Transferability in Deep Learning

Mingsheng Long

School of Software, Tsinghua University

MLA 2022





Outline

1	Introduction			
	1.1	Terminology	5	
	1.2	Overview	7	
2	Pre-	Training	9	
	2.1	Pre-Training Model	9	
	2.2	Supervised Pre-Training	12	
	2.3	Unsupervised Pre-Training	20	
3	Adaptation			
	3.1	Task Adaptation	31	
	3.2	Domain Adaptation	44	
4	Eva	luation	68	
	4.1	Datasets	68	
	4.2	Library	69	
	4.3	Benchmark	72	
5	Con	clusion	76	
6	Ack	nowledgements	77	
Re	eferer	ICes	78	

Transferability in Deep Learning

Suggested Citation: Junguang Jiang, Yang Shu and Mingsheng Long (2022), "Transfer- ability in Deep Learning", : Vol. xx, No. xx, pp 1–18. DOI: 10.1561/XXXXXXXX.
Junguang Jiang Tsinghua University jiangjunguang1123@outlook.com
Yang Shu Tsinghua University shu-y18@mails.tsinghua.edu.cn
Mingsheng Long Tsinghua University mingsheng@tsinghua.edu.cn

This article may be used only for the purpose of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval.





GPT: Large-scale Corpus Pre-training



BiT: General Visual Representation Learning





Pre-training and Adaptation

Pre-Training

Adaptation



 $Pre-training \rightarrow Adaptation$

A Paradigm for Deep Learning Application

Pre-training and Adaptation



All Hypotheses

All Hypotheses

Pre-training → Adaptation

A Paradigm for Deep Learning Application

Transferability in the Lifecycle



Pre-training



Supervised Pre-training



- Big Transfer (BiT) (Kolesnikov et al., 2020) emphasizes that training on larger datasets is vital for better transferability.
- Domain Adaptive Transfer (DAT) (Ngiam et al., 2018) uses importance weighting to carefully choose the pre-training data that are most relevant to the target task.





Task $i \in [1, ..., n]$ $\phi^* = \arg \max_{\phi} \sum_{i=1}^n \log P(\theta_i(\phi) | \mathcal{D}_i^{\text{ts}}),$ where $\theta_i(\phi) = \arg \max_{\theta} \log P(\theta | \mathcal{D}_i^{\text{tr}}, \phi).$

Model-Agnostic Meta-Learning (MAML)

for fast adaptation

$$\theta_i = \phi - \alpha \nabla_\phi L(\phi, \mathcal{D}_i^{\mathrm{tr}})$$

Causal Learning



Invariant Risk Minimization (IRM) for OOD generalization $\min_{\psi: \mathcal{X} \to \mathcal{Z}, h: \mathcal{Z} \to \mathcal{Y}} \sum_{e \in \mathcal{E}^{tr}} \epsilon^{e}(h \circ \psi),$

subject to $h \in \underset{\bar{h}: \mathcal{Z} \to \mathcal{Y}}{\operatorname{arg\,min}} \epsilon^{e}(\bar{h} \circ \psi)$, for all $e \in \mathcal{E}^{\operatorname{tr}}$

Essentially, this implies invariance to **data augmentation**!

• Causal learning seeks a model with causal mechanisms, and if the environment or distribution changes, only part of the causal mechanisms will be affected.

Generative Pre-training





12

Contrastive Pre-training



$$\min_{\psi} - \log \frac{\exp(\mathbf{q} \cdot \mathbf{k}_{+}/\tau)}{\sum_{j=0}^{K} \exp(\mathbf{q} \cdot \mathbf{k}_{j}/\tau)}$$

Supervised pre-training gains high-level semantic knowledge, while contrastive and generative pretraining gains mid-level & low-level representations



Remarks on Pre-training

Method	$\begin{array}{c} {\rm Modality} \\ {\rm Scalability}^1 \end{array}$	$\begin{array}{c} {\rm Task} \\ {\rm Scalability}^2 \end{array}$	Data Efficiency ³	$Label Efficiency^4$
Standard Pre-Training	***	**	***	*
Meta-Learning	***	*	*	*
Causal Learning	**	*	*	*
Generative Learning	**	***	***	***
Contrastive Learning	*	***	***	***

¹ Whether models can be pre-trained on various modalities, such as text, graph.

