

# Unsupervised Domain Adaptation Based on Source-guided Discrepancy

23<sup>th</sup> Sep.

Han Bao (The University of Tokyo / RIKEN AIP)



東京大学  
THE UNIVERSITY OF TOKYO



# Research interests

supervised learning + real-world constraints

- Learning theory: how to handle performance metrics for **class-imbalance**

[[Bao & Sugiyama 19](#)] (in submission)

- Reinforcement learning **with low-cost data**

[[WCBTS19](#)] (ICML2019) Imitation Learning from Imperfect Demonstration

- Domain adaptation: how to learn when **training  $\neq$  test**

today's topic

[[KCBHSS19](#)] (AAAI2019)

Unsupervised Domain Adaptation Based on Source-guided Discrepancy

- Weak supervision: how to learn **without labels**

[[BNS18](#)] (ICML2018)

Classification from Pairwise Similarity and Unlabeled Data

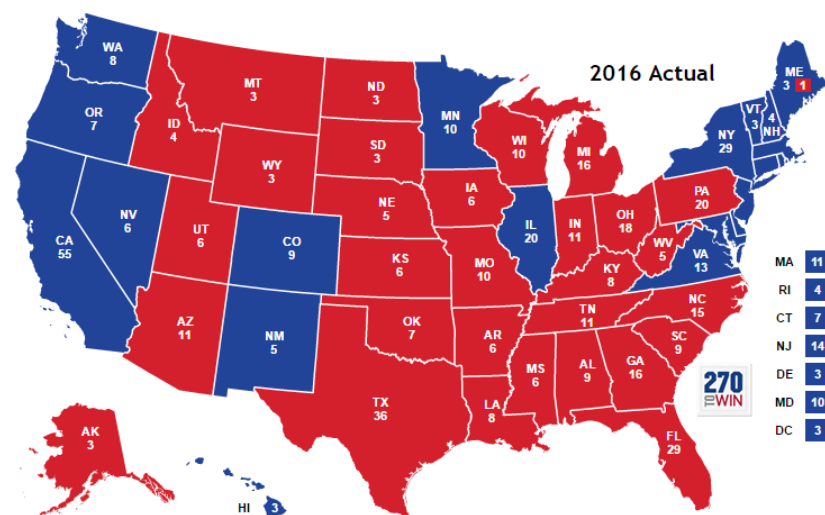
# Inference in Real-world

## ■ Prediction of President Election

[Brownback & Novotny 2018]

- ▶ cf. social desirability bias
- ▶ tend to answer in the ways “what others desire”
- ▶ unexpected results in 2016 US president election

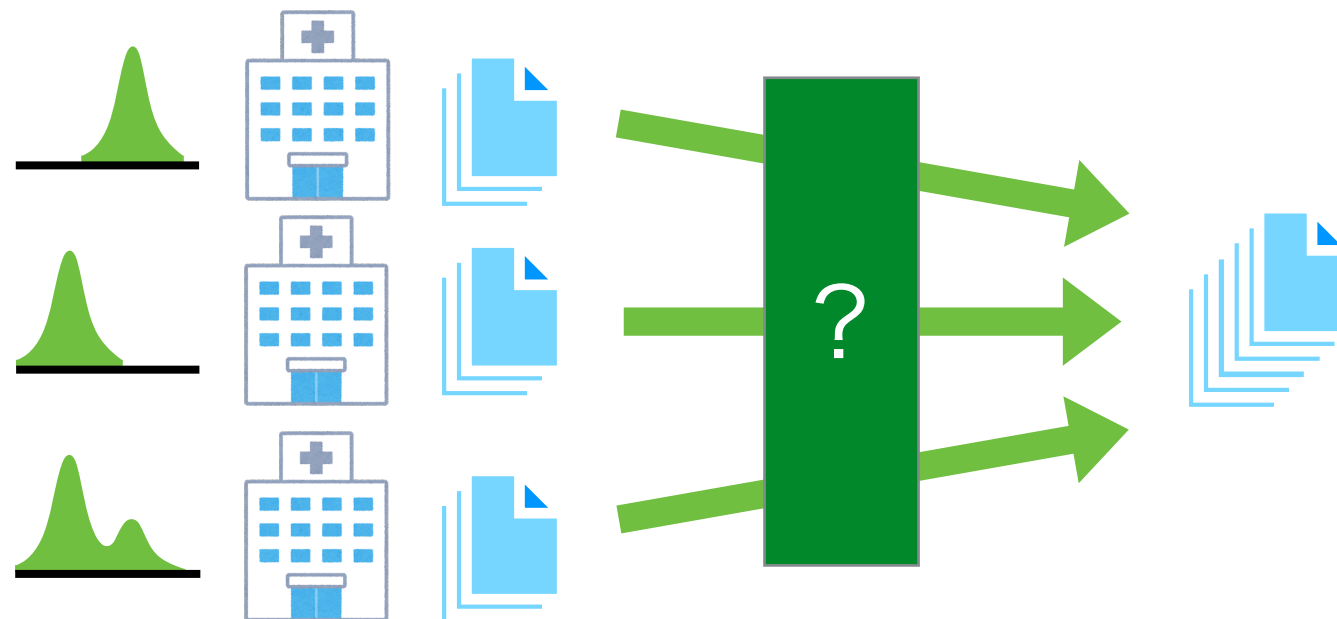
**Hard to obtain real answers!**



# Inference in Real-world

- Integration of hospital databases [Wachinger & Reuter 2016]
  - ▶ CAD (Computer-Aided Diagnosis) prevailing
  - ▶ each hospital has limited amount of data
  - ▶ want to unify among hospitals as much as possible

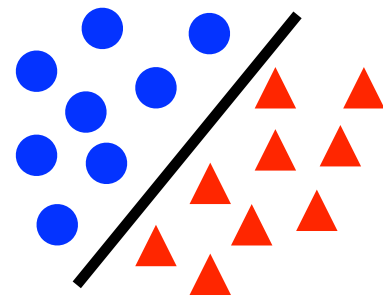
**Data distribution may differ!**



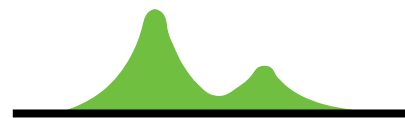
# What's transfer learning?

