# Lecture 7

- Muon .

# Recall:

Optimizer recepie.

$$\underset{\|\Delta w\| \leq \eta}{\text{argmin}} \quad \langle \nabla_w \mathcal{L}(w), \Delta w \rangle.$$

- $\nabla_w \mathcal{L}(w)$ — Gradient of Loss
- $\Delta w$ — change in $w$.
- $\|\Delta w\| \leq \eta$ — Choose appropriate norm.

Choose spectral norm.

$$\underset{\|\Delta w\|_2 \leq \eta}{\text{argmin}} \quad \langle \nabla_w \mathcal{L}(w), \Delta w \rangle = \eta \cdot U_r V_r^T$$

$$\text{if } \nabla_w \mathcal{L}(w) = U \Sigma V^T = U_r \Sigma_r V_r^T$$

$\Rightarrow$ Semi-orthogonal update

$$W_{t+1} = W_t - \eta \cdot U_r V_r^T$$

[ Key update of Shampoo
(without accumulation) ]

**Shampoo:**

$$L_t = L_{t-1} + G_t G_t^T$$

$$R_t = R_{t-1} + G_t^T G_t$$

$$\nabla_{w_t} \mathcal{L}(w_t) = G_t.$$

$$W_{t+1} = W_t - \eta \, L_t^{-1/4} G_t R_t^{-1/4}.$$

**Shampoo without accumulation:**

$$W_{t+1} = W_t - \eta \cdot (G_t G_t^T)^{-1/4} G_t (G_t^T G_t)^{-1/4}.$$

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

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

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

$$G_t \in \mathbb{R}^{d_{out} \times d_{in}} \quad \Rightarrow \quad \Sigma = \left[ \begin{array}{c|c} \Sigma_r & 0 \end{array} \right] \qquad \text{Assume full row rank.}$$

$$\Rightarrow \quad \Sigma\Sigma^T \underset{d_{out}}{=} \boxed{\begin{array}{c|c} \Sigma_r & 0 \end{array}} \; \boxed{\begin{array}{c} \overset{d_{out}}{\Sigma_r^T} \\ \hline \end{array}} = \Sigma_r \Sigma_r^T = \boxed{\begin{matrix} \sigma_1^2 & \\ & \sigma_2^2 \\ & & \ddots \end{matrix}}$$

$$\left(G_t G_t^T\right)^{-\frac{1}{4}} = U \boxed{\begin{matrix} \sigma_1^{-\frac{1}{2}} & \\ & \ddots \end{matrix}} U^T$$

$$\left(G_t^T G_t\right)^{-\frac{1}{4}} = V \boxed{\begin{matrix} \sigma_1^{-\frac{1}{2}} & \\ & \ddots \\ & & 0 \\ & & & 0 \\ & & & & 0 \end{matrix}} V^T$$

$$\left(G_t G_t^T\right)^{-\frac{1}{4}} G_t \left(G_t^T G_t\right)^{-\frac{1}{4}} = U \boxed{\begin{matrix} \sigma_1^{-\frac{1}{2}} & \\ & \ddots \end{matrix}} U^T U \Sigma V^T V \boxed{\begin{matrix} \sigma_1^{-\frac{1}{2}} & \\ & \ddots \end{matrix}} V^T$$

$$= U \boxed{\Sigma_r^{-\frac{1}{2}}} \boxed{\begin{array}{c|c} \Sigma_r & 0 \end{array}} \boxed{\begin{array}{c} \Sigma_r^{-\frac{1}{2}} \\ \hline 0 \end{array}} V^T = U_r V_r^T$$

Semi-orthogonal matrices:
$$\{ A \in \mathbb{R}^{m \times n} : AA^T = I_m \text{ or } A^TA = I_n \}$$

What if we had the RMS-RMS norm instead?     $W \in \mathbb{R}^{d_{out} \times d_{in}}$

$$\underset{\|\Delta W\|_{RMS \to RMS} \leq \eta}{\text{argmin}} \quad \langle \nabla_W \mathcal{L}(w), \Delta w \rangle$$

$$= \underset{\|\Delta W\|_2 \leq \sqrt{\frac{d_{out}}{d_{in}}} \cdot \eta}{\text{argmin}} \quad \langle \nabla_W \mathcal{L}(w), \Delta w \rangle$$

Muon key idea 1.

$$\Rightarrow \quad \Delta W^+ = -\eta \cdot \sqrt{\frac{d_{out}}{d_{in}}} \; U_r V_r^T$$

Why might semi-orthogonal matrices be a good idea?

$\longrightarrow$ Condition number $U_r V_r^T = 1.$     Uniform step in all directions.

$\longrightarrow$ Not dominated by largest singular values.

**But** computing $UV^T$ is expensive.

$\longrightarrow$ Needs SVD.

Muon solves this issue through two observations

(1) Getting the direction approximately correct is good enough.

(2) Newton-Schulz iterations. $\longleftarrow$ Muon key idea 2.

$$U\Sigma V^T \longrightarrow UV^T.$$

i.e. replace all the singular values by 1.

