

Question 1

- $\mathcal{H} = \{x \mapsto y = W_N \sigma(W_{N-1} \cdots W_2 \sigma(W_1 x) \cdots); W_n \in \mathbb{R}^{d_n, d_{n-1}}, n \in \mathbb{N}\}$

Let $L(W_1, \dots, W_N)$ be a continuously differentiable loss function that depends on W_1, \dots, W_N only through the input-output mapping of the network.

Assume that the global minimum of $L(\cdot)$ is not attained at $W_1=0, W_2=0, \dots, W_N=0$.

Assume also that $\sigma(\cdot)$ is continuously differentiable and $\sigma(0)=0$. Then, $L(\cdot)$ is non-convex.

Proof

Assume by contradiction that L is convex.

Since $L(\cdot)$ is convex, it is enough to show that $\nabla L(0, \dots, 0) = 0$, because if it's true than we can yield that $(W_1=0, \dots, W_N=0)$ is a global min of $L(\cdot)$ in contradiction.

[proof of: $\nabla L(\vec{0}) = 0 \Rightarrow \vec{0}$ is a global min of $L(\cdot)$]

Since L is convex we know from calculus 1 that

$$\begin{aligned} \forall x \quad L(x) &\geq L(\vec{0}) + \nabla L(\vec{0})^T (x - \vec{0}) \\ &= L(\vec{0}) + 0(x - \vec{0}) \\ &= L(\vec{0}) \end{aligned}$$

$\Rightarrow \vec{0} = (0, \dots, 0)$ is a global min of $L(\cdot)$

Let's assume we calculate the loss L on M training points x_1, \dots, x_M with respectively vector labels y_1^*, \dots, y_M^* .

Since L depends on w_1, \dots, w_N only through the input-output mapping of the network we can write L as:

$$L(w_1, \dots, w_N) = L(W_N w_N + \dots + w_1 x_1, y_1^*)$$

$$W_N w_{N-1} + \dots + w_1 x_2, y_2^* \quad \text{⊗}$$

$$\vdots$$

$$W_N w_{N-1} + \dots + w_1 x_M, y_M^*)$$

[L depends on w_1, \dots, w_N only through the network outputs on the inputs x_1, \dots, x_M].

We would like to prove that $\nabla L(o, \dots, o) = 0$ in order to get the required. Thus, we'll prove that for every

$$j \in \{1, \dots, N\} : \frac{\partial L}{\partial w_j}(o, \dots, o) = 0 : \quad + \text{functions of } h_i \text{ below}$$

Notice that from ⊗ and the chain rule, we get:

$$\begin{aligned} \frac{\partial L}{\partial w_j}(o, \dots, o) &= \frac{\partial L}{\partial h_1}(h_1(o, \dots, o), y_1^*, \dots, h_m(o, \dots, o), y_M^*) \cdot \frac{\partial h_1}{\partial w_j}(o, \dots, o) \\ &+ \frac{\partial L}{\partial h_2}(h_1(o, \dots, o), y_1^*, \dots, h_m(o, \dots, o), y_M^*) \cdot \frac{\partial h_2}{\partial w_j}(o, \dots, o) \\ &+ \vdots \\ &+ \frac{\partial L}{\partial h_m}(h_1(o, \dots, o), y_1^*, \dots, h_m(o, \dots, o), y_M^*) \cdot \frac{\partial h_m}{\partial w_j}(o, \dots, o) \end{aligned}$$

Where $\forall 1 \leq i \leq m \quad h_i(w_1, \dots, w_N) = W_N w_{N-1} + \dots + w_1 x_i$ [by definition].

Thus, in order to show that $\frac{\partial L}{\partial w_j}(o, \dots, o) = 0$, we can

Show that $\frac{\partial h_i}{\partial w_j}(o, \dots, o) = 0 \quad \forall 1 \leq i \leq m$ and to finish.

As defined $h_i(w_1, \dots, w_N) = w_N \sigma w_{N-1} \dots w_1 x_i$
 $= w_N \sigma \dots w_j \sigma \dots w_1 x_i$

We'll define: $f_i(w_i, \dots, w_j) = w_j \sigma \dots w_1 x_i$

$$g(v, w_1, \dots, w_N) = w_N \sigma w_{N-1} \dots w_1 \sigma(v)$$

Therefore by these definitions we yield

$$h_i(w_1, \dots, w_N) = g(f_i(w_1, \dots, w_j), w_{j+1}, \dots, w_N)$$

$$\Rightarrow \frac{\partial h_i}{\partial w_j}(o, \dots, o) = \frac{\partial g}{\partial f_i}(o, \dots, o) \cdot \frac{\partial f_i}{\partial w_j}(o, \dots, o)$$

Thus we can show that $\frac{\partial f_i}{\partial w_j}(o, \dots, o) = 0$ and finish.

$$f_i(w_1, \dots, w_j) = w_j \sigma \underbrace{w_{j-1} \dots w_1 x_i}_{\text{const}}$$

\Rightarrow When we calculate the derivative of f_i by w_j we can treat const as a const vector $y = \sigma w_{j-1} \dots w_1 x_i$:

$$\frac{\partial f_i}{\partial w_j}(w_1, \dots, w_j) = \frac{\partial (w_j y)}{\partial w_j}(w_1, \dots, w_j), \quad w_j \in \mathbb{R}^{d_j, d_{j-1}}, \quad y \in \mathbb{R}^{d_{j-1}} = \begin{pmatrix} y_1 \\ \vdots \\ y_{d_{j-1}} \end{pmatrix}$$

Let's sign $d_j = m$, $d_{j-1} = k$ for simplicity. Then:

$$w_j y = \begin{pmatrix} (w_j)_{11} y_1 + \dots + (w_j)_{1k} y_k \\ \vdots \\ (w_j)_{m1} y_1 + \dots + (w_j)_{mk} y_k \end{pmatrix}$$

Thus $\forall 1 \leq i \leq m, 1 \leq j \leq k$:

$$\frac{\partial w_j y}{\partial (w_i)_{ij}} = \begin{pmatrix} 0 \\ \vdots \\ y_j \\ \vdots \\ 0 \end{pmatrix} \leftarrow \text{1 zero}$$

As we proved $\forall 1 \leq i \leq m \quad \exists j \leq k$:

$$\frac{\partial w_i y}{\partial (w_j)_{ij}} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ y_j \\ 0 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i$$

$$\Rightarrow \frac{\partial w_i y}{\partial (w_j)_{ij}}(0, \dots, 0) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ y(0, \dots, 0)_j \\ 0 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i$$

Remember $y = \sigma^{-1} w_{j-1} \dots w_1 x_i$

$$\Rightarrow y(0, \dots, 0) = \sigma^{-1} 0 \dots x_i = \sigma(0) = 0$$

Therefore we've got $\frac{\partial w_i y}{\partial (w_j)_{ij}}(0, \dots, 0) = 0$

$$\Rightarrow \frac{\partial w_i y}{\partial w_j}(0, \dots, 0) = 0$$

$$\Rightarrow \frac{\partial f_i}{\partial w_j}(0, \dots, 0) = \frac{\partial w_i y}{\partial w_j}(0, \dots, 0) = 0$$

as required.

□

Foundations of Deep Learning – Homework Assignment #3

Adi Alm & Tomer Epshtain

Part 2: (1)

Question:

In this question we extend the convergence of stationary point result delivered in class from gradient descent to stochastic gradient descent.

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a twice continuously differentiable and β -smooth function that attains its global minimum $f^* = \min_{\vec{w} \in \mathbb{R}^d} f(\vec{w})$. Suppose we run stochastic gradient descent over $f(\cdot)$:

$$\vec{w}_{t+1} \leftarrow \vec{w}_t - \eta_t (\nabla f(\vec{w}_t) + \xi_t), t = 1, 2, 3, \dots$$

Where η_t is a (possibly time varying) learning rate, and ξ_t stands for a time-independent noise with zero mean ($\mathbb{E}[\xi_t] = 0$) and σ^2 variance ($\mathbb{E}[\|\xi_t\|^2] = \sigma^2$). Let $\epsilon > 0$.

Assume $\eta_t = \frac{1}{\beta}$ and derive an upper bound on the number of iterations needed for reaching an expected $(\epsilon + \sigma)$ -stationary point, i.e. on $T \in \mathbb{N}$ that will ensure:

$$\min_{t \in [T]} \mathbb{E}[\|\nabla f(\vec{w}_t)\|] \leq \epsilon + \sigma$$

Proof:

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a twice continuously differentiable and β -smooth function that attains its global minimum $f^* = \min_{\vec{w} \in \mathbb{R}^d} f(\vec{w})$.

Reminder: In class we proved the following lemma:

Lemma:

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be twice continuously differentiable and β -smooth. Then:

$$\forall \vec{w}_1, \vec{w}_2 \in \mathbb{R}^d: \left| f(\vec{w}_2) - (f(\vec{w}_1) + \underbrace{<\nabla f(\vec{w}_1), \vec{w}_2 - \vec{w}_1>}_{1^{st} \text{ order approx around } \vec{w}_1}) \right| \leq \frac{\beta}{2} \|\vec{w}_1 - \vec{w}_2\|^2$$

So, for any $t \geq 0$:

$$f(\vec{w}_{t+1}) \leq f(\vec{w}_t) + <\nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t> + \frac{\beta}{2} \|\vec{w}_{t+1} - \vec{w}_t\|^2$$

We'll take expected value over $\xi_1, \xi_2, \dots, \xi_t$:

Denote- $E_{t+1} := \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[f(\vec{w}_{t+1})]$, so

$$\begin{aligned} E_{t+1} &= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[f(\vec{w}_{t+1})] \leq \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}\left[f(\vec{w}_t) + <\nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t> + \frac{\beta}{2} \|\vec{w}_{t+1} - \vec{w}_t\|^2\right] = \\ &= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[f(\vec{w}_t)] + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[<\nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t>] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[\|\vec{w}_{t+1} - \vec{w}_t\|^2] = \end{aligned}$$

w_t is independent of ξ_t :

$$\begin{aligned} &= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}}[f(\vec{w}_t)] + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[<\nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t>] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[\|\vec{w}_{t+1} - \vec{w}_t\|^2] = \\ &= E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[<\nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t>] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t}[\|\vec{w}_{t+1} - \vec{w}_t\|^2] (*) = \end{aligned}$$

Plugging in the SGD step definition

$$\vec{w}_{t+1} \leftarrow \vec{w}_t - \eta_t (\nabla f(\vec{w}_t) + \xi_t), t = 1, 2, 3, \dots$$

We achieve:

$$\begin{aligned}
E_{t+1} &\leq E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t \rangle] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|w_{t+1} - w_t\|^2] = \\
&= E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), -\eta_t (\nabla f(\vec{w}_t) + \xi_t) \rangle] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\eta_t (\nabla f(\vec{w}_t) + \xi_t)\|^2] \stackrel{\eta_t \equiv \text{const}}{=} \\
&= E_t - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t)\|^2] - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \xi_t \rangle] + \frac{\beta}{2} \eta_t^2 \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2] = \\
&= E_t - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \xi_t \rangle] + \frac{\beta}{2} \eta_t^2 \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2] =
\end{aligned}$$

$\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \xi_t \rangle] = 0$, because $\mathbb{E}[\xi_t] = 0$, and linearity of expected value with inner product.

$$= E_t - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \frac{\beta}{2} \eta_t^2 \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2]$$

Now, let's focus on $\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2]$:

$$\begin{aligned}
\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2] &= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t)\|^2 + \|\xi_t\|^2 + 2 \langle \nabla f(\vec{w}_t), \xi_t \rangle] = \\
&= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t)\|^2] + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\xi_t\|^2] + 2 \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \xi_t \rangle] =
\end{aligned}$$

Here, once again, $\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \xi_t \rangle] = 0$. And $\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\xi_t\|^2] = \mathbb{E}_{\xi_t} [\|\xi_t\|^2] = \sigma^2$ by ξ_t 's distribution's definition.

$$= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \sigma^2$$

So we have:

$$\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2] = \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \sigma^2$$

Plugging in the result:

$$\begin{aligned}
E_{t+1} &\leq E_t - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \frac{\beta}{2} \eta_t^2 \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\nabla f(\vec{w}_t) + \xi_t\|^2] = \\
&= E_t - \eta_t \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \frac{\beta}{2} \eta_t^2 (\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \sigma^2)
\end{aligned}$$

Plugging in $\eta_t = \frac{1}{\beta}$:

$$E_{t+1} \leq E_t - \frac{1}{\beta} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \frac{1}{2\beta} (\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \sigma^2) = E_t - \frac{1}{2\beta} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \frac{\sigma^2}{2\beta}$$

Assume that for steps $t = 0, 1, \dots, T-1$ we didn't achieve an expected $(\epsilon + \sigma)$ -stationary point.

i.e. $\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|] > \epsilon + \sigma$ for all $t = 0, 1, \dots, T-1$.

