

# Modular Duality in Deep Learning

Zeyu Bian  
zebian@ucsd.edu

University of Science and Technology of China

May 11, 2025

- 1 Introduction
  - Intuition
  - Descent in Normed Spaces
- 2 Modular Norm
  - Invariant Optimal Step Size
  - Module and Modular Norm
  - Building Compound modules
- 3 Duality Map
  - Basic Modules
  - Computation
  - Concrete Examples
- 4 Summary and Takeaway
- 5 Appendix
  - Additional Modules
  - Omitted Proofs

- 1 Introduction
  - Intuition
  - Descent in Normed Spaces
- 2 Modular Norm
  - Invariant Optimal Step Size
  - Module and Modular Norm
  - Building Compound modules
- 3 Duality Map
  - Basic Modules
  - Computation
  - Concrete Examples
- 4 Summary and Takeaway
- 5 Appendix
  - Additional Modules
  - Omitted Proofs

## Shall we directly subtract the gradient from the weights?

The author argues we shouldn't, due to the reason that

- The geometry of the loss function may be non-isotropic, so we need to adjust the **size** and **direction** of the gradient.

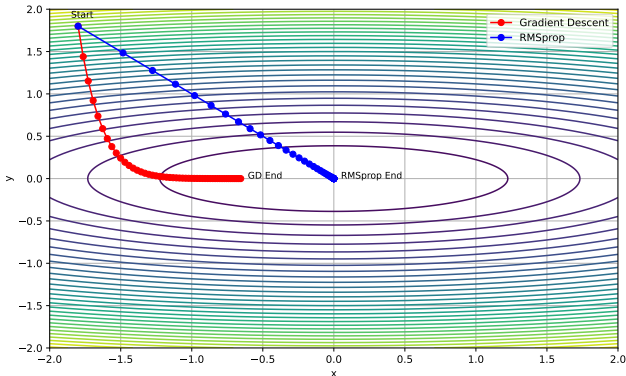


Figure 1: GD vs. RMSprop dynamic on  $Loss(x, y) = 0.1x^2 + y^2$

# Intuition From Dual Space Perspective

## Definition (Dual Space)

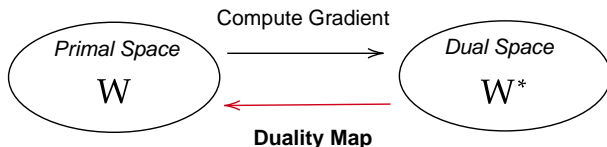
Given a vector space  $\mathcal{V}$ ,

- Say  $f : \mathcal{V} \rightarrow \mathbb{R}$  is a *linear functional* on  $\mathcal{V}$  if  $f$  is linear.
- Define the dual space  $\mathcal{V}^*$  to be the set of *linear functionals* on the vector space  $\mathcal{V}$ .

Consider the Taylor expansion of the loss function:

$$\mathcal{L}(\mathbf{w} + \Delta\mathbf{w}) = \mathcal{L}(\mathbf{w}) + \nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w})^\top \Delta\mathbf{w} + \text{higher-order terms} . \quad (1)$$

The gradient  $\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w})$  is a linear functional paired with  $\Delta\mathbf{w}$ , so that  $\nabla_{\mathbf{w}}\mathcal{L}(\mathbf{w}) \in \mathcal{W}^*$



# Descent in Normed Spaces


We use the diamond operator  $\diamond$  denotes tensor contraction. <sup>1</sup>

Suppose that we wish to use gradient descent to minimize a loss function  $\mathcal{L} : \mathcal{W} \rightarrow \mathbb{R}$  over a weight space  $\mathcal{W} = \mathbb{R}^N$ . Then we may desire the following three properties :

- 1 The loss function is differentiable, meaning that the gradient map  $\nabla_{\mathbf{w}} \mathcal{L} : \mathcal{W} \rightarrow \mathcal{W}$  exists;
- 2 The weight space  $\mathcal{W}$  carries a norm  $\| \cdot \| : \mathcal{W} \rightarrow \mathbb{R}$ .
- 3 The loss is Lipschitz smooth in the norm  $\| \cdot \|$ , with sharpness constant  $\lambda > 0$ , which implies:

$$\mathcal{L}(\mathbf{w} + \Delta \mathbf{w}) \leq \mathcal{L}(\mathbf{w}) + \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \diamond \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2. \quad (2)$$

---

<sup>1</sup>e.g. If  $\mathbf{x}$  and  $\mathbf{y}$  are compatible 4d tensors, then  $\mathbf{x} \diamond \mathbf{y} = \sum x_{ijkl} y_{ijkl}$  

# Descent in Normed Spaces

Under the above three properties, the weight update given by

$$\Delta \mathbf{w} = \arg \min \left[ \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}) \diamond \Delta \mathbf{w} + \frac{\lambda}{2} \|\Delta \mathbf{w}\|^2 \right]. \quad (3)$$

is guaranteed to reduce the loss (since the minimum is no greater than 0)  
We further give two definition to help discussions.

