# Weighted Multisource Tradaboost

João Antunes
Carnegie Mellon University
Instituto Superior Técnico
joaoa@andrew.cmu.edu

Alexandre Bernardino
Institudo Superior Técnico
alex@isr.tecnico.ulisboa.pt

Asim Smailagic
Carnegie Mellon University
asim@cs.cmu.edu

Daniel Siewiorek
Carnegie Mellon University
dps@cs.cmu.edu

## Abstract

*In this paper we propose an improved method for transfer learning that takes into account the balance between target and source data. This method builds on the state-of-the-art Multisource Tradaboost, but weighs the importance of each datapoint taking into account the amount of target and source data available. A comparative study is then presented exposing the performance of four transfer learning methods as well as the proposed Weighted Multisource Tradaboost. The experimental results show that the proposed method is able to outperform the base method as the number of target samples increase. These results are promising in the sense that source-target ratio weighing may be a path to improve current methods of transfer learning. However, against the asymptotic conjecture of [6], all transfer learning methods tested in this work get outperformed by a no-transfer SVM for large number on target samples.*

## 1. Introduction

Most machine learning techniques are based on the PAC (Probably Approximately Correct)[7] model, which states that while operating on a learning problem the samples used for training and the samples that we want to classify follow the same probability distribution. However, this assumption does not hold in a variety of cases. Frequently, the data used for training has become obsolete (e.g. due to changes on how data was collected) or simply that the data available is not enough to train a robust classifier. Insufficient data frequently occurs in classifiers that recognize a high number of classes (e.g. in object recognition systems routinely discriminate between $\approx 10^4$ categories)[5]. In this case, machine learning techniques give very little guarantees about the generalization error obtained. Transfer Learning is an approach to address the small dataset challenge. The intuition behind transfer learning is to mimic the way humans learn. The data we acquire from all our senses is stored in our memory along with concepts and inferences we make as to how to categorize this data. This makes it so that any new concept to be assimilated is not learned in isolation. Instead, we consider connections between what we already know and try to apply them to the new concept. The goal of transfer learning is to extract relevant information from data that does not need to come from the same probability distribution as the data to be classified by the final model. The ability to leverage more data during the learning process leads to more robust models since more information is used for training. In this paper, an improvement on a state-of-the-art transfer learning method is presented: Weighted Multisource Tradaboost.Our proposed approach incorporates the belief that if more target data is available, the contribution of the source data used in the model should gradually shift from model defining to fine-tuning. This is achieved with a re-weighing procedure. A comparative study is then provided between four state-of-the-art methods: Multisource Tradaboost[7], Task Tradaboost [7], Multi-KT [5], transfer learning decision forests [3], and Weighted Multisource Tradaboost. This study is evaluted on a subset of four classes of the Caltech-256 Dataset [4]. In turn, one of each of the four classes used is the target, while the other are used as sources. The classes chosen are dog, horse, leopard and zebra, chosen for empirically possessing a positive relationship with each other. The results show that our method can overcome the other methods in accuracy performance, but the higher asymptote assumption is still not achieved. This assumption states that a method employing transfer learning should outperform machine learning methods without transfer even when target data is abundant [6].

The contributions described in this paper are a comparative study exhibiting results not found in the literature, stat-

ing that the higher asymptote behavior theorized in [6] is not achieved by several state-of-the-art methods, and a novel approach for transfer learning that addresses this limitation of the methods studied. Although this limitation is not surpassed the method is showing a way to improve transfer learning approaches towards the theoretical asymptotic performance that can be applied in several transfer learning methods.

The rest of this document is organized as follows: Section 2 describes necessary concepts and notation introduced by transfer Learning. Section 3 describes the methods studied in this document. Section 4 describes the experiment ran, and the results obtained are discussed in Section 5. Finally, our conclusions and possible future research directions are presented in Section 6.

## 2. Transfer Learning

Transfer learning introduces several new concepts to machine learning. The definitions and notation described here will be used throughout this paper.

Standard machine learning tries to learn and then classify using one dataset for training and another one for testing. Both these datasets are assumed to come from the same distribution. In transfer learning information is leveraged from additional sources. The dataset that has the same distribution as the test data is called the target, and other(s) is(are) called source(s).

In this paper, the methods studied assume all datasets lie in the same feature space. This is called homogeneous transfer learning. If the feature space is different for at least one of the sources it is heterogeneous transfer learning (see MultiK-KT in [5]).

