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

- 4 @ > - 4 @ > - 4 @ >

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
 IDD: Independent and Differently Distributed (a case of Non-IID)
- How to effectively bound the generalization error on target domain?



Transfer Learning

- Transfer learning setups ($P \neq Q$): Feature Space X, Label Space Y
 - Domain Adaptation: common X, common Y, unlabeled T
 - Inductive Transfer Learning: common X, different Y, labeled T



Bias-Variance-Shift Tradeoff



Bridging Theory and Algorithm





• • • • • • • • • • • • •

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

2 Domain Adaptation

- $\mathcal{H}\Delta\mathcal{H}$ -Divergence
- MDD: Margin Disparity Discrepancy
- DEV: Deep Embedded Validation

Inductive Transfer Learning

A (10) A (10) A (10)

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.

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 $O(\sqrt{d \log n/n} + \sqrt{d \log m/m})$.

- 4 同 6 4 日 6 4 日 6

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) (周) (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)
$$\min_{\theta \in \Theta} \max_{k \in \mathcal{K}} \frac{1}{n_{a}} \sum_{i=1}^{n_{a}} L\left(\theta \left(\mathbf{x}_{i}^{a} \right), y_{i}^{a} \right) + \lambda \sum_{\ell=l_{1}}^{l_{2}} d_{k}^{2} \left(\widehat{P}_{\ell}, \widehat{Q}_{\ell} \right)$$
(7)

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

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

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

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

$$\leq \sup_{D \in \mathcal{H}_{D}} |\mathbb{E}_{P}[D(\mathbf{x}) = 1] + \mathbb{E}_{Q}[D(\mathbf{x}) = 0]|$$
(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. 🧃 🛌 🤊 🔍

Outline





Domain Adaptation

- $\mathcal{H}\Delta\mathcal{H}$ -Divergence
- MDD: Margin Disparity Discrepancy
- DEV: Deep Embedded Validation

Inductive Transfer Learning

3

< 回 ト < 三 ト < 三 ト

Towards Informative Margin Theory

- Towards a rigorous multiclass domain adaptation theory.
 - All existing theories are only applicable to binary classification.
 - Generalization bound with scoring functions has not been studied.
- Towards an informative margin theory.
 - Explore the idea of margin in measuring domain discrepancy.
 - Generalization bound with margin loss has not been studied.
- Towards a certain function class in the theoretical bound.
 - Eliminate approximation assumptions in all existing methods.
 - Computing the supremum in previous discrepancies requires an ergodicity over $\mathcal{H}\Delta\mathcal{H}$ that increases the difficulty of optimization.
- Towards bridging the existing gap between theories and algorithms.

(日) (同) (三) (三)

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:

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



(日) (周) (三) (三)

3

DD: Disparity Discrepancy

Definition (Disparity Discrepancy, DD)

Given a hypothesis space \mathcal{H} and a *specific classifier* $h \in \mathcal{H}$, the Disparity Discrepancy (DD) induced by $h' \in \mathcal{H}$ is defined by

$$d_{h,\mathcal{H}}(P,Q) = \sup_{h' \in \mathcal{H}} \left| \mathbb{E}_{Q}[h' \neq h] - \mathbb{E}_{P}[h' \neq h] \right|.$$
(11)

The supremum in the disparity discrepancy is taken **only over the hypothesis space** \mathcal{H} and thus can be optimized more easily.

Theorem

For every hypothesis $h \in \mathcal{H}$,

$$\epsilon_Q(h) \le \epsilon_P(h) + d_{h,\mathcal{H}}(P,Q) + \epsilon_{ideal},$$
(12)

where $\epsilon_{ideal} = \epsilon(\mathcal{H}, P, Q)$ is the ideal combined loss.

