Transfer Learning: Theories and Algorithms

Mingsheng Long

School of Software, Tsinghua University Research Center for Big Data, Tsinghua University National Engineering Laboratory for Big Data Software

mingsheng@tsinghua.edu.cn http://ise.thss.tsinghua.edu.cn/~mlong

(日) (同) (三) (三)

Outline

Transfer Learning

2 $\mathcal{H}\Delta\mathcal{H}$ -Divergence

- DAN: Deep Adaptation Network
- DANN: Domain Adversarial Neural Network
- MCD: Maximum Classifier Discrepancy
- 3 Margin Disparity Discrepancy
 - MDD: Margin Disparity Discrepancy
- 4) Transfer Model Selection
 - DEV: Deep Embedded Validation
- 5 Evaluation and Implementation

A (10) F (10)

Supervised Learning

Learner: $f: \mathbf{x} \to y$ Distribution: $(\mathbf{x}, y) \sim P(\mathbf{x}, y)$



complexity

Error Bound: $\epsilon_{\text{test}} \leq \hat{\epsilon}_{\text{train}}$

Transfer Learning

- Machine learning across domains of different distributions P \neq Q
 Independent and Differently Distributed (IDD)
- How to effectively bound the generalization error on target domain?



Bias-Variance-Shift Tradeoff



Bridging Theory and Algorithm





< □ > < 同 > < 三 > < 三

Everything should be made as simple as possible, but no simpler. —Albert Einstein There is nothing more practical than a good theory. —Vladimir Vapnik

Bridging Theory and Algorithm



Outline

Transfer Learning

$\mathcal{H}\Delta\mathcal{H}$ -Divergence

- DAN: Deep Adaptation Network
- DANN: Domain Adversarial Neural Network
- MCD: Maximum Classifier Discrepancy

Margin Disparity Discrepancy MDD: Margin Disparity Discrepancy

- Transfer Model Selection
 - DEV: Deep Embedded Validation
- 5 Evaluation and Implementation

Notations and Assumptions

- Source risk: $\epsilon_P(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P}[h(\mathbf{x}) \neq y], \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim P^n$
- Target risk: $\epsilon_Q(h) = \mathbb{E}_{(\mathbf{x}, y) \sim Q}[h(\mathbf{x}) \neq y], \{(\mathbf{x}_i, y_i)\}_{i=1}^m \sim Q^m$
- Source disparity: $\epsilon_P(h_1, h_2) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P}[h_1(\mathbf{x}) \neq h_2(\mathbf{x})]$
- Target disparity: $\epsilon_Q(h_1, h_2) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim Q}[h_1(\mathbf{x}) \neq h_2(\mathbf{x})]$
- Ideal joint hypothesis: $h^* = \arg \min_h \epsilon_P(h) + \epsilon_Q(h)$
- Assumption: ideal hypothesis has small risk $\epsilon_{ideal} = \epsilon_P(h^*) + \epsilon_Q(h^*)$





Ideal hypothesis with small error

イロト 不得下 イヨト イヨト 二日

Relating the Target Risk to the Source Risk

Theorem

Assuming small ϵ_{ideal} , the bound of the target risk $\epsilon_Q(h)$ of hypothesis $h \in \mathcal{H}$ is given by the source risk $\epsilon_P(h)$ plus the disparity difference:

 $\epsilon_{Q}(h) \leqslant \epsilon_{P}(h) + \left[\epsilon_{P}(h^{*}) + \epsilon_{Q}(h^{*})\right] + \left|\epsilon_{P}(h,h^{*}) - \epsilon_{Q}(h,h^{*})\right| \qquad (1)$

Proof.

Simply by using the triangle inequalities, we have

$$\epsilon_{Q}(h) \leq \epsilon_{Q}(h^{*}) + \epsilon_{Q}(h, h^{*})$$

$$\leq \epsilon_{Q}(h^{*}) + \epsilon_{P}(h, h^{*}) + \epsilon_{Q}(h, h^{*}) - \epsilon_{P}(h, h^{*})$$

$$\leq \epsilon_{Q}(h^{*}) + \epsilon_{P}(h, h^{*}) + |\epsilon_{Q}(h, h^{*}) - \epsilon_{P}(h, h^{*})|$$

$$\leq \epsilon_{P}(h) + [\epsilon_{P}(h^{*}) + \epsilon_{Q}(h^{*})] + |\epsilon_{P}(h, h^{*}) - \epsilon_{Q}(h, h^{*})|$$
(2)

How to Bound the Disparity Difference?

