

A survey on LoRA for Continual Learning

A simple survey

Yang, Junran¹ Deepseek²

¹School of AI
Tianjin Univeristy

²Deepseek Co., Ltd.

Group Meeting, March. 27th

Table of Contents

- 1 Gradient Projection Memory for Continual Learning
- 2 Gated Integration of Low-Rank Adaptation for Continual Learning of Large Language Models

Table of Contents

- 1 Gradient Projection Memory for Continual Learning
- 2 Gated Integration of Low-Rank Adaptation for Continual Learning of Large Language Models

Subspace

Lemma (spanning space)

n linearly independent vectors can span an n -dimensional space

The spanning space could be written as

$$\mathcal{F} = \{ \mathbf{u}_1\beta_1 + \mathbf{u}_2\beta_2 + \cdots + \mathbf{u}_n\beta_n \mid \beta_1, \beta_2, \cdots, \beta_n \in \mathbb{R} \} \quad (1)$$

The elements in the set are the vectors in the space, which is also the elements that can be linearly represented by $(\mathbf{u}_1, \mathbf{u}_2, \cdots, \mathbf{u}_n)$.

Projection Matrix

given a direction vector \mathbf{u} s.t. $\|\mathbf{u}\| = 1$, a random vector \mathbf{x} , the formula

$$\mathbf{proj}_{\mathbf{u}}(\mathbf{x}) = (\mathbf{u}^T \mathbf{x}) \mathbf{u}, \quad (2)$$

$$\text{where } \mathbf{u}^T \mathbf{x} \in \mathbb{R}, \mathbf{u} \in \mathbb{R}^n, \mathbf{x} \in \mathbb{R}^n \quad (3)$$

represents \mathbf{x} 's projection on \mathbf{u} . Since $(\mathbf{u}^T \mathbf{x})$ is scalar, Equation (2) is commutative and could be written as:

$$\mathbf{proj}_{\mathbf{u}}(\mathbf{x}) = \mathbf{u} \mathbf{u}^T \mathbf{x} \quad (4)$$

$$= \mathbf{P} \mathbf{x} \quad (5)$$

$$\text{where } \mathbf{P} = \mathbf{u} \mathbf{u}^T, \quad \mathbf{P} \in \mathbb{R}^{n \times n} \quad (6)$$

\mathbf{P} is the **Projection Matrix**

Projection Matrix

now $\mathbf{M} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n]$, $\mathbf{x} = [x_1, x_2, \dots, x_n]$, $\mathbf{M} \in \mathbb{R}^{m \times n}$, and \mathbf{u} are all orthonormal.

$$\mathbf{M}^T \mathbf{x} = \begin{bmatrix} \mathbf{u}_1^T \mathbf{x} = \alpha_1 \\ \mathbf{u}_2^T \mathbf{x} = \alpha_2 \\ \vdots \\ \mathbf{u}_n^T \mathbf{x} = \alpha_n \end{bmatrix} \quad (7)$$

means \mathbf{x} 's projection on every vectors in \mathbf{M} , then

$$\mathbf{M} \mathbf{M}^T \mathbf{x} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n] * \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} \quad (8)$$

$$= \mathbf{u}_1 \alpha_1 + \mathbf{u}_2 \alpha_2 + \dots + \mathbf{u}_n \alpha_n \quad (9)$$

which is a vector in the subspace spanned by the vectors in \mathbf{M} .

Projection Matrix

a simple example: 3-D vector projected into a plane spanned by 2 3-D vector.

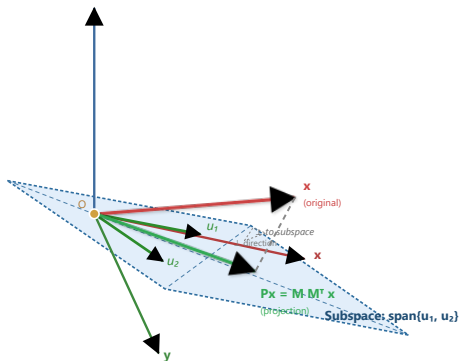


Figure: vector projection in subspace

Singular Value Decomposition (SVD)

Singular Value Decomposition (SVD) is a fundamental matrix factorization technique in linear algebra. For any real (or complex) matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, SVD decomposes it as:

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \quad (10)$$

where:

- $\mathbf{U} \in \mathbb{R}^{m \times m}$ is an orthogonal matrix whose columns are called left singular vectors,
- $\mathbf{V} \in \mathbb{R}^{n \times n}$ is an orthogonal matrix whose columns are called right singular vectors,
- $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ is a diagonal matrix containing non-negative singular values.