CDAN: Conditional Domain Adversarial Network⁴



Conditional adaptation of distributions over representation & prediction

$$\min_{G} \mathcal{E}(G) - \lambda \mathcal{E}(D, G)$$

$$\min_{D} \mathcal{E}(D, G),$$
(13)

 $\mathcal{E}(D,G) = -\mathbb{E}_{\mathbf{x}_{i}^{s} \sim \mathcal{D}_{s}} \log \left[D\left(\mathbf{f}_{i}^{s} \otimes \mathbf{g}_{i}^{s}\right) \right] - \mathbb{E}_{\mathbf{x}_{j}^{t} \sim \mathcal{D}_{t}} \log \left[1 - D\left(\mathbf{f}_{j}^{t} \otimes \mathbf{g}_{j}^{t}\right) \right]$ (14)

⁴Long et al. Conditional Adversarial Domain Adaptation. NIPS 2018. 🚓 🖘 📱 🔊 🔍

N/1.m	~~	hon	~	0.	20
	25	пеп	2		ı٢
	~		-		

MDD: Margin Disparity Discrepancy⁵

- 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

$$d_{f,\mathcal{F}}^{(\rho)}(P,Q) \triangleq \sup_{f' \in \mathcal{F}} \left(\epsilon_Q^{(\rho)}(f',f) - \epsilon_P^{(\rho)}(f',f) \right), \\ d_{f,\mathcal{F}}^{(\rho)}(\widehat{P},\widehat{Q}) \triangleq \sup_{f' \in \mathcal{F}} \left(\epsilon_{\widehat{Q}}^{(\rho)}(f',f) - \epsilon_{\widehat{P}}^{(\rho)}(f',f) \right).$$
(15)

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

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

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 $f \in \mathcal{F}$,

$$\epsilon_Q(f) \le \epsilon_P^{(\rho)}(f) + d_{f,\mathcal{F}}^{(\rho)}(P,Q) + \epsilon_{ideal}^{(\rho)}, \tag{16}$$

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

$$\epsilon_{ideal}^{(\rho)} = \min_{f^* \in \mathcal{F}} \{ \epsilon_P^{(\rho)}(f^*) + \epsilon_Q^{(\rho)}(f^*) \}.$$
(17)

- This upper bound has a similar form with previous bound.
 - $\epsilon_P^{(\rho)}(f)$ depicts the performance of f on source domain.
 - MDD bounds the performance gap caused by domain shift.
 - ϵ_{ideal} quantifies the inverse of "adaptability".

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

Mingsheng Long

Definitions

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

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

$$\Pi_1 \mathcal{F} = \{ x \mapsto f(x, y) | y \in \mathcal{Y}, f \in \mathcal{F} \}.$$
(18)

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

$$\Pi_{\mathcal{H}}\mathcal{F} \triangleq \{ x \mapsto f(x, h(x)) | h \in \mathcal{H}, f \in \mathcal{F} \}.$$
(19)

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

$$\widehat{\mathfrak{R}}_{\widehat{D}}(\mathcal{G}) = \mathbb{E}_{\sigma} \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_i g(z_i).$$
(20)

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

$$\mathfrak{R}_{n,D}(\mathcal{G}) = \mathbb{E}_{\widehat{D} \sim D^n} \widehat{\mathfrak{R}}_{\widehat{D}}(\mathcal{G}).$$
(21)

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 $f \in \mathcal{F}$:

$$\epsilon_{Q}(f) \leq \epsilon_{\widehat{P}}^{(\rho)}(f) + d_{f,\mathcal{F}}^{(\rho)}(\widehat{P},\widehat{Q}) + \epsilon_{ideal} + \frac{2k^{2}}{\rho} \mathfrak{R}_{n,P}(\Pi_{1}\mathcal{F}) + \frac{k}{\rho} \mathfrak{R}_{n,P}(\Pi_{\mathcal{H}}\mathcal{F}) + 2\sqrt{\frac{\log\frac{2}{\delta}}{2n}}$$
(22)
$$+ \frac{k}{\rho} \mathfrak{R}_{m,Q}(\Pi_{\mathcal{H}}\mathcal{F}) + \sqrt{\frac{\log\frac{2}{\delta}}{2m}}.$$

Mingsheng Long

(日) (同) (三) (三)

Rademacher Bound of Linear Classifier

We need to check the variation of $\mathfrak{R}_{n,D}(\Pi_{\mathcal{H}}\mathcal{F})$ with the growth of *n*. First, we include a simple example of binary linear classifiers.