## ■ Usual machine learning

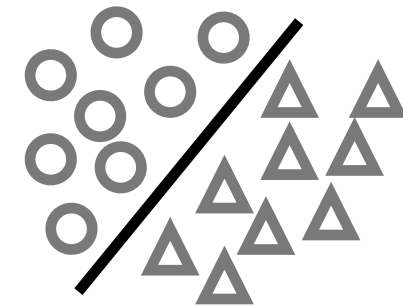
training data



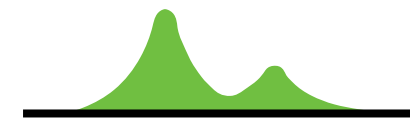
training  
distribution



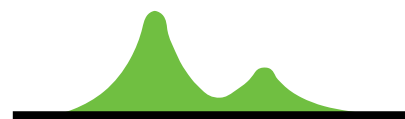
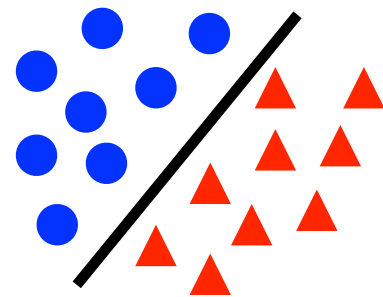
test data



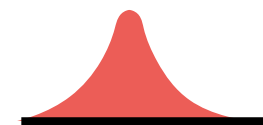
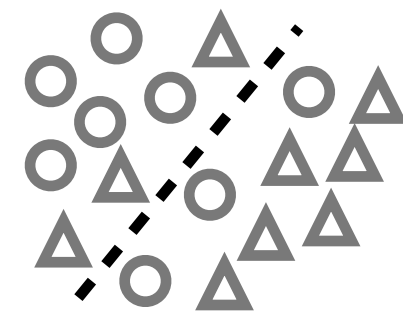
test  
distribution



## ■ Transfer learning



$\neq$

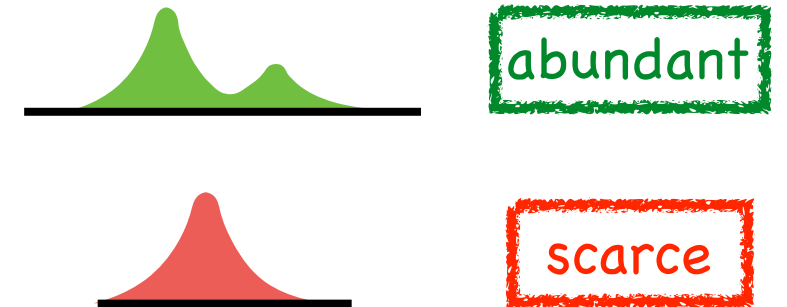


Many terminologies: transfer learning, covariate shift adaptation, domain adaptation, multi-task learning, etc.

# Unsupervised Domain Adaptation <sup>6</sup>

## Input

- ▶ training **labeled** data:  $\{x_i, y_i\} \sim p_S$   
(source)
- ▶ test **unlabeled** data:  $\{x'_j\} \sim p_T$   
(target)



## Goal

- ▶ obtain a predictor that performs well on test data

$$\operatorname{argmin}_g \operatorname{Err}_T(g) = \mathbb{E}_T[\underbrace{\ell(Y, g(X))}_{\text{no access}}]$$

- ▶ Q. How to estimate the target risk?

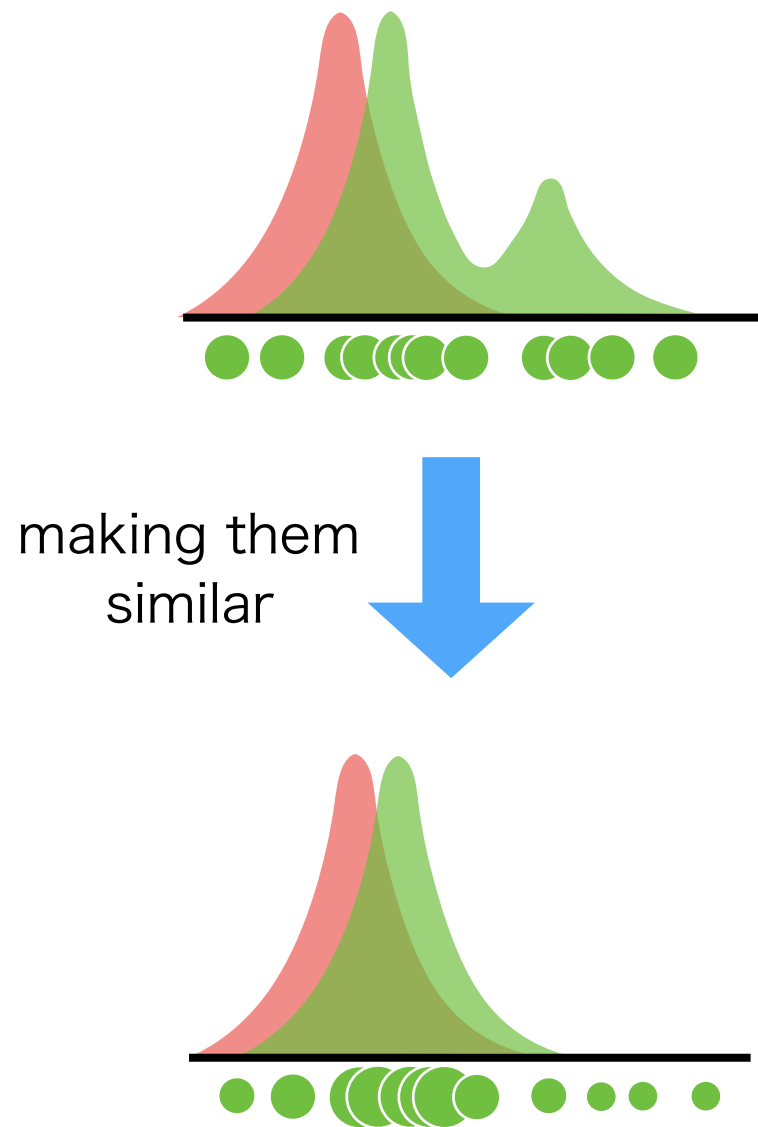
# Outline

- Introduction — Transfer Learning
- **History/Comparison of Existing Approaches**
- Proposed Method
- Experiments and Future Work

