

Parallelizing Stochastic Gradient Descent for Least Squares Regression: Mini-batching, Averaging, and Model Misspecification

Authors :-

- Prateek Jain, Praneeth Netrapalli
- Sham M.Kakade
- Rahul Kidambi
- Aaron Sidford

Angad Singh(22B1211)

PROBLEM SETUP

We consider the standard Linear Square Regression (LSR) problem with expected loss:

$$L(\boldsymbol{w}) = \frac{1}{2} \mathbb{E} [(y - \langle \boldsymbol{w}, \boldsymbol{x} \rangle)^2],$$

and let \boldsymbol{w}^* be the minimizer of $L(\boldsymbol{w})$. The corresponding population Hessian is $\boldsymbol{H} = \mathbb{E}[\boldsymbol{x}\boldsymbol{x}^\top]$, and the loss admits the quadratic form:

$$L(\boldsymbol{w}) = L(\boldsymbol{w}^*) + \frac{1}{2} (\boldsymbol{w} - \boldsymbol{w}^*)^\top \boldsymbol{H} (\boldsymbol{w} - \boldsymbol{w}^*).$$

ALGORITHM1:-MINIBATCH-TAIL AVERAGING-SGD

Input: Initial point \mathbf{w}_0 , stepsize γ , minibatch size b , initial iterations s , total samples n .

1: **for** $t = 1, 2, \dots, \lfloor \frac{n}{b} \rfloor$ **do**

2: Sample “ b ” tuples $\{(x_{ti}, y_{ti})\}_{i=1}^b \sim \mathcal{D}^b$

3: $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \frac{\gamma}{b} \sum_{i=1}^b \widehat{\nabla} L_{ti}(\mathbf{w}_{t-1})$

Output: $\bar{\mathbf{w}} = \frac{1}{\lfloor \frac{n}{b} \rfloor - s} \sum_{i>s} \mathbf{w}_i$

ASSUMPTIONS

- (A1) **Finite fourth moment:** The fourth moment tensor $\mathcal{M} = \mathbb{E} [\mathbf{x}^{\otimes 4}]$ exists and is finite.
- (A2) **Strong convexity:** The Hessian of $L(\cdot)$, $\mathbf{H} = \mathbb{E} [\mathbf{x}\mathbf{x}^\top]$ is positive definite i.e., $\mathbf{H} \succ 0$.

THEOREM

Theorem 1 Consider the general mis-specified case of the LSR problem 1. Running Algorithm 1 with a batch size $b \geq 1$, step size $\gamma \leq \gamma_{b,\max}/2$, number of unaveraged iterations s , total number of samples n , we obtain an iterate $\bar{\mathbf{w}}$ satisfying the following excess risk bound:

$$\mathbb{E} [L(\bar{\mathbf{w}})] - L(\mathbf{w}^*) \leq \frac{2}{\gamma^2 \mu^2} \cdot \frac{(1 - \gamma \mu)^s}{\left(\frac{n}{b} - s\right)^2} \cdot (L(\mathbf{w}_0) - L(\mathbf{w}^*)) + 4 \cdot \frac{\widehat{\sigma_{MLE}^2}}{b \cdot \left(\frac{n}{b} - s\right)}. \quad (4)$$

In particular, with $\gamma = \gamma_{b,\max}/2$, we have the following excess risk bound:

$$L(\bar{\mathbf{w}}) - L(\mathbf{w}^*) \leq \underbrace{\frac{2\kappa_b^2}{\left(\frac{n}{b} - s\right)^2} \exp\left(-\frac{s}{\kappa_b}\right) (L(\mathbf{w}_0) - L(\mathbf{w}^*))}_{\mathfrak{I}_1} + 4 \cdot \underbrace{\frac{\widehat{\sigma_{MLE}^2}}{b\left(\frac{n}{b} - s\right)}}_{\mathfrak{I}_2},$$

$$\text{with } \kappa_b = \frac{R^2 \cdot \rho_m + (b-1) \|\mathbf{H}\|_2}{b \lambda_{\min}(\mathbf{H})}.$$

THEOREM(ZEROth ORDER SGD)

$$\begin{aligned} & \mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N})] - L(\mathbf{w}^*) \\ & \leq \frac{2 \|\boldsymbol{\eta}_0\|^2}{\gamma^2 \lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} (1 - \gamma \lambda_{\min})^{2(s+1)} \\ & \quad + \frac{2(1+d) (\lambda_{\max}^2 + \sigma_H^2/b)}{\lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} \|\boldsymbol{\eta}_0\|^2 (1 - \gamma \lambda_{\min})^{2s} \\ & = \frac{2 \|\boldsymbol{\eta}_0\|^2 (1 - \gamma \lambda_{\min})^{2s}}{\lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} \left[\frac{(1 - \gamma \lambda_{\min})^2}{\gamma^2} + (1+d) \left(\lambda_{\max}^2 + \frac{\sigma_H^2}{b} \right) \right] \end{aligned}$$

Optimizer: Zeroth-Order SGD Update Rule

$$\boldsymbol{w}_t = \boldsymbol{w}_{t-1} - \gamma \hat{\boldsymbol{g}}_{t-1}$$

Zeroth-Order Gradient Estimator (Two-Point, Mini-Batch)

$$\hat{\boldsymbol{g}}_{t-1} = \frac{1}{2\delta} \left[L_{B_{t-1}}(\boldsymbol{w}_{t-1} + \delta \boldsymbol{u}_{t-1}) - L_{B_{t-1}}(\boldsymbol{w}_{t-1} - \delta \boldsymbol{u}_{t-1}) \right] \boldsymbol{u}_{t-1}$$

Where:

- γ is the learning rate.
- δ is a small positive smoothing radius.
- $\boldsymbol{u}_{t-1} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$
- B_{t-1} is a mini-batch of size b , sampled at iteration $t - 1$.
- $L_{B_{t-1}}(\boldsymbol{w}) = \frac{1}{2b} \sum_{i \in B_{t-1}} (y_i - \langle \boldsymbol{w}, \boldsymbol{x}_i \rangle)^2$

Derivation

Let $n_t = w_t - w^*$. Substitute $w = n + w^*$ into the update rule:

$$n_t = n_{t-1} - \gamma \hat{g}_{t-1}$$

Mini-batch Loss Terms in Gradient Estimator

We approximate the loss on the mini-batch B_{t-1} as:

$$L_{B_{t-1}}(w_{t-1} \pm \delta u_{t-1}) = L_{B_{t-1}}(w^*) + \frac{1}{2}(n_{t-1} \pm \delta u_{t-1})^\top \mathbf{H}_{B_{t-1}}(n_{t-1} \pm \delta u_{t-1})$$

where $\mathbf{H}_{B_{t-1}} = \frac{1}{b} \sum_{i \in B_{t-1}} \mathbf{x}_i \mathbf{x}_i^\top$ is the empirical Hessian over the mini-batch.

$$= L_{B_{t-1}}(w^*) + \frac{1}{2} [n_{t-1}^\top \mathbf{H}_{B_{t-1}} n_{t-1} \pm 2\delta n_{t-1}^\top \mathbf{H}_{B_{t-1}} u_{t-1} + \delta^2 u_{t-1}^\top \mathbf{H}_{B_{t-1}} u_{t-1}]$$

Loss Difference

$$\begin{aligned} & L_{B_{t-1}}(w_{t-1} + \delta u_{t-1}) - L_{B_{t-1}}(w_{t-1} - \delta u_{t-1}) \\ &= \delta \cdot 2n_{t-1}^\top \mathbf{H}_{B_{t-1}} u_{t-1} \end{aligned}$$

Substitute into \hat{g}_{t-1}

$$\hat{g}_{t-1} = \frac{1}{2\delta} \cdot (2\delta n_{t-1}^\top \mathbf{H}_{B_{t-1}} \mathbf{u}_{t-1}) \cdot \mathbf{u}_{t-1} = (n_{t-1}^\top \mathbf{H}_{B_{t-1}} \mathbf{u}_{t-1}) \mathbf{u}_{t-1}$$

Outer Product Form

$$\hat{g}_{t-1} = (\mathbf{u}_{t-1} \mathbf{u}_{t-1}^\top \mathbf{H}_{B_{t-1}}) \mathbf{n}_{t-1}$$

Update Rule for \mathbf{n}_t

$$\mathbf{n}_t = [\mathbf{I} - \gamma \mathbf{u}_{t-1} \mathbf{u}_{t-1}^\top \mathbf{H}_{B_{t-1}}] \mathbf{n}_{t-1}$$

Conditional Expectation

Assuming independence of \mathbf{u}_{t-1} and $\mathbf{H}_{B_{t-1}}$, and using $\mathbb{E}[\mathbf{u}_{t-1} \mathbf{u}_{t-1}^\top] = \mathbf{I}$, we get:

$$\mathbb{E}[\mathbf{n}_t | \mathbf{n}_{t-1}] = (\mathbf{I} - \gamma \mathbf{H}) \mathbf{n}_{t-1}$$

Noise Term

$$\text{Noise}_t = \gamma (\mathbf{I} - \mathbf{u}_{t-1} \mathbf{u}_{t-1}^\top) \mathbf{H}_{B_{t-1}} \mathbf{n}_{t-1}$$

Full Update Equation

$$\mathbf{n}_t = (\mathbf{I} - \gamma \mathbf{H}) \mathbf{n}_{t-1} + \gamma (\mathbf{I} - \mathbf{u}_{t-1} \mathbf{u}_{t-1}^\top) \mathbf{H}_{B_{t-1}} \mathbf{n}_{t-1}$$

Bias Evolution (Mean Error)

$$\mathbb{E}[\mathbf{n}_t] = (\mathbf{I} - \gamma \mathbf{H}) \mathbb{E}[\mathbf{n}_{t-1}]$$

Provided γ is chosen so that all eigenvalues of $(\mathbf{I} - \gamma \mathbf{H})$ lie within the unit circle, we have $\mathbb{E}[\mathbf{n}_t] \rightarrow 0$.

Variance Generation (Fluctuations)

The noise term:

$$\text{Noise}_t = \gamma(\mathbf{I} - \mathbf{u}_{t-1} \mathbf{u}_{t-1}^\top) \mathbf{H}_{B_{t-1}} \mathbf{n}_{t-1}$$

has zero conditional expectation given \mathbf{n}_{t-1} , and variance accumulates based on the empirical Hessian's fluctuations:

$$\text{Cov}(\mathbf{n}_t) = \mathbb{E}[\mathbf{n}_t \mathbf{n}_t^\top] - \mathbb{E}[\mathbf{n}_t] \mathbb{E}[\mathbf{n}_t]^\top$$

