

# Mind the class weight bias: weighted maximum mean discrepancy for unsupervised domain adaptation

Hongliang Yan<sup>1</sup>, Yukang Ding<sup>1</sup>, Peihua Li<sup>2</sup>, Qilong Wang<sup>2</sup>, Yong Xu<sup>3</sup>, Wangmeng Zuo<sup>1,\*</sup>

Hongliang Yan  
2017/06/21

# Domain Adaptation

## *Problem:*

Training and test sets are related but under different distributions.

## *Methodology:*

- Learn feature space that combine *discriminativeness* and *domain invariance*.

minimize **source error** + **domain discrepancy**

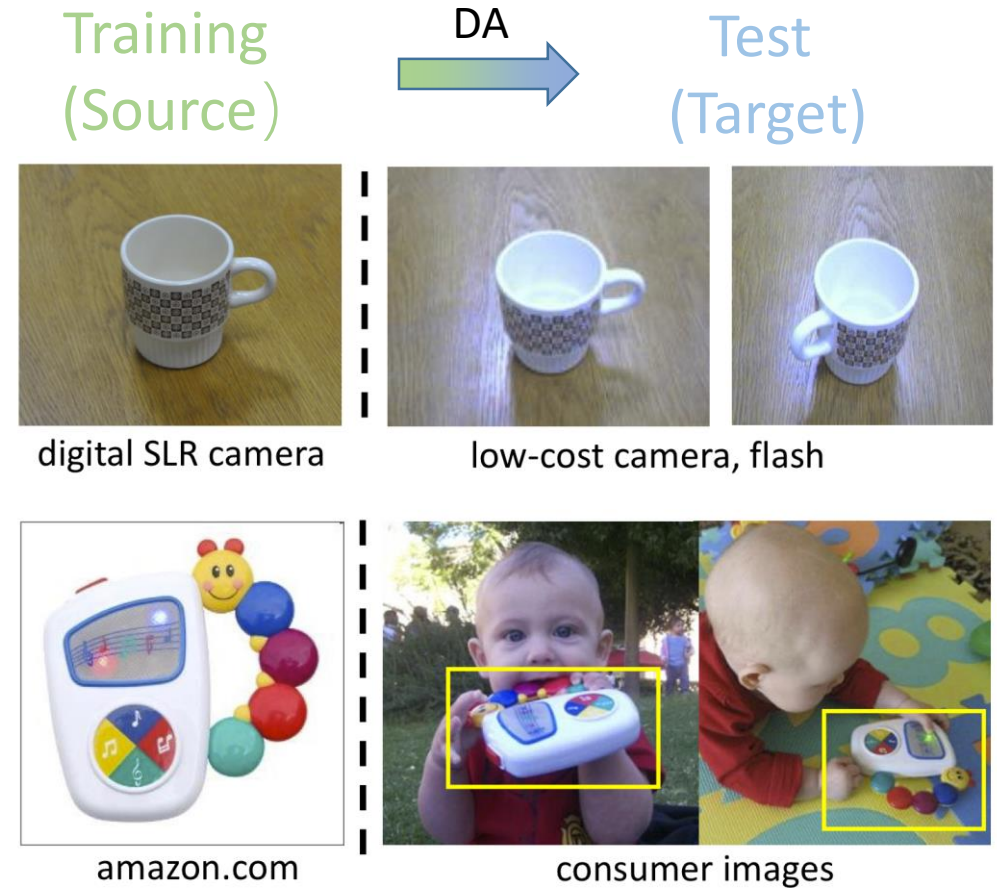


Figure 1. Illustration of dataset bias.

# Maximum Mean Discrepancy (MMD)

- representing distances between distributions as distances between mean embeddings of features

$$\text{MMD}^2(s, t) = \sup_{\|\phi\|_H \leq 1} \|E_{\mathbf{x}^s \sim s}[\phi(\mathbf{x}^s)] - E_{\mathbf{x}^t \sim t}[\phi(\mathbf{x}^t)]\|_H^2$$

- An empirical estimate

$$\text{MMD}^2(\mathbf{D}_s, \mathbf{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$

