

One-step full gradient suffices for low-rank fine-tuning, provably and efficiently

Fanghui Liu

fanghui.liu@warwick.ac.uk

Department of Computer Science, University of Warwick, UK
Centre for Discrete Mathematics and its Applications (DIMAP), Warwick
[joint work with Yuanhe Zhang (Warwick) and Yudong Chen (UW-Madison)]



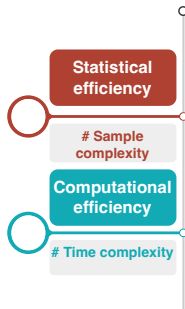
at HKUST@CS



My research

☐ Research interests

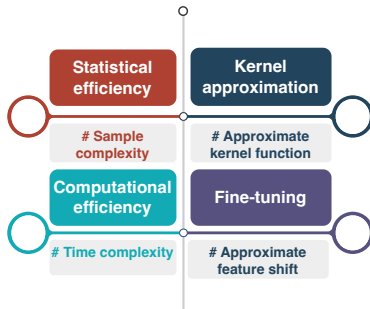
- Foundations of machine learning (ML)
 - Theory-grounded efficient algorithm design
 - Trustworthy ML



My research

☐ Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML



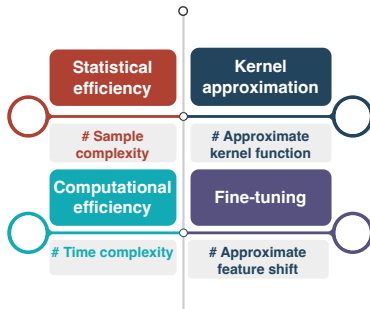
My research

❑ Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML

❑ Research goal

- characterize **learning efficiency** in theory
- contribute to practice



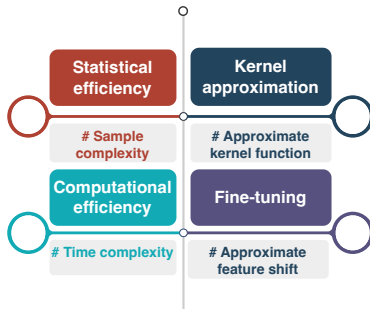
My research

❑ Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML

❑ Research goal

- characterize **learning efficiency** in theory
- contribute to practice



Learning efficiency (Curse of Dimensionality, CoD)

Machine learning works in **high dimensions** that can be a **curse**!

— David Donoho, 2000. (Richard E. Bellman, 1957)

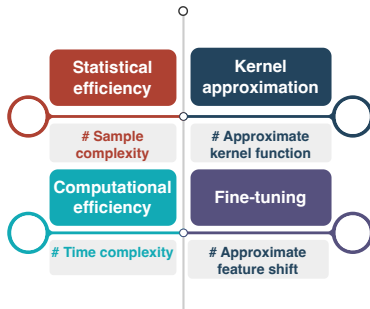
My research

❑ Research interests

- Foundations of machine learning (ML)
- Theory-grounded efficient algorithm design
- Trustworthy ML

❑ Research goal

- characterize **learning efficiency** in theory
- contribute to practice



Learning efficiency (Curse of Dimensionality, CoD)

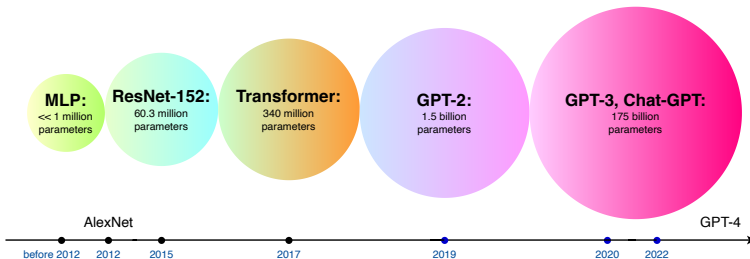
Machine learning works in **high dimensions** that can be a **curse**!

— David Donoho, 2000. (Richard E. Bellman, 1957)



In the era of machine learning (Pre-training)

relationship between data-centric, large model, huge compute resources

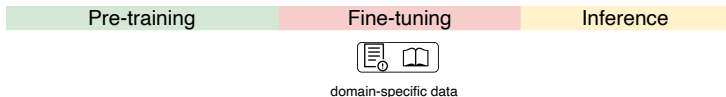


From pre-training to (parameter-efficient) fine-tuning

- GPT3: 175 billion parameters
- Llama3.1: > 400 billion parameters
- Deepseek-v3: > 600 billion parameters

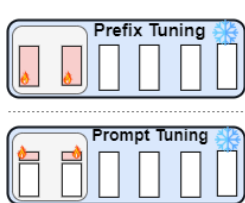
From pre-training to (parameter-efficient) fine-tuning

- GPT3: 175 billion parameters
- Llama3.1: > 400 billion parameters
- Deepseek-v3: > 600 billion parameters

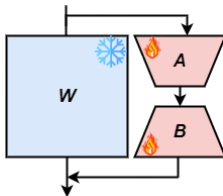


From pre-training to (parameter-efficient) fine-tuning

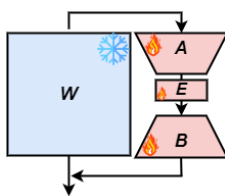
- GPT3: 175 billion parameters
- Llama3.1: > 400 billion parameters
- Deepseek-v3: > 600 billion parameters



(a) Prefix & Prompt



(b) LoRA



(c) LoRA variants

Low-rank adaption (LoRA) for fine-tuning [2]

$$\mathbf{W}^{\text{FT}} = \mathbf{W}^{\text{pre}} + \Delta \in \mathbb{R}^{d \times k}$$

- $\Delta \approx \mathbf{A}\mathbf{B}$ with $\mathbf{A} \in \mathbb{R}^{d \times r}$ and $\mathbf{B} \in \mathbb{R}^{r \times k}$
- initialization

$$[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \quad \text{and} \quad [\mathbf{B}_0]_{ij} = 0, \quad \alpha > 0. \quad (\text{LoRA-init})$$

- updated by gradient-based algorithms, e.g., SGD, AdamW
- obtain $(\mathbf{A}_t, \mathbf{B}_t)$

Low-rank adaption (LoRA) for fine-tuning [2]

$$\mathbf{W}^{\text{FT}} = \mathbf{W}^{\text{pre}} + \Delta \in \mathbb{R}^{d \times k}$$

- $\Delta \approx \mathbf{AB}$ with $\mathbf{A} \in \mathbb{R}^{d \times r}$ and $\mathbf{B} \in \mathbb{R}^{r \times k}$
- initialization

$$[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \quad \text{and} \quad [\mathbf{B}_0]_{ij} = 0, \quad \alpha > 0. \quad (\text{LoRA-init})$$

