

Scalable Greedy Algorithms for Transfer Learning

Ilja Kuzborskij^{*a,b,c}, Francesco Orabona^{†d}, and Barbara Caputo^{‡c,a}

^aIdiap Research Institute, Centre du Parc, Rue Marconi 19, 1920 Martigny, Switzerland

^bÉcole Polytechnique Fédérale de Lausanne (EPFL), Switzerland

^dYahoo Labs, 229 West 43rd Street, 10036 New York, NY, USA

^cUniversity of Rome La Sapienza, Dept. of Computer, Control and Management Engineering, Rome, Italy

September 11, 2018

Abstract

In this paper we consider the binary transfer learning problem, focusing on how to select and combine sources from a large pool to yield a good performance on a target task. Constraining our scenario to real world, we do not assume the direct access to the source data, but rather we employ the source hypotheses trained from them. We propose an efficient algorithm that selects relevant source hypotheses and feature dimensions simultaneously, building on the literature on the best subset selection problem. Our algorithm achieves state-of-the-art results on three computer vision datasets, substantially outperforming both transfer learning and popular feature selection baselines in a small-sample setting. We also present a randomized variant that achieves the same results with the computational cost independent from the number of source hypotheses and feature dimensions. Also, we theoretically prove that, under reasonable assumptions on the source hypotheses, our algorithm can learn effectively from few examples.

1 Introduction

Over the last few years, the visual recognition research landscape has been heavily dominated by Convolutional Neural Networks, thanks to their ability to leverage effectively over massive amount of training data [1]. This trend dramatically confirms the widely accepted truth that any learning algorithm performs better when trained on a lot of data. This is even more true when facing noisy or “hard” problems such as large-scale recognition [2]. However, when tackling large scale recognition problems, gathering substantial training data for all classes considered might be challenging, if not almost impossible. The occurrence of real-world objects follows a long tail distribution, with few objects occurring very often, and many with few instances. Hence, for the vast majority of visual categories known to human beings, it is extremely challenging to collect training data of the order of $10^4 - 10^5$ instances. The “long tail” distribution problem was noted and studied by Salakhutdinov *et al.* [3], who proposed to address it by leveraging on the prior knowledge available to the learner. Indeed, learning systems are often not trained from scratch: usually they can be built on previous knowledge acquired over time on related tasks [4]. The scenario of learning from few examples by *transferring* from what is already known to the learner is collectively known as Transfer Learning. The target domain usually indicates the task at hand and the source domain the prior knowledge of the learner.

*ilja.kuzborskij@idiap.ch

†francesco@orabona.com

‡caputo@dis.uniroma1.it

Most of the transfer learning algorithms proposed in the recent years focus on the object detection task (binary transfer learning), assuming access to the training data coming from both source and target domains [4]. While featuring good practical performance [5], they often demonstrate poor scalability w.r.t. the number of sources. An alternative direction, known as a Hypothesis Transfer Learning (HTL) [6, 7], consists in transferring from the *source hypotheses*, that is classifiers trained from them. This framework is practically very attractive [8, 9, 10], as it treats source hypotheses as black boxes without any regard of their inner workings.

The goal of this paper is to develop an HTL algorithm able to deal effectively and efficiently with a large number of sources, where our working definition of large is at least 10^3 . Note that this order of magnitude is also the current frontier in visual classification [2]. To this end, we cast Hypothesis Transfer Learning as a problem of *efficient selection* and *combination* of source hypotheses from a large pool. We pose it as a subset selection problem building on results from the literature [11, 12]. We present¹ a greedy algorithm, *GreedyTL*, which attains state of the art performance even with a very limited amount of data from the target domain. Moreover, we also present a randomized approximate variant of *GreedyTL*, called *GreedyTL-59*, that has a complexity *independent* from the number of sources, with no loss in performance. Our key contribution is a L_2 -regularized variant of the Forward Regression algorithm [14]. Since our algorithm can be viewed as a feature selection algorithm as well as an hypothesis transfer learning approach, we extensively evaluate it against popular feature selection and transfer learning baselines. We empirically demonstrate that *GreedyTL* dominates all the baselines in most small-sample transfer learning scenarios, thus proving the critical role of regularization in our formulation. Experiments over three datasets show the power of our approach: we obtain state of the art results in tasks with up to 1000 classes, totalling 1.2 million examples, with only 11 to 20 training examples from the target domain. We back our experimental results by proving generalization bounds showing that, under reasonable assumptions on the source hypotheses, our algorithm is able to learn effectively with very limited data.

The rest of the paper is organised as follows: after a review of the relevant literature in the field (section 2), we cast the transfer learning problem in the subset selection framework (section 3). We then define our *GreedyTL*, in section 4, deriving its formulation, analysing its computational complexity and its theoretical properties. Section 5 describes our experimental evaluation and discuss the related findings. We conclude with an overall discussion and presenting possible future research avenues.

2 Related Work

The problem of how to exploit prior knowledge when attempting to solve a new task with limited, if any, annotated samples is vastly researched. Previous work span from transfer learning [4] to domain adaptation [15, 16], and dataset bias [17]. Here we focus on the first. In the literature there are several transfer learning settings [16, 15, 5]. The oldest and most popular is the one assuming access to the data originating from both the source and the target domains [16, 5, 15, 18, 19, 20, 21]. There, one typically assumes that plenty of source data are available, but access to the target data is limited: for instance, we can have many unlabeled examples and only few labeled [22]. Here we focus on the Hypothesis Transfer Learning framework (HTL, [6, 7]). It requires to have access only to *source hypotheses*, that is classifiers or regressors trained on the source domains. No assumptions are made on how these source hypotheses are trained, or about their inner workings: they are treated as “black boxes”, in spirit similar to classifier-generated visual descriptors such as *Classemes* [23] or *Object-Bank* [24]. Several works proposed HTL for visual learning [8, 9, 25], some exploiting more explicitly the connection with *classemes*-like approaches [26, 27], demonstrating an intriguing potential. Although offering scalability, HTL-based approaches proposed so far have been tested on problems with less than a few hundred of sources [9], already showing some difficulties in selecting informative sources.

Recently, the growing need to deal with large data collections [2, 28] has started to change the focus

¹We build upon preliminary results presented in [13].

and challenges of research in transfer learning. Scalability with respect to the amount of data and the ability to identify and separate informative sources from those carrying noise for the task at hand have become critical issues. Some attempts have been made in this direction. For example, [29, 30] used taxonomies to leverage learning from few examples on the SUN09 dataset. In [29], authors attacked the transfer learning problem on the SUN09 dataset by using additional data from another dataset. Zero-shot approaches were investigated by [31] on a subset of the Imagenet dataset. Large-scale visual detection has been explored by [30]. However, all these approaches assume access to all source training data. A slightly different approach to transfer learning that aimed to circumvent this limitation, is reuse of a large convolutional neural network pre-trained on a large visual recognition dataset. The simplest approach is to use outputs of intermediate layers of such a network, such as DeCAF [1] or Caffe [32]. A more sophisticated way of reuse is fine-tuning, a kind of warm-start, that has been successfully exploited in visual detection [33] and domain adaptation [34, 35].

In many of these works the use of richer sources of information has been supported by an increase in the information available in the target domain as well. From an intuitive point of view, this corresponds to having more data points than dimensions. Of course, this makes the learning and selection process easier, but in many applications it is not a reasonable hypothesis. Also, none of the proposed algorithms has a theoretical backing.

While not explicitly mentioned before, the problem outlined above can also be viewed as a learning scenario where the number of features is by far larger than the number of training examples. Indeed, learning with classeme-like features [23, 24] when only few training examples are available can be seen as a Hypothesis Transfer Learning problem. Clearly, a pure empirical risk minimization would fail due to severe overfitting. In machine learning and statistics this is known as a feature selection problem, and is usually addressed by constraining or penalizing the solution with sparsity-inducing norms. One important sparsity constraint is a non-convex L_0 pseudo-norm constraint $\|\mathbf{w}\|_0 \leq k$, that simply corresponds to choosing up to k non-zero components of a vector \mathbf{w} . One usually resorts to the *subset selection* methods, and greedy algorithms for obtaining solutions under this constraint [11, 36, 12, 37]. However, in some problems introducing L_0 constraint might be computationally difficult. There, a computationally easier alternative is a convex relaxation of L_0 , the L_1 regularization. Empirical error minimization with L_1 penalty with various loss functions (for square loss is known as Lasso) has many favorable properties and is well studied theoretically [38]. Yet, L_1 penalty is known to suffer from several limitations, one of which is poor empirical performance when there are many correlated features. Perhaps the most famous way to resolve this issue is an *elastic net* regularization which is a weighted mixture of L_1 and squared L_2 penalties [14]. Since our work partially falls into the category of feature selection, we have extensively evaluated the aforementioned baselines in our task. As it will be shown below, none of them achieves competitive performances compared to our approach.

3 Transfer Learning through Subset Selection

Definitions. We will denote with small and capital bold letters respectively column vectors and matrices, e.g. $\mathbf{a} = [a_1, a_2, \dots, a_d]^T \in \mathbb{R}^d$ and $\mathbf{A} \in \mathbb{R}^{d_1 \times d_2}$. The subvector of \mathbf{a} with rows indexed by set S is \mathbf{a}_S , while the square submatrix of \mathbf{A} with rows and columns indexed by set S is \mathbf{A}_S . For $\mathbf{x} \in \mathbb{R}^d$, the *support* of \mathbf{x} is $\text{supp}(\mathbf{x}) = \{i \in \{1, \dots, d\} : x_i \neq 0\}$. Denoting by \mathcal{X} and \mathcal{Y} respectively the input and output space of the learning problem, the training set is $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$, drawn i.i.d. from the probability distribution p defined over $\mathcal{X} \times \mathcal{Y}$. We will focus on the binary classification problem so $\mathcal{Y} = \{-1, 1\}$, and, without loss of generality, $\mathcal{X} = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1, \mathbf{x} \in \mathbb{R}^d\}$.

To measure the accuracy of a learning algorithm, we have a non-negative *loss* function $\ell(h(\mathbf{x}), y)$, which measures the cost incurred predicting $h(\mathbf{x})$ instead of y . In particular, we will focus on the square loss, $\ell(h(\mathbf{x}), y) = (h(\mathbf{x}) - y)^2$, for its appealing computational properties. The *risk* of a hypothesis h , with respect to the probability distribution p , is then defined as $R(h) := \mathbb{E}_{(\mathbf{x}, y) \sim p}[\ell(h(\mathbf{x}), y)]$, while the *empirical risk* given a training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$ is $\hat{R}(h) := \frac{1}{m} \sum_{i=1}^m \ell(h(\mathbf{x}_i), y_i)$. Whenever the hypothesis

is a linear predictor, that is, $h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, we will also use risk notation as $R(\mathbf{w}) = R(h_{\mathbf{w}})$ and $\hat{R}(\mathbf{w}) = \hat{R}(h_{\mathbf{w}})$.

