Transferable Calibration with Lower Bias and Variance in Domain Adaptation

Ximei Wang, Mingsheng Long (⊠), Jianmin Wang, and Michael I. Jordan[♯]

School of Software, Tsinghua University National Engineering Laboratory for Big Data Software [#]University of California, Berkeley

wxm17@mails.tsinghua.edu.cn https://wxm17.github.io/ Neural Information Processing Systems (NeurIPS), 2020

Domain Adaptation (DA)

Transfer from a labeled source domain to an unlabeled target one.



Domain Adaptation Approaches

- Moment Matching
- Deep Adaptation Netwotk etc.

- Adversarial Training
- Domain Adversarial Neural Network etc.



Confidence Calibration in Deep Learning¹

• A model should output a probability reflecting the true frequency:

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

where \widehat{Y} is the class prediction and \widehat{P} is its associated confidence.

• Deep networks learn high accuracy at the expense of over-confidence.



Temperature Scaling for IID Calibration

• Calibration Metric: Expected Calibration Error (ECE)

$$\mathcal{L}_{\text{ECE}} = \sum_{m=1}^{B} \frac{|B_m|}{n} |\mathbb{A}(B_m) - \mathbb{C}(B_m)|$$

$$\mathbb{A}(B_m) = |B_m|^{-1} \sum_{i \in B_m} \mathbb{1}(\widehat{y}_i = y_i) \quad \text{(Accuracy)}$$

$$\mathbb{C}(B_m) = |B_m|^{-1} \sum_{i \in B_m} \max_k p(\widehat{y}_i^k | x_i, \theta) \quad \text{(Confidence)}$$
(2)

• IID Calibration: Temperature Scaling

$$T^* = \underset{T}{\arg\min} \ \mathsf{E}_{(\mathsf{x}_{\mathsf{v}},\mathsf{y}_{\mathsf{v}})\in\mathcal{D}_{\mathsf{v}}} \ \mathcal{L}_{\mathrm{NLL}} \left(\sigma(\mathsf{z}_{\mathsf{v}}/T),\mathsf{y}_{\mathsf{v}}\right)$$
(3)

 σ is the softmax function, $\mathcal{L}_{\rm NLL}$ is Negative Log-Likelihood loss.

• Transform logits z_{te} into calibrated probabilities $p_{te} = \sigma(z_{te}/T^*)$.

Dilemma of Accuracy vs Confidence in DA

• Transfer models yield high accuracy at the expense of over-confidence.



- Calibration in transfer learning is challenging due to the coexistence:
 - Domain shift ECE should be unbiased to the target domain
 - ${\scriptstyle \bullet }$ Unlabeled target ECE on the target domain is incomputable
- Bias-Variance-Shift Dilemma of confidence calibration in Transfer Learning

Transferable Calibration Framework

• $\mathsf{E}_{\mathsf{x}\sim p}\left[w(\mathsf{x})\mathcal{L}_{(\cdot)}(\phi(\mathsf{x}), \mathsf{y})\right]$ is an **unbiased** estimator of the target calibration error \mathbb{E}_q

$$\begin{aligned} \mathsf{E}_{\mathsf{x}\sim q} \left[\mathcal{L}_{(\cdot)}(\phi(\mathsf{x}), \mathsf{y}) \right] &= \int_{q} \mathcal{L}_{(\cdot)}(\phi(\mathsf{x}), \mathsf{y}) q(\mathsf{x}) \mathrm{d}\mathsf{x} \\ &= \int_{p} \frac{q(\mathsf{x})}{p(\mathsf{x})} \mathcal{L}_{(\cdot)}(\phi(\mathsf{x}), \mathsf{y}) p(\mathsf{x}) \mathrm{d}\mathsf{x} = \mathsf{E}_{\mathsf{x}\sim p} \left[w(\mathsf{x}) \mathcal{L}_{(\cdot)}(\phi(\mathsf{x}), \mathsf{y}) \right] \end{aligned}$$
(4)

