Transfer Learning Generalizing Deep Learning across Domains and Tasks

Mingsheng Long

School of Software National Engineering Lab for Big Data Software Tsinghua University

https://github.com/thuml Chinese Conference on Pattern Recognition and Computer Vision PRCV 2018

Mingsheng Long

May 18, 2019 1 / 31

A = A = A

Joint Work With:



Jianmin Wang Professor Tsinghua University jimwang@tsinghua.edu.cn



Mingsheng Long Associate Professor Tsinghua University mingsheng@tsinghua.edu.cn



Michael I. Jordan Professor UC Berkeley jordan@cs.berkeley.edu



Yuchen Zhang



Yue Cao



Han Zhu



Zhangjie Cao

Mingsheng Long

Transfer Learning

May 18, 2019 2 / 31

Outline

Transfer Learning

- 2 Problem I: P(X) ≠ Q(X)
 DAN: Deep Adaptation Network
- 3 Problem II: P(Y|X) ≠ Q(Y|X)
 CDAN: Conditional Domain Adversarial Network
- 4 Theoretical Analysis
- 5 Benchmarking

(日) (同) (三) (三)

Machine Learning

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



complexity

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

Transfer Learning

- Machine learning across domains of Non-IID distributions $P \neq Q$
- How to design models that effectively bound the generalization error?



Bias-Variance-Shift Tradeoff



Basic Approaches to Transfer Learning

Matching distributions across source and target domains s.t. P pprox Q

- Reduce marginal distribution mismatch: $P(\mathbf{X}) \neq Q(\mathbf{X})$
- Reduce conditional distribution mismatch: $P(Y|\mathbf{X}) \neq Q(Y|\mathbf{X})$
- Challenge: fail to align different domains of multimodal distributions



Kernel Embedding

Adversarial Learning

(日) (周) (三) (三)

Song et al. Kernel Embeddings of Conditional Distributions. **IEEE**, 2013. Goodfellow et al. Generative Adversarial Networks. **NIPS** 2014.

Mingsheng Long

May 18, 2019 7 / 31

Basic Guidelines to Algorithm Design



Everything should be made as simple as possible, but no simpler. —Albert Einstein

(日) (同) (三) (三)

Outline

Transfer Learning

2 Problem I: P(X) ≠ Q(X)
 • DAN: Deep Adaptation Network

Broblem II: P(Y|X) ≠ Q(Y|X) CDAN: Conditional Domain Adversarial Network

4 Theoretical Analysis

5 Benchmarking

(日) (同) (三) (三)

DAN: Deep Adaptation Network¹



Deep adaptation: match distributions in multiple domain-specific layers Optimal matching: maximize two-sample test power by multiple kernels

$$d_{k}^{2}(P,Q) \triangleq \left\| \mathbf{E}_{P} \left[\phi \left(\mathbf{x}^{s} \right) \right] - \mathbf{E}_{Q} \left[\phi \left(\mathbf{x}^{t} \right) \right] \right\|_{\mathcal{H}_{k}}^{2}$$
(1)

$$\min_{\theta \in \Theta} \max_{k \in \mathcal{K}} \frac{1}{n_a} \sum_{i=1}^{n_a} J(\theta(\mathbf{x}_i^a), y_i^a) + \lambda \sum_{\ell=l_1}^{l_2} d_k^2 \left(\mathcal{D}_s^{\ell}, \mathcal{D}_t^{\ell} \right)$$
(2)

 1 Long et al. Learning Transferable Features with Deep Adaptation Networks. IGML '15. and \sim

Mingsheng Long

DAN: MK-MMD

Multiple Kernel Maximum Mean Discrepancy (MK-MMD)

RKHS distance between kernel embeddings of distributions P_X and Q_X

$$d_{k}^{2}(P,Q) \triangleq \left\| \mathbf{E}_{P} \left[\phi \left(\mathbf{x}^{s} \right) \right] - \mathbf{E}_{Q} \left[\phi \left(\mathbf{x}^{t} \right) \right] \right\|_{\mathcal{H}_{k}}^{2}, \tag{3}$$

