







2D Edges



2D Keypoints

Depth



Image Reshading



3D Keypoints

Measuring Transferability in Transfer Learning 迁移学习中的可迁移性度量

Yang Li Dec 13,2021

中科院深圳先进技术研究院

TBSI 清华-伯克利深圳学院 Tsinghua-Berkeley Shenzhen Institute











Transfer Learning using a Pre-Trained Model Source coach Input S Task S table kitchen $\bullet \bullet \bullet$ Layer n Transfer frozen weights **Back-propagation** Target Mountain Beach Task T Input T Sea A A $\bullet \bullet \bullet$ fine-tuning **Back-propagation**



• Commonly used in computer vision, NLP etc



Transfer Learning with Pre-Trained Model

Improve target training efficiency, reduce number of target labeled data needed



Assumes represenation of S is transferable to T





Why is Knowing Task Transferability Important?

Source-task selection



e.g. Select the best word/sentence encoder for NLP tasks

• Learn more transferable features

• Task transfer policy learning



(Empirical) Transferability

(or accuracy) of the transferred network on target data



Given source data (X_s, Y_s) and source model (θ_s, h_s) , compute expected log loss



Related Works — Empirical Transferability

- Feature transferability in Neural Network (Yosinski 2014)
- Taskonomy (Zamir et. al 2018) for 2D/3D scene understanding tasks. Shape Inductive Biases (Feinman & Lake 2018) for 3D shapes

Limitation:

- minimum
- inefficient to compute

A. R. Zamir, A. Sax, W. Shen, L. Guibas, J. Malik and S. Savarese, "Taskonomy: Disentangling Task Transfer Learning," 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, 2018

• need to train the transfer network using gradient descend, subject to local

Can we estimate the transfer performance without any training of the target network?



Analytical Transferability Metric

• An analytical way to *estimate* empirical transferability

Algorithm	Cross-Task	Cross-Instance	Cross-Domain	
	$P(Y_S X_S) \neq P(Y_T X_T)$	$X_S \neq X_T$	$P(X_S) \neq P(X_T)$	
NCE (Tran et al. 2019) *	\checkmark	X	X	
H-Score (Bao et al. 2019)	\checkmark	\checkmark	X	first transferab work for task tra
LEEP (Nguyen et al. 2020)**	\checkmark	\checkmark	X	learning
OTCE (Tan et al. 2021)	\checkmark	\checkmark	\checkmark	
LogME (K. You 2021) ***	\checkmark	\checkmark	\checkmark	

* Anh T Tran, Cuong V Nguyen, and Tal Hassner. Transfer- ability and hardness of supervised classification tasks. ICCV, 2019. ** Cuong V Nguyen, Tal Hassner, Cedric Archambeau, and 952 Matthias Seeger. Leep: A new measure to evaluate transferability of learned representations.ICML, 2020.

*** K. You, Y. Liu, J. Wang, and M. Long, "Logme: Practical assessment of pre-trained models for transfer learning," in ICML. PMLR, 2021, pp. 12 133-12 143.





Talk Outline

- Analytical Transferability Metrics
 - H-score (efficient cross-task transferability metric)
 - OT-CE (cross-domain, cross-task transferability metric)
- Transferability-guided fine-tuning
- Measuring Transferability for Medical Segmentation





The Task Transferability Problem

Given:

- Input X, source task label Y_5 , target task label Y_T
- Trained source model with optimal feature $f_{s}(X)$



The Transfer Network

fs(**X**)

Task S

Ys YT Task T

coach

table ΤV

Transferability of $S \rightarrow T$: to what extent can fs help learning target task (X,Y_T)?

Assume same domain $P(X_s)=P(X_t)=P(X)$

Retrain-head setting



Task Transferability

Bao, Yajie, Yang Li et al. "An information-theoretic approach to transferability in task transfer learning." 2019 IEEE International Conference on Image Processing (ICIP). IEEE, 2019.

Transferability from Task S to Task T





$\mathfrak{T}(S,T) \triangleq \frac{\text{Target Performance of } f_S}{\text{Optimal Target Performance}}$

$$\begin{cases} \mathfrak{T}(S,T) = 1 & \textcircled{\vdots} \\ 0 \leq \mathfrak{T}(S,T) \leq 1 \\ \mathfrak{T}(S,T) = 0 & \textcircled{\vdots} \end{cases}$$

How to measure the performance of f_s(X) on target task (X,Y_T) ?

