

Quantifying and Improving Transferability in Domain Generalization

Guojun Zhang, Han Zhao, Yaoliang Yu and Pascal Poupart



UNIVERSITY OF ILLINOIS
URBANA-CHAMPAIGN

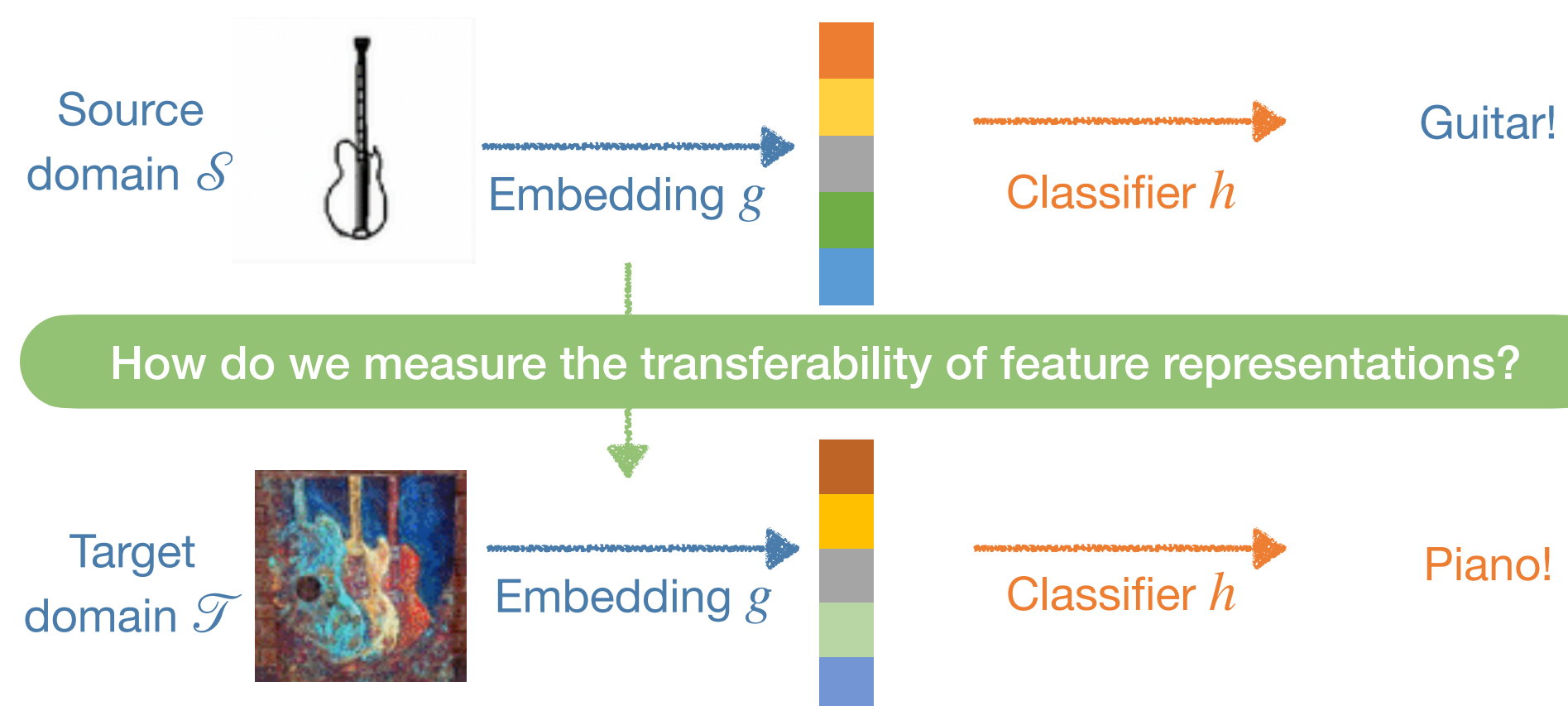
VECTOR INSTITUTE



Main result

We propose an evaluation metric for measuring the **transferability** in domain generalization. Based on this metric, we design algorithms to further improve out-of-domain generalization over SoTA methods.