## Definition (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:

$$\|\mathbf{g}\|^\dagger := \max_{\mathbf{t} \in \mathbb{R}^n: \|\mathbf{t}\|=1} \mathbf{g}^\top \mathbf{t}$$

## Example

The dual norm of  $\|\cdot\|_p$  is  $\|\cdot\|_q$  where  $(p,q)$  is conjugate pair,  $1 \leq p \leq \infty$

# Descent in Normed Spaces

## Definition (Duality map based on a norm)

Given a norm  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ , we consider the duality map:

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

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

We have the following expression for minimizing (3) in vector:

## Proposition (Steepest descent under a norm)

For any  $\mathbf{g} \in \mathbb{R}^n$ ,  $\lambda \geq 0$ , norm  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$  with dual norm  $\|\cdot\|^\dagger$  and duality map  $\text{dualize}_{\|\cdot\|}$ :

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

$\|\cdot\|$  influences the direction, while the sharpness  $\lambda$  influences the size.

We would ideally like the optimal step-size to remain invariant as we scale the width and the depth of the network. Thus, a fundamental problem is to design a norm such that:

- The upper bound (2) actually holds (and is not hopelessly lax).
- The corresponding sharpness constant  $\lambda$  is **invariant** to the relevant architectural dimensions.

We shall see later that the **modular norm** meets these requirements.

## Example (Duality Map)

- Euclidean norm: dualize  $\|\cdot\|_2 \mathbf{g} = \mathbf{g} / \|\mathbf{g}\|_2$
- Infinity norm: dualize  $\|\cdot\|_\infty \mathbf{g} = \text{sign}(\mathbf{g})$

# Norm Examples

## Definition (Induced operator norm)

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

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

For tensors, define the duality map via dualize  $\|\cdot\|_{\mathbf{G}} := \arg \max_{\|T\|=1} \mathbf{G} \diamond T$

## Definition (RMS norm)

For vector  $\mathbf{v} \in \mathbb{R}^d$ ,  $\|\mathbf{v}\|_{\text{RMS}} := \|\mathbf{v}\|_2 / \sqrt{d}$

## Example (The RMS $\rightarrow$ RMS operator norm)

Given a matrix  $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ , we have

$\|\mathbf{W}\|_{\text{RMS} \rightarrow \text{RMS}} = \sqrt{d_{\text{in}}/d_{\text{out}}} \times \|\mathbf{W}\|_*$ , where  $\|\cdot\|_*$  denotes the standard spectral norm.

# Duality Map Examples

Linear layers:

## Example (Duality map for the RMS $\rightarrow$ RMS operator norm)

For a matrix  $\mathbf{G} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$  with reduced singular value decomposition  $\mathbf{G} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ , the corresponding duality map is given by dualize  $\|\cdot\|_{\text{RMS} \rightarrow \text{RMS}} \mathbf{G} = \sqrt{d_{\text{out}}/d_{\text{in}}} \times \mathbf{U}\mathbf{V}^\top$ .

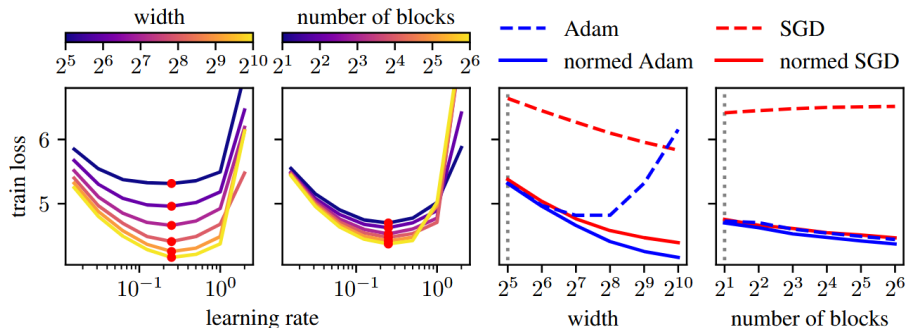
Embedding layers:

## Example (Duality map for the $\ell_1 \rightarrow$ RMS operator norm)

Given a matrix  $\mathbf{W} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ ,  $\|\mathbf{W}\|_{\ell_1 \rightarrow \text{RMS}} = \max_i \|\text{col}_i(\mathbf{W})\|_{\text{RMS}}$ .  
For a matrix  $\mathbf{G} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ , the corresponding duality map dualize  $\|\cdot\|_{\ell_1 \rightarrow \text{RMS}} \mathbf{G}$  simply normalizes each column of  $\mathbf{G}$  to have unit RMS norm:  $\text{col}_i(\mathbf{G}) \mapsto \text{col}_i(\mathbf{G}) / \|\text{col}_i(\mathbf{G})\|_{\text{RMS}}$  for each  $i = 1, \dots, d_{\text{in}}$ .