- updated by gradient-based algorithms, e.g., SGD, AdamW
- obtain $(\mathbf{A}_t, \mathbf{B}_t)$

Low-rank adaption (LoRA) for fine-tuning [2]

$$\mathbf{W}^{\text{FT}} = \mathbf{W}^{\text{pre}} + \Delta \in \mathbb{R}^{d \times k}$$

- $\Delta \approx \mathbf{AB}$ with $\mathbf{A} \in \mathbb{R}^{d \times r}$ and $\mathbf{B} \in \mathbb{R}^{r \times k}$
- initialization

$$[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2) \quad \text{and} \quad [\mathbf{B}_0]_{ij} = 0, \quad \alpha > 0. \quad (\text{LoRA-init})$$

- updated by gradient-based algorithms, e.g., SGD, AdamW
- obtain $(\mathbf{A}_t, \mathbf{B}_t)$

Motivation: non-linear dynamics and subspace alignment

- Even for linear model (pre-training and fine-tuning), **nonlinear dynamics**...

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^{\natural} \\ \eta_2 \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^{\natural} : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^{\natural} !

Target

- Q1: How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?
- Q2: How can our theoretical results contribute to algorithm design for LoRA in practice?

Motivation: non-linear dynamics and subspace alignment

- Even for linear model (pre-training and fine-tuning), **nonlinear dynamics**...

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^{\natural} \\ \eta_2 \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^{\natural} : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^{\natural} !

Target

- Q1: *How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?*
- Q2: *How can our theoretical results contribute to algorithm design for LoRA in practice?*

Motivation: non-linear dynamics and subspace alignment

- Even for linear model (pre-training and fine-tuning), **nonlinear dynamics**...

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^{\natural} \\ \eta_2 \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^{\natural} : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^{\natural} !

Target

- *Q1: How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?*
- *Q2: How can our theoretical results contribute to algorithm design for LoRA in practice?*

Motivation: non-linear dynamics and subspace alignment

- Even for linear model (pre-training and fine-tuning), **nonlinear dynamics**...

$$\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^{\natural} \\ \eta_2 \mathbf{G}^{\natural\top} & \mathbf{I}_k \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix} + \text{nonlinear term}.$$

- \mathbf{G}^{\natural} : one-step full gradient (from full fine-tuning)
- The dynamics $(\mathbf{A}_t, \mathbf{B}_t)$ heavily depends on \mathbf{G}^{\natural} !

Target

- Q1: *How to characterize low-rank dynamics of LoRA and the associated subspace alignment in theory?*
- Q2: *How can our theoretical results contribute to algorithm design for LoRA in practice?*

Alignment and theory-grounded algorithm

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^{\mathfrak{h}} \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^{\top} \mathbf{W}^{\mathfrak{h}})^{\top} \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^{\top} \mathbf{W}^{\mathfrak{h}})^{\top}] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^{\mathfrak{h}} := \mathbf{W}^{\mathfrak{h}} + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:

$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^{\top} \widetilde{\mathbf{W}}^{\mathfrak{h}})^{\top} \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^{\top} \widetilde{\mathbf{W}}^{\mathfrak{h}})^{\top}], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$

- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^{\mathfrak{h}} \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^{\top} \mathbf{W}^{\mathfrak{h}})^{\top} \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^{\top} \mathbf{W}^{\mathfrak{h}})^{\top}] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^{\mathfrak{h}} := \mathbf{W}^{\mathfrak{h}} + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:

$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^{\top} \widetilde{\mathbf{W}}^{\mathfrak{h}})^{\top} \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^{\top} \widetilde{\mathbf{W}}^{\mathfrak{h}})^{\top}], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$

- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^{\natural} \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^{\top} \mathbf{W}^{\natural})^{\top} \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^{\top} \mathbf{W}^{\natural})^{\top}] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^{\natural} := \mathbf{W}^{\natural} + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:

$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^{\top} \widetilde{\mathbf{W}}^{\natural})^{\top} \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^{\top} \widetilde{\mathbf{W}}^{\natural})^{\top}], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$

- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Problem setting and assumptions

- Pre-trained model: known $\mathbf{W}^{\mathfrak{h}} \in \mathbb{R}^{d \times k}$ and the ReLU activation σ

$$f_{\text{pre}}(\mathbf{x}) := \begin{cases} (\mathbf{x}^\top \mathbf{W}^{\mathfrak{h}})^\top \in \mathbb{R}^k & \text{linear} \\ \sigma[(\mathbf{x}^\top \mathbf{W}^{\mathfrak{h}})^\top] \in \mathbb{R}^k & \text{nonlinear} \end{cases}.$$

- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^{\mathfrak{h}} := \mathbf{W}^{\mathfrak{h}} + \Delta$
- $\text{Rank}(\Delta) = r^* < \min\{d, k\}$ with unknown r^*
- Downstream well-behaved data $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}_{i=1}^N$ for fine-tuning:

$$\tilde{\mathbf{y}} := \begin{cases} (\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^{\mathfrak{h}})^\top \in \mathbb{R}^k, \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \text{sub-Gaussian}, & \text{linear} \\ \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^{\mathfrak{h}})^\top], \quad \{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d) & \text{nonlinear} \end{cases}.$$

- We assume $N > d$, e.g., MetaMathQA, Code-Feedback, $d = 1,024$ and $N \sim 10^5$

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^{\natural}$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} \mathbf{W} - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} \mathbf{W}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} (\mathbf{W}^{\natural} + \mathbf{A}\mathbf{B}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} (\mathbf{W}^{\natural} + \mathbf{A}\mathbf{B})) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with different step-size, e.g., LoRA+ [1]

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta_1 \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta_2 \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}$: optimization and generalization!

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} \mathbf{W} - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} \mathbf{W}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A} \mathbf{B}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A} \mathbf{B})) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with different step-size, e.g., LoRA+ [1]

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta_1 \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta_2 \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}$: optimization and generalization!

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} \mathbf{W} - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} \mathbf{W}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A} \mathbf{B}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A} \mathbf{B})) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with different step-size, e.g., LoRA+ [1]

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta_1 \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta_2 \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}$: optimization and generalization!

Full fine-tuning and LoRA updates

- full fine-tuning (initialized at $\mathbf{W}_0 := \mathbf{W}^\natural$)

$$L(\mathbf{W}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} \mathbf{W} - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} \mathbf{W}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- LoRA update

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \begin{cases} \left\| \tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A} \mathbf{B}) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{linear} \\ \left\| \sigma(\tilde{\mathbf{X}} (\mathbf{W}^\natural + \mathbf{A} \mathbf{B})) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2 & \text{nonlinear} \end{cases}$$

- Gradient descent with different step-size, e.g., LoRA+ [1]

$$\mathbf{A}_{t+1} = \mathbf{A}_t - \eta_1 \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

$$\mathbf{B}_{t+1} = \mathbf{B}_t - \eta_2 \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)$$

- Evaluation by $\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}$: optimization and generalization!

