

Research Article

Least Square Regularized Regression for Multitask Learning

Yong-Li Xu,¹ Di-Rong Chen,² and Han-Xiong Li³

¹ Department of Mathematics, Beijing University of Chemical Technology, Beijing 100029, China

² Department of Mathematics, Beijing University of Aeronautics and Astronautics, Beijing 100091, China

³ Department of Systems Engineering and Engineering Management, City University of Hong Kong, Hong Kong

Correspondence should be addressed to Yong-Li Xu; xuyongli2312@sina.com

Received 11 October 2013; Accepted 13 November 2013

Academic Editor: Yiming Ying

Copyright © 2013 Yong-Li Xu et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The study of multitask learning algorithms is one of very important issues. This paper proposes a least-square regularized regression algorithm for multi-task learning with hypothesis space being the union of a sequence of Hilbert spaces. The algorithm consists of two steps of selecting the optimal Hilbert space and searching for the optimal function. We assume that the distributions of different tasks are related to a set of transformations under which any Hilbert space in the hypothesis space is norm invariant. We prove that under the above assumption the optimal prediction function of every task is in the same Hilbert space. Based on this result, a pivotal error decomposition is founded, which can use samples of related tasks to bound excess error of the target task. We obtain an upper bound for the sample error of related tasks, and based on this bound, potential faster learning rates are obtained compared to single-task learning algorithms.

1. Introduction

Multitask learning [1] is a learning paradigm which seeks to improve the generalization performance of a learning task with the help of some other related tasks. This learning paradigm is inspired by human learning activities in that people often apply the knowledge gained from previous learning tasks to help learn a new task [2]. Different multitask learning algorithms have been designed, such as multitask support vector machine (SVM) [3], multitask feature learning [4, 5], multitask clustering approach [6], multitask structure learning [7], and multitask gradients learning [8].

Multitask learning can be formulated under two different settings: symmetric and asymmetric [9]. The symmetric multitask learning tries to improve the performance of all tasks simultaneously, and the objective of asymmetric multitask learning tries to improve the performance of some target task using information from related tasks. The asymmetric multitask learning is related to transfer learning [10]. The major difference is that the source tasks are still learned simultaneously in asymmetric multitask learning while they are learned independently in transfer learning [2]. Much experimental work has achieved the target that is improving the

prediction performance of a learning task with the help of other related tasks [1, 11–15]. However, there has been relatively little progress on theoretical analysis for these results.

Baxter presented a general frame for model selection in multitask learning environment [16]. They showed that a hypothesis space that performs well on a sufficiently large number of training tasks will also perform well when learning novel tasks in the same environment. They proved that learning multiple tasks within an environment of related tasks can potentially give much better generalization than learning a single task. Ando and Zhang considered learning predictive structures on hypothesis spaces from multitask learning [17]. They presented a general framework in which the structural learning problem can be formulated and analyzed theoretically and related it to learning with unlabeled data. Ben-David and Borbely defined relatedness of tasks on the basis of similarity between the example generating distributions for classification [18], and then they gave precise conditions under which bounds guarantee generalization on the basis of smaller sample sizes than the standard single-task approach. Solnon et al. studied multitask regression, using penalization techniques [19]. They showed that the key element appearing for an optimal calibration is the covariance matrix of the

noise between the different tasks. They presented a new algorithm to estimate this covariance matrix and proved that this estimator converges towards the covariance matrix.

In this paper, we propose a least-square regularized regression algorithm for multitask learning with hypothesis space being the union of a sequence of Hilbert spaces. The relatedness of tasks is described by distributions that underlie these tasks and some property of the hypothesis space. We assume that the distributions are related by a set of transformations under which the norm of any Hilbert space in the hypothesis space is invariant. We design a multitask learning algorithm with two steps: firstly, samples of other related tasks are used to select an approximate optimal Hilbert space in the hypothesis space; secondly, in the optimal Hilbert space, we use standard least square regularized regression algorithm for the target task. It is proved that, under the above assumption, the optimal prediction function of every task is in the same Hilbert space. For error analysis, we decompose the excess error of prediction function in target task into regularization error and sample error in which the difference between error and empirical error of the prediction function in the target task is estimated by the average value of those in related tasks. This leads to a potential faster learning rate than that of standard regularized regression algorithm in single task.

