

Lecture 5

- Local linear perspective (Recap)
- Sign SGD
- Gradient Descent
- Shampoo
- Towards μ^p and μ_{non} .

Why do we care about optimizers?

- Want to train fast. Training time = \$\$\$.

- AdamW + tuned hyperparameters are default.

- Hyperparameter search is difficult for new optimizers.

But under Adam we see lazy training \rightarrow i.e. models barely change from their random initialization

\rightarrow Kernel / feature learning perspective \rightarrow NTK

\rightarrow There is no way random initialization should work
 \rightarrow we are not using all the power of these networks...

\rightarrow Initialization rules are nice.
 \rightarrow Before this everything was $N(0,1)$

$\vec{\theta}$: parameters.

$l(\vec{x}, \vec{\theta})$ is the loss.

\vec{x} : data.

Linearized perspective:

$$l(\vec{x}, \vec{\theta}_i + \Delta \vec{\theta}) \approx l(\vec{x}, \vec{\theta}_i) + \left. \frac{\partial l}{\partial \vec{\theta}} \right|_{\vec{\theta} = \vec{\theta}_i} \Delta \vec{\theta} + \dots$$

Gradient Descent takes a step.

→ Want this to be small enough that the approximation still holds

→ Want to converge fast

Follow perspective from Bernstein + Newhouse 2024.

↳ Build on many other works
average over data points.

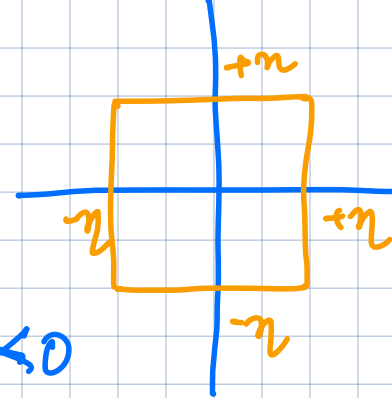
$$L(\vec{\theta}) = \frac{1}{n} \sum_{j=1}^n l(\vec{x}_j, \vec{\theta})$$

Want:

$$\underset{\|\Delta \vec{\theta}\| \leq \eta}{\operatorname{argmin}} \left(L(\vec{\theta}) + \langle \nabla_{\vec{\theta}} L(\vec{\theta}), \Delta \vec{\theta} \rangle \right)$$

$$= \underset{\|\Delta \vec{\theta}\| \leq \eta}{\operatorname{argmin}} \langle \nabla_{\vec{\theta}} L(\vec{\theta}), \Delta \vec{\theta} \rangle.$$

Choice 1: $\|\vec{\Delta\theta}\|_\infty \leq \eta \rightarrow$ Linear program



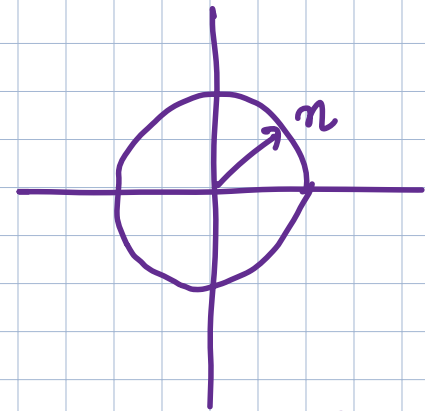
Solution:

$$\Delta\theta[j] = \begin{cases} +\eta & \text{if } \nabla_{\theta} \mathcal{L}(\vec{\theta})[j] \leq 0 \\ -\eta & \text{if } \nabla_{\theta} \mathcal{L}(\vec{\theta})[j] > 0 \end{cases}$$

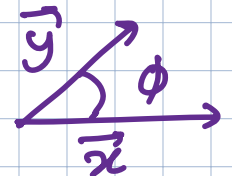
$$\Rightarrow \vec{\Delta\theta} = \underbrace{-\eta \cdot \text{sgn}(\nabla_{\theta} \mathcal{L}(\vec{\theta}))}_{\text{Sign SGD}}$$

$$\begin{cases} \underset{\|\vec{\Delta\theta}\|_\infty \leq \eta}{\text{argmin}} & \langle \nabla_{\theta} \mathcal{L}(\vec{\theta}), \vec{\Delta\theta} \rangle \\ \underset{\|\vec{y}\|_\infty \leq \eta}{\text{argmin}} & \vec{x}^T \vec{y} \end{cases}$$

Choice 2: $\|\vec{\Delta\theta}\|_2 \leq \eta$
 $\underset{\|\vec{\Delta\theta}\|_2 \leq \eta}{\text{argmin}} \quad \langle \nabla_{\theta} \mathcal{L}(\vec{\theta}), \vec{\Delta\theta} \rangle$



$$\vec{x}^T \vec{y} = \|\vec{x}\|_2 \|\vec{y}\|_2 \cos \phi$$



Cauchy-Schwartz: Inner product absolute value maximized when vectors are aligned.

$$\Rightarrow \vec{\Delta\theta} = -\eta \underbrace{\frac{\nabla_{\theta} \mathcal{L}(\vec{\theta})}{\|\nabla_{\theta} \mathcal{L}(\vec{\theta})\|_2}}_{\text{unit norm.}}$$

scale by η .