 $\frac{2}{2}$ Whether pre-trained models can be easily transferred to many downstream tasks.

³ Whether stronger transferability can be yielded from large-scale pre-training.

⁴ Whether pre-training relies on manual data labeling.

Foundation Model



[Data Universal]

Learn from various modalities

[Task Universal]

Adapt to a wide range of

downstream tasks

Bommasani et al. On the Opportunities and Risks of Foundation Models. Arxiv 2021.

General Relation Modeling in Transformers



Quadratic Complexity in Self-Attention



Quadratic Complexity in Self-Attention



 $(QK^T)V = Q(K^TV) \implies \mathcal{O}(n^2d) \rightarrow \mathcal{O}(nd^2)$

Recap: Softmax function

Softmax function is proposed as a differentiable generalization of the "winner-take-all" picking maximum operation.



Bridle et al. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. NeurIPS 1989.

Recap: Softmax function

Softmax function is proposed as a differentiable generalization of the "winner-take-all" picking maximum operation.



Bridle et al. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. NeurIPS 1989.

Flow Network Theory



Attention: A Flow Network View



Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, Mingsheng Long. Flowformer: Linearizing Transformers with Conservation Flows. ICML, 2022.

Attention: A Flow Network View



Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, Mingsheng Long. Flowformer: Linearizing Transformers with Conservation Flows. ICML, 2022.



- [Incoming Flow Conservation]:
 - Competition among Source tokens
- [Outgoing Flow Conservation]:
 - Competition among Sink tokens





Incoming flow: $I_i = \phi(Q_i) \sum_j \phi(K_j)^T$

Incoming flow conservation: $\frac{\phi(Q)}{I}$

Incoming flow:
$$\frac{\phi(Q_i)}{I_i} \sum_j \phi(K_j)^T = \frac{I_i}{I_i} = 1$$





Incoming flow: $I_i = \phi(Q_i) \sum_j \phi(K_j)^T$

Incoming flow conservation: $\frac{\phi(Q)}{I}$

Conserved outgoing flow: $\widehat{\mathbf{O}} = \phi(\mathbf{K}) \sum_{i} \frac{\phi(Q_i)^T}{I_i}$





Outgoing flow: $O_i = \phi(K_i) \sum_j \phi(Q_j)^T$

Outgoing flow conservation: $\frac{\phi(K)}{O}$

Outgoing flow:
$$\frac{\phi(K_i)}{O_i} \sum_j \phi(Q_j)^T = \frac{O_i}{O_i} = 1$$





Outgoing flow: $O_i = \phi(K_i) \sum_j \phi(Q_j)^T$

Outgoing flow conservation: $\frac{\phi(K)}{O}$

Conserved incoming flow: $\hat{I} = \phi(Q) \sum_{j} \frac{\phi(\kappa_{j})^{T}}{o_{j}}$



Competition: $\widehat{\mathbf{V}} = \operatorname{Softmax}(\widehat{\mathbf{O}}) \odot \mathbf{V}$ Aggregation: $\mathbf{A} = \frac{\phi(\mathbf{Q})}{\mathbf{I}} (\phi(\mathbf{K})^{\mathsf{T}} \widehat{\mathbf{V}})$ Allocation: $\mathbf{R} = \operatorname{Sigmoid}(\widehat{\mathbf{I}}) \odot \mathbf{A}$,

Successfully bring the **Competition Mechanism** Into Attention design to avoid trivial attention

Flowformer: Efficiency and Universality



[Efficiency]: All the calculations are in linear complexity.

