

Abstract

Model compression by way of parameter pruning, quantization, or distillation has recently gained popularity as an approach for reducing the computational requirements of modern deep neural network models for NLP. Inspired by previous work suggesting a connection between simpler, more generalizable models and those that lie within flat basins in the loss landscape, we hypothesize that optimizing for flat minima should lead to simpler parameterizations and thus more compressible models. We propose to combine sharpness-aware minimization with various task-specific model compression methods, including iterative magnitude pruning, structured pruning with a distillation objective, and post-training dynamic quantization. In our experiments, we show that optimizing for flat minima consistently leads to greater compressibility of parameters compared to standard Adam optimization when fine-tuning BERT models, leading to higher rates of compression with little to no loss in accuracy on the GLUE text classification and SQuAD question answering benchmarks.

Background

(BQ0): How do we train flat? Sharpness-Aware Minimization (SAM)

Cifar10 -			00
Cifar100 -			
Imagenet -	0.0000	•	
Finetuning - SVHN F-MNIST		••	٠
Noisy Cifar -		• •	
L	0	20	40

Error reduction (%)



Figure 1: (left) Error rate reduction obtained by switching to SAM. Each point is a different dataset / model / data augmentation. (middle) A sharp minimum to which a ResNet trained with SGD converged. (right) A wide minimum to which the same ResNet trained with SAM converged.



(BQ1): How do we compress? Many ways...

• Pruning

- We use unstructured iterative magnitude pruning (IMP); a proportion of individual parameters is pruned repeatedly based on absolute magnitude
- Lottery ticket hypothesis (Frankle and Carbin, 2019): Models are overparameterized... \exists sparse subnetworks that can be trained to reach the full model's accuracy in the same number of steps

• Distillation

- A smaller student model learns to mimic a teacher model's outputs
- We use a SoTA method that prunes coarse + fine-grained (CoFi) substructures from a student model using a distillation objective (Xia et al., 2022)

• Quantization

- Learn/calibrate transformations that map float32 params into smaller (e.g. int8) representations for inference
- We use a standard <u>PyTorch implementation</u> of Post-Training Dynamic Quantization

Models: pre-trained BERT models (BERT_{base}, + BERT_{large} & RoBERTa_{base}) *Tasks:* Text Classification (from GLUE), and Question Answering (SQuAD)

Optimizers: vanilla Adam (Kingma and Ba., 2014) vs. SAM (Foret et al., 2021) (with Adam base optimizer)

Metrics: accuracy metrics reported by related work at reference sparsity levels Hyperparameters

- $\rho = 0.05$ for size of epsilon-neighborhood in SAM
- Matched <u>Chen et al. (2020)</u> for unstructured IMP experiments
- Pruned 10% of weights between each of 10 fine-tuning iterations -> 90% sparsity
- We experiment both with and without *rewind* to BERT_{base} initialization
- Different: we set Adam w.d. = 0 in order to remove potential confound of L2 reg.
- Matched Xia et al. (2022) for structured/distillation (CoFi) experiments
- Different: we compare teacher models trained with SAM vs vanilla Adam

				Experime
Ŕ	Q2): How	does SAM	influence i	model compres
5A	AM learns	more com	pressible	winning tick
0.	Dataset	Ticket	Optim.	Accuracy
	RTE	Random	Adam	54.9 ± 1.3
	(60%)	Random	SAM	55.4 ± 1.6
		Adam	Adam	63.7 ± 0.9
		Adam	SAM	61.7 ± 2.4
		SAM	Adam	70.2 ± 1.9
		SAM	SAM	71.8 ± 1.7
-	MRPC	Random	Adam	70.8 ± 0.8
	(50%)	Random	SAM	70.1 ± 0.2
		Adam	Adam	84.2 ± 0.1
		Adam	SAM	85.3 ± 0.5
		SAM	Adam	85.3 ± 1.3
		SAM	SAM	85.7 ± 0.8

(RQ3): How does SAM affect compressed model generalization?

