

A Mathematical Framework for Quantifying Transferability in Multi-source Transfer Learning

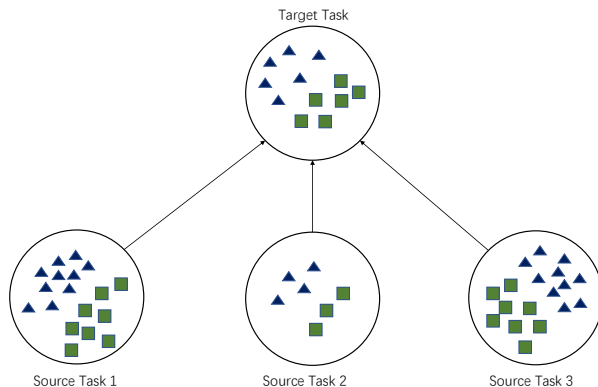
Xinyi Tong*, Xiangxiang Xu[†], Shao-Lun Huang*, Lizhong Zheng[†]

*Tsinghua-Berkeley Shenzhen Institute, Tsinghua University

[†]Massachusetts Institute of Technology

Tuesday 19th October, 2021

Which source task is more helpful for the target task?



We can intuitively claim that source task 1 helps more, i.e.,

Source task 1 has higher transferability.

Intuitively, the reasons are

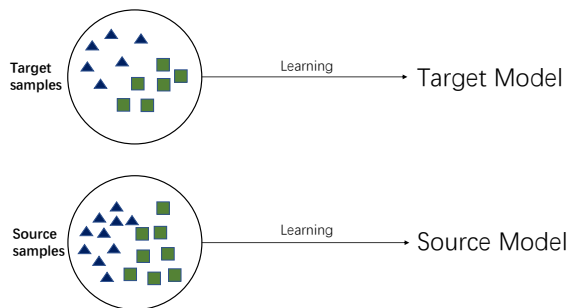
- Source task 2: less samples \rightarrow Sample size is important
- Source task 3: not like target task \rightarrow Similarity is important

What we hope to answer in this paper:

- Establish a mathematical framework to analyze transferability.
- Interpret all these factors.
- Apply the theoretical analyses to practical tasks.

Problem Formulation

First, how to establish the framework for transfer learning?





The final model using the source knowledge is **Optimal α ?**

$$(1 - \alpha) \cdot \text{Target Model} + \alpha \cdot \text{Source Model}$$

Discrete Case

Sample $x \in \mathcal{X}$ & Label $y \in \mathcal{Y}$

We hope to learn the target distribution $P_{XY}^{(0)}$

- Target samples $\{(x_\ell^{(0)}, y_\ell^{(0)})\}_{\ell=1}^{n_0}$  i.i.d from $P_{XY}^{(0)} \rightarrow \hat{P}_{XY}^{(0)}$
- Source samples $\{(x_\ell^{(1)}, y_\ell^{(1)})\}_{\ell=1}^{n_1}$  i.i.d from $P_{XY}^{(1)} \rightarrow \hat{P}_{XY}^{(1)}$

We use $(1 - \alpha)\hat{P}_{XY}^{(0)} + \alpha\hat{P}_{XY}^{(1)}$ to estimate $P_{XY}^{(0)}$.

Testing loss

How to evaluate the model? By the test data!

$$L_{\text{test}}^{(\alpha)} \triangleq \mathbb{E} \left[\chi^2 \left(P_{XY}^{(0)}, (1 - \alpha) \hat{P}_{XY}^{(0)} + \alpha \hat{P}_{XY}^{(1)} \right) \right] \quad (1)$$

Why not Log-loss?

What is optimal coefficient? \rightarrow **Transferability**

$$\alpha^* = \arg \min_{\alpha} L_{\text{test}}^{(\alpha)} \quad (2)$$

Testing loss

Theorem

The testing loss can be computed as

$$L_{\text{test}}^{(\alpha)} = \alpha^2 \chi^2 \left(P_{XY}^{(0)}, P_{XY}^{(1)} \right) + \frac{(1 - \alpha)^2}{n_0} V^{(0)} + \frac{\alpha^2}{n_1} V^{(1)}, \quad (3)$$

and the optimal α^* is

$$\alpha^* = \frac{\frac{1}{n_0} V^{(0)}}{\chi^2(P_{XY}^{(0)}, P_{XY}^{(1)}) + \frac{1}{n_0} V^{(0)} + \frac{1}{n_1} V^{(1)}}, \quad (4)$$

where $V^{(0)} = |\mathcal{X}||\mathcal{Y}| - 1$ and $V^{(1)} = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{P_{XY}^{(1)}(x, y) (1 - P_{XY}^{(1)}(x, y))}{P_{XY}^{(0)}(x, y)}.$

We claim that α^* is our **transferability measure**.




The affecting factors

- $\chi^2 \left(P_{XY}^{(0)}, P_{XY}^{(1)} \right) \rightarrow$ distance!
- n_0 & $n_1 \rightarrow$ sample size!
- $|\mathcal{X}||\mathcal{Y}| - 1 \rightarrow$ **task complexity!**

Consistent with our intuition

Multi-source Transfer Learning

We can extend this to the multi-source case.