The rest of the paper is organized as follows. In Section 2, we introduce some notions and definitions and then propose the least square regularized regression for multitask learning. In Section 3, we decompose excess error of target task into regularization error and sample error in which samples of other related tasks can be used to estimated difference between error and empirical error of the prediction function in target task. The main result is presented in Section 4. An upper bound for sample error of multiple tasks is given by Hoeffding's inequality and an estimation for covering number in multitask learning, and then, based on the upper bound, potential faster learning rates compared to single-task learning algorithms are obtained.

2. Preliminaries

To propose regularized regression algorithm for multitask learning, we introduce some definitions and notations. Let (X, d) be a compact metric space and let $Y = [-M, M]$. Let ρ be a probability distribution on $Z := X \times Y$. The regression function is defined as

$$f_\rho(x) = \int_Y y d\rho(y | x), \quad x \in X, \quad (1)$$

where $\rho(y | x)$ is the conditional probability measure at x induced by ρ . Knowing a set of samples from the probability distribution ρ , our goal is to find a good approximation of f_ρ .

For multitask learning, we define the relatedness of probability distribution of multiple tasks.

Definition 1. For a function $f : X \rightarrow X$, let $f[P]$ be the probability distribution over $X \times Y$ defined by $f[P](T) = \rho(\{(f(x), b) | (x, b) \in T\})$, for $T \subseteq X \times Y$. Let \mathcal{F} be a set of transformations $f : X \rightarrow X$ and let ρ_1 and ρ_2 be probability

distributions over $X \times Y$. We say that ρ_1 and ρ_2 are \mathcal{F} -related if there exists some $f \in \mathcal{F}$ such that $\rho_1 = f[\rho_2]$ or $\rho_2 = f[\rho_1]$.

Then, we describe the relatedness of the transformation set \mathcal{F} above and a hypothesis space.

Definition 2. Let \mathcal{F} be a set of transformations $f : X \rightarrow X$ and let \mathcal{H} be a set of functions $X \rightarrow Y$. We say that \mathcal{F} acts as a group over \mathcal{H} , if,

- (1) for every $f \in \mathcal{F}$ and every $h \in \mathcal{H}$, there holds $h \circ f \in \mathcal{H}$;
- (2) for every $f, g \in \mathcal{F}$, the inverse transformation f^{-1} and the composition $f \circ g$ are also members of \mathcal{F} .

Definition 3. Let \mathcal{H}_σ be a Hilbert space with norm $\|\cdot\|_\sigma$ and let \mathcal{F} act as a group over \mathcal{H}_σ . We say \mathcal{H}_σ is norm invariant under \mathcal{F} , if, for any $h \in \mathcal{H}_\sigma$ and any $f \in \mathcal{F}$, there holds

$$\|h\|_\sigma = \|h \circ f\|_\sigma. \quad (2)$$

To explain the above definitions, we give an example.

Example 4. Let $X = \mathbb{R}^2$, $Y = [-M, M]$, $\Gamma = \mathbb{R}$, and $K_\sigma(x, x') = \exp\{-\|x - x'\|_2^2/\sigma^2\}$, for $x, x' \in X$, and \mathcal{H}_σ be the closure of linear span of the set $\{K_{\sigma, x} := K_\sigma(x, \cdot) : x \in X\}$ with inner product $\langle K_{\sigma, x}, K_{\sigma, x'} \rangle_\sigma = K_\sigma(x, x')$, for any $\sigma \in \Gamma$. Let the norm $\|\cdot\|_\sigma$ of functions in \mathcal{H}_σ be induced by the inner product.