Theorem

Let $S \subseteq \mathcal{X} = \{\mathbf{x} \in \mathbb{R}^s | \|\mathbf{x}\|_2 \le r\}$ be a sample of size m and suppose

$$\begin{aligned} \mathcal{F} &= \big\{ f : \mathcal{X} \times \{ \pm 1 \} \to \mathbb{R} \ \big| \ f(\mathbf{x}, y) = \operatorname{sgn}(y) \ \mathbf{w} \cdot \mathbf{x}, \ \|\mathbf{w}\|_2 \leq \Lambda \big\}, \\ \mathcal{H} &= \big\{ h \mid h(\mathbf{x}) = \operatorname{sgn}(\mathbf{w} \cdot \mathbf{x}), \ \|\mathbf{w}\|_2 \leq \Lambda \big\}. \end{aligned}$$

Then the empirical Rademacher complexity of $\Pi_{\mathcal{H}}\mathcal{F}$ can be bounded as follows:

$$\widehat{\mathfrak{R}}_{\mathcal{S}}(\Pi_{\mathcal{H}}\mathcal{F}) \leq 2\Lambda r \sqrt{rac{d\lograc{em}{d}}{m}},$$

where d is the VC-dimension of \mathcal{H} .

(日) (周) (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 $f \in \mathcal{F}$:

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

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

$$(23)$$

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

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

$$f^* = \max_{f'} (\epsilon_{\psi(\widehat{Q})}^{(\rho)}(f', f) - \epsilon_{\psi(\widehat{P})}^{(\rho)}(f', f)).$$
(24)

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
DANN	82.0±0.4	96.9±0.2	$99.1{\pm}0.1$	79.7±0.4	68.2±0.4	$67.4 {\pm} 0.5$	82.2
JAN	85.4±0.3	97.4±0.2	99.8±0.2	84.7±0.3	68.6±0.3	$70.0{\pm}0.4$	84.3
MCD	88.6±0.2	98.5±0.1	100.0 ±.0	92.2±0.2	$69.5 {\pm} 0.1$	69.7±0.3	86.5
CDAN	94.1±0.1	98.6 ±0.1	100.0 ±.0	92.9±0.2	$71.0 {\pm} 0.3$	69.3±0.3	87.7
MDD (Proposed)	94.5 ±0.3	$98.4{\pm}0.1$	$100.0{\pm}.0$	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	CI-Pr	CI-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
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
JAN	45.9	61.2	68.9	50.4	59.7	61.0	45.8	43.4	70.3	63.9	52.4	76.8	58.3
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 (Proposed)	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

・ロト ・聞 ト ・ ヨト ・ ヨト ・ ヨー・

Results



Figure: Test accuracy and empirical values of $\sigma_{h_f} \circ f'$ (dashed line: $\frac{\gamma}{\gamma+1}$).

Table: Accuracy (%) on Office-31 by different margins.

Margin γ	1	2	3	4	5	6
$A\toW$	92.5	93.7	94.0	94.5	93.8	93.5
$D\toA$	72.4	73.0	73.7	74.6	74.3	74.2
Avg on Office-31	87.6	88.1	88.5	88.9	88.7	88.6

Mingsheng Long

э

э

A 🖓

Outline

Transfer Learning



Domain Adaptation

- $\mathcal{H} \Delta \mathcal{H}$ -Divergence
- MDD: Margin Disparity Discrepancy
- DEV: Deep Embedded Validation

Inductive Transfer Learning

3

< 回 ト < 三 ト < 三 ト

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 · 📑 🗠 🔍

Mingsheng Long	Transfer Learning	October 17, 2019	34 / 50

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.