Source Selection. Assume, that we are given a finite source hypothesis set $\{h_i^{\text{src}}\}_{i=1}^n$ and the training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$. As in previous works [39, 9, 26], we consider the target hypothesis to be of the form

$$h_{\mathbf{w}, \boldsymbol{\beta}}^{\text{trg}}(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + \sum_{i=1}^n \beta_i h_i^{\text{src}}(\mathbf{x}), \quad (1)$$

where \mathbf{w} and $\boldsymbol{\beta}$ are found by the learning procedure. The essential parameter here is $\boldsymbol{\beta}$, that is the one controlling the influence of each source hypothesis. Previous works in transfer learning have focused on finding $\boldsymbol{\beta}$ such that it minimizes the error on the training set, subject to some condition on $\boldsymbol{\beta}$. In particular, [9] proposed to minimize the leave-one-out error w.r.t. $\boldsymbol{\beta}$, subject to $\|\boldsymbol{\beta}\|_2 \leq \tau$, which is known to improve generalization for the right choice of τ [6]. A slightly different approach is to use $\|\boldsymbol{\beta}\|_1 \leq \tau$ regularization for this purpose [9], that induces solutions with most of the coefficients equal to 0, thus assuming that the optimal $\boldsymbol{\beta}$ is sparse.

In this work we embrace a weaker assumption, namely, there exist up to k sources that collectively improve the generalization on the target domain. Thus, we pose the problem of the Source Selection as a minimization of the regularized empirical risk on the target training set, while constraining the number of selected source hypotheses.

k -Source Selection. Given the training set $\{([\mathbf{x}_i^\top, h_1^{\text{src}}(\mathbf{x}_i), \dots, h_n^{\text{src}}(\mathbf{x}_i)]^\top, y_i)\}_{i=1}^m$ we have the optimal target hypothesis $h_{\mathbf{w}^*, \boldsymbol{\beta}^*}^{\text{trg}}$ by solving,

$$\begin{aligned} (\mathbf{w}^*, \boldsymbol{\beta}^*) &= \arg \min_{\mathbf{w}, \boldsymbol{\beta}} \left\{ \hat{R}(h_{\mathbf{w}, \boldsymbol{\beta}}^{\text{trg}}) + \lambda \|\mathbf{w}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \right\}, \\ \text{s.t. } &\|\mathbf{w}\|_0 + \|\boldsymbol{\beta}\|_0 \leq k. \end{aligned} \quad (2)$$

Notably, the problem (2) is a special case of the *Subset Selection* problem [11]: choose a subset of size k from the n observation variables, which collectively give the best prediction on the variable of interest. However, the Subset Selection problem is **NP-hard** [11]. In practice we can resort to algorithms generating approximate solutions, for many of which we have approximation guarantees. Hence, due to the extensive practical and theoretical results, we will treat the k -Source Selection as a Subset Selection problem, building atop of existing guarantees.

We note that our formulation, (2), differs from the classical subset selection for the fact that it is $L2$ -regularized. This technical modification makes an essential practical and theoretical difference and it is the crucial part of our algorithm. First, $L2$ regularization is known to improve the generalization ability of empirical risk minimization. Second, we show that regularization also improves the quality of the approximate solution in situations when the sources, or features, are correlated. At the same time, the experimental evaluation corroborates our theoretical findings: Our formulation substantially outperforms standard subset selection, feature selection algorithms, and transfer learning baselines.

4 Greedy Algorithm for k -Source Selection

In this section we state the algorithm proposed in this work, *GreedyTL*². In the following we will denote by $U = \{1, \dots, n + d\}$ the index set of all available source hypotheses and features, and by S , the index set of selected ones.

²Source code is available at <http://idiap.ch/~ikuzbor/>

GreedyTL. Let $\mathbf{X} \in \mathbb{R}^{m \times d}$ and $\mathbf{y} \in \{+1, -1\}^m$ be the zero-mean unit-variance training set, $\{h_i^{src}\}_{i=1}^n$, source hypothesis set, and k and λ , regularization parameters. Then, denote $\mathbf{C} = \mathbf{Z}^\top \mathbf{Z}$ and $\mathbf{b} = \mathbf{Z}^\top \mathbf{y}$, where $\mathbf{Z} = \begin{bmatrix} \mathbf{X} & h_1^{src}(\mathbf{x}_1) & \dots & h_n^{src}(\mathbf{x}_1) \\ \dots & \dots & \dots & \dots \\ h_1^{src}(\mathbf{x}_m) & \dots & \dots & h_n^{src}(\mathbf{x}_m) \end{bmatrix}$, and select set S of size k as follows: (I) Initialize $S \leftarrow \emptyset$ and $U \leftarrow \{1, \dots, n + d\}$. (II) Keep populating S with $i \in U$, that maximize $\mathbf{b}_S^\top ((\mathbf{C} + \lambda \mathbf{I})_S^{-1})^\top \mathbf{b}_S$, as long as $|S| \leq k$ and U is non-empty.

In this basic formulation, the algorithm requires to invert a $(d + n)$ -by- $(d + n)$ matrix at each iteration of a greedy search. Clearly, this naive approach gets prohibitive with the growth of the number of source hypotheses, feature dimensions, and desired subset size, since its computational complexity would be in $\mathcal{O}(k(d + n)^4)$. However, we note that in transfer learning one typically assumes that training set is much smaller than number of sources and feature dimension. For this reason we apply rank-one updates w.r.t. the dual solution of regularized subset selection, so that the size of the inverted matrix does not change. The computational complexity then improves to $\mathcal{O}(k(d + n)m^2)$. We present the pseudocode of such a variant of our algorithm, **GreedyTL with Rank-One Updates** in Algorithm 1. The computational complexity of the operations is shown at the end of each line.

Algorithm 1 GreedyTL with Rank-One Updates

Input: $\mathbf{Z} \in \mathbb{R}^{m \times (d+n)}$ – m examples formed from features and source predictions,

- 1: $\mathbf{y} \in \{-1, +1\}^m$ – labels,
- 2: $k \in \{1, \dots, d + n\}, \lambda \in \mathbb{R}_+$ – hyperparameters.

Output: \mathbf{w} – target predictor.

- 3: $U \leftarrow \{1, \dots, d + n\}$ ▷ All candidates
 - 4: $S \leftarrow \emptyset$ ▷ Selected sources and features
 - 5: $\mathbf{K} \leftarrow [\mathbf{0}, \dots, \mathbf{0}] \in \mathbb{R}^{m \times m}$
 - 6: $\mathbf{G} \leftarrow \lambda^{-1} \mathbf{I} \in \mathbb{R}^{m \times m}$
 - 7: **while** $U \neq \emptyset$ **and** $|S| \leq k$ **do**
 - 8:

$$i^* \leftarrow \arg \max_{i \in U} \left\{ \mathbf{y}^\top (\mathbf{K} + \mathbf{z}_i \mathbf{z}_i^\top) \mathbf{G}' \mathbf{y} \mid \mathbf{G}' \leftarrow \mathbf{G} - \frac{\mathbf{G} \mathbf{z}_i \mathbf{z}_i^\top \mathbf{G}}{1 + \mathbf{z}_i^\top \mathbf{G} \mathbf{z}_i} \right\}$$
▷ $\mathcal{O}((d + n)(m^2 + m))$
 - 9:
 - 10: **Computing** \mathbf{G}' : ▷ $\mathcal{O}(m^2 + m)$
 - 11: **Computing** score of i : ▷ $\mathcal{O}(m^2 + m)$
 - 12: $S \leftarrow S \cup \{i^*\}$
 - 13: $U \leftarrow U \setminus \{i^*\}$
 - 14: $\mathbf{K} \leftarrow \mathbf{K} + \mathbf{z}_{i^*} \mathbf{z}_{i^*}^\top$ ▷ $\mathcal{O}(m^2)$
 - 15: $\mathbf{G} \leftarrow \mathbf{G} - \frac{\mathbf{G} \mathbf{z}_{i^*} \mathbf{z}_{i^*}^\top \mathbf{G}}{1 + \mathbf{z}_{i^*}^\top \mathbf{G} \mathbf{z}_{i^*}}$ ▷ $\mathcal{O}(m^2 + m)$
 - 16:
 - 17: **end while** ▷ $\mathcal{O}(k(d + n)m^2)$
 - 18: $\mathbf{w} \leftarrow \mathbf{0} \in \mathbb{R}^{d+n}$
 - 19: $w_i \leftarrow \mathbf{z}_i^\top \mathbf{G} \mathbf{y}, \forall i \in S$
-

Derivation of the Algorithm. We derive GreedyTL by extending the well known Forward Regression (FR) algorithm [11], which gives an approximation to the subset selection problem, the problem of our interest. FR is known to find a good approximation as far as features are uncorrelated [11]. In the following, we build upon FR by introducing a Tikhonov (L_2) regularization into the formulation. The purpose of regularization is twofold: first, it improves the generalization ability of the empirical risk minimization, and second, it makes the algorithm more robust to the feature correlations, thus opting to find better approximate solution.

First, we briefly formalize the subset selection problem. In a subset selection problem one tries to achieve a good prediction accuracy on the *predictor* random variable Y , given a linear combination of a subset of the *observation* random variables $\{X_i\}_{i=1}^n$. The least squares subset selection then reads as

$$\min_{|S|=k, \mathbf{w} \in \mathbb{R}^k} \mathbb{E} \left[\left(Y - \sum_{i \in S} w_i X_i \right)^2 \right].$$

Now denote the covariance matrix of zero-mean unit-variance observation random variables by \mathbf{C} (a correlation matrix), and the correlations between Y and $\{X_i\}_{i=1}^n$ as \mathbf{b} . Note that the zero-mean unit-variance assumption will be necessary to prove the theoretical guarantees of our algorithm. By virtue of the analytic solution to least-squares and using the introduced notation, we can also state the equivalent *Subset Selection problem*: $\max_{|S|=k} \mathbf{b}_S^\top (\mathbf{C}_S^{-1})^\top \mathbf{b}_S$. However, our goal is to obtain the solution to (2), or a *L2-regularized* subset selection. Similarly to the unregularized subset selection, it is easy to get that (2) is equivalent to $\max_{|S|=k} \mathbf{b}_S^\top ((\mathbf{C}_S + \lambda \mathbf{I})^{-1})^\top \mathbf{b}_S$. As said above, the Subset Selection problem is **NP-hard**, however, there are several ways to approximate it in practice [36]. We choose FR for this task for its simplicity, appealing computational properties and provably good approximation guarantees. Now, to apply FR to our problem, all we have to do is to provide it with normalized matrix $(\mathbf{C} + \lambda \mathbf{I})^{-1}$ instead of \mathbf{C}^{-1} .

Approximated Randomized Greedy Algorithm. As mentioned above, the complexity of GreedyTL is linear in $d + n$, the number of features and the size of the source hypothesis set. In particular, the search in U for the index to add to S is responsible for the dependency on $d + n$. Here we show how to approximate this search with a randomized strategy. We will use the following Theorem.