- one-step full gradient: $\mathbf{G}^\natural \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}^\natural) = r^*$

$$\mathbf{G}^\natural := -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$

Theorem (Alignment between \mathbf{G}^\natural and B_t)

For the linear setting, consider the LoRA updates with (LoRA-init). We have

$$\left\| \mathbf{V}_{r^*, \perp}^\top \left(\mathbf{G}^\natural \right) \mathbf{V}_{r^*}(B_t) \right\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

Remark: $B_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}^\natural$ with $\text{Rank}(B_1) \leq r^*$

- one-step full gradient: $\mathbf{G}^\natural \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}^\natural) = r^*$

$$\mathbf{G}^\natural := -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$

Theorem (Alignment between \mathbf{G}^\natural and B_t)

For the linear setting, consider the LoRA updates with (LoRA-init). We have

$$\left\| \mathbf{V}_{r^*, \perp}^\top \left(\mathbf{G}^\natural \right) \mathbf{V}_{r^*}(\mathbf{B}_t) \right\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

Remark: $\mathbf{B}_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}^\natural$ with $\text{Rank}(\mathbf{B}_1) \leq r^*$

- one-step full gradient: $\mathbf{G}^\natural \in \mathbb{R}^{d \times k}$ and $\text{rank}(\mathbf{G}^\natural) = r^*$

$$\mathbf{G}^\natural := -\nabla_{\mathbf{W}} L(\mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top (\tilde{\mathbf{Y}} - \tilde{\mathbf{X}} \mathbf{W}^\natural) = \frac{1}{N} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \Delta.$$

Theorem (Alignment between \mathbf{G}^\natural and B_t)

For the linear setting, consider the LoRA updates with (LoRA-init). We have

$$\left\| \mathbf{V}_{r^*, \perp}^\top \left(\mathbf{G}^\natural \right) \mathbf{V}_{r^*}(\mathbf{B}_t) \right\|_{op} = 0, \quad \forall t \in \mathbb{N}_+.$$

Remark: $\mathbf{B}_1 = \eta_1 \mathbf{A}_0^\top \mathbf{G}^\natural$ with $\text{Rank}(\mathbf{B}_1) \leq r^*$

Theorem (Informal)

For $r \geq r^*$, recall $[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2)$ in (LoRA-init), for any $\epsilon \in (0, 1)$, choosing $\alpha = \mathcal{O}\left(\epsilon d^{-\frac{3}{4}\kappa^{\natural} - \frac{1}{2}}\right)$, running GD with $t^* \asymp \frac{\ln d}{\sqrt{\eta_1 \eta_2}}$, then we have

$$\left\| \mathbf{U}_{r^*, \perp}^{\top}(\mathbf{G}^{\natural}) \mathbf{U}_{r^*}(\mathbf{A}_{t^*}) \right\|_{op} \lesssim \epsilon, \text{ w.h.p.}$$

- small initialization: better alignment and better generalization performance
- imbalanced step-size finishes alignment earlier

Theorem (Informal)

For $r \geq r^*$, recall $[\mathbf{A}_0]_{ij} \sim \mathcal{N}(0, \alpha^2)$ in (LoRA-init), for any $\epsilon \in (0, 1)$, choosing $\alpha = \mathcal{O}\left(\epsilon d^{-\frac{3}{4}\kappa^{\natural} - \frac{1}{2}}\right)$, running GD with $t^* \asymp \frac{\ln d}{\sqrt{\eta_1 \eta_2}}$, then we have

$$\left\| \mathbf{U}_{r^*, \perp}^\top(\mathbf{G}^{\natural}) \mathbf{U}_{r^*}(\mathbf{A}_{t^*}) \right\|_{op} \lesssim \epsilon, \text{ w.h.p.}$$

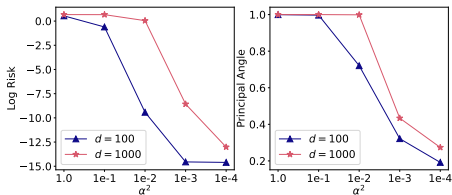


Figure 1: Left: the risk $\frac{1}{2} \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_F^2$. Right: the principal angle is defined as $\min_t \left\| \mathbf{U}_{r^*, \perp}^\top(\mathbf{G}^{\natural}) \mathbf{U}_{r^*}(\mathbf{A}_t) \right\|_{op}$.

- small initialization: better alignment and better generalization performance
- imbalanced step-size finishes alignment earlier

Key message: Algorithm design principle

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \sqrt{\gamma} \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]}.$$

$$\mathbf{B}_0 = \sqrt{\gamma} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top.$$

(Spectral-initialization)

Message

If we choose (Spectral-initialization), for both **linear/nonlinear** models, we can directly achieve the alignment at initialization.

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. \ 1 - C \exp(-\epsilon^2 N)$$

Key message: Algorithm design principle

- Take the SVD of \mathbf{G}^{\natural} : $\mathbf{G}^{\natural} = \tilde{\mathbf{U}}_{\mathbf{G}^{\natural}} \tilde{\mathbf{S}}_{\mathbf{G}^{\natural}} \tilde{\mathbf{V}}_{\mathbf{G}^{\natural}}^{\top}$

$$\mathbf{A}_0 = \sqrt{\gamma} \left[\tilde{\mathbf{U}}_{\mathbf{G}^{\natural}} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^{\natural}}^{1/2} \right]_{[1:r]}.$$

$$\mathbf{B}_0 = \sqrt{\gamma} \left[\tilde{\mathbf{S}}_{\mathbf{G}^{\natural}}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^{\natural}} \right]_{[:,1:r]}^{\top}.$$

(Spectral-initialization)

Message

If we choose (Spectral-initialization), for both linear/nonlinear models, we can directly achieve the alignment at initialization.

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \epsilon \|\Delta\|_{\text{op}}, \quad w.p. \ 1 - C \exp(-\epsilon^2 N)$$

Key message: Algorithm design principle

- Take the SVD of \mathbf{G}^\natural : $\mathbf{G}^\natural = \tilde{\mathbf{U}}_{\mathbf{G}^\natural} \tilde{\mathbf{S}}_{\mathbf{G}^\natural} \tilde{\mathbf{V}}_{\mathbf{G}^\natural}^\top$

$$\mathbf{A}_0 = \sqrt{\gamma} \left[\tilde{\mathbf{U}}_{\mathbf{G}^\natural} \right]_{[:,1:r]} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]}.$$

$$\mathbf{B}_0 = \sqrt{\gamma} \left[\tilde{\mathbf{S}}_{\mathbf{G}^\natural}^{1/2} \right]_{[1:r]} \left[\tilde{\mathbf{V}}_{\mathbf{G}^\natural} \right]_{[:,1:r]}^\top.$$

(Spectral-initialization)

Message

If we choose (Spectral-initialization), for both linear/nonlinear models, we can directly achieve the alignment at initialization.