Measuring Feature Effectiveness -Neural Network Perspective

- Classification using log-loss:
 - X, Y random variables; f(X) a zero-mean feature
 - **Expected log loss:** $L(f;\theta) = \mathbb{E}_{X,Y}[L(f(X),Y;\theta)]$

$$L(f, \theta^{\star}) = Const(X, Y) - H(f) + o(\epsilon^{2})$$

H-score of f(X)
$$Higher H-score => Better Performance$$
$$H(f) = tr(cov(f(X))^{-1}cov(\mathbb{E}_{P_{X|Y}}[f(X)|Y]))$$



By Local information geometry [Huang 2018], given feature f(X), the optimal loss is





Interpretation of $\mathcal{H}(f)$

Intuition in latent space

$$\mathcal{H}(f) = \operatorname{tr}(\operatorname{cov}(f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))^{-1}\operatorname{cov}(\mathbb{E}_{X|Y}[f(X))$$

Statistical operational meaning

- ratio predictor based on f(X)
- Higher H-score ↔ Faster error decay rate



H-score characterizes the asymptotic error probability of the likelihood



Computing Transferability

$$\mathfrak{T}(S,T) = \frac{\mathscr{H}_T(f_S)}{\mathscr{H}_T(f_T)}$$

H-score of source feature $\mathcal{H}_T(f_S)$

- Easy to compute
- O(mk²) time complexity
- Maximal H-score: $\mathscr{H}_T(f_T)$
 - Discrete X: Alternating Conditional Expectation (ACE) algorithm (Huang et. al. 2015)
 - Continuous X: Neural network formulation (Wang et al. 2017)

Python Code for H-Score

```
def Hscore(f,Y):
  Covf=np.cov(f)
  alphabetY=list(set(Y))
  g=np.zeros_like(f)
  for z in alphabetY:
    g[Y==y]=np.mean(f[Y==y,:], axis=0)
  Covg=np.cov(g)
  score=np.trace(np.dot(np.linalg.pinv(Covf,
         rcond=le-15), Covg))
  return score
```



Results: Image Classification Feature Selection

- Source task: ImageNet 1000 classification (ResNet50 features) from 6 layers 4a-5f)
- Target task: Cifar 100-class classification on 20,000 images







Results: Source Task Selection for 3D Scene Understanding



Query Image



2D Edges



3D (Occlusion) Edges

- 8 image-based tasks from Taskonomy dataset (Zamir et al. 2018) • 2 classification tasks: object-class, scene-class
- - 6 2D/3D image-to-image tasks: average H-score over all superpixels
- Source models: pre-trained task-specific models (4,000,000 training samples);
- Target model: linear feature transfer using 20,000 images (64 x 64)





Image Reshading



2D Keypoints

3D Keypoints

Depth

J.	
Ĩ	
-	7

Transferability Ranking



target task

source task



Task Relationships

Cluster the source task transferability scores for each target task.





Comparison with Task Af

- Reference metric: task affinity, an emp transferability score (Amir et al. 2018)
- Ranking comparison metric: Spearman's correlation and Discounted cumulative ga

Advantage of our approach:

- Efficiency: five times more efficient than Affinity
- Clear operational meaning based on statistics & information theory

finity					
		Spearman	DCG		_
oirical	edge2d	0.381	1.000		- ^ a
Jiiicat	keypoint2d	0.357	1.000		
	edge3d	0.429	0.851		
	keypoint3d	0.786	0.765		
ain (DCG)	reshade	0.810	0.998		- 0.0
	depth	0.738	0.996		- 0.5
	object class.	0.214	0.976		- 0.4
	scene class.	0.286	0.981		- 0.3
	F	Rank Co	mparis	30	n



Discussion on H-Score

- Proposed an efficient, easy-to-com operational meaning.
 - Theoretically proven for classification tasks.
 - Validated on image-processing, vision recognition applications
- Shrinkage based H-score (Ibrahim et. al. 2021) : improved the stability of covariance estimation in H-score computation

Ibrahim, S., Ponomareva, N., & Mazumder, R. (2021). Newer is not always better: Rethinking transferability metrics, their peculiarities, stability and performance. *arXiv preprint arXiv:2110.06893*.

