

NanoGPT Speedrun

A man simply must go *fast*



*But also for a reason: [1]

NanoGPT Speedrun

A man simply must go *fast*

#	Record time	Description	Date	Log	Contributors
1	45 minutes	llm.c baseline	05/28/24	log	@karpathy, llm.c contributors
2	31.4 minutes	Tuned learning rate & rotary embeddings	06/06/24	log	@kellerjordan0

⋮

22	2.990 minutes	Faster gradient all-reduce	05/24/25	log	@KonstantinWilleke, @alexrgilbert, @adricarda, @tuttyfrutye, @vdlad; The Enigma project
23	2.979 minutes	Overlap computation and gradient communication	05/25/25	log	@ryanyang0

NanoGPT Speedrun

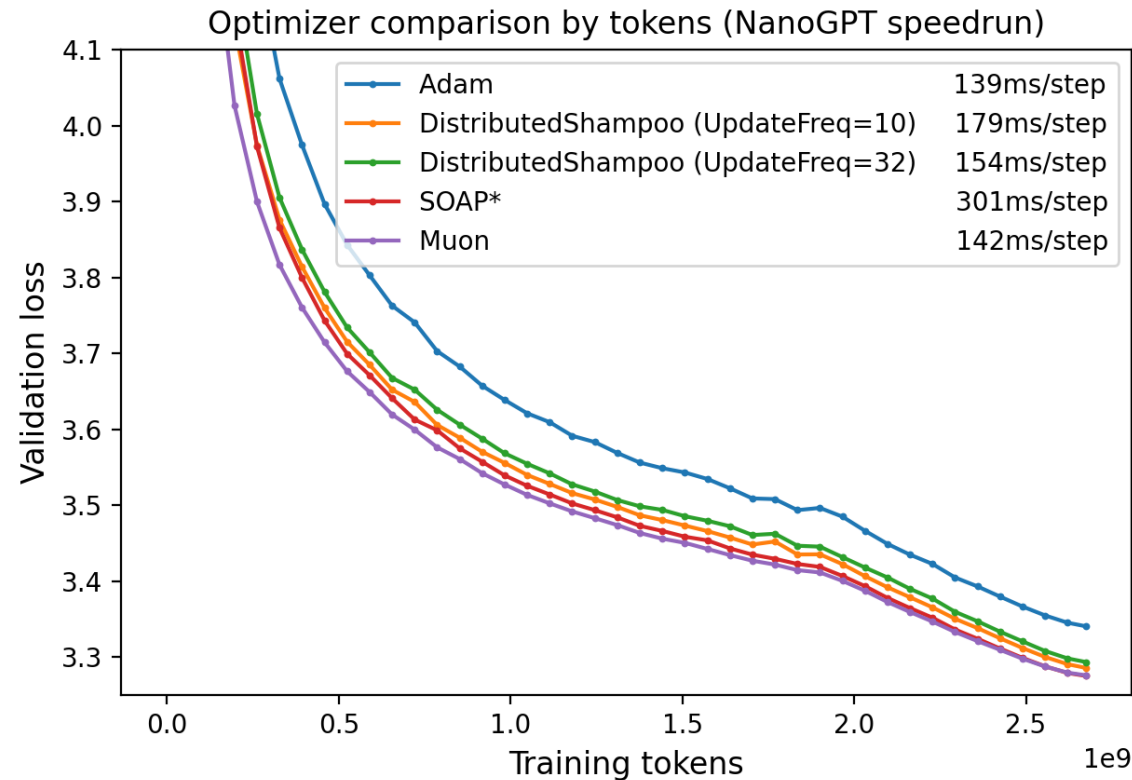
What's that?

#	Record time	Description	Date	Log	Contributors
1	45 minutes	llm.c baseline	05/28/24	log	@karpathy, llm.c contributors
2	31.4 minutes	Tuned learning rate & rotary embeddings	06/06/24	log	@kellerjordan0
3	24.9 minutes	Introduced the Muon optimizer	10/04/24	none	@kellerjordan0, @jxbz
4	22.3 minutes	Muon improvements	10/11/24	log	@kellerjordan0, @bozavlado
5	15.2 minutes	Pad embeddings, ReLU², zero-init projections, QK-norm	10/14/24	log	@Grad62304977, @kellerjordan0
6	13.1 minutes	Distributed the overhead of Muon	10/18/24	log	@kellerjordan0

~11min speedup

Faster Optimization

More sample-efficient than Adam but at the same speed?



Hyperparameter Transfer

Somehow makes scaling easier?

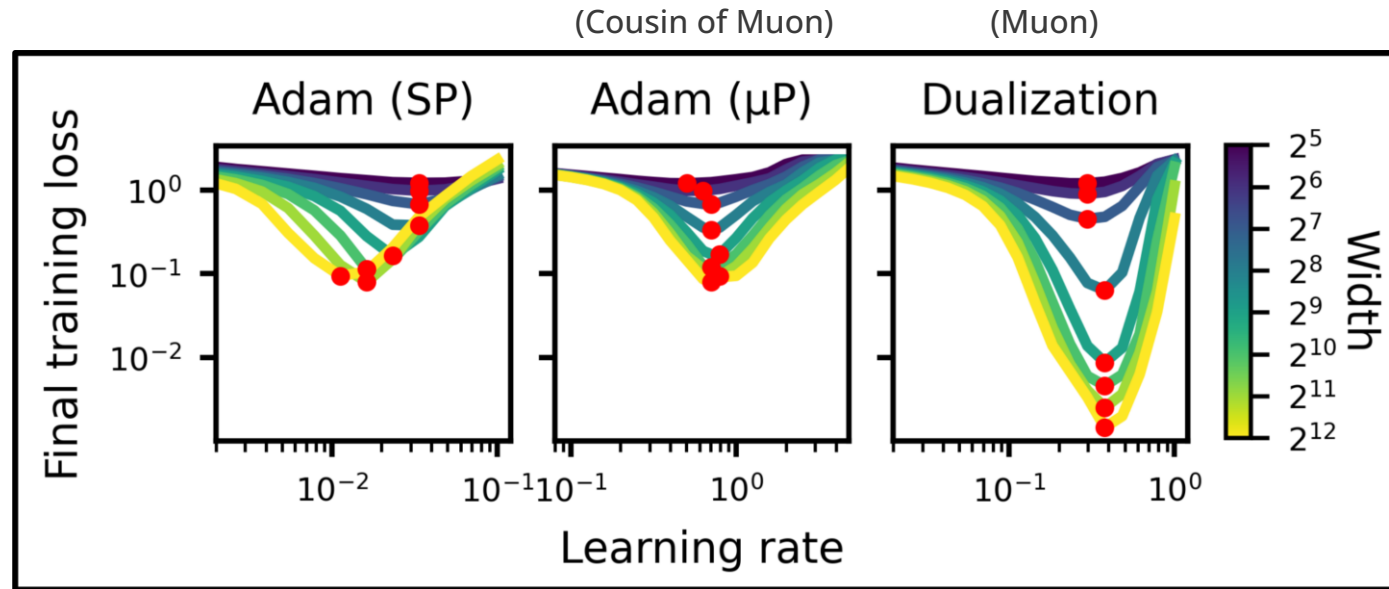


Table 1: **Hyperparameters That Can Be μ Transferred, Not μ Transferred, or μ Transferred Across**, with a few caveats discussed in Section 6.1. * means *empirically validated only on Transformers*, while all others additionally have theoretical justification.

μ Transferable	Not μ Transferable	μ Transferred Across
optimization related, init, parameter multipliers, etc	regularization (dropout, weight decay, etc)	width, depth*, batch size*, training time*, seq length*

Muon

Momentum + Orthogonal Gradients

- | | |
|---------------------|---|
| 1. Compute gradient | $G_t = \nabla_{\theta} L$ |
| 2. Update momentum | $B_t = \mu B_{t-1} + G_t$ |
| 3. Orthogonalize | $B_t = U \Sigma V^T \rightarrow O_t = UV^T$ |
| 4. Update | $\theta_t = \theta_{t-1} - \eta O_t$ |

Why would orthogonal gradients help?

Today's Goal

Understand that

muon is a steepest descent under spectral norm

and why it's a good idea.

Today's Goal

Understand that

*muon is a steepest descent under spectral norm
and why it's a good idea.*

I. Preliminary

Some Norms Are Induced by Others

Consider the linear transformation:

$$\mathbf{x} \xrightarrow{A} A\mathbf{x}$$

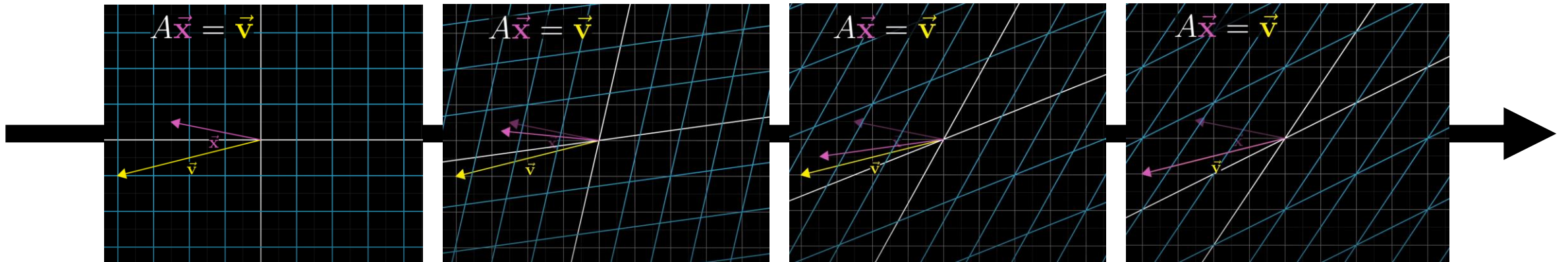
How do we measure the 'size' of A ?