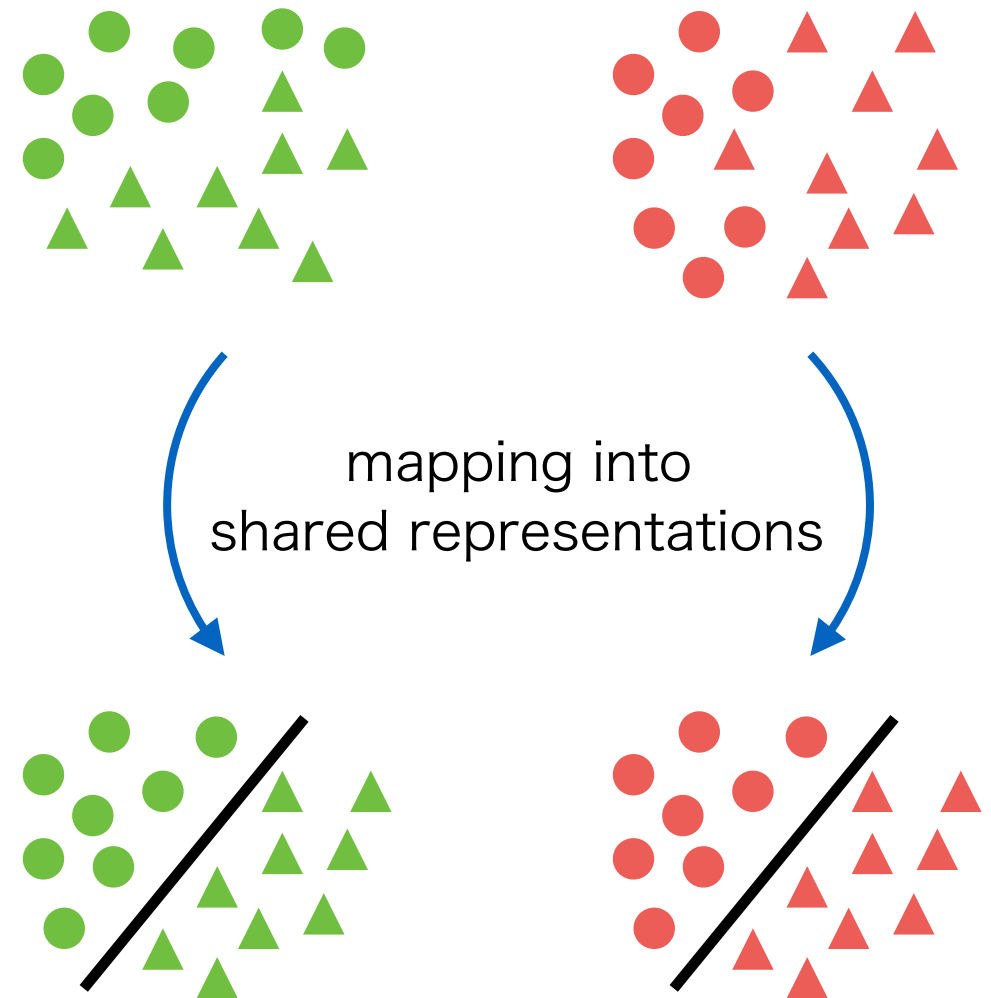
# Potential Solutions

8

## ■ Importance Weighting



## ■ Representation Learning

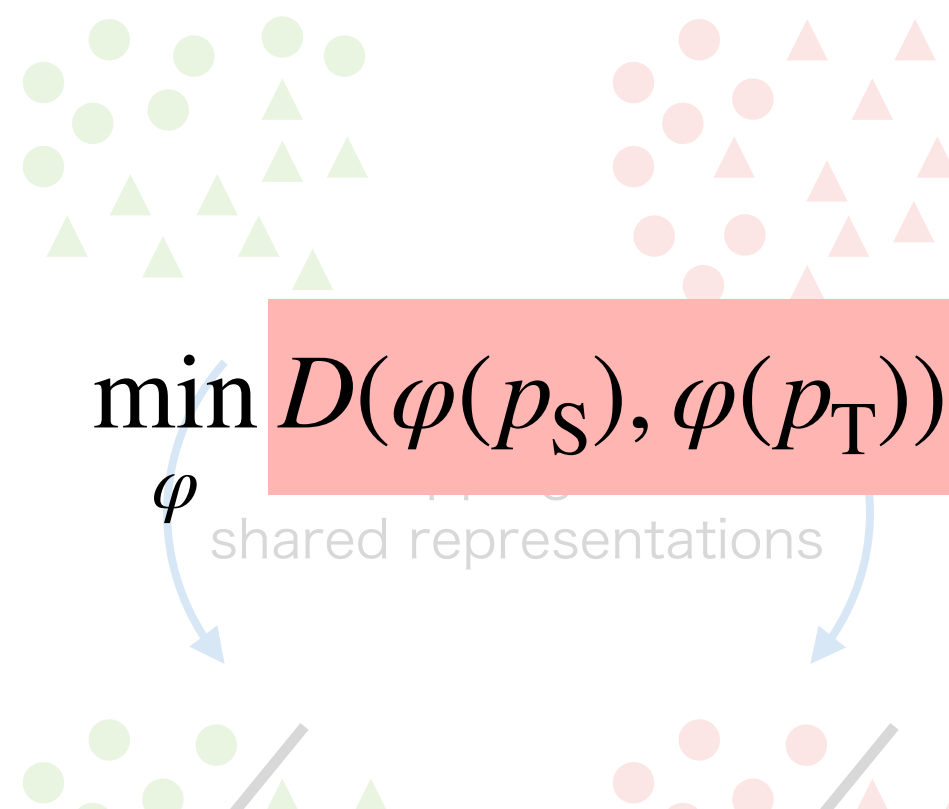
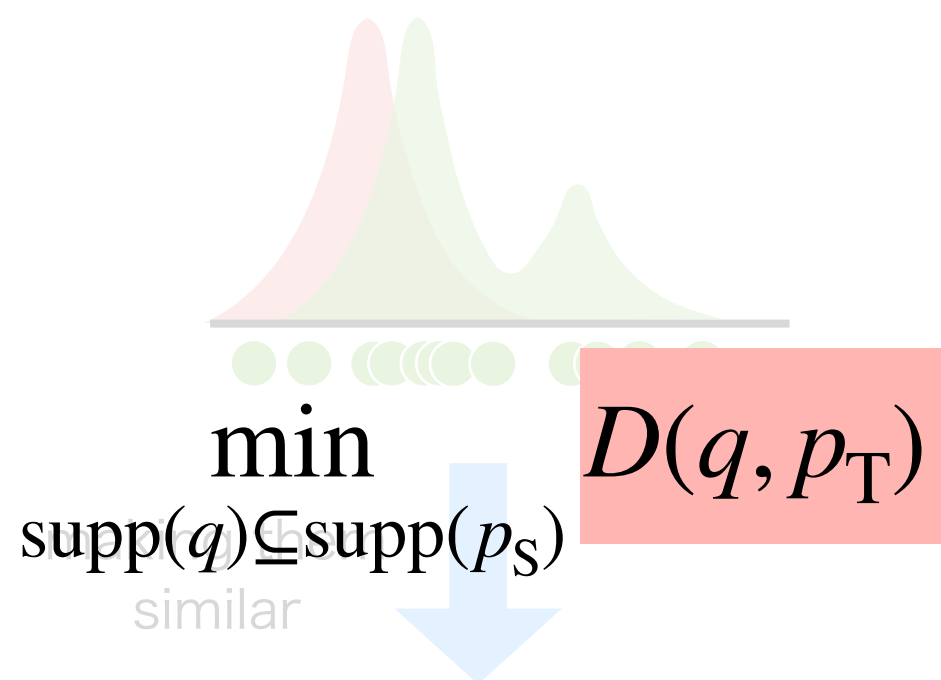