The success of transfer learning hinges on the inherent relationship between target and sources. In the case of a weak/non-existent connection between source and target the final classifier may actually be worse than its no-transfer counterpart. This phenomenon is known as negative transfer.

There are three measures by which transfer may improve learning [6] (see Fig. 1): Higher Start (better performance at the beginning of learning since source information is leveraged), Higher Slope (using the transferred knowledge the new task can be learned faster), and Higher Asymptote (since more information is being leveraged, the final system should have better performance). As shall be seen by the study presented in this paper, the Higher Asymptote hypothesis doesn't always hold true, even when a positive relationship between source and target can be established (See Section 5).

Finally, there is one more distinction between two types of transfer: instance transfer and task transfer. Instance transfer refers to scenarios in which some of the source data can actually be used to help train the new model. Multi-
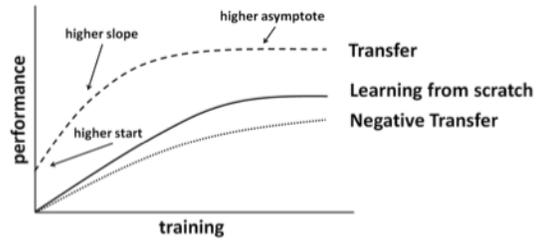


Figure 1. Three different ways in which transfer learning may improve traditional machine learning as a function of the number of target training samples: Higher Start, Higher Slope and Higher Asymptote (See Sec. 2.) Use of sources with no relation with the target may lead to the behavior described by the Negative Transfer curve. Figure adapted from [5].

SourceTradaboost is an example of this scenario. In task transfer, the source tasks are described explicitly by models trained *a priori*. Multi-KT and TaskTradaboost are examples of this type of transfer.

## 3. State-Of-The-Art

The methods compared in this paper comprise the recent state-of-the-art approached used in low data transfer learning. The methods are now described in detail.

### 3.1. Multi-KT - Support Vector Machines [5]

In 2014, Tommasi *et al.* [5] proposed a formulation for transfer learning using Support Vector Machines (SVM). Their problem setting was as follows: Assume that $j$ old (source) models (described by $\hat{w}_j$) are available *a priori*, and that these models can be expressed as a weighted sum of kernel functions (*e.g.* obtained *a priori* from an off-the-shelf SVM package). Then, in order to leverage the information already encoded in the other models, a simple framework is presented: change the cost function of an SVM solver to include a term imposing "model fidelity" (*i.e.* the cost function of the new model $w$ must be close to a weighted sum of the pre-existing j models) (see equation 1).

$$\operatorname*{argmin}_{w,b,\xi} \quad \frac{1}{2}\left\| w - \sum_{j=1}^{J} \beta_j \hat{w}_j \right\|^2 + C\sum_{i=1}^{n} \zeta_i \xi_i^2$$
$$\text{subject to} \quad y_i(w^T\phi(x_i) - b) \geq 1 - \xi_i,$$
$$\xi_i \geq 0$$

(1)

where $\zeta_i$ are used to balance the contribution of positive and negative samples, taking into account their proportion in the training set, $\beta_i$ are real numbers that control the influence that each *old* model should have over the new model (estimated *via* minimizing the leave-one-out error). The rest are

components that make up a standard SVM formulation(*i.e.* $\|w\|$ ensures margin maximization, $C \sum_{i=1}^{n} \zeta_i \xi_i^2$ encodes the trade-off between model fidelity and margin maximization and data fidelity; and the constraints ensure data fidelity).

## 3.2. Transfer Learning Decision Forests (TLDF)[3]

In 2014, Goussies *et al.*[3], proposed a method to do transfer learning using random decision forests. This is a method that uses data from several sources to shape the decision regions. Considering N+1 classification tasks, $T_0, \ldots, T_N$, the goal is to solve the classification task $T_0$, called the target task, using the knowledge of all tasks. By leveraging information from all datasets at once, the regions generated by the decision splits of each tree in the forest will construct a classifier with a higher classification accuracy, since more information is taken into account when shaping the decision regions (see Figure 2).
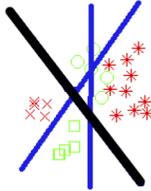


Figure 2. Consider two tasks (red and green) , each with two labels (stars and crosses for the red task, circles and squares for the green task. The red task is the target task, the green task is a source task. All three hyperplanes shown in the figure separate the target (red) dataset perfectly. The hyperplane represented in black, however, separates all the data from all the datasets simultaneously. According to the thinking presented in [3] the black hyperplane is preferable, and should be a better minimizer of the generalisation error. Image adapted from [3].