Some Norms Are Induced by Others

Operator Norm (Induced Norm)

$$\mathbf{x} \xrightarrow{A} A\mathbf{x}$$

A is an operator. A is defined by how it changes \mathbf{x}

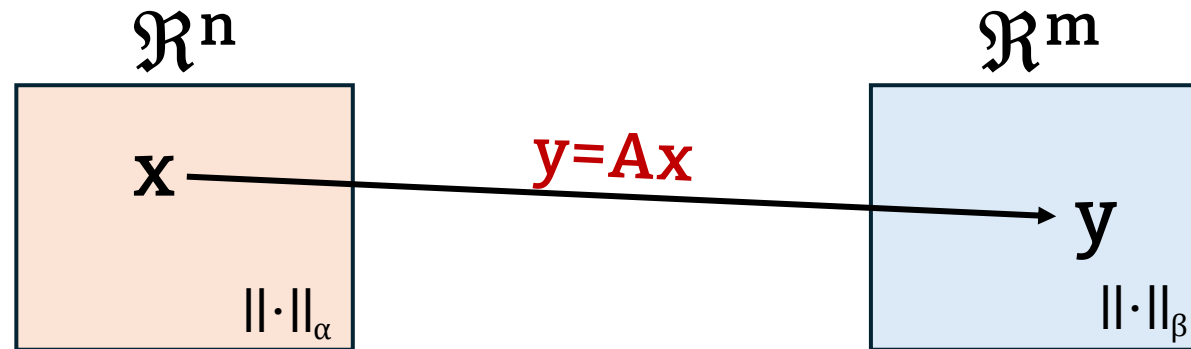


The 'size' of A should thus be defined by **how much it changes \mathbf{x}** .

Some Norms Are Induced by Others

Operator Norm (Induced Norm)

A is a linear operator that maps one space (equipped with $\|\cdot\|_\alpha$) to other space (with $\|\cdot\|_\beta$).
 α and β can be any type of norm of our choice.



Then the operator norm A is defined by the **maximum norm growth from input to output**:

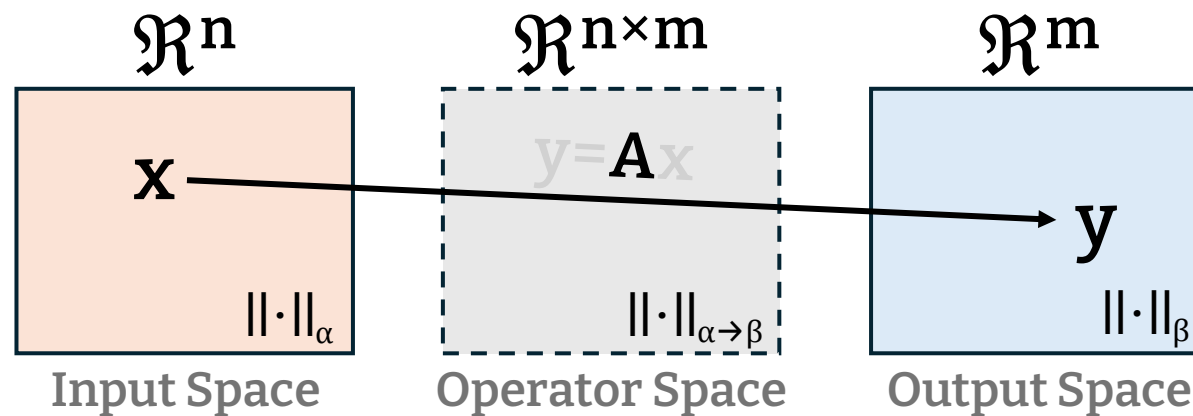
$$\|A\|_{\alpha \rightarrow \beta} \triangleq \sup_{x \neq 0} \frac{\|Ax\|_\beta}{\|x\|_\alpha} = \sup_{\|x\|_\alpha = 1} \|Ax\|_\beta$$

Foreshadowing: we should use this to measure the norm of weight matrices in ML/DL!

Some Norms Are Induced by Others

Operator Norm (Induced Norm)

Now we can draw the full diagram:



Example: Mapping from Euclidean to Euclidean (ℓ_2 -to- ℓ_2)

Then the operator norm of A is the **spectral norm**:

$$\|A\|_{\ell_2 \rightarrow \ell_2} \triangleq \sup_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = (\text{largest singular value of } A) = \|A\|_* = \|A\|_{S_\infty}$$

* ℓ_2 -to- ℓ_2 operator norm = Spectral norm = Schatten- ∞ norm

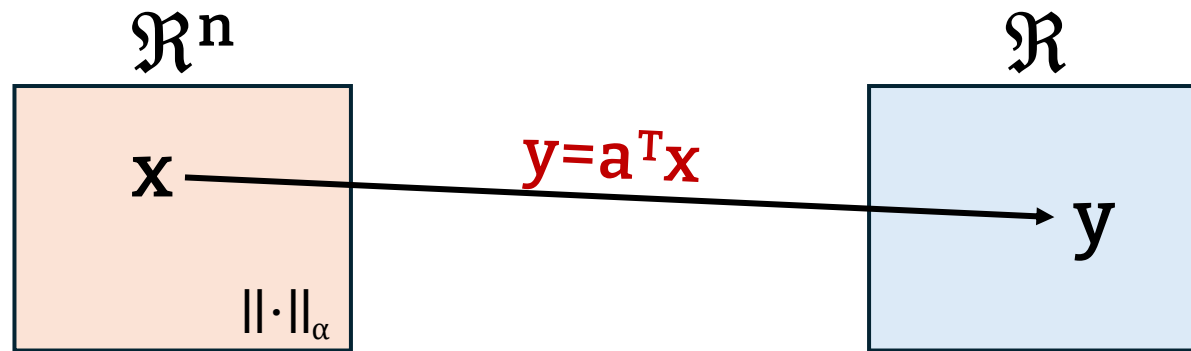
Some Norms Are Induced by Others

Dual Norm

A special case of operator norm, where \mathbf{A} is a vector rather than a matrix.

$$\mathbf{x} \xrightarrow{a} \mathbf{a}^T \mathbf{x}$$

Still an operator, but maps to a scalar rather than another vector space!

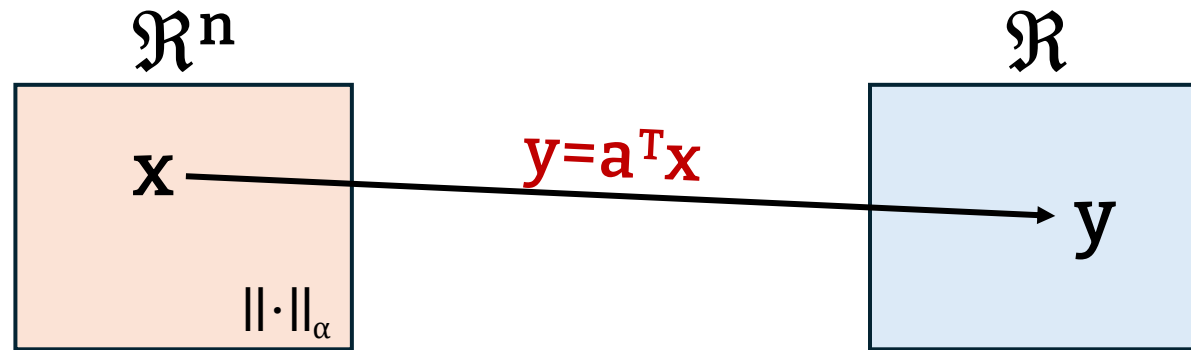


Some Norms Are Induced by Others

Dual Norm

\mathbf{a} is a linear operator that maps a vector space (equipped with $\|\cdot\|_\alpha$) to a scalar.

α can be any type of norm of our choice.