# Potential Solutions

■ Importance Weighting

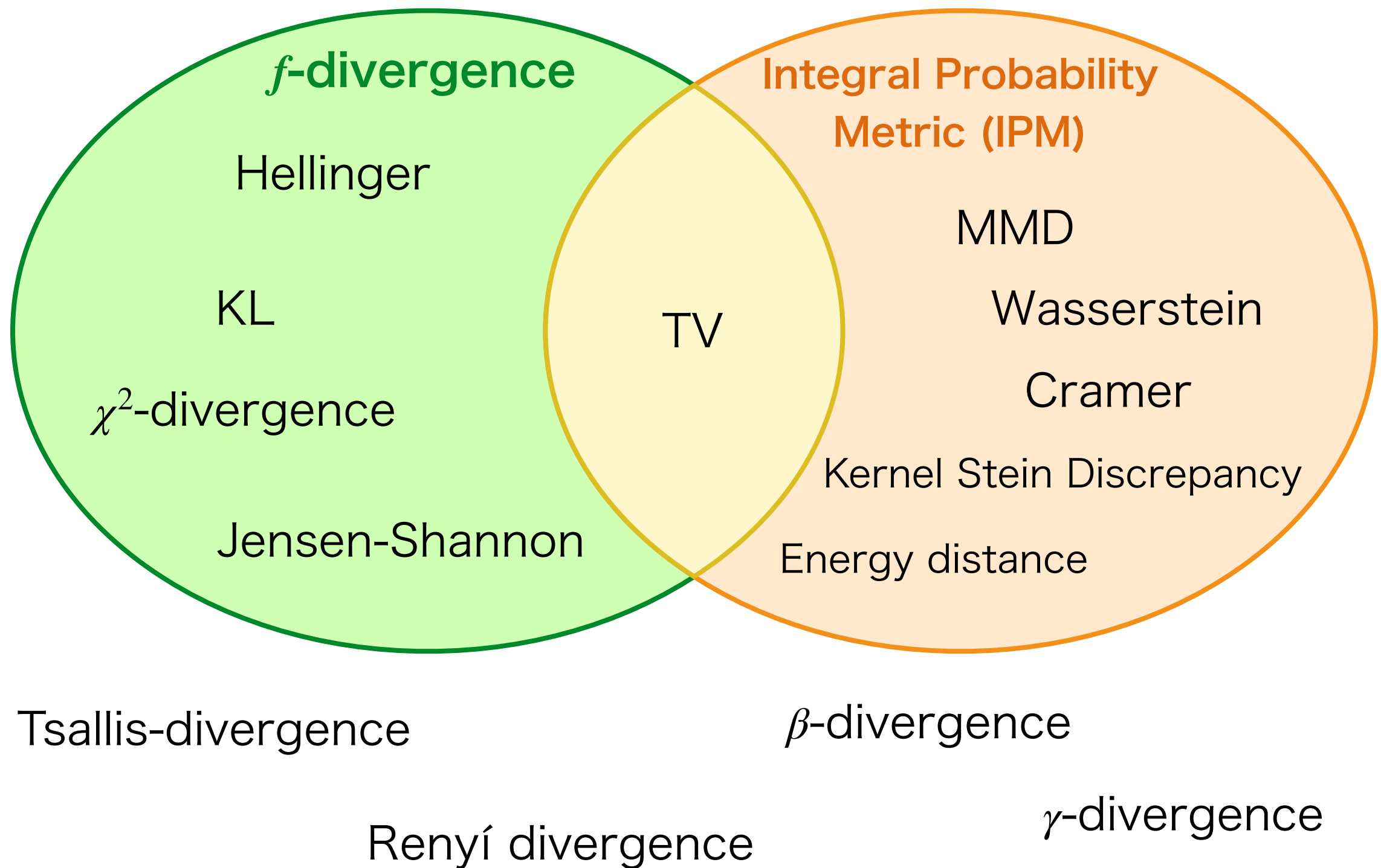
■ Representation Learning



It's important to measure closeness of distributions!

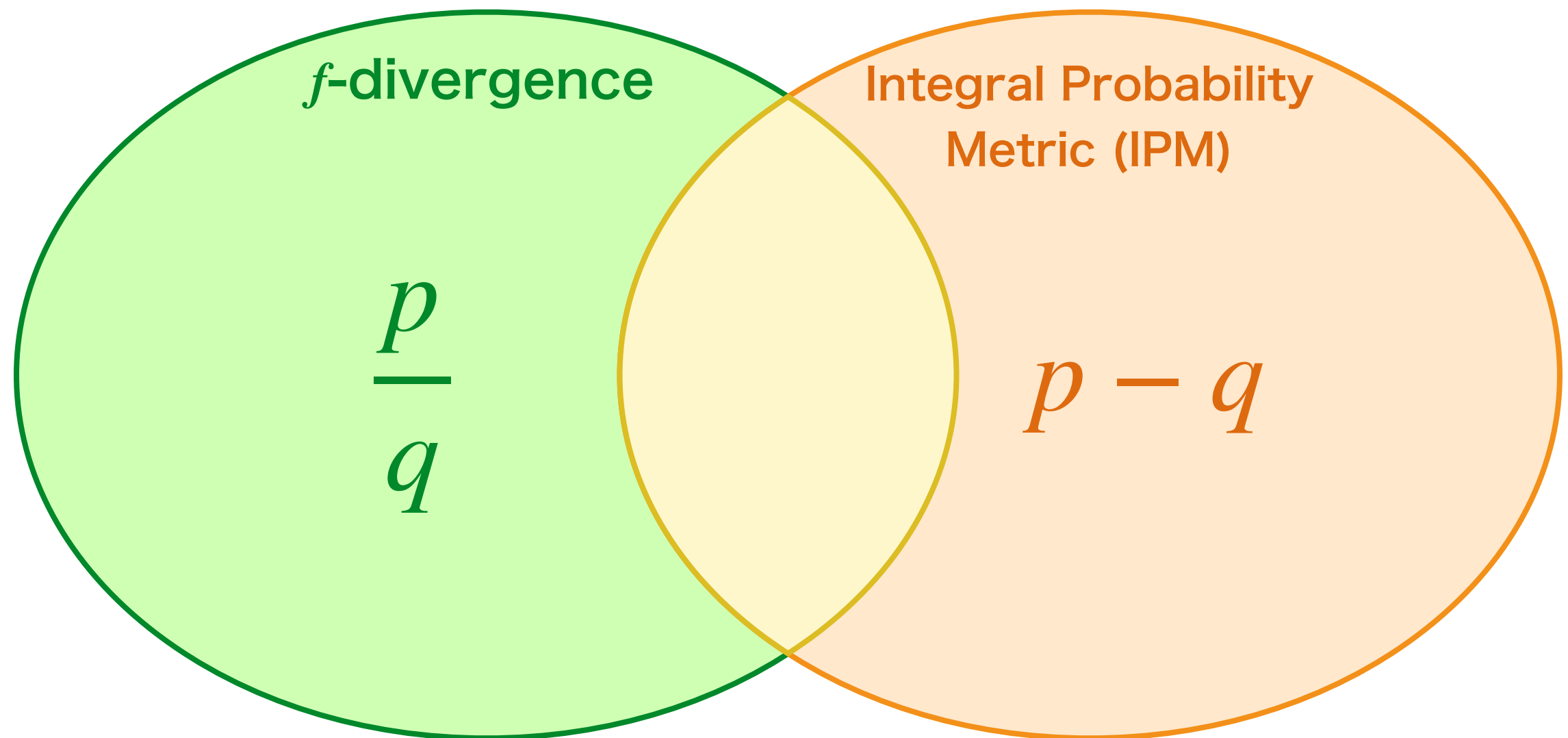
# Divergences

10



# Divergences

11



# What is a good measure?

- Postulate: classification risks should be closer if distances between distributions are small

$$\text{Err}_T(g) - \text{Err}_S(g) \leq D(p_T, p_S) + C$$

$$\mathbb{E}_T[\ell(g)] - \mathbb{E}_S[\ell(g)]$$

- IPM could be a more suitable family!