- 1 Introduction
  - Intuition
  - Descent in Normed Spaces
- 2 Modular Norm
  - Invariant Optimal Step Size
  - Module and Modular Norm
  - Building Compound modules
- 3 Duality Map
  - Basic Modules
  - Computation
  - Concrete Examples
- 4 Summary and Takeaway
- 5 Appendix
  - Additional Modules
  - Omitted Proofs

# Invariant Optimal Step Size



**Figure 2: Learning rate transfer in the modular norm.** We train GPT with context length 128 for 10k steps on OpenWebText.

**Left:** Learning rate sweeps for normed Adam across different width and number of blocks.

**Right:** For each method, the learning rate is tuned at the scale marked by the dotted line.

# Full Weight Space Norm

The advantage of the non-standard  $\|\cdot\|_{RMS \rightarrow RMS}$  matrix norm is that it allows one to estimate the amount of feature change induced by a gradient update:

$$\|\Delta \mathbf{W} \mathbf{x}\|_{RMS} \leq \|\Delta \mathbf{W}\|_{RMS \rightarrow RMS} \cdot \|\mathbf{x}\|_{RMS}, \quad (4)$$

Now we know how to metrize individual layers, can we combine layer-wise norms to produce a norm on the full weight space  $\mathcal{W} = \prod_k \mathcal{W}_k$  of the network?

The author propose to define the norm to be useful by the criteria 6. It ends up as a max (  $L^\infty$  combination) of scaled layer-wise norms  $\|\cdot\|_{\mathcal{W}_k}$  :

$$\|(\mathbf{w}_1, \dots, \mathbf{w}_L)\|_{\mathcal{W}} := \max \left( s_1 \|\mathbf{w}_1\|_{\mathcal{W}_1}, \dots, s_L \|\mathbf{w}_L\|_{\mathcal{W}_L} \right). \quad (5)$$

The positive scalar constants  $s_1, \dots, s_L$  are determined by both the architecture of the network and a set of user-specified "mass" parameters.

# Normed optimization

We define the following operation on weight updates

$\Delta \mathbf{w} = (\Delta \mathbf{w}_1, \dots, \Delta \mathbf{w}_L) \in \mathcal{W}$  :

$$\text{normalize}(\Delta \mathbf{w}) := \left( \frac{\Delta \mathbf{w}_1}{s_1 \|\Delta \mathbf{w}_1\|_{\mathcal{W}_1}}, \dots, \frac{\Delta \mathbf{w}_L}{s_L \|\Delta \mathbf{w}_L\|_{\mathcal{W}_L}} \right) \quad (6)$$

Use `normalize` as a wrapper, along with an explicit learning rate schedule for any base optimizer:

```
delta_w = optim(w.grad())
net.normalize(delta_w)
w -= eta(step) * delta_w
```

# get update from base optimizer  
# normalize update in the modular norm  
# apply update with learning rate eta

## Definition (Module)

Given input vector space  $\mathcal{X}$ , output vector space  $\mathcal{Y}$  and weight vector space  $\mathcal{W}$ , a module  $M$  is an object with the following four attributes:

- a *function*,  $M.forward : \mathcal{W} \times \mathcal{X} \rightarrow \mathcal{Y}$ , which maps an input and a weight vector to an output. we abbreviate this attribute to just  $M \equiv M.forward$ ;
- a *number*,  $M.mass \geq 0$ , which will turn out to set the **proportion of feature learning** that this module contributes to any supermodule;
- a *number*,  $M.sensitivity \geq 0$ , which estimates the module's sensitivity to input perturbations;
- a *norm* over the weight space,  $M.norm : \mathcal{W} \rightarrow \mathbb{R}_{\geq 0}$ , sometimes abbreviated to just  $\| \cdot \|_M$ .

# Modules

We care about three kinds of module in practice:

- *atomic modules*, whose attributes are hand-declared, and have weights. (e.g. linear modules, embedding modules, and convolution modules.)
- *bond modules*, whose attributes are hand-declared, but have no weights. Formally, their weight space is the zero vector space  $\mathcal{W} = 0$ . (e.g. ReLu, Add, Abs)
- *compound modules*, built out of other modules, with automatically inferred attributes.

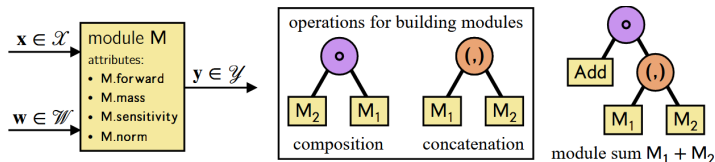


Figure 3: Modules and trees of modules

# Well-normed Module

We want a module to be good when its attributes are predictive of its behaviour.