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_F \leq \epsilon \|\Delta\|_{op}, \quad w.p. \ 1 - C \exp(-\epsilon^2 N)$$

Toy example (I)

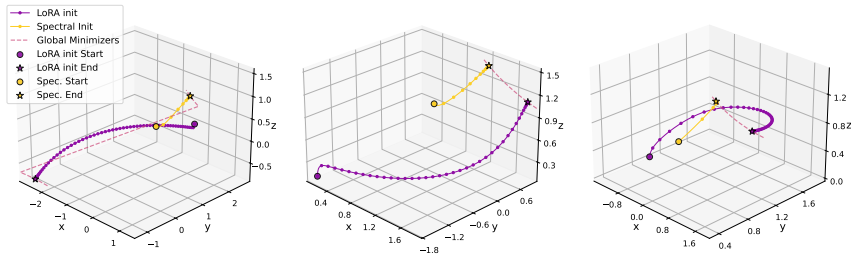
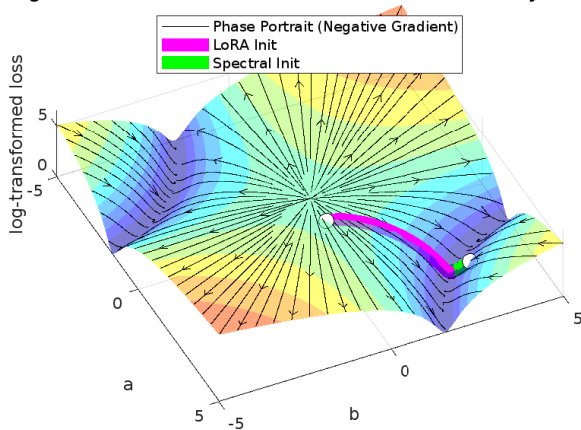


Figure 2: Comparison of the GD trajectories between LoRA and ours. $\mathbf{A} \in \mathbb{R}^2$ and $\mathbf{B} \in \mathbb{R}$. The set of global minimizers is $\{a_1^* = 2/t, a_2^* = 1/t, b^* = t \mid t \in \mathbb{R}\}$

Toy example (II): Phase portrait

Log-Transformed Surface with Phase Portrait and Trajectories



One-step full gradient may suffice for low-rank fine-tuning!

Dataset Size	MNLI 393k	SST-2 67k	CoLA 8.5k	QNLI 105k	MRPC 3.7k
Full	86.33 \pm 0.00	94.75 \pm 0.21	80.70 \pm 0.24	93.19 \pm 0.22	84.56 \pm 0.73
Pre-trained	-	89.79	59.03	49.28	63.48
One-step GD	-	90.48	73.00	69.13	68.38
LoRA ₈	85.30 \pm 0.04	94.04 \pm 0.09	72.84 \pm 1.25	93.02 \pm 0.07	68.38 \pm 0.01

Time cost

- CoLA LoRA: 47s, one-step: <1s
- MRPC LoRA: 25s, one-step: <1s

One-step full gradient may suffice for low-rank fine-tuning!

Dataset Size	MNLI 393k	SST-2 67k	CoLA 8.5k	QNLI 105k	MRPC 3.7k
Full	86.33 \pm 0.00	94.75 \pm 0.21	80.70 \pm 0.24	93.19 \pm 0.22	84.56 \pm 0.73
Pre-trained	-	89.79	59.03	49.28	63.48
One-step GD	-	90.48	73.00	69.13	68.38
LoRA ₈	85.30 \pm 0.04	94.04 \pm 0.09	72.84 \pm 1.25	93.02 \pm 0.07	68.38 \pm 0.01

Time cost

- **CoLA** LoRA: 47s, one-step: <1s
- **MRPC** LoRA: 25s, one-step: <1s

Clarification on gradient alignment based work

- Motivation: make LoRA's gradients align to full fine-tuning [5]
- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow [\tilde{\mathbf{U}}_{\mathbf{G}^\natural}]_{[:,1:r]}, \mathbf{B}_0 \leftarrow [\tilde{\mathbf{V}}_{\mathbf{G}^\natural}]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .

Clarification on gradient alignment based work

- Motivation: make LoRA's gradients align to full fine-tuning [5]
- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow [\tilde{\mathbf{U}}_{\mathbf{G}^\natural}]_{[:,1:r]}, \mathbf{B}_0 \leftarrow [\tilde{\mathbf{V}}_{\mathbf{G}^\natural}]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .

Clarification on gradient alignment based work

- Motivation: make LoRA's gradients align to full fine-tuning [5]
- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow [\tilde{\mathbf{U}}_{\mathbf{G}^\natural}]_{[:,1:r]}, \mathbf{B}_0 \leftarrow [\tilde{\mathbf{V}}_{\mathbf{G}^\natural}]_{[:,r+1:2r]}^\top. \quad (\text{LoRA-GA})$$

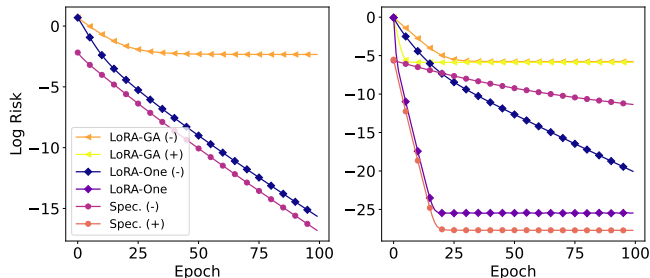
- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^\natural .

Clarification on gradient alignment based work

- Motivation: make LoRA's gradients align to full fine-tuning [5]
- best- $2r$ approximation: $\text{rank}(\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) + \text{rank}(\nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t)) \leq 2r$

$$\mathbf{A}_0 \leftarrow [\tilde{\mathbf{U}}_{\mathbf{G}^{\natural}}]_{[:,1:r]}, \mathbf{B}_0 \leftarrow [\tilde{\mathbf{V}}_{\mathbf{G}^{\natural}}]_{[:,r+1:2r]}^{\top}. \quad (\text{LoRA-GA})$$

- But! \mathbf{B}_t will align to the right-side rank- r^* singular subspace of \mathbf{G}^{\natural} .