Theorem 1 ([40](Theorem 6.33)). *Denote by $M := \{x_1, \dots, x_m\} \subset \mathbb{R}$ a set of cardinality m , and by $\tilde{M} \subset M$ a random subset of size \tilde{m} . Then the probability that $\max \tilde{M}$ is greater or equal than n elements of M is at least $1 - (\frac{n}{m})^{\tilde{m}}$.*

The surprising consequence is that, in order to approximate the maximum over a set, we can use a random subset of size $\mathcal{O}(1)$. In particular, if we want to obtain results in the $\frac{n}{m}$ percentile range with $1 - \eta$ confidence, we use³ $\tilde{m} = \frac{\log(\eta)}{\log \frac{n}{m}}$. Practically, if we desire values that are better than 95% of all other estimates with $1 - 0.05$ probability, then 59 samples are sufficient. This rule is commonly called the 59-trick and it has been widely used to speed-up a wide range of algorithms with negligible loss of accuracy, e.g. [41, 42]. Indeed, as we will show in Section 5.4, we virtually don't lose any accuracy using this strategy.

With the 59-trick, the search in U becomes a search for the maximum over a random set of size 59. So, the overall complexity is reduced to $\mathcal{O}(km^2)$, that is *independent* from all the quantities that are expected to be big.

Theoretical Guarantees. We now focus on the analysis of the generalization properties of GreedyTL for solving k -Source Selection problem (2). Throughout this paragraph we will consider a truncated target predictor $h_{\mathbf{w}, \beta}^{\text{trg}}(\mathbf{x}) := \tau(\mathbf{w}^\top \mathbf{x} + \sum_{i=1}^n \beta_i h_i^{\text{src}}(\mathbf{x}))$, with $\tau(a) := \min\{\max\{a, -1\}, 1\}$. We will also use big-O notation $\tilde{\mathcal{O}}$ to indicate the suppression of a logarithmic factor, in other words, $f(x) \in \tilde{\mathcal{O}}(g(x))$ is a short notation for $\exists n : f(x) \in \mathcal{O}(g(x) \log^n g(n))$. First we state the bound on the risk of an approximate solution returned by GreedyTL.⁴

³Note that the formula for \tilde{m} in [40] contains an error, the correct one is the one we report.

⁴Proofs for theorems can be found in the appendix.

Theorem 2. Let GreedyTL generate the solution $(\hat{\mathbf{w}}, \hat{\beta})$, given the training set (\mathbf{X}, \mathbf{y}) , source hypotheses $\{h_i^{\text{src}}\}_{i=1}^n$ with $\tau_\infty^{\text{src}} := \max_i \{\|h_i^{\text{src}}\|_\infty^2\}$, hyperparameters λ and k . Then with high probability,

$$R(h_{\hat{\mathbf{w}}, \hat{\beta}}^{\text{trg}}) - \hat{R}(h_{\hat{\mathbf{w}}, \hat{\beta}}^{\text{trg}}) \leq \tilde{\mathcal{O}}\left(\frac{1 + k\tau_\infty^{\text{src}}}{\lambda m} + \sqrt{\hat{R}^{\text{src}} \frac{1 + k\tau_\infty^{\text{src}}}{\lambda m}}\right),$$

where $\hat{R}^{\text{src}} := \frac{1}{m} \sum_{i=1}^m \ell\left(y_i, \mathbb{T}\left(\sum_{j \in \text{supp}(\hat{\beta})} \hat{\beta}_j h_j^{\text{src}}(\mathbf{x}_i)\right)\right)$.

This results in a generalization bound which tells us how close the performance of the algorithm on the test set will be to the one on the training set. The key quantity here is \hat{R}^{src} , which captures the quality of the sources selected by the algorithm. To understand its impact, assume that $\lambda = \mathcal{O}(1)$. The bound has two terms, a fast one of the order of $\tilde{\mathcal{O}}(k/m)$ and a slow one of the order $\tilde{\mathcal{O}}\left(\sqrt{\hat{R}^{\text{src}} k/m}\right)$. When m goes to infinity and $\hat{R}^{\text{src}} \neq 0$ the slow term will dominate the convergence rate, giving us a rate of the order of $\tilde{\mathcal{O}}\left(\sqrt{\hat{R}^{\text{src}} k/m}\right)$. If $\hat{R}^{\text{src}} = 0$ the slow term completely disappears, giving us a so called fast rate of convergence of $\tilde{\mathcal{O}}(k/m)$. On the other hand, for any finite m of the order of $\tilde{\mathcal{O}}(k/\hat{R}^{\text{src}})$, we still have a rate of the order of $\tilde{\mathcal{O}}(k/m)$. Hence, the quantity \hat{R}^{src} will govern the finite sample and asymptotic behavior of the algorithm, predicting a faster convergence in both regimes when it is small. In other words, when the source and target tasks are similar, TL facilitates a faster convergence of the empirical risk to the risk. A similar behavior was already observed in [6, 7].

However, one might ask what happens when the selected sources are providing bad predictions. Since $\hat{R}^{\text{src}} \leq 1$, due to truncation, the empirical risk converges to the risk at the standard rate $\tilde{\mathcal{O}}(\sqrt{k/m})$, the same one we would have without any transferring from the sources classifiers.

We now present another result that upper bounds the difference between the risk of solution of the algorithm and the empirical risk of the optimal solution to the k -Source Selection problem.

Theorem 3. In addition to conditions of Theorem 2, let (\mathbf{w}^*, β^*) be the optimal solution to (2). Given a sample correlation matrix $\hat{\mathbf{C}}$, assume that $\hat{C}_{i,j \neq i} \leq \gamma < \frac{1+\lambda}{6k}$, and $\epsilon := \frac{16(k+1)^2 \gamma}{1+\lambda}$. Then with high probability,

$$R(h_{\hat{\mathbf{w}}, \hat{\beta}}^{\text{trg}}) - \hat{R}(h_{\mathbf{w}^*, \beta^*}^{\text{trg}}) \leq (1 + \epsilon) \hat{R}_\lambda^{\text{src}} + \tilde{\mathcal{O}}\left(\frac{1 + k\tau_\infty^{\text{src}}}{\lambda m} + \sqrt{\hat{R}_\lambda^{\text{src}} \frac{1 + k\tau_\infty^{\text{src}}}{\lambda m}}\right),$$

where $\hat{R}_\lambda^{\text{src}} := \min_{|S| \leq k} \left\{ \frac{\lambda}{|S|} + \frac{1}{|S|} \sum_{i \in S} \hat{R}(h_i^{\text{src}}) \right\}$.

To analyze the implications of Theorem 3, let us consider few interesting cases. Similarly as done before, the quantity $\hat{R}_\lambda^{\text{src}}$ captures how well the source hypotheses are aligned with the target task and governs the asymptotic and finite sample regime. In fact, assume for any finite m that there is at least one source hypothesis with small empirical risk, in particular, in $\tilde{\mathcal{O}}(\sqrt{k/m})$, and set $\lambda = \tilde{\mathcal{O}}(\sqrt{k/m})$. Then we have that $R(h_{\hat{\mathbf{w}}, \hat{\beta}}^{\text{trg}}) - \hat{R}(h_{\mathbf{w}^*, \beta^*}^{\text{trg}}) = \tilde{\mathcal{O}}(\sqrt{k/m})$, that is we get the generalization bound as if we are able to solve the original NP-hard problem in (2). In other words, if there are useful source hypotheses, we expect our algorithm to perform similarly to the one that identifies the optimal subset. This might seem surprising, but it is important to note that we do not actually care about identifying the correct subset of source hypotheses. We only care about how well the returned solution is able to generalize. On the other hand, if not even one source hypothesis has low risk, selecting the best subset of k sources becomes meaningless. In this scenario, we expect the selection of any subset to perform in the same way. Thus the approximation guarantee does not matter anymore.

We now state the approximation guarantees of GreedyTL used to prove Theorem 3. In the following Corollary we show how far the optimal solution to the regularized subset selection is from the approximate one found by GreedyTL.

Corollary 1. Let $\lambda \in \mathbb{R}^+$ and $k \leq n$. Denote $\text{OPT} := \min_{\|\mathbf{w}\|_0=k} \left\{ \hat{R}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \right\}$. Assume that $\hat{\mathbf{C}}$ and $\hat{\mathbf{b}}$ are normalized, and $\hat{C}_{i,j \neq i} \leq \gamma < \frac{1+\lambda}{6k}$. Then, FR algorithm generates an approximate solution $\hat{\mathbf{w}}$ to the regularized subset selection problem that satisfies $\hat{R}(\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|_2^2 \leq \left(1 + \frac{16(k+1)^2 \gamma}{1+\lambda}\right) \text{OPT} - \frac{16(k+1)^2 \gamma \lambda}{(1+\lambda)^2}$.

Apart from being instrumental in the proof of Theorem 3, this statement also points to the secondary role of the regularization parameter λ : unlike in FR, we can control the quality of the approximate solution even if the features are correlated.

5 Experiments

In this section we present experiments comparing GreedyTL to several transfer learning and feature selection algorithms. As done previously, we considered the object detection task and, for all datasets, we left out one class considering it as the target class, while the remaining classes were treated as sources [9]. We repeated this procedure for every class and for every dataset at hand, and averaged the performance scores. In the following, we refer to this procedure as *leave-one-class-out*. We performed the evaluation for every class, reporting averaged class-balanced recognition scores.

We used subsets of Caltech-256 [43], Imagenet [2], SUN09 [28], SUN-397 [44]. The largest setting considered involves 1000 classes, totaling in 1.2M examples, where the number of training examples of the target domain varies from 11 to 20. Our experiments aimed at verifying three claims:

- I. L_2 -regularization is important when using greedy feature selection as a transfer learning scheme.
- II. In a small-sample regime GreedyTL is more robust than alternative feature selection approaches, such as L_1 -regularization.
- III. The approximated randomized greedy algorithm improves the computational complexity of GreedyTL with no significant loss in performance.

5.1 Datasets and Features

We used the whole Caltech-256, a public subset of Imagenet containing 10^3 classes, all the classes of SUN09 that have more than 1 example, which amounts to 819 classes, and the whole SUN-397 dataset containing 397 place categories. For Caltech-256 and Imagenet, we used as features the publicly-available 1000-dimensional SIFT-BOW descriptors, while for SUN09 we extracted 3400-dimensional PHOG descriptors. In addition, for Imagenet and SUN-397, we also ran experiments using convolutional features extracted from DeCAF neural network [1].

We composed a negative class by merging 100 held-out classes (*surrogate* negative class). We did so for each dataset, and we further split it into the *source* negative and the *target* negative class as 90% + 10% respectively, for training sources and the target. The source classifiers were trained for each class in the dataset, combining all the positive examples of that class and the source negatives. On average, each source classifier was trained using 10^4 examples for the Caltech-256, 10^5 for Imagenet and 10^3 for the SUN09 dataset. The training sets for the target task were composed by $\{2, 5, 10\}$ positive examples, and 10 negative ones. Following [9], the testing set contained 50 positive and 50 negative examples for Caltech-256, Imagenet, and SUN-397. For the skewed SUN09 dataset we took one positive and 10 negative training examples, with the rest left for testing. We drew each target training and testing set randomly 10 times, averaging the results over them.