## Definition (Well-normed)

Let  $M$  be a module on  $(\mathcal{X}, \mathcal{Y}, \mathcal{W})$ , where the input and output spaces have respective norms  $\|\cdot\|_{\mathcal{X}}$  and  $\|\cdot\|_{\mathcal{Y}}$ .  $M$  is **well-normed** if for all inputs  $\mathbf{x} \in \mathcal{X}$  and weights  $\mathbf{w} \in \mathcal{W}$ :

$$\begin{aligned} \|\nabla_{\mathbf{w}} M.\text{forward}(\mathbf{w}, \mathbf{x}) \diamond \Delta \mathbf{w}\|_{\mathcal{Y}} &\leq M.\text{norm}(\Delta \mathbf{w}) \quad \text{for all } \Delta \mathbf{w} \in \mathcal{W} \\ \|\nabla_{\mathbf{x}} M.\text{forward}(\mathbf{w}, \mathbf{x}) \diamond \Delta \mathbf{x}\|_{\mathcal{Y}} &\leq M.\text{sensitivity} * \|\Delta \mathbf{x}\|_{\mathcal{X}} \quad \text{for all } \Delta \mathbf{x} \in \mathcal{X} \end{aligned}$$

Establish Lipschitz-continuity over input and weights space.

# Building Upon Axioms

We could consider building new modules from old ones via the binary operations of **composition** and **concatenation**.

## Definition (Module composition)

Consider module  $M_1$  with input, output and weight space  $(\mathcal{X}_1, \mathcal{Y}_1, \mathcal{W}_1)$  and module  $M_2$  with input, output and weight space  $(\mathcal{X}_2, \mathcal{Y}_2, \mathcal{W}_2)$ .  $M_1$  and  $M_2$  are composable if  $\mathcal{X}_2 = \mathcal{Y}_1$ . Their composite  $M = M_2 \circ M_1$  lives on  $(\mathcal{X}_1, \mathcal{Y}_2, \mathcal{W}_1 \times \mathcal{W}_2)$  with attributes:

- $M.\text{forward}((\mathbf{w}_1, \mathbf{w}_2), \mathbf{x}) = M_2.\text{forward}(\mathbf{w}_2, M_1.\text{forward}(\mathbf{w}_1, \mathbf{x}))$ ;
- $M.\text{mass} = M_1.\text{mass} + M_2.\text{mass}$ ;
- $M.\text{sensitivity} = M_1.\text{sensitivity} * M_2.\text{sensitivity}$ ;
- $M.\text{norm}((\mathbf{w}_1, \mathbf{w}_2))$  given by:

$$\max \left( M_2.\text{sensitivity} * \frac{M.\text{mass}}{M_1.\text{mass}} * M_1.\text{norm}(\mathbf{w}_1), \frac{M.\text{mass}}{M_2.\text{mass}} * M_2.\text{norm}(\mathbf{w}_2) \right)$$

## Definition (Module concatenation)

Consider module  $M_1$  with input, output and weight space  $(\mathcal{X}_1, \mathcal{Y}_1, \mathcal{W}_1)$  and module  $M_2$  with input, output and weight space  $(\mathcal{X}_2, \mathcal{Y}_2, \mathcal{W}_2)$ . We say that  $M_1$  and  $M_2$  are concatenatable if their input spaces match:  $\mathcal{X}_1 = \mathcal{X}_2$ . The tuple  $M = (M_1, M_2)$  has input, output and weight space  $(\mathcal{X}_1, \mathcal{Y}_1 \times \mathcal{Y}_2, \mathcal{W}_1 \times \mathcal{W}_2)$  and attributes:

- $M.\text{forward}((\mathbf{w}_1, \mathbf{w}_2), \mathbf{x}) = (M_1.\text{forward}(\mathbf{w}_1, \mathbf{x}), M_2.\text{forward}(\mathbf{w}_2, \mathbf{x}))$ ;
- $M.\text{mass} = M_1.\text{mass} + M_2.\text{mass}$ ;
- $M.\text{sensitivity} = M_1.\text{sensitivity} + M_2.\text{sensitivity}$ ;
- $M.\text{norm}(\mathbf{w}_1, \mathbf{w}_2)$  given by:

$$\max \left( \frac{M.\text{mass}}{M_1.\text{mass}} * M_1.\text{norm}(\mathbf{w}_1), \frac{M.\text{mass}}{M_2.\text{mass}} * M_2.\text{norm}(\mathbf{w}_2) \right),$$

- 1 Introduction
  - Intuition
  - Descent in Normed Spaces
- 2 Modular Norm
  - Invariant Optimal Step Size
  - Module and Modular Norm
  - Building Compound modules
- 3 **Duality Map**
  - **Basic Modules**
  - **Computation**
  - **Concrete Examples**
- 4 Summary and Takeaway
- 5 Appendix
  - Additional Modules
  - Omitted Proofs

# Duality Maps Construction

We could first write down duality maps for atomic modules, then extend them to arbitrary compound modules by showing how these maps transfer. Typical procedure:

- 1 Select norms on the input and output spaces that respect the semantics<sup>2</sup> of  $M$ .forward.
- 2 Place a norm on the weight space to make  $M$  well-normed. Typically the induced operator norm.
- 3 Solve the duality maps for atomic modules, then extend to the whole network.