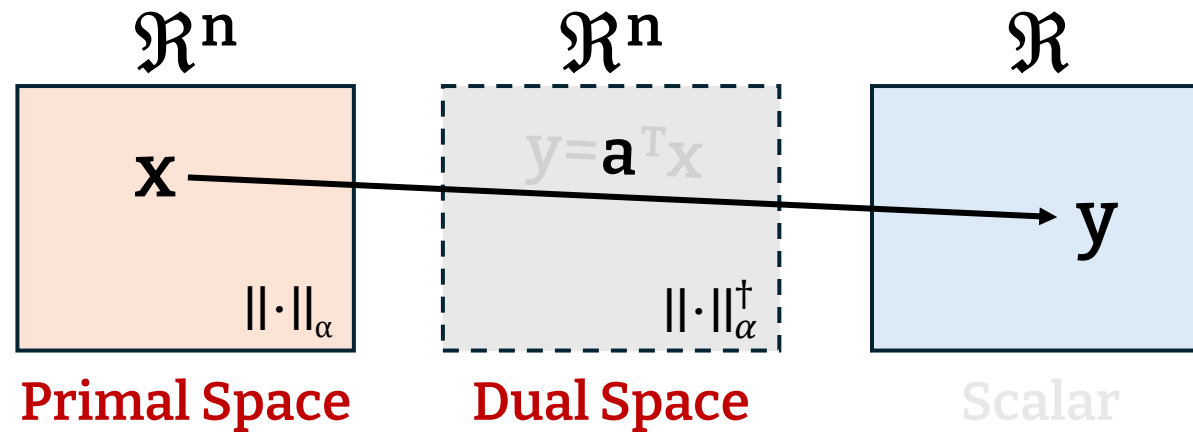
This defines the **dual norm of α** , indicated by the dagger(\dagger):

$$\|\mathbf{a}\|_\alpha^\dagger \triangleq \sup_{\mathbf{x} \neq 0} \frac{\mathbf{a}^T \mathbf{x}}{\|\mathbf{x}\|_\alpha} = \sup_{\|\mathbf{x}\|_\alpha = 1} \mathbf{a}^T \mathbf{x}$$

Some Norms Are Induced by Others

Dual Norm

We also have a special name for the spaces in this case:



Important note: Dual of ℓ_p is ℓ_q , where $1/p + 1/q = 1$

Example: Dual of ℓ_∞ is ℓ_1 , dual of ℓ_1 is ℓ_∞ .

Gradients Live in Dual Space

When I was a wee little undergraduate...

“Gradients tell us how fast the loss changes...
isn't it so weird that we directly subtract it from the parameter?”

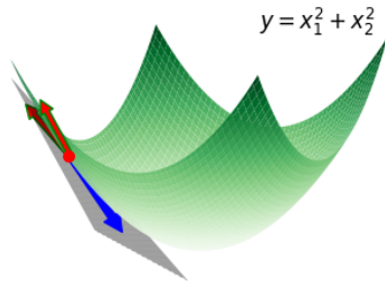
Intentional confusion

On higher dimension: moving x to minimize $y = f(x)$, $x \in \mathbb{R}^2$

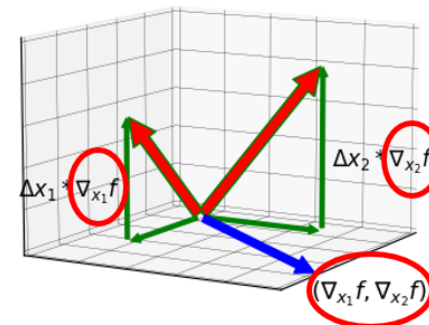
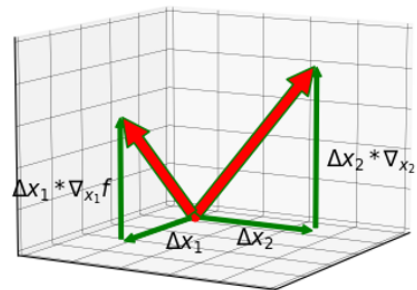
If x is currently at a , how much should x move? $-\nabla_x f(a)$!

Gradient $\nabla_x f(a) = \begin{bmatrix} \nabla_{x_1} f(a) \\ \nabla_{x_2} f(a) \end{bmatrix}$ consists of 2 partial derivatives to corresponding axes

Seems pretty “descent-y”



...but the blue arrow still seems to be out of nowhere!



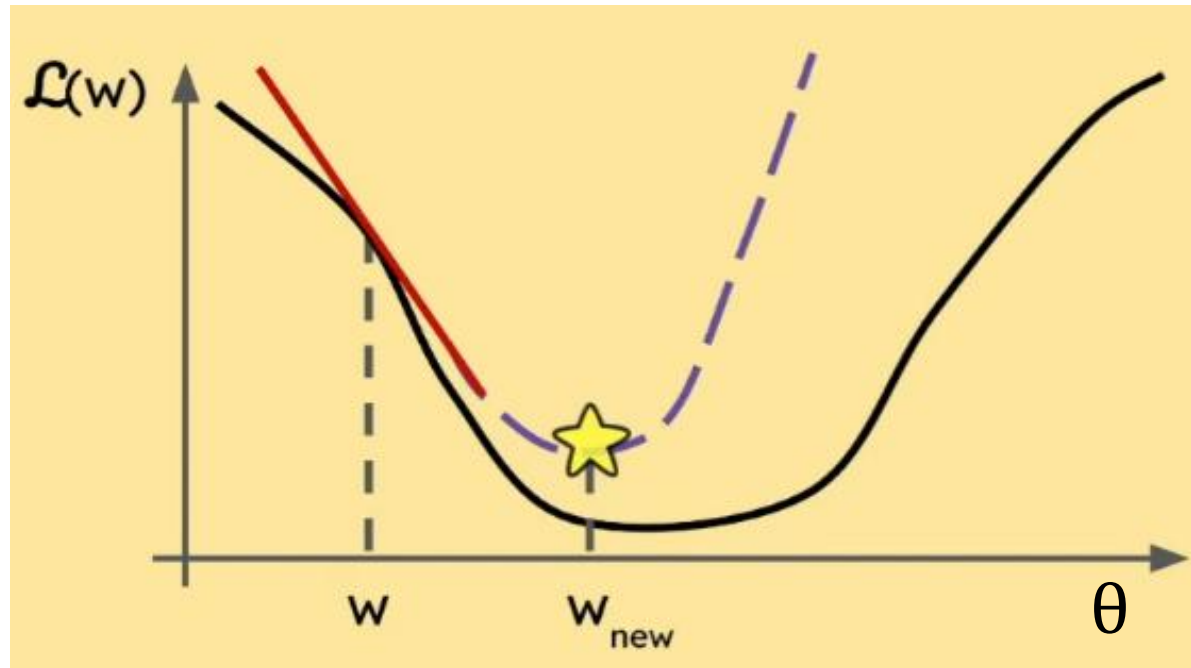
Gradients Live in Dual Space

↑↑↑↑↑ That's Because ↑↑↑↑↑

The gradient $\nabla_{\theta} L$ is a first-order approximation on how fast the loss L changes (near w).

$$L(w + \Delta w) \approx L(w) + g^T \Delta w$$

$$g^T \Delta w \approx L(w + \Delta w) - L(w) = \Delta L$$



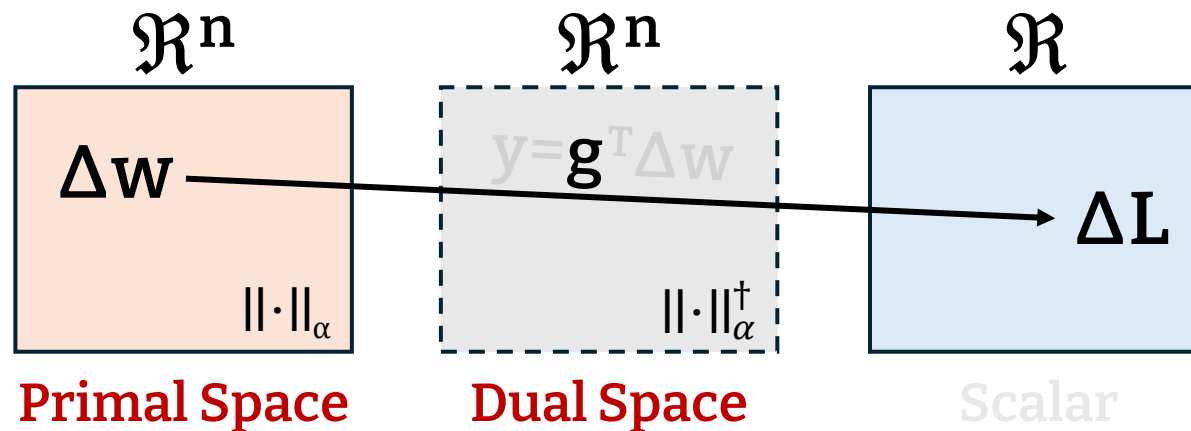
Gradients Live in Dual Space

↑↑↑↑↑ That's Because ↑↑↑↑↑

i.e., the gradient is a linear operator on Δw (that approximates ΔL).

$$\Delta w \xrightarrow{g} g^T \Delta w \quad g = \nabla_{\theta} L \quad g^T \Delta w \approx \Delta L$$

So even though Δw and g are both \mathbb{R}^n , they in fact live in different spaces: primal and dual!



Foreshadowing: everything gradient related will involve the dualization:

$$\sup_{\|x\|_{\alpha}=1} g^T x$$

Gradients Live in Dual Space

Why didn't we care about this the entire time?

Because we don't need to *if* we're using Euclidean norm!

1. Recall that: dual of ℓ_p is ℓ_q , where $1/p + 1/q = 1$

2. So, if we're using Euclidean norm (ℓ_2) on the parameters (primal space):

$$\|g\|_2^+ \triangleq \sup_{\|x\|_2=1} g^T x = \|g\|_2$$

The gradients (dual space) are also measured by Euclidean norm!

3. More importantly (and jumping a bit ahead):

$$\text{Steepest Descent} = \frac{\|g\|_2^+}{\lambda} \operatorname{argmax}_{\|x\|_2=1} g^T x = -\frac{\|g\|_2}{\lambda} \frac{g}{\|g\|_2} = -\frac{1}{\lambda} g$$

Gradient IS the steepest descent under Euclidean norm!