5.2 Baselines

We chose a linear SVM to train the source classifiers [45]. This allows us to compare fairly with relevant baselines (like Lasso) and is in line with recent trends in large scale visual recognition and trans-

fer learning [1]. The models were selected by 5-fold cross-validation having regularization parameter $C \in \{10^{-4}, 10^{-3}, \dots, 10^4\}$. In addition to trained source classifiers, for the Caltech-256, we also evaluated transfer from Clasesmes [23] and Object Bank [24], which are very similar in spirit to source classifiers. At the same time, for Imagenet, we evaluated transfer from the outputs of the final layers of the DeCAF convolutional neural network [1].

We divided the baselines into two groups - the linear transfer learning baselines that do not require access to the source data, and the feature selection baselines. We included the second group of baselines due to GreedyTL’s resemblance to a feature selection algorithm. We focus on the linear baselines, since we are essentially interested in the feature selection in high-dimensional spaces from few examples. In that scope, most feature selection algorithms, such as Lasso, are linear. In particular, amongst TL baselines we chose: *No transfer*: Regularized Least Squares (RLS) algorithm trained solely on the target data; *Best source*: indicates the performance of the best source classifier selected by its score on the testing set. This is a pseudo-indicator of what an HTL can achieve; *AverageKT*: obtained by averaging the predictions of all the source classifiers; *RLS src+feat*: RLS trained on the concatenation of feature descriptors and source classifier predictions; *MultiKT* $\|\cdot\|_2$: HTL algorithm by [9] selecting β in (1) by minimizing the leave-one-out error subject to $\|\beta\|_2 \leq \tau$; *MultiKT* $\|\cdot\|_1$: similar to previous, but applying the constraint $\|\beta\|_1 \leq \tau$; *DAM*: An HTL algorithm by [46], that can handle selection from multiple source hypotheses. It was shown to perform better than the well known and similar ASVM [47] algorithm. For the feature selection baselines we selected well-established algorithms involving sparsity assumption: *L1-Logistic*: Logistic regression with $L1$ penalty [14]; *Elastic-Net*: Logistic regression with mixture of $L1$ and $L2$ penalties [14]; *Forward-Reg*: Forward regression – a classical greedy feature selection algorithm. When comparing our algorithms to the baselines on large datasets, we also consider a Domain Adaptive *Dictionary Learning* baseline [48]. This baseline represents the family of dictionary learning methods for domain adaptation and transfer learning. In particular, it learns a dictionary on the source domain and adapts it to the target one. However, in our setup the only access to the source data is through the source hypotheses. Therefore, the only way to construct source features is by using the source hypotheses on the target data points.

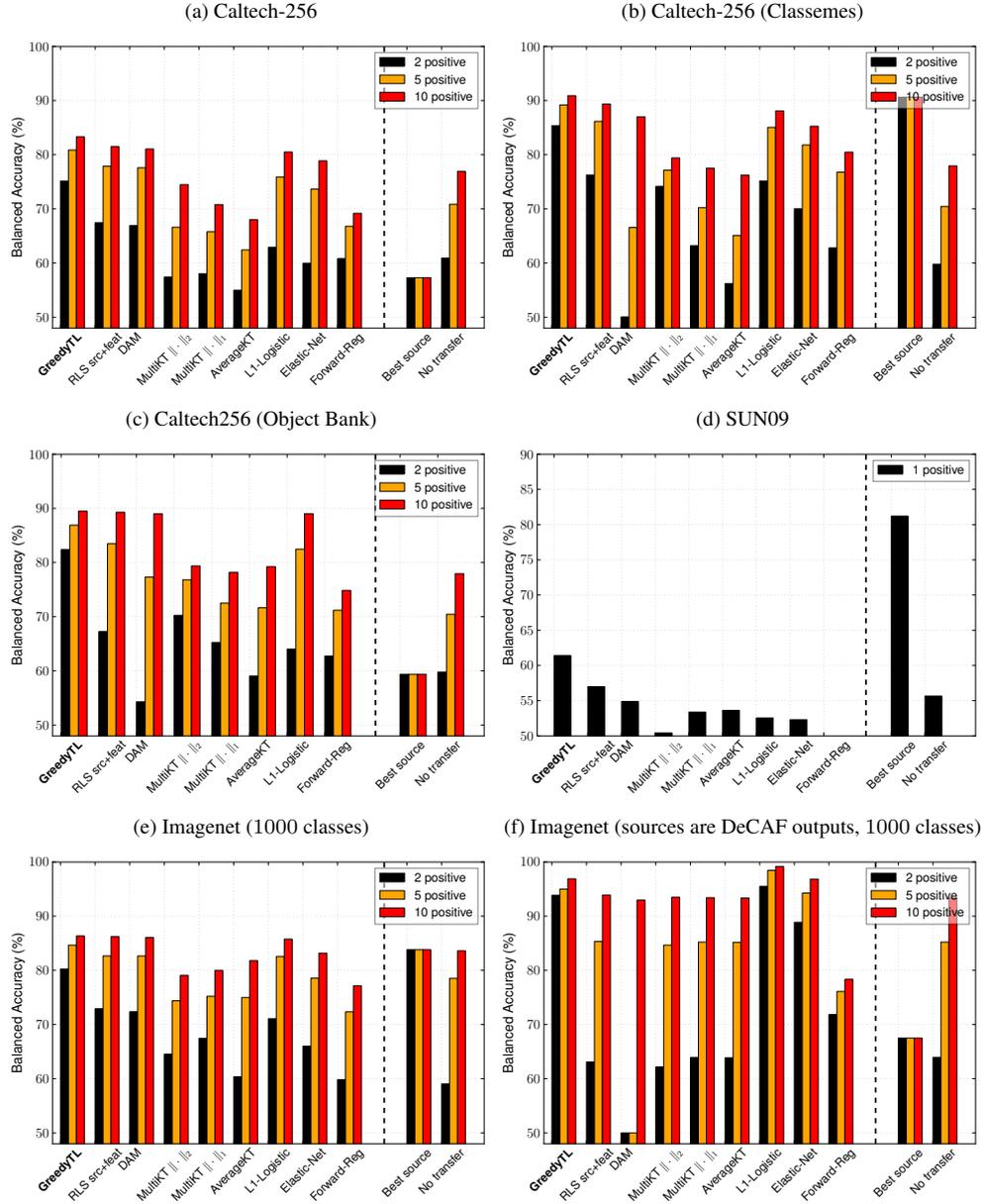
5.3 Results

Figure 1 shows the leave-one-class-out performance. In addition, Figures 1b, 1c, 1f show the performance when transferring from off-the-shelf classesmes, object-bank feature descriptors, and DeCAF neural network activations. Whenever any baseline algorithm has hyperparameters to tune, we chose the ones that minimize the leave-one-out error on the training set. In particular, we selected the regularization parameter $\lambda \in \{10^{-4}, 10^{-3}, \dots, 10^4\}$. MultiKT and DAM have an additional hyperparameter that we call τ with $\tau \in \{10^{-3}, \dots, 10^3\}$. Kernelized algorithms were supplied with a linear kernel. Model selection for GreedyTL involves two hyperparameters, that is k and λ . Instead of fixing k , we let GreedyTL select features as long as the regularized error between two consecutive steps is larger than δ . In particular, we set $\delta = 10^{-4}$, as in preliminary experiments we have not observed any gain in performance past that point. The λ is fixed to 1. Even better performance could be obtained tuning it.

We see that GreedyTL dominates TL and feature selection baselines throughout the benchmark, rarely appearing on-par, especially in the small-sample regime. In addition, on two datasets out of three, it manages to identify the source classifier subset that performs comparably or better than the Best source, that is the single best classifier selected by its performance on the testing set. The significantly stronger performance achieved by GreedyTL w.r.t. FR, on all databases and in all settings, confirms the importance of the regularization in our formulation.

Notably, GreedyTL outperforms RLS src+feat, which is equivalent to GreedyTL selecting all the sources and features. This observation points to the fact that GreedyTL successfully manages to discard irrelevant feature dimensions and sources. To investigate this important point further, we artificially add 10, 100 and 1000 dimensions of pure noise sampled from a standard distribution. Figure 2 compares feature selection methods to GreedyTL in robustness to noise. Clearly, in the small-sample setting, GreedyTL

Figure 1: Performance on the Caltech-256, subsets of Imagenet (1000 classes) and SUN09 (819 classes). Averaged class-balanced accuracies in the leave-one-class-out setting.



is tolerant to large amount of noise, while $L1$ and $L1/L2$ regularization suffer a considerable loss in performance. We also draw attention to the failure of $L1$ -based feature selection methods and MultiKT with $L1$ regularization to match the performance of GreedyTL.

Figure 2: Baselines and number of additional noise dimensions sampled from a standard distribution. Averaged class-balanced recognition accuracies in the leave-one-class-out setting.

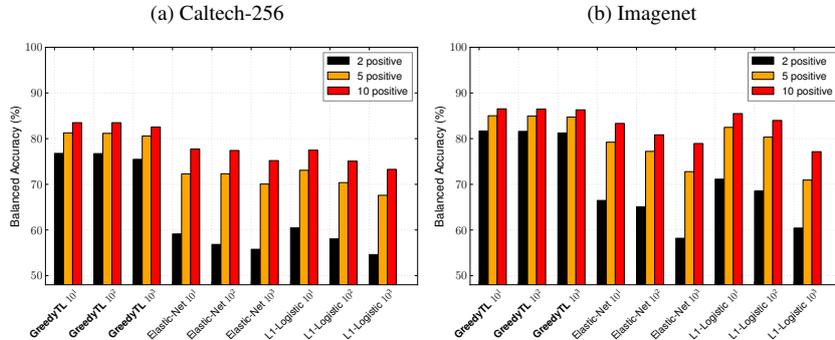


Table 1: Training time in seconds for transferring to a single target class. Results are averaged over 10 splits.