Proposed an efficient, easy-to-compute task transferability metric with clear

Talk Outline

- Analytical Transferability Metrics
 - H-score (efficient cross-task transferability metric)
 - OT-CE (cross-domain, cross-task transferability metric)
- Transferability-guided fine-tuning
- Measuring Transferability for Medical Segmentation





Transferability across tasks and domains



Domain Difference

Domain: Painting

transductive domain adaptation $P(X_S) \neq P(X_T)$





cross-domain cross-task transfer





Terminologies

- where $x_s^i, x_t^i \in \mathcal{X}$ $y_s^i \in \mathcal{Y}_s, y_t^i \in \mathcal{Y}_t$ • Source dataset: $D_s = \{(x_s^i, y_s^i)\}_{i=1}^m \sim P_s(x, y)$ where • Target dataset: $D_t = \{(x_t^i, y_t^i)\}_{i=1}^n \sim P_t(x, y)$
- Source model: (θ, h_s) where $\theta : \mathcal{X} \to \mathbb{R}^d$ $h_s : \mathbb{R}^d \to \mathcal{P}(\mathcal{Y}_s)$





OTCE Transferability Score

Yang Tan, Yang Li*, and Shao-Lun Huang. "OTCE: A Transferability Metric for Cross-Domain Cross-Task Representations." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15779-15788. 2021.Code and Supplementary Material:



joint work with Yang Tan

Compute Domain Difference

Optimal transport problem



$$\alpha \in \mathcal{P}(\mathcal{X}), \beta \in \mathcal{P}(\mathcal{X})$$
$$) \triangleq \{\pi \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) | P_{1\#}\pi = \alpha, P_{2\#}\pi = \beta\}$$
$$\alpha, \beta) \triangleq \min_{\pi \in \Pi(\alpha, \beta)} \int_{\mathcal{X} \times \mathcal{X}} c(x, y) d\pi(x, y)$$

When c is $c(x, y)^d$, OT is the pth-Wasserstain Distance W^p

Advantage of Wasserstein Distance

Observe the underlying geometry



• Provide a map about moving mass

$KL(p_1; p_2) = KL(p_1; p_3)$ $MMD(p_1, p_2) = MMD(p_1, p_3)$

 $W(p_1, p_2) < W(p_1, p_3)$





Compute Domain Difference

• Given two datasets

$$D_A = (x_A^i, y_A^i)_{i=1}^m \sim P_A(x, y) \qquad D_B = (x_B^j, y_B^j)_{j=1}^n \sim P_B(x, y)$$

• Regularized optimal transport (Sinkhorn Algorithm)

$$OT(D_A, D_B) \triangleq \min_{\pi} \int_{\mathcal{X} \times \mathcal{X}} c(x_A, x_B) d\pi(x_A) d\pi($$

• Domain difference as wasserstein distance

$$W(D_A, D_B) = \sum_{i,j=1}^{m,n} \pi^*(x_A^i, x_B^j) c(x_A^i, x_B^j)$$

 $(A, x_B) + \epsilon H(\pi)$

[1] Cuturi, Marco. "Sinkhorn distances: Lightspeed computation of optimal transport." Advances in neural information processing systems. 2013.



Task Difference: Theoretical Motivation

Given X and multi-task labels Y and Z

$$k_Y = \operatorname*{argmin}_{k \in K} \mathcal{L}_Y(w_Z, k)$$

 Approximat transferability by conditional entropy

$$\widetilde{\mathrm{Trf}}(T^Z \to T^Y) = l_Y(w_Z, k_Y)$$
$$\widetilde{\mathrm{Trf}}(T^Z \to T^Y) = \widetilde{\mathrm{Trf}}(T^Z \to T^Y)$$

Challenge: how to compute H(Y|Z) when source and target inputs are not the same?

