

Tuesday, May 5

Last class!

"Quiz" for specific feedback

↳ topics

↳ course components
(readings, activities, reviews,
self-grades, whiteboard lectures)

Student experience surveys

for more general feedback

Senior projects / talks due Thursday midnight

Presentations Wed., May 13 7-10pm

Only one more OH

(on Monday.)

Please reach out to meet!

Motivation

Previously, treated gradient descent as vector operation

But most weights are naturally in matrices (e.g., fully connected layers)

$W \in \mathbb{R}^{n \times d}$ is a weight matrix

$G = \nabla_W \mathcal{L} \in \mathbb{R}^{n \times d}$ is its gradient, $n \geq d$ WLOG

Goal: $\operatorname{argmax}_{\Delta W \in \mathbb{R}^{n \times d}} \mathcal{L}(W) - \mathcal{L}(W + \Delta W)$

Goal: $\operatorname{argmax}_{\Delta W \in \mathbb{R}^{n \times d}} \mathcal{L}(W) - \mathcal{L}(W + \Delta W)$

Taylor approximation: $\mathcal{L}(W + \Delta W) \approx \mathcal{L}(W) + \langle \nabla_W \mathcal{L}, \Delta W \rangle_F$

← sum of entry wise product

Rearranging,

$$\mathcal{L}(W) - \mathcal{L}(W + \Delta W) \approx - \langle \nabla_W \mathcal{L}, \Delta W \rangle_F$$

$$\begin{aligned} \langle A, B \rangle_F &= \sum_{j=1}^d \sum_{i=1}^n [A]_{i,j} [B]_{i,j} = \sum_{j=1}^d \sum_{i=1}^n [A^T]_{j,i} [B]_{i,j} \\ &= \sum_{j=1}^d [A^T B]_{j,j} = \operatorname{tr}(A^T B) \end{aligned}$$

$$\mathcal{L}(W) - \mathcal{L}(W + \Delta W) \approx - \operatorname{tr}(G^T \Delta W)$$

↑ max over ΔW st not too large

$$\mathcal{L}(W) - \mathcal{L}(W + \Delta W) \approx -\text{tr}(G^T \Delta W)$$

↑ max over ΔW st
 ΔW "not too large"

Frobenius norm: $\epsilon \geq \|\Delta W\|_F = \sqrt{\sum_{i,j} [\Delta W]_{i,j}^2} \Rightarrow$ gradient descent

Spectral norm $\epsilon \geq \|\Delta W\|_2 = \max_i \sigma_i(\Delta W) \Rightarrow$ Muon

$$\text{Let } G = \sum_{i=1}^d \sigma_i u_i v_i^T.$$

$$\text{Claim: } \underset{\Delta W: \|\Delta W\|_2 \leq \epsilon}{\text{argmin}} \text{tr}(G^T \Delta W) = -\epsilon \sum_{i=1}^d u_i v_i^T$$

Proof: Since $\{u_i\}_{i=1}^d$, $\{v_j\}_{j=1}^d$ orthonormal, we can write

$$\Delta W = \sum_{i=1}^n \sum_{j=1}^d c_{ij} u_i v_j^T$$

$$\text{tr}(G^T \Delta W) =$$

$$= \sum_{k=1}^d c_{kk} \sigma_k \quad \leftarrow \text{choose smallest } c_{kk} \text{ possible}$$

$$\|\Delta W\|_2 \leq \epsilon \Leftrightarrow \forall x, \|\Delta W x\|_2 \leq \epsilon$$

Choose $x = v_k$, then

$$\|\Delta W v_k\|_2^2 = \left\| \sum_{i=1}^n \sum_{j=1}^d c_{ij} u_i v_j^T v_k \right\|_2^2 = \left\| \sum_{i=1}^n c_{ik} u_i \right\|_2^2 = \sum_{i=1}^n c_{ik}^2 \leq \epsilon^2$$

$$c_{kk}^2 \leq \sum_{i=1}^n c_{ik}^2 \leq \epsilon^2$$

$$\Rightarrow |c_{kk}| < \epsilon$$

$$\Rightarrow \text{smallest } c_{kk} = -\epsilon$$

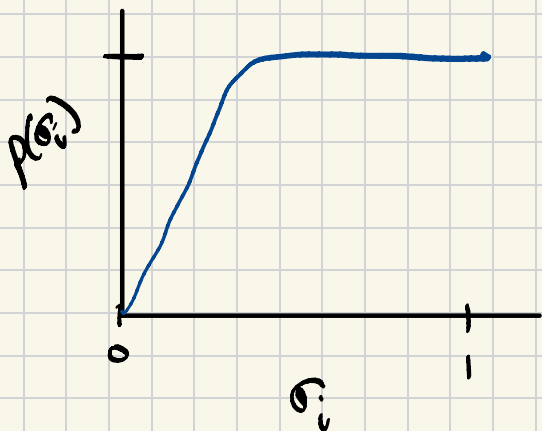
Then

$$c_{kk}^2 = \epsilon^2 \leq \sum_{i=1}^n c_{ik}^2 \leq \epsilon^2 \Leftrightarrow c_{ik} = 0 \quad \text{for } i \neq k$$

$$\Delta W = \sum_{i=1}^n -\epsilon \cdot u_i v_i^T$$

In practice, computing $\sum_{i=1}^m u_i v_i^T$ takes $O(nd^2)$ time.

Approximate via polynomial



Better polynomial \Rightarrow better optimizer

$$\frac{3}{2}x - \frac{1}{2}x^3$$

$$3.4445x - 4.775x^3 + 2.0315x^5$$

Polar Express