Experiments

Key features in our LoRA-One algorithm

Algorithm 1 LoRA-One training for a specific layer

Input: Pre-trained weight $\mathbf{W}^{\mathfrak{h}}$, batched data $\{\mathcal{D}_t\}_{t=1}^T$, LoRA rank r , LoRA alpha α , loss function L

Output: $\mathbf{W}^{\mathfrak{h}} + \frac{\alpha}{\sqrt{r}} \mathbf{A}_T \mathbf{B}_T$

Compute $\nabla_{\mathbf{W}} L(\mathbf{W}^{\mathfrak{h}})$ and $\mathbf{U}, \mathbf{S}, \mathbf{V} \leftarrow \text{SVD}(\nabla_{\mathbf{W}} L(\mathbf{W}^{\mathfrak{h}}))$

$$\mathbf{A}_0 \leftarrow \sqrt{\gamma} \cdot \mathbf{U}_{[:,1:r]}$$

$$\mathbf{B}_0 \leftarrow \sqrt{\gamma} \cdot \mathbf{V}_{[:,1:r]}^{\top}$$

$$\mathbf{W}^{\mathfrak{h}} \leftarrow \mathbf{W}^{\mathfrak{h}} - \frac{\alpha}{\sqrt{r}} \mathbf{A}_0 \mathbf{B}_0$$

for $t = 1, \dots, T$ **do**

$$\mathbf{G}_t^{\mathbf{A}} \leftarrow \nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_{t-1}, \mathbf{B}_{t-1}) \left(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^{\top} + \lambda \mathbf{I}_r \right)^{-1}$$

$$\mathbf{G}_t^{\mathbf{B}} \leftarrow \left(\mathbf{A}_{t-1}^{\top} \mathbf{A}_{t-1} + \lambda \mathbf{I}_r \right)^{-1} \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_{t-1}, \mathbf{B}_{t-1})$$

$$\text{Update } \mathbf{A}_t, \mathbf{B}_t \leftarrow \text{AdamW}(\mathbf{G}_t^{\mathbf{A}}, \mathbf{G}_t^{\mathbf{B}})$$

end

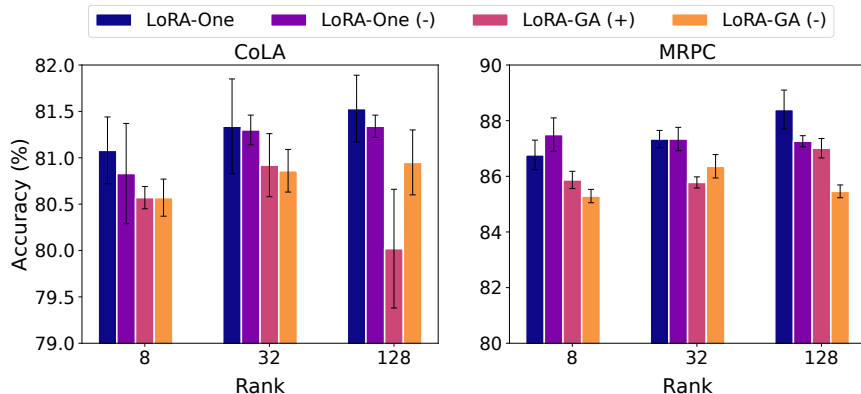
Experiments on NLP tasks from GLUE

Dataset Size	MNLI 393k	SST-2 67k	CoLA 8.5k	QNLI 105k	MRPC 3.7k
Full	86.33 \pm 0.00	94.75 \pm 0.21	80.70 \pm 0.24	93.19 \pm 0.22	84.56 \pm 0.73
Pre-trained	-	89.79	59.03	49.28	63.48
One-step GD	-	90.48	73.00	69.13	68.38
LoRA ₈ (Hu et al., 2022)	85.30 \pm 0.04	94.04 \pm 0.09	72.84 \pm 1.25	93.02 \pm 0.07	68.38 \pm 0.01
LoRA ₃₂	85.23 \pm 0.11	94.08 \pm 0.05	70.66 \pm 0.41	92.87 \pm 0.05	67.24 \pm 0.58
LoRA ₁₂₈	85.53 \pm 0.13	93.96 \pm 0.05	69.45 \pm 0.25	92.91 \pm 0.13	65.36 \pm 0.31
LoRA+ ₈ (Hayou et al., 2024)	85.81 \pm 0.09	93.85 \pm 0.24	77.53 \pm 0.20	93.14 \pm 0.03	74.43 \pm 1.39
LoRA+ ₃₂	85.88 \pm 0.16	94.15 \pm 0.25	79.29 \pm 0.96	93.25 \pm 0.08	79.49 \pm 0.64
LoRA+ ₁₂₈	86.07 \pm 0.15	94.08 \pm 0.30	78.59 \pm 0.73	93.06 \pm 0.23	78.76 \pm 0.12
P-LoRA ₈ (Zhang & Pilanci, 2024)	85.28 \pm 0.15	93.88 \pm 0.11	79.58 \pm 0.67	93.00 \pm 0.07	83.91 \pm 1.16
P-LoRA ₃₂	85.07 \pm 0.11	94.08 \pm 0.14	76.54 \pm 1.29	93.00 \pm 0.08	79.49 \pm 0.50
P-LoRA ₁₂₈	85.38 \pm 0.11	93.96 \pm 0.24	72.04 \pm 1.89	92.98 \pm 0.06	79.66 \pm 1.44
LoRA-GA ₈ (Wang et al., 2024a)	85.70 \pm 0.09	94.11 \pm 0.18	80.57 \pm 0.20	93.18 \pm 0.06	85.29 \pm 0.24
LoRA-GA ₃₂	83.32 \pm 0.10	94.49 \pm 0.32	80.86 \pm 0.23	93.06 \pm 0.14	86.36 \pm 0.42
LoRA-GA ₁₂₈	84.75 \pm 0.06	94.19 \pm 0.14	80.95 \pm 0.35	93.12 \pm 0.11	85.46 \pm 0.23
LoRA-One ₈ (Ours)	85.81 \pm 0.03	94.69 \pm 0.05	81.08 \pm 0.36	93.22 \pm 0.12	86.77 \pm 0.53
LoRA-One ₃₂	86.08 \pm 0.01	94.73 \pm 0.37	81.34 \pm 0.51	93.19 \pm 0.02	87.34 \pm 0.31
LoRA-One ₁₂₈	86.22 \pm 0.08	94.65 \pm 0.19	81.53 \pm 0.36	93.34 \pm 0.11	88.40 \pm 0.70

Experimental results on fine-tuning Llama 2-7B

	GSM8K	Human-eval
Full	59.36 \pm 0.85	35.31 \pm 2.13
LoRA ₈	46.89 \pm 0.05 (6.33h)	15.67 \pm 0.60 (6.75h)
LoRA ₃₂	47.44 \pm 0.74	16.02 \pm 0.85
LoRA ₁₂₈	47.33 \pm 0.32	15.57 \pm 0.75
LoRA-GA ₈	53.60 \pm 0.13	20.45 \pm 0.92
LoRA-GA ₃₂	55.12 \pm 0.30	20.18 \pm 0.19
LoRA-GA ₁₂₈	55.07 \pm 0.18	23.05 \pm 0.37
LoRA-One ₈	53.80 \pm 0.44 (+0.5h)	21.02 \pm 0.01 (+0.25h)
LoRA-One ₃₂	56.61 \pm 0.29	23.86 \pm 0.01
LoRA-One ₁₂₈	58.10 \pm 0.10	26.79 \pm 0.21

Ablation study



- (+): with preconditioners
- (-): no preconditioners

Theory and proof...