Alternatively, take the regularization / Lagrange multiplier perspective.

argmin.
 $\vec{\theta}$
 unconstrained.

$$\underbrace{\langle \nabla_{\theta} \mathcal{L}(\vec{\theta}), \Delta\theta \rangle + \lambda \|\Delta\theta\|_2^2}_{g(\Delta\theta). \text{ convex.}}$$

→ Sweeps out the same family of solutions as constrained version.

$$\nabla g(\Delta\theta) = \nabla_{\theta} \mathcal{L}(\vec{\theta}) + 2\lambda \vec{\Delta\theta} \rightarrow \text{Set to zero.}$$

$$\text{Minimize @ } \vec{\Delta\theta}^* = -\underbrace{\frac{1}{2\lambda}}_{\text{stepsize}} \underbrace{\nabla_{\theta} \mathcal{L}(\vec{\theta})}_{\text{gradient descent}}$$

Linearization around matrices?

$\vec{\theta}$: we've been thinking of this as a vector.

But neural network parameters are naturally organized in matrices.

$$\vec{h}_1 = \text{ReLU}[\vec{b}_1 + \boxed{W_0 \vec{x}}]$$

$$\vec{h}_j = \text{ReLU}[\vec{b}_j + W_{j-1} \cdot \vec{h}_{j-1}]$$

$$y = \text{ReLU}[\vec{b}_k + W_k \vec{h}_k]$$

Recall from Xavier initialization ... scaling problems come from the weight matrices.

So what if we think about change ΔW in "weight matrix space" instead.

$$\begin{array}{l} \text{argmin.} \\ \|\Delta W\| \leq \eta \end{array} \quad \underbrace{\langle \nabla_w \mathcal{L}(w), \Delta W \rangle}_{\text{matrix} \quad \text{matrix.}}$$

Recall: $\langle A, B \rangle_F = \text{trace}(A^T B)$

Consider: $\|\Delta W\|_2 \leq \eta$ ie. spectral norm.

Recall: $\|A\|_2 = \sigma_{\max}$ largest singular value of A .

$\max_{\|\vec{x}\|_2=1} \|A\vec{x}\|_2 \rightarrow$ what is the maximum scaling for a unit vector.

argmin
 $\|\Delta W\|_2 \leq \eta$ $\langle \nabla_w \mathcal{L}(w), \Delta W \rangle$.

Let $\nabla_w \mathcal{L}(w) = U \Sigma V^T$.

Then $\Delta W^* = -\eta U V^T$.

This gives us the Shampoo Optimizer (Gupta 2017/2018)

\rightarrow A variant of Shampoo (Dahl 2023) won AlgoPerf
 \rightarrow Training alg. competition.

$$\max_{\|B\|_2 \leq 1} \text{trace}(A^T B)$$

$$\text{trace}(ABC) \\ = \text{trace}(BCA)$$

$$\text{Let } A = U \Sigma V^T \\ = \sum_i \sigma_i \vec{u}_i \vec{v}_i^T$$

$$A = U \Sigma V^T$$

$$\max \text{trace}(A^T B)$$

$$\|B\|_2 \leq 1$$

$$= \max_{\|B\|_2 \leq 1} \text{trace}(V \Sigma^T U^T B)$$

$$\|B\|_2 \leq 1$$

$$= \max_{\|B\|_2 \leq 1} \text{trace} \left(\sum_i \sigma_i \vec{v}_i \vec{u}_i^T B \right) \quad \text{Trace is linear.}$$

$$= \max_{\|B\|_2 \leq 1} \sum_i \text{trace}(\sigma_i \vec{v}_i \vec{u}_i^T B)$$

$$= \max_{\|B\|_2 \leq 1} \sum_i \text{trace} \left(\sigma_i \underbrace{\vec{u}_i^T B \vec{v}_i}_{\text{Scalar}} \right)$$

$$= \max_{\|B\|_2 \leq 1} \sum_i \sigma_i \underbrace{\vec{u}_i^T B \vec{v}_i}_{\text{Scalar}}$$

$$\|B \vec{v}_i\|_2 \leq \|\vec{v}_i\|_2 \cdot \|B\|_2 \leq 1.$$

$$\|\vec{u}_i^T B \vec{v}_i\|_2 \leq 1. \quad \text{Cauchy-Sch.}$$

$$\leq \sum_i \sigma_i$$

→ To attain bound $B = U V^T$

$$\begin{aligned} & \text{trace}(V \Sigma^T U^T U V^T) \\ &= \text{trace}(\Sigma^T) \end{aligned}$$

Beinstein + Newhouse 2024

- Choose a norm
- Choose a stepsize.