 $k(\mathbf{x}^{s}, \mathbf{x}^{t}) = \langle \phi(\mathbf{x}^{s}), \phi(\mathbf{x}^{t}) \rangle \text{ is a convex combination of } m \text{ PSD kernels}$ $\mathcal{K} \triangleq \left\{ k = \sum_{u=1}^{m} \beta_{u} k_{u} : \sum_{u=1}^{m} \beta_{u} = 1, \beta_{u} \ge 0, \forall u \right\}.$ (4)

Theorem (Kernel Two-Sample Test (Gretton et al. 2012))

- P = Q if and only if $d_k^2(P,Q) = 0$ (In practice, $d_k^2(P,Q) < \epsilon$)
- $\max_{k \in \mathcal{K}} d_k^2(P, Q) \sigma_k^{-2} \Leftrightarrow \min Type \ II \ Error \ (d_k^2(P, Q) < \epsilon \ when \ P \neq Q)$

DAN: Feature Learning

Linear-Time Algorithm of MK-MMD (Streaming Algorithm)

$$\begin{array}{l} O(n^2): \ d_k^2(p,q) = \mathbf{E}_{\mathbf{x}^s \mathbf{x}'^s} k(\mathbf{x}^s, \mathbf{x}'^s) + \mathbf{E}_{\mathbf{x}^t \mathbf{x}'^t} k(\mathbf{x}^t, \mathbf{x}'^t) - 2\mathbf{E}_{\mathbf{x}^s \mathbf{x}^t} k(\mathbf{x}^s, \mathbf{x}^t) \\ O(n): \ d_k^2(p,q) = \frac{2}{n_s} \sum_{i=1}^{n_s/2} g_k(\mathbf{z}_i) \rightarrow \text{linear-time unbiased estimate} \\ \bullet \ \text{Quad-tuple } \mathbf{z}_i \triangleq (\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^s, \mathbf{x}_{2i-1}^t, \mathbf{x}_{2i}^t) \\ \bullet \ g_k(\mathbf{z}_i) \triangleq k(\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^s) + k(\mathbf{x}_{2i-1}^t, \mathbf{x}_{2i}^t) - k(\mathbf{x}_{2i-1}^s, \mathbf{x}_{2i}^t) - k(\mathbf{x}_{2i}^s, \mathbf{x}_{2i-1}^t) \end{array}$$

Stochastic Gradient Descent (SGD)

For each layer ℓ and for each quad-tuple $\mathbf{z}_i^\ell = \left(\mathbf{h}_{2i-1}^{s\ell}, \mathbf{h}_{2i}^{s\ell}, \mathbf{h}_{2i-1}^{t\ell}, \mathbf{h}_{2i}^{t\ell}\right)$

$$\nabla_{\Theta^{\ell}} = \frac{\partial J(\mathbf{z}_i)}{\partial \Theta^{\ell}} + \lambda \frac{\partial g_k(\mathbf{z}_i^{\ell})}{\partial \Theta^{\ell}}$$
(5)

DAN: Kernel Learning

Learning optimal kernel $k = \sum_{u=1}^{m} \beta_u k_u$

Maximizing test power \triangleq minimizing Type II error (Gretton et al. 2012)

$$\max_{k \in \mathcal{K}} d_k^2 \left(\mathcal{D}_s^\ell, \mathcal{D}_t^\ell \right) \sigma_k^{-2}, \tag{6}$$

where $\sigma_k^2 = \mathbf{E}_{\mathbf{z}} g_k^2 (\mathbf{z}) - [\mathbf{E}_{\mathbf{z}} g_k (\mathbf{z})]^2$ is the estimation variance.