Model	Results	Algorithm	Initialization	Conclusion
Linear	Theorem 3.1	GD	(LoRA-init)	Subspace alignment of \mathbf{B}_t
	Theorem 3.2	GD	(LoRA-init)	Subspace alignment of \mathbf{A}_t
	Proposition 3.3	GD	(Spectral-init)	$\ \mathbf{A}_0\mathbf{B}_0 - \Delta\ _F$ is small
	Theorem 3.5	GD	(Spectral-init)	Linear convergence of $\ \mathbf{A}_t\mathbf{B}_t - \Delta\ _F$
	Theorem 3.6	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
Nonlinear	Theorem 4.3	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
	Theorem C.15	Smoothed Precondition GD	(Spectral-init)	Better convergence performance with less assumptions

- subspace alignment
- global convergence

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^b \\ \eta_2 \mathbf{G}^b^\top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta_1 \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta_2 \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

- Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$
- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$ (decouple $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ and obtain the alignment to \mathbf{G}^b)
- Define the residual term $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in early stage $t < T_1 \sim \ln\left(\frac{\|\mathbf{G}^b\|_{op}}{\|\mathbf{A}_0\|_{op}^2}\right)$
- Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [4] (Stöger & Soltanolkotabi)
- $\|\mathbf{U}_{r^*, \perp}^\top (\mathbf{G}^b) \mathbf{U}_{r^*}(\mathbf{A}_t)\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top (\mathbf{P}_t^A) \mathbf{U}_{r^*}(\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op}$ is small, w.h.p.

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^b \\ \eta_2 \mathbf{G}^{b^\top} & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta_1 \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta_2 \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

- Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$
- Schur decomposition of \mathbf{H}
- obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$ (decouple $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ and obtain the alignment to \mathbf{G}^b)
- Define the residual term $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in early stage $t < T_1 \sim \ln\left(\frac{\|\mathbf{G}^b\|_{op}}{\|\mathbf{A}_0\|_{op}^2}\right)$

◦ Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [4] (Stöger & Soltanolkotabi)

$$\|\mathbf{U}_{r^*, \perp}^\top (\mathbf{G}^b) \mathbf{U}_{r^*}(\mathbf{A}_t)\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top (\mathbf{P}_t^A) \mathbf{U}_{r^*}(\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op} \text{ is small, w.h.p.}$$

Proof of sketch: Control the dynamics for alignment

$$\underbrace{\begin{bmatrix} \mathbf{A}_{t+1} \\ \mathbf{B}_{t+1}^\top \end{bmatrix}}_{:= \mathbf{Z}_{t+1}} = \underbrace{\begin{bmatrix} \mathbf{I}_d & \eta_1 \mathbf{G}^\natural \\ \eta_2 \mathbf{G}^\natural^\top & \mathbf{I}_k \end{bmatrix}}_{:= \mathbf{H}} \underbrace{\begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}}_{:= \mathbf{Z}_t} - \frac{1}{N} \begin{bmatrix} 0 & \eta_1 \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{A}_t \mathbf{B}_t \\ \eta_2 \mathbf{B}_t^\top \mathbf{A}_t^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{A}_t \\ \mathbf{B}_t^\top \end{bmatrix}.$$

- Approximated linear dynamical system $\mathbf{Z}_t^{\text{lin}} := \mathbf{H}^t \mathbf{Z}_0$
 - Schur decomposition of \mathbf{H}
 - obtain the dynamics of $\mathbf{Z}_t^{\text{lin}}$ (decouple $\mathbf{A}_t^{\text{lin}}$ and $\mathbf{B}_t^{\text{lin}}$ and obtain the alignment to \mathbf{G}^\natural)
 - Define the residual term $\mathbf{E}_t := \mathbf{Z}_t - \mathbf{Z}_t^{\text{lin}}$, control $\|\mathbf{E}_t\|_{op}$ in early stage $t < T_1 \sim \ln\left(\frac{\|\mathbf{G}^\natural\|_{op}}{\|\mathbf{A}_0\|_{op}^2}\right)$
 - Transfer the alignment from $\mathbf{A}_t^{\text{lin}}$ to \mathbf{A}_t [4] (Stöger & Soltanolkotabi)
- $\|\mathbf{U}_{r^*, \perp}^\top (\mathbf{G}^\natural) \mathbf{U}_{r^*}(\mathbf{A}_t)\|_{op} \lesssim \|\mathbf{U}_{r^*, \perp}^\top (\mathbf{P}_t^A) \mathbf{U}_{r^*}(\mathbf{P}_t^A \mathbf{A}_0 + \mathbf{E}_t)\|_{op}$ is small, *w.h.p.*

Global convergence on nonlinear models

Recall problem setting and assumptions for nonlinear model

- Pre-trained model $f_{\text{pre}}(\mathbf{x}) = \sigma[(\mathbf{x}^\top \mathbf{W}^\flat)^\top] \in \mathbb{R}^k$
- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\flat := \mathbf{W}^\flat + \Delta$ with $\text{Rank}(\Delta) = r^*$
- We assume $r = r^*$.
- Downstream well-behaved data $\tilde{\mathbf{y}} = \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\flat)^\top]$, $\{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$
- training loss

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \left\| \sigma(\tilde{\mathbf{X}}(\mathbf{W}^\flat + \mathbf{A}\mathbf{B})) - \tilde{\mathbf{Y}} \right\|_{\mathbb{F}}^2.$$

- gradient updates

$$\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top, \quad \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t},$$

where we define

$$\mathbf{J}_{\mathbf{W}_t} := \underbrace{\frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \widetilde{\mathbf{W}}^\flat) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} \mathbf{W}_t) \right]}_{\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}} \odot \sigma'(\tilde{\mathbf{X}} \mathbf{W}_t).$$

- GLM-tron style: [3, 6]

