



Rethinking Gradient Projection Continual Learning: Stability / Plasticity Feature Space Decoupling

Zhen Zhao¹, Zhizhong Zhang^{1,†}, Xin Tan¹, Jun Liu², Yanyun Qu³, Yuan Xie^{1,†}, Lizhuang Ma¹ ¹School of Computer Science and Technology East China Normal University, Shanghai, China ²Tencent Youtu Lab ³School of Informatics, Xiamen University, Fujian, China {51255901056}@stu.ecnu.edu.cn, {zzzhang,xtan,yxie,lzma}@cs.ecnu.edu.cn {junsenselee}@tencent.com, {yygu}@xmu.edu.cn

CVPR 2023



- As humans live, they acquire more knowledges through new experiences.
- In machine perspective, continual learning focuses on how to learn new knowledge from new incoming data with preserving previous knowledge.



Online vs Offline

- Online: allow only once to see the incoming data except for examples in the memory.
- Offline: allow to see the examples in current task and episodic memory many times.



Offline Continual Learning Setting



Online Continual Learning Setting

ParN C 模式识别与神经计算研究组 PAttern Recognition and NEural Computing

Previous Works for Continual Learning Methods

- Regularization based CL
 - **Trade-off:** keeping weights for maintaining previous knowledges vs. changing weights for learning new knowledges.
- Parameter isolation based CL
 - To prevent forgetting previous tasks, an intuition is to construct sufficiently large models, and for each task construct a subset of the larger model. This approach can be achieved by fixing Backbone and adding new branches for new tasks
- Rehearsal/Replay based CL
 - Assume that we can use small amount of data from previous tasks
 - Key Point: How to get informative data from the previous tasks?





Methods of Gradient Projection

Core idea: Project the gradient of the new task onto the orthogonal space of the input space of the old task



Former Methods

Low loss contour on previous tasks Low loss contour on current task

Methods of Gradient Projection

Core idea: Project the gradient of the new task into the **orthogonal space** of the **input space** of the old task

$$egin{aligned} oldsymbol{R}_1^l &= \left[ig(oldsymbol{X}_{1,1}^lig)^T,ig(oldsymbol{X}_{2,1}^lig)^T,\ldots,ig(oldsymbol{X}_{n_s,1}^lig)^T
ight] \ oldsymbol{R}_1^l &= oldsymbol{U}_1^l oldsymbol{\Sigma}_1^lig(oldsymbol{V}_1^lig)^T \ oldsymbol{M}^l &= ig[oldsymbol{u}_{1,1}^l,oldsymbol{u}_{2,1}^l,\ldots,oldsymbol{u}_{k,1}^lig] \
abla_{oldsymbol{W}_2^l}L_2 &=
abla_{oldsymbol{W}_2^l}L_2 - oldsymbol{M}^lig(oldsymbol{M}_1^lig)^Tig(
abla_{oldsymbol{W}_2^l}L_2ig) \end{aligned}$$

Addressing catastrophic forgetting

$$\begin{split} \boldsymbol{\theta}_{t}^{l} \boldsymbol{x}_{j,i}^{l} &= (\boldsymbol{\theta}_{t-1}^{l} + \Delta \boldsymbol{\theta}_{t-1}^{l}) \boldsymbol{x}_{j,i}^{l} \\ &= \boldsymbol{\theta}_{t-1}^{l} \boldsymbol{x}_{j,i}^{l} + \Delta \boldsymbol{\theta}_{t-1}^{l} \boldsymbol{x}_{j,i}^{l} \\ &= \boldsymbol{\theta}_{t-1}^{l} \boldsymbol{x}_{j,i}^{l}. \end{split}$$





Feature Space Continual Learning Paradigm





- OWM
- GPM
- TRGP
- Adam-NSCL
- AdNS

• • •



Stability-Plasticity Dilemma in Gradient Projection



 ϵ increases, stability increases, plasticity decreases



 ζ increases, plasticity increases, stability decreases

Motivation



Space Decoupling: Stability and Plasticity

Inspired by GPM: When updating the feature space, GPM directly **removes duplicate bases** between tasks.