Transfer learning

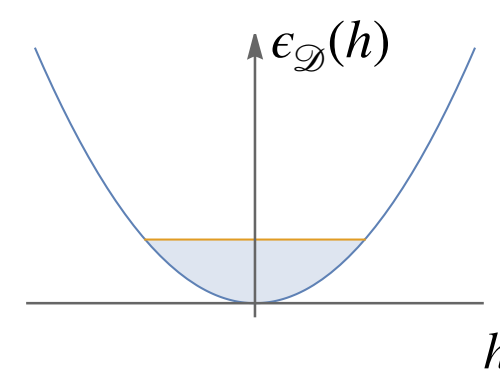


Transferable:
Every near-optimal source classifier is also near-optimal on the target

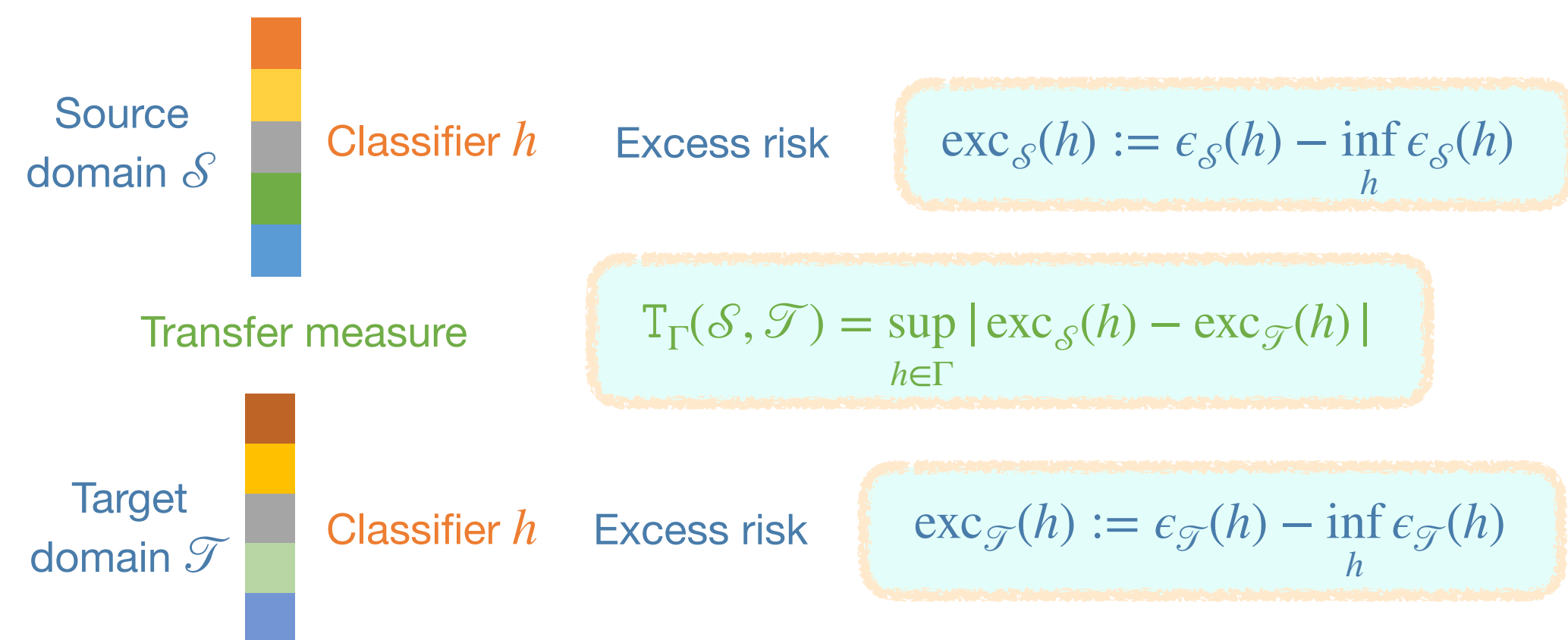
$$\text{argmin}(\epsilon_{\mathcal{D}}, \delta) := \{h \in \mathcal{H}, \epsilon_{\mathcal{D}}(h) \leq \inf_{h \in \mathcal{H}} \epsilon_{\mathcal{D}}(h) + \delta\} \quad \delta\text{-minimal set}$$

Definition: \mathcal{S} is $(\delta_{\mathcal{S}}, \delta_{\mathcal{T}})$ -transferable to \mathcal{T} if

$$\text{argmin}(\epsilon_{\mathcal{S}}, \delta_{\mathcal{S}}) \subseteq \text{argmin}(\epsilon_{\mathcal{T}}, \delta_{\mathcal{T}})$$



Use **Transfer Measures** to quantify transferability



Transferable \equiv Small transfer measure, if $\Gamma = \text{argmin}(\epsilon_{\mathcal{S}}, \delta)$.