---

<sup>2</sup>Select norms that describe both how large we would like the inputs and outputs to be, and in what geometry we would like the outputs to evolve.

# Duality Maps for Atomic Modules

## Example (The Linear module)

The Linear module sends inputs from  $\mathcal{X} = \mathbb{R}^{d_{in}}$  to outputs in  $\mathcal{Y} = \mathbb{R}^{d_{out}}$ . The weight space is given by the matrix space  $\mathcal{W} = \mathbb{R}^{d_{out} \times d_{in}}$ . We endow the Linear module with attributes:

- Linear.forward ( $\mathbf{W}, \mathbf{x}$ ) =  $\mathbf{W}\mathbf{x}$ , the matrix-vector product;
- Linear.sensitivity = 1;
- Linear.mass =  $\mu$ , where  $\mu \geq 0$  is a hyperparameter;
- Linear.norm ( $\mathbf{W}$ ) =  $\|\mathbf{W}\|_{\text{RMS} \rightarrow \text{RMS}}$ , the RMS  $\rightarrow$  RMS induced operator norm.

The duality map corresponding to Linear.norm is then given by:

Linear.dualize ( $\mathbf{G}$ ) =  $\sqrt{\frac{d_{out}}{d_{in}}} \times \mathbf{U}\mathbf{V}^T$ , where the gradient  $\mathbf{G} \in \mathbb{R}^{d_{out} \times d_{in}}$  has reduced SVD  $\mathbf{G} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ .

Well-normed if inputs and weights belong to the unit balls.

# Duality Maps for Atomic Modules

## Example (The Embed module)

The Embed module sends inputs from  $\mathcal{X} = \mathbb{R}^{d_{\text{in}}}$  to outputs in  $\mathcal{Y} = \mathbb{R}^{d_{\text{out}}}$ . The weight space is given by the matrix space  $\mathcal{W} = \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ . We endow the Embed module with attributes:

- Embed.forward ( $\mathbf{W}, \mathbf{x}$ ) =  $\mathbf{W}\mathbf{x}$ , the matrix-vector product;
- Embed.sensitivity = 1;
- Embed.mass =  $\mu$ , where  $\mu \geq 0$  is a hyperparameter;
- Embed.norm ( $\mathbf{W}$ ) =  $\|\mathbf{W}\|_{\ell_1 \rightarrow \text{RMS}}$ , the  $\ell_1 \rightarrow \text{RMS}$  induced operator norm.

The duality map for Embed.norm is:

Embed.dualize ( $\mathbf{G}$ ) performs the mapping  $\text{col}_j(\mathbf{G}) \mapsto \frac{\text{col}_j(\mathbf{G})}{\|\text{col}_j(\mathbf{G})\|_{\text{RMS}}}$  for each column index  $j = 1, \dots, d_{\text{in}}$ .

Well-normed if inputs and weights belong to the unit balls.

# Duality Maps for Atomic Modules

Finally, we consider a Conv2D module with a  $k \times k$  kernel. The calculations work by slicing up the weight tensor into a collection of  $k^2$  matrices.