Assume \mathcal{F} is a set of translation and rotation transformations on $\{\mathcal{H}_\sigma, \sigma \in \Gamma\}$: for any $f \in \mathcal{F}$ and $h \in \{\mathcal{H}_\sigma, \sigma \in \Gamma\}$, there holds $(h \circ f)(x) := h(x + x_f)$, for some $x_f \in X$, or $(h \circ f)(x) := h(A_f x)$, for

$$A_f = \begin{bmatrix} \cos \theta_f & -\sin \theta_f \\ \sin \theta_f & \cos \theta_f \end{bmatrix}, \quad (3)$$

where θ_f is an angle dependent on f . We can verify that \mathcal{F} acts as a group over \mathcal{H}_σ and \mathcal{H}_σ is norm invariant under \mathcal{F} , for any $\sigma \in \Gamma$.

Now we introduce standard least-square regularized regression algorithm. We denote error and empirical error of a function $f : X \rightarrow Y$ with squared loss as follows. For a distribution ρ on $X \times Y$, the error of f is defined as

$$E^\rho(f) = \int_Z (f(x) - y)^2 d\rho. \quad (4)$$

It is well known that the regression function minimizes the error. Indeed,

$$\|f - f_\rho\|_{\rho_X}^2 = \mathcal{E}(f) - \mathcal{E}(f_\rho), \quad (5)$$

where ρ_X is the marginal distribution of ρ on X and $\|f\|_{\rho_X}^2 = \int_X |f(x)|^2 d\rho_X$. The above difference is called the excess error of f , and for a sample set $S = \{(x_i, y_i), i = 1, 2, \dots, m\}$, independently drawn according to ρ , the empirical error of f is defined as

$$\hat{E}^S(f) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2. \quad (6)$$

In this paper, we consider a sequence of related learning tasks. The goal is to use information of related tasks to improve learning performance of one special learning task. Let \mathcal{F} be a transformation set and \mathcal{H}_σ be a Hilbert space with norm $\|\cdot\|_\sigma$, for any $\sigma \in \Gamma$, where Γ is an index set. We assume that there is $\kappa > 0$, such that

$$\|f\|_\infty \leq \kappa \|f\|_\sigma. \quad (7)$$

Let ρ_i denote the probability distribution of the i th task, for $i = 1, \dots, n$. We assume $\{\rho_i\}_{i=1}^n$ is pair-wise \mathcal{F} -related, \mathcal{F} acts as a group over \mathcal{H}_σ , and \mathcal{H}_σ is norm invariant under \mathcal{F} , for any $\sigma \in \Gamma$. Let $S_i = \{(x_i^j, y_i^j)\}_{j=1}^m$ be samples independently drawn according to ρ_i . Since the n tasks are related, we try to use samples $\mathbf{z} = \cup_{i=1}^n S_i$ of all the n tasks to improve the learning performance of the target task.

Standard least-square regularized regression algorithm associated with $\{\mathcal{H}_\sigma : \sigma \in \Gamma\}$ for the i th single task is defined as the minimizer

$$f_{S_i} = \arg \min_{f \in \mathcal{H}_\sigma, \sigma \in \Gamma} \{\hat{E}^{S_i}(f) + \lambda \|f\|_\sigma^2\}. \quad (8)$$

In the above optimization, there are two steps; firstly, for any fixed $\sigma \in \Gamma$, find the optimal function $f_{S_i, \sigma}$; secondly, find the global optimal function f_S , for $\sigma \in \Gamma$. Since the n tasks are related, we try to use samples $\mathbf{z} = \cup_{i=1}^n S_i$ of all the n tasks to improve the learning performance of the target task. Without loss of generality, we choose the first task as the target task.

Now, we propose least square regularized regression algorithm for multitask learning.