Theoretical results

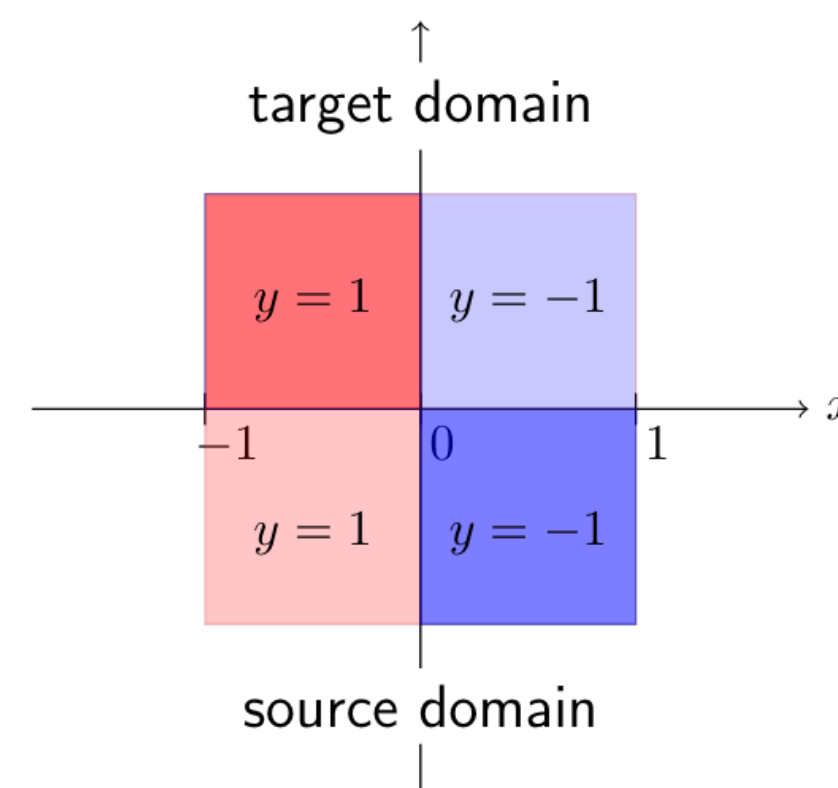
Generalization bound (tighter than the one using \mathcal{H} -divergence)

$$\epsilon_{\mathcal{T}}(h) \leq \epsilon_{\mathcal{S}}(h) + T_{\Gamma}(\mathcal{S}, \mathcal{T}) + \epsilon_{\mathcal{T}}^* - \epsilon_{\mathcal{S}}^*, \quad \epsilon_{\mathcal{D}}^* = \inf_{h \in \mathcal{H}} \epsilon_{\mathcal{D}}(h)$$

If the optimal errors are zero, transfer measure with 0-1 loss is equivalent to **total variation** when $\Gamma = \mathcal{H}_l$ includes all binary classifiers

$$\text{Total variation} \quad d_{TV}(\mathcal{S}, \mathcal{T}) = \sum_y \int |p_{\mathcal{S}}(x, y) - p_{\mathcal{T}}(x, y)| dx$$

Two domains can still be **transferable** even if the joint distributions are **dissimilar**



Invariance principles (error depends on both the marginal and conditional distributions)

- Transferability: invariance of excess risks (*joint*)
- H-divergence: invariance of feature distributions $p(z)$ (*marginal*)
- Invariant risk minimization (IRM): invariance of the optimal predictors (*conditional*)
- Generalized label shift: invariance of feature distributions within the same class $p(z|y)$ (*conditional*)

How to Compute Transfer Measures?

We make the following approximations

- We use **surrogate loss** to approximate 0-1 loss
- We use **the final loss** after training to approximate **the optimal loss**
- We use a **neighborhood in the parameter space** to approximate the **minimal set**
- The resulting transfer measure is a **lower bound**, but we can still use it to **refute** transferability

