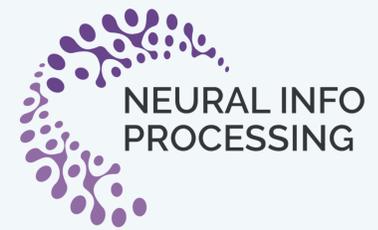




# On $f$ -Divergence Principled Domain Adaptation: An Improved Framework

Ziqiao Wang<sup>1</sup> Yongyi Mao<sup>2</sup>  
<sup>1</sup>Tongji University <sup>2</sup>University of Ottawa



## Unsupervised Domain Adaptation

- Unknown distributions  $\mu$  and  $\nu$
- Labeled source-domain sample  $\mathcal{S} = \{X_i, Y_i\}_{i=1}^n \sim \mu^{\otimes n}$
- Unlabelled target-domain sample  $\mathcal{T} = \{X_j\}_{j=1}^m \sim \nu^{\otimes m}$
- Goal: Efficiently transfer ML models between related domains at low cost  
 $\implies$  Find a hypothesis  $h \in \mathcal{H}$  “works well” on  $\nu$

## Limitations of Previous $f$ -Divergence-based DA Works

- Relying on a Weak Variational Representation (i.e. Eq. 1) 😞  
 $\implies$  Cannot recover Donsker and Varadhans representation of KL divergence.
- Slow Rate of Sample Complexity Bound 😞  $\implies$  e.g.,  $\mathcal{O}(\frac{1}{\sqrt{n}} + \frac{1}{\sqrt{m}})$ .
- Gap Between Theory and Algorithm 😞  
 $\implies$  e.g., Overestimation of the  $f$ -divergence in [1].

## Main Contributions

- We design a novel  $f$ -divergence-based domain discrepancy measure, termed  $f$ -DD, and derive an upper bound for the target error.
- To improve the convergence rate of our  $f$ -DD-based bound, we refine it using a localization technique.
- Our  $f$ -DD outperforms previous  $f$ -divergence-based algorithms on three popular UDA benchmarks.

## Background on $f$ -Divergence

- ( $f$ -Divergence) Let  $P$  and  $Q$  be two distributions on  $\Theta$ . Let  $\phi: \mathbb{R}_+ \rightarrow \mathbb{R}$  be a convex function with  $\phi(1) = 0$ . If  $P \ll Q$ , then

$$D_\phi(P||Q) \triangleq \mathbb{E}_Q \left[ \phi \left( \frac{dP}{dQ} \right) \right],$$

e.g., Total variation, KL,  $\chi^2$ , squared Hellinger, Jeffereys, Jensen-Shannon, etc.

- Variational Representation of  $f$ -divergence.
- Original Legendre Transformation:

$$D_\phi(P||Q) = \sup_{g \in \mathcal{G}} \mathbb{E}_{\theta \sim P} [g(\theta)] - \mathbb{E}_{\theta \sim Q} [\phi^*(g(\theta))]. \quad (1)$$

- Reparameterization of  $g \rightarrow g + \alpha$  (“Shift Transformation”) [2]

$$D_\phi(P||Q) = \sup_g \mathbb{E}_{\theta \sim P} [g(\theta)] - \inf_{\alpha \in \mathbb{R}} \{ \mathbb{E}_{\theta \sim Q} [\phi^*(g(\theta) + \alpha) - \alpha] \}. \quad (2)$$

Eq. (2) is point-wise “tighter” than Eq. (1)

$\Downarrow$

- Example: Donsker and Varadhans (DV) representation of KL divergence:  $\phi(x) = x \log x - x + 1$ , then  $\phi^*(y) = e^y - 1$

- By Eq. (1)

$$D_{\text{KL}}(P||Q) = \sup_{g \in \mathcal{G}} \mathbb{E}_P [g(\theta)] - \mathbb{E}_Q [e^{g(\theta)} - 1]. \quad (3)$$

- By Eq. (2)

$$D_{\text{KL}}(P||Q) = \sup_{g \in \mathcal{G}} \mathbb{E}_P [g(\theta)] - \log \mathbb{E}_Q [e^{g(\theta)}]. \quad (4)$$

Eq. (4) recovers the DV representation of KL, Eq. (4) is pointwise tighter than Eq. (3) by  $\log(x) \leq x - 1$  for  $x > 0$ .