Step 1. Use samples of other $n - 1$ tasks to select the approximate optimal $\mathcal{H}_{\hat{\sigma}}$ as follows:

$$\hat{\sigma} = \arg \min_{\substack{\sigma \in \Gamma, \\ h_2, \dots, h_n \in \mathcal{H}_\sigma}} \sum_{i=2}^n \{\hat{E}^{S_i}(h_i) + \lambda \|h_i\|_\sigma^2\}. \quad (9)$$

Step 2. In $\mathcal{H}_{\hat{\sigma}}$, search for the approximation \hat{h}_1 to f_{ρ_1} as follows:

$$(\hat{h}_1, \dots, \hat{h}_n) = \arg \min_{h_1, \dots, h_n \in \mathcal{H}_{\hat{\sigma}}} \sum_{i=1}^n \{\hat{E}^{S_i}(h_i) + \lambda \|h_i\|_{\hat{\sigma}}^2\}. \quad (10)$$

3. Error Decomposition

To estimate the bound of excess error $E^{\rho_1}(\hat{h}_1) - E^{\rho_1}(f_{\rho_1})$, we introduce some notations. Let σ_i^* be the optimal $\sigma \in \Gamma$ for the i th task,

$$(\sigma_i^*, \bar{h}_i^*) = \arg \min_{\substack{\sigma \in \Gamma, \\ h \in \mathcal{H}_\sigma}} \{E^{\rho_i}(h) + \lambda \|h\|_\sigma^2\}. \quad (11)$$

Lemma 5. Let ρ_i denote the probability distribution of the i th task for $i = 1, \dots, n$, let \mathcal{H}_σ be a Hilbert space with norm $\|\cdot\|_\sigma$ for $\sigma \in \Gamma$, and let \mathcal{F} be a transformation set. Assume $\{\rho_i\}_{i=1}^n$ is pair-wise \mathcal{F} -related, \mathcal{F} acts as a group over \mathcal{H}_σ , and \mathcal{H}_σ is norm invariant under \mathcal{F} , for any $\sigma \in \Gamma$. Then σ_i^* defined in (11) satisfies

$$\sigma_1^* = \sigma_2^* = \dots = \sigma_n^*. \quad (12)$$

Proof. By the assumption that ρ_i and ρ_j are \mathcal{F} -related, it is easy to show that

$$\inf_{h \in \mathcal{H}_\sigma} E^{\rho_i}(h) = \inf_{h \in \mathcal{H}_\sigma} E^{\rho_j}(h), \quad (13)$$

for any $i, j \in \{1, \dots, n\}$ and $\sigma \in \Gamma$. Notice that \mathcal{H}_σ is norm-invariant under \mathcal{F} for any $\sigma \in \Gamma$. Then, there holds

$$\inf_{h \in \mathcal{H}_\sigma} \{E^{\rho_i}(h) + \lambda \|h\|_\sigma^2\} = \inf_{h \in \mathcal{H}_\sigma} \{E^{\rho_j}(h) + \lambda \|h\|_\sigma^2\}, \quad (14)$$

for any $\sigma \in \Gamma$. Then the lemma follows. \square

To decompose the excess error, we give the following notations. Denote

$$(\bar{h}_1, \dots, \bar{h}_n) = \arg \min_{h_1, \dots, h_n \in \mathcal{H}_{\hat{\sigma}}} \sum_{i=1}^n \{E^{\rho_i}(h_i) + \lambda \|h_i\|_{\hat{\sigma}}^2\}, \quad (15)$$

$$(\hat{h}_1^*, \dots, \hat{h}_n^*) = \arg \min_{h_1, \dots, h_n \in \mathcal{H}_{\sigma^*}} \sum_{i=1}^n \{\hat{E}^{S_i}(h_i) + \lambda \|h_i\|_{\sigma^*}^2\}.$$