[Universality]: The whole design is based on flow network without specific inductive biases.

Flowformer Experiments



BENCHMARKS	ΤΑSΚ	VERSION	Length
LRA (2020C)	SEQUENCE	NORMAL	1000~4000
WIKITEXT (2017)	LANGUAGE	CAUSAL	512
IMAGENET (2009)	VISION	NORMAL	49~3136
UEA (2018)	TIME SERIES	NORMAL	29~1751
D4RL (2020)	OFFLINE RL	CAUSAL	60

- Extensive tasks (covering 5 mainstream tasks)
- Normal and causal versions
- Various sequence lengths (29-4000)
- Extensive baselines (20+)

Flowformer Experiments

Task	Metrics	Flowformer	Performer	Reformer	Vanilla Transformer
Long Sequence Modeling (LRA)	Avg Acc (%) ↑	56.48	51.41	50.67	OOM
Vision Recognization (ImageNet-1K)	Top-1 Acc (%) ↑	80.6	78.1	79.6	78.7
Language Modeling (WikiText-103)	Perplexity \downarrow	30.8	37.5	33.6	33.0
Time series classification (UEA)	Avg Acc (%) ↑	73.0	71.5	71.9	71.9
Offline RL (D4RL)	Avg Reward \uparrow Avg Deviation \downarrow	$ extsf{73.5} \pm extsf{2.9}$	63.8 ± 7.6	63.9 ± 2.9	72.2 ± 2.6

Strong performance on all five mainstream tasks within the linear complexity.

Flowformer

General Relation Modeling Quadratic Complexity



Model

Long Sequence

Model Efficiency

Big Model ⊗

Flowformer

Linear complexity w.r.t. sequence length

Based on flow network & without specific inductive biases Strong performance in Long Sequence, CV, NLP, Time Series, RL

Haixu Wu, Jialong Wu, Jiehui Xu, Jianmin Wang, Mingsheng Long. Flowformer: Linearizing Transformers with Conservation Flows. ICML, 2022.

Task Adaptation



Foundation Problems





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




Catastrophic Forgetting









Negative Transfer



Choose Pre-trained Models
 Penalize smallest singular values : L_{bs}

$$\sigma_{\rm s}(F) = \eta \sum_{i=1}^{\kappa} \sigma_{-i}^2$$

• LEEP, LogME

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

Negative Transfer

- Enhance Safe Transfer
 - BSS, Zoo-tuning



- Choose Pre-trained Models
 - LEEP, LogME

Pre-trained Model Hub

Various Models and Platforms



Avoid Heavy Pre-training



Plenty of Transferable Knowledge

IMAGENET SUP. MOCO PT. MaskRCNN PT. DeepLab PT. Keypoint pt.



Same architecture Pre-trained differently

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

Transferability Assessment by LogME

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



Theoretical Guarantee of LogME

- MacKay algorithm (1992) is a heuristic method for solving the evidence maximization procedure of empirical Bayesian learning (Bishop, 1995).
- We provide the theoretical guarantee for MacKay algorithm.

Algorithm 4 One iteration of evidence maximization in Algorithm 2.

- 1: Input: α, β ; Output: α', β' for the next iteration.
- 2: Compute $A = \alpha I + \beta F^T F, m = \beta A^{-1} F^T y, \gamma = \sum_{i=1}^{D} \frac{\beta \sigma_i^2}{\alpha + \beta \sigma_i^2}$

3: Return
$$\alpha' = \frac{\gamma}{m^T m}, \beta' = \frac{n-\gamma}{||Fm-y||_2^2}$$

Theorem 1 Algorithm 4 induces a scalar function (Equation 3) with $t = \frac{\alpha}{\beta}$ and $t' = \frac{\alpha'}{\beta'}$.