Recall problem setting and assumptions for nonlinear model

- Pre-trained model $f_{\text{pre}}(\mathbf{x}) = \sigma[(\mathbf{x}^\top \mathbf{W}^\flat)^\top] \in \mathbb{R}^k$
- Unknown low-rank feature shift Δ : $\widetilde{\mathbf{W}}^\flat := \mathbf{W}^\flat + \Delta$ with $\text{Rank}(\Delta) = r^*$
- We assume $r = r^*$.
- Downstream well-behaved data $\tilde{\mathbf{y}} = \sigma[(\tilde{\mathbf{x}}^\top \widetilde{\mathbf{W}}^\flat)^\top]$, $\{\tilde{\mathbf{x}}_i\}_{i=1}^N \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$
- training loss

$$\tilde{L}(\mathbf{A}, \mathbf{B}) := \frac{1}{2N} \left\| \sigma(\tilde{\mathbf{X}}(\mathbf{W}^\flat + \mathbf{A}\mathbf{B})) - \tilde{\mathbf{Y}} \right\|_{\text{F}}^2.$$

- gradient updates

$$\nabla_{\mathbf{A}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{J}_{\mathbf{W}_t} \mathbf{B}_t^\top, \quad \nabla_{\mathbf{B}} \tilde{L}(\mathbf{A}_t, \mathbf{B}_t) = -\mathbf{A}_t^\top \mathbf{J}_{\mathbf{W}_t},$$

where we define

$$\mathbf{J}_{\mathbf{W}_t} := \underbrace{\frac{1}{N} \tilde{\mathbf{X}}^\top \left[\sigma(\tilde{\mathbf{X}} \widetilde{\mathbf{W}}^\flat) - \frac{1}{N} \tilde{\mathbf{X}}^\top \sigma(\tilde{\mathbf{X}} \mathbf{W}_t) \right]}_{\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}} \odot \sigma'(\tilde{\mathbf{X}} \mathbf{W}_t).$$

- GLM-tron style: [3, 6]

Theorem (Linear convergence rate)

Under (**Spectral-initialization**) and $\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), w.h.p$$

- holds for standard gradient update, but requires more assumptions.
- $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \epsilon \|\Delta\|_{\text{op}}, w.h.p.$

A bit proof sketch at

- recover at initialization:

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \|\mathbf{A}_0 \mathbf{B}_0 - \gamma \mathbf{G}^{\natural}\|_{\text{F}} + \text{concentration on } \mathbf{G}^{\natural} + \rho \lambda_{r^*}^*, w.h.p$$

- $\mathbb{E}_{\tilde{\mathbf{x}}}[-\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)$ by Stein's lemma $\Rightarrow \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^{\natural}] = \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}\Delta$

- concentration:

$$\left\| \mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] \right\|_{\text{F}} \lesssim \sqrt{d} \epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, w.h.p. \Rightarrow \text{control } \mathbf{G}^{\natural}$$

Theorem (Linear convergence rate)

Under (**Spectral-initialization**) and $\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), w.h.p$$

- holds for standard gradient update, but requires more assumptions.
- $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \epsilon \|\Delta\|_{\text{op}}, w.h.p.$

A bit proof sketch at

- recover at initialization:

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \|\mathbf{A}_0 \mathbf{B}_0 - \gamma \mathbf{G}^{\natural}\|_{\text{F}} + \text{concentration on } \mathbf{G}^{\natural} + \rho \lambda_{r^*}^*, w.h.p$$

- $\mathbb{E}_{\tilde{\mathbf{x}}}[-\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)$ by Stein's lemma $\Rightarrow \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^{\natural}] = \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}\Delta$
- concentration:

$$\left\| \mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] \right\|_{\text{F}} \lesssim \sqrt{d} \epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, w.h.p. \Rightarrow \text{control } \mathbf{G}^{\natural}$$

Theorem (Linear convergence rate)

Under (**Spectral-initialization**) and $\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), w.h.p$$

- holds for standard gradient update, but requires more assumptions.
- $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \epsilon \|\Delta\|_{\text{op}}, w.h.p.$

A bit proof sketch at

- recover at initialization:

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \|\mathbf{A}_0 \mathbf{B}_0 - \gamma \mathbf{G}^{\natural}\|_{\text{F}} + \text{concentration on } \mathbf{G}^{\natural} + \rho \lambda_{r^*}^*, w.h.p$$

- $\mathbb{E}_{\tilde{\mathbf{x}}}[-\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)$ by Stein's lemma $\Rightarrow \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^{\natural}] = \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}\Delta$

- concentration:

$$\left\| \mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] \right\|_{\text{F}} \lesssim \sqrt{d} \epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, w.h.p. \Rightarrow \text{control } \mathbf{G}^{\natural}$$

Theorem (Linear convergence rate)

Under (**Spectral-initialization**) and $\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), w.h.p$$

- holds for standard gradient update, but requires more assumptions.
- $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \epsilon \|\Delta\|_{\text{op}}, w.h.p.$

A bit proof sketch at

- recover at initialization:

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \|\mathbf{A}_0 \mathbf{B}_0 - \gamma \mathbf{G}^{\natural}\|_{\text{F}} + \text{concentration on } \mathbf{G}^{\natural} + \rho \lambda_{r^*}^*, w.h.p$$

- $\mathbb{E}_{\tilde{\mathbf{x}}}[-\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)$ by Stein's lemma $\Rightarrow \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^{\natural}] = \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}\Delta$
- concentration:

$$\left\| \mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] \right\|_{\text{F}} \lesssim \sqrt{d} \epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, w.h.p. \Rightarrow \text{control } \mathbf{G}^{\natural}$$

Theorem (Linear convergence rate)

Under (**Spectral-initialization**) and $\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}$ for gradient update (adding preconditioners), choose constant step-size $\eta < 1$, we have

$$\|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}} \lesssim \left(1 - \frac{\eta}{4}\right)^t \lambda_{r^*}(\Delta), w.h.p$$

- holds for standard gradient update, but requires more assumptions.
- $\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \epsilon \|\Delta\|_{\text{op}}, w.h.p.$

A bit proof sketch at

- recover at initialization:

$$\|\mathbf{A}_0 \mathbf{B}_0 - \Delta\|_{\text{F}} \leq \|\mathbf{A}_0 \mathbf{B}_0 - \gamma \mathbf{G}^{\natural}\|_{\text{F}} + \text{concentration on } \mathbf{G}^{\natural} + \rho \lambda_{r^*}^*, w.h.p$$

- $\mathbb{E}_{\tilde{\mathbf{x}}}[-\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}(\mathbf{A}_t \mathbf{B}_t - \Delta)$ by Stein's lemma $\Rightarrow \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{G}^{\natural}] = \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] = \frac{1}{2}\Delta$
- concentration:

$$\left\| \mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - \mathbb{E}_{\tilde{\mathbf{x}}}[\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}}] \right\|_{\text{F}} \lesssim \sqrt{d} \epsilon \|\mathbf{A}_t \mathbf{B}_t - \Delta\|_{\text{F}}, w.h.p. \Rightarrow \text{control } \mathbf{G}^{\natural}$$