For any random variable X : $\mathbb{E}[X^2] = \mathbb{E}[X]^2 + \text{Var}(X) \geq \mathbb{E}[X]^2$.

So, for steps $t = 0, 1, \dots, T-1$:

$$\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] \geq \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|]^2 > (\epsilon + \sigma)^2$$

So:

$$E_{t+1} \leq E_t - \frac{1}{2\beta} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2] + \frac{\sigma^2}{2\beta} < E_t - \frac{1}{2\beta} (\epsilon + \sigma)^2 + \frac{\sigma^2}{2\beta} = E_t - \frac{\epsilon}{\beta} \left(\frac{\epsilon}{2} + \sigma \right)$$

This inequality holds for all $t = 0, 1, \dots, T-1$. So:

$$E_T < E_0 - T \frac{\epsilon}{\beta} \left(\frac{\epsilon}{2} + \sigma \right)$$

Since f^* is f 's global minimum, $f^* \leq E_T (= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{T-1}} [f(\vec{w}_T)])$

So we have

$$f^* < E_0 - T \frac{\epsilon}{\beta} \left(\frac{\epsilon}{2} + \sigma \right)$$

Where $E_0 = f(\vec{w_0})$, so:

$$f^* < f(\vec{w_0}) - T \frac{\epsilon}{\beta} \left(\frac{\epsilon}{2} + \sigma \right) \Rightarrow T < \frac{2\beta(f(\vec{w_0}) - f^*)}{\epsilon(\epsilon + 2\sigma)}$$

So, for any $\epsilon > 0$ an expected $(\epsilon + \sigma)$ -stationary point will be reached within a number of steps at most

$$T < \frac{2\beta(f(\vec{w_0}) - f^*)}{\epsilon(\epsilon + 2\sigma)}$$

Foundations of Deep Learning – Homework Assignment #3

Adi Almog & Tomer Epshteyn

Part 2: (2)

Question: (Bonus)

Show that by decaying learning rate it is possible to guarantee convergence to expected ϵ -stationary point. In particular, define a scheme for η_t and a respective upper bound on $T \in \mathbb{N}$ that will ensure:

$$\min_{t \in [T]} \mathbb{E}[\|\nabla f(\vec{w}_t)\|] \leq \epsilon$$

Proof:

Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a twice continuously differentiable and β -smooth function that attains its global minimum $f^* = \min_{\vec{w} \in \mathbb{R}^d} f(\vec{w})$.

We define an update scheme for $\eta_t := \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)}$

In part(2) 1 we derived

$$E_{t+1} \leq E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\langle \nabla f(\vec{w}_t), \vec{w}_{t+1} - \vec{w}_t \rangle] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} [\|\vec{w}_{t+1} - \vec{w}_t\|^2] \quad (*)$$

Plugging in SGD's update scheme $\vec{w}_{t+1} \leftarrow \vec{w}_t - \eta_t (\nabla f(\vec{w}_t) + \xi_t)$ where $\eta_t = \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)}$, we achieve:

$$\begin{aligned} E_{t+1} &\leq E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\langle \nabla f(\vec{w}_t), \vec{w}_t - \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} (\nabla f(\vec{w}_t) + \xi_t) - \vec{w}_t \rangle \right] \\ &\quad + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\left\| \vec{w}_t - \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} (\nabla f(\vec{w}_t) + \xi_t) - \vec{w}_t \right\|^2 \right] = \\ &= E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\langle \nabla f(\vec{w}_t), -\frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} (\nabla f(\vec{w}_t) + \xi_t) \rangle \right] \\ &\quad + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\left\| -\frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} (\nabla f(\vec{w}_t) + \xi_t) \right\|^2 \right] = \\ &= E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[-\frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \langle \nabla f(\vec{w}_t), \nabla f(\vec{w}_t) + \xi_t \rangle \right] \\ &\quad + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \|\nabla f(\vec{w}_t) + \xi_t\|^2 \right] = \\ &= E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[-\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} - \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \langle \nabla f(\vec{w}_t), \xi_t \rangle \right] \\ &\quad + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\underbrace{\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} (\|\nabla f(\vec{w}_t)\|^2 + 2 \langle \nabla f(\vec{w}_t), \xi_t \rangle + \|\xi_t\|^2)}_{(IV)} \right] = \\ &= E_t - \underbrace{\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \right]}_{(I)} - \underbrace{\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \langle \nabla f(\vec{w}_t), \xi_t \rangle \right]}_{(II)} \\ &\quad + \underbrace{\frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^6}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \right]}_{(III)} + \underbrace{\beta \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \langle \nabla f(\vec{w}_t), \xi_t \rangle \right]}_{(IV)} \\ &\quad + \underbrace{\frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \|\xi_t\|^2 \right]}_{(V)} \end{aligned}$$

We'll analyze each (I), (II), (III), (IV), (V) separately.

- (I) $= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \right] = \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \right]$
Because \vec{w}_t is independent of ξ_t
- (II) $= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} < \nabla f(\vec{w}_t), \xi_t \right] = \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[< \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \nabla f(\vec{w}_t), \xi_t \right] = 0$
Because \vec{w}_t is independent of ξ_t and $\mathbb{E}[\xi_t] = 0$, and from linearity of expected value with inner product, LHS is independent of ξ_t .
- (III) $= \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^6}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \right] = \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^6}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \right]$
Because \vec{w}_t is independent of ξ_t
- (IV) $= \beta \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} < \nabla f(\vec{w}_t), \xi_t \right] = \beta \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[< \frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \nabla f(\vec{w}_t), \xi_t \right] = 0$
Because \vec{w}_t is independent of ξ_t and $\mathbb{E}[\xi_t] = 0$, and from linearity of expected value with inner product, LHS is independent of ξ_t .
- (V) $= \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_t} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \|\xi_t\|^2 \right] = \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \right] \sigma^2$
Because \vec{w}_t is independent of ξ_t and $\mathbb{E}[\|\xi_t\|^2] = \sigma^2$

Bringing it all together we achieve

$$E_{t+1} \leq$$

$$\begin{aligned} E_t - \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \right] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^6}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \right] + \frac{\beta}{2} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta^2(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \right] \sigma^2 \\ = E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[-\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} + \frac{\|\nabla f(\vec{w}_t)\|^6}{2\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} + \frac{\|\nabla f(\vec{w}_t)\|^4}{2\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \sigma^2 \right] \\ = E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[-\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} + \frac{\|\nabla f(\vec{w}_t)\|^4}{2\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \|\nabla f(\vec{w}_t)\|^2 + \frac{\|\nabla f(\vec{w}_t)\|^4}{2\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} \sigma^2 \right] \\ = E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[-\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} + \frac{\|\nabla f(\vec{w}_t)\|^4}{2\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)^2} (\|\nabla f(\vec{w}_t)\|^2 + \sigma^2) \right] \\ = E_t + \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[-\frac{\|\nabla f(\vec{w}_t)\|^4}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} + \frac{\|\nabla f(\vec{w}_t)\|^4}{2\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)} \right] \\ = E_t - \frac{1}{2\beta} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\|\nabla f(\vec{w}_t)\|^2 + \sigma^2} \right] \end{aligned}$$

So

$$E_{t+1} \leq E_t - \frac{1}{2\beta} \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\|\nabla f(\vec{w}_t)\|^2 + \sigma^2} \right]$$

Assume that for steps $t = 0, 1, \dots, T-1$ we didn't achieve an expected ϵ -stationary point.

i.e. $\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|] > \epsilon$ for all $t = 0, 1, \dots, T-1$.

Define $g: [0, \infty) \rightarrow \mathbb{R}$ by $g(x) = \frac{x^4}{x^2 + \sigma^2}$. g is convex monotonically increasing.

So, for all $x > \epsilon$, $g(x) > \frac{\epsilon^4}{\epsilon^2 + \sigma^2}$.

g is convex, so by Jensen's inequality:

$$\mathbb{E}[X] > \epsilon \Rightarrow \mathbb{E}[g(X)] > g(\epsilon)$$

i.e.

$$\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\|\nabla f(\vec{w}_t)\|^2 + \sigma^2} \right] = \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [g(\|\nabla f(\vec{w}_t)\|)] > g(\epsilon) = \frac{\epsilon^4}{\epsilon^2 + \sigma^2}$$

So, for steps $t = 0, 1, \dots, T-1$:

$$\mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} \left[\frac{\|\nabla f(\vec{w}_t)\|^4}{\|\nabla f(\vec{w}_t)\|^2 + \sigma^2} \right] > \frac{\epsilon^4}{\epsilon^2 + \sigma^2}$$

So:

$$E_{t+1} < E_t - \frac{1}{2\beta} \frac{\epsilon^4}{\epsilon^2 + \sigma^2}$$

This inequality holds for all $t = 0, 1, \dots, T-1$:

$$E_T < E_0 - T \frac{1}{2\beta} \frac{\epsilon^4}{\epsilon^2 + \sigma^2}$$

Since f^* is f 's global minimum, $f^* \leq E_T (= \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{T-1}} [f(\vec{w}_T)])$

So,

$$f^* < E_0 - T \frac{1}{2\beta} \frac{\epsilon^4}{\epsilon^2 + \sigma^2}$$

Where $E_0 = f(\vec{w}_0)$, so:

$$f^* < f(\vec{w}_0) - T \frac{\epsilon^4}{2\beta(\epsilon^2 + \sigma^2)}$$

\Rightarrow

$$T < \frac{2\beta(\epsilon^2 + \sigma^2)(f(\vec{w}_0) - f^*)}{\epsilon^4}$$

Side note:

Where did the update rule for η_t , $\eta_t = \frac{\|\nabla f(\vec{w}_t)\|^2}{\beta(\|\nabla f(\vec{w}_t)\|^2 + \sigma^2)}$ come from?