		GreedyTL		
		2 + 10	5 + 10	10 + 10
Training examples pos.+neg.		2 + 10	5 + 10	10 + 10
Imagenet (SIFT-BOW)	1899 source+dim	1.541 ± 0.242	3.083 ± 0.486	5.291 ± 0.870
Imagenet (DECAF7)	4995 source+dim	3.481 ± 0.356	7.492 ± 0.655	13.408 ± 1.165
SUN-397 (Caffe-7)	4492 source+dim	3.245 ± 0.495	6.764 ± 1.051	11.282 ± 1.630

		GreedyTL-59		
		2 + 10	5 + 10	10 + 10
Training examples pos.+neg.		2 + 10	5 + 10	10 + 10
Imagenet (SIFT-BOW)	1899 source+dim	0.043 ± 0.005	0.088 ± 0.011	0.149 ± 0.021
Imagenet (DECAF7)	4995 source+dim	0.055 ± 0.006	0.114 ± 0.013	0.198 ± 0.020
SUN-397 (Caffe-7)	4492 source+dim	0.060 ± 0.021	0.120 ± 0.038	0.198 ± 0.055

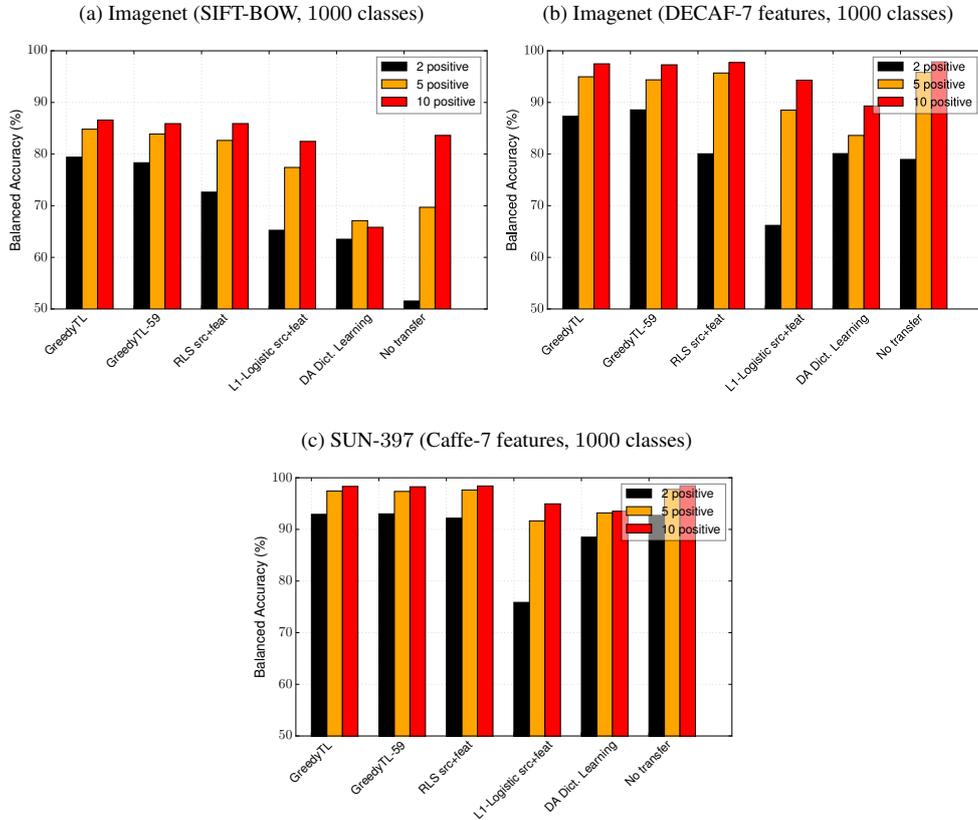
5.4 Approximated GreedyTL

As was discussed in Section 3, the computational complexity of GreedyTL is linear in the number of source hypotheses and feature dimensions. In this section we assess empirical performance of the approximated GreedyTL, which is *independent* from the number of source hypotheses, implemented through the approximated greedy algorithm described at the end of Section 3. In the following we refer to this version of an algorithm as GreedyTL-59. Instead of considering all the transfer learning and feature selection baselines, we restrict the performance comparison to the strongest competitors. To show the power of highly scalable approximated GreedyTL, we focus on the largest datasets in the number of source hypotheses and feature dimensions: Imagenet and SUN-397. In case of Imagenet, we consider standard SIFT-BOW features as in previous section and also DeCAF-7 convolutional features extracted from the seventh layer of the DeCAF neural network [1]. For the SUN-397, we use convolutional features of Caffe network trained on the Places-205 dataset [49], which was shown to perform particularly well in the scene recognition tasks. Figure 3 summarizes new results. Surprisingly, approximated GreedyTL performs on par with the version with exhaustive search over the candidate, maintaining dominant performance in the small-sample regime on the Imagenet dataset. Yet, training timings are dramatically improved as can be seen from Table 1. In the case of SUN-397 dataset, however, GreedyTL performs on par with the top competitors.

5.5 Selected Source Analysis

In this section we take a look at the source hypotheses selected by GreedyTL. In particular, we make a qualitative assessment with the goal to see if semantically related sources and targets are correlated, visualizing selected sources and the magnitude of their weights. We do so by grouping sources and targets semantically according to the WordNet [50] distance, and plotting them as matrices with columns corresponding to targets, rows to sources, and entries to the weights of the sources. Figure 4 shows such matrices

Figure 3: Comparison of the approximated GreedyTL: GreedyTL-59 to GreedyTL with exhaustive search and most competitive baselines on three largest datasets considered in our experiments.



for GreedyTL when evaluated on Imagenet with DECAF7 features and averaged over all splits, for 2 positive and 10 positive examples accordingly. First we note, that for certain supercategories there are clearly distinctive patterns, indicating cross-transfer within the same supercategory. We compare those matrices to the ones originating from the strongest RLS (*src+feat*) baseline, Figure 5. We notice a clear difference, as semantic patterns of GreedyTL are more distinctive in a small-sample setting (2+10), while the ones of RLS (*src+feat*) appear hazier. We argue that this is a consequence of greedy selection procedure implemented by GreedyTL, where sources are selected incrementally, thus many coefficients correspond to zeros. Due to the formulation of RLS (*src+feat*), however, even if a source is less relevant, its coefficient most likely will not be exactly equal to zero.

It is also instructive to compare exact GreedyTL to the approximated one. Figure 7 pictures semantic matrices for the approximated version. We note that approximated version appears to be slightly more conservative in a small-sample case (2+10), but in overall, semantic patterns seem to match, thus emphasizing the quality of the solution provided by the approximated version and empirically corroborating the theoretical motivation behind the randomized selection.

Finally, we take a closer look at some patterns of Figure 4a, that is in the case of learning from only 2 positive examples. This new analysis is shown in Figure 6. We notice that even at the smaller scale, there are emergent semantic patterns.

Figure 4: Semantic transferrability matrix for GreedyTL evaluated on Imagenet (DECAF7 features). Columns correspond to targets and rows to sources. Stronger color intensity means larger source weight. 4a corresponds to learning from 2 positive and 10 negative examples, while 4b, with 10 positive.

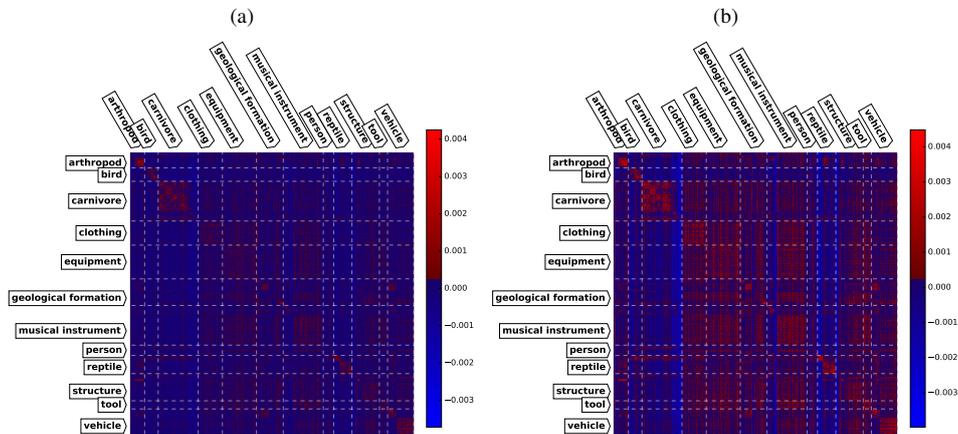
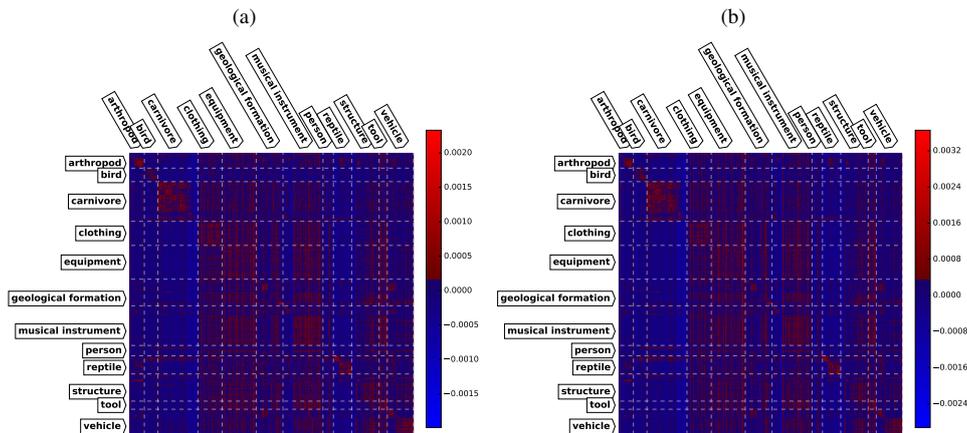


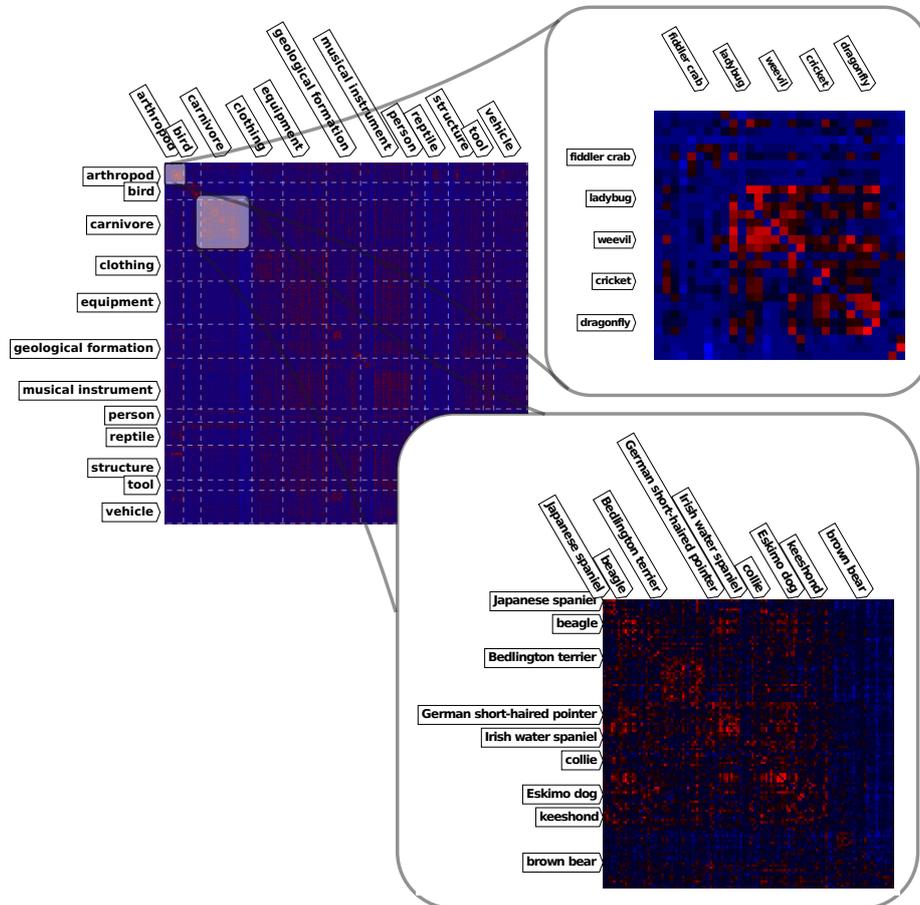
Figure 5: Semantic transferrability matrix for RLS (src+feat) evaluated on Imagenet (DECAF7 features).