Now if we take the tail averaging for eta -

$$\begin{aligned} \bar{\eta}_{s+1,N} &= \frac{1}{N} \sum_{t=s+1}^{s+N} \eta_t \\ &= \frac{1}{N} \sum_{t=s+1}^{s+N} (\eta_t^{\text{bias}} + \eta_t^{\text{variance}}) \\ &= \bar{\eta}_{s+1,N}^{\text{bias}} + \bar{\eta}_{s+1,N}^{\text{variance}}. \end{aligned}$$

$$\begin{aligned}
\mathbb{E} [\boldsymbol{\eta}_N \otimes \boldsymbol{\eta}_N] &= \mathbb{E} [(\boldsymbol{\eta}_N^{\text{bias}} + \boldsymbol{\eta}_N^{\text{variance}}) \otimes (\boldsymbol{\eta}_N^{\text{bias}} + \boldsymbol{\eta}_N^{\text{variance}})] \\
&\preceq 2 \cdot \left(\mathbb{E} [(\boldsymbol{\eta}_N^{\text{bias}} \otimes \boldsymbol{\eta}_N^{\text{bias}})] + \mathbb{E} [(\boldsymbol{\eta}_N^{\text{variance}} \otimes \boldsymbol{\eta}_N^{\text{variance}})] \right).
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} [L(\mathbf{w}_N)] - L(\mathbf{w}^*) &= \frac{1}{2} \langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_N \otimes \boldsymbol{\eta}_N] \rangle \\
&\leq \frac{1}{2} \langle \mathbf{H}, 2 \cdot (\mathbb{E} [\boldsymbol{\eta}_N^{\text{bias}} \otimes \boldsymbol{\eta}_N^{\text{bias}}] + \mathbb{E} [\boldsymbol{\eta}_N^{\text{variance}} \otimes \boldsymbol{\eta}_N^{\text{variance}}]) \rangle \\
&\leq 2 \cdot \left(\frac{1}{2} \langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_N^{\text{bias}} \otimes \boldsymbol{\eta}_N^{\text{bias}}] \rangle + \frac{1}{2} \langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_N^{\text{variance}} \otimes \boldsymbol{\eta}_N^{\text{variance}}] \rangle \right) \\
&= 2 \cdot \left((\mathbb{E} [L(\mathbf{w}_N^{\text{bias}})] - L(\mathbf{w}^*)) + (\mathbb{E} [L(\mathbf{w}_N^{\text{variance}})] - L(\mathbf{w}^*)) \right).
\end{aligned}$$

$$\mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}] \preceq 2 \cdot (\mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N}^{\text{bias}} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}^{\text{bias}}] + \mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N}^{\text{variance}} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}^{\text{variance}}])$$

$$\begin{aligned}
L(\bar{\mathbf{w}}_{s+1,N}) - L(\mathbf{w}^*) &= \frac{1}{2} \langle \mathbf{H}, \mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}] \rangle \\
&\leq \langle \mathbf{H}, \mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N}^{\text{bias}} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}^{\text{bias}}] \rangle + \langle \mathbf{H}, \mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N}^{\text{variance}} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}^{\text{variance}}] \rangle
\end{aligned}$$

$$\begin{aligned}
\langle \mathbf{H}, \mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}] \rangle &\leq \frac{1}{N^2} \cdot \sum_{l=s+1}^{s+N} \left(\langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] (\gamma \mathbf{H})^{-1} + (\gamma \mathbf{H})^{-1} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] \rangle \right. \\
&\quad \left. - \frac{1}{N^2} \cdot \sum_{l=s+1}^{s+N} \sum_{k=s+N+1}^{\infty} \left(\langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] (\mathbf{I} - \gamma \mathbf{H})^{k-l} \right. \right. \\
&\quad \left. \left. + (\mathbf{I} - \gamma \mathbf{H})^{k-l} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] \rangle \right) \right) \\
&\leq \frac{1}{N^2} \cdot \sum_{l=s+1}^{s+N} \left(\langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] (\gamma \mathbf{H})^{-1} + (\gamma \mathbf{H})^{-1} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] \rangle \right) \\
&= \frac{2}{\gamma N^2} \cdot \sum_{l=s+1}^{s+N} \text{Tr} \left(\mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] \right).
\end{aligned}$$

$$\begin{aligned}
&\langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] (\mathbf{I} - \gamma \mathbf{H})^{k-l} + (\mathbf{I} - \gamma \mathbf{H})^{k-l} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] \rangle \\
&= 2 \text{Tr} [\mathbf{H}(\mathbf{I} - \gamma \mathbf{H})^{k-l} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l]] \geq 0.
\end{aligned}$$

$$\begin{aligned} & \langle \mathbf{H}, \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] (\mathbf{I} - \gamma \mathbf{H})^{k-l} + (\mathbf{I} - \gamma \mathbf{H})^{k-l} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l] \rangle \\ &= 2 \operatorname{Tr} [\mathbf{H} (\mathbf{I} - \gamma \mathbf{H})^{k-l} \mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l]] \geq 0. \end{aligned}$$

$$L(\mathbf{w}) - L(\mathbf{w}^*) = \frac{1}{2} \operatorname{Tr} \left(\mathbf{H} \cdot (\boldsymbol{\eta} \otimes \boldsymbol{\eta}) \right), \text{ with } \boldsymbol{\eta} = \mathbf{w} - \mathbf{w}^*$$

$$\mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N})] - L(\mathbf{w}^*) = \frac{1}{2} \cdot \langle \mathbf{H}, \mathbb{E} [\bar{\boldsymbol{\eta}}_{s+1,N} \otimes \bar{\boldsymbol{\eta}}_{s+1,N}] \rangle$$

$$\begin{aligned} & \leq \frac{1}{\gamma N^2} \sum_{l=s+1}^{s+N} \operatorname{Tr} (\mathbb{E} [\boldsymbol{\eta}_l \otimes \boldsymbol{\eta}_l]) \\ & \leq \frac{2}{\gamma N^2} \cdot \sum_{l=s+1}^{s+N} \left(\operatorname{Tr} (\mathbb{E} [\boldsymbol{\eta}_l^{\text{bias}} \otimes \boldsymbol{\eta}_l^{\text{bias}}]) + \operatorname{Tr} (\mathbb{E} [\boldsymbol{\eta}_l^{\text{variance}} \otimes \boldsymbol{\eta}_l^{\text{variance}}]) \right) \end{aligned}$$

$$\mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N}^{\text{bias}})] - L(\mathbf{w}^*) \stackrel{\text{def}}{=} \frac{2}{\gamma N^2} \sum_{l=s+1}^{s+N} \operatorname{Tr} \left(\mathbb{E} [\boldsymbol{\eta}_l^{\text{bias}} \otimes \boldsymbol{\eta}_l^{\text{bias}}] \right).$$

$$\mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N}^{\text{variance}})] - L(\mathbf{w}^*) \stackrel{\text{def}}{=} \frac{2}{\gamma N^2} \sum_{l=s+1}^{s+N} \operatorname{Tr} \left(\mathbb{E} [\boldsymbol{\eta}_l^{\text{variance}} \otimes \boldsymbol{\eta}_l^{\text{variance}}] \right)$$