## Discrepancy-based DA Theory: Target Error Bound

- Additional Notations
  - Loss function  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_0^+$ . Assumptions:  $\begin{cases} \text{Triangle property} \\ \text{Bounded loss} \end{cases}$
  - Target error:  $R_\nu(h) \triangleq \mathbb{E}_{(X,Y) \sim \nu} [\ell(h(X), Y)]$ , Source error:  $R_\mu(h) \triangleq \mathbb{E}_{(X,Y) \sim \mu} [\ell(h(X), Y)]$ .
  - We use  $\ell(h, h')$  to denote  $\ell(h(x), h'(x))$ , i.e. the disagreement of  $h$  and  $h'$  on  $x$ .
- [1] defines:

$$\tilde{D}_\phi^{h, \mathcal{H}}(\mu||\nu) \triangleq \sup_{h' \in \mathcal{H}} |\mathbb{E}_\mu [\ell(h, h')] - \mathbb{E}_\nu [\phi^*(\ell(h, h'))]|$$

$\implies$  Additional absolute value function added. 😞

- Theory (Target Error Bound):

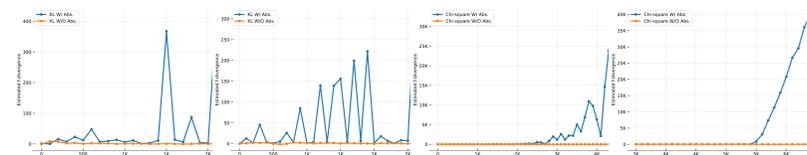
$$\text{Target Error} \leq \text{Source Error} + \tilde{D}_\phi^{h, \mathcal{H}}(\mu||\nu) + \text{Ideal Joint Error.}$$

$\implies$  Absolute value function is necessary for establishing this bound

- $f$ -Domain Adversarial Learning ( $f$ -DAL) Algorithm:

$$\min_h R_{\hat{\mu}}(h) + \max_{h'} \underbrace{\mathbb{E}_{\hat{\mu}} [\ell(h, h')] - \mathbb{E}_{\hat{\nu}} [\phi^*(\ell(h, h'))]}_{d(\hat{\mu}, \hat{\nu}; h)}.$$

$\implies d(\hat{\mu}, \hat{\nu}; h)$  drops the absolute value function compared with  $\tilde{D}_\phi^{h, \mathcal{H}}(\mu||\nu)$



(a) KL (Office31) (b) KL (OfficeHome) (c)  $\chi^2$  (Office31) (d)  $\chi^2$  (OfficeHome)

Figure 1. The  $y$ -axis is the estimated  $f$ -divergence and the  $x$ -axis is the # of iterations.

$\implies f$ -DAL algorithm fails if the absolute value function is added.

- Our  $f$ -DD:

$$D_\phi^{h, \mathcal{H}}(\nu||\mu) \triangleq \sup_{t \in \mathbb{R}, h' \in \mathcal{H}} \mathbb{E}_\nu [t\ell(h, h')] - \inf_{\alpha \in \mathbb{R}} \mathbb{E}_\mu [\phi^*(t\ell(h, h') + \alpha) - \alpha].$$

$\implies$  Introducing the scaling parameter  $t$  (i.e. “Affine Transformation”) 😞.

## Theorem (informal): $f$ -DD-based Bound

For any  $h \in \mathcal{H}$ ,

$$\text{Target Error} \leq \text{Source Error} + \inf_{t \geq 0} \frac{D_\phi^{h, \mathcal{H}}(\nu||\mu) + K_\mu(t)}{t} + \text{Ideal Joint Error,}$$

where  $K_\mu(t)$  is the upper bound for the “general CGF” for  $\mu$ .

- Ideal joint error can be  $\min_{h^* \in \mathcal{H}} R_\nu(h^*) + R_\mu(h^*)$  [3] or  $\min\{R_\nu(f_\mu), R_\mu(f_\nu)\}$  [5].
- If  $\phi$  is twice differentiable and  $\phi''$  is monotone, then  $\inf_{t \geq 0} \frac{D_\phi^{h, \mathcal{H}}(\nu||\mu) + K_\mu(t)}{t} = \sqrt{\frac{2}{\phi''(1)}} D_\phi^{h, \mathcal{H}}(\nu||\mu)$ , e.g.,  $\phi''(1) = 1$  for KL recovers [4, Theorem 4.2].
- Sample complexity bound: w.h.p.,

$$D_\phi^{h, \mathcal{H}}(\nu||\mu) \leq D_\phi^{h, \mathcal{H}}(\hat{\nu}||\hat{\mu}) + \text{Complexity Terms} + \mathcal{O}(1/\sqrt{n} + 1/\sqrt{m}).$$