- Discriminative density ratio estimation method: LogReg
- Use Bayesian formula to derive $\widehat{w}(x)$ from a logistic regression classifier

$$\widehat{w}(x) = \frac{q(x)}{p(x)} = \frac{v(x|d=0)}{v(x|d=1)} = \frac{P(d=1)}{P(d=0)} \frac{P(d=0|x)}{P(d=1|x)}$$
(5)

Transferable Calibration: Bias Reduction

- Importance-weighting for an unbiased estimate of target ECE if $\widehat{w}(x) = w(x)$
- The bias between the estimated ECE and the ground-truth ECE

$$\begin{aligned} \left| \mathsf{E}_{\mathsf{x}\sim q} \left[\mathcal{L}_{\mathrm{ECE}}^{\widehat{w}(\mathsf{x})} \right] - \mathsf{E}_{\mathsf{x}\sim q} \left[\mathcal{L}_{\mathrm{ECE}}^{w(\mathsf{x})} \right] \right| \\ &= \left| \mathsf{E}_{\mathsf{x}\sim p} \left[\widehat{w}(\mathsf{x}) \mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \ \mathsf{y}) \right] - \mathsf{E}_{\mathsf{x}\sim p} \left[w(\mathsf{x}) \mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \ \mathsf{y}) \right] \right| \\ &= \left| \mathsf{E}_{\mathsf{x}\sim p} \left[(w(\mathsf{x}) - \widehat{w}(\mathsf{x})) \mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \ \mathsf{y}) \right] \right|. \end{aligned}$$
(6)

• The bias of them can be further bounded by

$$\begin{aligned} &|\mathsf{E}_{\mathsf{x}\sim\rho}\left[\left(w(\mathsf{x})-\widehat{w}(\mathsf{x})\right)\mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \mathsf{y})\right]| \\ &\leq \sqrt{\mathsf{E}_{\mathsf{x}\sim\rho}\left[\left(w(\mathsf{x})-\widehat{w}(\mathsf{x})\right)^{2}\right]\mathsf{E}_{\mathsf{x}\sim\rho}\left[\left(\mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \mathsf{y})\right)^{2}\right]} & (\mathrm{Cachy}-\mathrm{Schwarz\ Inequality}) \quad (7) \\ &\leq \frac{1}{2}\left(\mathsf{E}_{\mathsf{x}\sim\rho}\left[\left(w(\mathsf{x})-\widehat{w}(\mathsf{x})\right)^{2}\right]+\mathsf{E}_{\mathsf{x}\sim\rho}\left[\left(\mathcal{L}_{\mathrm{ECE}}(\phi(\mathsf{x}), \mathsf{y})\right)^{2}\right]\right) & (\mathrm{AM}/\mathrm{GM\ Inequality}) \end{aligned}$$

Transferable Calibration: Bias Reduction

• For any x s.t. $P(d = 1|x) \neq 0$, the following inequality holds:

$$\frac{1}{M+1} \le P(d=1|\mathsf{x}) \le 1, \quad \text{since } w(\mathsf{x}) = \frac{P(d=0|\mathsf{x})}{P(d=1|\mathsf{x})} = \frac{1-P(d=1|\mathsf{x})}{P(d=1|\mathsf{x})} = \frac{1}{P(d=1|\mathsf{x})} - 1.$$
(8)

 \bullet The discrepancy between $\widehat{w}(\mathbf{x})$ and $w(\mathbf{x})$ can be bounded by

$$\mathbb{E}_{\mathbf{x}\sim\rho}\left[\left(w(\mathbf{x})-\widehat{w}(\mathbf{x})\right)^{2}\right] = \mathsf{E}_{\mathbf{x}\sim\rho}\left[\left(\frac{P(d=1|\mathbf{x})-\widehat{P}(d=1|\mathbf{x})}{P(d=1|\mathbf{x})\widehat{P}(d=1|\mathbf{x})}\right)^{2}\right]$$
$$\leq (M+1)^{4}\mathsf{E}_{\mathbf{x}\sim\rho}\left[\left(P(d=1|\mathbf{x})-\widehat{P}(d=1|\mathbf{x})\right)^{2}\right].$$
(9)

