# DaMSTF: Domain Adversarial Learning Enhanced Meta Self-Training for Domain Adaptation

**Menglong Lu**[1†], **Zhen Huang**[1†], **Yunxiang Zhao**[2∗], **Zhiliang Tian**[1∗],
**Yang Liu**[1] and **Dongsheng Li**[1]

[1]National Key Laboratory of Parallel and Distributed Computing,
National University of Defense Technology, China
[2] Beijing Institute of Biotechnology, China

{lumenglong, huangzhen, tianzhiliang, liuyang12a, dsli}@nudt.edu.cn,
zhaoyx1993@163.com

## Abstract

Self-training emerges as an important research line on domain adaptation. By taking the model's prediction as the pseudo labels of the unlabeled data, self-training bootstraps the model with pseudo instances in the target domain. However, the prediction errors of pseudo labels (label noise) challenge the performance of self-training. To address this problem, previous approaches only use reliable pseudo instances, i.e., pseudo instances with high prediction confidence, to retrain the model. Although these strategies effectively reduce the label noise, they are prone to miss the hard examples. In this paper, we propose a new self-training framework for domain adaptation, namely Domain adversarial learning enhanced Self-Training Framework (DaMSTF). Firstly, DaMSTF involves meta-learning to estimate the importance of each pseudo instance, so as to simultaneously reduce the label noise and preserve hard examples. Secondly, we design a meta constructor for constructing the meta validation set, which guarantees the effectiveness of the meta-learning module by improving the quality of the meta validation set. Thirdly, we find that the meta-learning module suffers from the training guidance vanishment and tends to converge to an inferior optimal. To this end, we employ domain adversarial learning as a heuristic neural network initialization method, which can help the meta-learning module converge to a better optimal. Theoretically and experimentally, we demonstrate the effectiveness of the proposed DaMSTF. On the cross-domain sentiment classification task, DaMSTF improves the performance of BERT with an average of nearly 4%.

## 1 Introduction

Domain adaptation, which aims to adapt the model trained on the source domain to the target domain, attracts much attention in Natural Language Processing (NLP) applications(Du et al., 2020; Chen et al., 2021; Lu et al., 2022). Since domain adaptation involves labeled data from the source domain and unlabeled data from the target domain, it can be regarded as a semi-supervised learning problem. From this perspective, self-training, a classical semi-supervised learning approach, emerges a prospective research direction on domain adaptation (Zou et al., 2019; Liu et al., 2021).

Self-training consists of a series of loops over the pseudo labeling phase and model retraining phase. In the pseudo labeling phase, self-training takes the model's prediction as the pseudo labels for the unlabeled data from the target domain. Based on these pseudo-labeled instances, self-training retrains the current model in the model retraining phase. The trained model can be adapted to the target domain by repeating these two phases. Due to the prediction errors, there exists label noise in pseudo instances, which challenges self-training approaches (Zhang et al., 2017).

Previous self-training approaches usually involve a data selection process to reduce the label noise, i.e., preserving the reliable pseudo instances and discarding the remaining ones. In general, higher prediction confidence implies higher prediction correctness, so existing self-training approaches prefer the pseudo instances with high prediction confidence (Zou et al., 2019; Shin et al., 2020). However, fitting the model on these easy pseudo instances cannot effectively improve the model, as the model is already confident about its prediction. On the contrary, pseudo instances with low prediction confidence can provide more information for improving the model, but contain more label noise at the same time.

To simultaneously reduce the label noise and preserve hard examples, we propose to involve in meta-learning to reweight pseudo instances. Within a learning-to-learn schema, the meta-learning mod-

---
† contributed equally to this work
∗ corresponding author