• We can illustrate the disparity difference $|\epsilon_{P}(h, h^{*}) - \epsilon_{Q}(h, h^{*})|$ as



- $\mathcal{H}\Delta\mathcal{H}$ -Divergence¹: $d_{\mathcal{H}\Delta\mathcal{H}}(P,Q) \triangleq \sup_{h,h'\in\mathcal{H}} |\epsilon_P(h,h') \epsilon_Q(h,h')|$
- Hypothesis-independent discrepancy—depending on hypothesis space.

¹Ben-David et al. A Theory of Learning from Different Domains. Machine Learning, 2010, and

Generalization Bound with $\mathcal{H}\Delta\mathcal{H}$ -Divergence

Theorem (Generalization Bound)

Denote by d the VC-dimension of hypothesis space \mathcal{H} . For any hypothesis $h \in \mathcal{H}$,

$$\epsilon_{Q}(h) \leq \epsilon_{\hat{P}}(h) + \frac{d_{\mathcal{H}\Delta\mathcal{H}}(\hat{P}, \hat{Q})}{n} + \epsilon_{ideal} + O(\sqrt{\frac{d\log n}{n}} + \sqrt{\frac{d\log m}{m}})$$
(3)

- $\epsilon_P(h)$ depicts the performance of h on source domain.
- $d_{\mathcal{H}\Delta\mathcal{H}}$ bounds the generalization gap caused by domain shift.
- ϵ_{ideal} quantifies the inverse of "adaptability" between domains.
- The order of the complexity term is $O(\sqrt{d \log n/n} + \sqrt{d \log m/m})$.

(人間) トイヨト イヨト

Approximating $\mathcal{H}\Delta\mathcal{H}$ -Divergence by Statistical Distance

For binary hypothesis h, the $\mathcal{H}\Delta\mathcal{H}$ -Divergence can be bounded by

$$d_{\mathcal{H}\Delta\mathcal{H}}(P,Q) \triangleq \sup_{\substack{h,h'\in\mathcal{H}\\ h,h'\in\mathcal{H}}} |\epsilon_{P}(h,h') - \epsilon_{Q}(h,h')|$$

$$= \sup_{\substack{h,h'\in\mathcal{H}\\ \delta\in\mathcal{H}\Delta\mathcal{H}}} |\mathbb{E}_{P}[|h(\mathbf{x}) - h'(\mathbf{x})| \neq 0] - \mathbb{E}_{Q}[|h(\mathbf{x}) - h'(\mathbf{x})| \neq 0]| \quad (4)$$

The last term takes the form of Integral Probability Metric (IPM):

$$d_{\mathcal{F}}(P,Q) = \sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathbf{x} \sim P} f(\mathbf{x}) - \mathbb{E}_{\mathbf{x} \sim Q} f(\mathbf{x})|$$
(5)

Assuming \mathcal{F} can be approximated by kernel functions in RKHS, $d_{\mathcal{F}}(P, Q)$ turns into Maximum Mean Discrepancy (MMD) (a statistical distance)

(日) (同) (三) (三)

DAN: Deep Adaptation Network²



Distribution matching: yield the upper-bound by multiple kernel learning

$$d_{k}^{2}(P,Q) \triangleq \left\| \mathbf{E}_{P} \left[\phi \left(\mathbf{x}^{s} \right) \right] - \mathbf{E}_{Q} \left[\phi \left(\mathbf{x}^{t} \right) \right] \right\|_{\mathcal{H}_{k}}^{2}$$
(6)
$$\min_{\theta \in \Theta} \max_{k \in \mathcal{K}} \frac{1}{n_{a}} \sum_{i=1}^{n_{a}} L\left(\theta \left(\mathbf{x}^{a}_{i} \right), y^{a}_{i} \right) + \lambda \sum_{\ell=l_{1}}^{l_{2}} d_{k}^{2} \left(\widehat{P}_{\ell}, \widehat{Q}_{\ell} \right)$$
(7)

²Long et al. Learning Transferable Features with Deep Adaptation Networks. *ICML* 2015. a.e.