► IPM:  $D_\Gamma(p, q) = \sup_{\gamma \in \Gamma} \left| \mathbb{E}_p[\gamma] - \mathbb{E}_q[\gamma] \right|$

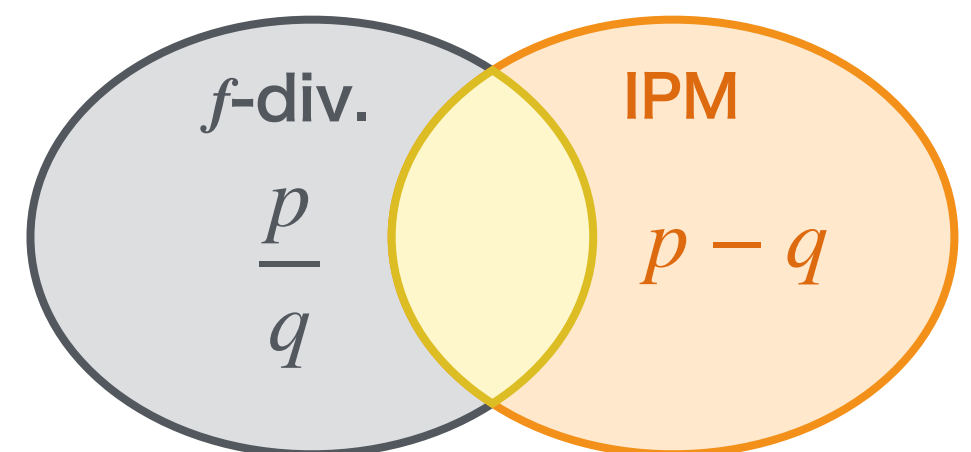
$\Gamma$  : real-valued function class  
(e.g. 1-Lipschitz for Wasserstein)

- represented in **difference of expectations**

expectation over marginal  
of source dist.

$$\text{Err}_S[g] = \underbrace{\mathbb{E}_{p_S}}_{\text{loss func.}} [\underbrace{\ell(g(X), f_S(X))}_{\text{labeling func.}}]$$

(parallel notation for target domain as well)



# Simple Approach: Total Variation <sup>13</sup>

[Kifer+ VLDB2004]

■ Total Variation  $D_{TV}(p, q) = 2 \sup_{A: \text{mes'ble}} |p(A) - q(A)|$

$p, q$  are distributions  
over  $\mathcal{X}$

■ classification risk bound

$$\text{Err}_T(g) - \text{Err}_S(g) \leq$$

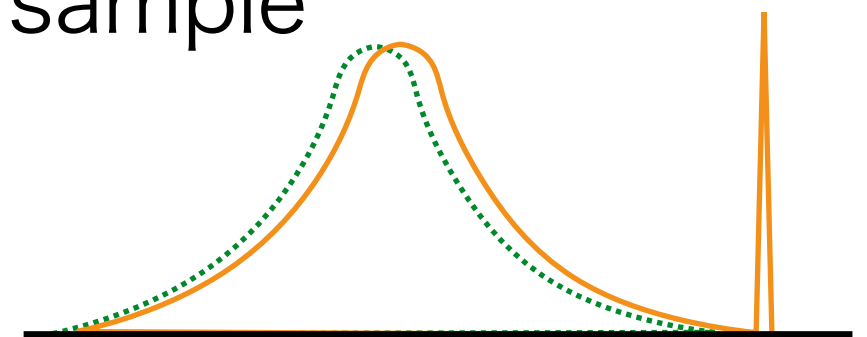
$$D_{TV}(p_S, p_T) + \min\{\mathbb{E}_S[|f_S - f_T|], \mathbb{E}_T[|f_S - f_T|]\}$$

■ Problems

► TV is overly pessimistic

► TV is hard estimate within finite sample

we can make a distribution  
with arbitrarily large TV



# First Attempt: $\mathcal{H} \Delta \mathcal{H}$ -divergence <sup>14</sup>

[Kifer+ VLDB2004; Blitzer+ NeurIPS2008]

## Definition ( $\mathcal{H}$ -divergence)

$$D_{\mathcal{H}}(p, q) = 2 \sup_{g \in \mathcal{H}} \left| p(g(X) = 1) - q(g(X) = 1) \right| \quad ; \mathcal{H} \subset \{\pm 1\}^{\mathcal{X}}$$

- ▶  $D_{\mathcal{H}}(p, q) \leq D_{\text{TV}}(p, q)$  by def.  $\Rightarrow$  could be less pessimistic
- ▶ estimator  $\hat{D}_{\mathcal{H}}(p, q)$  can be computed by ERM in  $\mathcal{H}$  (omitted)

## Lemma (finite-sample convergence)

Let  $d = \text{VCdim}(\mathcal{H})$ . Then, with prob. at least  $1 - \delta$ ,

$$D_{\mathcal{H}}(p_S, p_T) \leq \underbrace{\hat{D}_{\mathcal{H}}(p_S, p_T)}_{\text{empirical estimator}} + \tilde{O}_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right)$$

empirical estimator

$$\hat{D}_{\mathcal{H}}(p_S, p_T) = 2 \sup \left| \frac{1}{n_S} \sum_{x \in S} \mathbf{1}_{\{g(x)=1\}} - \frac{1}{n_T} \sum_{x \in T} \mathbf{1}_{\{g(x)=1\}} \right|$$

Kifer, D., Ben-David, S., & Gehrke, J. (2004, August). Detecting change in data streams. In *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30* (pp. 180-191). VLDB Endowment.

Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., & Wortman, J. (2008). Learning bounds for domain adaptation. In *Advances in neural information processing systems* (pp. 129-136).