Singular Value Decomposition (SVD)

Rewritten in vectors

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m], \quad \mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n], \quad \mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \\ 0 & 0 & \cdots & 0 \end{bmatrix} \quad (11)$$

$$\mathbf{A} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m] \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \\ 0 & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \quad (12)$$

, where the rank of $\mathbf{u}_i \mathbf{v}_i^T$ is 1

Low-rank Approximation

assume that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ the best rank- k approximation A_k for $\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is

$$\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T \quad (13)$$

Gradient of Matrix

$$L = \frac{1}{2} \|\mathbf{W}\mathbf{x} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{W}\mathbf{x} - \mathbf{y})^T (\mathbf{W}\mathbf{x} - \mathbf{y}) \quad (14)$$

$$\nabla_{\mathbf{W}} L = (\mathbf{W}\mathbf{x} - \mathbf{y})\mathbf{x}^T \quad (15)$$

highlights

- neither store old gradient directions
- nor store old examples for generating reference directions
- partitions the entire gradient space of the weights into two orthogonal subspaces: Core Gradient Space (CGS) and Residual Gradient Space (RGS)

Task 1

- 1 training one task 1 normally, obtain a set of parameters $\{\mathbb{W}'_l\}_{l=1}^L$
- 2 randomly samples n_s input data for each layers $[\mathbf{x}'_{1,1}, \mathbf{x}'_{2,1}, \dots, \mathbf{x}'_{n_s,1}]$
- 3 construct representation matrix $\mathbf{R}'_1 = [\mathbf{x}'_{1,1}, \mathbf{x}'_{2,1}, \dots, \mathbf{x}'_{n_s,1}]$, perform SVD on $\mathbf{R}'_1 = \mathbf{U}'_1 \Sigma'_1 (\mathbf{V}'_1)^T$ followed by its k-rank approximation $(\mathbf{R}'_1)_k$ according to the given threshold ϵ'_{th}

$$\|(\mathbf{R}'_1)_k\|_F^2 \geq \epsilon'_{th} \|\mathbf{R}'_1\|_F^2 \quad (16)$$

$(\mathbf{U}'_1)_k = [\mathbf{u}'_{1,1}, \mathbf{u}'_{2,1}, \dots, \mathbf{u}'_{k,1}] = \mathbf{M}'$, where \mathbf{M}' is the projection matrix would be used later.

Task 2 and successive tasks

On Initialization:

$$\mathbf{W}'_2 \leftarrow \mathbf{W}'_1 \quad (17)$$

Before taking gradient step, bases of the CGS are retrieved from GPM

$$\nabla_{\mathbf{W}'_2} L_2 \leftarrow \nabla_{\mathbf{W}'_2} L_2 - \nabla_{\mathbf{W}'_2} L_2 \mathbf{M}' (\mathbf{M}')^T \quad (18)$$

such that

$$\hat{\mathbf{W}}'_2 \mathbf{x}_{.,1} = (\mathbf{W}'_2 - \epsilon \nabla_{\mathbf{W}'_2} L_2) \mathbf{x}_{.,1} = \mathbf{W}'_1 \mathbf{x}_{.,1} \quad (19)$$

Task 2 and successive tasks

We sample from task 2 input data $\mathbf{x}_{\cdot,2}$, and Construct $\mathbf{R}'_2 = [\mathbf{x}'_{1,2}, \mathbf{x}'_{2,2}, \dots, \mathbf{x}'_{n_s,2}]$. Also, we perform SVD again to expand the subspace.