# Motivation

- *Class weight bias* cross domains remains unsolved but ubiquitous

$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$

# Motivation

- *Class weight bias* cross domains remains unsolved but ubiquitous

$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\left\| \sum_{c=1}^C \underbrace{w_c^s}_{\text{green circle}} E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C \underbrace{w_c^t}_{\text{blue circle}} E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

$$w_c^s = M_c / M \quad \text{and} \quad w_c^t = N_c / N$$

# Motivation

- *Class weight bias* cross domains remains unsolved but ubiquitous

$$\text{MMD}^2(D_s, D_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$

*Effect of class weight bias should be removed:*



① Changes in sample selection criteria

$$\left\| \sum_{c=1}^C w_c^s E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

$$w_c^s = M_c / M \quad \text{and} \quad w_c^t = N_c / N$$

# Motivation

- *Class weight bias* cross domains remains unsolved but ubiquitous

$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^N \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\left\| \sum_{c=1}^C w_c^s E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|$$

$$w_c^s = M_c / M \quad \text{and} \quad w_c^t = N_c / N$$

*Effect of class weight bias should be removed:*

① Changes in sample selection criteria

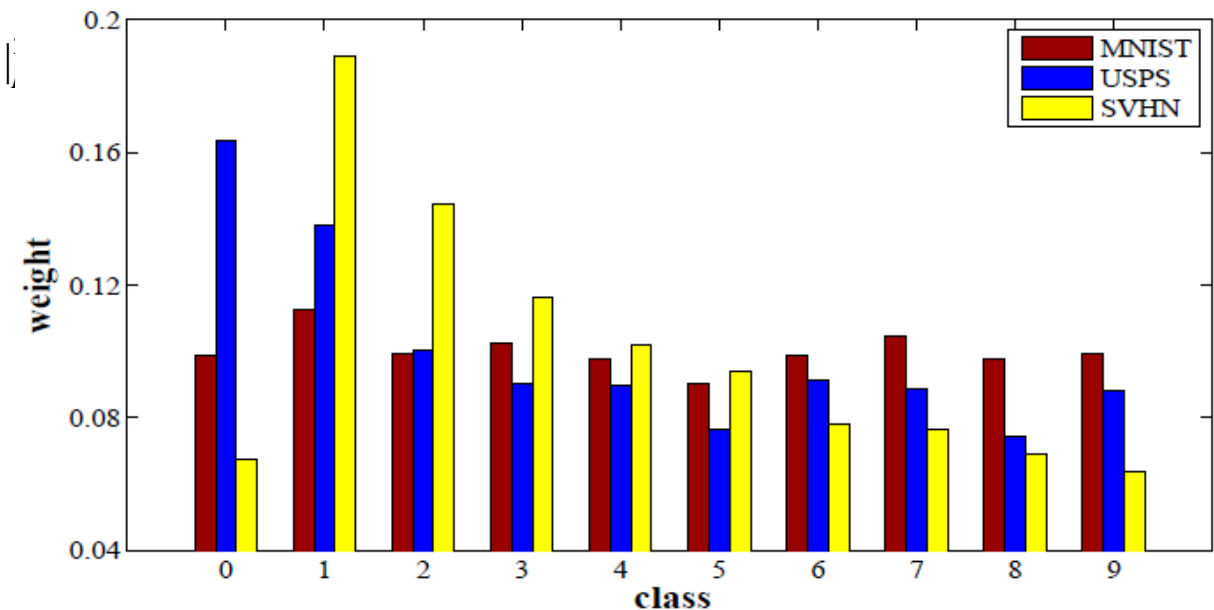


Figure 2. Class prior distribution of three digit recognition datasets.

# Motivation

- *Class weight bias* cross domains remains unsolved but ubiquitous

$$\text{MMD}^2(D_s, D_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\left\| \sum_{c=1}^C w_c^s E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

$$w_c^s = M_c / M \quad \text{and} \quad w_c^t = N_c / N$$

*Effect of class weight bias should be removed:*

- ① Changes in sample selection criteria
- ② Applications are not concerned with class prior distribution



# Motivation

- *Class weight bias* cross domains remains unsolved but ubiquitous

$$\text{MMD}^2(D_s, D_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\left\| \sum_{c=1}^C w_c^s E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

$$w_c^s = M_c / M \quad \text{and} \quad w_c^t = N_c / N$$

*Effect of class weight bias should be removed:*

① Changes in sample selection criteria

② Applications are not concerned with class prior distribution

MMD can be minimized by either learning domain invariant representation or preserving the class weights in source domain.