Tran, Anh T., Cuong V. Nguyen, and Tal Hassner. "Transferability and hardness of supervised classification tasks." Proceedings of the IEEE International Conference on Computer Vision. 2019.



$$) \ge l_Z(w_Z, h_Z) - H(Y|Z)$$





Computing Task Difference

- Build sample correspondence using the Coupling Matrix (joint distribution) π^* from the OT step
- Compute empirical distributions of labels

$$\hat{P}(y_s, y_t) = \sum_{i,j:y_s^i = y_s, y_t^j = y_t} \pi_{ij}^*, \ \hat{P}(y_s) = \sum_{y_t \in \mathcal{Y}_t} \hat{P}(y_s, y_t).$$
$$W_T = H(Y_t | Y_s) = -\sum_{y_t \in \mathcal{Y}_t} \sum_{y_s \in \mathcal{Y}_s} \hat{P}(y_s, y_t) \log \frac{\hat{P}(y_s, y_t)}{\hat{P}(y_s)}$$

Conditional Entropy

 $\pi^* = \begin{bmatrix} P(x_A^1, x_B^1) & P(x_A^1, x_B^2) & \dots & P(x_A^1, x_B^n) \\ P(x_A^2, x_B^1) & P(x_A^2, x_B^2) & \dots & P(x_A^2, x_B^n) \\ \vdots & \vdots & \ddots & \vdots \\ P(x_A^n, x_B^1) & P(x_A^n, x_B^2) & \dots & P(x_A^n, x_B^n) \end{bmatrix}$

$$\pi_{ij}^*, \ \hat{P}(y_s) = \sum_{y_t \in \mathcal{Y}_t} \hat{P}(y_s, y_t).$$

OTCE Score

Linear combination between Domain Difference and Task Difference

 $OTCE = \lambda_1$

• Learn coefficients from auxiliary tasks with known empirical transferability

Practical consideration: to avoid using auxiliary task, we can use task difference WT alone to characterize transferability $OTCE_{sim} = W_T$

$$_1\hat{W}_D + \lambda_2\hat{W}_T + b$$

Experiment Datasets (Classification)

• **DomainNet:** 345 categories in 5 domains

• Office31: 32 categories in 3 domains.

For a source domain D_{src} , randomly sample 100 classification tasks in each target domain $D_{tgt} \in D \setminus \{D_{src}\}$

С	Clipart		chu futebook	SCHOOL O-O	
Ρ	Painting	(A)			3
Q	Quickdraw	Ċ			
R	Real				
S	Sketch				

Α	Amazon	
D	DSLR	
W	Webcam	

Experiment Setup

• Transfer settings

Setting Names	target sample size	label size		
standard	all target task samples	1-50) tasks with	
few-shot	10 samples	1-50	diverse hardness	
fixed category size	category size all target task samples		tasks with similar	
		•	I I AI UI IESS	

- Evaluation metric: compare the rank and empirical transferability
 - Spearsman's rank correlation ρ
 - Kendall's rank τ