Goal

Derive bounds for the quantity

$$\text{tr} \left(\mathbb{E} \left[\eta_t^{\text{bias}} (\eta_t^{\text{bias}})^T \right] \right),$$

which represents the expected squared Euclidean norm of the bias component of the error vector at step t , i.e., $\mathbb{E}[\|\eta_t^{\text{bias}}\|^2]$.

Bias Recurrence

The bias evolves deterministically:

$$\eta_t^{\text{bias}} = (I - \gamma H) \eta_{t-1}^{\text{bias}},$$

where:

- γ is the effective learning rate.
- $H = \mathbb{E}[xx^T]$ is the Hessian, symmetric and positive definite with eigenvalues $\lambda_{\min} > 0$ and λ_{\max} .
- I is the identity matrix.

Outer Product Expansion

$$\begin{aligned}\eta_t^{\text{bias}}(\eta_t^{\text{bias}})^T &= (I - \gamma H)\eta_{t-1}^{\text{bias}}(\eta_{t-1}^{\text{bias}})^T(I - \gamma H)^T \\ &= (I - \gamma H)\eta_{t-1}^{\text{bias}}(\eta_{t-1}^{\text{bias}})^T(I - \gamma H) \quad (\text{since } H \text{ symmetric}).\end{aligned}$$

Expectation

Assuming deterministic η_0 :

$$\mathbb{E}[\eta_t^{\text{bias}}(\eta_t^{\text{bias}})^T] = (I - \gamma H)\mathbb{E}[\eta_{t-1}^{\text{bias}}(\eta_{t-1}^{\text{bias}})^T](I - \gamma H).$$

Let $M_{t-1} = \mathbb{E}[\eta_{t-1}^{\text{bias}}(\eta_{t-1}^{\text{bias}})^T]$ and $P = I - \gamma H$.

Then:

$$\mathbb{E}[\eta_t^{\text{bias}}(\eta_t^{\text{bias}})^T] = PM_{t-1}P.$$

Trace Bound

$$\begin{aligned}\mathrm{tr} \left(\mathbb{E}[\eta_t^{\mathrm{bias}} (\eta_t^{\mathrm{bias}})^T] \right) &= \mathrm{tr}(PM_{t-1}P) \\ &= \mathrm{tr}(P^2 M_{t-1}) \quad (\text{cyclic property of trace}) \\ &\leq \|P^2\|_2 \mathrm{tr}(M_{t-1}) \quad (\text{another property of trace}) \\ &= \|P\|_2^2 \mathrm{tr}(M_{t-1}).\end{aligned}$$

Spectral Norm of P

The eigenvalues of $P = I - \gamma H$ are $1 - \gamma \lambda_i(H)$. Thus,

$$\|P\|_2 = \max_i |1 - \gamma \lambda_i(H)|.$$

Define:

$$\rho(\gamma) := \|P\|_2 = \max(|1 - \gamma \lambda_{\min}|, |1 - \gamma \lambda_{\max}|).$$

Recursive Bound

$$\begin{aligned}\text{tr} \left(\mathbb{E}[\eta_t^{\text{bias}} (\eta_t^{\text{bias}})^T] \right) &\leq \rho(\gamma)^2 \text{tr} \left(\mathbb{E}[\eta_{t-1}^{\text{bias}} (\eta_{t-1}^{\text{bias}})^T] \right) \\ &\leq \rho(\gamma)^{2t} \text{tr}(\eta_0 \eta_0^T) \\ &= \rho(\gamma)^{2t} \|\eta_0\|^2,\end{aligned}$$

assuming η_0 is deterministic.

Conservative Learning Rate

If $0 < \gamma \leq \frac{1}{\lambda_{\max}}$, then all eigenvalues $1 - \gamma\lambda_i$ lie in $[0, 1]$. So:

$$\rho(\gamma) = 1 - \gamma\lambda_{\min},$$

and the bound becomes:

$$\mathbb{E}[\|\eta_t^{\text{bias}}\|^2] \leq (1 - \gamma\lambda_{\min})^{2t} \|\eta_0\|^2.$$

$$\begin{aligned}\mathbb{E} \left[L(\overline{\mathbf{w}}_{s+1,N}^{\text{bias}}) \right] - L(\mathbf{w}^*) &= \frac{2}{\gamma N^2} \sum_{t=s+1}^{s+N} \mathbb{E} \left[\|\eta_t\|^2 \right] \\ &\leq \frac{2}{\gamma N^2} \sum_{t=s+1}^{\infty} \mathbb{E} \left[\|\eta_t\|^2 \right]\end{aligned}$$

We start with the bound:

$$\mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N}^{\text{bias}})] - L(\mathbf{w}^*) \leq \frac{2}{\gamma N^2} \sum_{t=s+1}^{\infty} \mathbb{E} [\|\eta_t\|^2]$$