Goussies *et al.* go on to propose a mixed information gain formulation that formalizes the intuition described. For the k-th split:

$$\theta_k^* = \underset{\theta_k}{\operatorname{argmax}}(1 - \gamma)\mathcal{I}_0(\theta_k) + \gamma \sum_{n=1}^{N} \mathcal{I}_n(\theta_k) \qquad (2)$$

where $\mathcal{I}_0$ is the information gain on the target dataset (that stems from split $\theta_k$) and $\mathcal{I}_n, n = 1, \ldots, N$ are the information gains on the source datasets (stemming from the same split). $\gamma$ is a trade-off parameter that regulates the importance given to the information gain on the source and target datasets.

From this formulation a new problem arises: the leaf nodes of the tree are not required to have any datapoint belonging to the target. So, after creating a tree, a label propagation procedure is applied. For a given leaf node without a single target datapoint a distance vector is constructed with the distance from that node to all other leaf nodes that have at least one target datapoint. Then, the prediction made by the closest leaf node possessing at least one target datapoint is copied to the current node without target datapoints. The distance measure used is a Mahalanobis distance using the estimated mean and estimated covariance of the leaf nodes involved.

## 3.3. TrAdaboost

In 2007, Dai *et al.* proposed TrAdaboost[1], a transfer learning variant of AdaBoost. In 2010, Yao and Doretto[7], proposed two boosting models that perform transfer learning from multiple sources: MultiSource TrAdaboost and TaskTradaboost. Adaboost works as follows: At each iteration a weak classifier is trained. Then, the samples in the training set are re-weighted, increasing the weight of misclassified samples. This forces the next weak classifier trained to focus on getting the misclassified samples right. As such, expert models are being created for all the regions of the feature space of the dataset. Then, a final classifier is constructed by weighted majority voting of all the weak classifiers. The extensions proposed by Yao and Doretto in Tradaboost included:

### 3.3.1 MultiSource TrAdaboost [7]

For the MultiSource TrAdaboost model, proposed by Yao and Doretto[7] the availability of a very small target dataset is complemented by the availability of several larger datasets to be used as source. Information for all the datasets is leveraged by multiplexing between datasets in each iteration. When training one of the weak classifiers to boost, the target dataset is complemented by the source dataset that appears to be the most closely related to the target (*i.e.* the one that leads the weak classifier to the lowest error in the target dataset in the current iteration). Then, the weights of the datapoints in all the datasets are readjusted. However, unlike in Adaboost where misclassified points have their weight increased, the re-weighting procedure differs depending on which dataset is being used. Points in the target dataset have their weight increased if they are misclassified. On the other hand, misclassified points in any source dataset have their weight reduced. This is to express the belief that if a point in a source dataset is presenting conflicting information with the target dataset,then transfer from that datapoint should be avoided. The precise algorithm used is shown in Fig. 3, taken from [7].

### 3.3.2 Task TrAdaboost [7]

The TaskTrAdaboost performs transfer from previously available models, instead of from other datasets. It is divided in two phases.

**Algorithm 1:** *MultiSourceTrAdaBoost*

**Input:** Source training data $D_{S_1}, \cdots, D_{S_N}$, target training data $D_T$, and the maximum number of iterations $M$

**Output:** Target classifier function $\hat{f}_T : \mathcal{X} \to \mathcal{Y}$

1. Set $\alpha_S \doteq \frac{1}{2} \ln\left(1 + \sqrt{2 \ln \frac{n_S}{M}}\right)$, where $n_S \doteq \sum_k n_{S_k}$

2. Initialize the weight vector $(\mathbf{w}^{S_1}, \cdots, \mathbf{w}^{S_N}, \mathbf{w}^T)$, where $\mathbf{w}^{S_k} \doteq (w_1^{S_k}, \cdots, w_{n_{S_k}}^{S_k})$, and $\mathbf{w}^T \doteq (w_1^T, \cdots, w_{n_T}^T)$ to the desired distribution