• Evaluation metric: compare the ranking of target tasks ordered by OT-CE score

Experiment Result

Standard & Few-shot settings

		Source	Target			Spearma	n / Kendall		
Setting	Dataset	domain	domain	OTCE	OTCE _{sim}	LEEP [6]	NCE [4]	H-score [5]	LogME [15]
		С	P,Q,R,S	0.976 / 0.861	0.966 / 0.839	0.932 / 0.779	0.825 / 0.670	0.920 / 0.748	0.867 / 0.667
		Р	C,Q,R,S	0.977 / 0.868	<u>0.960</u> / <u>0.822</u>	0.906 / 0.743	0.849 / 0.686	0.937 / 0.777	0.929 / 0.761
	DomainNet	Q	C,P,R,S	<u>0.961</u> / <u>0.826</u>	0.963 / 0.832	0.953 / 0.810	0.943 / 0.793	0.942 / 0.784	0.912 / 0.744
		R	C,P,Q,S	0.975 / 0.863	<u>0.951</u> / <u>0.808</u>	0.910 / 0.747	0.872 / 0.707	0.942 / 0.786	0.855 / 0.670
Standard		S	C,P,Q,R	0.969 / 0.842	<u>0.967</u> / <u>0.839</u>	0.965 / 0.834	0.962 / 0.830	0.950 / 0.802	0.908 / 0.733
		А	D,W	0.848 / 0.643	<u>0.831</u> / <u>0.619</u>	0.817 / 0.606	0.817 / 0.604	0.601 / 0.417	0.660 / 0.459
	Office31	D	A,W	0.885 / 0.702	0.839 / 0.651	<u>0.862</u> / <u>0.676</u>	0.851 / 0.664	0.464 / 0.335	0.172 / 0.119
		W	A,D	0.859 / 0.676	0.801 / 0.614	<u>0.812</u> / <u>0.626</u>	0.805 / 0.616	0.524 / 0.371	0.470 / 0.318
			Average	0.931 / 0.785	<u>0.910</u> / <u>0.753</u>	0.895 / 0.728	0.866 / 0.696	0.785 / 0.628	0.722 / 0.559
		С	P,Q,R,S	0.926 / 0.756	0.909 / 0.729	0.836 / 0.640	0.745 / 0.576	0.762 / 0.567	0.731 / 0.524
		Р	C,Q,R,S	0.931 / 0.772	0.886 / 0.701	0.803 / 0.618	0.746 / 0.575	0.811 / 0.608	0.849 / 0.649
	DomainNet	Q	C,P,R,S	0.821 / 0.631	<u>0.829</u> / <u>0.636</u>	0.798 / 0.602	0.782 / 0.584	0.813 / 0.614	0.866 / 0.682
		R	C,P,Q,S	0.929 / 0.769	<u>0.853</u> / <u>0.666</u>	0.770 / 0.589	0.728 / 0.559	0.845 / 0.652	0.774 / 0.574
Few-shot		S	C,P,Q,R	0.914 / 0.742	<u>0.895</u> / <u>0.710</u>	0.872 / 0.680	0.872 / 0.679	0.838 / 0.645	0.867 / 0.684
		А	D,W	0.859 / 0.662	0.845 / 0.640	0.818 / 0.609	0.811 / 0.602	0.651 / 0.456	0.659 / 0.460
	Office31	D	A,W	0.929 / 0.773	<u>0.925</u> / <u>0.766</u>	<u>0.925</u> / 0.764	0.924 / 0.765	0.429 / 0.308	0.002 / 0.021
		W	A,D	<u>0.927</u> / <u>0.765</u>	0.929 / 0.767	0.916 / 0.749	0.919 / 0.752	0.316 / 0.235	0.250 / 0.181
			Average	0.905 / 0.734	<u>0.884</u> / <u>0.702</u>	0.842 / 0.656	0.816 / 0.637	0.683 / 0.511	0.625 / 0.472

Experiment Result

• Same category set size setting: similar target task complexity

Sourc	ce Target			Spearr	nan / Kendall		
doma	in domair	D OTCE	$OTCE_{sim}$	LEEP [6]	NCE [4]	H-score [5]	LogME [15]
С	P,Q,R,S	0.701 / 0.500	<u>0.687</u> / <u>0.487</u>	0.685 / 0.486	0.666 / 0.472	-0.438 / -0.290	-0.222 / -0.151
Р	C,Q,R,S	0.670 / 0.485	<u>0.631</u> / <u>0.448</u>	0.630 / 0.446	0.612 / 0.430	-0.529 / -0.371	-0.043 / -0.039
Q	C,P,R,S	0.341 / 0.225	0.316 / 0.211	0.210 / 0.136	0.291 / 0.191	-0.256 / -0.172	0.066 / 0.037
R	C,P,Q,S	0.637 / 0.455	0.598 / 0.415	0.587 / 0.407	0.586 / 0.406	-0.094 / -0.063	-0.382 / -0.252
S	C,P,Q,R	0.428 / 0.292	$\overline{0.436}$ / $\overline{0.299}$	0.404 / 0.277	<u>0.432</u> / <u>0.298</u>	-0.247 / -0.164	0.027 / 0.006
	Average	0.555 / 0.391	<u>0.534</u> / <u>0.372</u>	0.503 / 0.350	0.517 / 0.359	-0.313 / -0.212	-0.111 / -0.080
				1			

Transferability between different task pairs less distinguishable => More challenging scenario

Experiment Results

Visualization of transferability experiment results on DomainNet



Experiments: Transferability Applications

Source Model Selection



Multi-Source Feature Fusion



Table 2. Quantitative comparisons of source model selection accu-
racy (%) among transferability metrics on DomainNet.

