Adaptive Transfer Learning from Pre-trained Models

Mingsheng Long

School of Software, Tsinghua University

August 24, 2022





Supervised Training

Training images, scarce



Heavy parameters for large capacity



Related data, unlabeled

Adaptation

Training images, scarce



Fast convergence, better performance



Related images, labeled



Pre-training and Adaption

Pre-Training

Adaptation



 $Pre-training \rightarrow Adaption$

A Paradigm for Deep Learning Application

Pre-trained Models



GPT: Large-scale Corpus Pre-training



BiT: General Visual Representation Learning



Adaptation: Foundation Problems





Overcome Catastrophic Forgetting

Regularization Tuning

Loss Function: $\min_{\theta} \sum_{i=1}^{m} L(h_{\theta}(x_i), y_i) + \lambda \cdot \Omega(\theta)$

Regularization term





Overcome Catastrophic Forgetting





Overcome Negative Transfer





Overcome Negative Transfer

• Enhance Safe Transfer
• BSS, Zoo-tuning

$$\operatorname{err}_{p}(g) \leq \operatorname{err}_{p}^{\gamma}(f) + O\left(\sqrt{\frac{p_{g} \log_{2} r_{g}}{n}}\right)$$
 $\operatorname{Penalize smallest singular values}$
 $\operatorname{Lhss}(F) = n^{\sum_{k} \sigma^{2}}$

r challes shallest singular values r

$$\sigma_{\mathrm{oss}}(F) = \eta \sum_{i=1}^k \sigma_{-i}^2$$

• LEEP, LogME

Xinyang Chen, Sinan Wang, Bo Fu, , Jianmin Wang, Mingsheng Long . Catastrophic Forgetting Meets Negative Transfer: Batch Spectral Shrinkage for Safe Transfer Learning, NeurIPS 2019

Overcome Negative Transfer

- Enhance Safe Transfer
 - BSS, Zoo-tuning



- Choose Pre-trained Models
 - LEEP, LogME

Adaptive Transfer of Model Hubs



Pre-trained Model Hub

Various Models and Platforms



PyTorch Hugging Face TensorFlow Hub Pytorch Image Models

Avoid Heavy Pre-training



Plenty of Transferable Knowledge

IMAGENET SUP. MOCO PT. MASKRCNN PT. DEEPLAB PT. KEYPOINT PT.



Same architectue Pre-trained differently

- Adapt one model
- Which one is the best?
- Adapt multiple models
- How to aggregate transferable knowledge?

Transferability Assessment

Estimate adaption performance of PTM on given dataset without finetuning.

LogME Approach





- Fixed PTM (as feature extractor).
- P(y | F): Graphical modeling

between extracted features and GT label.

- Parameterize P(y | F) by prior α, β .
- Maximize evidence $P(y | F, \alpha, \beta)$.
 - MacKay algorithm with guarantee!

Effectiveness of LogME

General and Accurate



Effectiveness of LogME

Computation Efficient --- MacKay algorithm with theoretical guarantee!

Algorithm 2 Evidence Maximization by MacKay's Algorithm
1: Input: Extracted features $F \in \mathbb{R}^{n \times D}$ and corresponding labels $y \in \mathbb{R}^n$
2: Output: Logarithm of Maximum Evidence (LogME)
3: Note: F has been pre-decomposed into $F = U\Sigma V^T$
4: Initialize $\alpha = 1, \beta = 1$
5: while α, β not converge do
6: Compute $\gamma = \sum_{i=1}^{D} \frac{\beta \sigma_i^2}{\alpha + \beta \sigma_i^2}, \Lambda = \text{diag}\{(\alpha + \beta \sigma^2)\}$
7: Naïve : $A = \alpha I + \beta F^T F, m = \beta A^{-1} F^T y$
9: Update $\alpha \leftarrow \frac{\gamma}{m^T m}, \beta \leftarrow \frac{n-\gamma}{\ Fm-y\ _2^2}$

```
10: end while
11: Compute and return \mathcal{L} = \frac{1}{n}\mathcal{L}(\alpha,\beta) using Equation 2
```