At the end of the task 2 training, we update the GPM with new task-specific bases (of CGS). To obtain such bases, we construct $\mathbf{R}_2^l = [\mathbf{x}_{1,2}^l, \mathbf{x}_{2,2}^l, ..., \mathbf{x}_{n_s,2}^l]$ using data from task 2 only. However, before performing SVD and subsequent k-rank approximation, from \mathbf{R}_2^l we eliminate the common directions (bases) that are already present in the GPM so that newly added bases are unique and orthogonal to the existing bases in the memory. To do so, we perform the following step :

$$\hat{R}_{2}^{l} = R_{2}^{l} - M^{l} (M^{l})^{T} (R_{2}^{l}) = R_{2}^{l} - R_{2,Proj}^{l}.$$
(8)

The authors think Task-Shared bases are more important than Task-Specific bases



Space Decoupling Algorithm To Decouple The Feature Space



Space Decoupling Algorithm To Decouple The Feature Space

子空间的交和直和:

$$\mathcal{P} \cap \mathcal{Q} = \{ \alpha \mid \alpha \in \mathcal{P}, \alpha \in \mathcal{Q} \}$$

 $\mathcal{P} + \mathcal{Q} = \{ \alpha + \beta \mid \alpha \in \mathcal{P}, \beta \in \mathcal{Q} \}$
稳定性空间:
 $\mathcal{I}(t) = \sum_{1 < i \le t} \overline{S}(i-1) \cap S_i$
(特征子空间的交的直和) → Task-Shared
 $\mathbf{R}(t) = \overline{S}(t) - \operatorname{Proj}_{\mathcal{I}(t)}(\overline{S}(t))$

(稳定性空间在特征空间下的正交补) → Task-Specific



Empirical Evidence

Defining Gradient Update Influence Score:

• The degree of perturbation of the gradient with respect to the feature subspace of the old task.

$$\omega(\boldsymbol{g}) = \sum_{j=1}^{t} \left\| \operatorname{Proj}_{\mathcal{S}_{j}}(\boldsymbol{g}) \right\|_{F}^{2}$$

Further define the Influence Score of the feature subspace:

• The sum of the Influence Scores of the gradient components in that subspace.



Figure 5. Mean knowledge interference of $\mathcal{I}(t)$ and $\mathcal{R}(t)$.

 $\Omega(\mathcal{I}(t)) = \frac{\sum_{i} \omega\left(\boldsymbol{g}_{i}^{\mathcal{I}}\right)}{\dim(\mathcal{I}(t))} \qquad \Omega(\mathcal{R}(t)) = \frac{\sum_{i} \omega\left(\boldsymbol{g}_{i}^{\mathcal{R}}\right)}{\dim(\mathcal{R}(t))}$



Space Decoupling Algorithm To Decouple The Feature Space



 $\hat{\boldsymbol{I}}(t) = \mathcal{A}\left(\boldsymbol{I}(t); \boldsymbol{\epsilon}^{\mathcal{I}}\right)$ $\hat{\boldsymbol{R}}(t) = \mathcal{A}\left(\boldsymbol{R}(t); \boldsymbol{\epsilon}^{\mathcal{R}}\right)$

 $\nabla_{\boldsymbol{\theta}} \mathcal{L}_{t+1} = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{t+1}$ $- \nabla_{\boldsymbol{\theta}} \mathcal{L}_{t+1} \left(1 - \zeta^{\mathcal{I}} \right) \hat{\boldsymbol{I}}(t) (\hat{\boldsymbol{I}}(t))^{T}$ $- \nabla_{\boldsymbol{\theta}} \mathcal{L}_{t+1} \left(1 - \zeta^{\mathcal{R}} \right) \hat{\boldsymbol{R}}(t) (\hat{\boldsymbol{R}}(t))^{T}$

