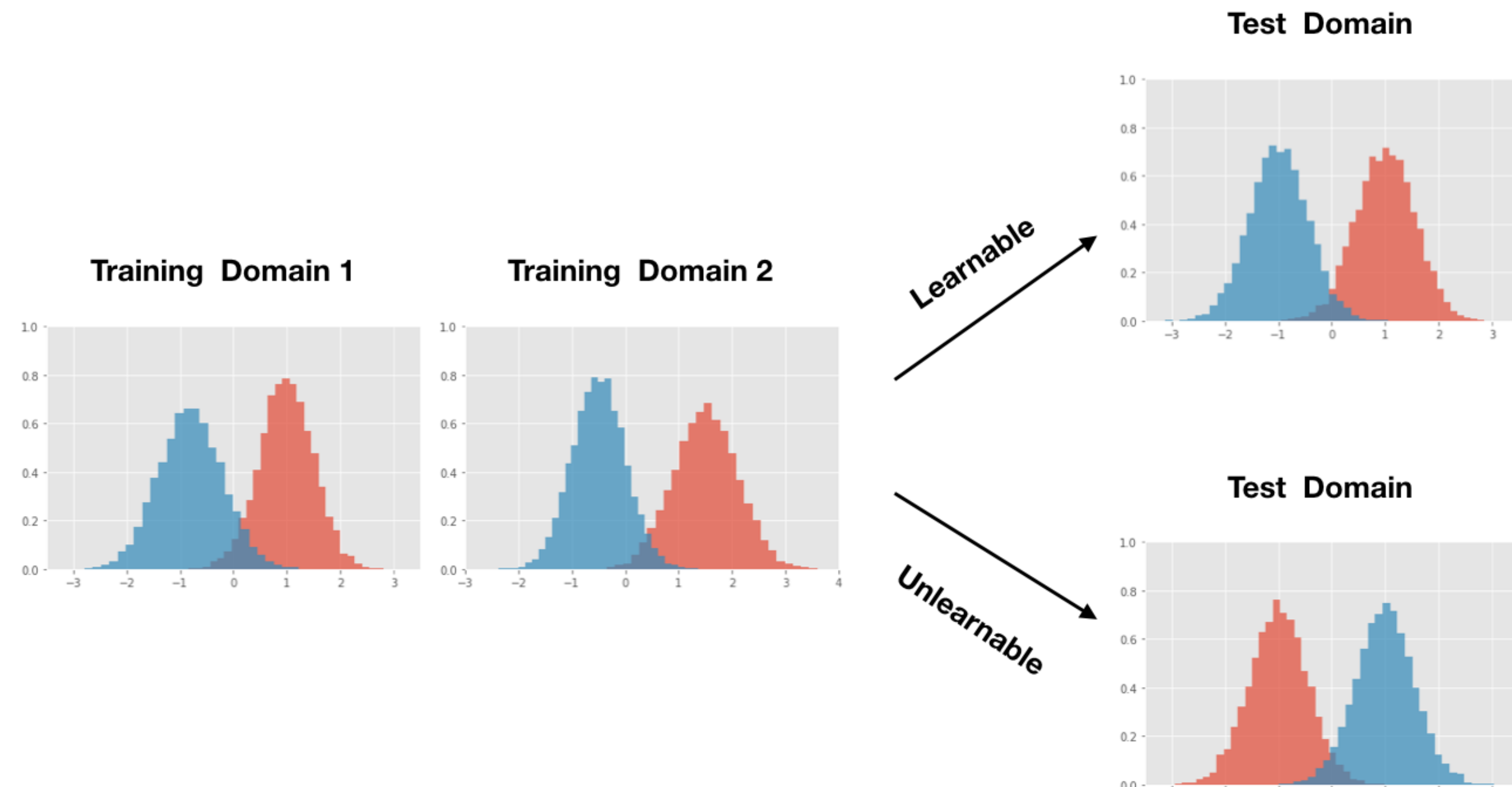


Towards a Theoretical Framework of Out-of-Distribution Generalization

Haotian Ye¹, Chuanlong Xie², Tianle Cai¹, Ruichen Li¹, Zhenguo Li², Liwei Wang¹

¹Peking University; ²Huawei Noah's Ark Lab

Takeaways



- Out-of-Distribution (OOD) generalization is to train on sets of domains, and generalize to unseen domains.
- As the above figure shows, generalization to any domains is impossible. We lack theories of **when we can guarantee generalization (OOD – learnable)**.
- Intuition: the *extent* to which the *invariance* of *informative* feature is preserved, determines OOD learnability. **We build a framework to describe OOD learnability, bound OOD generalization error based on it, and propose an effective model selection algorithm.**

Our Framework (Informal)

- $\mathcal{E}_{avail} \subset \mathcal{E}_{all}$ are two sets of domains. We train on \mathcal{E}_{avail} , and (hope to) generalize to \mathcal{E}_{all} .
- Hypothesis Space: $f = g \circ h$, where h is d -dimensional feature extractors, and g is top classifier.
- For 1-dim feature ϕ , define $P_y^e = \text{Prob}(\phi(X^e|Y = y))$.
- Goal: minimizes maximum loss of f in all domains, i.e., $\mathcal{L}(f, \mathcal{E}) \triangleq \max_{e \in \mathcal{E}} \ell(f, e)$.