DEV: Deep Embedded Validation

Algorithm 1 Deep Embedded Validation (DEV)

- 1: Input: Candidate model $g(\mathbf{x}) = \mathcal{T}(F(\mathbf{x}))$ Training set $\mathcal{D}_{tr} = \{(\mathbf{x}_i^{tr}, y_i^{tr})\}_{i=1}^{n_{tr}}$ Validation set $\mathcal{D}_v = \{(\mathbf{x}_i^{v}, y_i^v)\}_{i=1}^{n_v}$ Test set $\mathcal{D}_{ts} = \{(\mathbf{x}_i^{ts})\}_{i=1}^{n_{ts}}$
- 2: **Output:** DEV Risk $\mathcal{R}_{DEV}(g)$ of model g
- 3: Compute features and predictions using model g:

$$\mathcal{F}_{tr} = \{\mathbf{f}_{i}^{tr}\}_{i=1}^{n_{tr}}, \ \mathcal{F}_{ts} = \{\mathbf{f}_{i}^{ts}\}_{i=1}^{n_{ts}}, \ \mathcal{F}_{v} = \{\mathbf{f}_{i}^{v}\}_{i=1}^{n_{v}}, \ \mathcal{Y}_{v} = \{\hat{y}_{i}^{v}\}_{i=1}^{n_{v}}$$
4: Train a two-layer logistic regression model M to classify \mathcal{F}_{tr} and \mathcal{F}_{ts}
(label \mathcal{F}_{tr} as 1 and \mathcal{F}_{ts} as 0)
5: Compute $w_{\mathbf{f}}(\mathbf{x}_{i}^{v}) = \frac{n_{tr}}{n_{ts}} \frac{1-M(f_{i}^{v})}{M(f_{i}^{v})}, \ W = \{w_{\mathbf{f}}(\mathbf{x}_{i}^{v})\}_{i=1}^{n_{v}}$
6: Compute weighted loss $L = \{w_{\mathbf{f}}(\mathbf{x}_{i}^{v})\ell(\hat{y}_{i}^{v}, y_{i}^{v})\}_{i=1}^{n_{v}}$
7: Estimate coefficient $\eta = -\frac{\widehat{Cov}(L,W)}{\widehat{Var}[W]}$
8: Compute DEV Risk $\mathcal{R}_{\text{DEV}}(g) = \text{mean}(L) + \eta \text{mean}(W) - \eta$

(a)

Outline

Transfer Learning

Domain Adaptation

- $\mathcal{H} \Delta \mathcal{H}$ -Divergence
- MDD: Margin Disparity Discrepancy
- DEV: Deep Embedded Validation

Inductive Transfer Learning

3

(日) (同) (三) (三)

Inductive Transfer Learning¹¹

Successful of transfer learning: Pre-train a model on a large-scale source dataset, and use the parameters as initialization for training a target task.

Compared to training from scratch:

- Generalization: better accuracy
- Optimization: faster convergence

How to understand the transferability of deep representations?



 11 Liu et al. Towards Understanding the Transferability of Deep Representations, arXiv, 2019 $_{\odot}$ $_{\odot}$

M	ings	heng	ong	
			- ong	

Transferred Parameters Induce Better Generalization

We can quantify how pre-trained knowledge is preserved when transferring to the target dataset with $\frac{1}{\sqrt{n}} \|\mathbf{W}_Q - \mathbf{W}_P\|_F$.