   **for** $t \leftarrow 1$ **to** $M$ **do**

3.    Empty the set of candidate weak classifiers, $\mathcal{F} \leftarrow \emptyset$

4.    Normalize to 1 the weight vector $(\mathbf{w}^{S_1}, \cdots, \mathbf{w}^{S_N}, \mathbf{w}^T)$

      **for** $k \leftarrow 1$ **to** $N$ **do**

5.       Find the candidate weak classifier $h_t^k : \mathcal{X} \to \mathcal{Y}$ that minimizes the classification error over the combined set $D_{S_k} \bigcup D_T$, weighted according to $(\mathbf{w}^{S_k}, \mathbf{w}^T)$

6.       Compute the error of $h_t^k$ on $D_T$:
$$\epsilon_t^k \doteq \sum_j \frac{w_j^T [y_j^T \neq h_t^k(\mathbf{x}_j^T)]}{\sum_i w_i^T}$$

7.       $\mathcal{F} \leftarrow \mathcal{F} \cup (h_t^k, \epsilon_t^k)$

8.    Find the weak classifier $h_t : \mathcal{X} \to \mathcal{Y}$ such that
$$(h_t, \epsilon_t) \doteq \arg\min_{(h,\epsilon) \in \mathcal{F}} \epsilon$$

9.    Set $\alpha_t \doteq \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$, where $\epsilon_t < 1/2$

10.   Update the weight vector
$$w_i^{S_k} \leftarrow w_i^{S_k} e^{-\alpha_S |h_t(\mathbf{x}_i^{S_k}) - y_i^{S_k}|}$$
$$w_i^T \leftarrow w_i^T e^{\alpha_t |h_t(\mathbf{x}_i^T) - y_i^T|}$$

**return** $\hat{f}_T(\mathbf{x}) \doteq \text{sign}\left(\sum_t \alpha_t h_t(\mathbf{x})\right)$

Figure 3. MultiSource Tradaboost algorithm. Taken from [7]

Phase I consists of training off-the-shelf Adaboost models on each of the source datasets available.

Phase II mimicks Adaboost by boosting several weak classifiers on a weighted dataset. However, in TaskTradaboost the weak models used are the Adaboost models trained on the source datasets. The weight update step in this algorithm is identical to the one in Adaboost.

### 3.3.3 Weighted Multisource Tradaboost

When using transfer learning, information from both target and source datasets is leveraged. Naturally, most strategies have some way to weigh the data according to the prior belief of how similar the target data's and the source data's distribution is (*i.e.:* the $\beta_j$ in Multi-KT, the $\gamma$ parameter in TLDF's, and the weight vectors $w_i^{S_k}$ and $w_i^T$ in Tradaboost). However, to our knowledge, no method incorporates the proportion of target and source data available as prior knowledge in the mixing of target and source information in the learning stage.

We believe this approach to be sound because, if more target data is available, the contribution of the source data used in the model should gradually shift from model defining to fine-tuning. We postulate this is the case because as

more target data becomes available, the model built using only target data becomes more and more robust. In that case, forcing the model to acommodate source data can actually be detrimental to the model's performance. We shall prove this with a comparative study in the results section.

Our approach follows the general method described by tradaboost but replaces the weight update rules defined in step 10 of the algorithm (see Fig. 3) to take into account the proportion of target and source data available.

Instead, we propose:

$$w_i^{S_k} \leftarrow w_i^{S_k} e^{-\eta \alpha_S |h_t(x_i^{S_k}) - y_i^{S_k}|}$$
$$or \qquad (3)$$
$$w_i^{S_k} \leftarrow w_i^{S_k} \eta e^{-\alpha_S |h_t(x_i^{S_k}) - y_i^{S_k}|}$$

$$w_i^T \leftarrow w_i^T e^{\eta \alpha_T |h_t(x_i^T) - y_i^T|}$$
$$or \qquad (4)$$
$$w_i^T \leftarrow w_i^T \eta e^{\alpha_T |h_t(x_i^T) - y_i^T|}$$

where $\eta$ is a term that depends on the amount of target and source data available for training. The same term can be used for both target and source datapoints because the weight update step (shown in equations 3 and 4)inverts the signal of the exponent when switching dataset. Strategies for how to define this quantity are discussed in Sec. 4.