• Use λ (0 $\leq \lambda \leq$ 1) to control the bound *M* of the importance weights

$$T^* = \arg\min_{T,\lambda} \mathsf{E}_{\mathsf{x}_{\mathsf{v}} \sim p} \left[\widetilde{w}(\mathsf{x}_{\mathsf{v}}) \mathcal{L}_{\mathrm{ECE}}(\sigma(\phi(\mathsf{x}_{\mathsf{v}})/T), \mathsf{y}) \right], \quad \widetilde{w}(\mathsf{x}_{\mathsf{v}}^{i}) = \left[\widehat{w}(\mathsf{x}_{\mathsf{v}}^{i}) \right]^{\lambda}.$$
(10)

Control Variate Method

• (a) Feature adaptation reduces distribution discrepancy $d_{\alpha+1}(q||p)$

$$\begin{aligned} \operatorname{Var}_{\mathsf{x}\sim\rho}\left[\mathcal{L}_{\mathrm{ECE}}^{\mathsf{w}}\right] &= \mathsf{E}_{\mathsf{x}\sim\rho}\left[\left(\mathcal{L}_{\mathrm{ECE}}^{\mathsf{w}}\right)^{2}\right] - (\mathsf{E}_{\mathsf{x}\sim\rho}\left[\mathcal{L}_{\mathrm{ECE}}^{\mathsf{w}}\right])^{2} \\ &\leq d_{\alpha+1}(q||\rho)(\mathsf{E}_{\mathsf{x}\sim\rho}\mathcal{L}_{\mathrm{ECE}}^{\mathsf{w}})^{1-\frac{1}{\alpha}} - (\mathsf{E}_{\mathsf{x}\sim\rho}\mathcal{L}_{\mathrm{ECE}}^{\mathsf{w}})^{2}, \quad \forall \alpha > 0. \end{aligned}$$
(11)

• (b) Control variate explicitly reduces the variance ²

²Lemieux. Control variates. In Wiley StatsRef: Statistics Reference Online, American Cancer Society, 2017.

Transferable Calibration: Variance Reduction

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

$$u^{*} = u + \eta_{1}(t_{1} - \tau_{1})$$

$$u^{**} = u^{*} + \eta_{2}(t_{2} - \tau_{2})$$
(12)

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

$$\mathbb{E}_{q}^{*}(\widehat{\mathbf{y}},\mathbf{y}) = \widetilde{\mathbb{E}}_{q}(\widehat{\mathbf{y}},\mathbf{y}) - \frac{1}{n_{s}} \frac{\operatorname{Cov}(\mathcal{L}_{\mathrm{ECE}}^{\widetilde{w}},\widetilde{w}(\mathbf{x}))}{\operatorname{Var}[\widetilde{w}(\mathbf{x})]} \sum_{i=1}^{n_{s}} [\widetilde{w}(\mathbf{x}_{s}^{i}) - 1].$$
(13)

• Second, use the prediction correctness $r(x_s)$ as another control variate

$$\mathbb{E}_{q}^{**}(\widehat{\mathbf{y}},\mathbf{y}) = \mathbb{E}_{q}^{*}(\widehat{\mathbf{y}},\mathbf{y}) - \frac{1}{n_{s}} \frac{\operatorname{Cov}(\mathcal{L}_{\mathrm{ECE}}^{\widetilde{w}*},r(\mathbf{x}))}{\operatorname{Var}[r(\mathbf{x})]} \sum_{i=1}^{n_{s}} [r(\mathbf{x}_{s}^{i}) - c],$$
(14)