 $\mathcal{O}(nCD^2 + CD^3)$ Biquadrate complexity



$$\mathcal{O}(nD^2 + nCD + CD^2 + D^3)$$

Cubic complexity

Algorithm 3 Evidence Maximization by Optimized Fixed Point Iteration1: Input: Extracted features $F \in \mathbb{R}^{n \times D}$ and corresponding labels $y \in \mathbb{R}^n$ 2: Output: Logarithm of Maximum Evidence (LogME)3: Require: Truncated SVD of $F: F = U_r \Sigma_r V_r^T$, with $U_r \in \mathbb{R}^{n \times r}, \Sigma_r \in \mathbb{R}^{r \times r}, V_r \in \mathbb{R}^{D \times r}.$ 4: Compute the first r entries of $z = U_r^T y$ 5: Compute the sum of remaining entries $\Delta = \sum_{i=r+1}^n z_i^2 = \sum_{i=1}^n y_i^2 - \sum_{i=1}^r z_i^2$ 6: Initialize $\alpha = 1, \beta = 1, t = \frac{\alpha}{\beta} = 1$ 7: while t not converge do8: Compute $m^T m = \sum_{i=1}^r \frac{\sigma_i^2 z_i^2}{(t + \sigma_i^2)^2}, \gamma = \sum_{i=1}^r \frac{\sigma_i^2}{t + \sigma_i^2}, ||Fm - y||_2^2 = \sum_{i=1}^r \frac{z_i^2}{(1 + \sigma_i^2/t)^2} + \Delta$ 9: Update $\alpha \leftarrow \frac{\gamma}{m^T m}, \beta \leftarrow \frac{n - \gamma}{||Fm - y||_2^2}, t = \frac{\alpha}{\beta}$ 10: end while11: Compute $m = V_r \Sigma' z$, where $\Sigma'_{ii} = \frac{\sigma_i}{t + \sigma_i^2} (1 \le i \le r)$.12: Compute and return $\mathcal{L} = \frac{1}{n} \mathcal{L}(\alpha, \beta)$ using Equation 2

$\mathcal{O}(nD^2 + nCD)$

Cubic complexity with fewer terms

memory footprint
r bound) $161000s$ fine-tune (upper bound) 6.3 GB
(lower bound) 37s extract feature (lower bound) 43 MB
43s LogME 53 MB
$3700 \uparrow $ benefit $120 \uparrow$
r bound) 100200s fine-tune (upper bound) 88 GB
(lower bound) 1130s extract feature (lower bound) 1.2 GB
1136s LogME 1.2 GB
$88\uparrow$ benefit $73\uparrow$
<u>פ</u>





LEEP, NCE, LogME...

Top-K: Heuristic but Effective

- Architecture heterogeneity
- Dimensionality of features
- Always challenging part...

B-Tuning

Consider simple Knowledge Distillation (KD):

$$L_{\text{KD}} = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{K} \sum_{k=1}^{K} |\phi_k(x_i) - W_k \phi_t(x_i)|_2^2$$

- Needs additional learnable projection W_k for each teacher model.
- Treats all teacher models as equal:
 - No adaptive mechanism to transfer only useful knowledge.

Kaichao You, Yong Liu, Jianmin Wang, Michael I Jordan, Mingsheng Long . Ranking and Tuning Pre-trained Models: A New Paradigm of Exploiting Model Hubs, JMLR 2022

B-Tuning



- Project teacher features into a common output space by LogME.
- Transfer them to target model with weighting from their LogME score.

Intuition: encourage the target model to behave like the best top-K teachers.

Kaichao You, Yong Liu, Jianmin Wang, Michael I Jordan, Mingsheng Long . Ranking and Tuning Pre-trained Models: A New Paradigm of Exploiting Model Hubs, JMLR 2022

Experiments

Reduced burden of Selection and Adaptation.

- Exhaustively fine-tune 10 times: 84.41% accuracy.
- Rank by LogME and fine-tune once: 84.29% accuracy.



Experiments

Fully utilization of transferable knowledge in Model Hub.



- Just fine-tune the most popular model is sub-optimal.
- The ranking and B-Tuning paradigm brings 3%~5% accuracy gain.

Adaption from Model Zoo

Considering models with same architecture but different knowledge.



Zoo-Tuning

Adaptively aggregate source model parameters to derive target model.



Adaptive Aggregation



Adaptive Aggregation



Channel alignment

• Channels in different pre-trained models may have different semantic meanings.

$$\widetilde{\mathbf{W}}_i^l = \mathbf{T}_i^l \ast \mathbf{W}_i^l$$

Transformed parameter

Adaptive Aggregation



Experiments

- Adaptive transfer from multiple models \rightarrow Better accuracy.
- Adaptive aggregation of model parameters \rightarrow More efficient than ensemble.