Quadratic Program (QP), scaling linearly to sample size: $O(m^2n + m^3)$

$$\min_{\mathbf{d}^{\mathsf{T}}\boldsymbol{\beta}=1,\boldsymbol{\beta}\geq\mathbf{0}}\boldsymbol{\beta}^{\mathsf{T}}\left(\mathbf{Q}+\epsilon\mathbf{I}\right)\boldsymbol{\beta},\tag{7}$$

(日) (同) (三) (三)

where $\mathbf{d} = (d_1, d_2, \dots, d_m)^T$, and each d_u is MMD using base kernel k_u .

DANN: Domain Adversarial Neural Network²



Adversarial adaptation: learning features indistinguishable across domains

$$E\left(\theta_{f},\theta_{y},\theta_{d}\right) = \sum_{\mathbf{x}_{i}\in\mathcal{D}_{s}}L_{y}\left(G_{y}\left(G_{f}\left(\mathbf{x}_{i}\right)\right),y_{i}\right) - \lambda\sum_{\mathbf{x}_{i}\in\mathcal{D}_{s}\cup\mathcal{D}_{t}}L_{d}\left(G_{d}\left(G_{f}\left(\mathbf{x}_{i}\right)\right),d_{i}\right)$$
(8)

$$(\hat{\theta}_{f},\hat{\theta}_{y}) = \arg\min_{\theta_{f},\theta_{y}} E\left(\theta_{f},\theta_{y},\theta_{d}\right) \quad (\hat{\theta}_{d}) = \arg\max_{\theta_{d}} E\left(\theta_{f},\theta_{y},\theta_{d}\right) \tag{9}$$

²Ganin et al. Domain Adversarial Training of Neural Networks. JMLR '16 🕨 🧃 🐑 🧝 🗠 🔍

Outline

Transfer Learning

2 Problem I: P(X) ≠ Q(X)
 • DAN: Deep Adaptation Network

Il: P(Y|X) ≠ Q(Y|X) CDAN: Conditional Domain Adversarial Network

4 Theoretical Analysis

Benchmarking

- **(())) (())) ())**

CDAN: Conditional Domain Adversarial Network³

Main Idea of This Work: Distribution Embeddings with Statistics

- Capture cross-covariance statistics across multiple random vectors
 - Concatenation: $\mathbb{E}_{XY}[X \oplus Y] = \mathbb{E}_{X}[X] \oplus \mathbb{E}_{Y}[Y]$
 - Multilinear: $\mathbb{E}_{XY}[X \otimes Y] = \mathbb{E}_{X}[X|Y = 1] \oplus \ldots \oplus \mathbb{E}_{X}[X|Y = C]$



³Long et al. Conditional Adversarial Domain Adaptation. NIPS '18. . . .

Mingsheng Long

May 18, 2019 16 / 31

CDAN: Multilinear Conditioning



Conditional adaptation of distributions over representation & prediction $\min_{G} E_{G} - \lambda E_{D,G}$ $\min_{D} E_{D,G}$ (10)

$$E_{D,G} = -\frac{1}{n_s} \sum_{i=1}^{n_s} \log\left(D\left(\mathbf{f}_i^s \otimes \mathbf{g}_i^s\right)\right) - \frac{1}{n_t} \sum_{j=1}^{n_t} \log\left(1 - D\left(\mathbf{f}_j^t \otimes \mathbf{g}_j^t\right)\right) \quad (11)$$

-

- ₹ 🗦 🕨

CDAN: Randomized Multilinear Conditioning



Conditional adaptation of distributions over representation & prediction

$$\mathcal{T}_{\otimes}\left(\mathbf{f},\mathbf{g}\right) = \mathbf{f} \otimes \mathbf{g} \tag{12}$$

$$T_{\odot}(\mathbf{f}, \mathbf{g}) = \frac{1}{\sqrt{d}} \left(\mathbf{R}_{\mathbf{f}} \mathbf{f} \right) \odot \left(\mathbf{R}_{\mathbf{g}} \mathbf{g} \right)$$
(13)

$$\mathcal{T}(\mathbf{h}) = \begin{cases} \mathcal{T}_{\otimes}(\mathbf{f}, \mathbf{g}) & \text{if } d_f \times d_g \leqslant 4096 \\ \mathcal{T}_{\odot}(\mathbf{f}, \mathbf{g}) & \text{otherwise} \end{cases}$$
(14)