• Reduce bias, variance, and shift all-in-one for Transferable Calibration

Ximei Wang

TransCal Algorithm

Algorithm 1 Transferable Calibration in Domain Adaptation

- 1: Input: Labeled source dataset $S = \{ (\mathbf{x}_s^i, \mathbf{y}_s^i) \}_{i=1}^{n_s}$ and unlabled target dataset $\mathcal{T} = \{ (\mathbf{x}_t^i) \}_{i=1}^{n_t}$
- 2: **Parameter:** Temperature T and learnable meta parameter λ
- 3: Partition \mathcal{S} into $\mathcal{S}_{tr} = \{ (\mathbf{x}_{tr}^i, \mathbf{y}_{tr}^i) \}_{i=1}^{n_{tr}}$ and $\mathcal{S}_v = \{ (\mathbf{x}_v^i, \mathbf{y}_v^i) \}_{i=1}^{n_v}$.
- 4: Train a DA model $\phi(\mathbf{x}) = G(F(\mathbf{x}))$ on \mathcal{S}_{tr} and \mathcal{T} via any DA method until convergy
- 5: Randomly upsample the source or the target dataset to make $n_{tr} = n_t$
- 6: Fix the DA model and compute features $\mathcal{F}_{tr} = \{f_{tr}^i\}_{i=1}^{n_{tr}}, \mathcal{F}_v = \{f_v^i\}_{i=1}^{n_v}, \mathcal{F}_t = \{f_t^i\}_{i=1}^{n_t}$
- 7: Train a logistic regression model H to discriminate the features \mathcal{F}_{tr} and \mathcal{F}_t until converge
- 8: Compute $\widehat{w}(\mathbf{x}_v^i) = \left[1 H(f_v^i)\right] / H(f_v^i)$ and $\widetilde{w}(\mathbf{x}_v^i) = \left[\widehat{w}(\mathbf{x}_v^i)\right]^{\lambda}$
- 9: Compute $\mathsf{E}_{\mathbf{x}\sim p}\mathcal{L}^{\widetilde{w}}_{\mathrm{ECE}}, \mathbb{E}^*_a(\widehat{\mathbf{y}}, \mathbf{y})$ and $\mathbb{E}^{**}_a(\widehat{\mathbf{y}}, \mathbf{y})$ as in Eq. 9, Eq. 11 and Eq. 13 respectively
- 10: Jointly optimize the transferable calibration objective as $T^* = \arg \min \mathbb{E}_a^{**} (\sigma(\phi(\mathbf{x}_v)/T), \mathbf{y}_v))$ 11: Calibrate the logit vectors on the target domain by $\hat{\mathbf{y}}_t = \sigma(\phi(\mathbf{x}_t)/T^*)$

- 3

Experiments and Results

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

Metric	Cal. Method	A→C	$A \rightarrow P$	$A \rightarrow R$	$C {\rightarrow} A$	$C \rightarrow P$	$C \rightarrow R$	$R \rightarrow A$	$R \rightarrow C$	$R \rightarrow P$	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 21.2 <u>21.3</u>	23.5 20.5 25.4 15.4 <u>13.2</u> 29.0	21.9 13.2 24.1 11.5 8.2 5.6	24.8 25.6 38.1 27.3 26.0 19.9	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	21.7 31.2 <u>22.9</u> 5.8	10.8 16.4 9.3 8.1	5.8 6.5 5.1 4.8	27.6 31.1 21.7 10.0	9.2 14.7 14.0 7.7	6.0 16.1 6.4 4.2	27.4 27.5 <u>21.6</u> 5.5	5.2 4.1 <u>4.5</u> 3.9	16.9 20.0 15.6 6.2	$\frac{ \underline{14.5} }{ \underline{18.6} }$

Experiments and Results



Figure 2: Reliability diagrams from Clipart to Product with CDAN [25] before and after calibration.



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

< A

Summary

- 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.
- A Transferable Calibration (TransCal) method, achieving more accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework.
- Extensive experiments on various DA methods, datasets, and calibration metrics, while the effectiveness of our method has been justified both theoretically and empirically.
- Code will be available @ github.com/thuml/TransCal

Future Work

- 1. Design DA methods based on our ECE-Accuracy-Dilemma observation
- 2. TransCal may still fall short under the following circumstances:
- The domain gap is extremely large even after applying domain adaptation methods
- The source or the target dataset is too small to estimate importance weights
- TransCal is based on the covariate shift assumption and it remains unclear whether it can still perform well under label shift, especially when we meet with a long-tailed distribution.