Denote $g_t := \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_{t-1}} [\|\nabla f(\vec{w}_t)\|^2]$, In part 2(1) we showed:

$$E_{t+1} \leq E_t + \left(-\eta_t g_t + \frac{\beta}{2} \eta_t^2 (g_t + \sigma^2) \right)$$

We would like to select an update scheme for η_t such that the above bound $(-\eta_t g_t + \frac{\beta}{2} \eta_t^2 (g_t + \sigma^2))$ is minimal.

Let $t \in \mathbb{N}$, define $h_t: \mathbb{R} \rightarrow \mathbb{R}$, $h_t(\eta_t) := \frac{\beta}{2} (g_t + \sigma^2) \eta_t^2 - g_t \eta_t$.

h_t is a parabola that attains its minimum at $\eta_t = \frac{g_t}{\beta(g_t + \sigma^2)}$

(This can easily be seen noticing h_t is a parabola of the form $h_t(x) = ax^2 + bx$ with $a > 0$ so its minimum is achieved at $x = -\frac{b}{2a}$)

Selecting this update scheme gives us the desired bound

Foundations of Deep Learning – Homework Assignment #3

Adi Alm & Tomer Epshtain

Part 2: (3)

Question:

Let $\ell: \mathbb{R}^{d,d} \rightarrow \mathbb{R}$ be a twice continuously differentiable convex loss, overparameterized by a depth N linear neural network with hidden widths d :

$$\phi: \mathbb{R}^{d,d} \times \mathbb{R}^{d,d} \times \cdots \times \mathbb{R}^{d,d} \rightarrow \mathbb{R}, \quad \phi(W_1, W_2, \dots, W_N) = \ell(W_N W_{N-1} \cdots W_1)$$

Prove the following claims:

- A. If we restrict the domain of $\phi(\cdot)$ to $B^{d,d} \times B^{d,d} \times \cdots \times B^{d,d}$, where $B^{d,d} := \{x \in \mathbb{R}^{d,d}: \|x\|_F \leq 1\}$, then we obtain a smooth (Lipschitz gradient) function.
- B. If we do not restrict the domain of $\phi(\cdot)$, the function is smooth if and only if at least one of the following conditions hold:
 - a. $\ell(\cdot)$ is constant
 - b. $\ell(\cdot)$ is affine and the depth is $N = 2$.

Proof (A):

Let $(W_1, W_2, \dots, W_N), (\tilde{W}_1, \tilde{W}_2, \dots, \tilde{W}_N) \in B^{d,d} \times B^{d,d} \times \cdots \times B^{d,d}$.

Denote by $W_{j:j'} = W_j \cdots W_{j'}$ for $1 \leq j \leq j' \leq N$ and Id matrix otherwise.

Let's take a look at the expression:

$$\nabla \phi(\tilde{W}_1, \tilde{W}_2, \dots, \tilde{W}_N) - \nabla \phi(W_1, W_2, \dots, W_N) \in \underbrace{\mathbb{R}^{d,d} \times \mathbb{R}^{d,d} \times \cdots \times \mathbb{R}^{d,d}}_{N \text{ times}}$$

We saw in class:

For $i \in [N]$, let's look at $(\nabla \phi(W_1, W_2, \dots, W_N))_i \in \mathbb{R}^{d,d}$

$$(\nabla \phi(W_1, W_2, \dots, W_N))_i = \frac{\partial}{\partial W_i} \phi(W_1, W_2, \dots, W_N) = W_{i+1:N}^T \nabla l(W_{1:N}) W_{1:i-1}^T$$

Yielding

$$\begin{aligned} & (\nabla \phi(\tilde{W}_1, \tilde{W}_2, \dots, \tilde{W}_N) - \nabla \phi(W_1, W_2, \dots, W_N))_i = \\ &= \frac{\partial}{\partial \tilde{W}_i} \phi(\tilde{W}_1, \tilde{W}_2, \dots, \tilde{W}_N) - \frac{\partial}{\partial W_i} \phi(W_1, W_2, \dots, W_N) = \\ &= \tilde{W}_{i+1:N}^T \nabla l(\tilde{W}_{1:N}) \tilde{W}_{1:i-1}^T - W_{i+1:N}^T \nabla l(W_{1:N}) W_{1:i-1}^T \end{aligned}$$

We'll add and subtract terms...

$$\begin{aligned} &= \underbrace{\tilde{W}_{i+1:N}^T \nabla l(\tilde{W}_{1:N}) \tilde{W}_{1:i-1}^T - W_{i+1:N}^T \nabla l(W_{1:N}) \tilde{W}_{1:i-1}^T}_{(A)} + \\ &+ \underbrace{W_{i+1:N}^T \nabla l(\tilde{W}_{1:N}) \tilde{W}_{1:i-1}^T - W_{i+1:N}^T \nabla l(W_{1:N}) \tilde{W}_{1:i-1}^T}_{(B)} + \\ &+ \underbrace{W_{i+1:N}^T \nabla l(W_{1:N}) \tilde{W}_{1:i-1}^T - W_{i+1:N}^T \nabla l(W_{1:N}) W_{1:i-1}^T}_{(C)} = \\ &= \underbrace{(\tilde{W}_{i+1:N}^T - W_{i+1:N}^T) \nabla l(\tilde{W}_{1:N}) \tilde{W}_{1:i-1}^T}_{(A)} + \underbrace{W_{i+1:N}^T (\nabla l(\tilde{W}_{1:N}) - \nabla l(W_{1:N})) \tilde{W}_{1:i-1}^T}_{(B)} + \underbrace{W_{i+1:N}^T \nabla l(W_{1:N}) (\tilde{W}_{1:i-1}^T - W_{1:i-1}^T)}_{(C)} \end{aligned}$$

So:

$$\|\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N)\|_2^2 = \sum_{i=1}^N \left\| (\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N))_i \right\|_F^2 \leq$$

Denote $j := \operatorname{argmax}_{i \in [N]} \left(\left\| (\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N))_i \right\|_F^2 \right)$, so

$$\begin{aligned} & \leq \sum_{i=1}^N \left\| (\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N))_j \right\|_F^2 = N \cdot \left\| (\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N))_j \right\|_F^2 = \\ & = N \cdot \left\| \underbrace{(\widetilde{W}_{j+1:N}^T - W_{j+1:N}^T) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T}_{(A)} + \underbrace{W_{j+1:N}^T (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_{1:j-1}^T}_{(B)} + \underbrace{W_{j+1:N}^T \nabla l(W_{1:N}) (\widetilde{W}_{1:j-1}^T - W_{1:j-1}^T)}_{(C)} \right\|_F^2 \stackrel{\Delta}{\leq} \\ & N \cdot \left(\underbrace{\left\| (\widetilde{W}_{j+1:N}^T - W_{j+1:N}^T) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T \right\|_F^2}_{(A)} + \underbrace{\left\| W_{j+1:N}^T (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_{1:j-1}^T \right\|_F^2}_{(B)} \right. \\ & \quad \left. + \underbrace{\left\| W_{j+1:N}^T \nabla l(W_{1:N}) (\widetilde{W}_{1:j-1}^T - W_{1:j-1}^T) \right\|_F^2}_{(C)} \right) \end{aligned}$$

We will bound each (A), (B) & (C) separately.

Propositions:

1. $\forall k \in \mathbb{N}$, and $A_1, \dots, A_k, \tilde{A}_1, \dots, \tilde{A}_k \in B^{d,d}$:

$$\|\tilde{A}_{1:k} - A_{1:k}\|_F^2 \leq \sum_{i=1}^k \|\tilde{A}_i - A_i\|_F^2$$

Notice that if we denote $(A_1, \dots, A_k), (\tilde{A}_1, \dots, \tilde{A}_k) \in B^{d,d} \times B^{d,d} \times \dots \times B^{d,d}$ we obtain:

$$\|\tilde{A}_{1:k} - A_{1:k}\|_F^2 \leq \|(\tilde{A}_1, \dots, \tilde{A}_k) - (A_1, \dots, A_k)\|_2^2$$

2. $\forall k \in \mathbb{N}$, and $A_1, \dots, A_k \in B^{d,d}$:

$$A_1 \cdots A_k \in B^{d,d}$$

We will use these propositions to complete our proof. We'll prove them later.

Note:

- i. $\ell: \mathbb{R}^{d,d} \rightarrow \mathbb{R}$ is a twice continuously differentiable convex function, $B^{d,d}$ is compact and convex and therefore $\ell|_{B^{d,d}}: B^{d,d} \rightarrow \mathbb{R}$ is Lipschitz gradient. I.e there exists $\beta > 0$, such that $\forall A, B \in B^{d,d}$:

$$\|\nabla\ell(B) - \nabla\ell(A)\|_F \leq \beta \|B - A\|_F$$

(This is true by claim in course's recap notes on convexity and Lipschitzness)

- ii. $\ell: \mathbb{R}^{d,d} \rightarrow \mathbb{R}$ is a twice continuously differentiable function, $B^{d,d}$ is compact, therefor $\nabla\ell|_{B^{d,d}}: B^{d,d} \rightarrow \mathbb{R}$ is bounded on $B^{d,d}$. I.e. there exists $M > 0$ such that $\forall A \in B^{d,d}$:

$$\|\nabla\ell(A)\|_F \leq M$$

Bound on (A):

$$(A) = \|(\tilde{W}_{j+1:N}^T - W_{j+1:N}^T) \nabla l(\tilde{W}_{1:N}) \tilde{W}_{1:j-1}^T\|_F^2 \leq \|\tilde{W}_{j+1:N}^T - W_{j+1:N}^T\|_F^2 \|\nabla l(\tilde{W}_{1:N})\|_F^2 \|\tilde{W}_{1:j-1}^T\|_F^2$$

- By proposition 1:

$$\|\tilde{W}_{j+1:N}^T - W_{j+1:N}^T\|_F^2 \leq \sum_{i=j+1}^N \|\tilde{W}_i^T - W_i^T\|_F^2$$

- $\tilde{W}_{1:N} \in B^{d,d}$ because $\tilde{W}_{1:N} = \tilde{W}_N \cdots \tilde{W}_1$ where $\tilde{W}_i \in B^{d,d} \forall i \in [N]$ (from proposition 2), so, by note (ii)

$$\|\nabla l(\tilde{W}_{1:N})\|_F^2 \leq M^2$$

- $\tilde{W}_{1:N}^T \in B^{d,d}$ because $\tilde{W}_{1:N}^T = (W_N \cdots W_1)^T = W_1^T \cdots W_N^T$ where $W_i^T \in B^{d,d} \forall i \in [N]$ (from proposition 2), so

$$\|\tilde{W}_{1:N}^T\|_F \leq 1$$

Bringing it together:

$$\begin{aligned} (A) &\leq \|\tilde{W}_{j+1:N}^T - W_{j+1:N}^T\|_F^2 \|\nabla l(\tilde{W}_{1:N})\|_F^2 \|\tilde{W}_{1:j-1}^T\|_F^2 \leq M^2 \cdot \sum_{i=j+1}^N \|\tilde{W}_i^T - W_i^T\|_F^2 \leq \\ &\leq M^2 \cdot \sum_{i=1}^N \|\tilde{W}_i^T - W_i^T\|_F^2 = M^2 \cdot \|(\tilde{W}_1^T, \dots, \tilde{W}_N^T) - (W_1^T, \dots, W_N^T)\|_2^2 = \\ &= M^2 \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2 \end{aligned}$$