# First Attempt: $\mathcal{H} \Delta \mathcal{H}$ -divergence <sup>15</sup>

[Kifer+ VLDB2004; Blitzer+ NeurIPS2008]

## Definition (symmetric difference hypothesis $\mathcal{H} \Delta \mathcal{H}$ )

$$g \in \mathcal{H} \Delta \mathcal{H} \iff g = h \oplus h' \text{ for some } h, h' \in \mathcal{H} \quad (\oplus : \text{XOR})$$

## Theorem (domain adaptation bound)

Let  $d = \text{VCdim}(\mathcal{H})$ . Then, with prob. at least  $1 - \delta$ , for any  $g$ ,

$$\text{Err}_T(g) \leq \text{Err}_S(g) + \frac{1}{2} \hat{D}_{\mathcal{H} \Delta \mathcal{H}}(p_S, p_T) + \tilde{O}_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right) + \lambda$$

where  $\lambda = \min_{h \in \mathcal{H}} \text{Err}_S(h) + \text{Err}_T(h)$  (joint minimizer)

## Issues

- ▶  $\hat{D}_{\mathcal{H} \Delta \mathcal{H}}$  is **intractable**; though  $\hat{D}_{\mathcal{H}}$  is tractable
- ▶  $\lambda$  is intrinsically **impossible to estimate**; assume to be small

( $\because \text{Err}_T$  cannot be accessed)

# Extension: discrepancy measure <sup>16</sup>

[Mansour+ COLT2009]

## Definition (discrepancy)

$$D_{\text{disc},\ell}(p, q) = \sup_{g, g' \in \mathcal{H}} \left| \text{Err}_p(g, g') - \text{Err}_q(g, g') \right| ; \text{Err}(g, g') = \underbrace{\int \ell(g(X), g'(X)) dp}_{\text{loss is generalized}}$$

- ▶ intuition: seeking for potential labelings maximizing diff. of losses
- ▶  $\hat{D}_{\text{disc},\ell}$  : empirical estimator of  $D_{\text{disc},\ell}$ ;  $\hat{D}_{\text{disc},\ell}(p, q) = \sup_{g, g' \in \mathcal{H}} \left| \widehat{\text{Err}}_p(g, g') - \widehat{\text{Err}}_q(g, g') \right|$

## Lemma (finite-sample convergence)

Let Rademacher averages of  $\mathcal{H}$  on the distribution  $p_S$  ( $p_T$  resp.) are bounded by  $O_p(n_S^{-1/2})$  ( $O_p(n_T^{-1/2})$  resp.). Assume  $\ell$  is Lipschitz cont. Then, with prob. at least  $1 - \delta$ ,

$$D_{\text{disc},\ell}(p_S, p_T) \leq \hat{D}_{\text{disc},\ell}(p_S, p_T) + O_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right)$$



# Extension: discrepancy measure <sup>17</sup>

[Mansour+ COLT2009]

## Theorem (domain adaptation bound)

Let Rademacher averages of  $\mathcal{H}$  on the distribution  $p_S$  ( $p_T$  resp.) are bounded by  $O_p(n_S^{-1/2})$  ( $O_p(n_T^{-1/2})$  resp.). Assume  $\mathcal{H}$  is symmetric. Then, with prob. at least  $1 - \delta$ , for any  $g$ ,

$$\text{Err}_T(g, f_T) - \underbrace{\text{Err}_T^*}_{= \text{Err}_T(g_S^*, f_T)} \leq \widehat{\text{Err}}_S(g, g_S^*) + \hat{D}_{\text{disc},01}(p_S, p_T) + O_p\left(\frac{1}{\sqrt{\min\{n_S, n_T\}}}\right) + \lambda$$

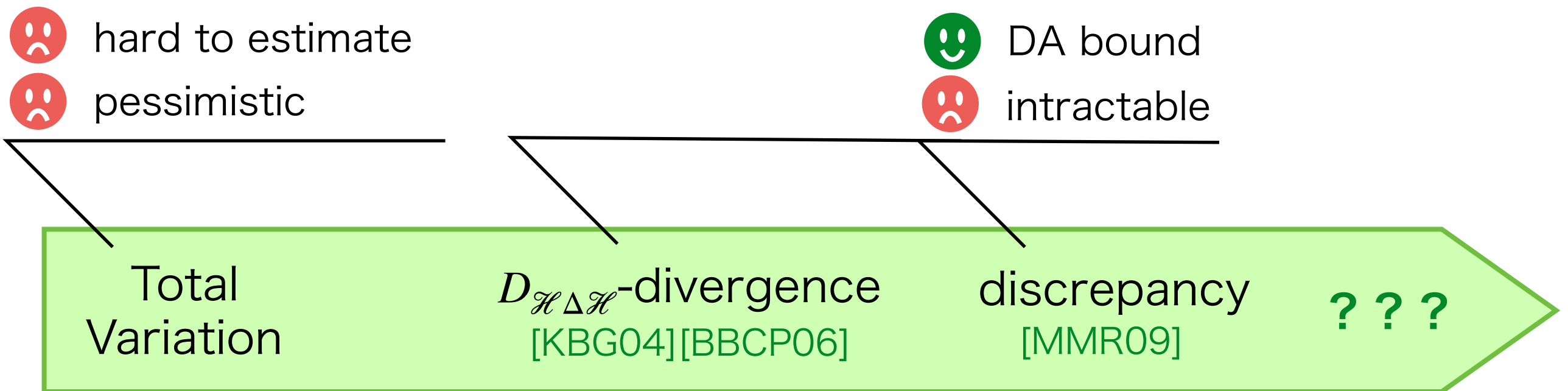
where  $\lambda = \text{Err}_T(g_S^*, g_T^*)$  (joint minimizer)

## Issues

- ▶  $\hat{D}_{\text{disc},\ell}$  is generally **intractable**; needs joint sup of  $g$  and  $g'$   
(tractable in simple cases)
- ▶  $\lambda$  is intrinsically **impossible to estimate**; assume to be small

# Comparison of Existing Measures<sup>18</sup>

Q. Can we construct a tractable/tighter measure?



# Outline


- Introduction — Transfer Learning
- History/Comparison of Existing Approaches
- **Proposed Method**
- Experiments and Future Work

# Proposed: Source-guided Discrepancy <sup>20</sup>

**Idea:** supremum with one variable should be tractable

## Definition (Source-guided Discrepancy)

$$D_{\text{sd},\ell}(p, q) = \sup_{g \in \mathcal{H}} \left| \text{Err}_p(g, g_S^*) - \text{Err}_q(g, g_S^*) \right| ; \text{Err}(g, g') = \int \ell(g(X), g'(X)) dp$$

 fix one function

where  $g_S^* = \operatorname{argmin}_{g \in \mathcal{H}} \text{Err}_S(g)$  (source risk minimizer)

cf. (discrepancy)

$$D_{\text{disc},\ell}(p, q) = \sup_{g, g' \in \mathcal{H}} \left| \text{Err}_p(g, g') - \text{Err}_q(g, g') \right|$$

►  $D_{\text{sd},\ell}(p, q) \leq D_{\text{disc},\ell}(p, q)$  by definition (S-disc is finer)

# S-disc Estimator = ERM

- Consider binary classification (loss function:  $\ell_{01}$ )
  - ▶ assume  $\mathcal{H}$  is symmetric:  $g \in \mathcal{H} \implies -g \in \mathcal{H}$

**Theorem**  $\hat{D}_{\text{sd},01}(p_S, p_T) = 1 - \min_{g \in \mathcal{H}} J_{\ell_{01}}(g)$

where  $J_{\ell}(g) = \frac{1}{n_S} \sum_{i=1}^{n_S} \ell(g(x_i^S), \underbrace{g_S^*(x_i^S)}_{\text{source: labeled by } g_S^*}) + \frac{1}{n_T} \sum_{j=1}^{n_T} \ell(g(x_j^T), \underbrace{-g_S^*(x_j^T)}_{\text{target: labeled by } -g_S^*})$  (cost-sensitive risk)