- For more similar target datasets, $\frac{1}{\sqrt{n}} \| \mathbf{W}_Q \mathbf{W}_P \|_F$ is smaller
- For more similar target datasets, generalization error is smaller
- Is $\frac{1}{\sqrt{n}} \| \mathbf{W}_Q \mathbf{W}_P \|_F$ implicitly bounded? (we will formally study this)



Transferred Parameters Induce Better Generalization

Staying close to the transferred parameters benefits generalization

- Even for the same target dataset, different pre-trained parameters lead to significantly different solutions
- At the convergence point, pre-trained networks stay in the original flat region, leading to flatter minima than random initialization



(a) t-SNE of parameters



(b) Randomly initialized



(c) ImageNet pre-trained

(日) (同) (三) (三)

Transferred Parameters Enable Faster Optimization

Modern neural networks are equipped with Batch Normalization (BN) and skip connections to enable better loss landscapes

• However, at the initialization point, the loss landscapes are still very messy even in the presence of Batch-Norm and residual connections



- 4 回 ト - 4 回 ト

Transferred Parameters Enable Faster Optimization

Pre-trained parameters help smoothen the loss landscape and accelerate training in the early stages

• The landscapes can be described with the Lipschitzness of the loss function, i.e. the magnitude of gradient



Transferred Parameters Enable Better Optimization

Why is the magnitude of gradient better controlled with the pre-trained representations?

• The gradient is computed through back-prop, $\frac{\partial L}{\partial \mathbf{x}_i^{k-1}} = \mathbf{W}_k \mathbb{I}_i^k \left(\frac{\partial L}{\partial \mathbf{x}_i^k} \right)$. Pre-trained weight matrices provide more stable scaling factors.



∃ → (∃ →

Feasibility of Transfer Learning

Varying input with fixed labels.

• Choosing a model pre-trained on more similar inputs yields a larger performance gain.

Varying labels with fixed input.

• The similarity of the input (images) is just one point. Another factor of similarity is the relationship between the nature of tasks (labels).



Feasibility of Transfer Learning

Choices of pre-training epochs.

 Although the test accuracy on the pre-training dataset continues increasing, the test accuracy on the target dataset starts to decline.



Theoretical Analysis

• Two-layer ReLU network of *m* hidden units $f_{\mathbf{W},\mathbf{a}}(\mathbf{x}) = \frac{1}{\sqrt{m}} \mathbf{a}^{\top} \sigma(\mathbf{W}^{\top} \mathbf{x})$.

•
$$\mathbf{x} \in \mathbb{R}^d$$
, $\mathbf{W} = (\mathbf{w}_1, \cdots, \mathbf{w}_m) \in \mathbb{R}^{d imes m}$.

- $\mathbf{w}_r(0) \sim \mathcal{N}(\mathbf{0}, \kappa^2 \mathbf{I}), a_r \sim \text{unif} (\{-1, 1\}).$
- $L(\mathbf{W}) = \frac{1}{2}(\mathbf{y} f_{\mathbf{W},\mathbf{a}}(\mathbf{X}))^{\top}(\mathbf{y} f_{\mathbf{W},\mathbf{a}}(\mathbf{X})).$

We first pre-train the model on $\{\mathbf{x}_{P,i}, y_{P,i}\}_{i=1}^{n_P}$ drawn i.i.d from P to obtain $\mathbf{W}(P)$, then train on the target dataset $\{\mathbf{x}_{Q,i}, y_{Q,i}\}_{i=1}^{n_Q}$ drawn i.i.d from Q.

Definition (Gram matrix of P and Q)

$$\mathbf{H}_{P,ij}^{\infty} = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0},\mathbf{I})}[\mathbf{x}_{P,i}^{\top}\mathbf{x}_{P,j}\mathbb{I}\{\mathbf{w}^{\top}\mathbf{x}_{P,i} \geq 0, \ \mathbf{w}^{\top}\mathbf{x}_{P,j} \geq 0\}].$$
(25)

$$\mathbf{H}_{Q,ij}^{\infty} = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0},\mathbf{I})}[\mathbf{x}_{Q,i}^{\top}\mathbf{x}_{Q,j}\mathbb{I}\{\mathbf{w}^{\top}\mathbf{x}_{Q,i} \geq 0, \ \mathbf{w}^{\top}\mathbf{x}_{Q,j} \geq 0\}].$$
 (26)