# Weighted MMD

Main idea: reweighting classes in source domain so that they have the same class weights as target domain

- Introducing an auxiliary weight  $\alpha_c$  for each class  $c$  in source domain

$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\alpha_c = w_c^t / w_c^s$$

$$\left\| \sum_{c=1}^C w_c^s E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

# Weighted MMD

Main idea: reweighting classes in source domain so that they have the same class weights as target domain

- Introducing an auxiliary weight  $\alpha_c$  for each class  $c$  in source domain

$$\text{MMD}^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\left\| \sum_{c=1}^C w_c^s E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

$$\alpha_c = w_c^t / w_c^s$$

$$\text{MMD}_w^2(\mathcal{D}_s, \mathcal{D}_t) = \left\| \frac{1}{M} \sum_{i=1}^M \alpha_{y_i^s} \phi(\mathbf{x}_i^s) - \frac{1}{N} \sum_{j=1}^t \phi(\mathbf{x}_j^t) \right\|_H^2$$



$$\left\| \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^s)) - \sum_{c=1}^C w_c^t E_c(\phi(\mathbf{x}^t)) \right\|_H^2$$

# Weighted DAN

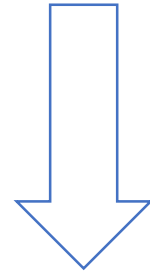
1. Replace MMD with weighted MMD item in DAN[4]:

$$\min_{\mathbf{W}} \frac{1}{M} \sum_{i=1}^M \ell(\mathbf{x}_i^s, \mathbf{y}_i^s; \mathbf{W}) + \lambda \sum_{l \in \{l_1, \dots, l_L\}} \text{MMD}_l(\mathbf{D}_s^l, \mathbf{D}_t^l)$$

# Weighted DAN

1. Replace MMD with weighted MMD item in DAN[4]:

$$\min_{\mathbf{W}} \frac{1}{M} \sum_{i=1}^M \ell(\mathbf{x}_i^s, y_i^s; \mathbf{W}) + \lambda \sum_{l \in \{l_1, \dots, l_L\}} \text{MMD}_l(\mathbf{D}_s^l, \mathbf{D}_t^l)$$



$$\min_{\mathbf{W}, \alpha} \frac{1}{M} \sum_{i=1}^M \ell(\mathbf{x}_i^s, y_i^s; \mathbf{W}) + \lambda \sum_{l \in \{l_1, \dots, l_L\}} \text{MMD}_{l,w}(\mathbf{D}_s^l, \mathbf{D}_t^l)$$

[4] Long M, Cao Y, Wang J. Learning Transferable Features with Deep Adaptation Networks[J]., 2015.

# Weighted DAN

1. Replace MMD with Weighted MMD item in DAN[4]:

$$\min_{\mathbf{W}} \frac{1}{M} \sum_{i=1}^M \ell(\mathbf{x}_i^s, y_i^s; \mathbf{W}) + \lambda \sum_{l \in \{l_1, \dots, l_L\}} \text{MMD}_l(\mathbf{D}_s^l, \mathbf{D}_t^l)$$

2. To further exploit the unlabeled data in target domain, empirical risk is considered as semi-supervised model in [5]:

$$\min_{\mathbf{W}, \{\hat{y}_j\}_{j=1}^N, \alpha} \frac{1}{M} \sum_{i=1}^M \ell(\mathbf{x}_i^s, y_i^s; \mathbf{W}) + \gamma \frac{1}{N} \sum_{j=1}^N \ell(\mathbf{x}_j^t, \hat{y}_j^t; \mathbf{W}) + \lambda \sum_{l \in \{l_1, \dots, l_L\}} \text{MMD}_{l,w}(\mathbf{D}_s^l, \mathbf{D}_t^l)$$

[4] Long M, Cao Y, Wang J. Learning Transferable Features with Deep Adaptation Networks[J]., 2015.

[5] Amini, Massih-Reza, and Patrick Gallinari. "Semi-supervised logistic regression." *Proceedings of the 15th European Conference on Artificial Intelligence*. IOS Press, 2002.

# Optimization: an extension of CEM[6]

Parameters to be estimated including three parts, i.e.,  $\mathbf{W}, \boldsymbol{\alpha}, \{\hat{y}_j^t\}_{j=1}^N$

The model is optimized by alternating between three steps :

- E-step:

Fixed  $\mathbf{W}$ , estimating the class posterior probability  $p(y_j^t = c | \mathbf{x}_j^t)$  of target samples:

$$p(y_j^t = c | \mathbf{x}_j^t) = g(\mathbf{x}_j^t, \mathbf{W})$$

# Optimization: an extension of CEM[6]

Parameters to be estimated including three parts, i.e.,  $\mathbf{W}, \boldsymbol{\alpha}, \{\hat{y}_j^t\}_{j=1}^N$

The model is optimized by alternating between three steps :

- E-step:

Fixed  $\mathbf{W}$ , estimating the class posterior probability  $p(y_j^t = c | \mathbf{x}_j^t)$  of target samples:

$$p(y_j^t = c | \mathbf{x}_j^t) = g(\mathbf{x}_j^t, \mathbf{W})$$

- C-step:

① Assign the pseudo labels  $\{\hat{y}_j^t\}_{j=1}^N$  on target domain:  $\hat{y}_j^t = \arg \max_c p(y_j^t = c | \mathbf{x}_j^t)$

② update the auxiliary class-specific weights  $\boldsymbol{\alpha}$  for source domain:

$$\alpha_c = \hat{w}_c^t / w_c^s \quad \text{where} \quad \hat{w}_c^t = \sum_j \mathbf{1}_c(\hat{y}_j^t) / N$$

$\mathbf{1}_c(x)$  is an indicator function which equals 1 if  $x = c$ , and equals 0 otherwise.

[7] Celeux, Gilles, and Gérard Govaert. "A classification EM algorithm for clustering and two stochastic versions." *Computational statistics & Data analysis* 14.3 (1992): 315-332.



# Optimization: an extension of CEM[6]

Parameters to be estimated including three parts, i.e.,  $\mathbf{W}, \boldsymbol{\alpha}, \{\hat{y}_j^t\}_{j=1}^N$

The model is optimized by alternating between three steps :

- M-step:

Fixed  $\{\hat{y}_j^t\}_{j=1}^N$  and  $\boldsymbol{\alpha}$ , updating  $\mathbf{W}$ . The problem is reformulated as:

$$\min_{\mathbf{W}} \frac{1}{M} \sum_{i=1}^M \ell(\mathbf{x}_i^s, y_i^s; \mathbf{W}) + \gamma \sum_{j=1}^N \ell(\mathbf{x}_i^t, y_i^t; \mathbf{W}) + \lambda \sum_{l \in \{l_1, \dots, l_L\}} \text{MMD}_{l,w}(\mathbf{D}_s^l, \mathbf{D}_t^l)$$

The gradient of the three items is computable and  $\mathbf{W}$  can be optimized by using a mini-batch SGD.

[7] Celeux, Gilles, and Gérard Govaert. "A classification EM algorithm for clustering and two stochastic versions." *Computational statistics & Data analysis* 14.3 (1992): 315-332.

# Experimental results

- Comparison with state-of-the-arts

Method	A→C	W→C	D→C	C→A	C→W	C→D	Avg.
AlexNet [21]	84.0±0.3	77.9±0.4	81.0±0.4	91.3±0.2	83.2±0.3	89.1±0.2	84.0
LapCNN (AlexNet) [38]	83.6±0.6	77.8±0.5	80.6±0.4	92.1±0.3	81.6±0.4	87.8±0.4	83.9
DDC (AlexNet) [36]	84.3±0.5	76.9±0.4	80.5±0.2	91.3±0.3	85.5±0.3	89.1±0.3	84.6
DAN (AlexNet) [25]	86.0±0.5	81.5±0.3	82.0±0.4	92.0±0.3	92.6±0.4	90.5±0.2	87.3
WDAN (AlexNet)	<b>86.9±0.1</b>	<b>84.1±0.2</b>	<b>83.9±0.1</b>	<b>93.1±0.2</b>	<b>93.6±0.2</b>	<b>93.4±0.2</b>	<b>89.2</b>
WDAN* (AlexNet)	87.1±0.2	85.1±0.3	85.2±0.2	93.2±0.1	93.5±0.3	94.5±0.2	89.8
GoogLeNet [34]	91.3±0.2	88.2±0.3	88.9±0.3	95.2±0.1	92.5±0.2	94.7±0.3	91.8
DDC (GoogLeNet) [36]	91.4±0.2	88.7±0.3	89.0±0.4	95.3±0.2	93.0±0.1	94.9±0.4	92.1
DAN (GoogLeNet) [25]	91.4±0.3	89.7±0.2	89.1±0.4	95.5±0.2	93.1±0.3	95.3±0.1	92.3
WDAN (GoogLeNet)	<b>92.2±0.2</b>	<b>91.0±0.5</b>	<b>89.8±0.3</b>	<b>95.5±0.3</b>	<b>95.4±0.2</b>	<b>95.5±0.5</b>	<b>93.2</b>
VGGnet-16 [32]	89.6±0.4	88.1±0.4	85.4±0.5	93.7±0.2	94.3±0.2	93.7±0.2	90.8
DAN (VGGnet-16) [25]	91.2±0.2	90.6±0.3	87.1±0.4	95.7±0.2	95.3±0.3	94.7±0.1	92.4
WDAN (VGGnet-16)	<b>91.4±0.2</b>	<b>91.0±0.2</b>	<b>89.0±0.3</b>	<b>95.7±0.1</b>	<b>95.8±0.2</b>	<b>95.9±0.3</b>	<b>93.1</b>