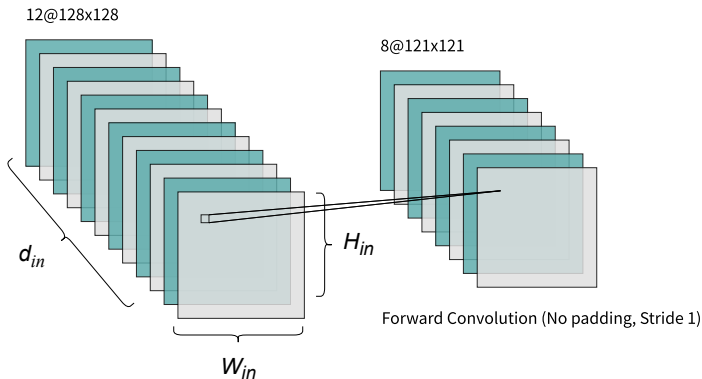


Figure 4: Illustration of forward passing in Conv2D

# Duality Maps for Atomic Modules

## Example (The Conv2D module)

The Conv2D module sends inputs from  $\mathcal{X} = \mathbb{R}^{W_{\text{in}} \times H_{\text{in}} \times d_{\text{in}}}$  to outputs in  $\mathcal{Y} = \mathbb{R}^{W_{\text{out}} \times H_{\text{out}} \times d_{\text{out}}}$ . We think of this as mapping an input image of width  $W_{\text{in}}$ , height  $H_{\text{in}}$  and with  $d_{\text{in}}$  color channels to an output image of width  $W_{\text{out}}$ , height  $H_{\text{out}}$  and with  $d_{\text{out}}$  color channels. The weight space is given by the tensor space  $\mathcal{W} = \mathbb{R}^{d_{\text{out}} \times d_{\text{in}} \times k \times k}$ , where  $k$  is the kernel size. We endow Conv2D with attributes:

- Conv2D.forward ( $\mathbf{W}, \mathbf{x}$ ) =  $\mathbf{W} \circledast \mathbf{x}$ , where  $\circledast$  denotes 2D convolution;
- Conv2D.sensitivity = 1;
- Conv2D.mass =  $\mu$ , where  $\mu \geq 0$  is a hyperparameter;
- Conv2D.norm ( $\mathbf{W}$ ) =  $k^2 \max_{i,j=1}^k \|\mathbf{W}_{..ij}\|_{\text{RMS} \rightarrow \text{RMS}}$ , the max RMS  $\rightarrow$  RMS norm over kernel indices.

# Duality Maps for Atomic Modules

- We would like pixel intensities in the inputs and outputs to be order one and undergo order one change.
- We formalize this by taking the input and output norms to be the spatial maximum of the RMS norms of all the color channel vectors:  
$$\|\mathbf{x}\|_{\mathcal{X}} = \max_{w=1}^{W_{\text{in}}} \max_{h=1}^{H_{\text{in}}} \|\mathbf{x}_{wh\cdot}\|_{\text{RMS}} \text{ and}$$
$$\|\mathbf{y}\|_{\mathcal{Y}} = \max_{w=1}^{W_{\text{out}}} \max_{h=1}^{H_{\text{out}}} \|\mathbf{y}_{wh\cdot}\|_{\text{RMS}}.$$
- Since the duality map for a max of norms decouples into one duality map per sub-norm, the duality map corresponding to Conv2D.norm is given by:  $\mathbf{G} \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}} \times k \times k}$   
Conv2D.dualize ( $\mathbf{G}$ ) does  $\mathbf{G}_{..ij} \mapsto \frac{1}{k^2} \sqrt{\frac{d_{\text{out}}}{d_{\text{in}}}} \times \mathbf{U}_{ij} \mathbf{V}_{ij}^{\top}$ , where  $\mathbf{G}_{..ij}$  has reduced SVD  $\mathbf{U}_{ij} \mathbf{\Sigma}_{ij} \mathbf{V}_{ij}^{\top}$ .

# Duality Maps for Bond Modules and Compound Modules

- Since the weight space of Bound Module B is  $\{0\}$ , then  
 $B.norm = 0 \mapsto 0$ ,  $B.dualize = 0 \mapsto 0$
- For compound modules,
  - Composition:

M.dualize ( $\mathbf{g}_1, \mathbf{g}_2$ ) =

$$\left( \frac{1}{M_2.sensitivity} * \frac{M_1.mass}{M.mass} * M_1.dualize(\mathbf{g}_1), \frac{M_2.mass}{M.mass} * M_2.dualize(\mathbf{g}_2) \right) \quad (7)$$

- Concatenation:

M.dualize ( $\mathbf{g}_1, \mathbf{g}_2$ ) =

$$\left( \frac{M_1.mass}{M.mass} * M_1.dualize(\mathbf{g}_1), \frac{M_2.mass}{M.mass} * M_2.dualize(\mathbf{g}_2) \right) \quad (8)$$

# Rectangular Newton-Schulz Iteration

According to [Higham, 2008], the method works by first normalizing the matrix  $\mathbf{G}$  according to  $\mathbf{X}_0 = \mathbf{G} / \|\mathbf{G}\|_{\ell_2 \rightarrow \ell_2}$  and then iterating:

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

then as  $t \rightarrow \infty$ , the sequence  $\mathbf{X}_t \rightarrow \mathbf{UV}^\top$ .