→ The right norm may depend on geometry / architecture of network.

- But we have a recipe

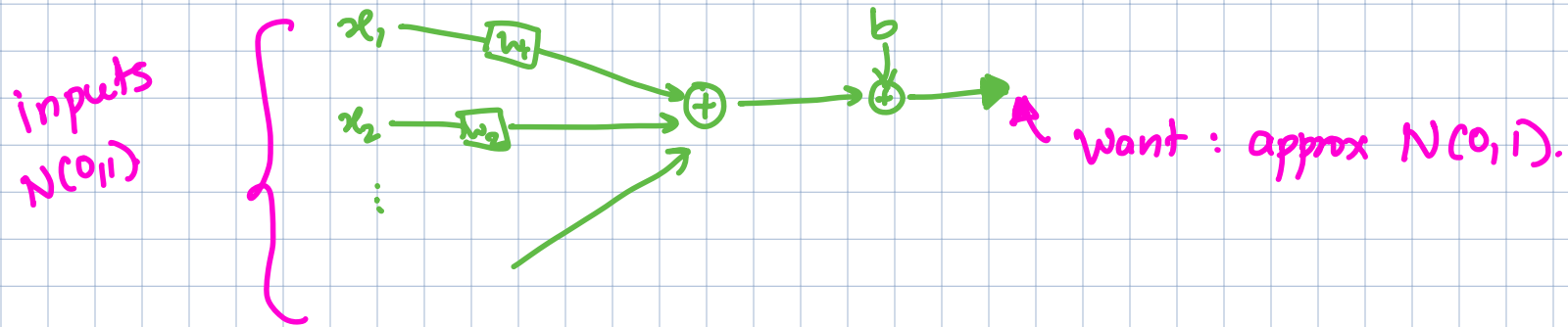
→ Given a norm, make an optimizer.

→ examples for $\|\cdot\|_2$, $\|\cdot\|_0$, spectral norm.

→ Why do we care about this recipe?

So what might be a good norm for deep learning?

Recap Xavier initialization.



Can we rewrite this in terms of norms?

What norm are we trying to preserve?

Say $\|\vec{x}\|_2 = 1$ Xavier does not imply $\|\vec{h}\|_2 = 1$.

$$x_i \sim \frac{1}{\sqrt{d_{in}}}$$

$$h_i \sim \frac{1}{\sqrt{d_{out}}}$$

Consider RMS-norm:

$$\|\vec{x}\|_{\text{RMS}} = \left(\frac{1}{\sqrt{d_{in}}} \right) \|\vec{x}\|_2 = \frac{1}{\sqrt{d_{in}}} \sqrt{\sum_i x_i^2} \rightarrow \text{Average size of each entry 1!}$$

Root mean square.

Hmm... So maybe RMS norm is interesting for DNNs.
Is there a matrix version? (Recall Shampoo...)

Recall: induced Matrix-norm

$$\|A\|_{\alpha \rightarrow \beta} = \max_{\|\vec{x}\|_{\alpha}=1} \|A\vec{x}\|_{\beta}$$

RMS \rightarrow RMS norm?

$$\begin{aligned} \|A\|_{\text{RMS} \rightarrow \text{RMS}} &= \max_{\|\vec{x}\|_{\text{RMS}}=1} \|A\vec{x}\|_{\text{RMS}} \\ A \quad \begin{array}{|c|} \hline d_1 \\ \hline \end{array} \quad \begin{array}{|c|} \hline d_2 \\ \hline \end{array} &= \max_{\|\vec{x}\|_2 = \sqrt{d_1}} \frac{1}{\sqrt{d_2}} \cdot \|A\vec{x}\|_2 \end{aligned}$$

$$= \frac{\sqrt{d_1}}{\sqrt{d_2}} \cdot \max_{\|\vec{x}\|_2=1} \|A\vec{x}\|_2$$

Spectral norm!

So can we use this induced RMS \rightarrow RMS norm in our optimizer recipe?