Mingsheng Long

Approximating $\mathcal{H}\Delta\mathcal{H}$ -Divergence by Domain Discriminator

For binary hypothesis h, the $\mathcal{H}\Delta\mathcal{H}$ -Divergence can be bounded by

$$d_{\mathcal{H}\Delta\mathcal{H}}(P,Q) \triangleq \sup_{\substack{h,h'\in\mathcal{H}\\\delta\in\mathcal{H}\Delta\mathcal{H}}} |\epsilon_P(h,h') - \epsilon_Q(h,h')|$$

$$= \sup_{\delta\in\mathcal{H}\Delta\mathcal{H}} |\mathbb{E}_P[\delta(\mathbf{x})\neq 0] - \mathbb{E}_Q[\delta(\mathbf{x})\neq 0]|$$

$$\leq \sup_{\substack{D\in\mathcal{H}_D}} |\mathbb{E}_P[D(\mathbf{x})=1] + \mathbb{E}_Q[D(\mathbf{x})=0]|$$
(8)

This upper-bound can be yielded by training a domain discriminator $D(\mathbf{x})$



DANN: Domain Adversarial Neural Network³



Adversarial adaptation: learning features indistinguishable across domains

$$E(\theta_{f},\theta_{y},\theta_{d}) = \sum_{\mathbf{x}_{i}\sim\widehat{P}} L_{y}(G_{y}(G_{f}(\mathbf{x}_{i})),y_{i}) - \lambda \sum_{\mathbf{x}_{i}\sim\widehat{P}\cup\widehat{Q}} L_{d}(G_{d}(G_{f}(\mathbf{x}_{i})),d_{i})$$
(9)

$$(\hat{\theta}_{f}, \hat{\theta}_{y}) = \arg\min_{\theta_{f}, \theta_{y}} E\left(\theta_{f}, \theta_{y}, \theta_{d}\right) \quad (\hat{\theta}_{d}) = \arg\max_{\theta_{d}} E\left(\theta_{f}, \theta_{y}, \theta_{d}\right) \tag{10}$$

³Ganin et al. Domain Adversarial Training of Neural Networks. JMLR 2016. 🧃 👘 🚊 🔊 🔍

Approximating $\mathcal{H}\Delta\mathcal{H}$ -Divergence by Classifier Consistency⁴



- Use two classifiers G_1, G_2 to approximate $\sup_{h,h' \in \mathcal{H}} |\epsilon_P(h, h') \epsilon_Q(h, h')|$
- Assume $G_1 = h$ and $G_2 = h'$ should agree on source domain.
- Use L₁-loss of two classifiers' outputs to approximate disagreement:

$$\min_{\phi} \{ \min_{G_1, G_2} \mathbb{E}_{\widehat{P}}[L(G_1(\mathbf{x}), y) + L(G_2(\mathbf{x}), y)] + \max_{G_1, G_2} \mathbb{E}_{\widehat{Q}}|G_1(\mathbf{x}) - G_2(\mathbf{x})| \}$$
(11)

Outline

Transfer Learning

2 $\mathcal{H}\Delta\mathcal{H}$ -Divergence

- DAN: Deep Adaptation Network
- DANN: Domain Adversarial Neural Network
- MCD: Maximum Classifier Discrepancy

Margin Disparity Discrepancy MDD: Margin Disparity Discrepancy

Transfer Model Selection

- DEV: Deep Embedded Validation
- 5 Evaluation and Implementation

< 67 ▶

Towards Informative Margin Theory⁵

- Towards a rigorous multiclass domain adaptation theory.
 - All existing theories are only applicable to binary classification.
- Towards an informative margin theory.
 - Explore the idea of margin in measuring domain discrepancy.
- Towards a certain function class in the theoretical bound.
 - Eliminate approximation assumptions in all existing methods.
- Towards bridging the existing gap between theories and algorithms.