Using the bound:

$$\mathbb{E}[\|\eta_t^{\text{bias}}\|^2] \leq (1 - \gamma \lambda_{\min})^{2t} \|\eta_0\|^2,$$

we substitute:

$$\begin{aligned} \mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N}^{\text{bias}})] - L(\mathbf{w}^*) &\leq \frac{2}{\gamma N^2} \sum_{t=s+1}^{\infty} (1 - \gamma \lambda_{\min})^{2t} \|\eta_0\|^2 \\ &= \frac{2}{\gamma N^2} \|\eta_0\|^2 \sum_{t=s+1}^{\infty} (1 - \gamma \lambda_{\min})^{2t} \\ &= \frac{2}{\gamma N^2} \|\eta_0\|^2 (1 - \gamma \lambda_{\min})^{2(s+1)} \sum_{t=0}^{\infty} (1 - \gamma \lambda_{\min})^{2t} \\ &= \frac{2}{\gamma N^2} (1 - \gamma \lambda_{\min})^{2(s+1)} \|\eta_0\|^2 \cdot \frac{1}{1 - (1 - \gamma \lambda_{\min})^2} \\ &= \frac{2}{\gamma N^2} \cdot \frac{(1 - \gamma \lambda_{\min})^{2(s+1)}}{1 - (1 - 2\gamma \lambda_{\min} + \gamma^2 \lambda_{\min}^2)} \|\eta_0\|^2 \\ &= \frac{2}{\gamma N^2} \cdot \frac{(1 - \gamma \lambda_{\min})^{2(s+1)}}{2\gamma \lambda_{\min} - \gamma^2 \lambda_{\min}^2} \|\eta_0\|^2 \\ &= \frac{2}{\gamma^2 \lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} (1 - \gamma \lambda_{\min})^{2(s+1)} \|\eta_0\|^2 \end{aligned}$$

Let η_t^{var} denote the noise term representing the fluctuation around the mean update direction at step t . Based on the decomposition derived from the mini-batch Zeroth-Order (ZO) update rule, this term is given by:

$$\eta_t^{\text{var}} = \gamma(\mathbf{I} - \mathbf{u}_{t-1}\mathbf{u}_{t-1}^\top)\mathbf{H}_{B_{t-1}}\eta_{t-1},$$

where:

- γ is the learning rate parameter (which may incorporate dimension d).
- $\mathbf{u}_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is an isotropic random direction vector (standard Gaussian).
- $\mathbf{H}_{B_{t-1}} = \frac{1}{b} \sum_{i \in B_{t-1}} \mathbf{x}_i \mathbf{x}_i^\top$ is the random empirical Hessian calculated from the mini-batch B_{t-1} of size b . Note that $\mathbb{E}_{B_{t-1}}[\mathbf{H}_{B_{t-1}}] = \mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$.
- $\eta_{t-1} = \mathbf{w}_{t-1} - \mathbf{w}^*$ is the total error vector at step $t - 1$.

We are interested in bounding the expected squared Euclidean norm of this noise term:

$$\mathbb{E}[\|\eta_t^{\text{var}}\|^2] = \text{tr} \left(\mathbb{E}[\eta_t^{\text{var}}(\eta_t^{\text{var}})^\top] \right).$$

The outer expectation \mathbb{E} is taken over all sources of randomness, which include the random direction \mathbf{u} and the mini-batch sampling B .

Step 1: Conditional Second Moment

We first compute the expectation conditioned on the state η_{t-1} and the specific mini-batch B_{t-1} . The only remaining randomness is $\mathbf{u} := \mathbf{u}_{t-1}$. Let $\mathbf{H}_B := \mathbf{H}_{B_{t-1}}$ for simpler notation in this step.

The conditional mean is:

$$\mathbb{E}_{\mathbf{u}}[\eta_t^{\text{var}} \mid \eta_{t-1}, B_{t-1}] = \gamma \mathbb{E}_{\mathbf{u}}[(\mathbf{I} - \mathbf{u}\mathbf{u}^\top) \mathbf{H}_B \eta_{t-1}] = \gamma(\mathbf{I} - \mathbb{E}_{\mathbf{u}}[\mathbf{u}\mathbf{u}^\top]) \mathbf{H}_B \eta_{t-1} = \gamma(\mathbf{I} - \mathbf{I}) \mathbf{H}_B \eta_{t-1} = \mathbf{0}.$$

Since the noise term is conditionally zero-mean, its conditional second moment equals its conditional covariance:

$$\mathbb{E}_{\mathbf{u}}[\eta_t^{\text{var}} (\eta_t^{\text{var}})^\top \mid \eta_{t-1}, B_{t-1}] = \text{Cov}_{\mathbf{u}}(\eta_t^{\text{var}} \mid \eta_{t-1}, B_{t-1}).$$

Let $\mathbf{M}_B = \mathbf{H}_B \eta_{t-1} \eta_{t-1}^\top \mathbf{H}_B$. Note that \mathbf{M}_B is fixed when conditioning on η_{t-1} and B_{t-1} .

We compute the conditional covariance:

$$\begin{aligned} \text{Cov}_{\mathbf{u}}(\eta_t^{\text{var}} \mid \eta_{t-1}, B_{t-1}) &= \mathbb{E}_{\mathbf{u}} \left[\left(\gamma(\mathbf{I} - \mathbf{u}\mathbf{u}^\top) \mathbf{H}_B \eta_{t-1} \right) \left(\gamma(\mathbf{I} - \mathbf{u}\mathbf{u}^\top) \mathbf{H}_B \eta_{t-1} \right)^\top \right] \\ &= \gamma^2 \mathbb{E}_{\mathbf{u}} \left[(\mathbf{I} - \mathbf{u}\mathbf{u}^\top) \mathbf{H}_B \eta_{t-1} \eta_{t-1}^\top \mathbf{H}_B (\mathbf{I} - \mathbf{u}\mathbf{u}^\top)^\top \right] \\ &= \gamma^2 \mathbb{E}_{\mathbf{u}} \left[(\mathbf{I} - \mathbf{u}\mathbf{u}^\top) \mathbf{M}_B (\mathbf{I} - \mathbf{u}\mathbf{u}^\top) \right] \quad (\text{since } \mathbf{I} - \mathbf{u}\mathbf{u}^\top \text{ and } \mathbf{M}_B \text{ are symmetric}) \\ &= \gamma^2 (\mathbf{M}_B + \text{tr}(\mathbf{M}_B) \mathbf{I}) \quad (\text{using Gaussian identity } \mathbb{E}[\mathbf{u}\mathbf{u}^\top \mathbf{M} \mathbf{u}\mathbf{u}^\top] = 2\mathbf{M} + \text{tr}(\mathbf{M}) \mathbf{I}) \\ &= \gamma^2 (\mathbf{H}_B \eta_{t-1} \eta_{t-1}^\top \mathbf{H}_B + \text{tr}(\mathbf{H}_B \eta_{t-1} \eta_{t-1}^\top \mathbf{H}_B) \mathbf{I}) \\ &= \gamma^2 (\mathbf{H}_B \eta_{t-1} \eta_{t-1}^\top \mathbf{H}_B + \text{tr}(\mathbf{H}_B^2 \eta_{t-1} \eta_{t-1}^\top) \mathbf{I}) \quad (\text{using cyclic property of trace}) \end{aligned}$$

Step 2: Total Expectation

Now, we take the expectation over the remaining randomness in $\boldsymbol{\eta}_{t-1}$ (inherited from previous steps) and the mini-batch sampling B_{t-1} . Let $\mathbb{E}_{\boldsymbol{\eta}, B}$ denote this expectation.

$$\begin{aligned}\mathbb{E}[\boldsymbol{\eta}_t^{\text{var}}(\boldsymbol{\eta}_t^{\text{var}})^\top] &= \mathbb{E}_{\boldsymbol{\eta}_{t-1}, B_{t-1}} \left[\mathbb{E}_{\mathbf{u}}[\boldsymbol{\eta}_t^{\text{var}}(\boldsymbol{\eta}_t^{\text{var}})^\top \mid \boldsymbol{\eta}_{t-1}, B_{t-1}] \right] \\ &= \mathbb{E}_{\boldsymbol{\eta}_{t-1}, B_{t-1}} \left[\gamma^2 \left(\mathbf{H}_B \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \mathbf{H}_B + \text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top) \mathbf{I} \right) \right] \\ &= \gamma^2 \left(\mathbb{E}[\mathbf{H}_B \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \mathbf{H}_B] + \mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] \mathbf{I} \right) \quad (\text{linearity of } \mathbb{E}) \\ &= \gamma^2 \left(\mathbb{E}[\mathbf{H}_B \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \mathbf{H}_B] + \text{tr}(\mathbb{E}[\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top]) \mathbf{I} \right). \quad (\text{linearity of trace})\end{aligned}$$

Step 3: Apply Trace and Bound

We apply the trace operator to get the expected squared norm:

$$\begin{aligned}\mathbb{E}[\|\boldsymbol{\eta}_t^{\text{var}}\|^2] &= \text{tr}(\mathbb{E}[\boldsymbol{\eta}_t^{\text{var}}(\boldsymbol{\eta}_t^{\text{var}})^\top]) \\ &= \text{tr} \left(\gamma^2 \left[\mathbb{E}[\mathbf{H}_B \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \mathbf{H}_B] + \text{tr}(\mathbb{E}[\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top]) \mathbf{I} \right] \right) \\ &= \gamma^2 \left[\text{tr}(\mathbb{E}[\mathbf{H}_B \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \mathbf{H}_B]) + \text{tr} \left(\text{tr}(\mathbb{E}[\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top]) \mathbf{I} \right) \right] \\ &= \gamma^2 \left[\mathbb{E}[\text{tr}(\mathbf{H}_B \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \mathbf{H}_B)] + d \cdot \text{tr}(\mathbb{E}[\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top]) \right] \quad (\text{trace properties}) \\ &= \gamma^2 \left[\mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] + d \cdot \mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] \right] \quad (\text{cyclic property}) \\ &= \gamma^2(1 + d) \mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)].\end{aligned}$$

Now we bound the remaining expectation term. Let \mathbb{E}_B denote expectation only over the mini-batch B .

$$\begin{aligned}\mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] &= \mathbb{E}_{\boldsymbol{\eta}_{t-1}} \left[\mathbb{E}_B[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top) \mid \boldsymbol{\eta}_{t-1}] \right] \\ &= \mathbb{E}_{\boldsymbol{\eta}_{t-1}} \left[\text{tr}(\mathbb{E}_B[\mathbf{H}_B^2 \mid \boldsymbol{\eta}_{t-1}] \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top) \right].\end{aligned}$$

We use the identity for the expected squared empirical Hessian:

$$\mathbb{E}_B[\mathbf{H}_B^2] = \mathbf{H}^2 + \frac{1}{b} \mathbf{V}_H, \quad \text{where } \mathbf{V}_H = \mathbb{E}[(\mathbf{x}\mathbf{x}^\top - \mathbf{H})^2] = \mathbb{E}[(\mathbf{x}\mathbf{x}^\top)^2] - \mathbf{H}^2.$$

Note that \mathbf{V}_H represents the variance of the Hessian estimator. Substitute this back:

$$\mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] = \mathbb{E}_{\boldsymbol{\eta}_{t-1}} \left[\text{tr} \left(\left(\mathbf{H}^2 + \frac{1}{b} \mathbf{V}_H \right) \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \right) \right].$$