## 4. Experimental Design

We compare all the methods described in Section 3 with a subset of the Caltech-256 Dataset [4]. This is a dataset composed of 256 classes, with images as datapoints. The images range from high-quality pictures to poor drawings of the subject of the class. A subset of 4 classes was chosen from those available: dog, horse, leopard and zebra as well as the background class. These classes were shown to test positive transfer from empirically related classes: 4-legged animals. For these classes we downloaded the Scale Invariant Feature Transform (SIFT) features from [2] (See [2] for details). These features have a dimension of 300.

The results presented are averaged over 5 tests done with random permutations of the data. For each test, the results are averaged over 4 runs, each with a different 4-legged animal as target. As such, all experiments are averaged over twenty runs. Finally, the tests are run with the number of target points available ranging from 1 to 10.

For comparison with the no-transfer scenario, an off-the-shelf SVM classifier is trained exclusively on the target data.

### 4.1. Method Hyperparameters

For Multi-KT the $C$ parameter (see Equation 1) is chosen via cross-validation on the source data. The $\beta_j$ parameters are chosen by minimizing the leave-one-out error. In [5]

feature fusion is used. For fairness of comparison with the other methods only SIFT features were used.

For the Random Forests methods the parameters were decided according to the values found in the literature instead of chosen by testing different values for the parameters. This was due to the long time needed for each run of this method. The parameters used were: $\gamma = 0.8$ (Controls the influence of sources and target when calculating splits), Maximum tree depth = 10, number of trees in a forest = 3.

For MultiSource Tradaboost and Task Tradaboost the only hyperparameter is the number of iterations to run. This value was set at 50 due to computational limitations.

For Weighted Multisource Tradaboost 2 different values were empirically chosen for testing:

- $\eta = \frac{N_T * 100}{N_S}$

- $\eta = \frac{N_T^2 * 100}{N_S}$

where $N_T$ is the number of target datapoints and $N_S$ is the number of source datapoints. The factor of 100 inserted in the numerator describes the belief that in most transfer learning settings target data will be scarce while source data will be abundant. So both these terms enforce that the influence of source samples will be greatly diminished as more target samples become available. Each of these values was tested on both variants shown in Equations 3 and 4, resulting in four different tests.

For comparison with the no-transfer scenario, an off-the-shelf SVM classifier is trained exclusively on the target data. The hyperparameters for this model are the same as those used for Multi-KT but setting all the $\beta_j$ to 0.

## 5. Results

Running the comparative study of the methods described in Section 3 on the Caltech-256 Dataset, the graph on Figure 4 was drawn.

As can be observed in Figure 4, the no transfer scenario outperforms all other approaches when 10 target samples are available. None of the methods studied are able to achieve the higher asymptote behaviour[6] (see Section 2).

All methods studied outperform the no-transfer approach in scenarios where the number of target samples available is very limited, up to 7 target samples.

The fact that all methods get overtaken when more target data is available suggests that once "high-quality" (target) data is available in sufficient quantity the methods are unable to extract information from the sources in a way that is not conflicting with the targets. This implies that further protection from negative transfer is required.

The TLDF method shows very unstable performance. Also, results have been found in the literature stating that this method outperforms no-transfer in cases where more
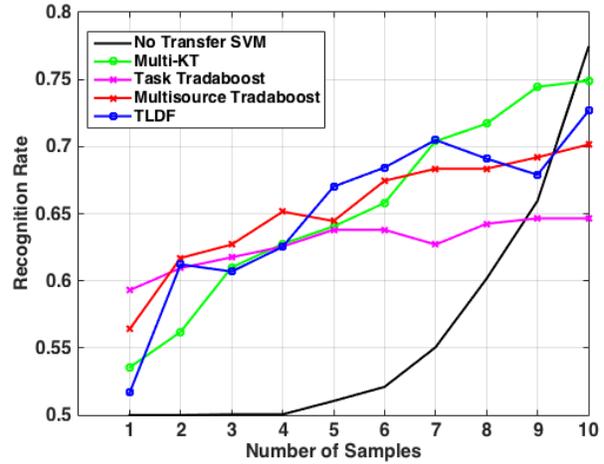


Figure 4. Results obtained running the experiment described in Section 4. Each point represents the average over 20 runs with random sample and target selection.

target data is available. Only 1-10 target datapoints are available in our experiment, and this amount of data is not enough to populate a feature space with a dimension of 300. Since the feature space is sparsely populated, during the label propagation step the distances between leaf nodes with no target datapoints and the closest leaf node with a target can be immense, which could justify the instability found.