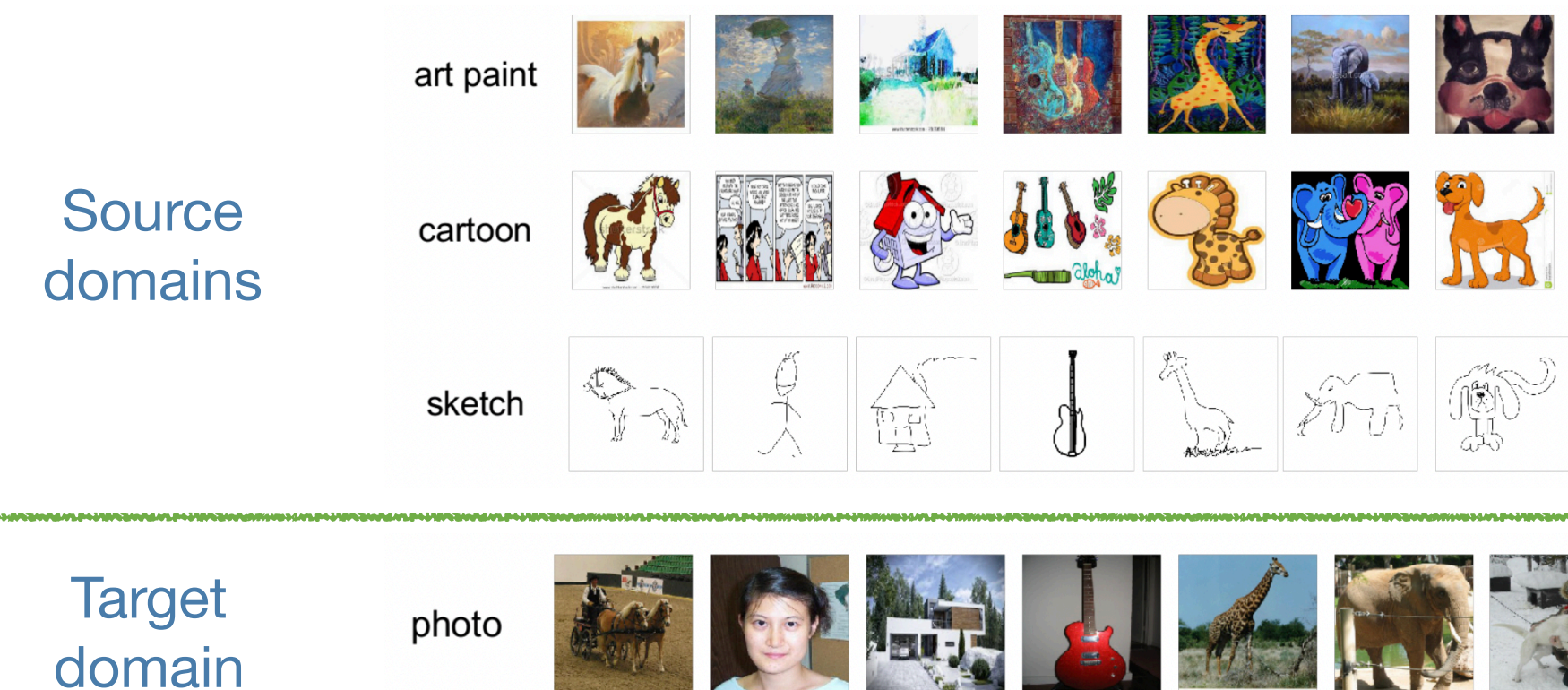
$$\text{Final result as a lower bound of transfer measure} \quad \sup_{\|\theta - \hat{\theta}^*\| \leq \delta} \epsilon_{\mathcal{T}}(h) - \epsilon_{\mathcal{S}}(h) - \epsilon_{\mathcal{T}}(\hat{h}^*)$$

θ : parameter of classifier $\hat{h}^* / \hat{\theta}^*$: the learned classifier after training

If the lower bound is large, then it is not transferable

Algorithms

Domain Generalization (DG) learns feature embeddings and classifiers from **several** source domains



Suppose we have n source domains: $\mathcal{S}_1, \dots, \mathcal{S}_n$ and a target domain $\mathcal{T} = \mathcal{S}_0$

$$\text{Evaluate transferability} \quad \max_{i, j \in [0, n]} \max_{\|h - \hat{h}^*\| \leq \delta} \text{exc}_{\mathcal{S}_i}(h \circ g) - \text{exc}_{\mathcal{S}_j}(h \circ g)$$