SAM-learned tickets transfer better

Figure 8: Heatmaps indicating the absolute difference in target task performance between SAM and Adam tickets when transferring tickets across tasks, fine-tuned with either SAM (left) or Adam (right) optimizers during IMP. Values > 0 indicate the extent to which SAM tickets transferred better than Adam tickets; Values < 0 indicate where Adam tickets transferred better than SAM. Overall, SAM tickets transfer better regardless of the final fine-tuning optimizer.

SAM vs	1 vs Adam Ticket Transfer Performance, w/ SAM Optimizer SAM vs Adam Ticket Transfer Performance, w/ Adam Optimizer																						
ilum	0.1	0.3	0.9	0.8	2.2	0.0	-0.6	1.3	0.6		- 5	uli m	0.2	-0.0	0.5	0.2	0.7	2.5	0.2	-1.6	0.3		- 2.5
ddb	-0.3	0.1	1.9	0.1	5.2	1.8	-0.1	2.4	1.4		- 4	dbb	-0.5	0.2	0.5	-1.3	3.4	1.1	0.8	-6.3	-0.3		- 2.0
stsb	0.0	0.4	0.1	-0.3	0.2	3.2	-0.1	2.3	0.7		- 3	stsb -	-0.3	0.0	0.1	-0.3	-1.2	0.7	0.2	-0.4	-0.1		- 1.5
Task qnli	0.2	0.2	0.6	0.3	-1.2	-1.4	0.3	6.1	0.6		- 7	Task qnli	-0.2	0.0	-0.1	0.3	0.2	0.4	0.5	0.6	0.2		- 1.0
ource mrpc	-0.4	-0.1	0.3	-0.7	1.0	1.8	0.0	-0.1	0.2		2	Source mrpc	-0.2	0.1	0.2	0.1	1.0	1.4	-0.3	-0.8	0.2		- 0.5
rte S	0.5	-0.1	0.7	0.1	3.9	10.5	-0.1	-2.7	1.6		- 1	t a	-0.3	-0.1	0.4	0.0	2.0	7.2	0.2	1.0	1.3		- 0.0
sst2	0.1	-0.2	0.5	0.1	-0.7	0.0	0.5	-3.6	-0.4		- 0	a sst2	0.2	-0.1	0.4	-0.4	0.5	1.8	0.2	0.6	0.4		0.5
j cola	0.1	-0.1	-0.1	-0.5	0.5	-1.1	0.1	-0.2	-0.1		1	- cola	-0.1	0.0	0.1	-0.0	-1.2	-1.1	0.5	-0.6	0.0		1.0
Avg	mnli		U.6	-0.0	1.4	rte	u.u	cola	U.6			Ave	mnli	qqp	stsb	qnli	mrpc	rte	sst2	cola	Avg		
		996		Та	irget T	ask										Та	rget Ta	ask					
SAM vs	a Adar	n Opt	imize	er Tra	nsfer	Perfo	rman	ce, w	/ SAM	Tick	ket	SAM vs	Adar	n Opt	imize	r Trai	nsfer	Perfo	rman	ce, w/	Adar	n Ti	cket
iln -	0.6	0.3	0.6	0.1	-8.1	-7.2	1.1	2.2	-1.3		- 2	il -	0.7	-0.1	0.1	-0.6	-9.6	-4.7	1.9	-0.7	-1.6		- 1
dbb	0.4	0.3	1.1	-0.8	-1.5	1.8	0.7	5.9	1.0		- 1	dbb	0.2	0.4	-0.3	-2.1	-3.2	1.1	1.6	-2.8	-0.6		- 0
stsb	0.6	0.1	0.5	-0.8	0.7	1.1	0.5	1.7	0.5		- 0	stsb	0.2	-0.2	0.4	-0.8	-0.7	-1.4	0.8	-1.0	-0.3		1
Task qnli	0.2	-0.0	0.1	0.6	-4.4	-3.2	1.3	-1.8	-0.9		1	Task qnli	-0.3	-0.2	-0.6	0.6	-2.9	-1.4	1.4	-7.3	-1.4		1
ource ⁻ mrpc	0.2	-0.0	0.3	-0.9	1.5	0.4	1.3	0.5	0.4		2	ource mrpc	0.4	0.2	0.2	-0.1	1.5	0.0	0.9	-0.3	0.3		2
rte -	0.5	0.2	0.6	0.1	-1.2	-2.5	0.3	-1.4	-0.4		3	- rt S	-0.3	0.1	0.2	-0.0	-3.2	-5.8	0.7	2.4	-0.7		3
sst2	0.1	0.1	0.2	-0.2	-3.7	-5.4	0.8	-1.9	-1.2		4	sst2	0.2	0.2	0.1	-0.6	-2.5	-3.6	0.6	2.2	-0.4		4
cola	0.1	-0.2	-0.1	-0.4	2.7	-6.1	0.9	-1.1	-0.5		5	cola	-0.1	-0.1	-0.0	-0.0	1.0	-6.1	1.3	1.0	-0.4		5
Avg	0.3	0.1	0.4	-0.3	-1.8	-2.7	0.9	0.5	-0.3		6	Avg	0.1	0.0	0.0	-0.5	-2.5	-2.8	1.1	-0.8	-0.6		6
	mnli	qqp	stsb	qnli Ta	mrpc rget Ta	rte ask	sst2	cola	Avg				mnli	qqp	stsb	qnli Ta	mrpc rget Ta	rte ask	sst2	cola	Avg		