Summary

Operator norm $\alpha \rightarrow \beta$ is defined by the maximum norm growth from input to output:

$$\|A\|_{\alpha \rightarrow \beta} \triangleq \sup_{x \neq 0} \frac{\|Ax\|_{\beta}}{\|x\|_{\alpha}} = \sup_{\|x\|_{\alpha}=1} \|Ax\|_{\beta}$$

We will use this to measure the norm of weight matrices in neural networks.

The **dual norm of α** is a special case of operator norm where the output is a scalar

$$\|a\|_{\alpha}^{\dagger} \triangleq \sup_{x \neq 0} \frac{a^T x}{\|x\|_{\alpha}} = \sup_{\|x\|_{\alpha}=1} a^T x$$

Gradients live in dual space (of parameters), so it will always involve the dualization:

$$\sup_{\|x\|_{\alpha}=1} g^T x$$

So far...

Understand that

*muon is a steepest descent under spectral **norm***

and why it's a good idea.

II. Steepest Descent

Optimization Algorithms

Gradient-based optimizers are basically...

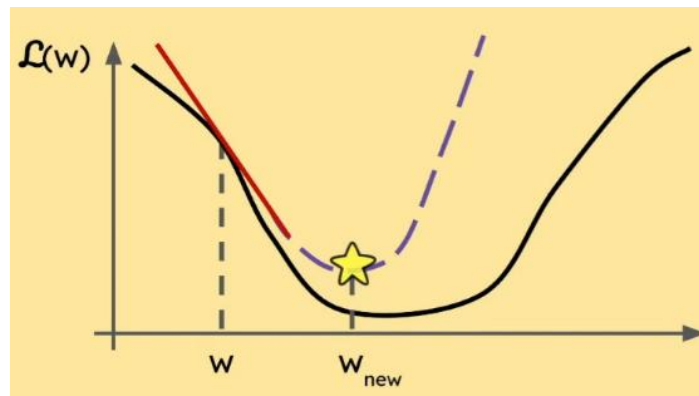
1. Taylor Expansion on the loss function

$$L(w+\Delta w) = L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T H \Delta w + \dots$$

2. Approximate high-order terms

$$L(w+\Delta w) \approx L(w) + g^T \Delta w + D(w, w+\Delta w)$$

3. Minimize (find Δw that minimize RHS)



Optimization Algorithms

Different higher-order modeling = Different optimization

$$L(w + \Delta w) = L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T H \Delta w + \dots$$

ex1) Euclidean Norm \rightarrow Gradient Descent (1st order method)

$$\approx L(w) + g^T \Delta w + \frac{1}{2} \lambda^* \|\Delta w\|_2^2$$

Optimization Algorithms

Different higher-order modeling = Different optimization

$$L(w + \Delta w) = L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T H \Delta w + \dots$$

ex1) Euclidean Norm \rightarrow Gradient Descent (1st order method)

$$\approx L(w) + g^T \Delta w + \frac{1}{2} \lambda \|\Delta w\|_2^2$$

ex2) Hessian Matrix \rightarrow Newton's Method (2nd order)

$$\approx L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T H \Delta w$$

Optimization Algorithms

Different higher-order modeling = Different optimization

$$L(w + \Delta w) = L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T H \Delta w + \dots$$

ex1) Euclidean Norm \rightarrow Gradient Descent (1st order method)

$$\approx L(w) + g^T \Delta w + \frac{1}{2} \lambda \|\Delta w\|_2^2$$

ex2) Hessian Matrix \rightarrow Newton's Method (2nd order)

$$\approx L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T H \Delta w$$

ex3) Some distance function D on output f (e.g., TRPO: KL divergence)
 \rightarrow Natural Gradient Descent (2nd order)

$$\approx L(w) + g^T \Delta w + D(f, f + \Delta f)$$

$$\approx L(w) + g^T \Delta w + \frac{1}{2} \Delta w^T (\nabla_w f^T \nabla_f^2 D \nabla_w f) \Delta w$$

Optimization Algorithms

A lot to argue about which is better, but let's just say that...



Same Steepest Descent, Different Norms

...and we will see that these all use the same formula as SGD



...but with different norms!

Same Steepest Descent, Different Norms



Same Steepest Descent, Different Norms

The first-order method:

$$\begin{aligned} L(\mathbf{w} + \Delta \mathbf{w}) &= L(\mathbf{w}) + \mathbf{g}^\top \Delta \mathbf{w} + \frac{1}{2} \Delta \mathbf{w}^\top \mathbf{H} \Delta \mathbf{w} + \dots \\ &\leq L(\mathbf{w}) + \mathbf{g}^\top \Delta \mathbf{w} + \frac{1}{2} \lambda^* \|\Delta \mathbf{w}\|^2 \end{aligned}$$

Sharpness parameter Square of *any* norm
e.g., L2 → Gradient descent

This is also called a 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 $\|\cdot\|^\dagger$:

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

Same Steepest Descent, Different Norms

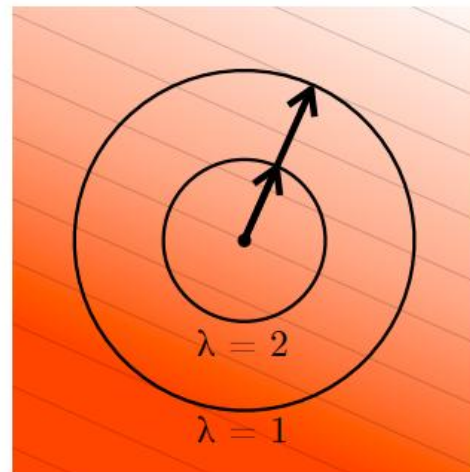
Simple intuition: λ decides step size, $\|\cdot\|$ decides direction

$$\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right] = \underbrace{-\frac{\|\mathbf{g}\|^\dagger}{\lambda}}_{\text{Step size}} \cdot \underbrace{\arg \max_{\|t\|=1} \mathbf{g}^\top t}_{\text{Step direction}}.$$

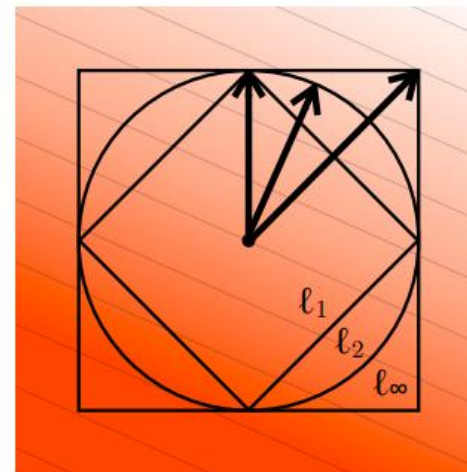
Step size Step direction

If the landscape is sharp, take smaller steps

1. Draw a unit ball of selected norm
2. Find the direction that changes L the most



a) varying sharpness λ



b) varying choice of norm $\|\cdot\|$

Same Steepest Descent, Different Norms

Extra intuition: gradients are always ‘dualized’

$$\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right] = - \underbrace{\frac{\|\mathbf{g}\|^\dagger}{\lambda}}_{\text{Step size}} \cdot \underbrace{\arg \max_{\|\mathbf{t}\|=1} \mathbf{g}^\top \mathbf{t}}_{\text{Step direction}}.$$

which uses dual norm on \mathbf{g} also called a duality map on \mathbf{g}

Recall: \mathbf{g} is an operator on \mathbf{w} , so any operation on \mathbf{g} must depend on \mathbf{w} (or \mathbf{t}).

Definition 1 (Dual norm). Given a norm $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$, the dual norm $\|\cdot\|^\dagger$ of a vector $\mathbf{g} \in \mathbb{R}^n$ is given by:

$$\|\mathbf{g}\|^\dagger := \max_{\mathbf{t} \in \mathbb{R}^n : \|\mathbf{t}\|=1} \mathbf{g}^\top \mathbf{t}. \quad (5)$$

Definition 2 (Duality map based on a norm). Given a norm $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$, we consider the duality map:

$$\text{dualize}_{\|\cdot\|} \mathbf{g} := \arg \max_{\mathbf{t} \in \mathbb{R}^n : \|\mathbf{t}\|=1} \mathbf{g}^\top \mathbf{t}, \quad (6)$$

where, if the $\arg \max$ is not unique, $\text{dualize}_{\|\cdot\|}$ returns any maximizer.