Bound on (B):

$$(B) = \|W_{j+1:N}^T (\nabla l(\tilde{W}_{1:N}) - \nabla l(W_{1:N})) \tilde{W}_{1:j-1}^T\|_F^2 \leq \|W_{j+1:N}^T\|_F^2 \|\nabla l(\tilde{W}_{1:N}) - \nabla l(W_{1:N})\|_F^2 \|\tilde{W}_{1:j-1}^T\|_F^2$$

- $\tilde{W}_{j+1:N}^T \in B^{d,d}$ because $\tilde{W}_{j+1:N}^T = (W_N \cdots W_{j+1})^T = W_{j+1}^T \cdots W_N^T$ where $W_i^T \in B^{d,d} \forall i \in [N]$ (from proposition 2), so

$$\|\tilde{W}_{j+1:N}^T\|_F^2 \leq 1$$

- by note (i):

$$\|\nabla l(\tilde{W}_{1:N}) - \nabla l(W_{1:N})\|_F^2 \leq \beta^2 \cdot \|\tilde{W}_{1:N} - W_{1:N}\|_F^2 \leq$$

Now, by proposition 1:

$$\leq \beta^2 \cdot \sum_{i=1}^N \|\tilde{W}_i^T - W_i^T\|_F^2 = \beta^2 \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2$$

- $\tilde{W}_{1:j-1}^T \in B^{d,d}$ because $\tilde{W}_{1:j-1}^T = (W_{j-1} \cdots W_1)^T = W_1^T \cdots W_{j-1}^T$ where $W_i^T \in B^{d,d} \forall i \in [N]$ (from proposition 2), so

$$\|\tilde{W}_{1:j-1}^T\|_F \leq 1$$

Bringing it together:

$$(B) \leq \|W_{j+1:N}^T\|_F^2 \|\nabla l(\tilde{W}_{1:N}) - \nabla l(W_{1:N})\|_F^2 \|\tilde{W}_{1:j-1}^T\|_F^2 \leq \beta^2 \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2$$

Bound on (C):

$$(C) = \|W_{j+1:N}^T \nabla l(W_{1:N})(\tilde{W}_{1:j-1}^T - W_{1:j-1}^T)\|_F^2 \leq \|W_{j+1:N}^T\|_F^2 \|\nabla l(W_{1:N})\|_F^2 \|\tilde{W}_{1:j-1}^T - W_{1:j-1}^T\|_F^2$$

- $\tilde{W}_{j+1:N}^T \in B^{d,d}$ because $\tilde{W}_{j+1:N}^T = (W_N \cdots W_{j+1})^T = W_{j+1}^T \cdots W_N^T$ where $W_i^T \in B^{d,d} \forall i \in [N]$ (from proposition 2), so

$$\|\tilde{W}_{j+1:N}^T\|_F \leq 1$$

- $W_{1:N} \in B^{d,d}$ because $W_{1:N} = W_N \cdots W_1$ where $W_i \in B^{d,d} \forall i \in [N]$ (from proposition 2), so, by note (ii)

$$\|\nabla l(W_{1:N})\|_F^2 \leq M^2$$

- By proposition 1:

$$\|\tilde{W}_{1:j-1}^T - W_{1:j-1}^T\|_F^2 \leq \sum_{i=1}^{j-1} \|\tilde{W}_i^T - W_i^T\|_F^2$$

Bringing it together:

$$\begin{aligned} (C) &\leq \|W_{j+1:N}^T\|_F^2 \|\nabla l(W_{1:N})\|_F^2 \|\tilde{W}_{1:j-1}^T - W_{1:j-1}^T\|_F^2 \leq M^2 \cdot \sum_{i=1}^{j-1} \|\tilde{W}_i^T - W_i^T\|_F^2 \leq \\ &\leq M^2 \cdot \sum_{i=1}^N \|\tilde{W}_i^T - W_i^T\|_F^2 = M^2 \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2 \end{aligned}$$

So:

$$\begin{aligned} \|\nabla \phi(\tilde{W}_1, \tilde{W}_2, \dots, \tilde{W}_N) - \nabla \phi(W_1, W_2, \dots, W_N)\|_2^2 &\leq \\ &\leq N \cdot ((A) + (B) + (C)) \leq \\ &\leq N \cdot (M^2 \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2 + \beta^2 \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2 + M^2 \\ &\quad \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2) = \\ &= (2 \cdot N \cdot M^2 + N \cdot \beta^2) \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2 \end{aligned}$$

Overall:

$$\|\nabla \phi(\tilde{W}_1, \tilde{W}_2, \dots, \tilde{W}_N) - \nabla \phi(W_1, W_2, \dots, W_N)\|_2^2 \leq (2 \cdot N \cdot M^2 + N \cdot \beta^2) \cdot \|(\tilde{W}_1, \dots, \tilde{W}_N) - (W_1, \dots, W_N)\|_2^2$$

So $\phi: B^{d,d} \times B^{d,d} \times \dots \times B^{d,d} \rightarrow \mathbb{R}$ is a smooth ρ -Lipschitz-gradient function with $\rho = \sqrt{2 \cdot N \cdot M^2 + N \cdot \beta^2}$.

■

We will now prove the two remaining propositions:

1. $\forall k \in \mathbb{N}$, and $A_1, \dots, A_k, \tilde{A}_1, \dots, \tilde{A}_k \in B^{d,d}$:

$$\|\tilde{A}_{1:k} - A_{1:k}\|_F^2 \leq \sum_{i=1}^k \|\tilde{A}_i - A_i\|_F^2$$

Notice that if we denote $(A_1, \dots, A_k), (\tilde{A}_1, \dots, \tilde{A}_k) \in B^{d,d} \times B^{d,d} \times \dots \times B^{d,d}$ we obtain:

$$\|\tilde{A}_{1:k} - A_{1:k}\|_F^2 \leq \|(\tilde{A}_1, \dots, \tilde{A}_k) - (A_1, \dots, A_k)\|_2^2$$

2. $\forall k \in \mathbb{N}$, and $A_1, \dots, A_k \in B^{d,d}$:

$$A_1 \cdots A_k \in B^{d,d}$$

Proof of proposition 1:

Proof by induction:

- $k = 1$:
Immediate.

$$\|\tilde{A}_{1:1} - A_{1:1}\|_F^2 = \|\tilde{A}_1 - A_1\|_F^2 = \sum_{i=1}^1 \|\tilde{A}_i - A_i\|_F^2$$

- $(k-1) \rightarrow k$:

Let $A_1, \dots, A_k, \tilde{A}_1, \dots, \tilde{A}_k \in B^{d,d}$.

$$\begin{aligned} \|\tilde{A}_{1:k} - A_{1:k}\|_F^2 &= \|\tilde{A}_{1:k-1}\tilde{A}_k - A_{1:k-1}A_k\|_F^2 = \|\tilde{A}_{1:k-1}\tilde{A}_k - \tilde{A}_{1:k-1}A_k + \tilde{A}_{1:k-1}A_k - A_{1:k-1}A_k\|_F^2 = \\ &= \|\tilde{A}_{1:k-1}(\tilde{A}_k - A_k) + (\tilde{A}_{1:k-1} - A_{1:k-1})A_k\|_F^2 \stackrel{\Delta}{\leq} \|\tilde{A}_{1:k-1}(\tilde{A}_k - A_k)\|_F^2 + \|(\tilde{A}_{1:k-1} - A_{1:k-1})A_k\|_F^2 \leq \\ &\leq \underbrace{\|\tilde{A}_{1:k-1}\|_F^2}_{\leq 1; prop. 2} \|\tilde{A}_k - A_k\|_F^2 + \underbrace{\|\tilde{A}_{1:k-1} - A_{1:k-1}\|_F^2}_{(*)-induction} \underbrace{\|A_k\|_F^2}_{\leq 1} \leq \|\tilde{A}_k - A_k\|_F^2 + \sum_{i=1}^{k-1} \|\tilde{A}_i - A_i\|_F^2 = \\ &= \sum_{i=1}^k \|\tilde{A}_i - A_i\|_F^2 \end{aligned}$$

Overall:

$$\|\tilde{A}_{1:k} - A_{1:k}\|_F^2 \leq \sum_{i=1}^k \|\tilde{A}_i - A_i\|_F^2$$

Proof of proposition 2:

Proof by induction

- $k = 1$:
Immediate.

- $(k-1) \rightarrow k$:

Let $A_1, \dots, A_k \in B^{d,d}$:

$$\|A_1 \cdots A_{k-1} A_k\|_F \leq \underbrace{\|A_1 \cdots A_{k-1}\|_F}_{\leq 1; induction} \underbrace{\|A_k\|_F}_{\leq 1} \leq 1$$

So $A_1 \cdots A_k \in B^{d,d}$

Proof (B):

Assume that $\phi(\cdot)$ is Lipschitz smooth on $\mathbb{R}^{d,d} \times \dots \times \mathbb{R}^{d,d}$, i.e.

$\forall (W_1, \dots, W_N), (\widetilde{W}_1, \dots, \widetilde{W}_N) \in \mathbb{R}^{d,d} \times \dots \times \mathbb{R}^{d,d}$:

$$\|\nabla\phi(\widetilde{W}_1, \dots, \widetilde{W}_N) - \nabla\phi(W_1, \dots, W_N)\|_2 \leq \rho \|(\widetilde{W}_1, \dots, \widetilde{W}_N) - (W_1, \dots, W_N)\|_2$$

Let $(W_1, \dots, W_N), (\widetilde{W}_1, \dots, \widetilde{W}_N) \in \mathbb{R}^{d,d} \times \dots \times \mathbb{R}^{d,d}$, so:

$$\rho \|(\widetilde{W}_1, \dots, \widetilde{W}_N) - (W_1, \dots, W_N)\|_2 \geq$$

$$\begin{aligned} & \geq \|\nabla\phi(\widetilde{W}_1, \dots, \widetilde{W}_N) - \nabla\phi(W_1, \dots, W_N)\|_2 = \\ & = \sqrt{\sum_{i=1}^N \left\| \left(\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N) \right)_i \right\|_F^2} \geq \\ & \geq \left\| \left(\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N) \right)_j \right\|_F; \quad \forall j \in [N] \end{aligned}$$

So, $\forall j \in [N]$

$$\rho \|(\widetilde{W}_1, \dots, \widetilde{W}_N) - (W_1, \dots, W_N)\|_2 \geq$$

$$\begin{aligned} & \geq \left\| \left(\nabla\phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla\phi(W_1, W_2, \dots, W_N) \right)_j \right\|_F \geq \\ & \geq \left\| \widetilde{W}_{j+1:N}^T \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T - W_{j+1:N}^T \nabla l(W_{1:N}) W_{1:j-1}^T \right\|_F = \end{aligned}$$

We'll add and subtract terms:

$$= \left\| \underbrace{\left(\widetilde{W}_{j+1:N}^T - W_{j+1:N}^T \right) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T}_{(A)} + \underbrace{W_{j+1:N}^T \left(\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N}) \right) \widetilde{W}_{1:j-1}^T}_{(B)} + \underbrace{W_{j+1:N}^T \nabla l(W_{1:N}) \left(\widetilde{W}_{1:j-1}^T - W_{1:j-1}^T \right)}_{(C)} \right\|_F$$

Overall:

$\forall (W_1, \dots, W_N), (\widetilde{W}_1, \dots, \widetilde{W}_N) \in \mathbb{R}^d \times \dots \times \mathbb{R}^d$:

$$\begin{aligned} & \rho \|(\widetilde{W}_1, \dots, \widetilde{W}_N) - (W_1, \dots, W_N)\|_2 \stackrel{(*)}{\geq} \\ & \geq \left\| \underbrace{\left(\widetilde{W}_{j+1:N}^T - W_{j+1:N}^T \right) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T}_{(A)} + \underbrace{W_{j+1:N}^T \left(\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N}) \right) \widetilde{W}_{1:j-1}^T}_{(B)} + \underbrace{W_{j+1:N}^T \nabla l(W_{1:N}) \left(\widetilde{W}_{1:j-1}^T - W_{1:j-1}^T \right)}_{(C)} \right\|_F \end{aligned}$$

First, we assume $N = 2$. We would like to prove $l(\cdot)$ is affine.

