

# Towards more Reliable Transfer Learning (Supplementary Material)

Zirui Wang and Jaime Carbonell

Language Technologies Institute  
Carnegie Mellon University, Pittsburgh PA, USA  
`{ziruiw,jgc}@cs.cmu.edu`

## A Proof of Theorem 1

First, we rewrite the definition of *ideal hypothesis* [1] between the  $i^{th}$  source and the target as:

$$h_i^* = \arg \min_{h \in \mathcal{H}} \epsilon_T(h) + \epsilon_{S_i}(h)$$

We then denote the combined risk of the  $i^{th}$  ideal hypothesis as:

$$\lambda_i = \epsilon_T(h_i^*) + \epsilon_{S_i}(h_i^*)$$

Notice that the  $i^{th}$  ideal hypothesis is a single hypothesis that performs well on *both* domains and it explicitly defines the transferability between the  $i^{th}$  source and the target. When  $\lambda_i$  is large, we cannot expect to transfer knowledge from the  $i^{th}$  source to the target by minimizing source error.

Similarly, for the multi-source setting, we define the ideal hypothesis on domains weighted by  $\alpha$  and  $\mu$  as:

$$h_{\alpha,\mu}^* = \arg \min_{h \in \mathcal{H}} \epsilon_T(h) + \sum_i^K (\alpha_i \mu + (1 - \alpha_i) \frac{1 - \mu}{K - 1}) \epsilon_{S_i}(h)$$

The corresponding the combined risk is then defined as:

$$\lambda_{\alpha,\mu} = \epsilon_T(h_{\alpha,\mu}^*) + \sum_i^K (\alpha_i \mu + (1 - \alpha_i) \frac{1 - \mu}{K - 1}) \epsilon_{S_i}(h_{\alpha,\mu}^*)$$

Similar to the definition above,  $\lambda_{\alpha,\mu}$  defines the adaptability between weighted sources and the target, where  $\alpha$  specifies individual source importance and  $\mu$  controls the importance of inter-source relationships.

Now we can proof the bound in Theorem 1, assuming the linearity of risk:

$$\begin{aligned}
\epsilon_T(\hat{h}) &= \epsilon_T\left(\sum_i^K \alpha_i \hat{h}_i\right) = \sum_i^K \alpha_i \epsilon_T(\hat{h}_i) = \sum_i^K \alpha_i \left[ \mu \epsilon_T(\hat{h}_i) + \frac{1-\mu}{K-1} \sum_{j \neq i}^K \epsilon_T(\hat{h}_j) \right] \\
&\leq \sum_i^K \alpha_i \left[ \mu \left[ \epsilon_{S_i}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_i}, U_T) + 4 \sqrt{\frac{2d \log(2m') + \log(\frac{4}{\delta})}{m'}} + \lambda_i \right] + \right. \\
&\quad \left. \frac{1-\mu}{K-1} \sum_{j \neq i}^K \left[ \epsilon_{S_j}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_j}, U_T) + 4 \sqrt{\frac{2d \log(2m') + \log(\frac{4}{\delta})}{m'}} + \lambda_j \right] \right] \\
&\quad \text{(Theorem 1 of Blitzer)}
\end{aligned}$$

For the combined risks of ideal hypotheses  $\lambda_i$ , we can rearrange and show that:

$$\begin{aligned}
&\sum_i^K \alpha_i \left[ \mu \lambda_i + \frac{1-\mu}{K-1} \sum_{j \neq i} \lambda_j \right] \\
&= \sum_i^K \alpha_i \left[ \mu (\epsilon_T(h_i^*) + \epsilon_{S_i}(h_i^*)) + \frac{1-\mu}{K-1} \sum_{j \neq i} (\epsilon_T(h_j^*) + \epsilon_{S_j}(h_j^*)) \right] \\
&\leq \sum_i^K \alpha_i \left[ \mu (\epsilon_T(h_{\alpha, \mu}^*) + \epsilon_{S_i}(h_{\alpha, \mu}^*)) + \frac{1-\mu}{K-1} \sum_{j \neq i} (\epsilon_T(h_{\alpha, \mu}^*) + \epsilon_{S_j}(h_{\alpha, \mu}^*)) \right] \\
&= \epsilon_T(h_{\alpha, \mu}^*) + \sum_i^K (\alpha_i \mu + (1-\alpha_i) \frac{1-\mu}{K-1}) \epsilon_{S_i}(h_{\alpha, \mu}^*) \\
&= \lambda_{\alpha, \mu}
\end{aligned}$$

Then, the risk bound becomes:

$$\begin{aligned}
&\leq \left[ \sum_i^K \alpha_i \left[ \mu \left[ \epsilon_{S_i}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_i}, U_T) \right] + \frac{1-\mu}{K-1} \sum_{j \neq i}^K \left[ \epsilon_{S_j}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_j}, U_T) \right] \right] \right] + \\
&\quad 4 \sqrt{\frac{2d \log(2m') + \log(\frac{4}{\delta})}{m'}} + \lambda_{\alpha, \mu} \\
&\leq \left[ \sum_i^K \alpha_i \left[ \mu \left[ \hat{\epsilon}_{S_i}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_i}, U_T) \right] + \frac{1-\mu}{K-1} \sum_{j \neq i}^K \left[ \hat{\epsilon}_{S_j}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_j}, U_T) \right] \right] \right. \\
&\quad \left. + \sqrt{\frac{\mu^2}{\beta_i} + (\frac{1-\mu}{K-1})^2 \sum_{j \neq i} \frac{1}{\beta_j} \sqrt{\frac{d \log(2m) - \log(\delta)}{2m}}} \right] + 4 \sqrt{\frac{2d \log(2m') + \log(\frac{4}{\delta})}{m'}} + \lambda_{\alpha, \mu} \\
&\quad \text{(Lemma 4 of Blitzer)} \\
&= \left[ \sum_i^K \alpha_i \left[ \mu \left[ \hat{\epsilon}_{S_i}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_i}, U_T) \right] + \frac{1-\mu}{K-1} \sum_{j \neq i}^K \left[ \hat{\epsilon}_{S_j}(\hat{h}_i) + \frac{1}{2} \hat{d}_{\mathcal{H} \triangle \mathcal{H}}(U_{S_j}, U_T) \right] \right] \right] + \\
&\quad \left( \sum_i^K \alpha_i \sqrt{\frac{\mu^2}{\beta_i} + (\frac{1-\mu}{K-1})^2 \sum_{j \neq i} \frac{1}{\beta_j}} \right) \sqrt{\frac{d \log(2m) - \log(\delta)}{2m}} + 4 \sqrt{\frac{2d \log(2m') + \log(\frac{4}{\delta})}{m'}} + \lambda_{\alpha, \mu}
\end{aligned}$$

■

We also note that if  $\lambda_{\alpha,\mu}$  is small, we can set the hypothesis space  $\mathcal{H}$  to be the set containing all consistent hypotheses (i.e. the version space) and the weighted  $\mathcal{H}\Delta\mathcal{H}$ -divergence would also be small.

## References

1. Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., Wortman, J.: Learning bounds for domain adaptation. In: NIPS. pp. 129-136 (2008)