Table 1. Experimental results on office-10+Caltech-10

# Experimental results

- Empirical analysis

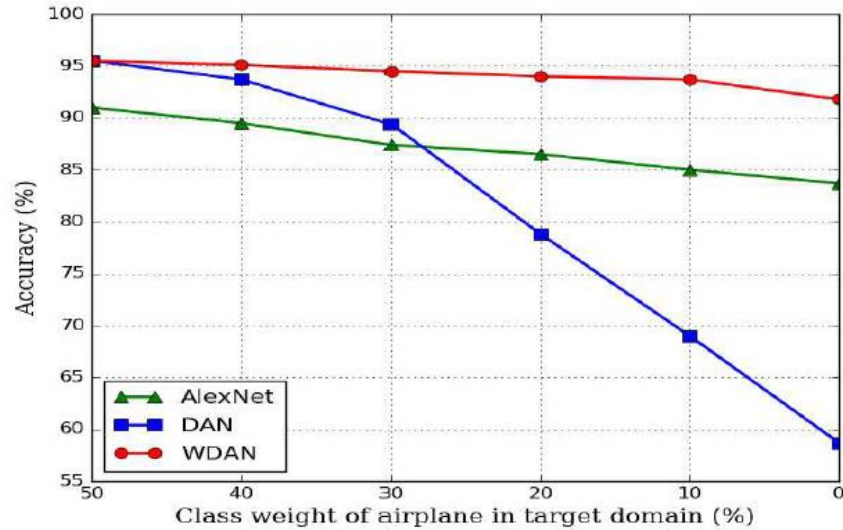


Figure 3. Performance of various model under different class weight bias.

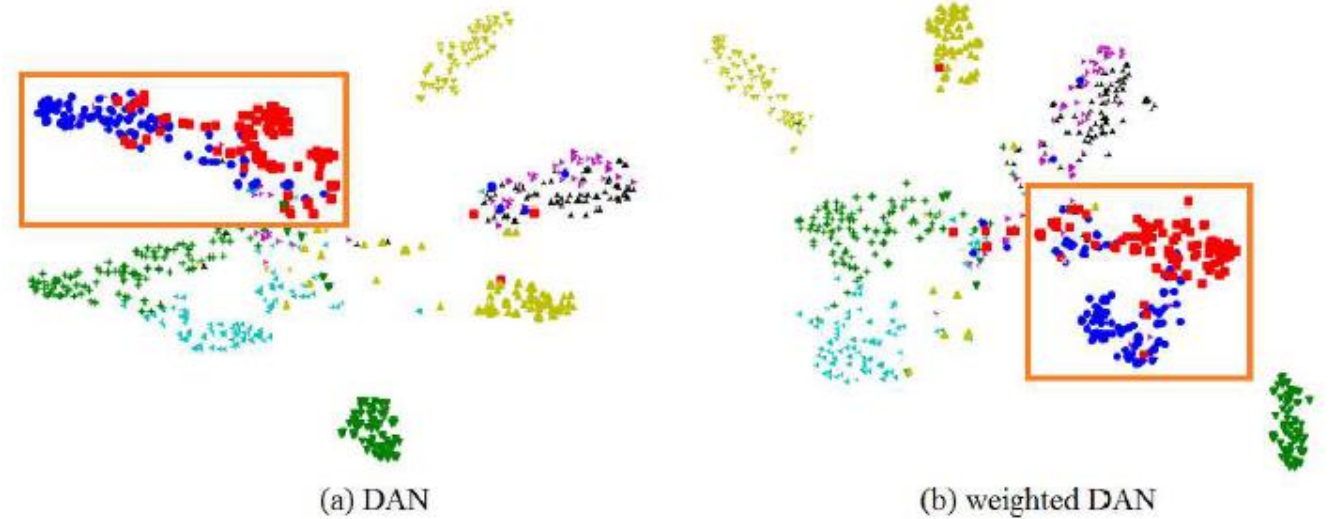


Figure 4. Visualization of the learned features of DAN and weighted DAN.

# Summary

- Introduce class-specific weight into MMD to reduce the effect of class weight bias cross domains.
- Develop WDAN model and optimize it in an CEM framework.
- Weighted MMD can be applied to other scenarios where MMD is used for distribution distance measurement, e.g., image generation

# Thanks!

Paper & code are available

Paper: <https://arxiv.org/abs/1705.00609>

Code: <https://github.com/yhldhit/WMMMD-Caffe>