Experiment



Comparison with SOTA CL Methods

Model	Venue	20-split-MiniImageNet		20-split-CIFAR-100		10-split-CIFAR-100	
		ACC(%)	BWT(%)	ACC(%)	BWT(%)	ACC(%)	BWT (%)
LWF [22]	PAMI'17	57.63	-8.72	74.38	-9.11	70.7	-6.27
EWC [17]	PANS'17	52.01	-12	71.66	-3.72	70.77	-2.83
MAS [2]	ECCV'18	50.12	-5.82	63.84	-6.29	66.93	-4.03
MUC-MAS [24]	ECCV'20	46.24	-3.79	67.22	-5.72	63.73	-3.38
GEM [25]	NIPS'17	-	_	68.89	-1.2	49.48	2.77
A-GEM [5]	ICLR'18	57.24	-12	61.91	-6.88	49.57	-1.13
*AdNS [18]	ECCV'22	60.82	-4.24	77.33	-3.25	77.21	-2.32
OWM [41]	NMI'19	47.48	-8.57	68.47	-3.37	68.89	-1.88
GPM [31]	ICLR'21	60.41 ± 0.61	-0.7 ± 0.4	77.53 ± 0.83	-0.97 ± 0.59	$72.48 {\pm} 0.4$	-0.9±0
Adam-NSCL [38]	CVPR'21	59.07±1.1	-4.9 ± 1.32	75.81±0.93	$-3.98 {\pm} 0.85$	74.97 ± 1.15	-2.64 ± 0.91
TRGP [23]	ICLR'22	63.51±0.74	-0.76 ± 0.25	80.68±0.7	$-0.87 {\pm} 0.46$	$74.46 {\pm} 0.32$	-0.9 ± 0.01
GPM+SD		$62.39 {\pm} 0.56$	-0.61 ± 0.11	$80.71 {\pm} 0.82$	-0.73 ± 0.27	73.53 ± 0.44	-0.83 ± 0.31
Adam-NSCL+SD		$60.38 {\pm} 0.75$	-4.81 ± 1	76.5 ± 1.02	-3.99 ± 0.96	75.97±0.66	$-2.88 {\pm} 0.89$
TRGP+SD		65.8±0.16	$-0.49 {\pm} 0.08$	$83.84{\pm}0.12$	-0.72 ± 0.2	$75.5 {\pm} 0.35$	$-0.96 {\pm} 0.09$

Experiment



Model	I/R-A	\mathcal{I}/\mathcal{R} -P	ACC(%)	BWT(%)
			77.53 ± 0.83	-0.97±0.59
GPM	\checkmark		$79.99 {\pm} 0.48$	-0.7 ± 0.45
		\checkmark	78.41±1.3	-1.44 ± 0.65
	\checkmark	\checkmark	80.71±0.82	-0.73 ± 0.27
4.			$80.68 {\pm} 0.7$	-0.87±0.46
TRGP	\checkmark		83.21±0.36	-0.4 ± 0.02
		\checkmark	81.15±1.22	-1.24 ± 0.71
	\checkmark	\checkmark	83.84 ± 0.12	-0.72 ± 0.2

Computational Complexity Analysis

Detecto	Methods						
Datasets	GPM	GPM+SD	TRGP	TRGP+SD	Adam-NSCL	Adam-NSCL+SD	
10-split-CIFAR-100	0.25	0.27	0.41	0.45	2.71	2.93	
20-split-CIFAR-100	0.31	0.36	0.48	0.53	4.14	4.53	
20-split-MiniImageNet	0.58	0.66	0.81	0.92	11.28	12.95	

Comparison of feature space dimension





Experiment



Forcing GPM to have the same feature dimension with GPM+SD.



Comparisons with different subspaces construction.



Stability and Plasticity analysis.



Figure 8. Stability and Plasticity analysis. Experiments are conducted by GPM+SD on 10-split-CIFAR-100. Red and blue values denote the comparison between \mathcal{R} and \mathcal{I} .

Thanks