Using the trace inequality $\text{tr}(XY) \leq \|X\|_2 \text{tr}(Y)$ for symmetric X and PSD $Y = \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top$ (note $\text{tr}(Y) = \|\boldsymbol{\eta}_{t-1}\|^2$):

$$\text{tr} \left(\left(\mathbf{H}^2 + \frac{1}{b} \mathbf{V}_H \right) \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top \right) \leq \left\| \mathbf{H}^2 + \frac{1}{b} \mathbf{V}_H \right\|_2 \cdot \text{tr}(\boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top) = \left\| \mathbf{H}^2 + \frac{1}{b} \mathbf{V}_H \right\|_2 \cdot \|\boldsymbol{\eta}_{t-1}\|^2.$$

We commonly bound the spectral norm using the triangle inequality and assuming a bound on the Hessian variance, $\|\mathbf{V}_H\|_2 \leq \sigma_H^2$:

$$\left\| \mathbf{H}^2 + \frac{1}{b} \mathbf{V}_H \right\|_2 \leq \|\mathbf{H}^2\|_2 + \frac{1}{b} \|\mathbf{V}_H\|_2 \leq \lambda_{\max}(\mathbf{H})^2 + \frac{\sigma_H^2}{b}.$$

Applying this bound:

$$\mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] \leq \mathbb{E}_{\boldsymbol{\eta}_{t-1}} \left[\left(\lambda_{\max}(\mathbf{H})^2 + \frac{\sigma_H^2}{b} \right) \|\boldsymbol{\eta}_{t-1}\|^2 \right] = \left(\lambda_{\max}(\mathbf{H})^2 + \frac{\sigma_H^2}{b} \right) \mathbb{E}[\|\boldsymbol{\eta}_{t-1}\|^2].$$

Here, σ_H^2 is a constant bounding $\|\mathbb{E}[(\mathbf{x}\mathbf{x}^\top)^2] - \mathbf{H}^2\|_2$.

Combining the results, we obtain the final bound on the expected squared norm of the variance term:

$$\begin{aligned} \mathbb{E}[\|\boldsymbol{\eta}_t^{\text{var}}\|^2] &= \gamma^2(1+d) \mathbb{E}[\text{tr}(\mathbf{H}_B^2 \boldsymbol{\eta}_{t-1} \boldsymbol{\eta}_{t-1}^\top)] \\ &\leq \gamma^2(1+d) \left(\lambda_{\max}(\mathbf{H})^2 + \frac{\sigma_H^2}{b} \right) \mathbb{E}[\|\boldsymbol{\eta}_{t-1}\|^2]. \end{aligned}$$

This bound shows that the variance introduced by the ZO estimation noise at each step:

- Is proportional to the square of the learning rate, γ^2 .
- Grows approximately linearly with the dimension d .
- Depends on the maximum curvature of the loss landscape, $\lambda_{\max}(\mathbf{H})^2$.
- Includes a term related to the variance of the mini-batch Hessian estimator, σ_H^2/b , which decreases as the mini-batch size b increases.
- Is scaled by the expected squared total error from the previous step, $\mathbb{E}[\|\boldsymbol{\eta}_{t-1}\|^2]$, highlighting its multiplicative nature.

We are interested in bounding the expected optimization error (excess loss) attributed to the accumulated variance, averaged over N iterates starting from step $s+1$. Let $\bar{\mathbf{w}}_{s+1,N}$ represent some form of averaged iterate over $\mathbf{w}_{s+1}, \dots, \mathbf{w}_{s+N}$. Under certain analysis frameworks (e.g., for Polyak-Ruppert averaging or related analyses), the contribution of the variance terms to the excess loss can be related to the sum of the expected squared norms of the noise terms.

Let's assume the relevant quantity to bound is proportional to the sum of the expected squared norms of the noise terms:

$$\text{Error}_{\text{variance}} \propto \sum_{l=s+1}^{s+N} \mathbb{E} [\|\boldsymbol{\eta}_l^{\text{variance}}\|^2]$$

We need to bound $\sum_{l=s+1}^{s+N} \mathbb{E} [\|\boldsymbol{\eta}_l^{\text{variance}}\|^2]$.

From the previous analysis of the noise term $\boldsymbol{\eta}_l^{\text{var}} = \gamma(\mathbf{I} - \mathbf{u}_{l-1}\mathbf{u}_{l-1}^\top)\mathbf{H}_{B_{l-1}}\boldsymbol{\eta}_{l-1}$ in the mini-batch setting (assuming Gaussian \mathbf{u}), we derived the bound:

$$\mathbb{E}[\|\boldsymbol{\eta}_l^{\text{variance}}\|^2] \leq \gamma^2(1+d) \left(\lambda_{\max}(\mathbf{H})^2 + \frac{\sigma_H^2}{b} \right) \mathbb{E}[\|\boldsymbol{\eta}_{l-1}\|^2].$$

where $\boldsymbol{\eta}_{l-1} = \mathbf{w}_{l-1} - \mathbf{w}^*$ is the *total* error vector at step $l-1$, d is the dimension, b is the mini-batch size, $\lambda_{\max}(\mathbf{H})$ is the max eigenvalue of the population Hessian, and σ_H^2 relates to the variance of the Hessian estimator (e.g., $\sigma_H^2 \approx \|\mathbb{E}[(\mathbf{x}\mathbf{x}^\top)^2] - \mathbf{H}^2\|_2$).