$$t' = f(t) = \left(\frac{n}{n - \sum_{i=1}^{D} \frac{\sigma_i^2}{t + \sigma_i^2}} - 1\right) t^2 \frac{\sum_{i=1}^{n} \frac{z_i^2}{(t + \sigma_i^2)^2}}{\sum_{i=1}^{n} \frac{\sigma_i^2 z_i^2}{(t + \sigma_i^2)^2}}.$$
(3)

Theorem 2 If r < n and $\sum_{1 \le i,j \le n} (z_i^2 - z_j^2) (\sigma_i^2 - \sigma_j^2) > 0$, then f(t) has a fixed point and thus MacKay's algorithm will converge.

Efficiency of LogME

 \mathbb{R}^{n}

Computation Efficient --- MacKay algorithm with improved complexity.

Algorithm 2 Evidence Maximization by MacKay's Algorithm
1: Input: Extracted features $F \in \mathbb{R}^{n \times D}$ and corresponding labels $y \in$
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 $\gamma = \sum_{i=1}^{D} \frac{\beta \sigma_i^2}{\alpha + \beta \sigma^2}, \Lambda = \text{diag}\{(\alpha + \beta \sigma^2)\}$
7: Naïve : $A = \alpha I + \beta F^T F, m = \beta A^{-1} F^T y$
9: Update $\alpha \leftarrow \frac{\gamma}{m^T m}, \beta \leftarrow \frac{n - \gamma}{ Fm - \nu _2^2}$
10 and while

```
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 $\Sigma'_{i=1} = \frac{\sigma_i}{t+\sigma_i^2} (1 \le i \le r)$.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

wall-clock time		memory footprint	
fine-tune (upper bound)	161000s	fine-tune (upper bound)	6.3 GB
extract feature (lower bou	nd) 37s	extract feature (lower bound)	$43 \mathrm{MB}$
LogME	43s	LogME	$53 \mathrm{MB}$
benefit	$3700\uparrow$	benefit	$120\uparrow$
fine-tune (upper bound)	100200s	fine-tune (upper bound)	88 GB
extract feature (lower bou	nd) $1130s$	extract feature (lower bound)	$1.2~\mathrm{GB}$
LogME	1136s	LogME	$1.2 \ \mathrm{GB}$
benefit	$88\uparrow$	benefit	$73\uparrow$
	wall-clock time fine-tune (upper bound) extract feature (lower bou LogME benefit fine-tune (upper bound) extract feature (lower bou LogME benefit	wall-clock timefine-tune (upper bound)161000sextract feature (lower bound)37sLogME43sbenefit3700 \uparrow fine-tune (upper bound)100200sextract feature (lower bound)1130sLogME1136sbenefit88 \uparrow	wall-clock timememory footprintfine-tune (upper bound)161000sfine-tune (upper bound)extract feature (lower bound)37sextract feature (lower bound)LogME43sLogMEbenefit3700 \uparrow benefitfine-tune (upper bound)100200sfine-tune (upper bound)extract feature (lower bound)1130sextract feature (lower bound)LogME1136sLogMEbenefit88 \uparrow benefit





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.
- Violates the "Many could be better than all" theorem (Zhou et al. 2002).

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

Effectiveness of B-Tuning

Reduced burden of Selection and Adaptation.

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



Effectiveness of B-Tuning

Fully utilization of transferable knowledge in Model Hub.



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

Homogeneous Model Zoo

Considering models with same architecture but different knowledge.



Zoo-Tuning

Adaptively aggregate source model parameters to derive target model.