		Target domain						
N	lethod	С	Р	Q	R	S	average	
L	EEP[24]	31.1	26.7	5.6	97.8	100.0	52.2	
Ν	[CE[38]	41.1	94.4	2.2	100.0	100.0	67.5	
<u> </u>	TCE	41.1	93.3	97.8	100.0	100.0	86.4	
1								





Semantic Segmentation Task

Input image



• A fundamental task in Autonomous driving, medical image analysis etc

Semantic label



Which feature to transfer?

• U-Net architecture commonly used for semantic segmentation



Challenge in computing transferability: Output Y is high-dimensional

Observation on prediction error map



Input image



Predicted error map



Semantic label



Extract edge from semantic label

OTCE for Semantic Segmenation

- Solution:
 - Sample N pixels from all images

 Y_{s}

 Y_t

• Compute OTCE_{sim} over the feature set of the sampled labels

Sample Heuristic: sample pixels near segmentation boundaries





Experiment Datasets (Segmentation)

Cityscapes (real captured)











BDD100K (real captured)

GTA5 (computer game)





Results: Intra-dataset transfer

- Two cities in Cityscapes
- Comparing OT-CE_{sim} with LEEP and LogME



Results: cross-dataset transfer

- 6 model architectures: Fcn8s, UNet, SegNet, PspNet, FrrnA, FrrnB
- Comparing OT-CE_{sim} with transfer accuracy

source: BDD100k target: Cityscapes



et, SegNet, PspNet, FrrnA, FrrnB r accuracy

source: GTA5 target: Cityscapes



Comparison with Other Works

Spearman Kendall								
Transfer setting	Target task	$OTCE_{sim}$	LEEP [6]	LogME [15]	$OTCE_{sim}$	LEEP [6]	LogME [15]	
	aachen	0.774	0.627	-0.005	0.642	0.484	-0.074	
	cologne	0.750	0.639	0.620	0.579	0.459	0.474	
Intra-dataset Transfer	erfurt	0.750	0.585	0.565	0.575	0.432	0.400	
	jena	<u>0.735</u>	0.868	0.561	<u>0.579</u>	0.695	0.400	
	strasbourg	<u>0.791</u>	0.680	0.838	<u>0.632</u>	0.505	0.684	
	aachen	0.771	0.371	0.371	0.600	0.333	0.200	
Inter-dataset Transfer	cologne	0.657	0.371	0.600	0.467	0.333	0.333	
(source: BDD100K)	erfurt	0.086	0.714	0.257	<u>0.200</u>	0.600	0.067	
	jena	<u>0.600</u>	0.314	0.657	0.467	0.200	0.467	
	strasbourg	0.657	0.429	0.657	0.467	0.333	0.467	
	aachen	0.200	0.314	0.429	0.067	0.200	0.333	
Inter-dataset Transfer	cologne	0.829	0.429	<u>0.771</u>	0.733	0.200	0.600	
(source: GTA5)	erfurt	<u>0.600</u>	0.543	0.943	<u>0.467</u>	0.333	0.867	
	jena	0.714	0.257	0.886	0.467	0.200	0.733	
	strasbourg	0.886	-0.029	<u>0.429</u>	0.733	0.067	0.200	
	Average	0.653	0.474	0.572	0.512	0.358	<u>0.410</u>	
bold der otes the best result and <u>underline</u> denotes the 2 nd best result.								

OT-CE Discussion

- Shown to be effective for both classification and semantic segmentation tasks
- OT-CE and OT-CEsim can out-perform other state-of-the-art transferability metrics in crossdomain cross-task settings
- Only a small number of auxiliary tasks are needed for OT-CE



Talk Outline

- Analytical Transferability Metrics
 - H-score (efficient cross-task transferability metric)
 - OT-CE (cross-domain, cross-task transferability metric)
- Transferability-guided fine-tuning for few-shot transfer learning
- Measuring Transferability for Medical Segmentation

OT-CE Score based Fine Tuning

• Step 1: Minimize task difference



• Step 2: Fine-tune on target data



Learning transferable feature Source data Target data $D_s = \{ (x_s^i, y_s^i) \}_{i=1}^m \sim P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) \qquad D_t = \{ (x_t^i) \}_{i=1}^m \geq P_s(x, y) > D_t = P_s(x, y)$