 5 Zhang et al. Bridging Theory and Algorithm for Domain Adaptation. ICML 2019. 🚊 🗠 🔍

Notations

- Scoring function: $f \in \mathcal{F} : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$
- Labeling function induced by $f: h_f : x \mapsto \arg \max_{y \in \mathcal{Y}} f(x, y)$
- Labeling function class: $\mathcal{H} = \{h_f | f \in \mathcal{F}\}$
- Margin of a hypothesis:

$$\rho_f(x, y) = \frac{1}{2} (f(x, y) - \max_{y' \neq y} f(x, y'))$$

• Margin Loss:

$$\Phi_{\rho}(x) = \begin{cases} 0 & \rho \leq x \\ 1 - x/\rho & 0 \leq x \leq \rho \\ 1 & x \leq 0 \end{cases}$$



(日) (周) (三) (三)

3

MDD: Margin Disparity Discrepancy

- Margin risk: $\epsilon_D^{(\rho)}(f) = \mathbb{E}_{(x,y)\sim D} \left[\Phi_{\rho}(\rho_f(x,y)) \right]$
- Margin disparity: $\epsilon_D^{(\rho)}(f', f) \triangleq \mathbb{E}_{x \sim D_X}[\Phi_{\rho}(\rho_{f'}(x, h_f(x)))]$

Definition (Margin Disparity Discrepancy, MDD)

With above definitions, we define Margin Disparity Discrepancy (MDD) and its empirical version by

$$d_{f,\mathcal{F}}^{(\rho)}(P,Q) \triangleq \sup_{f' \in \mathcal{F}} \left(\epsilon_Q^{(\rho)}(f',f) - \epsilon_P^{(\rho)}(f',f) \right), \\ d_{f,\mathcal{F}}^{(\rho)}(\widehat{P},\widehat{Q}) \triangleq \sup_{f' \in \mathcal{F}} \left(\epsilon_{\widehat{Q}}^{(\rho)}(f',f) - \epsilon_{\widehat{P}}^{(\rho)}(f',f) \right).$$
(12)

MDD satisfies $d_{f,\mathcal{F}}^{(\rho)}(P,P) = 0$ as well as nonnegativity and subadditivity.

(日) (周) (三) (三)

Bounding the Target Risk by MDD

Theorem

Let $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X} \times \mathcal{Y}}$ be a hypothesis set with label set $\mathcal{Y} = \{1, \dots, k\}$ and $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be the corresponding \mathcal{Y} -valued labeling function class. For every scoring function $f \in \mathcal{F}$,

$$\epsilon_Q(f) \le \epsilon_P^{(\rho)}(f) + d_{f,\mathcal{F}}^{(\rho)}(P,Q) + \epsilon_{ideal}^{(\rho)}, \tag{13}$$

where $\epsilon_{ideal}^{(\rho)}$ is the margin error of ideal joint hypothesis f^* :

$$\epsilon_{ideal}^{(\rho)} = \min_{f^* \in \mathcal{F}} \{ \epsilon_P^{(\rho)}(f^*) + \epsilon_Q^{(\rho)}(f^*) \}.$$

$$(14)$$

• Main proof difficulties: margin loss does not satisfy triangle inequality.

• Solution: One-sided triangle inequality for the margin loss.

• A new tool for analyzing transfer learning with margin theory.

Mingsheng Long	Transfer Learning		August	21, 2019		22 / 37
				2 E 7	-	+) Q (+

Definitions

Definition (Function Class $\Pi_1 \mathcal{F}$)

Given a class of scoring functions \mathcal{F} , $\Pi_1 \mathcal{F}$ is defined as

$$\Pi_1 \mathcal{F} = \{ x \mapsto f(x, y) | y \in \mathcal{Y}, f \in \mathcal{F} \}.$$
(15)

Definition (Function Class $\Pi_{\mathcal{H}}\mathcal{F}$)

Given a class of scoring functions \mathcal{F} and a class of the induced labeling functions \mathcal{H} , we define $\Pi_{\mathcal{H}}\mathcal{F}$ as

$$\Pi_{\mathcal{H}}\mathcal{F} \triangleq \{ x \mapsto f(x, \mathbf{h}(x)) | h \in \mathcal{H}, f \in \mathcal{F} \}.$$
(16)

By applying the margin error over each entry in $\Pi_{\mathcal{H}}\mathcal{F}$, we obtain the "scoring" version of $\mathcal{H}\Delta\mathcal{H}$ (symmetric difference hypothesis space)