Adaptive Aggregation



Channel alignment

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

Data-dependent gating



Experiments

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

	GEN	ERAL	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.



Heterogeneous Model Hub

Design data-dependent pathways throughout the Model Hub.



Hub-Pathway

- 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.



Experiments

• Data dependent pathways \rightarrow General for heterogenous models.

Modal	General		Fi	ne-Grain	ed	Spec	Ava	
MOUEI	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 Pathways



Keypoint MoCo DeepLab MaskRCNN ImageNet

Keypoint MoCo DeepLab MaskRCNN ImageNet

Keypoint MoCo DeepLab MaskRCNN ImageNet

Remarks on Task Adaptation

	$egin{array}{c} { m Adaptation} \\ { m Accuracy}^1 \end{array}$	$\begin{array}{c} {\rm Data} \\ {\rm Efficiency}^2 \end{array}$	$\begin{array}{c} \text{Parameter} \\ \text{Efficiency}^3 \end{array}$	$\begin{array}{c} {\rm Modality} \\ {\rm Scalability}^4 \end{array}$	Task Scalability ⁵
Feature Transfer	*	**	***	***	***
Vanilla Fine-tuning	***	*	*	***	***
Domain Adaptive Tuning	***	**	*	**	***
Regularization Tuning	***	**	*	***	*
Residual Tuning	**	**	**	**	**
Parameter Difference Tuning	**	**	**	***	***
Metric Learning	*	***	***	***	*
Prompt Learning	**	***	***	*	*

¹ Accuracy when there are large-scale labeled data in downstream tasks.

² Accuracy when there are only small-scale labeled data in downstream tasks.

³ Whether parameters can be controlled when the number of downstream tasks increases.

⁴ Whether pre-trained model can be adapted to various modalities, such as text, graph.

⁵ Whether pre-trained model can be adapted to different downstream tasks, such as detection.

Domain Adaptation



Domain Adaptation



- How to measure the discrepancy between P and Q?
- Can we control target error $\epsilon_Q(h)$?

HAH-Divergence

• HAH-Divergence

$$d_{\mathcal{H}\Delta\mathcal{H}}(P,Q) = \sup_{\substack{h,h'\in\mathcal{H}}} \left| \epsilon_Q(h,h') - \epsilon_P(h,h') \right|$$



HAH-Divergence

• HAH-Divergence

$$d_{\mathcal{H}\Delta\mathcal{H}}(P,Q) = \sup_{\substack{h,h'\in\mathcal{H}}} \left| \epsilon_Q(h,h') - \epsilon_P(h,h') \right|$$

- Theorem (Generalization Bound with HΔH-Divergence)
- Denote by d the VC-dimension of hypothesis space $\mathcal H$ and ideal joint error

 $\epsilon_{ideal} = \epsilon_P(h^*) + \epsilon_Q(h^*)$. We have

$$\epsilon_{Q}(h) \leq \epsilon_{\hat{P}}(h) + \frac{d_{\mathcal{H}\Delta\mathcal{H}}(\hat{P},\hat{Q})}{m} + \epsilon_{ideal} + O\left(\sqrt{\frac{d\log n}{n}} + \sqrt{\frac{d\log m}{m}}\right)$$

Domain Adversarial Learning

