

# Transferable Calibration with Lower Bias and Variance in Domain Adaptation

#### Summary

- A Transferable Calibration (TransCal) method, achieving more accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework.
- ► A dilemma in the open problem of Calibration in DA: existing domain adaptation models learn higher classification accuracy at the expense of well-calibrated probabilities.
- Extensive experiments on various DA methods, datasets, and calibration metrics, while the effectiveness of our method has been justified both theoretically and empirically.
- Code available @ github.com/thuml/TransCal

## **Domain Adaptation (DA)**

▶ Deep learning across domains:  $(P \neq Q)$ Non independent and identically distributed distributions (Non-IID)



## Mainstream Approaches to DA

- Numerous deep DA methods can be mainly grouped into two categories: moment matching and adversarial training.
- Most of DA methods focus on improving the accuracy in the target domain but fail to estimate the predictive uncertainty, falling short of a miscalibration problem.





(a) Moment Matching: DAN

(b) Adversarial Training: DANN

## School of Software - Tsinghua University - China

Ximei Wang, Mingsheng Long ( $\boxtimes$ ), Jianmin Wang, and Michael I. Jordan<sup>‡</sup>

School of Software, KLiss, BNRist, Tsinghua University <sup>#</sup>University of California, Berkeley

### **Confidence Calibration in Deep Learning**

A model should output a prediction probability reflecting the true frequency of an event:

$$\mathbb{P}(\widehat{Y} = Y | \widehat{P} = c) = c, \forall c \in [0, 1]$$
(1)

where  $\widehat{Y}$  is the class prediction and  $\widehat{P}$  is its confidence. ► DNNs learn high accuracy at the cost of over-confidence.



#### Calibration Metric

Expected Calibration Error (ECE)  $\mathcal{L}_{ ext{ECE}} = \sum_{m=1}^{B} \frac{|B_m|}{n} |\mathbb{A}(B_m) - \mathbb{C}(B_m)|$  $\mathbb{A}(B_m) = |B_m|^{-1} \sum \mathbb{1}(\widehat{y}_i = y_i) \quad (Accuracy)$ (2)  $\mathbb{C}(B_m) = |B_m|^{-1} \sum_{i} \max_{i} p(\widehat{y}_i^k | x_i, \theta) \quad (Confidence)$ 

### **Temperature Scaling for IID Calibration**

- Fix the neural model trained on the training set  $\mathcal{D}_{tr}$
- $\blacktriangleright$  Attain the optimal temperature  $T^*$  by minimizing

$$T^* = \arg\min_{\tau} E_{(x_v, y_v) \in \mathcal{D}_v} \mathcal{L}_{NLL}(\sigma(z_v/T), y_v)$$
(3)

 $\sigma$  is the *softmax* function,  $\mathcal{L}_{NLL}$  is Negative Log-Likelihood. Transform  $z_{te}$  into calibrated probabilities  $p_{te} = \sigma(z_{te}/T^*)$ .

## Dilemma of Accuracy vs Confidence in DA

► DA models yield high acc at the cost of poorly-calibration.



Calibration in DA is challenging due to the existence of domain shift and the lack of target label

### **Transferable Calibration Framework**

Estimate the target ECE by importance weighting

$$E_{\mathbf{x}\sim q} \left[ \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y}) \right] = \int_{q} \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y}) q(\mathbf{x}) d\mathbf{x}$$
$$= \int_{p} \frac{q(\mathbf{x})}{p(\mathbf{x})} \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y}) p(\mathbf{x}) d\mathbf{x} = E_{\mathbf{x}\sim p} \left[ w(\mathbf{x}) \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y}) \right],$$
(4)

Estimate density ratio from a logistic regression classifier  $\widehat{w}(\mathbf{x}) = \frac{q(\mathbf{x})}{p(\mathbf{x})} = \frac{v(\mathbf{x}|d=0)}{v(\mathbf{x}|d=1)} = \frac{P(d=1)P(d=0|\mathbf{x})}{P(d=0)P(d=1|\mathbf{x})}, \quad (5)$ 

## Transferable Calibration: Bias Reduction

Bias between the estimated ECE and the ground-truth one  $\left|\mathsf{E}_{\mathsf{x}\sim q}\left[\mathcal{L}_{\mathrm{ECE}}^{\widehat{w}(\mathsf{x})}\right] - \mathsf{E}_{\mathsf{x}\sim q}\left[\mathcal{L}_{\mathrm{ECE}}^{w(\mathsf{x})}\right]\right|$  $= |\mathsf{E}_{\mathsf{x}\sim p} \left[ \widehat{w}(\mathsf{x}) \mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \mathsf{y}) \right] - \mathsf{E}_{\mathsf{x}\sim p} \left[ w(\mathsf{x}) \mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \mathsf{y}) \right]|$  $= |\mathsf{E}_{\mathsf{x}\sim \rho} \left[ (\mathsf{w}(\mathsf{x}) - \widehat{\mathsf{w}}(\mathsf{x})) \mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \mathsf{y}) \right]|.$ (6)