(日) (周) (三) (三)

Definitions

Definition (Rademacher Complexity)

The empirical Rademacher complexity of function class ${\cal G}$ with respect to the sample \widehat{D} is defined as

$$\widehat{\mathfrak{R}}_{\widehat{D}}(\mathcal{G}) = \mathbb{E}_{\sigma} \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_i g(z_i).$$
(17)

where σ_i 's are independent uniform random variables taking values in $\{-1, +1\}$. The Rademacher complexity is

$$\mathfrak{R}_{n,D}(\mathcal{G}) = \mathbb{E}_{\widehat{D} \sim D^n} \widehat{\mathfrak{R}}_{\widehat{D}}(\mathcal{G}).$$
(18)

Definition (**Covering Number**)

(Informal) A covering number $\mathcal{N}_2(\tau, \mathcal{G})$ is the minimal number of \mathcal{L}_2 balls of radius $\tau > 0$ needed to cover a class \mathcal{G} of bounded functions $g : \mathcal{X} \to \mathbb{R}$

Generalization Bound with Rademacher Complexity

Theorem (Generalization Bound with Rademacher Complexity)

Let $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X} \times \mathcal{Y}}$ be a hypothesis set with label set $\mathcal{Y} = \{1, \dots, k\}$ and $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be the corresponding \mathcal{Y} -valued labeling function class. Fix $\rho > 0$. For all $\delta > 0$, with probability $1 - 3\delta$ the following inequality holds for all hypothesis $f \in \mathcal{F}$:

$$\epsilon_{Q}(f) \leq \epsilon_{\widehat{P}}^{(\rho)}(f) + d_{f,\mathcal{F}}^{(\rho)}(\widehat{P},\widehat{Q}) + \epsilon_{ideal} + \frac{2k^{2}}{\rho} \mathfrak{R}_{n,P}(\Pi_{1}\mathcal{F}) + \frac{k}{\rho} \mathfrak{R}_{n,P}(\Pi_{\mathcal{H}}\mathcal{F}) + 2\sqrt{\frac{\log\frac{2}{\delta}}{2n}}$$
(19)
$$+ \frac{k}{\rho} \mathfrak{R}_{m,Q}(\Pi_{\mathcal{H}}\mathcal{F}) + \sqrt{\frac{\log\frac{2}{\delta}}{2m}}.$$

Mingsheng Long

(日) (同) (三) (三)

Generalization Bound with Covering Numbers

Theorem (Generalization Bound with Covering Numbers)

Let $\mathcal{F} \subseteq \mathbb{R}^{\mathcal{X} \times \mathcal{Y}}$ be a hypothesis set with label set $\mathcal{Y} = \{1, \cdots, k\}$ and $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be the corresponding \mathcal{Y} -valued labeling function class. Suppose $\Pi_1 \mathcal{F}$ is bounded in \mathcal{L}_2 by L. Fix $\rho > 0$. For all $\delta > 0$, with probability $1 - 3\delta$ the following inequality holds for all hypothesis $f \in \mathcal{F}$:

$$\epsilon_{Q}(f) \leq \epsilon_{\widehat{P}}^{(\rho)}(f) + d_{f,\mathcal{F}}^{(\rho)}(\widehat{P},\widehat{Q}) + \epsilon_{ideal} + 2\sqrt{\frac{\log\frac{2}{\delta}}{2n}} + \sqrt{\frac{\log\frac{2}{\delta}}{2m}} + \frac{16k^{2}\sqrt{k}}{\rho} \inf_{\epsilon \geq 0} \left\{ \epsilon + 3\left(\frac{1}{\sqrt{n}} + \frac{1}{\sqrt{m}}\right) \right.$$

$$\left. \left(\int_{\epsilon}^{L} \sqrt{\log\mathcal{N}_{2}(\tau,\Pi_{1}\mathcal{F})} \mathrm{d}\tau + L \int_{\epsilon/L}^{1} \sqrt{\log\mathcal{N}_{2}(\tau,\Pi_{1}\mathcal{H})} \mathrm{d}\tau \right) \right\}.$$

$$(20)$$

3

イロト イポト イヨト イヨト

MDD: Margin Disparity Discrepancy



Minimax game: Adversarial learning induced by informative margin theory

$$\min_{\substack{f,\psi \\ \psi(\widehat{P})}} \epsilon_{\psi(\widehat{P})}^{(\rho)}(f) + (\epsilon_{\psi(\widehat{Q})}^{(\rho)}(f^*, f) - \epsilon_{\psi(\widehat{P})}^{(\rho)}(f^*, f)),$$

$$f^* = \max_{f'} (\epsilon_{\psi(\widehat{Q})}^{(\rho)}(f', f) - \epsilon_{\psi(\widehat{P})}^{(\rho)}(f', f)).$$
(21)