So $j \in \{1,2\}$. Choose $j = 1$

- Define a series of matrices $k \in \mathbb{N}$:

$$W_1^{(k)} := \frac{1}{k} W_1, \widetilde{W}_1^{(k)} := \frac{1}{k} \widetilde{W}_1$$

Where $W_1, \widetilde{W}_1 \in \mathbb{R}^{d,d}$ are arbitrary matrices (const w.r.t k)

- Define a series of matrices $k \in \mathbb{N}$:

$$W_2^{(k)} := k \cdot I_d, \widetilde{W}_2^{(k)} := k \cdot I_d$$

Where $I_d \in \mathbb{R}^{d,d}$ is the identity matrix.

For every $k \in \mathbb{N}$, inequality (*) holds for $(W_1^{(k)}, W_2^{(k)}), (\widetilde{W}_1^{(k)}, \widetilde{W}_2^{(k)})$

So:

- $LHS = \rho \left\| (\widetilde{W}_1^{(k)}, \widetilde{W}_2^{(k)}) - (W_1^{(k)}, W_2^{(k)}) \right\|_2 = \rho \left\| \widetilde{W}_1^{(k)} - W_1^{(k)} \right\|_F = \frac{\rho}{k} \left\| \widetilde{W}_1 - W_1 \right\|_F$
- $(A) = (\widetilde{W}_2^{(k)^T} - W_2^{(k)^T}) (\nabla l(\widetilde{W}_2^{(k)} \widetilde{W}_1^{(k)})) = (kI_d - kI_d) (\nabla l(\widetilde{W}_2 \widetilde{W}_1)) = 0$
- $(C) = W_2^{(k)^T} \nabla l(W_2^{(k)} W_1^{(k)}) (I_d - I_d) = 0$
- $(B) = W_2^{(k)^T} (\nabla l(\widetilde{W}_2^{(k)} \widetilde{W}_1^{(k)}) - \nabla l(W_2^{(k)} W_1^{(k)})) = k \cdot \left(\nabla l(k \cdot \frac{1}{k} \cdot \widetilde{W}_1) - \nabla l(k \cdot \frac{1}{k} \cdot W_1) \right) = k \cdot (\nabla l(\widetilde{W}_1) - \nabla l(W_1))$

So:

$$\underbrace{\frac{\rho}{k} \left\| \widetilde{W}_1 - W_1 \right\|_F}_{\xrightarrow[k \rightarrow \infty]{}} \geq \underbrace{\sum_{k=1}^{\infty} \left\| \nabla l(\widetilde{W}_1) - \nabla l(W_1) \right\|_F}_{\xrightarrow[k \rightarrow \infty]{}}$$

$\Rightarrow \left\| \nabla l(\widetilde{W}_1) - \nabla l(W_1) \right\|_F \xrightarrow[k \rightarrow \infty]{} 0$, but $\left\| \nabla l(\widetilde{W}_1) - \nabla l(W_1) \right\|_F$ is constant w.r.t k

so $\left\| \nabla l(\widetilde{W}_1) - \nabla l(W_1) \right\|_F = 0 \forall \widetilde{W}_1, W_1 \in \mathbb{R}^{d,d}$

$\Rightarrow \nabla l(\cdot) \equiv \text{const} \Rightarrow l(\cdot) \text{ is affine.}$

Now, we assume $N > 2$. We would like to prove $\ell(\cdot)$ is constant.

$j \in [N]$. Choose $j = 2$

$$\begin{aligned} \rho \|(\widetilde{W}_1, \dots, \widetilde{W}_N) - (W_1, \dots, W_N)\|_2 &\stackrel{(*)}{\geq} \\ &\geq \left\| \underbrace{(\widetilde{W}_{j+1:N}^T - W_{j+1:N}^T) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T}_{(A)} + \underbrace{W_{j+1:N}^T (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_{1:j-1}^T}_{(B)} + \underbrace{W_{j+1:N}^T \nabla l(W_{1:N}) (\widetilde{W}_{1:j-1}^T - W_{1:j-1}^T)}_{(C)} \right\|_F = \\ &= \left\| \underbrace{(\widetilde{W}_{3:N}^T - W_{3:N}^T) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_1^T}_{(A)} + \underbrace{W_{3:N}^T (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_1^T}_{(B)} + \underbrace{W_{3:N}^T \nabla l(W_{1:N}) (\widetilde{W}_1^T - W_1^T)}_{(C)} \right\|_F \end{aligned}$$

- Define a series of matrices $k \in \mathbb{N}$:

$$\widetilde{W}_1^{(k)} = k^{N-2} I_d \quad ; \quad W_1^{(k)} = k^{N-2} I_d$$

Where I_d is the $d \times d$ identity matrix

- Define a series of matrices $k \in \mathbb{N}$:

$$\widetilde{W}_2^{(k)} = W_2 \quad ; \quad W_2^{(k)} = W_2$$

Where $W_2 \in \mathbb{R}^{d,d}$ is an arbitrary matrix

- For all $3 \leq i \leq N$, define a series of matrices $k \in \mathbb{N}$:

$$\widetilde{W}_i^{(k)} = \frac{1}{k} I_d \quad ; \quad W_i^{(k)} = 0$$

For every $k \in \mathbb{N}$, inequality $(*)$ holds for $(W_1^{(k)}, \dots, W_N^{(k)})$, $(\widetilde{W}_1^{(k)}, \dots, \widetilde{W}_N^{(k)})$

So:

- $LHS = \rho \|(\widetilde{W}_1^{(k)}, \dots, \widetilde{W}_N^{(k)}) - (W_1^{(k)}, \dots, W_N^{(k)})\|_2 =$
 $= \rho \|((0, 0, \widetilde{W}_3^{(k)} - W_3^{(k)}, \dots, \widetilde{W}_N^{(k)} - W_N^{(k)})\|_2 = \rho \|(0, 0, \frac{1}{k} I_d, \dots, \frac{1}{k} I_d)\|_2 = \frac{\rho}{k} (N-2) \|I_d\|_F =$
 $= \frac{\rho(N-2)d}{k}$
- $(A) = (\widetilde{W}_{3:N}^{(k)T} - W_{3:N}^{(k)T}) \nabla l(\widetilde{W}_{1:N}^{(k)}) \widetilde{W}_1^{(k)T} = \left(\frac{1}{k^{N-2}} I_d - 0\right) \nabla l(\widetilde{W}_{3:N}^{(k)} \widetilde{W}_2^{(k)} \widetilde{W}_1^{(k)}) k^{N-2} I_d =$
 $= \nabla l\left(\frac{1}{k^{N-2}} I_d W_2 k^{N-2} I_d\right) = \nabla l(W_2)$
- $(B) = W_{3:N}^T (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_1^T = 0 (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_1^T = 0$
- $(C) = W_{3:N}^T \nabla l(W_{1:N}) (\widetilde{W}_1^T - W_1^T) = 0 \nabla l(W_{1:N}) (\widetilde{W}_1^T - W_1^T) = 0$

So,

$$\underbrace{\frac{\rho(N-2)d}{k}}_{\substack{k \rightarrow \infty \\ \rightarrow 0}} \geq \|\nabla l(W_2)\|_F$$

$\Rightarrow \|\nabla l(W_2)\|_F \xrightarrow{k \rightarrow \infty} 0$, but $\|\nabla l(W_2)\|_F$ is constant w.r.t k so $\|\nabla l(W_2)\|_F = 0 \forall W_2 \in \mathbb{R}^{d,d}$
 $\Rightarrow \|\nabla l(\cdot)\|_F \equiv 0 \Rightarrow l(\cdot)$ is const.

It remains to prove the opposite direction:

- (1) $l(\cdot)$ is constant $\Rightarrow \phi(\cdot)$ is smooth
- (2) $l(\cdot)$ is affine and $N = 2 \Rightarrow \phi(\cdot)$ is smooth

Proof:

- (1) $l(\cdot) = \text{const} \Rightarrow \nabla l(\cdot) \equiv 0$
- (2) $l(\cdot)$ is affine $\Rightarrow \nabla l(\cdot) \equiv 0$

In either case $\nabla l(\cdot) \equiv 0$. This will be the only assumption we need.

We saw in “Proof (A)”:

$$\forall (W_1, \dots, W_N), (\widetilde{W}_1, \dots, \widetilde{W}_N) \in \mathbb{R}^{d,d} \times \dots \times \mathbb{R}^{d,d}:$$

$$\begin{aligned} & \| \nabla \phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla \phi(W_1, W_2, \dots, W_N) \|_2^2 \leq \\ & \leq N \cdot \left(\underbrace{\| (\widetilde{W}_{j+1:N}^T - W_{j+1:N}^T) \nabla l(\widetilde{W}_{1:N}) \widetilde{W}_{1:j-1}^T \|_F^2}_{(A)} + \underbrace{\| W_{j+1:N}^T (\nabla l(\widetilde{W}_{1:N}) - \nabla l(W_{1:N})) \widetilde{W}_{1:j-1}^T \|_F^2}_{(B)} \right. \\ & \quad \left. + \underbrace{\| W_{j+1:N}^T \nabla l(W_{1:N}) (\widetilde{W}_{1:j-1}^T - W_{1:j-1}^T) \|_F^2}_{(C)} \right) = \\ & = N \cdot (0 + 0 + 0) = 0 \end{aligned}$$

So in particular, selecting $\rho = 1$:

$$\forall (W_1, \dots, W_N), (\widetilde{W}_1, \dots, \widetilde{W}_N) \in \mathbb{R}^{d,d} \times \dots \times \mathbb{R}^{d,d}:$$

$$(0 =) \| \nabla \phi(\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - \nabla \phi(W_1, W_2, \dots, W_N) \|_2^2 \leq \rho^2 \| (\widetilde{W}_1, \widetilde{W}_2, \dots, \widetilde{W}_N) - (W_1, W_2, \dots, W_N) \|_2^2$$

$\Rightarrow \phi(\cdot)$ is ρ -smooth on $\mathbb{R}^{d,d} \times \dots \times \mathbb{R}^{d,d}$.