## The discrepancy between $\widehat{w}(x)$ and w(x) can be bounded by $\mathbb{E}_{\mathsf{x}\sim p}\left[\left(w(\mathsf{x})-\widehat{w}(\mathsf{x})\right)^{2}\right] \leq (M+1)^{4}\mathsf{E}_{\mathsf{x}\sim p}\left[\left(P(d=1|\mathsf{x})-\widehat{P}(d=1|\mathsf{x})\right)^{2}\right].$

▶ Use  $\lambda$  ( $0 \le \lambda \le 1$ ) to control the bound M of  $\widehat{w}(x)$  $T^* = \arg\min \mathsf{E}_{\mathsf{x}_v \sim p} \left[ \widetilde{w}(\mathsf{x}_v) \mathcal{L}_{\mathrm{ECE}}(\sigma(\phi(\mathsf{x}_v)/T), \mathsf{y}) \right], \quad \widetilde{w}(\mathsf{x}_v^i) = \left[ \widehat{w}(\mathsf{x}_v^i) \right]^{\lambda}.$ (8)

## National Engineering Lab for Big Data System Software



## **Transferable Calibration: Variance Reduction**

Serial Control Variate: 
$$\operatorname{Var}[u^{**}] \leq \operatorname{Var}[u^*] \leq \operatorname{Var}[u]$$

$$u^* = u + \eta_1(t_1 - \tau_1)$$

$$u^{**} = u^* + \eta_2(t_2 - \tau_2)$$
(9)

First, use importance weight  $\widetilde{w}(x_s)$  as a control covariate

$$\mathbb{E}_{q}^{*}(\widehat{\mathbf{y}},\mathbf{y}) = \widetilde{\mathbb{E}}_{q}(\widehat{\mathbf{y}},\mathbf{y}) - \frac{1}{n_{s}} \frac{\operatorname{Cov}(\mathcal{L}_{\mathrm{ECE}}^{\widetilde{w}},\widetilde{w}(\mathbf{x}))}{\operatorname{Var}[\widetilde{w}(\mathbf{x})]} \sum_{i=1}^{n_{s}} [\widetilde{w}(\mathbf{x}_{s}^{i}) - 1].$$
(10)

Second, use the prediction correctness  $r(x_s)$  as another one

$$\mathbb{E}_{q}^{**}(\widehat{\mathbf{y}},\mathbf{y}) = \mathbb{E}_{q}^{*}(\widehat{\mathbf{y}},\mathbf{y}) - \frac{1}{n_{s}} \frac{\operatorname{Cov}(\mathcal{L}_{\mathrm{ECE}}^{\widetilde{w}*},r(\mathbf{x}))}{\operatorname{Var}[r(\mathbf{x})]} \sum_{i=1}^{n_{s}} [r(\mathbf{x}_{s}^{i}) - c],$$
(11)

#### **Experiments and Results**

Table 2: ECE (%) vs. Acc (%) via various calibration methods on *Office-Home* with CDAN

Metric	Cal. Method	A→C	$A \rightarrow P$	$A \rightarrow R$	C→A	$C \rightarrow P$	$C \rightarrow R$	$R \rightarrow A$	$R \rightarrow C$	$R \rightarrow P$	Avg
Acc	Before Cal. MC-dropout [12] TransCal (ours)	49.4 47.2 49.4	68.4 66.2 68.4	75.5 71.4 75.5	57.6 57.1 57.6	70.1 65.7 70.1	70.4 70.6 70.4	68.9 68.3 68.9	54.4 53.6 54.4	81.2 80.7 81.2	68.3 66.7 68.3
ECE	Before Cal. MC-dropout [12] Matrix Scaling Vector Scaling Temp. Scaling CPCS [38]	40.2 33.1 44.7 34.7 28.3 35.0	26.4 21.3 28.8 18.0 17.6 29.4	17.8 15.0 19.7 11.3 10.1 8.3	35.8 24.2 36.1 23.4 <b>21.2</b> <u>21.3</u>	$23.5 \\ 20.5 \\ 25.4 \\ 15.4 \\ \underline{13.2} \\ 29.0$	<ul> <li>21.9</li> <li>13.2</li> <li>24.1</li> <li>11.5</li> <li>8.2</li> <li><b>5.6</b></li> </ul>	24.8 25.6 38.1 27.3 26.0 <b>19.9</b>	36.4 14.2 15.7 8.5 8.8 9.1	14.5 22.4 29.5 20.0 18.1 20.3	26.8 19.6 29.1 18.9 16.8 19.8
	TransCal (w/o Bias) TransCal (w/o Variance) TransCal (ours) Oracle	<b>21.7</b> 31.2 22.9 5.8	<u>10.8</u> 16.4 <b>9.3</b> 8.1	<u>5.8</u> 6.5 <b>5.1</b> 4.8	27.6 31.1 21.7 10.0	<b>9.2</b> 14.7 14.0 7.7	<u>6.0</u> 16.1 6.4 4.2	27.4 27.5 <u>21.6</u> 5.5	5.2 4.1 4.5 3.9	16.9         20.0         15.6         6.2	$   \begin{vmatrix}     14.5 \\     18.6 \\     13.5 \\     6.2   \end{vmatrix} $





Figure 3: The estimated calibration error with respect to different values of temperature T and meta parameter  $\lambda$  (both are *learnable*), showing that different models achieve optimal values at different  $\lambda$ .

Mail: wxm17@mails.tsinghua.edu.cn