

Singh & Jaggi, NeurIPS 2020

DLCT Reading Group: 25th September, 2020



Sidak Pal Singh





Problem Setup

- Given K independent neural nets (models), with each trained on a possibly different training set, i.e., $(M_1, S_1), \dots, (M_K, S_K)$.
- Aim: Fuse their capabilities in one single model *M*.
- Question: How to obtain the fusion operator?

$$\mathscr{M} \leftarrow f(M_1, \cdots, M_K) \ heta_\mathscr{M} \leftarrow f(heta_1, \cdots, heta_K)$$

• (Personal opinion: fundamental yet relatively unexplored)

How hard is it?

- As a start, let's assume, same architecture and size!
- Highly non-convex parameter space.
- Neural networks possess inherent symmetries.
- No apriori guarantee that a suitable solution exists.
- Additional constraints to care about (data free; no stacking of params)

Why is it important?

- Make predictions more robust: Ensemble
- Combine models learned on different datasets/tasks:
 - Federated learning
 - Continual learning

•

• "Skill/Knowledge" transfer to a smaller model

Things to care about

Resource costs (fusion procedure, storage)

□ Inference costs (communication, latency)

Privacy COStS (sharing training data)

Related Work



Umm, what if we average parameters one-to-one?

Input Models





- No direct correspondence between the neurons across two models.
- The models could have been initialized differently or trained on separate data.
- Does not really work in general.

Layer-wise alignment of neurons & then averaging!

(to simplify discussion, let's consider only two models A and B)

Input Models

Aligned Models





Output Model







Opening the black-box

Optimal Transport

Algorithm:

Layerwise alignment & average!

For each layer

$$\begin{split} & \mathbf{I}. \quad \widehat{\boldsymbol{W}}_{A}^{(\ell,\ell-1)} \leftarrow \boldsymbol{W}_{A}^{(\ell,\ell-1)} \boldsymbol{T}^{(\ell-1)} \operatorname{diag} \left(\frac{1}{\beta^{(\ell-1)}} \right) \qquad // \text{ Align ind} \\ & \mathbf{2}. \quad \boldsymbol{\alpha}^{(\ell)}, \quad \boldsymbol{\beta}^{(\ell)} \leftarrow \mathbf{1}_{n_{A}^{(\ell)}} / n_{A}^{(\ell)}, \quad \mathbf{1}_{n_{B}^{(\ell)}} / n_{B}^{(\ell)} \qquad // \text{ Initialise} \\ & \mathbf{3}. \quad \mu_{A}^{(\ell)}, \quad \mu_{B}^{(\ell)} \leftarrow \left(\boldsymbol{\alpha}^{(\ell)}, \boldsymbol{S}_{A}[\ell] \right), \quad \left(\boldsymbol{\beta}^{(\ell)}, \boldsymbol{S}_{B}[\ell] \right) \qquad // \text{ Write as} \\ & \mathbf{4}. \quad D_{S}^{(\ell)} \leftarrow \left\| \boldsymbol{S}_{A}[\ell]^{\top} - \boldsymbol{S}_{B}[\ell]^{\top} \right\|_{2} \qquad // \text{ Form gradients} \\ & \mathbf{5}. \quad \boldsymbol{T}^{(\ell)}, \mathcal{W}_{2}^{(\ell)} \leftarrow \operatorname{OT} \left(\boldsymbol{\mu}_{A}^{(\ell)}, \boldsymbol{\mu}_{B}^{(\ell)}, D_{S}^{(\ell)} \right) \qquad // \text{ Compute the set of the set$$

coming edges for M_A weights probability mass values for neurons s empirical measures ound metric te OT map and distance ation of neurons in A

// Average model weights

End for

Similar extension for multiple models, pick one model to align against!