Figure 9: Heatmaps indicating the absolute difference in target task performance between SAM and Adam optimizers during fine-tuning when transferring tickets across tasks. Values greater than 0 indicate the extent to which SAM optimizer worked better than Adam: Values less than 0 indicate where Adam optimizer worked better than SAM: Values close to 0 indicate

AM vs	vs Adam Ticket Transfer Performance, w/ SAM Optimizer SAM vs Adam Ticket Transfer Performance, w/ Adam Optimizer																						
ilum.	0.1	0.3	0.9	0.8	2.2	0.0	-0.6	1.3	0.6		- 5	ilum -	0.2	-0.0	0.5	0.2	0.7	2.5	0.2	-1.6	0.3		- 2.5
ddb	-0.3	0.1	1.9	0.1	5.2	1.8	-0.1	2.4	1.4		- 4	ddb	-0.5	0.2	0.5	-1.3	3.4	1.1	0.8	-6.3	-0.3		- 2.0
stsb	0.0	0.4	0.1	-0.3	0.2	3.2	-0.1	2.3	0.7		- 3	stsb	-0.3	0.0	0.1	-0.3	-1.2	0.7	0.2	-0.4	-0.1		- 1.5
Task qnli	0.2	0.2	0.6	0.3	-1.2	-1.4	0.3	6.1	0.6		- 7	Task qnli	-0.2	0.0	-0.1	0.3	0.2	0.4	0.5	0.6	0.2		- 1.0
ource mrpc	-0.4	-0.1	0.3	-0.7	1.0	1.8	0.0	-0.1	0.2		2	Source	-0.2	0.1	0.2	0.1	1.0	1.4	-0.3	-0.8	0.2		- 0.5
rt o	0.5	-0.1	0.7	0.1	3.9	10.5	-0.1	-2.7	1.6		- 1	te to	-0.3	-0.1	0.4	0.0	2.0	7.2	0.2	1.0	1.3		- 0.0
sst2	0.1	-0.2	0.5	0.1	-0.7	0.0	0.5	-3.6	-0.4		- 0	i sst2	0.2	-0.1	0.4	-0.4	0.5	1.8	0.2	0.6	0.4		0.5
- cola	0.1	-0.1	-0.1	-0.5	0.5	-1.1	0.1	-0.2	-0.1		1	g cola	-0.1	0.0	0.1	-0.0	-1.2	-1.1	0.5	1.9	0.0		1.0
Avg	0.0		U.6	-0.0	1.4	1.9	0.0	0.7				Avg	mnli	aap	stsb	-0.2	mrpc	rte	sst2	cola	v.s Ava		
		ЧЧР	3(3)	Та	irget T	ask	3312	cold	Avg					-1-11-		Та	rget Ta	ask			J		
AM vs	a Adar	n Opt	timize	er Tra	nsfer	Perfo	rman	ce, w	/ SAM	Tic	ket	SAM vs	Adar	n Opt	imize	r Trai	nsfer	Perfo	rman	ce, w/	/ Adar	n Ti	cket
iln '	0.6	0.3	0.6	0.1	-8.1	-7.2	1.1	2.2	-1.3		- 2	mnli.	0.7	-0.1	0.1	-0.6	-9.6	-4.7	1.9	-0.7	-1.6		- 1
ddb	0.4	0.3	1.1	-0.8	-1.5	1.8	0.7	5.9	1.0		- 1	ddb	0.2	0.4	-0.3	-2.1	-3.2	1.1	1.6	-2.8	-0.6		- 1
stsb	0.6	0.1	0.5	-0.8	0.7	1.1	0.5	1.7	0.5		- 0	stsb	0.2	-0.2	0.4	-0.8	-0.7	-1.4	0.8	-1.0	-0.3		- 0
rask qnli	0.2	-0.0	0.1	0.6	-4.4	-3.2	1.3	-1.8	-0.9		1	Task qnli	-0.3	-0.2	-0.6	0.6	-2.9	-1.4	1.4	-7.3	-1.4		1
nrpc	0.2	-0.0	0.3	-0.9	1.5	0.4	1.3	0.5	0.4		2	burce mrpc	0.4	0.2	0.2	-0.1	1.5	0.0	0.9	-0.3	0.3		2