Outline

Transfer Learning

2 $\mathcal{H} \Delta \mathcal{H}$ -Divergence

- DAN: Deep Adaptation Network
- DANN: Domain Adversarial Neural Network
- MCD: Maximum Classifier Discrepancy
- Margin Disparity Discrepancy
 MDD: Margin Disparity Discrepancy

Transfer Model Selection DEV: Deep Embedded Validation

Evaluation and Implementation

___ ▶

Model Selection in Domain Adaptation

Supervised Learning



- Semi-Supervised Learning (SSL)?
- Unsupervised Domain Adaptation (UDA)?



IWCV: Importance-Weighted Cross-Validation⁶

- Covariate shift assumption: $P(y|\mathbf{x}) = Q(y|\mathbf{x})$
- Model selection by estimating Target Risk $\epsilon_Q(h) = \mathbb{E}_Q[h(\mathbf{x}) \neq y]$
- Importance-Weighted Cross-Validation (IWCV)

$$\mathbb{E}_{P}w(\mathbf{x})\cdot[h(\mathbf{x})\neq y] = \mathbb{E}_{P}\frac{Q(\mathbf{x})}{P(\mathbf{x})}\cdot[h(\mathbf{x})\neq y] = \mathbb{E}_{Q}[h(\mathbf{x})\neq y] = \epsilon_{Q}(h)$$

- The estimation is unbiased but the variance is unbounded
- Density ratio is not accessible due to unknownness of P and Q



⁶Covariate shift adaptation by importance weighted cross validation JMLR'2007 > 💿 🤄 🦿

DEV: Deep Embedded Validation¹⁰

• Variance of IWCV (bounded by Rényi divergence)⁷:

 $\operatorname{Var}_{\mathbf{x}\sim P}[w(\mathbf{x})\cdot[h(\mathbf{x})\neq y]] \leq d_{\alpha+1}(Q\|P)\epsilon_Q(h)^{1-\frac{1}{\alpha}} - \epsilon_Q(h)^2$

- Density ratio $w(\mathbf{x}) = \frac{Q(\mathbf{x})}{P(\mathbf{x})}$ is estimated by discriminative learning⁸
- Feature adaptation reduces distribution discrepancy $d_{\alpha+1}(Q||P)^9$
- Control variate explicitly reduces the variance of $\mathbb{E}_P w(\mathbf{x}) \cdot [h(\mathbf{x}) \neq y]$
 - $\mathbb{E}[z] = \zeta, \mathbb{E}[t] = \tau$
 - $z^{\star} = z + \eta(t \tau)$
 - $\mathbb{E}[z^*] = \mathbb{E}[z] + \eta \mathbb{E}[t-\tau] = \zeta + \eta (\mathbb{E}[t] \mathbb{E}[\tau]) = \zeta.$
 - $\operatorname{Var}[z^{\star}] = \operatorname{Var}[z + \eta(t \tau)] = \eta^2 \operatorname{Var}[t] + 2\eta \operatorname{Cov}(z, t) + \operatorname{Var}[z]$
 - min Var[z^*] = $(1 \rho_{z,t}^2)$ Var[z], when $\hat{\eta} = -\frac{\text{Cov}(z,t)}{\text{Var}[t]}$

⁷Learning Bounds for Importance Weighting, NeurIPS 2010

⁸Discriminative learning for differing training and test distributions, ICML 2007

⁹Conditional Adversarial Domain Adaptation, NeurIPS 2018

¹⁰ You et al. Towards Accurate Model Selection in Deep Unsupervised Domain Adaptation. ICML 2019.