CDAN: Entropy Conditioning



Control the uncertainty of classifier prediction to guarantee transferability

$$\max_{D} \frac{1}{n_{s}} \sum_{i=1}^{n_{s}} e^{-H(\mathbf{g}_{i}^{s})} \log \left[D\left(T\left(\mathbf{h}_{i}^{s}\right)\right) \right] + \frac{1}{n_{t}} \sum_{j=1}^{n_{t}} e^{-H\left(\mathbf{g}_{j}^{t}\right)} \log \left[1 - D\left(T\left(\mathbf{h}_{j}^{t}\right)\right) \right]$$

$$(15)$$

CDAN: Minimax Optimization Problem

Principled approaches: Conditional Domain Adversarial Networks (CDAN)

- Multilinear Conditioning: capture the cross-covariance between feature representation & classifier prediction to boost discriminability
- Entropy Conditioning: control the uncertainty of classifier prediction to guarantee transferability (entropy minimization principle)

A B F A B F

Outline

Transfer Learning

2 Problem I: P(X) ≠ Q(X)
 • DAN: Deep Adaptation Network

3 Problem II: P(Y|X) ≠ Q(Y|X) • CDAN: Conditional Domain Adversarial Network

Theoretical Analysis

Benchmarking

(日) (同) (三) (三)

Notations and Assumptions

- Source risk: $\epsilon_{P}(G) = \mathbb{E}_{(\mathbf{f}, \mathbf{y}) \sim P}[G(\mathbf{f}) \neq \mathbf{y}]$
- Target risk: $\epsilon_Q(G) = \mathbb{E}_{(\mathbf{f}, \mathbf{y}) \sim Q}[G(\mathbf{f}) \neq \mathbf{y}]$
- Disagreement on source: $\epsilon_{P}(G_{1}, G_{2}) = \mathbb{E}_{(\mathbf{f}, \mathbf{y}) \sim P}[G_{1}(\mathbf{f}) \neq G_{2}(\mathbf{f})]$
- Disagreement on target: $\epsilon_Q(G_1, G_2) = \mathbb{E}_{(\mathbf{f}, \mathbf{y}) \sim Q}[G_1(\mathbf{f}) \neq G_2(\mathbf{f})]$
- Idea hypothesis: $G^* = \arg \min_G \epsilon_P(G) + \epsilon_Q(G)$
- Assumption: idea hypothesis has small risk $\epsilon_{ideal} = \epsilon_P(G^*) + \epsilon_Q(G^*)$





Ideal hypothesis with small error

Generalization Bound

Theorem

The probabilistic bound of the target risk $\epsilon_Q(G)$ of hypothesis G is given by the source risk $\epsilon_P(G)$ plus the distribution discrepancy:

$$\epsilon_{Q}(G) \leq \epsilon_{P}(G) + [\epsilon_{P}(G^{*}) + \epsilon_{Q}(G^{*})] + |\epsilon_{P}(G, G^{*}) - \epsilon_{Q}(G, G^{*})|$$
(18)

Proof.

By using the triangle inequalities, we have

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

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

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

$$\leq \epsilon_{P}(G) + [\epsilon_{P}(G^{*}) + \epsilon_{Q}(G^{*})] + |\epsilon_{P}(G, G^{*}) - \epsilon_{Q}(G, G^{*})|$$
(19)

Joint Distribution Discrepancy

Define the proxies of the joint distributions $P(\mathbf{x}, \mathbf{y})$ and $Q(\mathbf{x}, \mathbf{y})$

•
$$P_G = (\mathbf{f}, G(\mathbf{f}))_{\mathbf{f} \sim P(\mathbf{f})}, \ Q_G = (\mathbf{f}, G(\mathbf{f}))_{\mathbf{f} \sim Q(\mathbf{f})}$$