6 Conclusions

In this work we studied the transfer learning problem involving hundreds of sources. The kind of transfer learning scenario we consider assumes no access to the source data directly, but through the use of the source hypotheses induced by them. In particular, we focused on the efficient source hypothesis selection and combination, improving the performance on the target task. We proposed a greedy algorithm, GreedyTL, capable of selecting relevant sources and feature dimensions at the same time. We verified these claims by obtaining the best results among the competing feature selection and TL algorithms, on the Imagenet, SUN09 and Caltech-256 datasets. At the same time, comparison against the non-regularized version of the algorithm clearly show the power of our intuition. We support our empirical findings by showing theoretically that under reasonable assumptions on the sources, the algorithm can learn effectively from few

Figure 6: GreedyTL evaluated on Imagenet (DECAF7 features): a closer look at some strongly related sources and targets.



target examples.

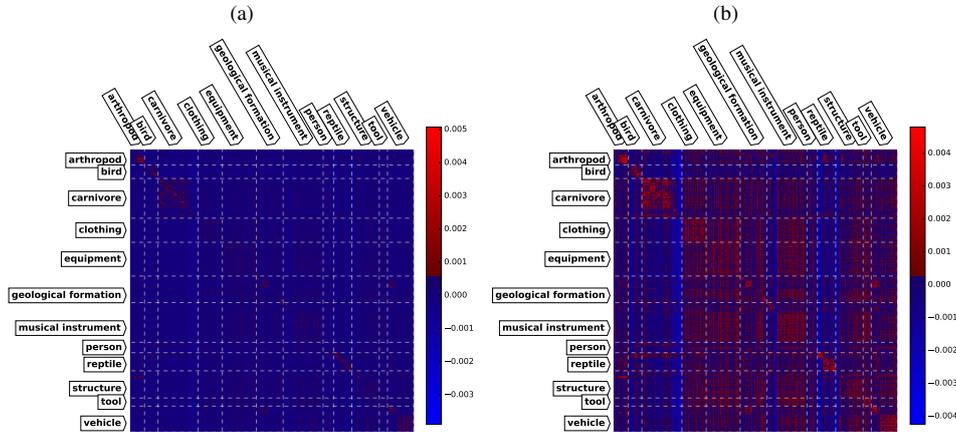
Acknowledgments

This work was partially supported by the ERC grant 367076 -RoboExNovo (B.C. and I. K.).

References

- [1] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *Proceedings of The 31st International Conference on Machine Learning*, pages 647–655, 2014.
- [2] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. IEEE, 2009.
- [3] R. Salakhutdinov, A. Torralba, and J. Tenenbaum. Learning to share visual appearance for multiclass object detection. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1481–1488. IEEE, 2011.

Figure 7: Semantic transferrability matrix for the approximated GreedyTL evaluated on Imagenet (DECAF7 features).



- [4] S. J. Pan and Q. Yang. A survey on transfer learning. *Knowledge and Data Engineering, IEEE Transactions on*, 22(10):1345–1359, 2010.
- [5] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2066–2073. IEEE, 2012.
- [6] I. Kuzborskij and F. Orabona. Stability and hypothesis transfer learning. In *Proceedings of the 30th International Conference on Machine Learning*, pages 942–950, 2013.
- [7] S. Ben-David and R. Uner. Domain adaptation as learning with auxiliary information. *New Directions in Transfer and Multi-Task - Workshop @ NIPS, 2013*.
- [8] Y. Aytar and A. Zisserman. Tabula rasa: Model transfer for object category detection. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2252–2259. IEEE, 2011.
- [9] 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.
- [10] I. Kuzborskij, F. Orabona, and B. Caputo. From N to N+1: Multiclass Transfer Incremental Learning. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 3358–3365. IEEE, 2013.
- [11] A. Das and D. Kempe. Algorithms for subset selection in linear regression. In *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 45–54. ACM, 2008.
- [12] T. Zhang. Adaptive forward-backward greedy algorithm for sparse learning with linear models. In *Advances in Neural Information Processing Systems*, pages 1921–1928, 2009.
- [13] I. Kuzborskij, F. Orabona, and B. Caputo. Transfer learning through greedy subset selection. In *Image Analysis and Processing - ICIAP 2015 - 18th International Conference, Proceedings, Part I*, pages 3–14, 2015.
- [14] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements Of Statistical Learning*. Springer, 2009.
- [15] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *Proceedings of the European Conference on Computer Vision*, pages 213–226. Springer, 2010.
- [16] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J.W. Vaughan. A theory of learning from different domains. *Machine Learning*, 79(1):151–175, 2010.
- [17] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, 2011.
- [18] L. Duan, D. Xu, and S.-F. Chang. Exploiting web images for event recognition in consumer videos: A multiple source domain adaptation approach. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 1338–1345. IEEE, 2012.

- [19] C.-W. Seah, I. W.-H. Tsang, and Y.-S. Ong. Healing sample selection bias by source classifier selection. In *Data Mining (ICDM), 2011 IEEE 11th International Conference on*, pages 577–586. IEEE, 2011.
- [20] T. Tommasi and B. Caputo. Frustratingly easy nbnn domain adaptation. In *Computer Vision (ICCV), IEEE International Conference on*, 2013.
- [21] I. Kuzborskij, F. M. Carlucci, and B. Caputo. When Naïve Bayes Nearest Neighbours Meet Convolutional Neural Networks. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on (in press)*, 2016.
- [22] V. M. Patel, R. Gopalan, R. Li, and R. Chellappa. Visual domain adaptation: A survey of recent advances. *Signal Processing Magazine, IEEE*, 32(3):53–69, 2015.
- [23] A. Bergamo and L. Torresani. Classemes and other classifier-based features for efficient object categorization. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, PP(99):1–1, 2014.
- [24] L.-J. Li, H. Su, L. Fei-Fei, and E. P. Xing. Object bank: A high-level image representation for scene classification & semantic feature sparsification. In *Advances in neural information processing systems*, pages 1378–1386, 2010.
- [25] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1717–1724. IEEE, 2014.
- [26] L. Jie, T. Tommasi, and B. Caputo. Multiclass transfer learning from unconstrained priors. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 1863–1870. IEEE, 2011.
- [27] N. Patricia and B. Caputo. Learning to learn, from transfer learning to domain adaptation: A unifying perspective. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1442–1449. IEEE, 2014.
- [28] M. J. Choi, J. J. Lim, A. Torralba, and A. S. Willsky. Exploiting hierarchical context on a large database of object categories. In *Computer vision and pattern recognition (CVPR), 2010 IEEE conference on*, pages 129–136. IEEE, 2010.
- [29] J. J. Lim, A. Torralba, and R. Salakhutdinov. Transfer learning by borrowing examples for multiclass object detection. In *Advances in Neural Information Processing Systems 24*, pages 118–126, 2011.
- [30] A. Vezhnevets and V. Ferrari. Associative embeddings for large-scale knowledge transfer with self-assessment. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 1987–1994. IEEE, 2014.
- [31] M. Rohrbach, M. Stark, and B. Schiele. Evaluating knowledge transfer and zero-shot learning in a large-scale setting. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1641–1648. IEEE, 2011.
- [32] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *ACM International Conference on Multimedia*, pages 675–678. ACM, 2014.
- [33] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Region-based convolutional networks for accurate object detection and segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, PP(99):1–1, 2015.
- [34] Y. Ganin and V. S. Lempitsky. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015*, pages 1180–1189, 2015.
- [35] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 97–105, 2015.
- [36] A. Das and D. Kempe. Submodular meets spectral: Greedy algorithms for subset selection, sparse approximation and dictionary selection. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 1057–1064. ACM, 2011.
- [37] T. Zhang. On the consistency of feature selection using greedy least squares regression. In *Journal of Machine Learning Research*, pages 555–568, 2009.
- [38] P. Bühlmann and S. Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- [39] Y. Mansour, M. Mohri, and A. Rostamizadeh. Domain Adaptation with Multiple Sources. In *Advances in neural information processing systems*, volume 21, pages 1041–1048, 2009.
- [40] A. Smola and B. Schölkopf. *Learning with Kernels*. MIT press, Cambridge, MA, USA, 2002.

- [41] C. Domingo and O. Watanabe. MadaBoost: A modification of AdaBoost. In *Proceedings of the Thirteenth Annual Conference on Computational Learning Theory (COLT 2000), June 28 - July 1, 2000, Palo Alto, California*, pages 180–189, 2000.
- [42] A. J. Smola and B. Schölkopf. Sparse greedy matrix approximation for machine learning. In *International Conference on Machine Learning, ICML '00*, pages 911–918, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc.
- [43] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. Technical report, Caltech, 2007.
- [44] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Computer vision and pattern recognition (CVPR), 2010 IEEE conference on*, pages 3485–3492. IEEE, 2010.
- [45] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [46] L. Duan, I. W. Tsang, D. Xu, and T.-S. Chua. Domain adaptation from multiple sources via auxiliary classifiers. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 289–296. ACM, 2009.
- [47] J. Yang, R. Yan, and A. G. Hauptmann. Cross-domain video concept detection using adaptive svms. In *Proceedings of the 15th international conference on Multimedia*, pages 188–197. ACM, 2007.
- [48] Q. Qiu, V. M. Patel, P. Turaga, and R. Chellappa. Domain adaptive dictionary learning. In *Proceedings of the European Conference on Computer Vision*, pages 631–645. Springer, 2012.
- [49] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *Advances in Neural Information Processing Systems*, pages 487–495, 2014.
- [50] G. A. Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [51] N. Srebro, K. Sridharan, and A. Tewari. Smoothness, low noise and fast rates. In *Advances in Neural Information Processing Systems*, pages 2199–2207, 2010.
- [52] M. Mohri, A. Rostamizadeh, and A. Talwalkar. *Foundations of machine learning*. The MIT Press, 2012.