Same Steepest Descent, Different Norms

Example: Vanilla Gradient Descent is steepest under ℓ_2

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

$$\|\mathbf{g}\|_2^\dagger \triangleq \sup_{\|\mathbf{x}\|_2=1} \mathbf{g}^\top \mathbf{x} = \|\mathbf{g}\|_2$$

$$\arg \max_{\|\mathbf{x}\|_2=1} \mathbf{g}^\top \mathbf{x} = \frac{\mathbf{g}}{\|\mathbf{g}\|_2}$$

$$\frac{\|\mathbf{g}\|_2^\dagger}{\lambda} \arg \max_{\|\mathbf{x}\|_2=1} \mathbf{g}^\top \mathbf{x} = \frac{\|\mathbf{g}\|_2}{\lambda} \frac{\mathbf{g}}{\|\mathbf{g}\|_2} = \frac{1}{\lambda} \mathbf{g}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{\lambda} \mathbf{g}$$

Same Steepest Descent, Different Norms

Example: Adam without EMA is steepest under ℓ_∞

Adam's update rule:

$$\begin{aligned} \mathbf{m}_t &= \beta_1 \cdot \mathbf{m}_{t-1} + (1 - \beta_1) \cdot \mathbf{g}_t, \\ \mathbf{v}_t &= \beta_2 \cdot \mathbf{v}_{t-1} + (1 - \beta_2) \cdot \mathbf{g}_t^2, \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \cdot \mathbf{m}_t / \sqrt{\mathbf{v}_t}, \end{aligned}$$

Adam without EMA ($\beta_1 = \beta_2 = 0$) is a sign descent:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta \cdot \mathbf{g}_t / \sqrt{\mathbf{g}_t^2} \\ &= \mathbf{w}_t - \eta \cdot \text{sign}(\mathbf{g}_t). \end{aligned}$$

Note: Same applies to any sign descent algorithms such as RMSProp and Lion.

Note: "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"

Old Optimizer, New Norm: An Anthology., Bernstein & Newhouse.

Same Steepest Descent, Different Norms

Example: Adam without EMA is steepest under ℓ_∞

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

$$\|g\|_\infty^\dagger \triangleq \sup_{\|x\|_\infty=1} g^\top x = \|g\|_1$$

$$\arg \max_{\|x\|_\infty=1} g^\top x = \text{sign}(g)$$

$$\frac{\|g\|_\infty^\dagger}{\lambda} \arg \max_{\|x\|_\infty=1} g^\top x = \frac{\|g\|_1}{\lambda} \text{sign}(g)$$

$$w_{t+1} = w_t - \frac{\|g\|_1}{\lambda} \text{sign}(g)$$

Summary

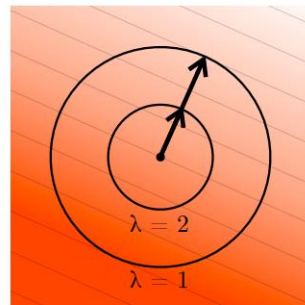
The first-order method:

$$\begin{aligned} L(\mathbf{w} + \Delta \mathbf{w}) &= L(\mathbf{w}) + \mathbf{g}^\top \Delta \mathbf{w} + \frac{1}{2} \Delta \mathbf{w}^\top \mathbf{H} \Delta \mathbf{w} + \dots \\ &\leq L(\mathbf{w}) + \mathbf{g}^\top \Delta \mathbf{w} + \frac{1}{2} \lambda^* \|\Delta \mathbf{w}\|^2 \end{aligned}$$

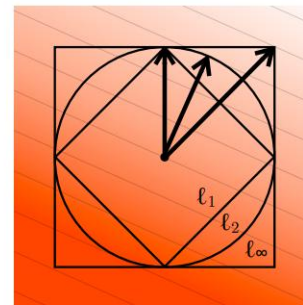
↓
Square of *any* norm

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

and the **choice of norm** results in a completely different algorithm!



a) varying sharpness λ



b) varying choice of norm $\|\cdot\|$

Next Question

Which norm should be use?

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

$\ell_1, \ell_2, \dots, \ell_\infty$?

Hold on...

w is not a vector, they are matrices with structure...

Can't we use matrix norms here?

So far...

Understand that

*muon is a **steepest descent** under spectral **norm***

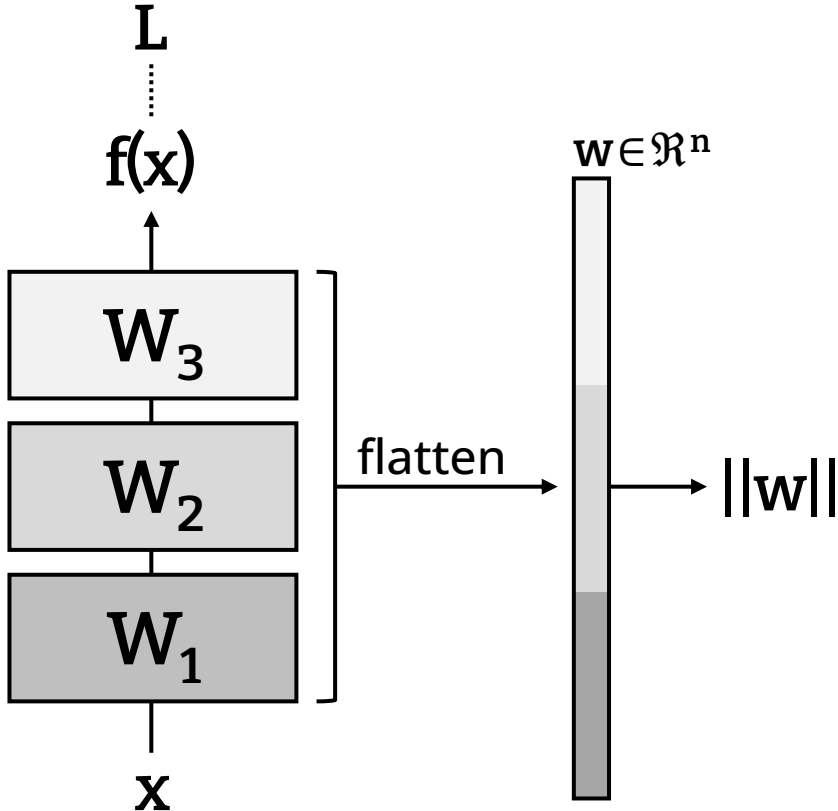
and why it's a good idea.

III. Modular Steepest Descent

Norms Should be Defined Layerwise

What we've been doing so far:

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



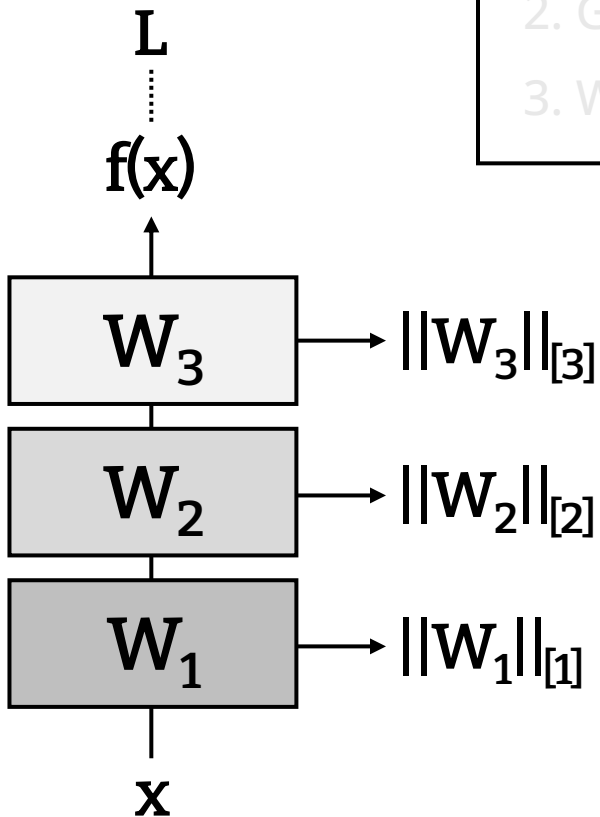
Flattening completely erases the structure of the network!

- Number of layers (Depth)
- Size of each layer (layerwise width)
- Layer type (CNN, MLP, Attention, ...)
- **Weights are matrices, not vectors!**
- **Weights are operators, not just matrices!**

Norms Should be Defined Layerwise

What we should be doing:

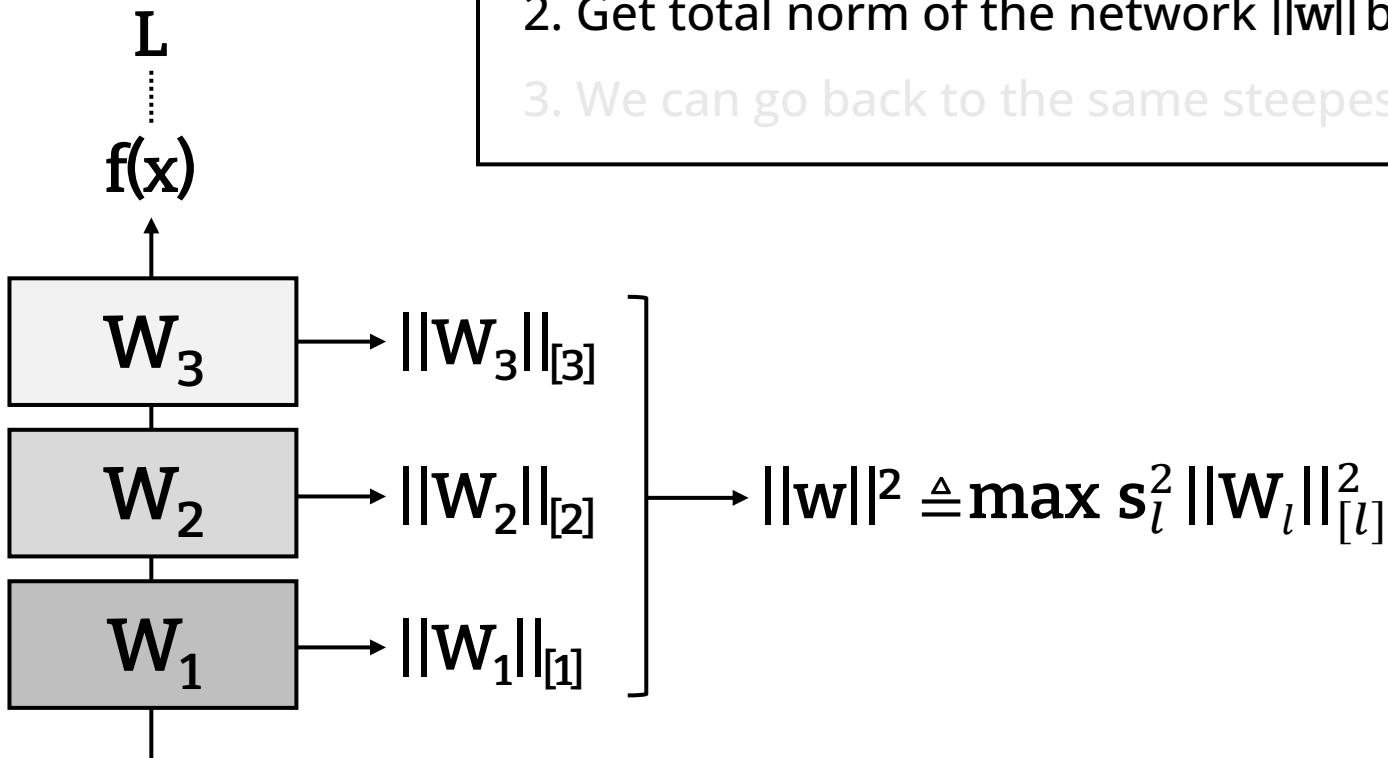
1. Measure each matrix individually: $\|\cdot\|_{[1]}, \|\cdot\|_{[2]}, \dots, \|\cdot\|_{[L]}$
2. Get total norm of the network $\|w\|$ by weighted max.
3. We can go back to the same steepest descent.



Norms Should be Defined Layerwise

What we should be doing:

1. Measure each matrix individually: $\|\cdot\|_{[1]}, \|\cdot\|_{[2]}, \dots, \|\cdot\|_{[L]}$
2. Get total norm of the network $\|w\|$ by weighted max.
3. We can go back to the same steepest descent.



Note: The weights s seems to connect to 'sensitivity' in later works (for now they're just set to 1).

Note: The max operation kinda makes sense when you actually derive the results (although I'm still unsure of the intuition). See Appendix C.

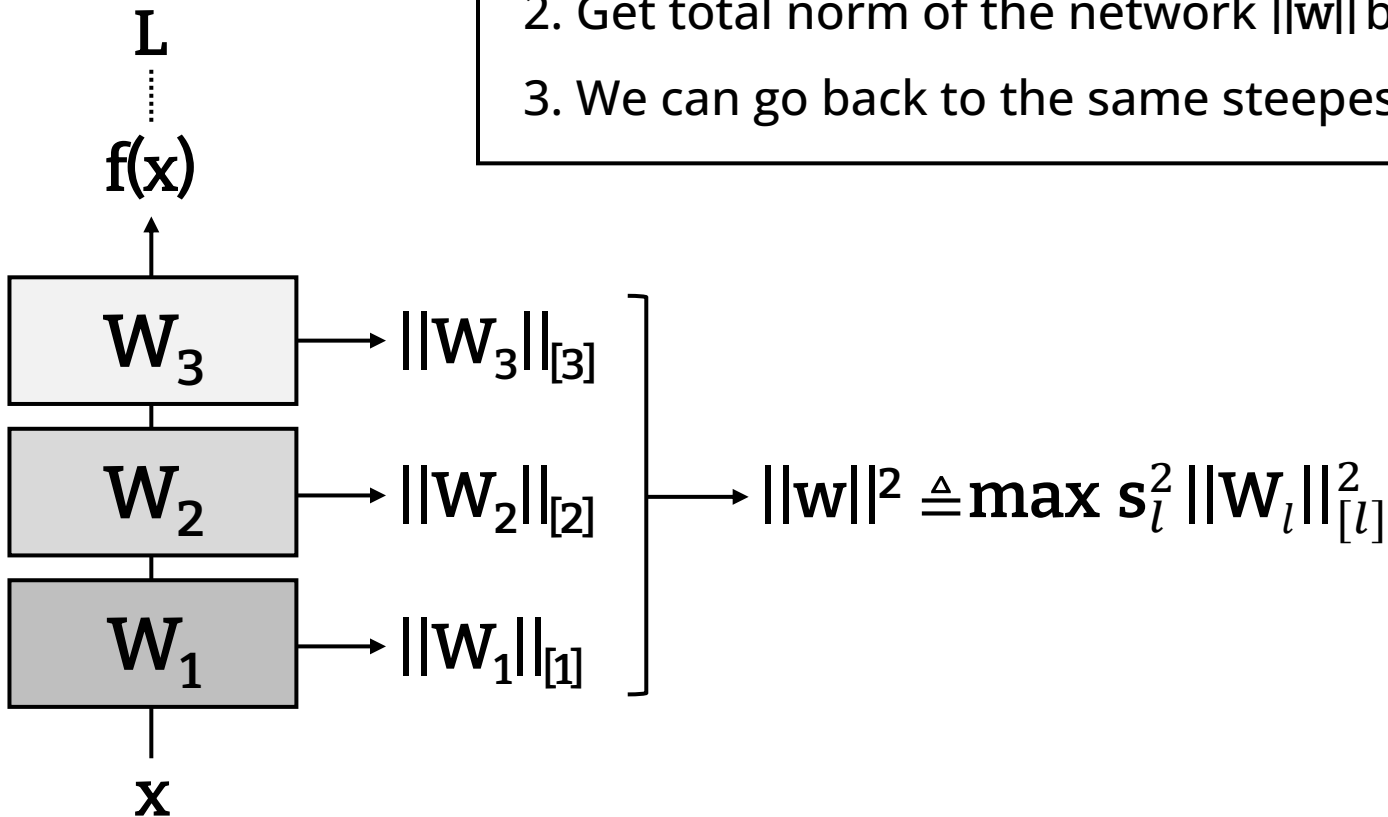
■ Scalable Optimization in the Modular Norm., Large et al.

■ Modular Duality in Deep Learning., Bernstein & Newhouse.

Norms Should be Defined Layerwise

What we should be doing:

1. Measure each matrix individually: $\|\cdot\|_{[1]}, \|\cdot\|_{[2]}, \dots, \|\cdot\|_{[L]}$
2. Get total norm of the network $\|w\|$ by weighted max.
3. We can go back to the same steepest descent.



Steepest Descent

$$\arg \min_{\Delta w \in \mathbb{R}^n} \left[g^\top \Delta w + \frac{\lambda}{2} \|\Delta w\|^2 \right]$$

Modular Steepest Descent

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

Norms Should be Defined Layerwise

No big difference other than “is it layerwise or not”.

Steepest Descent

$$\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right]$$

$$\Delta \mathbf{w} = -\frac{\|\mathbf{g}\|^\dagger}{\lambda} \operatorname{argmax}_{\|t\|=1} \mathbf{g}^\top t$$

Norm on flattened vector

Can use only vector norm

All layers are jointly solved

Modular Steepest Descent

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

$$\Delta \mathbf{W}_i = -\frac{\eta}{s_i} \operatorname{argmax}_{\|T_i\|=1} \langle \mathbf{G}_i, T_i \rangle$$

Norm on each matrix

Can use matrix norm

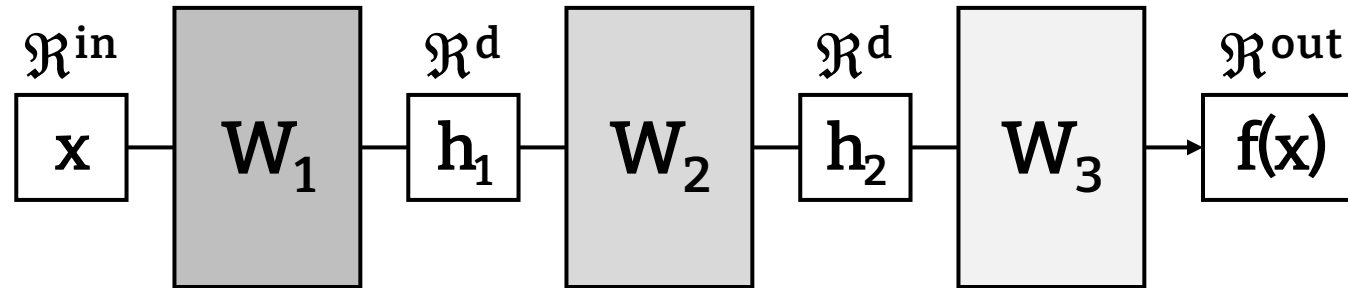
Each layer is individually solved

Now the question is: “Which matrix norm do we use?”

Which Norm Should We Use?

We need to first ask: what are the matrices doing?