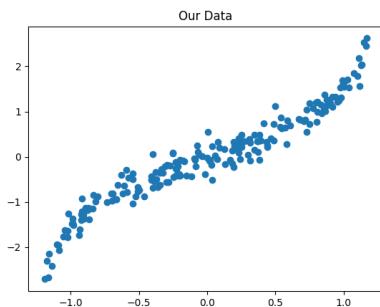
■

Question 4:

We used in our experiment a LNN network with the following properties:

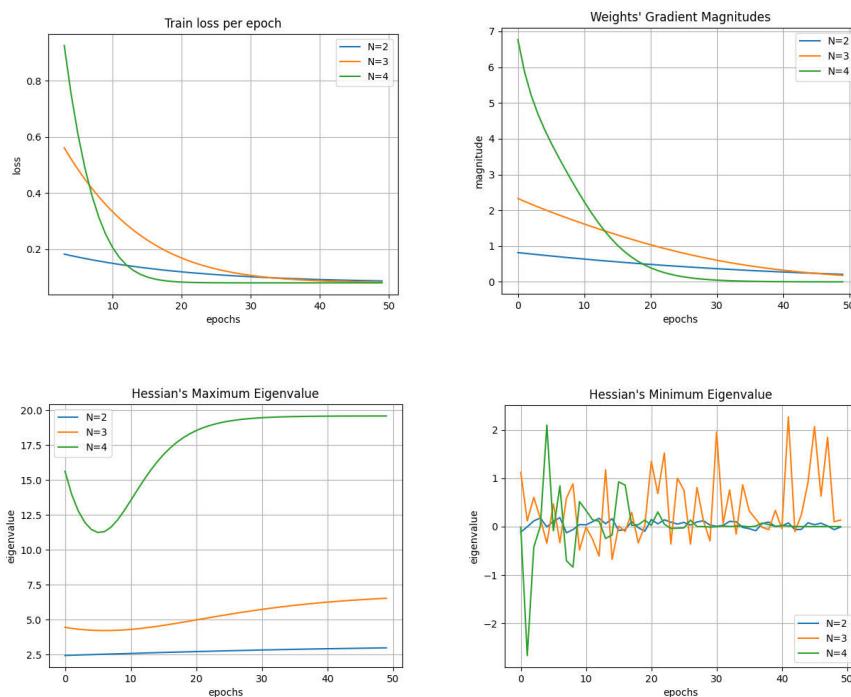
1. Input dimension = 1
2. Output dimensions = 1
3. Width of hidden layers = 3
4. Depth of the LNN: N=2,3,4

We ran our experiments on the following dataset:



We chose uniformly 200 random points between -1.2 and 1.2, we calculate for each point the function $\tan(x)$ and added noise which was chosen normally around 0 with standard deviation of 0.2.

The results of our experiment are the followings:



Those are the results in each epoch while running of train loss, magnitude of weights' gradient, maximal value of Hessian of the weights and minimum value of Hessian of the weights.

Explanation of the results:

1. No bad local minima – as we proved in the lecture of optimization 1, for LNN networks without a “bottleneck” layers (which is the case in our case) every minima is a global one for the objective function which induced by the loss of the network.
We can clearly see that result by the figure of train loss per epoch as all the values of the loss (for N=2, 3, 4 networks) converge to the same value (and not stay stuck in a local minima which is not global).
2. From the proposition in class where we vectorized the end-to-end matrix and wrote the dynamics with the preconditioned PSD matrix $P_{W_{1:N}(t)}$, we concluded that the end-to-end dynamics are in fact an instance of preconditioned GF, Thus overparameterizing the loss with a LNN “promotes movement in directions already taken” – which means stretches the gradient in singular directions of large singular values and attenuates it in directions of small singular values.
 - a. We can clearly see the results we derived in class here in our results - Note that in the figure of maximal eigenvalue of the Hessian, as N is bigger; the maximal value of the Hessian is bigger. Thus, it implies on the movement of the gradient such that as N is bigger the weights converge **faster** to the global min, as we can see in the figure of the train loss (as N is bigger, it converge faster). In addition the magnitude of the gradient decreases faster as N is bigger as we can see in our results.
 - b. In addition we can see (in the figure with the maximal and minimal eigenvalues of the Hessian) that in convergence the minimal eigenvalue is bigger than 0 (very close but still bigger), and thus we converge to a Hessian which is PSD so we really get closer to a local minima point (also a global as we mentioned in 1) as the point is a stationary point.
3. From the lecture of optimization 1, we proved that for LNN with $depth = 2$ any stationary point which is not a global minimizer - a strict saddle point.
In addition we proved that for LNNs with $depth \geq 3$ the objective function doesn't have any non-strict saddle, and that GD escapes a strict saddle point if it arrives to it.
We clearly can see from the figure of the training loss that since all the networks converge to the same value - it strengthens the proof of escaping strict saddle points for the networks such that they will arrive to a global min.

Part 3 - Trajectory approach

Question 1

We have to show that under the conditions of the respective theorem:

$$W_{1:j}(t)^T W_{1:j}(t) = [W_{1:N}(t)^T W_{1:N}(t)]^{\frac{j}{N}}$$

$$\forall t \in \mathbb{R}_{\geq 0}, j \in [N]$$

Proof

By definition: $W_{1:j}(t) = W_j(t) \cdots W_1(t)$

As we proved in lecture of optimization 2:

$\forall j \in [N-1] \exists V_j \in \mathbb{R}^{d \times d}$ orthogonal and $\exists D_j \in \mathbb{R}^{d \times d}$ diagonal matrix whose diagonal values are in $\{+1, -1\}$

such that $W_j(t) = V_{j+1} D_j \sum V_j^T$.

For fixed $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_d) \in \mathbb{R}^{d \times d}$
 $\forall i \sigma_i \geq 0$

Therefore,

$$W_{1:j}(t) = W_j(t) \cdots W_1(t) = \underbrace{V_{j+1} D_j \sum V_j^T}_{\Sigma} \underbrace{V_j D_{j-1} \sum V_{j-1}^T}_{\Sigma} \cdots \underbrace{V_2 D_1 \sum V_1^T}_{\Sigma}$$

\Rightarrow

$$W_{1:j}(t)^T W_{1:j}(t) = (V_{j+1} D_j \sum V_j^T \cdots V_2 D_1 \sum V_1^T)^T (V_{j+1} D_j \sum V_j^T \cdots V_2 D_1 \sum V_1^T)$$

$$= (V_1 \sum D_1^T V_2^T \cdots V_j \sum D_j^T V_{j+1}^T) (V_{j+1} D_j \sum V_j^T \cdots V_2 D_1 \sum V_1^T)$$

Since the matrices Σ, D_1, \dots, D_N are diagonal, it means

$$\forall j \quad D_j^T = D_j \quad , \quad \Sigma^T = \Sigma.$$

$$\Rightarrow W_{1:j}(t)^T W_{1:j}(t) = \\ (V_1 \Sigma D_1 V_2^T V_2 \Sigma D_2 V_3^T V_3 \cdots V_j^T V_j \Sigma D_j V_{j+1}^T) (V_{j+1} D_j \Sigma V_j^T \cdots V_2^T V_2 D_1 \Sigma V_1^T)$$

$$= (V_1 \Sigma D_1 \Sigma D_2 \cdots \Sigma D_j) (D_j \Sigma D_{j-1} \cdots D_1 \Sigma V_1^T)$$

$$\forall j \quad V_j^T V_j = I$$

$$= V_1 \sum^{2j} D_1^2 \cdots D_j^2 V_1^T$$

D_1, \dots, D_j, Σ
are diagonal and therefore commutative w.r.t. matrix multiplication

$$\Rightarrow W_{1:j}(t)^T W_{1:j}(t) = V_1 \sum^{2j} D_1^2 \cdots D_j^2 V_1^T = V_1 \sum^{2j} V_1^T$$

(since $D_i^2 = I$ $\forall i$ because its diagonal matrix with values ± 1 on diagonal)

$$\text{In particular } W_{1:N}(t)^T W_{1:N}(t) = V_1 \sum^{2N} V_1^T$$

$\Rightarrow \forall \epsilon \in \mathbb{R}_{\geq 0}, \quad j \in [N]:$

$$W_{1:j}(t)^T W_{1:j}(t) = \left[W_{1:N}(t)^T W_{1:N}(t) \right]^{\frac{j}{N}}$$

□

Part 3 - question 2

Let $\varrho: \mathbb{R}^{d,d} \rightarrow \mathbb{R}$ be a continuously differentiable loss parameterized by: $\phi: \mathbb{R}^{d,d} \rightarrow \mathbb{R}$

$$\phi(U) = \varrho(UU^\top)$$

$$W(\epsilon) := U(\epsilon)U(\epsilon)^\top.$$

Develop an expression for the dynamics of $W(\epsilon)$ that are induced by gradient flow over $\phi(\cdot)$:

$$\dot{U}(t) = \frac{d}{dt} U(t) = -\nabla \phi(U(t))$$

Solution

We'll calculate the dynamics for $W(\epsilon)$:

$$\begin{aligned}\dot{W}(\epsilon) &= \frac{d}{d\epsilon} W(\epsilon) = \frac{d}{d\epsilon} (U(\epsilon)U(\epsilon)^\top) = (\dot{U}(t)U(t)^\top + U(t)\dot{U}(t)^\top) \\ &= \dot{U}(t)U(t)^\top + U(t)\dot{U}(t)^\top = \\ &= -\nabla \phi(U(t))U(t)^\top + U(t) \cdot (-\nabla \phi(U(t)))^\top \\ &= -\nabla \phi(U(t))U(t)^\top - U(t)\nabla \phi(U(t))^\top \quad (*)\end{aligned}$$

Now, we would like to continue develop the expression by develop the term $\nabla \phi(U(t))$ using the first order Taylor expansions of ϕ and ϱ .

As we know: $\phi(U) = \epsilon(UU^T)$

First order Taylor expansion of ϕ around U :

$$\phi(U + \Delta) = \phi(U) + \langle \nabla \phi(U), \Delta \rangle + o(\|\Delta\|_{Fro})$$

We also know that:

$$\begin{aligned} \phi(U + \Delta) &= \epsilon((U + \Delta)(U + \Delta)^T) = \\ &= \epsilon(U + \Delta)(U^T + \Delta^T) = \\ &= \epsilon(UU^T + U\Delta^T + \Delta U^T + \Delta\Delta^T) \end{aligned}$$

Develop the first order Taylor expansion of ϵ around UU^T will yield:

$$\begin{aligned} \epsilon(UU^T + \underline{U\Delta^T + \Delta U^T + \Delta\Delta^T}) &= \\ &= \epsilon(UU^T) + \langle \nabla \epsilon(UU^T), U\Delta^T + \Delta U^T + \Delta\Delta^T \rangle + o(\|\Delta\|_{Fro}) \\ &= \epsilon(UU^T) + \langle \nabla \epsilon(UU^T), U\Delta^T + \Delta U^T \rangle \\ &\quad + \langle \nabla \epsilon(UU^T), \Delta\Delta^T \rangle + o(\|\Delta\|_{Fro}) \end{aligned}$$

We've got: $\phi(U + \Delta) = \phi(U) + \langle \nabla \phi(U), \Delta \rangle + o(\|\Delta\|_{Fro})$

$$\begin{aligned} \epsilon(UU^T + U\Delta^T + \Delta U^T + \Delta\Delta^T) &= \\ &= \epsilon(UU^T) + \langle \nabla \epsilon(UU^T), U\Delta^T + \Delta U^T \rangle \\ &\quad + \langle \nabla \epsilon(UU^T), \Delta\Delta^T \rangle + o(\|\Delta\|_{Fro}) \end{aligned}$$