- Estimation Algorithm
  - ▶ train a classifier only using source ( $g_S^*$ )
  - ▶ minimize cost-sensitive risk  $J_{\ell}$

Similar idea to  $\mathcal{H}$ -divergence, but we don't need to use  $\mathcal{H} \Delta \mathcal{H}$

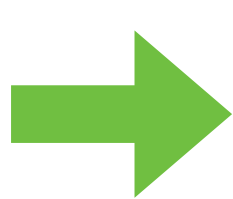
# Finite-Sample Consistency

22

## Theorem

Let Rademacher averages of  $\mathcal{H} \otimes \mathcal{H}$  on the distribution  $p_S$  ( $p_T$  resp.) are bounded by  $O_p(n_S^{-1/2})$  ( $O_p(n_T^{-1/2})$  resp.). Then, with prob. at least  $1 - \delta$ ,

$$D_{\text{sd},\ell}(p_S, p_T) \leq \hat{D}_{\text{sd},\ell}(p_S, p_T) + O_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right)$$



empirical S-disc is tractable



consistent (as well as  $D_{\mathcal{H}}$ ,  $D_{\text{disc}}$ )

- ▶  $\mathcal{H} \otimes \mathcal{H} = \{g \cdot g' \mid g, g' \in \mathcal{H}\}$
- ▶  $\text{Rad}(\mathcal{H}) = O_p(n^{-1/2}) \implies \text{Rad}(\mathcal{H} \otimes \mathcal{H}) = O_p(n^{-1/2})$

# Domain Adaptation Bound

## Theorem (domain adaptation bound)

Let Rademacher averages of  $\mathcal{H} \otimes \mathcal{H}$  on the distribution  $p_S$  ( $p_T$  resp.) are bounded by  $O_p(n_S^{-1/2})$  ( $O_p(n_T^{-1/2})$  resp.).

Assume the loss  $\ell$  satisfies the triangle inequality.

Then, with prob. at least  $1 - \delta$ , for any  $g$ ,

$$\text{Err}_T(g, f_T) - \text{Err}_T^* \leq \widehat{\text{Err}}_S(g, g_S^*) + \hat{D}_{\text{sd}, \ell}(p_S, p_T) + O_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right) + \lambda$$

where  $\lambda = \text{Err}_T(g_S^*, g_T^*)$  (joint minimizer)

😊  $\hat{D}_{\text{sd}, \ell}$  is tractable

😊  $D_{\text{sd}, \ell} \leq D_{\text{disc}, \ell}$  (always tighter bound)

😞  $\lambda$  is impossible to estimate

# Summary

## Source-guided Discrepancy

$$D_{\text{sd},\ell}(p, q) = \sup_{g \in \mathcal{H}} \left| \text{Err}_p(g, g_S^*) - \text{Err}_q(g, g_S^*) \right|$$

fix one function

## DA bound

$$\text{Err}_T(g, f_T) - \text{Err}_T^* \leq \widehat{\text{Err}}_S(g, g_S^*) + \hat{D}_{\text{sd},\ell}(p_S, p_T) + O_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right) + \lambda$$

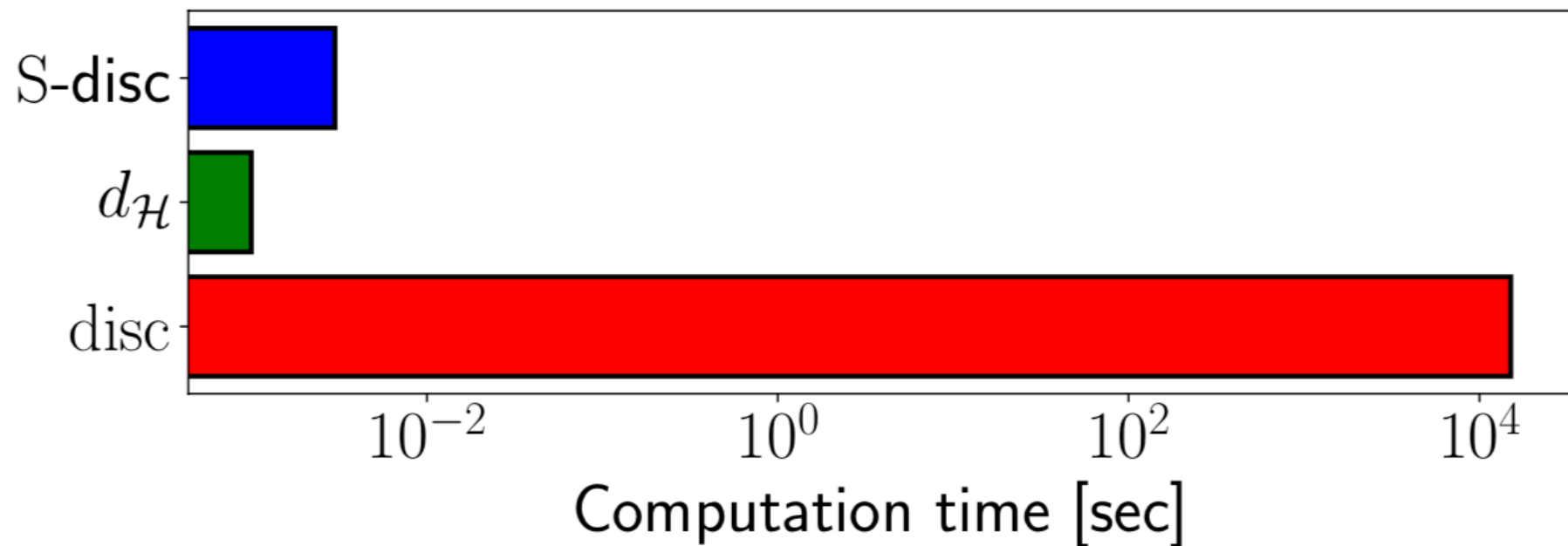
- 😊 Tractable estimator: can be computed by ERM
- 😊 Tighter measure
- 😞 DA bound, but still  $\lambda$  (impossible term) exists



# Outline

- Introduction — Transfer Learning
- History/Comparison of Existing Approaches
- Proposed Method
- **Experiments and Future Work**

# Computational Time



- $d = 2$ , 200 synthetic examples for both source and target
- $d_{\mathcal{H}}$  is an approximator of  $D_{\mathcal{H} \Delta \mathcal{H}}$ 
  - faster, but does not entail DA bound
- discrepancy is computed via approximation
  - resorted to semi-definite relaxation