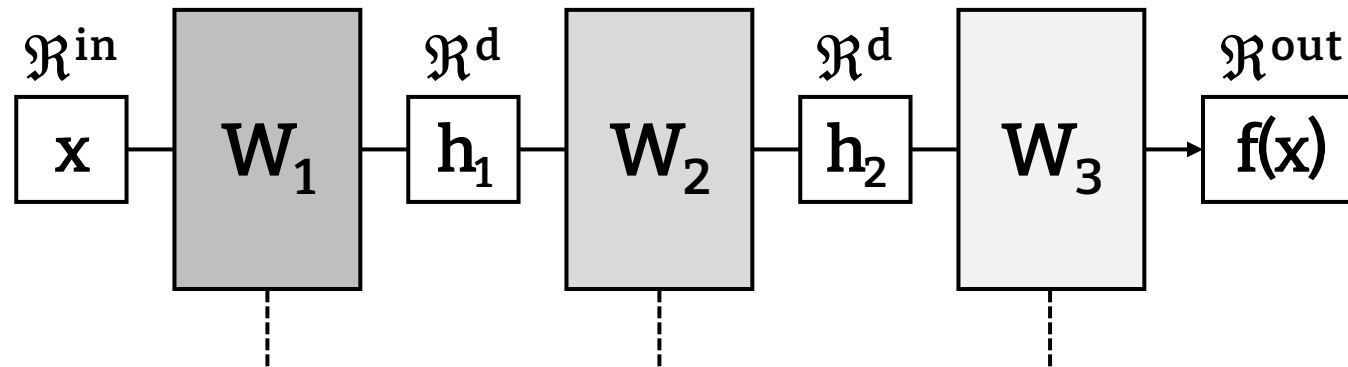
They are linear operators that map (each of their) input space to output space!



Which Norm Should We Use?

We need to first ask: what are the matrices doing?

They are linear operators that map (each of their) input space to output space!



Which means their norms are induced by the norm of the features they work on!

$$\|W_1\|_{x \rightarrow h_1}$$

$$\|W_2\|_{h_1 \rightarrow h_2}$$

$$\|W_3\|_{h_2 \rightarrow y}$$

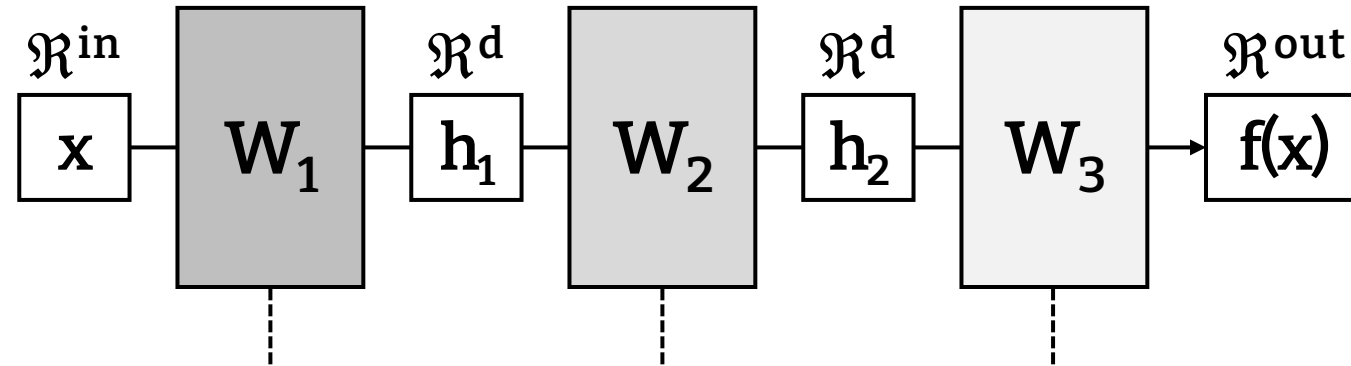
Recall:

$$\|A\|_{\alpha \rightarrow \beta} \triangleq \sup_{x \neq 0} \frac{\|Ax\|_{\beta}}{\|x\|_{\alpha}} = \sup_{\|x\|_{\alpha}=1} \|Ax\|_{\beta}$$

Which Norm Should We Use?

We need to first ask: what are the matrices doing?

They are linear operators that map (each of their) input space to output space!



Which means their norms are induced by the norm of the features they work on!

$$\|W_1\|_{x \rightarrow h_1} \quad \|W_2\|_{h_1 \rightarrow h_2} \quad \|W_3\|_{h_2 \rightarrow y}$$

Now the question becomes: “Which feature norm do we use?”

Which Feature Norm Do We Use?

Example: Adam without EMA is steepest under max $\ell_1 \rightarrow \ell_\infty$ norm.

Interestingly, Adam can also be thought as using $\ell_1 \rightarrow \ell_\infty$ norm on every weight matrix.

First, the $\ell_1 \rightarrow \ell_\infty$ norm is simply the largest entry of the matrix.

$$\|A\|_{\ell_1 \rightarrow \ell_\infty} \triangleq \sup_{x \neq 0} \frac{\|Ax\|_\infty}{\|x\|_1} = \sup_{\|x\|_1=1} \|Ax\|_\infty = \max_{i,j} |A_{i,j}|$$

*Maximized when $x = \text{one-hot}(j)$ where max entry is at A_{ij}

Then, ℓ_∞ is (coincidentally) the maximum of $\ell_1 \rightarrow \ell_\infty$ norms (a.k.a. max-of-max norm)!

$$\|w\|_\infty = \max_l \max_r \|\text{row}_r(W_l)\|_\infty = \max_l \|W_l\|_{\ell_1 \rightarrow \ell_\infty}$$

$$\arg \min_{\Delta w \in \mathbb{R}^n} \left[g^\top \Delta w + \frac{\lambda}{2} \|\Delta w\|^2 \right]$$

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

Which Feature Norm Do We Use?

Example: Adam without EMA is steepest under max $\ell_1 \rightarrow \ell_\infty$ norm.

Is $\ell_1 \rightarrow \ell_\infty$ natural? Probably not...

Still, they ARE doing modular steepest descent, which may explain why they work so well.

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

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](#).

Still still, we can probably do better than this!

Summary

Norms should be measured layer-by-layer

Steepest Descent

$$\arg \min_{\Delta \mathbf{w} \in \mathbb{R}^n} \left[\mathbf{g}^\top \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right]$$

$$\Delta \mathbf{w} = -\frac{\|\mathbf{g}\|^\dagger}{\lambda} \operatorname{argmax}_{\|t\|=1} \mathbf{g}^\top t$$

Modular Steepest Descent

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

$$\Delta \mathbf{W}_i = -\frac{\eta}{s_i} \operatorname{argmax}_{\|T_i\|=1} \langle \mathbf{G}_i, T_i \rangle$$

Question: Which matrix norm should we use?

Answer: Matrices in NNs are operators, so it should be an **operator norm** ($\alpha \rightarrow \beta$).

New Question: Which **feature norm** (α, β) then?

$$\|\mathbf{A}\|_{\alpha \rightarrow \beta} \triangleq \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|_\beta}{\|\mathbf{x}\|_\alpha} = \sup_{\|\mathbf{x}\|_\alpha=1} \|\mathbf{A}\mathbf{x}\|_\beta$$

So far...

Understand that

muon is **a steepest descent under modular norm**

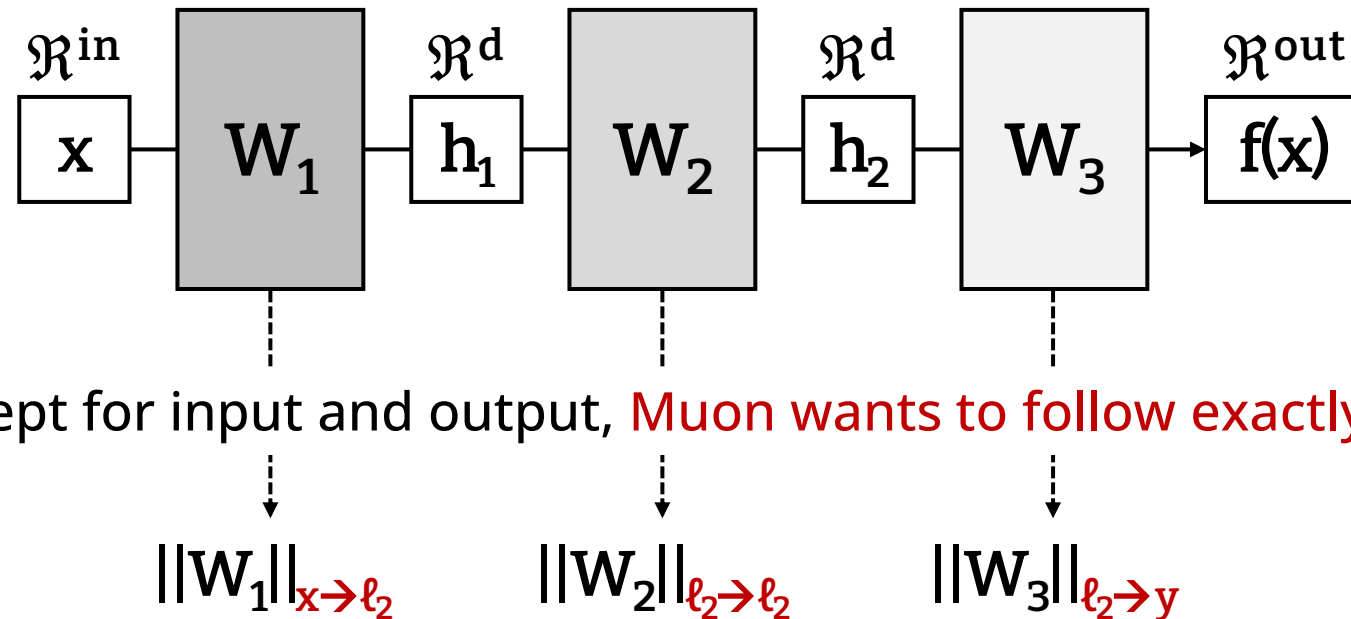
and why it's a good idea.