$$\hat{\mathbf{R}}'_2 = \mathbf{R}'_2 - \mathbf{M}'(\mathbf{M}')^T(\mathbf{R}'_2) = \mathbf{R}'_2 - \mathbf{R}'_{2,\text{Proj}} \quad (20)$$

$$\hat{\mathbf{R}}'_2 = (\hat{\mathbf{U}}'_2 \hat{\Sigma}'_2 (\hat{\mathbf{V}}'_2)^T) \quad (21)$$

and update \mathbf{M}'

$$\mathbf{M}' \leftarrow [\mathbf{M}', \mathbf{u}'_{1,2}, \mathbf{u}'_{2,2}, \dots, \mathbf{u}'_{k,2}] \quad (22)$$

Experimental Setup

Benchmarks. PMNIST (10 tasks), Split CIFAR-100 (10 tasks), Split minilmageNet (20 tasks), and 5-Datasets.

Backbones. 2-layer MLP for PMNIST; 5-layer AlexNet for Split CIFAR-100; reduced ResNet18 for minilmageNet and 5-Datasets.

Protocol. PMNIST uses *single-head* inference without task identity; other benchmarks use *multi-head* inference.

Baselines. OGD, OWM, GEM, A-GEM, ER-Reservoir, EWC, HAT, plus a multitask upper bound.

Training & Metric. Models are trained with SGD. Main metrics are average accuracy (ACC) and backward transfer (BWT), where more negative BWT means more forgetting.

Results and Discussion

Method	PMNIST		CIFAR-100		minilmaGeNet		5-Datasets	
	ACC	BWT	ACC	BWT	ACC	BWT	ACC	BWT
OWM	90.71	-0.01	50.94	-0.30	–	–	–	–
EWC	89.97	-0.04	68.80	-0.02	52.01	-0.12	88.64	-0.04
HAT	–	–	72.06	-0.00	59.78	-0.03	91.32	-0.01
A-GEM	83.56	-0.14	63.98	-0.15	57.24	-0.12	84.04	-0.12
ER Res	87.24	-0.11	71.73	-0.06	58.94	-0.07	88.31	-0.04
GPM	93.91	-0.03	72.48	-0.00	60.41	-0.00	91.22	-0.01

- GPM achieves the best ACC on PMNIST, CIFAR-100, and minilmaGeNet, while staying competitive on 5-Datasets.
- Forgetting is consistently very small: nearly zero BWT on CIFAR-100 and minilmaGeNet.
- The paper also reports lower memory use than replay-based methods and fast training; e.g., total training time is 245s on PMNIST, 770s on CIFAR-100, 3387s on minilmaGeNet, and 5008s on 5-Datasets.

Table of Contents

- 1 Gradient Projection Memory for Continual Learning
- 2 Gated Integration of Low-Rank Adaptation for Continual Learning of Large Language Models

Novelty

- expand new LoRA branch to learn new task
- introduce gating modules to integrate new and old LoRA branches
- minimize the influence from the new LoRA branch to old tasks

Architecture of Gating Module The formula for the gating module of the i^{th} task could be written as:

$$\mathbf{p}_{i,0} = \mathbf{p}_0 = \text{Pool}(\text{Token}(\mathbf{x})), \quad (23)$$

$$\mathbf{p}_{i,l} = \sigma(\mathbf{G}_{i,l}\mathbf{p}_{i,l-1}), \quad l \in \{1, 2, \dots, L\}, \quad (24)$$

$$g_i(\mathbf{x}) = f(\mathbf{G}_{i,L+1}\mathbf{p}_{i,L}). \quad (25)$$

$$\mathbf{G}_{i,l} \in \mathbb{R}^{n \times n}, \quad \mathbf{G}_{i,L+1} \in \mathbb{R}^n, \quad \mathbf{p} \in \mathbb{R}^n, \quad f: \mathbb{R} \rightarrow [0, 1]$$

We hope that for $\mathbf{x} \in \mathcal{T}_i$, $g_i(\cdot) = 1$, for $\mathbf{x} \notin \mathcal{T}_i$, $g_i(\cdot) = 0$

Pipeline

- 1 training on task 1 freely.
- 2 initialize the gating module on task 2 with some constraints, update the gating module also with some constraints.
- 3 repeat the procedure in task 2 for the following tasks.