- Target task:  i.i.d from $P_{XY}^{(0)} \rightarrow \hat{P}_{XY}^{(0)}$
- Source task 1:  i.i.d from $P_{XY}^{(1)} \rightarrow \hat{P}_{XY}^{(1)}$
- ...
- Source task k :  i.i.d from $P_{XY}^{(k)} \rightarrow \hat{P}_{XY}^{(k)}$

Multi-source Transfer Learning

We use $\alpha_0 \hat{P}_{XY}^{(0)} + \alpha_1 \hat{P}_{XY}^{(1)} + \cdots + \alpha_k \hat{P}_{XY}^{(k)}$ to estimate $P_{XY}^{(0)}$.

We have the testing loss

$$L_{\text{test}} = \chi^2 \left(P_{XY}^{(0)}, \sum_{i=0}^k \alpha_i P_{XY}^{(i)} \right) + \sum_{i=0}^k \frac{\alpha_i^2}{n_i} V^{(i)}, \quad (5)$$

and we can find the optimal coefficients.

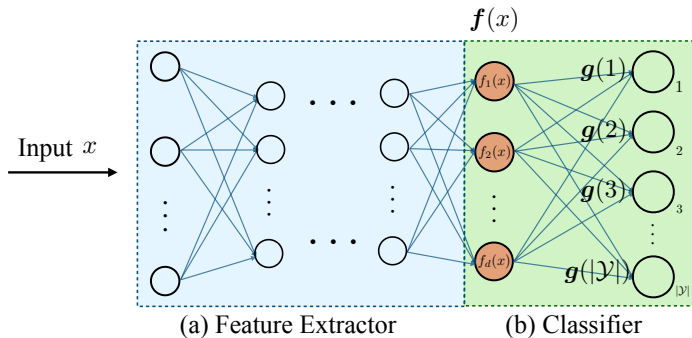
*We request $\sum_{i=0}^k \alpha_i = 1$ here.

Is our theory practical? **Yes**

However, there are 2 things we need to solve

- How the neural network models the distribution
- How to avoid the high dimensionality $|\mathcal{X}||\mathcal{Y}|$

Parametric Model



Discriminative Model




$$\tilde{P}_{Y|X}^{(f,g)}(y|x) \triangleq P_Y^{(0)}(y) (1 + \mathbf{f}^T(x) \mathbf{g}(y)) , \quad (6)$$

Parametric Model

When f is fixed, we can train a classifier g by samples under the referenced χ^2 -loss.

$$\hat{g}_i = \arg \min_g \chi_R^2(\hat{P}_{XY}^{(i)}, P_X^{(0)} \tilde{P}_{Y|X}^{(f,g)}),$$

where $\chi_R^2(P, Q) \triangleq \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \frac{(P(x, y) - Q(x, y))^2}{P_X^{(0)}(x) P_Y^{(0)}(y)}$.

- Target task:  $\rightarrow \hat{g}_0 \quad (\mathbb{E}[\hat{g}_0] = g_0)$
- Source task 1:  $\rightarrow \hat{g}_1 \quad (\mathbb{E}[\hat{g}_1] = g_1)$
- ...
- Source task k :  $\rightarrow \hat{g}_k \quad (\mathbb{E}[\hat{g}_k] = g_k)$

We use

$$\alpha_0 \tilde{P}_{Y|X}^{(\mathbf{f}, \hat{\mathbf{g}}_0)} + \alpha_1 \tilde{P}_{Y|X}^{(\mathbf{f}, \hat{\mathbf{g}}_1)} + \cdots + \alpha_k \tilde{P}_{Y|X}^{(\mathbf{f}, \hat{\mathbf{g}}_k)} \quad (7)$$

as our estimation. We have the testing loss

$$\begin{aligned} L_{\text{test}} = & \chi_R^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_i)} \right) + \sum_{i=0}^k \frac{\alpha_i^2}{n_i} \tilde{V}^{(i)} \\ & + \chi_R^2 \left(P_{XY}^{(0)}, P_X^{(0)} \tilde{P}_{Y|X}^{(\mathbf{f}, \mathbf{g}_0)} \right) \end{aligned}$$

Consistent with the theory in the discrete case

Transferability

$$(\alpha_0^*, \alpha_1^*, \dots, \alpha_k^*) = \arg \min_{(\alpha_0^*, \alpha_1^*, \dots, \alpha_k^*): \sum_{i=0}^k \alpha_i = 1} L_{\text{test}} \quad (8)$$

Let's see what's the affecting factors in the parametric model

- $\chi_R^2 \left(P_X^{(0)} \tilde{P}_{Y|X}^{(f, g_0)}, \sum_{i=0}^k \alpha_i P_X^{(0)} \tilde{P}_{Y|X}^{(f, g_i)} \right) \rightarrow \text{distance!}$
- $n_i \rightarrow \text{sample size!}$
- $\tilde{V}^{(i)} \rightarrow \text{task complexity!}$

Algorithm and Experimental Results

We have an **iterative algorithm**:

- $(f, g) \leftarrow$ Training Loss with given $\alpha_0, \alpha_1, \dots, \alpha_k$
- $(\alpha_0, \alpha_1, \dots, \alpha_k) \leftarrow$ Testing Loss with given f, g
- Until Converge

We made experiments on CIFAR-10, Office-31 and Office-Caltech datasets.

What's the contributions of our work?

- A theoretical analysis for transferability covering distance, sample sizes, task complexity at the same time
- An extension to continuous data
- A consistent algorithm that works