• $\epsilon_P(G, G^*) = \epsilon_{P_G}(G^*), \ \epsilon_Q(G, G^*) = \epsilon_{Q_G}(G^*)$

Proof.

$$\epsilon_{P}\left(G,G^{*}\right) = \mathbb{E}_{\left(\mathbf{f},\mathbf{y}\right)\sim P}\left[G\left(\mathbf{f}\right)\neq G^{*}\left(\mathbf{f}\right)\right] = \mathbb{E}_{\left(\mathbf{f},\mathbf{g}\right)\sim P_{G}}\left[\mathbf{g}\neq G^{*}\left(\mathbf{f}\right)\right] = \epsilon_{P_{G}}\left(G^{*}\right)$$

How to bound the distribution discrepancy $|\epsilon_P(G, G^*) - \epsilon_Q(G, G^*)|$?



Upper-Bounding the Distribution Discrepancy

The distribution discrepancy $|\epsilon_{P}(G, G^{*}) - \epsilon_{Q}(G, G^{*})|$ is bounded by

 $\begin{aligned} |\epsilon_{P}\left(G,G^{*}\right)-\epsilon_{Q}\left(G,G^{*}\right)| &= \left|\mathbb{E}_{\left(\mathbf{f},\mathbf{g}\right)\sim P_{G}}\left[\mathbf{g}\neq G^{*}\left(\mathbf{f}\right)\right]-\mathbb{E}_{\left(\mathbf{f},\mathbf{g}\right)\sim Q_{G}}\left[\mathbf{g}\neq G^{*}\left(\mathbf{f}\right)\right]\right| \\ &\leq \sup_{G^{*}\in\mathcal{H}}\left|\mathbb{E}_{\left(\mathbf{f},\mathbf{g}\right)\sim P_{G}}\left[|\mathbf{g}-G^{*}\left(\mathbf{f}\right)|\neq0\right]-\mathbb{E}_{\left(\mathbf{f},\mathbf{g}\right)\sim Q_{G}}\left[|\mathbf{g}-G^{*}\left(\mathbf{f}\right)|\neq0\right]\right| \end{aligned}$

$$\leqslant \sup_{\delta \in \Delta} \left| \mathbb{E}_{(\mathbf{f}, \mathbf{g}) \sim P_{G}} \left[\delta\left(\mathbf{f}, \mathbf{g}\right) \neq 0 \right] - \mathbb{E}_{(\mathbf{f}, \mathbf{g}) \sim Q_{G}} \left[\delta\left(\mathbf{f}, \mathbf{g}\right) \neq 0 \right] \right|$$

$$\leqslant \sup_{D \in \mathcal{H}_{D}} \left| \mathbb{E}_{(\mathbf{f}, \mathbf{g}) \sim P_{G}} \left[D\left(\mathbf{f}, \mathbf{g}\right) \neq 0 \right] - \mathbb{E}_{(\mathbf{f}, \mathbf{g}) \sim Q_{G}} \left[D\left(\mathbf{f}, \mathbf{g}\right) \neq 0 \right] \right|$$

The upper-bound can be yielded by training the domain discriminator D.



Outline

Transfer Learning

2 Problem I: P(X) ≠ Q(X)
 • DAN: Deep Adaptation Network

3 Problem II: P(Y|X) ≠ Q(Y|X) • CDAN: Conditional Domain Adversarial Network

4 Theoretical Analysis



(日) (同) (三) (三)