Outline

Transfer Learning

2 $\mathcal{H}\Delta\mathcal{H}$ -Divergence

- DAN: Deep Adaptation Network
- DANN: Domain Adversarial Neural Network
- MCD: Maximum Classifier Discrepancy
- 3 Margin Disparity Discrepancy
 - MDD: Margin Disparity Discrepancy
- Transfer Model Selectior
 - DEV: Deep Embedded Validation

5 Evaluation and Implementation

A (1) > A (2) > A

Datasets



VisDA Challenge 2017

3

イロト イヨト イヨト イヨト

Results

Method	$A\toW$	$D\toW$	$W\toD$	$A\toD$	$D\toA$	$W\toA$	Avg
ResNet-50	68.4±0.2	96.7±0.1	99.3±0.1	68.9±0.2	62.5±0.3	60.7±0.3	76.1
DAN	$80.5 {\pm} 0.4$	97.1±0.2	$99.6{\pm}0.1$	78.6±0.2	63.6±0.3	62.8±0.2	80.4
DANN	82.0±0.4	96.9±0.2	99.1±0.1	79.7±0.4	68.2±0.4	67.4±0.5	82.2
CDAN	93.0±0.2	98.4±0.2	100.0 ±.0	89.2±0.3	$70.2{\pm}0.4$	$69.4{\pm}0.4$	86.7
CDAN+E	$93.1{\pm}0.1$	98.6 ±0.1	100.0 ±.0	93.4±0.2	$71.0{\pm}0.3$	$70.3{\pm}0.3$	87.7
MDD	94.5 ±0.3	$98.4{\pm}0.1$	$\boldsymbol{100.0}{\pm}.0$	93.5 ±0.2	74.6 ±0.3	72.2 ±0.1	88.9

Table: Accuracy (%) on Office-31 for unsupervised domain adaptation

Table: Accuracy (%) on Office-Home for unsupervised domain adaptation

Method	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	/ Rw→Ar	Rw→Cl	Rw→Pr	Avg
ResNet-50	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1
DAN	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3
DANN	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6
CDAN	50.7	70.6	76.0	57.6	70.0	70.0	57.4	50.9	77.3	70.9	56.7	81.6	65.8
MDD	54.9	73.7	77.8	60.0	71.4	71.8	61.2	53.6	78.1	72.5	60.2	82.3	68.1

イロン 不聞 とくほど 不良とう ヨート

-

-

Results

Method	$Synthetic \to Real$
MCD	69.2
GTA	69.5
CDAN	70.0
MDD	74.6

Table: Accuracy (%) on VisDA-2017 (ResNet-50)

Table: Accuracy (%) of MCD by different validation methods on VisDA-2017

Method	plane	bcycl	bus	car	horse	knife	mcycl	person	plant	sktbrd	train	truck	mean
Original	87.00	60.90	83.70	64.00	88.90	79.60	84.70	76.90	88.60	40.30	83.00	25.80	71.90
Source Risk	84.39	54.11	69.15	46.37	80.49	80.45	85.04	65.24	87.22	36.86	78.04	28.91	66.36
IWCV	81.21	60.95	76.00	56.53	82.83	72.06	84.05	68.65	86.85	44.37	69.29	23.81	67.22
DEV (w/o control variate)	84.21	63.95	79.00	59.53	85.83	75.06	87.05	71.65	89.85	47.37	72.29	26.81	70.22
DEV	81.83	53.48	82.95	71.62	89.16	72.03	89.36	75.73	97.02	55.48	71.19	29.17	72.42
Target Risk (Upper Bound)	81.95	53.60	83.07	72.02	89.25	72.15	89.55	75.83	97.10	55.57	71.19	29.27	72.55

3

(日) (同) (三) (三)

Transfer Learning Systems

Tsinghua Dataway Big Data Software Stack



3

National Engineering Lab for Big Data Software



Jiaguang Sun (孙家广**)** Director Tsinghua University



Jianmin Wang (王建民**)** Deputy Director Dean, School of Software Tsinghua University



Michael I. Jordan Academic Committee Chair UC Berkeley

Mingsheng Long (龙明盛)

Machine Learning Group @ NELBDS



Yuchen Zhang



Yue Cao



Han Zhu

(日) (同) (三) (三)



Zhangjie Cao

August 21, 2019 37 / 37