So Tte	0.5	0.2	0.6	0.1	-1.2	-2.5	0.3	-1.4	-0.4		3	rte i	-0.3	0.1	0.2	-0.0	-3.2	-5.8	0.7	2.4	-0.7		3
sst2	0.1	0.1	0.2	-0.2	-3.7	-5.4	0.8	-1.9	-1.2		4	sst2	0.2	0.2	0.1	-0.6	-2.5	-3.6	0.6	2.2	-0.4		4
cola	0.1	-0.2	-0.1	-0.4	2.7	-6.1	0.9	-1.1	-0.5		5	cola	-0.1	-0.1	-0.0	-0.0	1.0	-6.1	1.3	1.0	-0.4		5
Avg	0.3	0.1	0.4	-0.3	-1.8	-2.7	0.9	0.5	-0.3		6	Avg	0.1	0.0	0.0	-0.5	-2.5	-2.8	1.1	-0.8	-0.6		6
	mnli	qqp	stsb	qnli Ta	mrpc rget Ta	rte ask	sst2	cola	Avg				mnli	qqp	stsb	qnli Ta	mrpc irget Ta	rte ask	sst2	cola	Avg		

Key findings

- compared to vanilla Adam optimization
- even a vanilla Adam optimizer can learn on successfully:
- over a given ticket

(DQ0): What's the catch?

- Standard SAM is slower... bilevel optimization can be costly
- Lottery ticket-style IMP is expensive (fine-tuning x *n* iterations)
- compressed networks

(DQ1): What's next?

(e) SQuAD, IMP w/ RoBERTa

- flatness (e.g. stochastic weight averaging)?
- modalities?
- What role could pre-training with SAM play?



nts (cont'd)

ssibility?

tet" structures

Table 2: For RTE and MRPC at their reference sparsity values, we fine-tune using 1) SAM and 2) Adam optimizers from pre-trained BERT-base initializations using only the remaining weights based on a) a random mask, b) an Adam-learned mask, and c) a SAM-learned mask.

across	tasks
across	tasks

Discussion

• Fine-tuning with SAM improves compressibility of BERT models across a variety of settings,

• Specifically, SAM helps us *find winning lottery tickets* during iterative magnitude pruning that

• The *structure* of a winning ticket learned by SAM is more important than the *optimization*

• The winning tickets learned by SAM *generalize better* to other tasks • On the other hand, SAM does not always help us optimize a given winning ticket

• Hardware does not necessarily support inference or training speedups proportional to size in

• SAM works, but what about more efficient alternatives, or different methods of obtaining

• Do these results hold in much larger models? Other architectures? Other types of tasks or