A Simple Finetuning Strategy Based on Bias-Variance Ratios of Layer-Wise Gradients



¹ Institute of Science Tokyo, Japan, ² Denso IT Laboratory, Japan, ³ National Institute of Informatics, Japan

1. Background

Institute of

SCIENCE

ΤΟΚΥΟ

Finetuning is a common deep learning technique for adapting a pretrained model to a new, related tasks.



2. Which layers should we finetune?

- Performance strongly depends on the choice of tuning layers. ullet
- Which layers should be tuned depends on the target task. \bullet



Task			
Pet	15.1 %	10.3 %	8.2 %
Food	18.0 %	15.8 %	35.5 %
CIFAR100	22.1 %	23.1 %	56.9 %

To propose a simple finetuning strategy that is independent of the target Our aim task by identifying the layers that should be tuned based on some score.

3. Proposed method < BVG-LS: Bias-Variance Guided Layer Selection >







(Test error rate)



4. Evaluation



- Pretrained model: WideResNet-50-2, on ImageNet
- We finetuned pretrained model on 7 small dataset

* WRN-50-2 consists of six layers:



A) Test error rate

Dataset	Full finetuning	L2-6	Partial fi L3-6	inetuning L4-6	L5-6	Linear probing	BVG-LS (ours)	_	2.0			Full finetuning BVG-LS (ours)	1.0		D	TD		
Flowers Pets DTD	$ \begin{array}{c c} 7.65 \\ 15.14 \\ 40.37 \end{array} $	$ \begin{array}{c c} 6.67 \\ 13.24 \\ 38.78 \end{array} $	$6.37 \\ 11.96 \\ 37.23$	$\frac{5.88}{10.29}$ 34.57	4.61 8.33 31.60	$ \begin{array}{r} 12.84 \\ \underline{8.22} \\ 35.32 \end{array} $	6.18 6.26 29.20	_	1.8 % 1.6				- 8.0 - 8.0					
Aircraft Food SUN CIFAR100	$ 19.34 \\ 18.01 \\ 35.54 \\ 22.07 $	$18.32 \\ 17.05 \\ 35.18 \\ 21.65$	16.28 16.52 34.23 21.89	$\frac{16.22}{15.84}$ 33.21 23.06	17.72 16.97 31.77 31.93	$49.40 \\ 35.45 \\ 36.72 \\ 56.87$	$14.60 \\ 14.19 \\ 30.23 \\ 17.25$		1.4 1.2				0.4 0.2					
Average	22.59	21.55	20.64	<u>19.87</u>	20.42	33.55	16.84	_		0	10 20 3 Epoch	40 50	0.0	0 10	20 Ep	30 och	40	50

B) Test loss (SUN)



Layer selection frequencies **C**)



- BVG-LS outperforms the best performance of typical finetuning on 6/7 dataset. **A**)
 - BVG-LS significantly outperforms full finetuning on all 7/7 datasets, although both methods do update all layers in the end.
- B) Despite fewer updates per layer, BVG-LS undergoes faster loss descent and achieves smaller test loss than full finetuning.
- **C**) The last layer undergoes dominant updates in early stage.
 - Surprisingly, the selection frequency of the input layer is quite high (roughly, the second highest).

This work is supported by DENSO IT LAB Recognition and Learning Algorithm Collaborative Research Chair (Science Tokyo.)