Discussion

- How to design the ground metric (i.e., compare neurons)?
 - Activation-based alignment
 - Weight-based alignment
- We still have a greedy procedure:
 - Ideal procedure is combinatorially hard
- Runtime Efficiency (even with exact OT):
 - 15 seconds to fuse 6VGG11 models 1 Nvidia V100 GPU

Fusion for heterogeneous data and tasks

Applications: Skill Transfer, Federated Learning

Different Initialization

Model A is a "specialist", Model B is a "generalist". How to get the best?

Over here, it is just one-shot averaging!

Applications: Skill Transfer, Federated Learning

Fusion for different sized models

(OT can also provide soft-maps)

Applications: Adapting sizes across client-server

Model size at client != Model size at server

Applications: Post-processing tool for pruning

Fraction of channels removed

Data free!

VGG11, CIFAR10 (similar for CIFAR100)

Fusion for efficient ensembling

Applications: Efficient ensemble

DATASET +	M_A	M_B	PREDICTION	VANILLA	OT	FINETUNING	
	00.21	00.50	AVG.	AVG.	AVG.		00.72
VGG11	90.31 90.50		91.34	17.02	85.98	90.39	90.73
	1 ×		1 ×	2 ×	2 ×	2 ×	2 ×
CIFAR10+	93.11	93.20	93.89	18.49	77.00	93.49	93.78
ResNet18	1	×	1 ×	2 ×	2 ×	2 ×	2 ×

Two model case: CIFARIO

Applications: Efficient ensemble

CIFAR10+	INDIVIDUAL MODELS	PREDICTION	VANILLA	OT	Finetun	ΓUNING	
VGG11	INDIVIDUAL WIODELS	AVG.	AVG.	AVG.	VANILLA	OT	
Accuracy	[90.31, 90.50, 90.43, 90.51]	91.77	10.00	73.31	12.40	90.91	
Efficiency	$1 \times$	1 ×	4 ×	4 ×	4 ×	4 ×	
Accuracy	[90.31, 90.50, 90.43, 90.51, 90.49, 90.40]	91.85	10.00	72.16	11.01	91.06	
Efficiency	1 ×	1 ×	6 ×	6 ×	6 ×	6 ×	

Multiple models case: CIFARIO

Applications: Efficient ensemble

CIFAR100 +	INDIVIDUAL MODELS	PREDICTION	<u>Fine</u> tuning	
VGG11	INDIVIDUAL MODELS	AVG.	VANILLA	ОТ
Accuracy	[62.70, 62.57, 62.50, 62.92]	66.32	4.02	64.29± 0.26
Efficiency	1 ×	1 ×	4 ×	4 ×
Accuracy	[62.70, 62.57, 62.50, 62.92, 62.53, 62.70]	66.99	0.85	$\textbf{64.55} \pm \textbf{0.30}$
Efficiency	1 ×	1 ×	6 ×	6 ×
Accuracy	[62.70, 62.57, 62.50, 62.92, 62.53, 62.70, 61.60, 63.20]	67.28	1.00	65.05±0.53
Efficiency	1 ×	1 ×	8 ×	8 ×

Table 3: Efficient alternative to ensembling via OT fusion on CIFAR100 for VGG11. Vanilla average fails to retrain. Results shown are mean \pm std. deviation over 5 seeds.

Multiple models case: CIFAR 100

Conclusion

- Fuse multiple models (of possibly different widths):
 - Layerwise soft alignment of neurons via OT & then average (+ mild fine-tuning)
- Significantly outperforms vanilla averaging (accuracy) and ensembling (efficiency)
- Can be used alongside Distillation (as an initialization)
- Overall, this results in
 - (a) Successful one-shot transfer of knowledge without sharing data
 - (b) Data free and algorithm independent post-processing
 - (c) Combine parameters of different sized models
 - (d) Keep one model rather than an ensemble of models

Future work

- Open problem: How to fuse models with different depths?
- Adjust the size of the fused model
- Fusing generative models: GANs
- Continual learning
- Toolkit to compare & interpolate between models

Thank you!

sidakpal.com

https://github.com/modelfusion/otfusion