# Source Selection

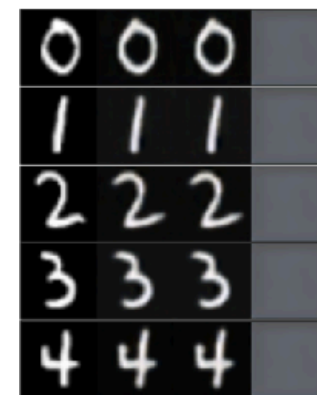
## ■ Domains

- ▶ source: 5 clean MNIST-M, 5 noisy MNIST-M
- ▶ target: MNIST  
(clean MNIST-M is known to be useful for MNIST)

## ■ Setup

- ▶ measure the distance between target and each source
- ▶ sort in ascending order

➡ 5 clean MNIST-M should admit smaller distance than noisy ones

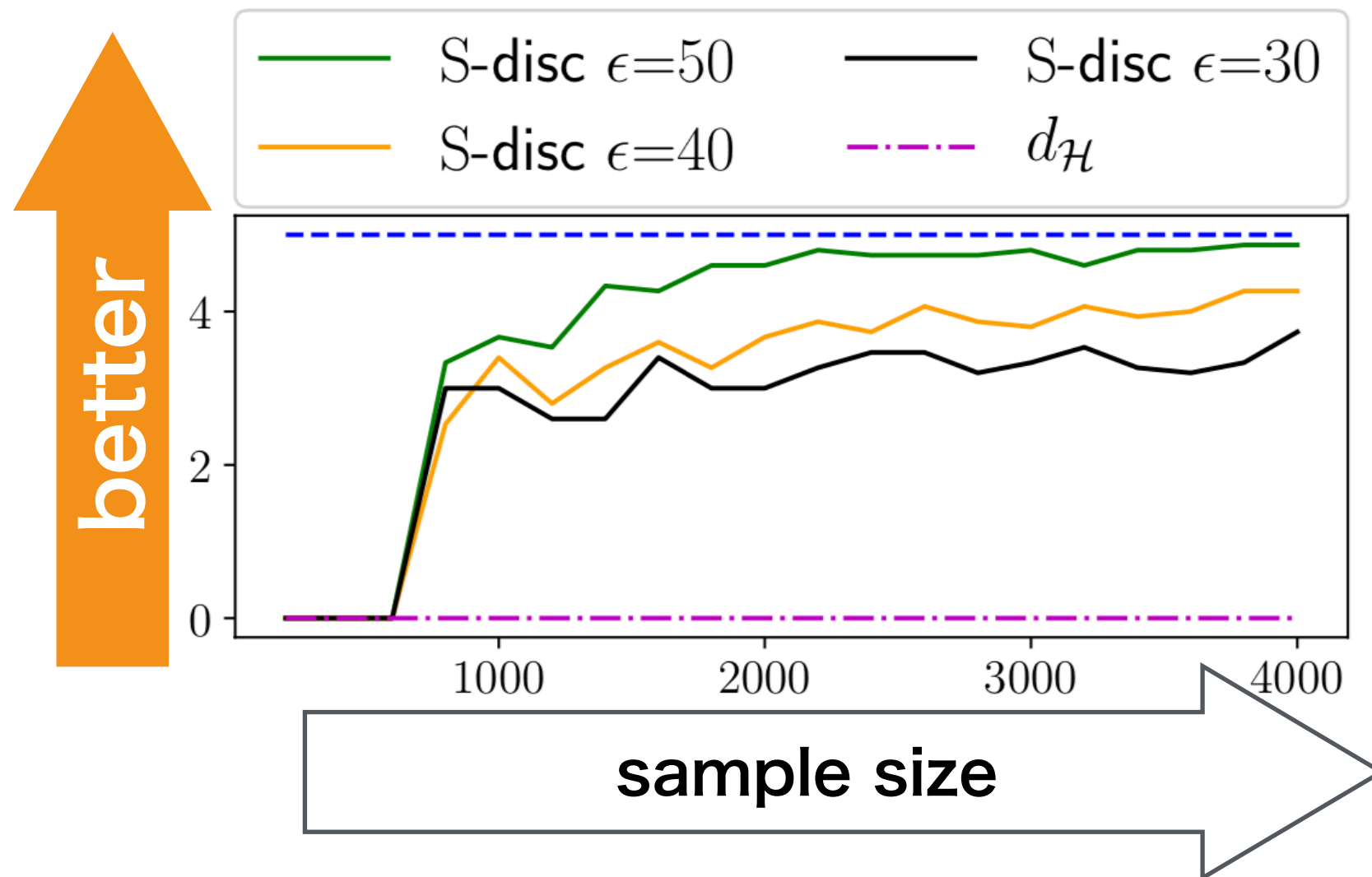


MNIST



MNIST-M

# Source Selection



Vertical-axis: # of clean MNIST-M domains in top 5

- S-disc successfully capture the difference between clean and noisy MNIST-M

# Following Work

[Zhang+ ICML2019]

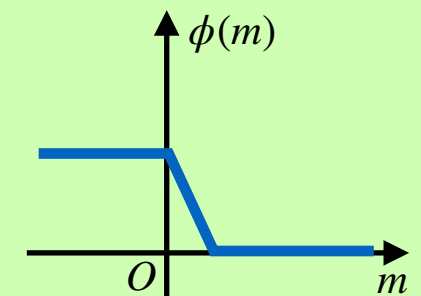
## Source-guided Discrepancy

$$D_{sd,\ell}(p, q) = \sup_{g \in \mathcal{H}} \left| \text{Err}_p(g, \underbrace{g_S^*}_{\text{fix source-risk minimizer}}) - \text{Err}_q(g, \underbrace{g_S^*}_{\text{fix source-risk minimizer}}) \right|$$

## Definition: Margin Disparity Discrepancy

② limited to margin loss

$$D_{MDD,f,\ell}(p, q) = \sup_{g \in \mathcal{H}} \left| \text{Err}_p(g, \underbrace{f}_{\text{① fix an arbitrary}}) - \text{Err}_q(g, \underbrace{f}_{\text{① fix an arbitrary}}) \right|$$



## DA bound based on MDD

$$\text{Err}_T(g, f_T) \leq \widehat{\text{Err}}_S(g, g_S^*) + \hat{D}_{MDD,g,\ell}(p_S, p_T) + O_p \left( \frac{1}{\sqrt{\min\{n_S, n_T\}}} \right) + \lambda$$

⇒ extended to **multi-class (one-vs-all)** case

# Conclusion

## ■ Discrepancy measure is important in domain adaptation

- ▶ IPM is a nice family; can be connected to DA bound
- ▶ “**fixing one function**” would be a good idea

$$D_{\text{sd},\ell}(p, q) = \sup_{g \in \mathcal{H}} \left| \text{Err}_p(g, \underbrace{g_S^*}_{\text{fix source-risk minimizer}}) - \text{Err}_q(g, \underbrace{g_S^*}_{\text{fix source-risk minimizer}}) \right|$$

## ■ Potential directions

- ▶ remove the unestimable term in DA bound ( $\lambda$ )
- ▶ any “optimality” in DA bound?

rethinking DA framework (adaptation algorithms, available supervision)  
might be needed...