1) Solve the Optimal Transport problem via Sinkhorn Algorithm

$$OT(D_s, D_t) \triangleq \min_{\pi \in \Pi(D_s, D_t)} \sum_{i,j=1}^{m,n} c(\theta(x_s^i), \theta(x_t^j))\pi_{ij} + \epsilon H(\pi)$$

$$\hat{P}(y_s, y_t) = \sum_{i, j: y_s^i = y_s, y_t^j = y_t} \pi_{ij}^*, \ \hat{P}(y_s) = \sum_{y_t \in \mathcal{Y}_t} \hat{P}(y_s, y_t).$$

2) Minimize the conditional entropy (task difference) $\min_{a} H(Y_t|Y_s) = -\min_{a}$

Then we can finetune (theta, h t) on the target training data and evaluate the transfer performance

Source Model

$$\{x_t^i, y_t^i\}_{i=1}^n \sim P_t(x, y)$$
 (θ, h_s)

$$\sum_{y_t \in \mathcal{Y}_t} \sum_{y_s \in \mathcal{Y}_s} \hat{P}(y_s, y_t) \log \frac{\hat{P}(y_s, y_t)}{\hat{P}(y_s)}$$



Few-Shot Transfer Learning Dataset

Dataset visualization

MNIST (number 0-9)



USPS (number 0-9)



Omniglot (characters of 1,623 classes)

Few-Shot Transfer Learning Results

source: **MNIST** (number **0-6**) target: **USPS** (number **7-9**) (3-way-5-shot task)

Method	Target testing acc
Train from scratch	81.16%
Vanilla finetune	86.43%
OTCE-based finetune	91.29% (+4.86%)



source: **MNIST** (number **0-9**) target: **Omniglot** (randomly sample 100 5-way-5shot tasks)

Method	Target testing acc		
Vanilla finetune	86.11%		
OTCE-based finetune	90.52% (+4.41%)		



Few-Shot Transfer Learning Results



vanilla fine-tuning.

On average. OTCE-based fine tuning has higher transferability score and higher accuracy than

Comparison with other few-shot learning methods

Model	Method	$MNIST \rightarrow Omniglot$	Caltech101 \rightarrow MiniImageNet	
FewshotNet	MAML [29] MatchingNet [27] ProtoNet [30] RelationNet [28] Vanilla finetune OTCE-based finetune	$\begin{array}{c} 88.60 \pm 1.14\% \\ 87.92 \pm 1.10\% \\ 83.11 \pm 1.34\% \\ 69.35 \pm 1.62\% \\ 91.30 \pm 0.95\% \\ \textbf{92.32} \pm \textbf{0.87\%} \end{array}$	$28.23 \pm 0.44\% \\ 44.75 \pm 1.30\% \\ 50.40 \pm 1.35\% \\ 29.55 \pm 0.61\% \\ 49.49 \pm 1.27\% \\ \mathbf{51.36 \pm 1.33\%}$	meta learni based meth
LeNet	Vanilla finetune OTCE-based finetune	$\begin{array}{c} 86.11 \pm 1.10\% \\ \textbf{90.52} \pm \textbf{0.94}\% \end{array}$	_	
ResNet-18	Vanilla finetune OTCE-based finetune	_	$\begin{array}{c} 48.48 \pm 1.39\% \\ 50.02 \pm \mathbf{1.34\%} \end{array}$	



Fine-Tuning Discussion

- setting
- performance on target task.