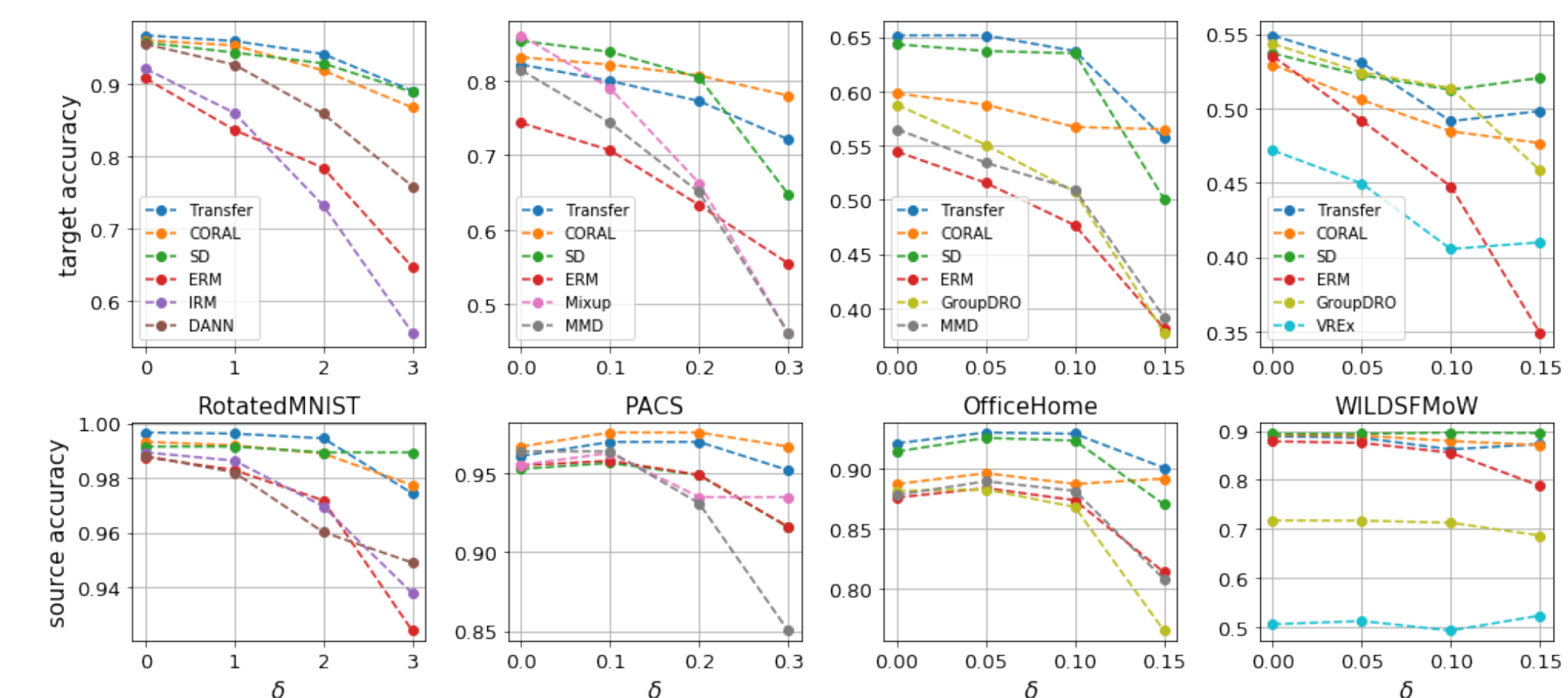
Improve transferability

$$\min_{g, h} \frac{1}{n} \sum_{i=1}^n \epsilon_{\mathcal{S}_i}(h \circ g) + \max_{i, j \in [1, n]} \max_{\|h' - h\| \leq \delta} \text{exc}_{\mathcal{S}_i}(h' \circ g) - \text{exc}_{\mathcal{S}_j}(h' \circ g)$$

ERM

adversarial training

Experiments



- Many popular algorithms are not learning transferable features: a near-optimal source classifier is **not** near-optimal on the target
- Our **Transfer** algorithm learns more **transferable** features than ERM, DANN, IRM, GroupDRO, etc.
- Best performers: **Transfer** (ours), **CORrelation ALignment (CORAL)**, **Spectral Decomposition (SD)**

References

- Arjovsky et al. "Invariant risk minimization." *arXiv preprint arXiv:1907.02893* (2019).
- Ben-David et al. "A theory of learning from different domains." *Machine Learning*, 2010
- Koltchinskii, "Rademacher complexities and bounding the excess risk in active learning," *JMLR* 2010
- Tachet des Combes et al. "Domain adaptation with conditional distribution matching and generalized label shift." *NeurIPS* 2020

Full talk: <https://www.youtube.com/watch?v=Ce3PyHA54Gj>

Code: <https://github.com/Gordon-Guojun-Zhang/Transferability-NeurIPS2021>