How?

$$f(U\Sigma V^T) \approx UV^T$$

# Newton- Schulz.

① Odd polynomials commute with the SVD.

$$p(x) = a_0 X + a_1 X X^T X + a_2 (X X^T)^2 X + \cdots + a_n (X X^T)^n \cdot X.$$

e.g. $\underline{p(x)} = \frac{3}{2} \cdot X - \frac{1}{2} X X^T X.$

$$p(U \Sigma Y^T) = \frac{3}{2} U \Sigma V^T - \frac{1}{2} (U \Sigma V^T)(V \Sigma^T U^T) \cdot U \Sigma V^T$$

$$= \frac{3}{2} U \Sigma V^T - \frac{1}{2} U \Sigma \Sigma^T \Sigma V^T$$

$$= U \left[ \frac{3}{2} \Sigma - \frac{1}{2} \Sigma \Sigma^T \Sigma \right] V^T$$

$$= U \, p(\Sigma) \cdot V^T$$

So $p(x) = U \underline{p(\Sigma)} V^T.$

So can apply to a matrix without changing the singular vectors,

② Can we find $p(x)$ such that

$$p(x) \to 1 \quad \text{for } x > 0$$
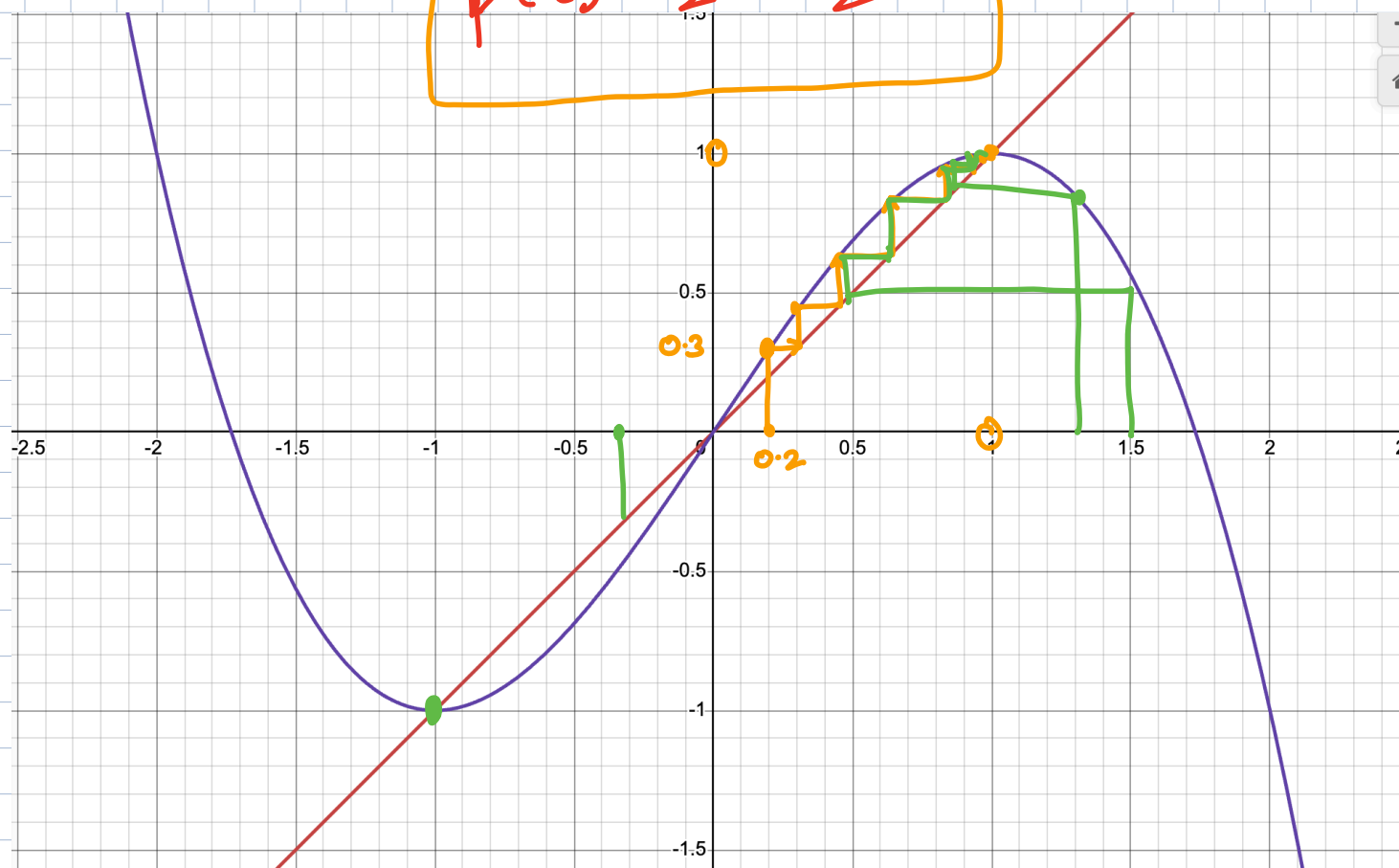
$$P(X) \to \bigcup P(\Sigma) \vee T$$
$$P(P(X)) \to \bigcup P(P(\Sigma)) \vee T$$
$$\vdots$$

Then, iteratively apply $p$.

$$p(x) = \frac{3}{2}x - \frac{1}{2}x^3$$



$$P(P(P(P(0)))) \to \quad p^n(0) \to 1$$
$$n \to \infty$$

$$D(x) = \frac{3}{2}x - \frac{1}{2}x^2 \qquad \text{Zoomed out}$$



→ Need to have singular values $\in (0,1]$ before applying!

→ How? Normalize by the Frobenius norm : $\dfrac{G}{\|G\|_F}$.

# Muon  (_Momentum _Orthogonalized by _Newton Schulz)

$$B_t = \mu \cdot B_{t-1} + \nabla_w \mathcal{L}(w)$$

$$O_t = NewtonSchulz(B_t)$$

$$W_t = W_{t-1} - \eta \cdot O_t.$$

How to choose $p(x)$?

$$p(x) = ax + bx^3 + cx^5 + \dots$$

Can "tune" co-efficients for specific characteristics.

• Higher-order may converge faster, but each step is more expensive.

Do you have to converge?

NanoGPT speedrun :  $f(x) = 3.4442x - 4.7750x^3 + 2.0315x^5$

$$f(1) \neq 1$$