IV. Deriving Muon

Which Feature Norm Does Muon Use?

Muon asks: which feature norm **are** we using?

Since we love LayerNorm sooo much, we're using ℓ_2 norm* almost everywhere!



So except for input and output, **Muon wants to follow exactly that!**

We can define input and output norms, but let's only think about hidden features.
In fact, Muon just uses Adam on input and output layers.

Muon

Muon w/o momentum is steepest under $\ell_2 \rightarrow \ell_2$ norm

- | | |
|-------------------------|---|
| 1. Compute gradient | $G_t = \nabla_{\theta} L$ |
| 2. Update momentum | $B_t = \mu B_{t-1} + G_t$ |
| 3. Orthogonalize | $B_t = U \Sigma V^T \rightarrow O_t = UV^T$ |
| 4. Update | $\theta_t = \theta_{t-1} - \eta O_t$ |

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

$$\Delta W_i = -\frac{\eta}{s_l} \operatorname{argmax}_{\|T_i\|=1} \langle G_i, T_i \rangle$$

$$\operatorname{argmax}_{\|T\|_{\ell_2 \rightarrow \ell_2} = 1} \langle G, T \rangle = \mathbf{UV}^T$$

Muon

Muon w/o momentum is steepest under $\ell_2 \rightarrow \ell_2$ norm

$$\operatorname{argmax}_{\|T\|_{\ell_2 \rightarrow \ell_2} = 1} \langle G, T \rangle = \mathbf{UV}^T$$

$$\Delta W_i = -\frac{\eta}{s_l} \operatorname{argmax}_{\|T_i\| = 1} \langle G_i, T_i \rangle$$

$$\begin{aligned} \max_{\|X\|_{l_2 \rightarrow l_2} = 1} \langle G, X \rangle &= \max_{\|X\|_{l_2 \rightarrow l_2} = 1} \operatorname{tr}(G^T X) \\ &= \max_{\|X\|_{l_2 \rightarrow l_2} = 1} \operatorname{tr}(V \Sigma U^T X) \quad (\text{Spectral decomposition } G = U \Sigma V^T) \\ &= \max_{\|X\|_{l_2 \rightarrow l_2} = 1} \operatorname{tr}(\Sigma U^T X V) \quad (\text{Cycle property of trace: } \operatorname{tr}(ABC) = \operatorname{tr}(BCA) = \operatorname{tr}(CAB)) \\ &= \operatorname{tr}(\Sigma U^T U V^T V) \quad \text{Maximized by: } X = UV^T \\ &= \operatorname{tr}(\Sigma) \end{aligned}$$

Orthogonalization via Newton-Schulz

Should we perform SVD every time we update?

- | | |
|-------------------------|---|
| 1. Compute gradient | $G_t = \nabla_{\theta} L$ |
| 2. Update momentum | $B_t = \mu B_{t-1} + G_t$ |
| 3. Orthogonalize | $B_t = U \Sigma V^T \rightarrow O_t = UV^T$ |
| 4. Update | $\theta_t = \theta_{t-1} - \eta O_t$ |

People have found a much faster way to find **UV^T** without performing SVD!

The Newton-Schulz iteration:

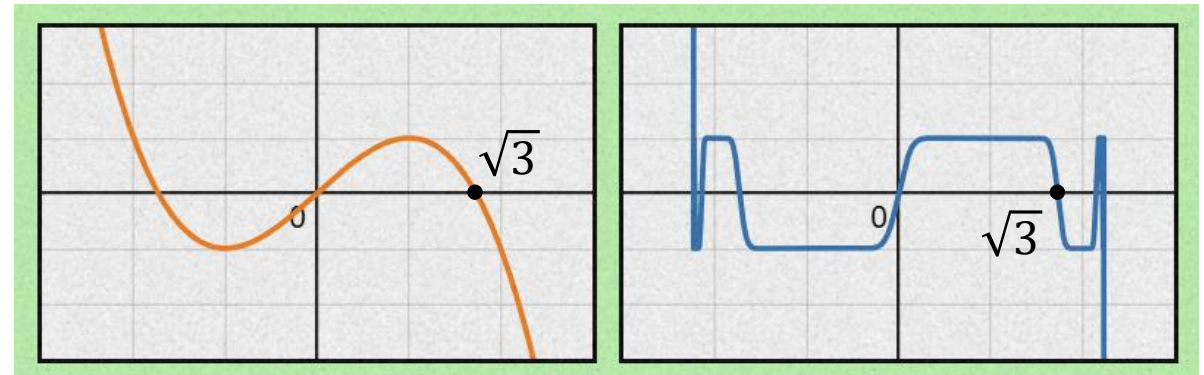
$$X_{t+1} = \frac{2}{3}X_t - \frac{1}{2}X_t X_t^T X_t$$

...that's it. X gets closer to orthogonal after every iter.

Orthogonalization via Newton-Schulz

I can't believe it's that easy!

$$\begin{aligned}X_{t+1} &= \frac{2}{3}X_t - \frac{1}{2}X_tX_t^TX_t \\&= \frac{2}{3}U\Sigma V^T - \frac{1}{2}U\Sigma V^T \cdot V\Sigma^TU^T \cdot U\Sigma V^T \\&= \frac{2}{3}U\Sigma V^T - \frac{1}{2}U\Sigma^3V^T \\&= U \underbrace{\left(\frac{2}{3}\Sigma - \frac{1}{2}\Sigma^3\right)}_{f(\sigma) \text{ on each entry}} V^T\end{aligned}$$



$f(\sigma)$

$fffff(\sigma)$

As long as $0 \leq \sigma \leq \sqrt{3}$, every iteration gets σ closer to 1!

Note: Muon uses higher polynomials with better coefficients in practice, but they basically do the same thing but faster.

Note: The $0 \leq \sigma \leq \sqrt{3}$ condition can be satisfied by normalizing the matrix by its Frobenius norm first (which doesn't affect the final output).

So far...

Understand that

muon is a steepest descent under spectral norm

and why it's a good idea.

V. Unfinished Business

Unanswered Questions

(And Where to Find the Answers)

Q1. Why is spectral descent a good idea?

A1. This blog (3min): <https://jeremybernste.in/writing/deriving-muon>
tl;dr – Because it takes the largest improvement within a safe range.

For the third step in the derivation, we consider choosing a weight update to maximize the linear improvement in loss \mathcal{L} while maintaining a bound on the amount that the outputs can change in response. The rationale is that if the weight update makes the layer outputs change too much, this could destabilize the overall network. In symbols, we would like to solve:

$$\min_{\Delta W} \langle \nabla_W \mathcal{L}, \Delta W \rangle \quad \text{subject to} \quad \|\Delta y\|_{\text{RMS}} \leq \eta.$$

change to directly controlling the size of the weight update itself. If the input has size $\|x\|_{\text{RMS}} \leq 1$, we obtain the following problem as a proxy:

$$\min_{\Delta W} \langle \nabla_W \mathcal{L}, \Delta W \rangle \quad \text{subject to} \quad \|\Delta W\|_{\text{RMS} \rightarrow \text{RMS}} \leq \eta. \quad (\dagger)$$

$$\Delta W = -\eta \times \sqrt{\frac{\text{fan-out}}{\text{fan-in}}} \times UV^\top.$$

Also recommended: [Spectral Condition](#) / [Appendix J of Tensor Programs V](#)

Unanswered Questions

(And Where to Find the Answers)

Q2. Why should we use L2 or RMS norm?

A2. We don't need to. There is little known about what is the best norm.
(The modular norm paper came out late 2024, sooo...)
In fact, there are words that EMA allows optimizers like Adam or Shampoo to adjust to the 'right' norm for each layer.

 <https://x.com/leloykun/status/1847919153589735705>

Unanswered Questions

(And Where to Find the Answers)

Q3. What will happen in the future?

A3. Bernstein is expanding on the idea of modular optimization. He believes that every layer can be designed like lego blocks, which will make it easier to understand what's going on inside the NNs.

 <https://jeremybernste.in/writing/deriving-muon> Conclusion

 <https://docs.modula.systems/>

 Scalable Optimization in the Modular Norm., Large et al.

 Modular Duality in Deep Learning., Bernstein, Newhouse.

Unanswered Questions

(And Where to Find the Answers)

Q4. Has Muon been applied to larger LLMs?

A4. Yes.

 Muon is Scalable for LLM Training., Liu et al.

Q5. Has Muon been applied to RL?

A5. Not yet. Our attempt while building V-Simba was unsuccessful. Could not outperform Adam, but I may not have been cautious enough with the implementation details (e.g., on convolution).

 Modular Duality in Deep Learning., Bernstein, Newhouse.

 <https://x.com/jxbz/status/1846188906733044029>

 <https://x.com/tianyin/status/1896542545557262606>

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