❖ **Variation** of ϕ across domain set \mathcal{E} :

$$V(\phi, \mathcal{E}) = \max_{y \in Y} \max_{e, e' \in \mathcal{E}} \text{dist}(P_y^e, P_y^{e'})$$

❖ **Informativeness** of ϕ across domain set \mathcal{E} :

$$I(\phi, \mathcal{E}) = \text{avg}_{y \neq y'} \min_{e \in \mathcal{E}} \text{dist}(P_y^e, P_{y'}^e)$$

❖ **OOD Learnability**: $(\mathcal{E}_{avail}, \mathcal{E}_{all})$ is $(s(\cdot), \delta)$ -learnable, if for all $\phi \in \Phi$ satisfying $I(\phi, \mathcal{E}_{avail}) \geq \delta$, we have $s(V(\phi, \mathcal{E}_{avail})) \geq V(\phi, \mathcal{E}_{all})$.

Here $s(\cdot)$ is a special type of monotonically increasing function, which we call **expansion function**.

Theorems and Model Selections

- We prove that OOD generalization error can be bounded with our framework.
- **Main Theorem**: Under some conditions, if the problem is $(s(\cdot), I(h, \mathcal{E}_{avail}))$ -learnable, then
$$\mathcal{L}(f, \mathcal{E}_{all}) - \mathcal{L}(f, \mathcal{E}_{avail}) \leq O\left(s(V(h, \mathcal{E}_{avail}))^{\frac{\alpha^2}{(\alpha+d)^2}}\right)$$
- **Lower Bound**: Exists a $(s(\cdot), \delta)$ -learnable problem, where the optimal classifier f with $V(h, \mathcal{E}_{avail}) = \varepsilon$ has
$$\mathcal{L}(f, \mathcal{E}_{all}) - \mathcal{L}(f, \mathcal{E}_{avail}) \geq \Omega\left(s(V(h, \mathcal{E}_{avail}))\right)$$
- We also propose a **model selection algorithm**, and it outperforms other selection methods:

Algorithm 1: Model Selection

Input: available dataset $\mathcal{X}_{avail} = (\mathcal{X}_{train}, \mathcal{X}_{val})$, candidate models set \mathcal{M} , var_acc_rate r_0 .

for $f = g \circ h$ **in** \mathcal{M} **do**

for i **in** $[d]$ **do**

$\hat{V}_i \leftarrow \max_{y \in \mathcal{Y}, \mathcal{X}^e \neq \mathcal{X}^{e'} \in \mathcal{X}_{avail}} \text{Total Variation}(\mathbb{P}(\phi_i^e|y), \mathbb{P}(\phi_i^{e'}|y));$ ▷ Use GPU KDE

end

$\mathcal{V}_f \leftarrow \text{mean}_{i \in [d]} \hat{V}_i$

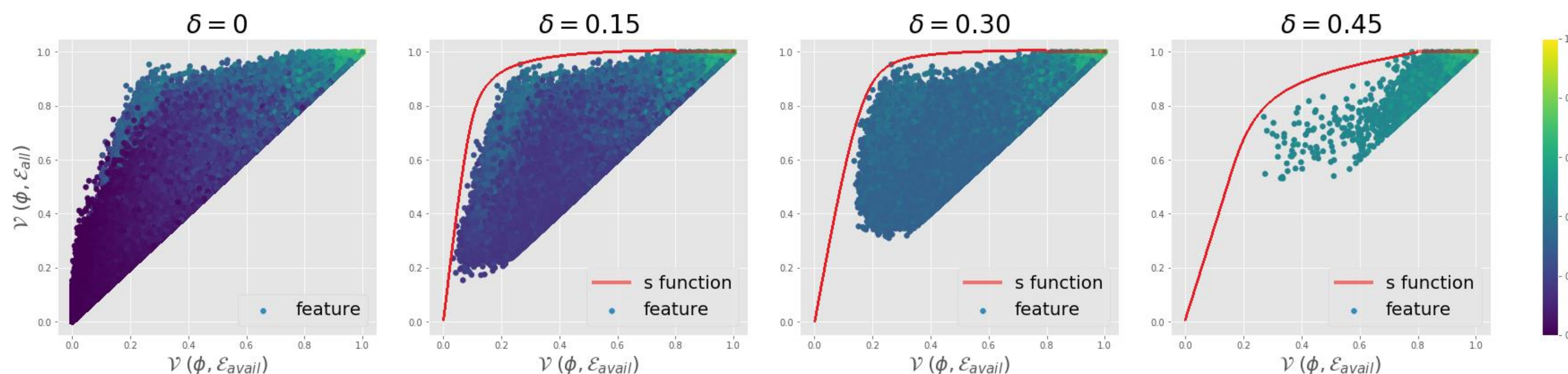
$\text{Acc}_f \leftarrow \text{compute validation accuracy of } f \text{ using } \mathcal{X}_{val}$

end

Return $\text{argmax}_{f \in \mathcal{M}} (\text{Acc}_f - r_0 \mathcal{V}_f)$

Features' V and I (lightness) in Office-Home [1].

Larger δ corresponds to flatter $s(\cdot)$.



PACS	Env	A	C	P	S	avg	acc inc
	Val	85.20%	80.42%	96.17%	77.86%	84.91%	-
	Ours	88.72%	81.74%	96.83%	79.00%	86.57%	1.66%↑
OfficeHome	Env	A	C	P	R	avg	acc inc
	Val	61.85%	55.56%	74.72%	76.25%	67.09%	-
	Ours	65.76%	55.07%	75.20%	76.31%	68.09%	1.00%↑
VLCS	Env	C	L	S	V	avg	acc inc
	Val	97.46%	64.83%	69.50% ⁶	70.97%	75.69%	-
	Ours	97.81%	66.98%	69.50%	70.97%	76.32%	0.63%↑

[1] H Venkateswara, et al. Deep hashing network for unsupervised domain adaptation.