Proof of sketch on $\mathbf{A}_t \mathbf{B}_t - \Delta$

$$\begin{aligned} \|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_F &\lesssim \|\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - c_H(\mathbf{A}_t \mathbf{B}_t - \Delta)\|_F \text{ [concentration+Hermite]} \\ &\quad + (1 - \eta) \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right\|_F \\ &\quad + \left\| \left(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \right) (\mathbf{A}_t \mathbf{B}_t - \Delta) \left(\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right) \right\|_F \\ &\quad + \text{other projections} \end{aligned}$$

$$\mathbf{L} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t} & \mathbf{0}_{d \times r} \\ \mathbf{0}_{k \times r} & \mathbf{V}_{\mathbf{B}_t} \end{bmatrix} \in \mathbb{R}^{(d+k) \times 2r},$$

then $\mathbf{L} \mathbf{L}^\top$ is a projection matrix, $\mathbf{I}_{d+k} - \mathbf{L} \mathbf{L}^\top = \mathbf{L}_\perp \mathbf{L}_\perp^\top$

◦ transformed to lower bound $\left\| \mathbf{L}_\perp^\top \Delta \mathbf{L} \right\|_F^2$

◦ upper bound $\left\| \mathbf{L}_\perp^\top \mathbf{U} \right\|_{op} < 1$ by Wedin's sin- θ theorem

Proof of sketch on $\mathbf{A}_t \mathbf{B}_t - \Delta$

$$\begin{aligned}
 \|\mathbf{A}_{t+1} \mathbf{B}_{t+1} - \Delta\|_F &\lesssim \|\mathbf{J}_{\mathbf{W}_t}^{\text{GLM}} - c_H(\mathbf{A}_t \mathbf{B}_t - \Delta)\|_F \text{ [concentration+Hermite]} \\
 &\quad + (1 - \eta) \left\| \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top (\mathbf{A}_t \mathbf{B}_t - \Delta) \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right\|_F \\
 &\quad + \left\| \left(\mathbf{I}_d - \mathbf{U}_{\mathbf{A}_t} \mathbf{U}_{\mathbf{A}_t}^\top \right) (\mathbf{A}_t \mathbf{B}_t - \Delta) \left(\mathbf{I}_k - \mathbf{V}_{\mathbf{B}_t} \mathbf{V}_{\mathbf{B}_t}^\top \right) \right\|_F \\
 &\quad + \text{other projections}
 \end{aligned}$$

$$\mathbf{L} = \begin{bmatrix} \mathbf{U}_{\mathbf{A}_t} & \mathbf{0}_{d \times r} \\ \mathbf{0}_{k \times r} & \mathbf{V}_{\mathbf{B}_t} \end{bmatrix} \in \mathbb{R}^{(d+k) \times 2r},$$

then $\mathbf{L} \mathbf{L}^\top$ is a projection matrix, $\mathbf{I}_{d+k} - \mathbf{L} \mathbf{L}^\top = \mathbf{L}_\perp \mathbf{L}_\perp^\top$

- transformed to lower bound $\left\| \mathbf{L}_\perp^\top \Delta \mathbf{L} \right\|_F^2$
- upper bound $\left\| \mathbf{L}_\perp^\top \mathbf{U} \right\|_{op} < 1$ by Wedin's sin- θ theorem

Takeaway messages

- [arXiv: 2502.01235](#) and [code](#)

Model	Results	Algorithm	Initialization	Conclusion
Linear	Theorem 3.1	GD	(LoRA-init)	Subspace alignment of B_t
	Theorem 3.2	GD	(LoRA-init)	Subspace alignment of A_t
	Proposition 3.3	GD	(Spectral-init)	$\ A_0 B_0 - \Delta\ _F$ is small
	Theorem 3.5	GD	(Spectral-init)	Linear convergence of $\ A_t B_t - \Delta\ _F$
	Theorem 3.6	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
Nonlinear	Theorem 4.3	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
	Theorem C.15	Smoothed Precondition GD	(Spectral-init)	Better convergence performance with less assumptions

- subspace alignment: G^\natural and $(A_t, B_t) \Rightarrow$ theory-grounded algorithm design
- clarification on gradient alignment based algorithms

Target

- How to handle **nonlinearity** at a theoretical level (e.g., training dynamics)
- How to precisely and efficiently approximate **nonlinearity** at a practical level under theoretical guidelines

Takeaway messages

- [arXiv: 2502.01235](#) and [code](#)

Model	Results	Algorithm	Initialization	Conclusion
Linear	Theorem 3.1	GD	(LoRA-init)	Subspace alignment of B_t
	Theorem 3.2	GD	(LoRA-init)	Subspace alignment of A_t
	Proposition 3.3	GD	(Spectral-init)	$\ A_0 B_0 - \Delta\ _F$ is small
	Theorem 3.5	GD	(Spectral-init)	Linear convergence of $\ A_t B_t - \Delta\ _F$
	Theorem 3.6	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
Nonlinear	Theorem 4.3	Precondition GD	(Spectral-init)	Linear convergence rate independent of $\kappa(\Delta)$
	Theorem C.15	Smoothed Precondition GD	(Spectral-init)	Better convergence performance with less assumptions

- subspace alignment: G^\natural and $(A_t, B_t) \Rightarrow$ theory-grounded algorithm design
- clarification on gradient alignment based algorithms

Target

- How to handle **nonlinearity** at a theoretical level (e.g., training dynamics)
- How to precisely and efficiently approximate **nonlinearity** at a practical level under theoretical guidelines



Soufiane Hayou, Nikhil Ghosh, and Bin Yu.

LoRA+: Efficient low rank adaptation of large models.

arXiv preprint arXiv:2402.12354, 2024.



Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.

LoRA: Low-rank adaptation of large language models.

In *International Conference on Learning Representations*, 2022.



Sham M Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai.

Efficient learning of generalized linear and single index models with isotonic regression.

In *Advances in Neural Information Processing Systems*, 2011.



Dominik Stöger and Mahdi Soltanolkotabi.

**Small random initialization is akin to spectral learning:
Optimization and generalization guarantees for overparameterized
low-rank matrix reconstruction.**

In *Advances in Neural Information Processing Systems*, pages
23831–23843, 2021.



Shaowen Wang, Linxi Yu, and Jian Li.

LoRA-GA: Low-rank adaptation with gradient approximation.

In *Advances in Neural Information Processing Systems*, 2024.



Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Sham M Kakade.

Finite-sample analysis of learning high-dimensional single relu neuron.

In *International Conference on Machine Learning*, pages 37919–37951, 2023.