Meta learning doesn't work as well as Vanilla fine-tuning in cross-domain

• Minimizing the task difference characterized by OTCE is an effective way to obtain a better initial source weights so that we can achieve higher transfer

Talk Outline

- Analytical Transferability Metrics
 - H-score (efficient cross-task transferability metric)
 - OT-CE (cross-domain, cross-task transferability metric)
- Transferability-guided fine-tuning for few-shot transfer learning
- Measuring Transferability for Medical Segmentation

Transfer Learning in MRI Segmentation Tasks

 An important problem in federated learning and distributed meidcal AI systems

iSeg2019: brain matter segmentation

3 tasks, 2 modalities

T1w



Ground Truth







FeTS Challenge: brain tumor segmentation 4 tasks, 4 modalities



ground truth



Empirical Study on Transferability for Brain MRI Segmentation What affects



6 modality-task combinations

12 modality-task combinations

Finding #1: Transferring from related tasks is easier then unrelated tasks

ground truth label

brain tumor recognition



tumor core

brain matter recognition



white matter

transfer from related source



whole tumor -> tumor core



gray matter -> white matter

transfer from unrelated source



gray matter -> tumor core



whole tumor -> white matter

Finding #1: Transferring from related tasks is easier then unrelated tasks

many-to-one transfer experiments





Finding #2: Cross-task transfer is easier than cross-modality transfer

		Target Task					
		ET,T1CE	ED,T1CE	NCR,T1CE	WM,T1	GM,T1	CSF,T1
Source	same task different modality	0.755	0.731	0.726	0.864	0.881	0.935
Tasks	different task same modality	0.821	0.786	0.782	0.877	0.892	0.934

Mean Transfer Accuracy (dice loss)





Mutually exclusive case: "complementary tasks" are more transferable than "non-complementary tasks"

CSF

GM





Mutually exclusive case: "complementary tasks" are more transferable than "non-complementary tasks"

CSF

GM





Subset case: "smaller task" is more transferable than "bigger task"





Subset case: "smaller task" is more transferable than "bigger task"





Subset case: "smaller task" is more transferable than "bigger task"

Source Task Selection



Source(s) Selection Pipeline



1 Target Task: ET-22-T1CE, Correlation: 0.371204



Source Task Selection

- Reduced the size of candidate source task pool to 16
- (smaller rank displacement)

Target	Method	Top 1	Top 2	Top 3	Top 4
	H-score w/o F	5	10	22	27
ЕТ 22 Т2	H-score w/ F	4	5	6	7
121-22-12	OTCE w/o F	2	2	4	12
	OTCE w/F	2	2	4	7
	H-score w/o F	14	24	30	40
ET 20 T1	H-score w/ F	0	9	9	13
L1-20-11	OTCE w/o F	2	14	17	23
	OTCE w/F	2	11	13	17

Analytical transferability metrics are more accurate on the filtered candidate set

Spearsman's Footrule (Total Rank Displacement)

still work in progress

Future Works on Medical Image Transfer Learning

Inspired by how neuroradiologist uses expert knowledge in diagnosis

- Use transferability to automatically decide which tasks to be learned together in multi-• task learning
- Design a multi-task transfer curriculum to perform hierachical transfer

Brain MRI Image Understainding Tasks



low level processing tasks

structure

tumor

white matter lesion

Alzhemer's

mid-level detection tasks

high-level diagnosis

Summary

Estimating transferability is important in practical transfer learning

- For same-domain transferability, H-score is very efficient and theoretically proven in information theory (Shrinkage-based H-Score resolves numerical issues)
- For cross-domain case, OT-CE and OT-CE_{sim} are robust to many challenging scenarios, with less data assumption
- Transferability metrics can lead to cross task/domain feature learning algorithms (HGR-regression and OT-CE based fine tuning)

design based on transferability

Future work: efficient transferability for regression problems, transfer strategy

Acknowledgement



黄绍伦 (TBSI)



Leonidas Guibas (Stanford)







鲍亚捷 (TBSI, MS 2016) University of Georgia





awesome collaborators

Thank You! Q&A



Related Papers and Code

- (ICIP), pp. 2309-2313. IEEE, 2019.
- Code and Supplementary Material:
 - http://yangli-feasibility.com/home/ttl.html
 - <u>https://github.com/tanyang1231/OTCE_Transferability_CVPR21</u>

• Yajie Bao, Yang Li, Shao-Lun Huang*, Lin Zhang, Lizhong Zheng, Amir Zamir, and Leonidas Guibas. "An information-theoretic approach to transferability in task transfer learning." In 2019 IEEE International Conference on Image Processing

• Yang Tan, Yang Li*, and Shao-Lun Huang. "OTCE: A Transferability Metric for Cross-Domain Cross-Task Representations." In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15779-15788. 2021.