Datasets



VisDA Challenge 2017

May 18, 2019 27 / 31

3

Results

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

Method	$A\toW$	$D\toW$	$W\toD$	$A\toD$	$D\toA$	$W\toA$	Avg
AlexNet	$61.6 {\pm} 0.5$	95.4±0.3	99.0±0.2	$63.8{\pm}0.5$	$51.1 {\pm} 0.6$	49.8±0.4	70.1
DAN	$68.5{\pm}0.5$	$96.0 {\pm} 0.3$	99.0±0.3	67.0±0.4	$54.0{\pm}0.5$	$53.1 {\pm} 0.5$	72.9
RTN	73.3±0.3	96.8±0.2	$99.6{\pm}0.1$	$71.0 {\pm} 0.2$	50.5 ± 0.3	$51.0 {\pm} 0.1$	73.7
DANN	$73.0 {\pm} 0.5$	96.4±0.3	99.2±0.3	$72.3 {\pm} 0.3$	$53.4{\pm}0.4$	$51.2 {\pm} 0.5$	74.3
ADDA	$73.5{\pm}0.6$	96.2±0.4	98.8±0.4	$71.6{\pm}0.4$	$54.6{\pm}0.5$	$53.5{\pm}0.6$	74.7
JAN	$74.9 {\pm} 0.3$	96.6±0.2	99.5±0.2	$71.8 {\pm} 0.2$	58.3±0.3	$55.0 {\pm} 0.4$	76.0
CDAN-RM	77.9±0.3	96.9±0.2	100.0 ±.0	74.6 ±0.2	$55.1 {\pm} 0.3$	57.5 ±0.4	77.0
CDAN-M	77.6±0.2	97.2 ±0.1	100.0 ±.0	$73.0{\pm}0.1$	57.3±0.2	$56.1 {\pm} 0.3$	76.9
ResNet-50	68.4±0.2	96.7±0.1	99.3±0.1	68.9±0.2	$62.5 {\pm} 0.3$	60.7±0.3	76.1
DAN	80.5±0.4	97.1±0.2	$99.6{\pm}0.1$	78.6±0.2	$63.6{\pm}0.3$	62.8±0.2	80.4
RTN	84.5±0.2	$96.8 {\pm} 0.1$	$99.4{\pm}0.1$	77.5±0.3	$66.2 {\pm} 0.2$	64.8±0.3	81.6
DANN	82.0±0.4	96.9±0.2	$99.1 {\pm} 0.1$	79.7±0.4	68.2±0.4	67.4±0.5	82.2
ADDA	$86.2 {\pm} 0.5$	96.2±0.3	98.4±0.3	77.8±0.3	$69.5 {\pm} 0.4$	$68.9{\pm}0.5$	82.9
JAN	85.4±0.3	97.4±0.2	99.8±0.2	84.7±0.3	$68.6{\pm}0.3$	$70.0{\pm}0.4$	84.3
JAN-A	92.6±0.2	98.2±0.1	99.8±0.2	86.3±0.1	$71.4{\pm}0.2$	$72.4{\pm}0.1$	86.8
CDAN-RM	93.0±0.2	98.4±0.2	100.0 ±.0	89.2±0.3	70.2±0.4	69.4±0.4	86.7
CDAN-M	93.1 ±0.1	98.6 ±0.1	$\boldsymbol{100.0}{\pm}.0$	93.4 ±0.2	71.0 ±0.3	70.3±0.3	87.7

<ロ> (日) (日) (日) (日) (日)

Mingsheng Long

Results



Mingsheng Long

May 18, 2019 29 / 31

3

<ロ> (日) (日) (日) (日) (日)

Analysis

Ming



Figure: T-SNE on features by (a) ResNet, (b) DANN, (c) CDAN-f, (d) CDAN-fg.



sheng	Long	Transfer Learning	
-------	------	-------------------	--

Open Problems

- Conditional Shift: $P(Y^s | \mathbf{X}^s) \neq Q(Y^t | \mathbf{X}^t)$
- Simulation-to-Real: $P(\mathbf{X}_{\mathsf{low-level}}^s) \neq Q(\mathbf{X}_{\mathsf{low-level}}^t)$
- Open-Set/Zero-Shot (auxiliary info): $\mathbf{Y}^{s} \neq \mathbf{Y}^{t}$
- Heterogeneous (almost impossible): $\mathbf{X}^{s} \neq \mathbf{X}^{t}$
- Learning Transferable Architectures: BN, Skip-connection, etc.
- Xlearn library is available: https://github.com/thuml/Xlearn

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ののの