## Shaper Bound: Localization Technique

- Restricted Hypothesis Space (Rashomon set):  $\mathcal{H}_r \triangleq \{h \in \mathcal{H} | R_\mu(h) \leq r\}$
- Localized  $f$ -DD: For a given  $h \in \mathcal{H}_{r_1}$

$$D_\phi^{h, \mathcal{H}_r}(\nu||\mu) \triangleq \sup_{h' \in \mathcal{H}_r, t \geq 0} \mathbb{E}_\nu [t\ell(h, h')] - \inf_{\alpha \in \mathbb{R}} \mathbb{E}_\mu [\phi^*(t\ell(h, h') + \alpha) - \alpha].$$

## Theorem (informal): Localized KL-DD-based Bound

For any  $h \in \mathcal{H}_{r_1}$ , w.h.p.

$$\text{Target Error} \leq \text{Source Train Error} + \mathcal{O}(D_{\text{KL}}^{h, \mathcal{H}_r}(\hat{\nu}||\hat{\mu})) + \mathcal{O}(1/n + 1/m) \\ + \mathcal{O}\left(\sqrt{(r+r_1)/n} + \sqrt{r/m}\right) + (r, r_1)\text{-Related Terms} + \text{Complexity.}$$

Small  $r, r_1 \implies$  fast decaying rate (i.e.  $\mathcal{O}(\frac{1}{n} + \frac{1}{m})$ ) 😊.

## Algorithms and Experiments

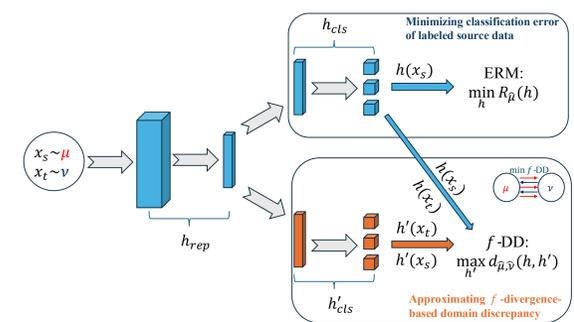


Figure 2. Overview of  $f$ -DD.

- Three specific discrepancy measures:
  - KL-DD,  $\chi^2$ -DD, the weighted Jeffereys-DD:  $\gamma_1 D_{\text{KL}}(\hat{\mu}||\hat{\nu}) + \gamma_2 D_{\text{KL}}(\hat{\nu}||\hat{\mu})$
- Objective Function:

Bounded  $\ell \rightarrow$  Unbounded  $\hat{\ell}$  (Optimizing over  $t$  may not be necessary) 😞

$$\min_h R_{\hat{\mu}}(h) + \max_{h'} \left\{ \mathbb{E}_{\hat{\mu}} [\hat{\ell}(h, h')] - \inf_{\alpha} \mathbb{E}_{\hat{\nu}} [\phi^*(\hat{\ell}(h, h') + \alpha) - \alpha] \right\}.$$

Table 1. Accuracy (%) on UDA Classification Tasks

Method	Office-31	Office-Home	Digits
$f$ -DAL [1]	89.5	68.5	96.3
Ours (KL-DD)	89.8	69.4	96.9
Ours ( $\chi^2$ -DD)	89.7	69.2	96.4
Ours (Jeffereys-DD)	<u>90.1</u>	<u>70.2</u>	<u>97.1</u>

## Key Takeways

- Utilizing stronger variational formulations of  $f$ -divergences can obtain improved results in both theoretical analysis and algorithmic performance.
- The best performance is achieved by Jeffereys-DD, suggesting the value of manually adjusting the asymmetric properties in domain adaptation.

## References

- David Acuna, Guojun Zhang, Marc T Law, and Sanja Fidler.  $f$ -domain adversarial learning: Theory and algorithms. In International Conference on Machine Learning, pages 66–75. PMLR, 2021.
- Rohit Agrawal and Thibaut Horel. Optimal bounds between  $f$ -divergences and integral probability metrics. In International Conference on Machine Learning, pages 115–124. PMLR, 2020.
- Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. Advances in neural information processing systems, 19, 2006.
- Ziqiao Wang and Yongyi Mao. Information-theoretic analysis of unsupervised domain adaptation. In International Conference on Learning Representations, 2023.
- Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In International Conference on Machine Learning, pages 7523–7532. PMLR, 2019.