A Proofs

In this section we present proofs of theorems. For brevity, we define $\mathbf{h}^{\text{src}}(\mathbf{x}) := [h_1^{\text{src}}(\mathbf{x}), \dots, h_n^{\text{src}}(\mathbf{x})]^\top$, and we will consider a truncated target predictor

$$h_{\mathbf{w}, \boldsymbol{\beta}}^{\text{trg}}(\mathbf{x}) := \tau(\mathbf{w}^\top \mathbf{x} + \boldsymbol{\beta}^\top \mathbf{h}^{\text{src}}(\mathbf{x})),$$

with $\tau(a) := \min\{\max\{a, -1\}, 1\}$. That said, we will assume that

$$\hat{R}(h_{\mathbf{w}, \boldsymbol{\beta}}^{\text{trg}}) \leq \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + \boldsymbol{\beta}^\top \mathbf{h}^{\text{src}}(\mathbf{x}_i) - y_i)^2,$$

in other words, empirical risk of truncated predictor cannot be greater, since all the labels belong to $\{-1, 1\}$.

To prove Theorem 2 we need the following supplementary lemmas.

Lemma 1. *Let GreedyTL generate solution $(\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}})$, given the training set (\mathbf{X}, \mathbf{y}) , source hypotheses $\{h_i^{\text{src}}\}_{i=1}^n$, and hyperparameters λ and k . Then we have that,*

$$\lambda \|\hat{\mathbf{w}}\|^2 + \lambda \|\hat{\boldsymbol{\beta}}\|^2 + \hat{R}(h_{\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}}^{\text{trg}}) \leq \min_{|S| \leq k} \left\{ \frac{1}{|S|} \sum_{j \in S} \hat{R}(h_j^{\text{src}}) + \frac{\lambda}{|S|} \right\},$$

$$\lambda \|\hat{\mathbf{w}}\|^2 + \lambda \|\hat{\boldsymbol{\beta}}\|^2 + \hat{R}(h_{\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}}^{\text{trg}}) \leq \hat{R}(h_{\mathbf{0}, \hat{\boldsymbol{\beta}}}^{\text{trg}}).$$

and also,

$$\lambda \|\hat{\mathbf{w}}\|^2 + \lambda \|\hat{\boldsymbol{\beta}}\|^2 + \hat{R}(h_{\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}}^{\text{trg}}) \leq 1.$$

Proof. Define $J(\mathbf{w}, \boldsymbol{\beta}) := \hat{R}(h_{\mathbf{w}, \boldsymbol{\beta}}^{\text{trg}}) + \lambda \|\mathbf{w}\|^2 + \lambda \|\boldsymbol{\beta}\|^2$. For any $\boldsymbol{\alpha} \in \left\{0, \frac{1}{p}\right\}^n$ such that $\|\boldsymbol{\alpha}\|_0 = p$ we have,

$$\begin{aligned} J(\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}) &\leq J(\mathbf{0}, \boldsymbol{\alpha}) = \frac{1}{m} \sum_{i=1}^m \ell \left(y_i, \frac{1}{p} \sum_{j \in \text{supp}(\boldsymbol{\alpha})} h_j^{\text{src}}(\mathbf{x}_i) \right) + \frac{\lambda}{p} \\ &\leq \frac{1}{p} \sum_{j \in \text{supp}(\boldsymbol{\alpha})} \hat{R}(h_j^{\text{src}}) + \frac{\lambda}{p}. \end{aligned} \quad (3)$$

We have the last inequality due to Jensen's inequality. The fact that (5) holds for any $p \in \{1, \dots, k\}$ proves the first statement.

We have the second statement from,

$$\begin{aligned} \hat{R}(h_{\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}}^{\text{trg}}) + \lambda \|\hat{\mathbf{w}}\|^2 + \lambda \|\hat{\boldsymbol{\beta}}\|^2 &\leq \hat{R}(h_{\mathbf{0}, \boldsymbol{\beta}}^{\text{trg}}) + \lambda \|\hat{\boldsymbol{\beta}}\|^2 \\ \Rightarrow \hat{R}(h_{\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}}^{\text{trg}}) &\leq \hat{R}(h_{\hat{\mathbf{w}}, \hat{\boldsymbol{\beta}}}^{\text{trg}}) + \lambda \|\hat{\mathbf{w}}\|^2 \leq \hat{R}(h_{\mathbf{0}, \hat{\boldsymbol{\beta}}}^{\text{trg}}). \end{aligned}$$

The last statement comes from,

$$\lambda \|\hat{\mathbf{w}}\|^2 + \lambda \|\hat{\boldsymbol{\beta}}\|^2 \leq J(\mathbf{0}, \mathbf{0}) \leq 1. \quad (4)$$

□

Lemma 2. Let $(\mathbf{w}^*, \boldsymbol{\beta}^*)$ be the optimal solution to (3), given the training set (\mathbf{X}, \mathbf{y}) , source hypotheses $\{h_i^{\text{src}}\}_{i=1}^n$, and hyperparameters λ and k . Then, the following holds,

$$\begin{aligned} &\lambda \|\mathbf{w}^*\|^2 + \lambda \|\boldsymbol{\beta}^*\|^2 + \hat{R}(h_{\mathbf{w}^*, \boldsymbol{\beta}^*}^{\text{trg}}) \\ &\leq \min_{|S| \leq k} \left\{ \frac{1}{|S|} \sum_{j \in S} \hat{R}(h_j^{\text{src}}) + \frac{\lambda}{|S|} \right\}. \end{aligned}$$

Proof. Define $J(\mathbf{w}, \boldsymbol{\beta}) := \hat{R}(h_{\mathbf{w}, \boldsymbol{\beta}}^{\text{trg}}) + \lambda \|\mathbf{w}\|^2 + \lambda \|\boldsymbol{\beta}\|^2$. For any $\boldsymbol{\alpha} \in \left\{0, \frac{1}{p}\right\}^n$ such that $\|\boldsymbol{\alpha}\|_0 = p$ we have,

$$\begin{aligned} J(\mathbf{w}^*, \boldsymbol{\beta}^*) &\leq J(\mathbf{0}, \boldsymbol{\alpha}) = \frac{1}{m} \sum_{i=1}^m \ell \left(y_i, \frac{1}{p} \sum_{j \in \text{supp}(\boldsymbol{\alpha})} h_j^{\text{src}}(\mathbf{x}_i) \right) + \frac{\lambda}{p} \\ &\leq \frac{1}{p} \sum_{j \in \text{supp}(\boldsymbol{\alpha})} \hat{R}(h_j^{\text{src}}) + \frac{\lambda}{p}. \end{aligned} \quad (5)$$

We have the last inequality due to Jensen's inequality. The fact that (5) holds for any $p \in \{1, \dots, k\}$ proves the statement.

□

Proof of Theorem 2. To prove the statement we will use the optimistic rate Rademacher complexity bounds of [51]. In particular, we will have to do two things: upper-bound the worst-case Rademacher complexity of the hypothesis class of GreedyTL, and upper-bound the empirical risk of members of that hypothesis class. Before proceeding, we spend a moment to define the loss class of GreedyTL, assuring that it is consistent with the definition by [51],

$$\mathcal{L} := \left\{ (\mathbf{x}, y) \mapsto \frac{1}{2} (h(\mathbf{x}) - y)^2 : h \in (\mathcal{T} \circ \mathcal{H}), \hat{R}(h) \leq r \right\}. \quad (6)$$

Here, $(\mathsf{T} \circ \mathcal{H})$ is the class of truncated hypotheses, \mathcal{H} is the hypothesis class of `GreedyTL` and r is the mentioned bound on the empirical risk. We define the hypothesis class as,

$$\mathcal{H} := \left\{ \mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} + \boldsymbol{\beta}^\top \mathbf{h}^{\text{src}}(\mathbf{x}) : \|\mathbf{w}\|_2^2 + \|\boldsymbol{\beta}\|_2^2 \leq \frac{1}{\lambda} \right\}.$$

In this definition we have used the fact shown in Lemma 1, that is the constraint on $\|\mathbf{w}\|_2^2 + \|\boldsymbol{\beta}\|_2^2$, which translates into a constraint on the hypothesis class. Now we are ready to analyze its complexity.

Recall that the worst case Rademacher complexity is defined as,

$$\mathfrak{R}(\mathcal{F}) := \sup_{\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{X}} \left\{ \mathbb{E}_\sigma \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{m} \sum_{i=1}^m \sigma_i f(\mathbf{x}_i) \right\} \right] \right\},$$

where σ_i is r.v., such that $\mathbb{P}(\sigma_i = 1) = \mathbb{P}(\sigma_i = -1) = \frac{1}{2}$.

Let us focus on the analysis of empirical Rademacher complexity $\hat{\mathfrak{R}}(\mathsf{T} \circ \mathcal{H})$, that is the part inside the outer supremum. The truncation $\mathsf{T}(\cdot)$ is 1-Lipschitz, therefore by Talagrand's contraction lemma [52] we have that $\hat{\mathfrak{R}}(\mathsf{T} \circ \mathcal{H}) \leq \hat{\mathfrak{R}}(\mathcal{H})$. Hence, now we proceed with an upper-bound on $\hat{\mathfrak{R}}(\mathcal{H})$. Define $\boldsymbol{\iota} \in \{0, 1\}^n$ such that $\iota_i := \begin{cases} 1, & i \in \text{supp}(\boldsymbol{\beta}) \\ 0, & \text{otherwise} \end{cases}$. Then we have that,

$$\hat{\mathfrak{R}}(\mathsf{T} \circ \mathcal{H}) \leq \hat{\mathfrak{R}}(\mathcal{H}) \tag{7}$$

$$\begin{aligned} &= \mathbb{E}_\sigma \left[\sup_{\|\mathbf{w}\|_2^2 + \|\boldsymbol{\beta}\|_2^2 \leq \frac{1}{\lambda}} \frac{1}{m} \sum_{i=1}^m \sigma_i (\mathbf{w}^\top \mathbf{x}_i + \boldsymbol{\beta}^\top \mathbf{h}^{\text{src}}(\mathbf{x}_i)) \right] \\ &= \frac{1}{m\sqrt{\lambda}} \mathbb{E}_\sigma \left[\left\| \sum_{i=1}^m \sigma_i \begin{bmatrix} \mathbf{x}_i \\ \boldsymbol{\iota} \circ \mathbf{h}^{\text{src}}(\mathbf{x}_i) \end{bmatrix} \right\| \right] \end{aligned} \tag{8}$$

$$\leq \sqrt{\frac{1}{m^2 \lambda} \sum_{i=1}^m \|\mathbf{x}_i\|^2 + \|\boldsymbol{\iota} \circ \mathbf{h}^{\text{src}}(\mathbf{x}_i)\|^2} \tag{9}$$

$$\leq \sqrt{\frac{1 + k \|\mathbf{h}^{\text{src}}\|_\infty^2}{\lambda m}}. \tag{10}$$

To obtain (8) we have applied Cauchy-Schwartz inequality on the inner product of $[\mathbf{w}^\top \boldsymbol{\beta}^\top]^\top$ and $[\mathbf{x}_i^\top \mathbf{h}^{\text{src}}(\mathbf{x}_i)^\top]^\top$, then upper-bounding norms with constraints given by definition of a class \mathcal{H} . To get (9) we have applied Jensen's inequality w.r.t. $\mathbb{E}[\cdot]$, along with the fact that $\mathbb{E}[\sigma_i \sigma_{j \neq i}] = 0$ and $\mathbb{E}[\sigma_i \sigma_i] = 1$. Next, we have bounded the $L2$ norms of features and sources, recalling that by assumption, $\|\mathbf{x}_i\|^2 \leq 1$. Finally, taking supremum over (10) w.r.t. data, we obtain,

$$\mathfrak{R}(\mathsf{T} \circ \mathcal{H}) \leq \sqrt{\frac{1 + k \|\mathbf{h}^{\text{src}}\|_\infty^2}{\lambda m}}.$$

Next, we upper bound the empirical risk of the members of \mathcal{H} by Lemma 1. By plugging the bound on the $\mathfrak{R}(\mathcal{H})$, and the bound on the empirical risk of (6) into Theorem 1 in [51] we have the statement. \square

Next we prove the approximation guarantee of a Regularized Subset Selection (RSS), Corollary 1, that is needed for proof of Theorem 3. First we note that the solution returned by FR enjoys the following guarantees in solving the Subset Selection.