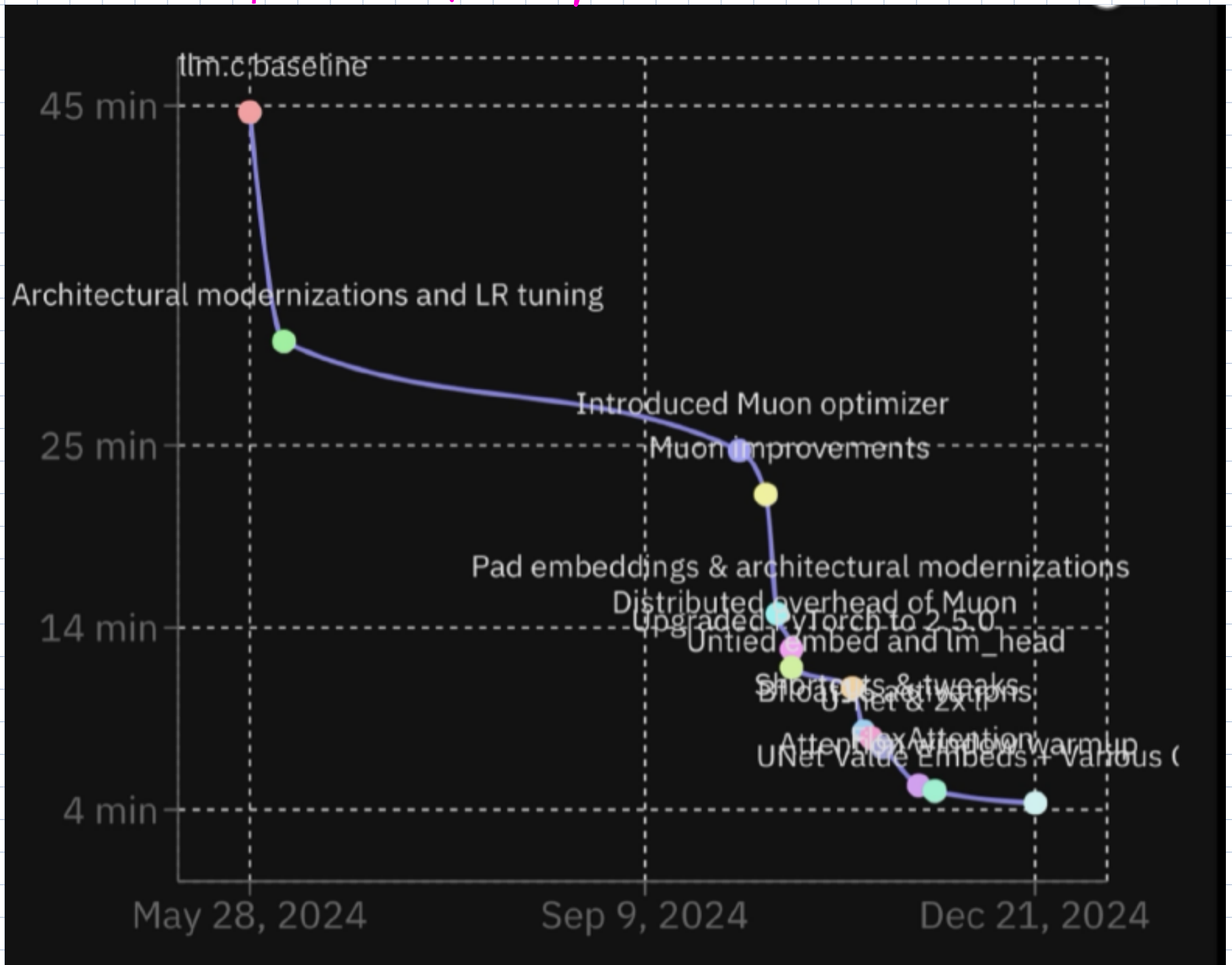
$$\|\Delta W\|_{\text{RMS}^*} \leq \eta?$$

Key observation: If we choose this norm, we can choose only one hyperparameter η across all layers, but the effect will be layer specific learning rates, due to fan-in and fan-out.

$$\begin{aligned}\|\Delta W\|_{\text{RMS}^*} \leq \eta &\Rightarrow \sqrt{\frac{d_1}{d_2}} \|\Delta W\|_2 \leq \eta. \\ &\Rightarrow \|\Delta W\|_2 \leq \eta \sqrt{\frac{d_1}{d_2}}\end{aligned}$$

\rightarrow Essential idea behind mup \rightarrow maximal update parameterization.

Nano GPT speedrun.



MUON IS SCALABLE FOR LLM TRAINING

TECHNICAL REPORT

Jingyuan Liu¹ Jianlin Su¹ Xingcheng Yao² Zhejun Jiang¹ Guokun Lai¹ Yulun Du¹
Yidao Qin¹ Weixin Xu¹ Enzhe Lu¹ Junjie Yan¹ Yanru Chen¹ Huabin Zheng¹
Yibo Liu¹ Shaowei Liu¹ Bohong Yin¹ Weiran He¹ Han Zhu¹ Yuzhi Wang¹
Jianzhou Wang¹ Mengnan Dong¹ Zheng Zhang¹ Yongsheng Kang¹ Hao Zhang¹
Xinran Xu¹ Yutao Zhang¹ Yuxin Wu¹ Xinyu Zhou¹ Zhilin Yang¹

¹ Moonshot AI ² UCLA



22 February 2025