	GENERAL FINE-GRAINED SPECIALIZED			TRAIN		INFERENCE						
Model	CIFAR-100	COCO-70	AIRCRAFT	CARS	INDOORS	DMLAB	EuroSAT	Avg. Acc	GFLOPS	Params	GFLOPS	Params
IMAGENET SUP.	81.18	81.97	84.63	89.38	73.69	74.57	98.43	83.41	4.12	23.71M	4.12	23.71M
MOCO PT.	75.31	75.66	83.44	85.38	70.98	75.06	98.82	80.66	4.12	23.71M	4.12	23.71M
MASKRCNN PT.	79.12	81.64	84.76	87.12	73.01	74.73	98.65	82.72	4.12	23.71M	4.12	23.71M
DEEPLAB PT.	78.76	80.70	84.97	88.03	73.09	74.34	98.54	82.63	4.12	23.71M	4.12	23.71M
Keypoint pt.	76.38	76.53	84.43	86.52	71.35	74.58	98.34	81.16	4.12	23.71M	4.12	23.71M
ENSEMBLE	82.26	82.81	87.02	91.06	73.46	76.01	98.88	84.50	20.60	118.55M	20.60	118.55M
DISTILL	82.32	82.44	85.00	89.47	73.97	74.57	98.95	83.82	24.72	142.28M	4.12	23.71M
KNOWLEDGE FLOW	81.56	81.91	85.27	89.22	73.37	75.55	97.99	83.55	28.83	169.11M	4.12	23.71M
LITE ZOO-TUNING	83.39	83.50	85.51	89.73	75.12	75.22	99.12	84.51	4.53	130.43M	4.12	23.71M
ZOO-TUNING	83.77	84.91	86.54	90.76	75.39	75.64	99.12	85.16	4.53	130.43M	4.18	122.54M

Applied to RL tasks

- Reinforcement Learning: Atari Games.
- Pre-trained Models: Models trained from other games.



Towards A More General Framework

- Parameter aggregation: efficient and adaptive.
- Not general enough for diverse models or even heterogeneous architectures.



Hub-Pathway Framework

Design data-dependent pathways throughout the Model Hub.



Hub-Pathway Framework

- Input level: route different data to different models.
- Output level: aggregate transferred knowledge to make predictions.
- Pathway flow: control training and inference costs with Top-K activation.



Exploration and Exploitation



- How to explore better pathway configurations?
- How to fully exploit and transfer knowledge in pre-trained models?

Enhance Exploration

• Direct training may activate only one or a few models \rightarrow hub collapse.



• A noisy pathway generator to add randomness.

Enhance Exploration

• Direct training may activate only one or a few models \rightarrow hub collapse.



• Inspired from RL: maximum-entropy regularization to encourage exploration.

Enhance Exploitation

• Additional tuning with specific data to enhance knowledge transfer.



Experiments

• Data dependent pathways \rightarrow General for heterogenous models.

Model	General		Fi	ne-Grain	ed	Spec	Δυσ	
WIOUEI	CIFAR	COCO	Aircraft	Cars	Indoors	DMLab	EuroSAT	Avg.
MaskRCNN	79.12	81.64	84.76	87.12	73.01	74.73	98.65	82.72
MobileNetV3	83.14	83.28	80.26	86.37	75.09	70.09	98.95	82.45
EffNet-B3	87.28	86.97	83.99	89.34	78.16	72.69	99.13	85.37
Swin-T	84.37	84.12	80.82	89.10	73.39	72.22	98.69	83.24
ConvNeXt-T	86.96	87.15	84.23	90.67	81.66	73.80	98.65	86.16
Ensemble	87.72	88.04	87.11	92.68	82.79	74.86	99.23	87.49
Distill	87.33	88.09	85.26	91.39	81.51	74.75	99.24	86.80
Hub-Pathway	89.01	89.14	88.12	92.93	84.40	74.80	99.26	88.24

• Control the costs with top-k activation \rightarrow More efficient than ensemble.

Model	Acc (%)	Params (M)	FLOPs (G)	Train (iters/s)	Inference (samples/s)
ImageNet	83.41	23.71	4.11	10.87	484.92
Ensemble	84.50	118.55	20.55	2.30	98.64
Hub-Pathway	85.63	128.43	9.11	4.68	240.48

Adaptive Transfer with Generated Pathways



Thank You!



树扬





刘雍



龙明盛



王建民

大数据系统软件国家工程研究中心



游凯超