Constraint on Initialization for Task 2 and the Successive

$$\text{Init}(\mathbf{G}_{t,l}) \leftarrow \mathbf{G}_{t-1,l}, l \in [1, L] \quad (26)$$

$$\text{Init}(\mathbf{G}_{t,L+1}) \leftarrow \text{Init}(\mathbf{G}_{t,L+1}) - \mathbf{M}_{t,L+1} \mathbf{M}_{t,L+1}^T \text{Init}(\mathbf{G}_{t,L+1}) \quad (27)$$

such that

$$\mathbf{M}_{t,L+1}^T (\text{Init}(\mathbf{G}_{t,L+1}) - \mathbf{M}_{t,L+1} \mathbf{M}_{t,L+1}^T \text{Init}(\mathbf{G}_{t,L+1})) \quad (28)$$

$$= \mathbf{M}_{t,L+1}^T (\mathbf{I} - \mathbf{M}_{t,L+1} \mathbf{M}_{t,L+1}^T) \text{Init}(\mathbf{G}_{t,L+1}) \quad (29)$$

$$= (\mathbf{I} - \mathbf{M}_{t,L+1}^T \mathbf{M}_{t,L+1}) \mathbf{M}_{t,L+1}^T \text{Init}(\mathbf{G}_{t,L+1}) \quad (30)$$

$$= \mathbf{0} \quad (31)$$

Constraint on Updating

$$\Delta \mathbf{G}_{t,l} \leftarrow \Delta \mathbf{G}_{t,l} - \mathbf{M}_{t,l} \mathbf{M}_{t,l}^T \Delta \mathbf{G}_{t,l}, \quad l \in [1, L + 1] \quad (32)$$

No Constraint on LoRA parameters, only on gating modules.

Benchmarks and Settings

- Evaluated on two continual learning benchmarks: **SuperNI** and **Long Sequence**.
- Each benchmark contains **15 sequential tasks**, and two different task orders are considered.
- Evaluation metrics are **Average Performance (AP)** and **Forgetting (FT)**.
- Experiments are conducted on multiple backbones, including **T5-Large**, **T5-XL**, **Llama-2-7B/13B**, and **Llama-3-8B**.

Main Results

- GainLoRA consistently outperforms strong LoRA-based CL baselines such as **O-LoRA** and **InfLoRA**.
- On **T5-Large**, GainLoRA substantially improves AP and reduces FT across all task orders.
- The same advantage remains when scaling to larger backbones, showing good **robustness across model sizes**.

Experiment

Method	Order 1		Order 2		Order 3		Order 4	
	AP \uparrow	FT \downarrow	AP \uparrow	FT \downarrow	AP \uparrow	FT \downarrow	AP \uparrow	FT \downarrow
O-LoRA	26.37	19.15	32.83	11.99	70.98	3.69	71.21	4.03
GainLoRA (O-LoRA)	47.84	2.26	46.84	2.91	73.37	3.02	76.01	2.49
InfLoRA	39.78	7.64	39.57	8.93	75.15	4.19	75.79	3.47
GainLoRA (InfLoRA)	46.21	2.40	46.44	2.61	78.01	0.77	77.54	1.25

Observation: GainLoRA consistently improves *average performance* and reduces *forgetting* over the corresponding LoRA-based continual learning baselines.

Takeaway

- GainLoRA achieves **better overall continual learning performance** than existing state-of-the-art methods.
- Its advantages are consistent across **different task orders, benchmarks, and backbone scales**.
- Therefore, the paper shows that **how to integrate LoRA branches** is a key factor in continual learning of LLMs.

Limitations

- The imposed constraints may **accumulate as the number of tasks increases**, which can hinder future learning.
- The method is mainly evaluated on **non-overlapping task settings**; its behavior in more complex scenarios with **task overlap or blurry boundaries** still needs further study.