Let $C_{\text{var}} = \gamma^2(1 + d) \left(\lambda_{\max}(\mathbf{H})^2 + \frac{\sigma_u^2}{b} \right)$. Then:

$$\mathbb{E}[\|\boldsymbol{\eta}_l^{\text{variance}}\|^2] \leq C_{\text{var}} \cdot \mathbb{E}[\|\boldsymbol{\eta}_{l-1}\|^2].$$

Summing over the iterations:

$$\begin{aligned} \sum_{l=s+1}^{s+N} \mathbb{E}[\|\boldsymbol{\eta}_l^{\text{variance}}\|^2] &\leq \sum_{l=s+1}^{s+N} C_{\text{var}} \cdot \mathbb{E}[\|\boldsymbol{\eta}_{l-1}\|^2] \\ &= C_{\text{var}} \sum_{k=s}^{s+N-1} \mathbb{E}[\|\boldsymbol{\eta}_k\|^2]. \end{aligned}$$

To proceed, we need a bound on the sum of the total expected squared errors $\sum_{k=s}^{s+N-1} \mathbb{E}[\|\boldsymbol{\eta}_k\|^2]$. Assuming the algorithm converges, the total error $\mathbb{E}[\|\boldsymbol{\eta}_k\|^2]$ typically decreases (possibly to a non-zero noise floor if the learning rate is constant). A common, simple bound, especially if the error is decreasing or stable around iteration s , is:

$$\sum_{k=s}^{s+N-1} \mathbb{E}[\|\boldsymbol{\eta}_k\|^2] \leq N \cdot \sup_{k \geq s} \mathbb{E}[\|\boldsymbol{\eta}_k\|^2] \quad \text{or perhaps} \quad \sum_{k=s}^{s+N-1} \mathbb{E}[\|\boldsymbol{\eta}_k\|^2] \leq N \cdot \mathbb{E}[\|\boldsymbol{\eta}_s\|^2] \quad (\text{if error is decreasing}).$$

Alternatively, if we know $\mathbb{E}[\|\boldsymbol{\eta}_k\|^2]$ converges to some value \mathcal{E}_∞ , then for large s , the sum is approximately $N \cdot \mathcal{E}_\infty$.

Using the simpler bound $\sum_{k=s}^{s+N-1} \mathbb{E}[\|\boldsymbol{\eta}_k\|^2] \leq N \cdot \mathbb{E}[\|\boldsymbol{\eta}_s\|^2]$ (assuming error doesn't significantly increase after step s):

$$\sum_{l=s+1}^{s+N} \mathbb{E}[\|\boldsymbol{\eta}_l^{\text{variance}}\|^2] \leq C_{\text{var}} \cdot N \cdot \mathbb{E}[\|\boldsymbol{\eta}_s\|^2].$$

Substituting the definition of C_{var} :

$$\sum_{l=s+1}^{s+N} \mathbb{E}[\|\boldsymbol{\eta}_l^{\text{variance}}\|^2] \leq N\gamma^2(1+d) \left(\lambda_{\max}(\boldsymbol{H})^2 + \frac{\sigma_H^2}{b} \right) \mathbb{E}[\|\boldsymbol{\eta}_s\|^2].$$

If we plug this into the specific error expression provided in the original template (assuming its correctness for the analysis context):

$$\begin{aligned} \mathbb{E} \left[L(\bar{\mathbf{w}}_{s+1,N}^{\text{variance}}) \right] - L(\mathbf{w}^*) &\triangleq \frac{2}{\gamma N^2} \sum_{l=s+1}^{s+N} \mathbb{E} \left[\|\boldsymbol{\eta}_l^{\text{variance}}\|^2 \right] \\ &\leq \frac{2}{\gamma N^2} \left[N\gamma^2(1+d) \left(\lambda_{\max}(\boldsymbol{H})^2 + \frac{\sigma_H^2}{b} \right) \mathbb{E}[\|\boldsymbol{\eta}_s\|^2] \right] \\ &= \frac{2\gamma(1+d)}{N} \left(\lambda_{\max}(\boldsymbol{H})^2 + \frac{\sigma_H^2}{b} \right) \mathbb{E}[\|\boldsymbol{\eta}_s\|^2]. \end{aligned}$$

This bound highlights that the average error contribution from variance over N steps:

- Decreases as $1/N$ (the averaging effect).
- Increases linearly with the learning rate γ .
- Increases with dimension d .
- Depends on the Hessian properties $(\lambda_{\max}, \sigma_H^2)$ and mini-batch size b .
- Depends on the magnitude of the total error $\mathbb{E}[\|\boldsymbol{\eta}_s\|^2]$ at the beginning of the averaging window.

First order

$$\mathbb{E} [L(\bar{\mathbf{w}})] - L(\mathbf{w}^*) \leq \frac{2}{\gamma^2 \mu^2} \cdot \frac{(1 - \gamma \mu)^s}{\left(\frac{n}{b} - s\right)^2} \cdot (L(\mathbf{w}_0) - L(\mathbf{w}^*)) + 4 \cdot \frac{\widehat{\sigma_{MLE}^2}}{b \cdot \left(\frac{n}{b} - s\right)}.$$

Zeroth order

$$\begin{aligned} & \mathbb{E} [L(\bar{\mathbf{w}}_{s+1,N})] - L(\mathbf{w}^*) \\ & \leq \frac{2 \|\boldsymbol{\eta}_0\|^2}{\gamma^2 \lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} (1 - \gamma \lambda_{\min})^{2(s+1)} \\ & \quad + \frac{2(1+d) (\lambda_{\max}^2 + \sigma_H^2/b)}{\lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} \|\boldsymbol{\eta}_0\|^2 (1 - \gamma \lambda_{\min})^{2s} \\ & = \frac{2 \|\boldsymbol{\eta}_0\|^2 (1 - \gamma \lambda_{\min})^{2s}}{\lambda_{\min} N^2 (2 - \gamma \lambda_{\min})} \left[\frac{(1 - \gamma \lambda_{\min})^2}{\gamma^2} + (1+d) \left(\lambda_{\max}^2 + \frac{\sigma_H^2}{b} \right) \right] \end{aligned}$$



THANK YOU