Learning representation ϕ to minimize $d_{\mathcal{H}\Delta\mathcal{H}}(\phi(P), \phi(Q))$: • $\min_{\phi,h} \left\{ \mathbb{E}_{(x,y)\sim P} L(h(\phi(x)), y) + \max_{D} \left(\mathbb{E}_{P} L(D(\phi(x)), 1) + \mathbb{E}_{Q} L(D(\phi(x)), 0) \right) \right\}$ Supervised Learning on source Minimize Upper bound of $d_{\mathcal{H}A\mathcal{H}}$ ∂L_y loss L_y $\partial \theta_u$ $\overline{\partial \theta}_f$ h features class label ylabel predictor $G_u(\cdot; \theta_u)$ $\lambda \frac{\partial L_d}{\partial \theta_f}$ φ domain classifier $G_d(\cdot; \theta_d)$ Sradient Jreversal feature extractor $G_f(\cdot; \theta_f)$ layer **\bigcirc** domain label dD loss $\overline{\partial \theta}$. backprop (and produced derivatives) forwardprop

Ganin et al. Domain Adversarial Training of Neural Networks. JMLR 2016.

Domain Adversarial Learning



Theory vs. Practice



- Binary Classification vs. Multiclass Classification
- Discrete Classifier vs. Classifier with Scoring Function
- HΔH is excessively large that is hard to estimate and optimize

Disparity Discrepancy

• Disparity Discrepancy

$$d_{h,\mathcal{H}}(P,Q) = \sup_{h' \in \mathcal{H}} \left(\epsilon_Q(h,h') - \epsilon_P(h,h') \right)$$



Disparity Discrepancy

Disparity Discrepancy

$$d_{h,\mathcal{H}}(P,Q) = \sup_{h' \in \mathcal{H}} \left(\epsilon_Q(h,h') - \epsilon_P(h,h') \right)$$

- Theorem (Generalization Bound with Disparity Discrepancy)
- For any $\delta > 0$ and binary classifier $h \in \mathcal{H}$, with probability $1 3\delta$, we have $\epsilon_Q(h) \le \epsilon_{\hat{P}}(h) + d_{h,\mathcal{H}}(\hat{P},\hat{Q}) + \epsilon_{ideal} + 2\Re_{n,P}(\mathcal{H})$

$$+2\Re_{n,P}(\mathcal{H}\Delta\mathcal{H})+2\Re_{m,Q}(\mathcal{H}\Delta\mathcal{H})+2\sqrt{\frac{\log\frac{2}{\delta}}{2n}}+\sqrt{\frac{\log\frac{2}{\delta}}{2m}}$$

Yuchen Zhang & Mingsheng Long et al. Bridging Theory and Algorithm for Domain Adaptation. ICML 2019.

Margin Disparity Discrepancy

 $\begin{array}{c} & & \\$

Margin Disparity Discrepancy

$$d_{f,\mathcal{F}}^{(\rho)}(P,Q) = \sup_{f'\in\mathcal{F}} \left(\epsilon_Q^{(\rho)}(f',f) - \epsilon_P^{(\rho)}(f',f) \right)$$
 Margin Loss



 $\mathcal{H}\Delta\mathcal{H}$ -Divergence

Supremum over single space



Disparity Discrepancy

Margin Disparity Discrepancy



Margin Disparity Discrepancy

$$d_{f,\mathcal{F}}(P,Q) = \sup_{f' \in \mathcal{F}} \left(\epsilon_Q^{(\rho)}(f',f) - \epsilon_P^{(\rho)}(f',f) \right)$$
 Margin Loss

- Theorem (Generalization Bound with Margin Disparity Discrepancy)
- For any $\delta > 0$ and scoring classifier $f \in \mathcal{F}$, with probability $1 3\delta$, we have

$$\begin{split} \epsilon_{Q}(f) &\leq \epsilon_{\hat{p}}^{(\rho)}(f) + d_{h,\mathcal{H}}^{(\rho)}(\hat{p},\hat{Q}) + \epsilon_{ideal} + \frac{2k^{2}}{\rho} \Re_{n,P}(\Pi_{1}\mathcal{F}) \\ &+ \frac{k}{\rho} \Re_{n,P}(\Pi_{\mathcal{H}}\mathcal{F}) + \frac{k}{\rho} \Re_{m,Q}(\Pi_{\mathcal{H}}\mathcal{F}) + 2\sqrt{\frac{\log\frac{2}{\delta}}{2n}} + \sqrt{\frac{\log\frac{2}{\delta}}{2m}} \end{split}$$

Yuchen Zhang & Mingsheng Long et al. Bridging Theory and Algorithm for Domain Adaptation. ICML 2019.
Margin Disparity Discrepancy

• Margin Disparity Discrepancy

$$\min_{\psi,f} \max_{f'} \epsilon_{\hat{p}}^{(\rho)}(f) + \left(\epsilon_{\hat{Q}}^{(\rho)}(f',f) - \epsilon_{\hat{p}}^{(\rho)}(f',f)\right)$$