Definition (Gram matrix of transfer learning)

$$\mathbf{H}_{PQ,ij}^{\infty} = \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0},\mathbf{l})}[\mathbf{x}_{P,i}^{\top}\mathbf{x}_{Q,j}\mathbb{I}\{\mathbf{w}^{\top}\mathbf{x}_{P,i} \ge 0, \ \mathbf{w}^{\top}\mathbf{x}_{Q,j} \ge 0\}].$$
(27)

Theoretical Analysis

Definition (Transformed labels from source label set to target label set) $\mathbf{y}_{P \to Q} \triangleq \mathbf{H}_{PQ}^{\infty} {}^{\top} \mathbf{H}_{P}^{\infty-1} \mathbf{y}_{P}.$ (28)

Theorem (Improved Lipschitzness) Denote by \mathbf{X}^1 the activations in the target dataset. If $m \ge \text{poly}(n_P, n_Q, \delta^{-1}, \lambda_P^{-1}, \lambda_Q^{-1}, \kappa^{-1}), \ \kappa = O\left(\frac{\lambda_P^2 \delta}{n_P^2 n_Q^2}\right)$, with probability no less than $1 - \delta$ over the random initialization,

$$\|\frac{\partial L(\mathbf{W}(P))}{\partial \mathbf{X}^{1}}\|^{2} = \|\frac{\partial L(\mathbf{W}(0))}{\partial \mathbf{X}^{1}}\|^{2} - \mathbf{y}_{Q}^{\top}\mathbf{y}_{Q} + (\mathbf{y}_{Q} - \mathbf{y}_{P \to Q})^{\top}(\mathbf{y}_{Q} - \mathbf{y}_{P \to Q}) + \frac{\operatorname{poly}(n_{P}, n_{Q}, \delta^{-1}, \lambda_{P}^{-1}, \kappa^{-1})}{m^{\frac{1}{4}}} + O\left(\frac{n_{P}^{2}n_{Q}^{\frac{1}{2}}\kappa}{\lambda_{P}^{2}\delta}\right).$$
(29)

Theoretical Analysis

Theorem (Improved generalization) Suppose $m \ge \text{poly}(n_P, n_Q, \delta^{-1}, \lambda_P^{-1}, \lambda_Q^{-1}, \kappa^{-1})$, $\kappa = O\left(\frac{\lambda_P^2 \lambda_Q^2 \delta}{n_P^2 n_Q^2}\right)$, with probability

no less than $1 - \delta$ over the random initialization.

$$\|\mathbf{W}(Q) - \mathbf{W}(P)\|_{F} \leq \sqrt{(\mathbf{y}_{Q} - \mathbf{y}_{P \to Q})^{\top} \mathbf{H}_{Q}^{\infty^{-1}}(\mathbf{y}_{Q} - \mathbf{y}_{P \to Q})} + O\left(\frac{n_{P}n_{Q}^{\frac{1}{2}}\kappa^{\frac{1}{2}}}{\lambda_{P}\lambda_{Q}\delta^{\frac{1}{2}}}\right) + \frac{\operatorname{poly}(n_{P}, n_{Q}, \delta^{-1}, \lambda_{P}^{-1}, \lambda_{Q}^{-1}, \kappa^{-1})}{m^{\frac{1}{4}}}.$$
(30)

Lemma (Arora et al., 2019)

Under the same conditions as (30), with probability no less than $1 - \delta$,

$$\mathbb{E}_{Q}(L(f(\mathbf{x}))) \leq \sqrt{\frac{2}{n_{Q}}} \|\mathbf{W}(Q) - \mathbf{W}(P)\|_{F} + O\left(\sqrt{\frac{1}{n_{Q}}}\right)$$

Mingsheng Long

Transfer Learning System

Tsinghua Dataway Big Data Software Stack



Mingsheng Long

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

- 4 同 1 4 回 1 4 回 1



Zhangjie Cao

October 17, 2019 50 / 50