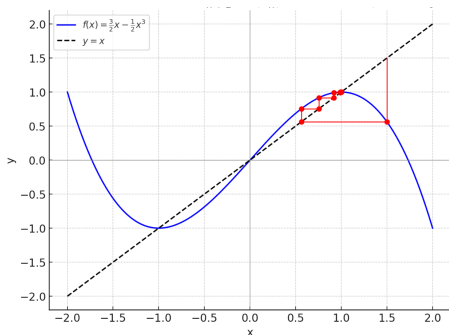


Figure 5: 1d illustration of iteration ( $x_0 = 1.5$ )

## Definition (Blocks and deep networks)

A deep neural network is a module  $M$  formed by a composition

$$M = \text{OutputLayer} \circ \text{Block}_L \circ \dots \circ \text{Block}_1 \circ \text{InputLayer}$$

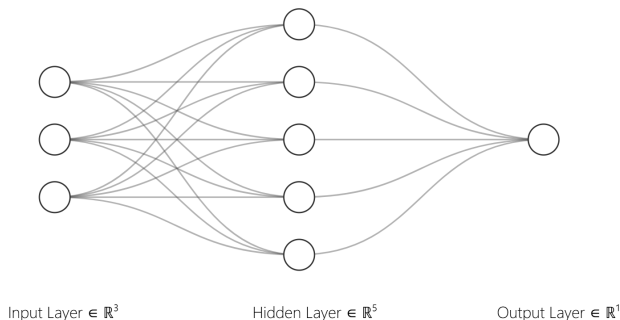
where  $\text{InputLayer}$ ,  $\text{Block}_1, \dots, \text{Block}_L$ ,  $\text{OutputLayer}$  are modules.

Typically, each of  $\text{Block}_1, \dots, \text{Block}_L$  will be copies of the same module, so that the network can be written as an iterated composition

$$M = \text{OutputLayer} \circ \text{Block}^L \circ \text{InputLayer}.$$

Usually  $\text{InputLayer}$  is often some form of embedding module, and  $\text{OutputLayer}$  is usually a linear module.

# Example: Two-layer network



$$M := \text{Linear}(1,5) \circ \left( \text{ReLU} \circ \text{Linear}(5,3) \circ \text{LayerNorm} \right) \circ \text{Embed}$$

We show *bond module* in cyan to mark no dualization.

# Example: Multi-layer network

We denote  $M(d) = \text{ReLU} \circ \text{Linear}(d, d) \circ \text{LayerNorm}$  for intermediate layers.

Then  $M = \text{OutputLayer} \circ M(d)^L \circ \text{Embed}$

- The *bound module* still take effects in its sensitivity and specified mass. (e.g. ReLu layer has sensitivity  $1/\sqrt{2}$ )

# Example: Multi-head Attention

## Definition (Module broadcasting)

Suppose  $M$  is a module with inputs  $\mathcal{X}$ , outputs  $\mathcal{Y}$  and weights  $\mathcal{W}$ . Then for any  $h \geq 1$ , the  $h$ -times-broadcast of  $M$  is the module  $M^{(h)}$  with the same weight space  $\mathcal{W}$ , mass, sensitivity and norm as  $M$ , but inputs the Cartesian power  $\mathcal{X}^h = \mathcal{X} \times \dots \times \mathcal{X}$  and outputs  $\mathcal{Y}^h = \mathcal{Y} \times \dots \times \mathcal{Y}$ , and forward function

$$(\mathbf{w}, (\mathbf{x}_1, \dots, \mathbf{x}_h)) \mapsto (M.\text{forward}(\mathbf{w}, \mathbf{x}_1), \dots, M.\text{forward}(\mathbf{w}, \mathbf{x}_h)).$$

The attention module have both inputs and outputs.  $\mathcal{X} = \mathbb{R}^{\ell \times d}$  where  $\ell$  is the context length and  $d$  is the embedding dimension. The attention module itself will depend on three additional dimensional arguments:

- $h$ , the number of heads;
- $d_Q$ , the key/query dimension;
- $d_V$ , the value dimension;

## Example: Multi-head Attention

Also an  $\ell \times \ell$  matrix mask, which we usually take to be either

$$\text{mask}_{ij} = \begin{cases} 0 & \text{if } i \geq j \\ -\infty & \text{otherwise} \end{cases}$$

for causal attention, and  $\text{mask} = 0$  for non-causal attention.

### Definition (FuncAttention Module)

Inputs  $\mathcal{X} = \mathbb{R}^{\ell \times d_Q} \times \mathbb{R}^{\ell \times d_Q} \times \mathbb{R}^{\ell \times d_V}$ , outputs  $\mathcal{Y} = \mathbb{R}^{\ell \times d_V}$ , and forward function

$$\text{FuncAttention.forward}(\mathbf{q}, \mathbf{k}, \mathbf{v}) = \text{softmax}\left(\frac{\mathbf{q}\mathbf{k}^T}{d_Q} + \text{mask}\right)\mathbf{v}.$$

# Example: Multi-head Attention

For positive integers  $\ell, d, h, d_Q, d_V$  and a choice of mask, take four instances of the linear module:

- Query = Linear ( $h * d_Q, d$ )
- Key = Linear ( $h * d_Q, d$ )
- Value = Linear ( $h * d_V, d$ )
- Exit = Linear ( $d, h * d_V$ )

which by broadcasting we consider to have inputs of shape  $\mathbb{R}^{\ell \times d}$ . The multi-headed attention MultiHeadAttention module is then the composition:

$$\text{MultiHeadAttention} = \text{Exit} \circ \frac{1}{3} * \text{FuncAttention}^{(h)} \circ (\text{Query}, \text{Key}, \text{Value})$$

- 1 Introduction
  - Intuition
  - Descent in Normed Spaces
- 2 Modular Norm
  - Invariant Optimal Step Size
  - Module and Modular Norm
  - Building Compound modules
- 3 Duality Map
  - Basic Modules
  - Computation
  - Concrete Examples
- 4 Summary and Takeaway
- 5 Appendix
  - Additional Modules
  - Omitted Proofs

# Summary and Takeaway

- The modular norm is proposed to let optimal learning rate transfer when scaling up width and depth using different "normed" base optimizers.
- We could build modules and duality maps for whole network through composition and concatenation of atomic and bond modules.
- We could compute the duality map (essentially the SVD) with Newton-Schulz Iteration.

Module	Weight Space $\mathcal{W}$	Module.norm	Module.dualize
Linear	$\mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$	$\mathbf{W} \mapsto \ \mathbf{W}\ _{\text{RMS} \rightarrow \text{RMS}}$	$\mathbf{G} \mapsto \sqrt{\frac{d_{\text{out}}}{d_{\text{in}}}} \times \mathbf{UV}^\top$
Embed	$\mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$	$\mathbf{W} \mapsto \ \mathbf{W}\ _{\ell_1 \rightarrow \text{RMS}}$	$\text{col}_j(\mathbf{G}) \mapsto \frac{\text{col}_j(\mathbf{G})}{\ \text{col}_j(\mathbf{G})\ _{\text{RMS}}}$
Conv2D	$\mathbb{R}^{d_{\text{out}} \times d_{\text{in}} \times k \times k}$	$\mathbf{W} \mapsto k^2 \max_{i,j=1}^k \ \mathbf{W}_{\cdot ij}\ _{\text{RMS} \rightarrow \text{RMS}}$	$\mathbf{G}_{\cdot ij} \mapsto \frac{1}{k^2} \sqrt{\frac{d_{\text{out}}}{d_{\text{in}}}} \times \mathbf{U}_{ij} \mathbf{V}_{ij}^\top$

Figure 6: Summary of atomic duality maps

- 1 Introduction
  - Intuition
  - Descent in Normed Spaces
- 2 Modular Norm
  - Invariant Optimal Step Size
  - Module and Modular Norm
  - Building Compound modules
- 3 Duality Map
  - Basic Modules
  - Computation
  - Concrete Examples
- 4 Summary and Takeaway
- 5 Appendix
  - Additional Modules
  - Omitted Proofs

## Example (Bond module ReLU)

In any dimension  $d$ , we define the "rectified linear unit" bond module ReLU to have input space  $\mathcal{X} \subset \mathbb{R}^d$ , output space  $\mathcal{Y} = \mathbb{R}^d$ , forward function

ReLU.forward :  $(x_1, \dots, x_d) \mapsto (\max(0, x_i))_{i=1, \dots, d}$ . and sensitivity  $1/\sqrt{2}$ .

## Example (LayerNorm)

We define

MeanSubtract.forward :  $(x_1, \dots, x_d) \mapsto (x_1 - \bar{x}, \dots, x_d - \bar{x})$

and

RMSDivide.forward :  $\mathbf{x} \mapsto \frac{\mathbf{x}}{\|\mathbf{x}\|_{\text{RMS}}} = \frac{\sqrt{d}\mathbf{x}}{\|\mathbf{x}\|_2}$ .

Such that: LayerNorm = RMSDivide  $\circ$  MeanSubtract.

# Duality Map for Composition and Concatenation

We prove for (8), the proof for (7) is similar by substituting the norm.


Denote  $(\mathbf{q}_1, \mathbf{q}_2) = M.dualize(\mathbf{g}_1, \mathbf{g}_2)$ .

Then  $(q_1, q_2) = \underset{M.norm(q_1, q_2)=1}{argmax} q_1 g_1 + q_2 g_2$

According to the norm constraints, we could assume WLOG that

$\frac{M.mass}{M_1.mass} \times M_1.norm(q_1) = 1$ ,  $\frac{M.mass}{M_2.mass} \times M_2.norm(q_2) \leq 1$ , then  
optimizing the  $q_1 g_1$  term yields  $q_1 \sim \alpha M_1.dualize(\mathbf{g}_1)$ , similarly for the  
second term  $q_2 \sim \beta M_2.dualize(\mathbf{g}_2)$

We then scale up the size of two terms to satisfy the norm constraints.

-  Higham, N. J. (2008).  
*Functions of matrices: theory and computation.*  
SIAM.