Theory vs. Practice



- A common observation is that difficulty of transfer is asymmetric.
- Previous bounds will remain unchanged after switching P and Q.
- Previous discrepancies are supremum over the whole hypothesis space.

Localized Disparity Discrepancy

Pre-train on source

Margin Disparity Discrepancy

$$d_{h,\mathcal{H}_r}(P,Q) = \sup_{\substack{h' \in \mathcal{H}_r \coloneqq \{h \in \mathcal{H} | \mathbb{E}_P L(h(x), y) \le r\}}} \left(\epsilon_Q(h,h') - \epsilon_P(h,h') \right)$$

Supremum over localized space



Localized $\mathcal{H}\Delta\mathcal{H}$ -Divergence

Supremum over localized space



Localized Disparity Discrepancy

Foundation Model

Localized Disparity Discrepancy

Pre-train on source

Localized Disparity Discrepancy

$$d_{h,\mathcal{H}_r}(P,Q) = \sup_{\substack{h' \in \mathcal{H}_r \coloneqq \{h \in \mathcal{H} \mid \mathbb{E}_P L(h(x), y) \le r\}}} \left(\epsilon_Q(h,h') - \epsilon_P(h,h') \right)$$

- Theorem (Generalization Bound with Localized Disparity Discrepancy)
- For any $\delta > 0$ and binary classifier $h \in \mathcal{H}$, with probability 1δ , we have

$$\begin{aligned} \epsilon_{Q}(h) &\leq \epsilon_{\hat{P}}(f) + d_{h,\mathcal{H}_{r}}\left(\hat{P},\hat{Q}\right) + \epsilon_{ideal} + O\left(\frac{d\log n}{n} + \frac{d\log m}{m}\right) \\ &+ O\left(\sqrt{\frac{(\epsilon_{\hat{P}}(h) + r)d\log n}{n}} + \sqrt{\frac{(\epsilon_{\hat{P}}(h) + d_{h,\mathcal{H}_{r}}\left(\hat{P},\hat{Q}\right) + r)d\log m}{m}}\right) \end{aligned}$$

Yuchen Zhang & Mingsheng Long et al. Localized Discrepancies for Domain Adaptation. arXiv 2021.

Foundation Model

Remarks on Domain Adaptation

	$\begin{array}{c} \text{Adaptation} \\ \text{Accuracy}^2 \end{array}$	Data Efficiency ²	2 Modality 2 Scalability 3	${ m Task} \ { m Scalability}^4$	Theory Guarantee ⁵
Statistics Matching	*	***	***	**	***
Domain Adversarial Learning	**	**	***	**	***
Hypothesis Adversarial Learning	***	**	***	**	***
Domain Translation	**	*	*	***	*
Semi-Supervised Learning	**	**	**	*	*

¹ Accuracy when there are large-scale data in source and target domains.

 2 Accuracy when there are only small-scale data in source and target domains.

³ Whether the model can be adapted to various modalities, such as text, time series.

 $\frac{4}{2}$ Whether the model can be adapted to different tasks, such as regression, detection.

⁵ Whether the generalization error of target domain can be theoretically bounded in adaptation.

Transfer Learning Library

Library Public Public

Transfer Learning Library for Domain Adaptation, Task Adaptation, and Domain Generalization





Thank You!



Mingsheng Long (龙明盛) Tsinghua University mingsheng@tsinghua.edu.cn



Jianmin Wang (王建民) Tsinghua University jimwang@tsinghua.edu.cn



Michael I. Jordan (迈克尔・欧文・乔丹) UC Berkeley jordan@cs.berkeley.edu



Yuchen Zhang Zhangjie Cao (张育宸) (曹张杰)



(朱晗)

Yue Cao (曹越)



(游凯超)

(江俊广)



(王希梅)



(陈新阳)

Yang Shu (树扬)