To address this limitation of failing to achieve the higher asymptote behaviour, the Weighted Multisource Tradaboost method was applied to the same dataset in the same conditions. The results obtained are shown in Figure 5.

In this figure, WMS-Exponential and WMS - Multiplicative refer to using the linear $\eta$ weight shown in Section 4.1 witht the weight update rules defined in Equations 3 and 4 respectively. WMS-Squared exponential and WMS-Squared Multiplicative correspond to the same weight update rules using the squared $\eta$ weight shown in Section 4.1. As can be seen in Figure 5 all attempts outperform Multisource Tradaboost except for the Squared Outside attempt. However, the failure to achieve the higher asymptote behavior still eludes us. More research is needed on this topic.

## 6. Conclusion

All methods studied outperform the no-transfer approach when very little target data is available only to get outperformed by it when more target data is accessible. This failure to achieve the higher asymptote behaviour theorized in [6] is an unpublished result, and is one of the contributions of this work. Our proposed attempt to improve Multisource Tradaboost to achieve the higher asymptote behaviour did not come to fruition, albeit managing to improve the classification performance. Further research is required to determine when to expect the no-transfer approach to be prefer-
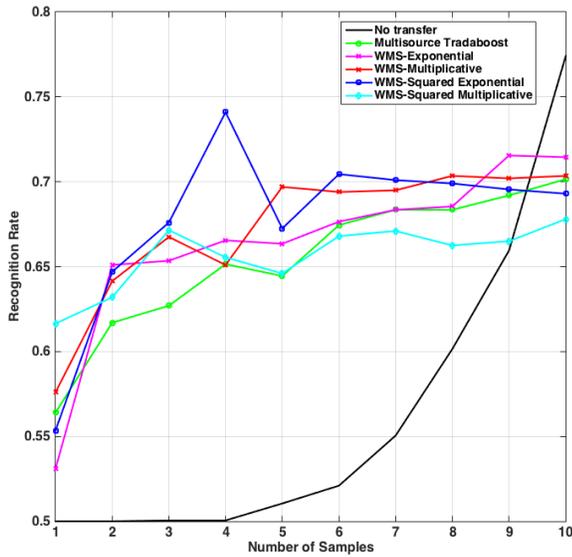
Figure 5. Results obtained running the experiment described in Section 4 for Weighted Multisource Tradaboost. Each point represents the average over 20 runs with random sample and target selection.

able, and how to achieve the higher asymptote. The strategies employed in this paper can be ported to other transfer learning methods, namely the Multi-KT method. Our strategy to improve transfer learning approaches towards the theoretical asymptotic performance predicted in [6] is our other contribution.

## 6.1. Future Work

According to the results obtained the following future research is suggested:

- Add a regularization term in the $\beta_j$ calculation steps for Multi-KT that takes into account how much target data is available. This would hopefully lead the method to not be outperformed by no-transfer approaches, fusing the best of both worlds

- Further testing with the Random Forests approaches is needed in order to evaluate the performance of these methods in situations where more target data is available. Results found in the literature indicate that this research is promising.

- Test all the methods in a scenario where the classes are unrelated and test specifically for resilience against negative transfer

## References

[1] W. Dai, Q. Yang, G.-R. Xue, and Y. Yu. Boosting for transfer learning. In *Proceedings of the 24th international conference on Machine learning*, pages 193–200. ACM, 2007.

[2] P. Gehler and S. Nowozin. On feature combination for multiclass object classification. In *Computer Vision, 2009 IEEE 12th International Conference on*, pages 221–228. IEEE, 2009.

[3] N. A. Goussies, S. Ubalde, and M. Mejail. Transfer Learning Decision Forests for Gesture Recognition. *Journal of Machine Learning Research*, 15:3667–3690, 2014.

[4] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. 2007.

[5] T. Tommasi, F. Orabona, and B. Caputo. Learning categories from few examples with multi model knowledge transfer. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 36(5):928–941, 2014.

[6] L. Torrey and J. Shavlik. Transfer learning. *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques*, 1:1–2, 2009.

[7] Y. Yao and G. Doretto. Boosting for transfer learning with multiple sources. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 1855–1862. IEEE, 2010.