Merging and MoErging for Compositional Generalization Colin Raffel



Tasks can be considered as a composition of skills



Multitask models can generalize to new tasks



From "Multitask Prompted Training Enables Zero-Shot Task Generalization" by Sanh et al.









Model merging



"Vanilla" merging is just parameter averaging



From "Weight Averaging for Neural Networks and Local Resampling Schemes" by Utans

Model merging as an optimization problem



Parameter averaging makes an isotropic Gaussian assumption



Fisher merging uses the Laplace approximation



Fisher merging can combine the capabilities of different models



Merging models based on loss landscape curvature



Solving a general merging problem with the conjugate gradient method

$$oldsymbol{ heta}^{oldsymbol{st}} = igg(\sum_{m=1}^M oldsymbol{C}_migg)^{-1}igg(\sum_{m=1}^M oldsymbol{C}_moldsymbol{ heta}_migg)$$



MaTS allows flexibly combining merging objectives and initializations...



... and achieved state-of-the-art resulted across settings



Merging models with task vectors



TIES Merging resolves interference between task vectors



From "Resolving Interference When Merging Models" by Yadav et al.

TIES helps retain specialist model performance



From "Resolving Interference When Merging Models" by Yadav et al.

Is merging actually doing what we want?



NLP Task \downarrow / Language $ ightarrow$	English	Arabic	Thai	German	Korean
Question-Answering (SQuaD/XQuaD)	M	1	Ø	Ø	
Natural Language Inference (XNLI)	and the second se	5	and the second s	and the second s	
Summarization (WikiLingua)	100 Marine	100 miles	3	100 Cart	(and
Word Sense Disambiguation (WiC/XLWiC)	and the second s			5	and the second s
Is question answerable? (TyDiQA)	1	(Jacobian Contraction of the con	(Jaci		3

From "Realistic Evaluation of Model Merging for Compositional Generalization" by Tam et al.

... not really.



From "Realistic Evaluation of Model Merging for Compositional Generalization" by Tam et al.

Merging methods also have different practical requirements...

	Prerequisites				Computational cost (FLOPs)		
	$\theta_{\mathbf{p}}$	Stats	Data	Hparams?	Merging	Statistics	
Average					Mdk		
SLERP					$\mathcal{O}((5M-2)dk)$		
Task Arith.	X		X	\checkmark	(2M+1)dk		
DARE	X		×	\checkmark	(6M+1)dk		
TIES	X		X *	\checkmark^*	(4M+1)dk	$\mathcal{O}(MKdk)$	
Fisher		X^*	\boldsymbol{X}^*		(3M-1)dk	$4MTd^2k$	
RegMean		X *	X *	\checkmark	$\mathcal{O}((M+2)d^2k)$	MTd^2k	
MaTS		X *	X *	\checkmark	$\mathcal{O}((M+N)d^2k)$	$4MTd^2k$	

From "Realistic Evaluation of Model Merging for Compositional Generalization" by Tam et al.

... and different hyperparameter sensitivity



From "Realistic Evaluation of Model Merging for Compositional Generalization" by Tam et al.

Multitask performance is still poor as you the number of models...



From "Realistic Evaluation of Model Merging for Compositional Generalization" by Tam et al.

... but the picture for generalization is better.



From "Realistic Evaluation of Model Merging for Compositional Generalization" by Tam et al.

An alternative: MoErging?



Differentiable routing between specialist submodels with SMEAR



From "Soft Merging of Experts with Adaptive Routing" by Muqeeth et al.

SMEAR is pareto-optimal across different routing strategies



From "Soft Merging of Experts with Adaptive Routing" by Muqeeth et al.

Experts specialize and are shared across different datasets



From "Soft Merging of Experts with Adaptive Routing" by Muqeeth et al.

Post-Hoc Adaptive Tokenwise Gating Over an Ocean of Specialized Experts



PHATGOOSE outperforms prior routing methods and multitask training



PHATGOOSE learns nontrivial routing strategies



The MoErging design space



From "A Survey on Model MoErging: Recycling and Routing Among Specialized Experts for Collaborative Learning" by Yadav et al.

Thanks. Please give me feedback: <u>http://bit.ly/colin-talk-feedback</u> <u>craffel@qmail.com</u>