Theorem 4 ([11]). *Assume that \mathbf{C} and \mathbf{b} are normalized, and $C_{i,j \neq i} \leq \gamma < \frac{1}{6k}$ for subset size $k \leq n$. Then, the FR algorithm generates an approximate solution $\hat{\mathbf{w}}$ to the Subset Selection such that, $R(\hat{\mathbf{w}}) \leq (1 + 16(k+1)^2 \gamma) \min_{\|\mathbf{w}\|_0=k} R(\mathbf{w})$.*

This theorem is instrumental in stating our corollary.

Proof of Corollary 1. In addition to the sample covariance matrix $\hat{\mathbf{C}}$, define also correlations $\mathbf{b} := \frac{1}{m} \mathbf{X}^\top \mathbf{y}$. Denote $\hat{\mathbf{C}}' = \frac{\hat{\mathbf{C}} + \lambda \mathbf{I}}{1 + \lambda}$. Now, suppose that $\hat{\mathbf{w}}_S$ is the solution found by the forward regression algorithm, given the input $(\hat{\mathbf{C}}', \hat{\mathbf{b}}, k)$. So, the empirical risk that the algorithm attains is $1 - \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S$, as follows from the analytic solution to empirical risk minimization for given S . In fact, we can upper-bound it right away using Theorem 4. But, recall that our goal is to upper-bound the quantity $\hat{R}(\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|^2 = 1 - \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}_S + \lambda \mathbf{I})^{-1} \hat{\mathbf{b}}_S$, that is the regularized empirical risk of the approximation $\hat{\mathbf{w}}_S$ to the RSS. This quantity is obtained via the unnormalized covariance matrix, therefore we cannot analyze it directly by Theorem 4. For this reason we rewrite it as $\hat{R}(\hat{\mathbf{w}}) + \lambda \|\hat{\mathbf{w}}\|^2 = 1 - \frac{1}{1 + \lambda} \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S$. From Theorem 4 assumptions we then have $(\hat{\mathbf{C}}_S)'_{i,j \neq i} \leq \gamma' \leq \frac{1}{6k}$, denote $\epsilon = 16(k+1)^2 \gamma'$, and let S^* be the optimal subset of size k . Now we plug $1 - \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S$ into Theorem 4, and proceed with algebraic transformations,

$$\begin{aligned} 1 - \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S &\leq (1 + \epsilon) (1 - \hat{\mathbf{b}}_{S^*}^\top (\hat{\mathbf{C}}'_{S^*})^{-1} \hat{\mathbf{b}}_{S^*}) \\ &\Rightarrow \frac{1}{1 + \lambda} (1 - \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S) \leq \frac{1 + \epsilon}{1 + \lambda} (1 - \hat{\mathbf{b}}_{S^*}^\top (\hat{\mathbf{C}}'_{S^*})^{-1} \hat{\mathbf{b}}_{S^*}) \\ &\Rightarrow 1 - \frac{1}{1 + \lambda} \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S \end{aligned} \quad (11)$$

$$\begin{aligned} &\leq (1 + \epsilon) \left(\frac{1}{1 + \lambda} - \frac{1}{1 + \lambda} \hat{\mathbf{b}}_{S^*}^\top (\hat{\mathbf{C}}'_{S^*})^{-1} \hat{\mathbf{b}}_{S^*} \right) + \frac{\lambda}{1 + \lambda} \\ &\Rightarrow 1 - \frac{1}{1 + \lambda} \hat{\mathbf{b}}_S^\top (\hat{\mathbf{C}}'_S)^{-1} \hat{\mathbf{b}}_S \end{aligned} \quad (12)$$

$$\begin{aligned} &\leq (1 + \epsilon) \left(1 - \frac{1}{1 + \lambda} \hat{\mathbf{b}}_{S^*}^\top (\hat{\mathbf{C}}'_{S^*})^{-1} \hat{\mathbf{b}}_{S^*} \right) - \frac{\epsilon \lambda}{1 + \lambda}. \end{aligned}$$

The last step is to relate γ' to γ . The fact $(\hat{\mathbf{C}}_S)'_{i,j \neq i} \leq \gamma' \leq \frac{1}{6k}$ is equivalent to $\frac{(\hat{\mathbf{C}}_S)_{i,j \neq i}}{1 + \lambda} \leq \gamma' \leq \frac{1}{6k}$. Therefore we can set $\gamma = \gamma'(1 + \lambda)$ and obtain $(\hat{\mathbf{C}}_S)_{i,j \neq i} \leq \gamma \leq \frac{1 + \lambda}{6k}$. This concludes the proof. \square

Proof of Theorem 3. The proof follows the composition of Theorem 2, Corollary 1 and Lemma 2. In particular, we upper-bound the empirical risk of Theorem 2 with an approximation given by Corollary 1, ignoring the negative term. Next, we upper-bound $\epsilon(\lambda \|\mathbf{w}^*\|^2 + \lambda \|\beta^*\|^2 + \hat{R}(h_{\mathbf{w}^*, \beta^*}^{\text{reg}})) + \lambda \|\mathbf{w}^*\|^2 + \lambda \|\beta^*\|^2$ by Lemma 2. \square

The following proposition is used to derive the GreedyTL in Section 4.

Proposition 1. Define the regularized accuracy as,

$$\hat{A}^\lambda(\mathbf{w}) := 1 - \left(\frac{1}{m} \|\mathbf{X}^\top \mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \right).$$

We are given $\mathbf{X} \in \mathbb{R}^{n \times m}$, $\mathbf{y} \in \mathbb{R}^m$, $S \subseteq \{1, \dots, n\}$, and $\lambda \in \mathbb{R}^+$. Furthermore, assume that $\frac{\|\mathbf{y}\|_2^2}{m} = 1$, and let $\hat{\mathbf{X}}$ be the submatrix of \mathbf{X} , selecting rows indexed by S . Then we have that,

$$\max_{\mathbf{w}, \text{supp}(\mathbf{w})=S} \left\{ \hat{A}^\lambda(\mathbf{w}) \right\} = \frac{1}{m} \mathbf{y}^\top \hat{\mathbf{X}}^\top (\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m \lambda \mathbf{I})^{-1} \hat{\mathbf{X}} \mathbf{y} \quad (13)$$

$$= \frac{1}{m} \mathbf{y}^\top (\hat{\mathbf{X}}^\top \hat{\mathbf{X}} + m \lambda \mathbf{I})^{-1} \hat{\mathbf{X}}^\top \hat{\mathbf{X}} \mathbf{y}. \quad (14)$$

Proof. Expanding the $\|\cdot\|^2$ in $\hat{A}^\lambda(\mathbf{w})$ and using the fact that $\frac{\|\mathbf{y}\|^2}{m} = 1$, gives us

$$\hat{A}^\lambda(\mathbf{w}) = \frac{2}{m} \mathbf{w}^\top \hat{\mathbf{X}} \mathbf{y} - \frac{1}{m} \mathbf{w}^\top (\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m\lambda \mathbf{I}) \mathbf{w}.$$

Now we have that $\frac{\partial \hat{A}^\lambda(\mathbf{w})}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = (\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m\lambda \mathbf{I})^{-1} \hat{\mathbf{X}} \mathbf{y}$. Denote $\mathbf{G} = (\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m\lambda \mathbf{I})^{-1}$ and set optimal solution $\mathbf{w}^* = \mathbf{G} \hat{\mathbf{X}} \mathbf{y}$. By putting \mathbf{w}^* into the objective we have,

$$\begin{aligned} \hat{A}^\lambda(\mathbf{w}^*) &= \frac{2}{m} \mathbf{y}^\top \hat{\mathbf{X}}^\top \mathbf{G}^\top \hat{\mathbf{X}} \mathbf{y} - \frac{1}{m} \mathbf{y}^\top \hat{\mathbf{X}}^\top \mathbf{G}^\top \mathbf{G}^{-1} \mathbf{G} \hat{\mathbf{X}} \mathbf{y} \\ &= \frac{1}{m} \mathbf{y}^\top \hat{\mathbf{X}}^\top \mathbf{G}^\top \hat{\mathbf{X}} \mathbf{y}. \end{aligned}$$

This proves the first statement.

Now we turn to the second statement, that is solution in the dual variables. By using dual variable identity $(\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m\lambda \mathbf{I})^{-1} \hat{\mathbf{X}} = \hat{\mathbf{X}} (\hat{\mathbf{X}}^\top \hat{\mathbf{X}} + m\lambda \mathbf{I})^{-1}$ [52], we write solution w.r.t. \mathbf{w} as $\mathbf{w} = \hat{\mathbf{X}} (\hat{\mathbf{X}}^\top \hat{\mathbf{X}} + m\lambda \mathbf{I})^{-1} \mathbf{y}$. Denoting $\mathbf{G} = (\hat{\mathbf{X}}^\top \hat{\mathbf{X}} + m\lambda \mathbf{I})^{-1}$, setting optimal solution $\mathbf{w}^* = \hat{\mathbf{X}} \mathbf{G} \mathbf{y}$, and putting \mathbf{w}^* into the objective we have,

$$\begin{aligned} \hat{A}^\lambda(\mathbf{w}^*) &= \frac{2}{m} \mathbf{y}^\top \mathbf{G}^\top \hat{\mathbf{X}}^\top \hat{\mathbf{X}} \mathbf{y} \\ &\quad - \frac{1}{m} \mathbf{y}^\top \mathbf{G}^\top \hat{\mathbf{X}}^\top (\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m\lambda \mathbf{I}) \hat{\mathbf{X}} \mathbf{G} \mathbf{y} = \frac{1}{m} \mathbf{y}^\top \mathbf{G} \hat{\mathbf{X}}^\top \hat{\mathbf{X}} \mathbf{y}. \end{aligned}$$

The last fact comes from the observation that $\hat{\mathbf{X}} \mathbf{G} = (\hat{\mathbf{X}} \hat{\mathbf{X}}^\top + m\lambda \mathbf{I})^{-1} \hat{\mathbf{X}}$ by dual variable identity. This concludes the proof of the second statement. \square