Since $\phi(U) = \epsilon(UU^T)$, $\phi(U + \Delta) = \epsilon(UU^T + U\Delta^T + \Delta U^T + \Delta\Delta^T)$

and the uniqueness of first order Taylor expansion
we conclude:

$$\langle \nabla \phi(U), \Delta \rangle = \langle \nabla \epsilon(UU^T), U\Delta^T + \Delta U^T \rangle$$

$$\text{We've got: } \langle \nabla \phi(u), \Delta \rangle = \langle \nabla e(uu^T), u\Delta^T + \Delta u^T \rangle \quad (\star)$$

$$\Leftrightarrow \langle \nabla e(uu^T), u\Delta^T \rangle + \langle \nabla e(uu^T), \Delta u^T \rangle =$$

$$= \text{tr}(\nabla e(uu^T)^T u\Delta^T) + \text{tr}(\nabla e(uu^T)^T \Delta u^T) =$$

$$\text{tr}(AB^T) = \text{tr}(A^T B)$$

$$\leftarrow = \text{tr}((\nabla e(uu^T)^T u)^T \Delta) + \text{tr}(\nabla e(uu^T)^T \Delta u^T)$$

$$= \text{tr}(u^T \nabla e(uu^T) \Delta) + \text{tr}(\nabla e(uu^T)^T \Delta u^T) =$$

$$\text{tr}(ABC) = \text{tr}(ACB)$$

$$\leftarrow = \text{tr}(u^T \nabla e(uu^T) \Delta) + \text{tr}(\nabla e(uu^T)^T u^T \Delta)$$

$$\leftarrow = \text{tr}(\nabla e(uu^T) u^T \Delta) + \text{tr}(\nabla e(uu^T)^T u^T \Delta)$$

$$= \text{tr}((u \nabla e(uu^T))^T \Delta) + \text{tr}((u \nabla e(uu^T))^T \Delta)$$

$$= \langle u \nabla e(uu^T)^T, \Delta \rangle + \langle u \nabla e(uu^T), \Delta \rangle$$

$$= \langle u \nabla e(uu^T) + u \nabla e(uu^T)^T, \Delta \rangle$$

$$= \langle u(\nabla e(uu^T) + \nabla e(uu^T)^T), \Delta \rangle$$

$$\Rightarrow \nabla \phi(u) = u(\nabla e(uu^T) + \nabla e(uu^T)^T)$$

Plugging it in (\star) :

$$\dot{W}(t) = -\nabla \phi(u(t)) u(t)^T - u(t) \cdot \nabla \phi(u(t))^T =$$

$$= -u[\nabla e(uu^T) + \nabla e(uu^T)^T]u^T - u[\nabla e(uu^T) + \nabla e(uu^T)^T]^T$$

$$[u=u(t)]$$

for simplicity

$$= -u[\nabla e(uu^T) + \nabla e(uu^T)^T]u^T - u[\nabla e(uu^T) + \nabla e(uu^T)^T]^T u^T$$

$$= -u[\nabla e(uu^T) + \nabla e(uu^T)^T]u^T - u[\nabla e(uu^T) + \nabla e(uu^T)^T]u^T$$

$$= -2 \cdot u[\nabla e(uu^T) + \nabla e(uu^T)^T]u^T$$

$$\Rightarrow \dot{W}(t) = -2u(t)[\nabla e(u(t)u(t)^T) + \nabla e(u(t)u(t)^T)^T]u(t)^T$$

□

Question 3

Simplify the end-to-end dynamics delivered in class for the special case of $d_N = 1$. Explain how this resonates with our interpretation of the end-to-end dynamics "promoting movement in direction already taken"

Solution

As we've already showed the end-to-end dynamics are:

$$\textcircled{X} \quad \dot{W}_{1:N}(t) = - \sum_{j=1}^N [W_{1:N}(t) W_{1:N}(t)^T]^{\frac{j-1}{N}} \nabla (W_{1:N}(t)) [W_{1:N}(t)^T W_{1:N}(t)]^{\frac{N-j}{N}}$$

$t \in \mathbb{R} \geq 0$

By the assumption $d_N = 1$, thus $W_{1:N}(t) = W_N(t) \cdots W_1(t) \in \mathbb{R}^{1,d_0}$
 $\dim g: (1, d_{N-1}), \dots, (d_1, d_0)$

① Therefore, for each $j \in [N]$:

$$W_{1:N}(t) W_{1:N}(t)^T = \langle W_{1:N}(t), W_{1:N}(t) \rangle = \|W_{1:N}(t)\|_2^2 \in \mathbb{R}$$

$$[W_{1:N}(t) W_{1:N}(t)^T]^{\frac{j-1}{N}} = \|W_{1:N}(t)\|_2^{2 \cdot \frac{j-1}{N}}$$

② On the other hand, $[W_{1:N}(t)^T W_{1:N}(t)]^{\frac{N-j}{N}} \in \mathbb{R}^{d_0 \times d_0}$, thus:

- For $j=N$: $[W_{1:N}(t)^T W_{1:N}(t)]^{\frac{N-j}{N}} = [W_{1:N}(t)^T W_{1:N}(t)]^0 = I_{d,d}$

↑

By definition

- Claim

For $j=1, \dots, N-1$:

$$[W_{1:N}(t)^T W_{1:N}(t)]^{\frac{N-j}{N}} = \|W_{1:N}(t)\|_2^{2 \cdot \frac{N-j}{N}} \cdot \frac{W_{1:N}(t)^T}{\|W_{1:N}(t)\|_2} \cdot \frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2}$$

Proof of claim:

First, we know that $W_{1:N}(t)^T W_{1:N}(t)$ is PSD matrix and thus all its singular values are ≥ 0 . Notice that $\text{rank}(W_{1:N}(t)) = 1 \Rightarrow \text{rank}(W_{1:N}(t)^T W_{1:N}(t)) \leq 1$.

In addition $W_{1:N}(t)^T$ is an eigenvector with eigenvalue

$$\|W_{1:N}(t)\|_2^2 :$$

$$\begin{aligned} [W_{1:N}(t)^T W_{1:N}(t)] W_{1:N}(t)^T &= W_{1:N}(t)^T \cdot \|W_{1:N}(t)\|_2^2 \\ &= \|W_{1:N}(t)\|_2^2 \cdot W_{1:N}(t)^T \end{aligned}$$

Thus $\sigma_1 := \|W_{1:N}(t)\|_2^2 \neq 0$ is the only non-zero eigenvalue of $W_{1:N}(t)^T W_{1:N}(t)$. [Suppose $\|W_{1:N}(t)\|_2 = 0$]

We'll write the SVD composition of $W_{1:N}(t)^T W_{1:N}(t)$:

$$W_{1:N}(t)^T W_{1:N}(t) = U \Sigma V^T = U \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \end{pmatrix} V^T$$

Thus for every $x \geq 0$:

$$\begin{aligned} [W_{1:N}(t)^T W_{1:N}(t)]^x &= U \Sigma^x V^T = U \begin{pmatrix} \sigma_1^x & 0 & \dots & 0 \end{pmatrix} V^T \\ &= \sigma_1^x \cdot \underbrace{U \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \end{pmatrix} V^T}_{\sqrt{\sigma_1} \quad \sqrt{\sigma_1}} = \cancel{U \cancel{\Sigma} \cancel{V^T}}^{\cancel{\sigma_1^x}} \cancel{\sigma_1^x} \\ &= \|W_{1:N}(t)\|_2^{2x} \cdot \frac{W_{1:N}(t)^T W_{1:N}(t)}{\sqrt{\sigma_1} \cdot \sqrt{\sigma_1}} \\ &= \|W_{1:N}(t)\|_2^{2x} \cdot \left(\frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2} \right)^T \cdot \left(\frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2} \right) \end{aligned}$$

We can ~~also~~ use $x = \frac{N-j}{N}$ and get the required.

Now, we can use the equations ① and ② into ⑧ and we'll get:

$$\begin{aligned}
 W_{1:N}(t) &= -\sum_{j=1}^N [W_{1:N}(t) W_{1:N}(t)^T]^{\frac{j-1}{N}} \nabla \varphi(W_{1:N}(t)) [W_{1:N}(t)^T W_{1:N}(t)]^{\frac{N-j}{N}} \\
 &= -\sum_{j=1}^{N-1} \|W_{1:N}(t)\|_2^{2 \cdot \frac{j-1}{N}} \nabla \varphi(W_{1:N}(t)) \cdot \|W_{1:N}(t)\|_2^{2 \cdot \frac{N-j}{N}} \frac{W_{1:N}(t)^T}{\|W_{1:N}(t)\|_2} \cdot \frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2} \\
 &\quad - \|W_{1:N}(t)\|_2^{2 \cdot \frac{N-1}{N}} \nabla \varphi(W_{1:N}(t)) \cdot I_{d,d} \\
 &= -\sum_{j=1}^{N-1} \|W_{1:N}(t)\|_2^{2 \cdot \frac{j-1}{N}} \cdot \nabla \varphi(W_{1:N}(t)) \frac{W_{1:N}(t)^T}{\|W_{1:N}(t)\|_2} \cdot \frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2} \\
 &\quad - \|W_{1:N}(t)\|_2^{2 \cdot \frac{N-1}{N}} \cdot \nabla \varphi(W_{1:N}(t)) \\
 &= -\|W_{1:N}(t)\|_2^{2-\frac{2}{N}} \left[\nabla \varphi(W_{1:N}(t)) + (N-1) \cdot \nabla \varphi(W_{1:N}(t)) \frac{W_{1:N}(t)^T}{\|W_{1:N}(t)\|_2} \cdot \frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2} \right]
 \end{aligned}$$

Where we assumed $W_{1:N}(t) \neq \vec{0}$ and $N \geq 2$ on our proof.

We've got: $\dot{W}_{1:N}(t) = -\|W_{1:N}(t)\|_2^{2-\frac{2}{N}} \left[\nabla \varphi(W_{1:N}(t)) + (N-1) \cdot \nabla \varphi(W_{1:N}(t)) \frac{W_{1:N}(t)^T}{\|W_{1:N}(t)\|_2} \cdot \frac{W_{1:N}(t)}{\|W_{1:N}(t)\|_2} \right]$

Thus, besides the multiplication by $\|W_{1:N}(t)\|_2^{2-\frac{2}{N}}$, we also add, to $\nabla \varphi(W_{1:N}(t))$ the projection of itself on $W_{1:N}(t)$.

Recall that we view $W_{1:N}(t)$ not only as the parameters but also as the overall movement of the system.

Thus by adding the projection of the gradient on $W_{1:N}(t)$ we can think about it as promoting movements in directions which already taken.



Foundations of Deep Learning – Homework Assignment #3

Adi Almog & Tomer Epshtain

Part 3: (4)

Question:

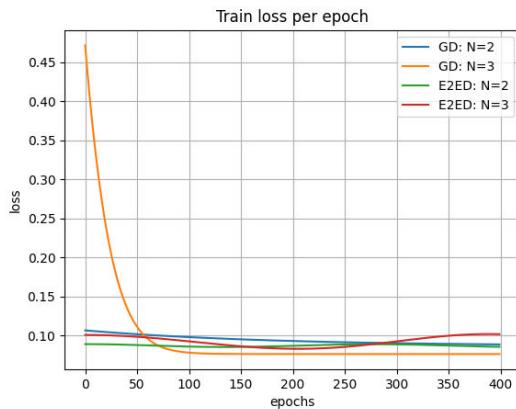
On a scalar regression dataset of your choice, train a depth N linear neural network (with hidden widths no smaller than the minimum between input dimension and output dimension) by minimizing ℓ_2 loss via (full batch) gradient descent with small learning rate and initialization close to zero. Compare the trajectory taken by the end-to-end matrix to that obtained by directly applying the (discrete version of the) end-to-end dynamics to a linear model:

$$W_{t+1} \leftarrow W_t - \eta \sum_{j=1}^N [W_t W_t^T]^{\frac{j-1}{N}} \nabla l(W_t) [W_t^T W_t]^{\frac{N-j}{N}} ; \quad t = 1, 2, 3, \dots$$

Where η is the learning rate used for gradient descent over the linear neural network. Repeat the experiment with depths $N = 2$ and 3

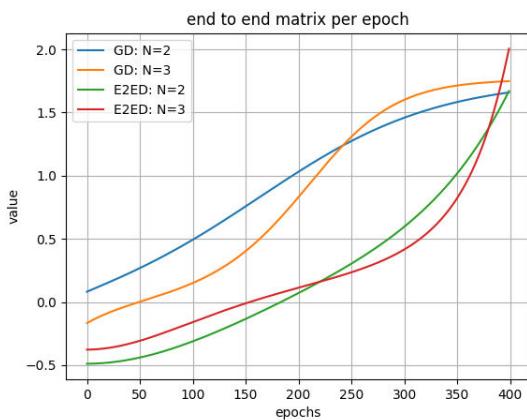
Solution:

The training procedure can be viewed in the following plots:



Convergence of all four models: GD with $N = 2, 3$ and E2ED (End-to-end dynamics) with $N = 2, 3$ is achieved.

Additionally, because our network is a Linear Neural Network with input dimension 1 and output dimension 1, the end-to-end matrix is a 1×1 matrix. I.e. a single real value. We can see convergence to global minimum, with different learning rates, in the following graph, where all four models converge to an end-to-end matrix $\approx (1.7)_{1 \times 1}$



Part 4 - Ultra wide neural networks

Question 1

$$\dot{u}(t) = -H^*(u(t) - y) \quad t \in \mathbb{R}_{\geq 0}$$

Where:

- $u(t) \in \mathbb{R}^m$ holds the neural network predictions at time t .
- $H^* \in \mathbb{R}^{m \times m}$ is the Neural Tangent Kernel's Gram matrix PD with eigenvalues $\geq \lambda_{\min} > 0$
- $y \in \mathbb{R}^m$ holds training labels.

We have to show that $u(t) \rightarrow y$ exponentially fast

Without using the change of variables as done in class.

Solution

We'll develop the expression $\frac{d}{dt} \|u(t) - y\|^2$:

$$\begin{aligned}\frac{d}{dt} \|u(t) - y\|^2 &= 2 \cdot (u(t) - y)^T \cdot \frac{d}{dt} u(t) = \\ &= 2 \cdot (u(t) - y)^T \cdot \dot{u}(t) \\ &= -2 (u(t) - y)^T H^* (u(t) - y)\end{aligned}$$

[Notice that $\lambda_{\min}(H^*) = \min_{x \neq 0 \in \mathbb{R}^m} \frac{x^T H^* x}{\|x\|^2} \Rightarrow \forall x \quad x^T H^* x \geq \lambda_{\min} \cdot \|x\|^2$]

$$\Rightarrow \frac{d}{dt} \|u(t) - y\|^2 = -2 (u(t) - y)^T H^* (u(t) - y) \leq -2 \lambda_{\min} \cdot \|u(t) - y\|^2$$

$[x = u(t) - y]$

$$\Rightarrow \frac{\frac{d}{dt} \|u(t) - y\|^2}{\|u(t) - y\|^2} \leq -2 \lambda_{\min}$$

Thus, we can write it as:

$$\frac{d}{dt} \ln(\|u(t) - y\|^2) = \frac{\frac{d}{dt} \|u(t) - y\|^2}{\|u(t) - y\|^2} \leq -2\lambda_{\min}$$

∴

$$\int_0^t \frac{d}{ds} \ln(\|u(s) - y\|^2) \leq \int_0^t -2\lambda_{\min}$$

∴

$$\ln(\|u(s) - y\|^2) [s=t] - \ln(\|u(s) - y\|^2) [s=0] \leq -2\lambda_{\min} t$$

$$\ln\left(\frac{\|u(t) - y\|^2}{\|u(0) - y\|^2}\right) \leq -2\lambda_{\min} t \quad / e^x \text{ is monotonic}$$

∴

$$\frac{\|u(t) - y\|^2}{\|u(0) - y\|^2} \leq e^{-2\lambda_{\min} t}$$

∴

$$\|u(t) - y\|^2 \leq \|u(0) - y\|^2 \cdot e^{-2\lambda_{\min} t}$$

⇒ $u(t) \rightarrow y$ exponentially fast as required.

□

Scanned with CamScanner

Question 2

The dynamics in question (1) (Part 4) are actually an idealization of: $\dot{u}(t) = -H(t)(u(t) - y) \quad t \in \mathbb{R}_{\geq 0} \quad \textcircled{4}$

Where $H(t) \in \mathbb{R}^{mm}$ satisfies $\|H(t) - H^*\|_{\text{Frob}} \leq \epsilon$.

Assuming identical initialization $u(0)$, show that under $\textcircled{4}$

the value of $u(t)$ is at most $O(\sqrt{\epsilon t})$ away (in Euclidean distance) from what it would be under question (1) dynamics.

Solution

① Recall that we proved in question (1) that under the dynamics $\dot{u}_1(t) = -H^*(u_1(t) - y)$ we got:

$$\|u_1(t) - y\|^2 \leq \|u_1(0) - y\|^2 \cdot e^{-2\lambda_{\min} t}$$

$$\Rightarrow \|u_1(t) - y\|^2 \leq \|u_1(0) - y\|^2 \cdot e^{-2\lambda_{\min} t}$$

② We'll prove that under the dynamics $\dot{u}_2(t) = -H(t)(u_2(t) - y)$ we'll get:

$$\|u_2(t) - y\|^2 \leq \|u_2(0) - y\|^2 \cdot e^{-2(\lambda_{\min} + \epsilon)t}$$

Proof

$$\text{Note: } \frac{d}{dt} \|u_2(t) - y\|^2 = (u_2(t) - y)^T \frac{d}{dt} (u_2(t)) =$$

$$= -(u_2(t) - y)^T H(t) (u_2(t) - y)$$

$$\Rightarrow \frac{d}{dt} \|u_2(t) - y\|^2 = -(u_2(t) - y)^T H(t) (u_2(t) - y)$$

By the assumption $\forall \epsilon \quad \|H(\epsilon) - H^*\|_{\text{spec}} \leq \epsilon$.

Since $\|\cdot\|_{\text{spec}}$ is a norm then we conclude:

$$\|H(\epsilon)\|_{\text{spec}} = \|H(\epsilon) - H^* + H^*\|_{\text{spec}}$$

$$\Delta \text{ inequality} \quad \geq \|H^*\|_{\text{spec}} - \|H(\epsilon) - H^*\|_{\text{spec}}$$

$$\text{by definition } H^* \leq \lambda_{\min} = \|H(\epsilon) - H^*\|_{\text{spec}}$$

PSD so $\|H^*\|_{\text{spec}} \geq \lambda_{\min}$

$$\geq \lambda_{\min} - \epsilon$$

\Rightarrow

$$\frac{d}{dt} \|u_2(t) - y\|^2 = -2(u_2(t) - y)^T H(\epsilon)(u_2(t) - y) \quad [\forall X^T A X \leq \|A\|_{\text{spec}} \cdot \|X\|^2]$$

$$\leq -2\|H(\epsilon)\|_{\text{spec}} \cdot \|u_2(t) - y\|^2$$

$$\leq -2(\lambda_{\min} - \epsilon) \|u_2(t) - y\|^2$$

$$\Rightarrow \frac{\frac{d}{dt} \|u_2(t) - y\|^2}{\|u_2(t) - y\|^2} \leq -2(\lambda_{\min} - \epsilon)$$

\Rightarrow With exactly the same steps as in question 1

$$\text{where we proved } \frac{\frac{d}{dt} \|u_2(t) - y\|^2}{\|u_2(t) - y\|^2} \leq -2\lambda_{\min} \Rightarrow \|u_2(t) - y\|^2 \leq \|u_2(0) - y\|^2 \cdot e^{-2\lambda_{\min} t}$$

we'll get the required:

$$\|u_2(t) - y\|^2 \leq \|u_2(0) - y\|^2 \cdot e^{2(\lambda_{\min} - \epsilon)t} = \|u_2(0) - y\|^2 \cdot e^{2(\lambda_{\min} - \epsilon)t}$$

as we claimed.

Therefore, from ① and ② we've got:

$$\|U_1(t) - y\|^2 \leq \|U(0) - y\|^2 \cdot e^{-2\lambda_{\min} t}$$

$$\|U_2(t) - y\|^2 \leq \|U(0) - y\|^2 \cdot e^{-2(\lambda_{\max} - \varepsilon)t}$$

Now:

$$\|U_1(t) - U_2(t)\|_2^2 = \|U_1(t) - y - (U_2(t) - y)\|^2$$

$$\begin{aligned} \left[\|U - V\|^2 \leq 2\|U\|^2 + 2\|V\|^2 \right] &\leq 2 \cdot \|U_1(t) - y\|^2 + 2 \cdot \|U_2(t) - y\|^2 \\ &\leq 2 \cdot \|U(0) - y\|^2 \cdot e^{-2\lambda_{\min} t} + 2 \cdot \|U(0) - y\|^2 \cdot e^{-2(\lambda_{\max} - \varepsilon)t} \\ &= 2 \|U(0) - y\|^2 \cdot (e^{-2\lambda_{\min} t} + e^{-2(\lambda_{\max} - \varepsilon)t}) \end{aligned}$$

$$e^{-2\lambda_{\min} t} \leq e^0 = 1 \leq 2 \cdot \|U(0) - y\|^2 \cdot (1 + e^{-2(\lambda_{\max} - \varepsilon)t})$$

$\forall t \geq 0$ since $\lambda_{\min} > 0$

$$e^{-2(\lambda_{\max} - \varepsilon)t} \leq 1 \quad \leftarrow \leq 2 \cdot \|U(0) - y\|^2 \cdot (1 + 1)$$

for $\varepsilon < \lambda_{\min}$

$$\varepsilon t \geq 1 \quad \leftarrow \leq 2 \cdot \|U(0) - y\|^2 \cdot (\varepsilon t + \varepsilon t) = 4 \cdot \|U(0) - y\|^2 \varepsilon t$$

for $t \geq \frac{1}{\varepsilon}$

$$\Rightarrow \text{for } \varepsilon < \lambda_{\min}(H^*) \text{, } t \geq \frac{1}{\varepsilon} = t_0$$

$$\|U_1(t) - U_2(t)\| \leq \sqrt{\mu \cdot \|U(0) - y\|^2} \cdot \varepsilon t = 2 \cdot \|U(0) - y\| \cdot \sqrt{\varepsilon t}$$

$$\Rightarrow \|U_1(t) - U_2(t)\| \leq \mathcal{O}(\sqrt{\varepsilon t})$$

as required.

□