Proposition 6. Let \hat{h}_i be defined in (10). Then under the assumption of Lemma 5, $E^{\rho_1}(\bar{h}_1) - E^{\rho_1}(f_{\rho_1}) \leq E^{\rho_1}(\hat{h}_1) - E^{\rho_1}(f_{\rho_1}) + \lambda \|\hat{h}_1\|_{\hat{\sigma}}^2$ can be bounded by

$$\begin{aligned} & \{E^{\rho_1}(\bar{h}_1^*) - E^{\rho_1}(f_{\rho_1}) + \lambda \|\bar{h}_1^*\|_{\sigma^*}^2\} + \{E^{\rho_1}(\hat{h}_1) - \hat{E}^{S_1}(\hat{h}_1)\} \\ & + \{\hat{E}^{S_1}(\bar{h}_1) - E^{\rho_1}(\bar{h}_1)\} + \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(\bar{h}_i^*) - \hat{E}^{S_i}(\bar{h}_i^*)\} \\ & + \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(\hat{h}_i) - \hat{E}^{S_i}(\hat{h}_i)\}. \end{aligned} \quad (16)$$

Proof. Write the regularization error as

$$\begin{aligned} & E^{\rho_1}(\hat{h}_1) - E^{\rho_1}(f_{\rho_1}) + \lambda \|\hat{h}_1\|_{\hat{\sigma}}^2 \\ & = \{E^{\rho_1}(\hat{h}_1) - \hat{E}^{S_1}(\hat{h}_1)\} \\ & + \left\{ \left(\hat{E}^{S_1}(\hat{h}_1) + \lambda \|\hat{h}_1\|_{\hat{\sigma}}^2 \right) - \left(\hat{E}^{S_1}(\bar{h}_1) + \lambda \|\bar{h}_1\|_{\hat{\sigma}}^2 \right) \right\} \\ & + \{\hat{E}^{S_1}(\bar{h}_1) - E^{\rho_1}(\bar{h}_1)\} \\ & + \left\{ \left(E^{\rho_1}(\bar{h}_1) + \lambda \|\bar{h}_1\|_{\hat{\sigma}}^2 \right) - \left(E^{\rho_1}(\bar{h}_1^*) + \lambda \|\bar{h}_1^*\|_{\sigma^*}^2 \right) \right\} \\ & + \left\{ E^{\rho_1}(\bar{h}_1^*) - E^{\rho_1}(f_{\rho_1}) + \lambda \|\bar{h}_1^*\|_{\sigma^*}^2 \right\}. \end{aligned} \quad (17)$$

By Lemma 5, for any $i, j \in \{1, \dots, n\}$, there holds $E^{\rho_i}(\bar{h}_i) + \lambda \|\bar{h}_i\|_{\hat{\sigma}}^2 = E^{\rho_j}(\bar{h}_j) + \lambda \|\bar{h}_j\|_{\hat{\sigma}}^2$. Then we obtain

$$\begin{aligned}
& E^{\rho_1}(\bar{h}_1) + \lambda \|\bar{h}_1\|_{\hat{\sigma}}^2 \\
&= \frac{1}{n-1} \sum_{i=2}^n \left\{ E^{\rho_i}(\bar{h}_i) + \lambda \|\bar{h}_i\|_{\hat{\sigma}}^2 \right\} \\
&= \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(E^{\rho_i}(\bar{h}_i) + \lambda \|\bar{h}_i\|_{\hat{\sigma}}^2 \right) - \left(E^{\rho_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right) \right\} \\
&\quad + \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(E^{\rho_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right) - \left(\hat{E}^{S_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right) \right\} \\
&\quad + \frac{1}{n-1} \sum_{i=2}^n \left\{ \hat{E}^{S_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right\}, \\
& E^{\rho_1}(\bar{h}_1^*) + \lambda \|\bar{h}_1^*\|_{\sigma^*}^2 \\
&= \frac{1}{n-1} \sum_{i=2}^n \left\{ E^{\rho_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right\} \\
&= \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(E^{\rho_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right) \right. \\
&\quad \left. - \left(\hat{E}^{S_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right) \right\} \\
&\quad + \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(\hat{E}^{S_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right) \right. \\
&\quad \left. - \left(\hat{E}^{S_i}(\hat{h}_i^*) + \lambda \|\hat{h}_i^*\|_{\sigma^*}^2 \right) \right\} \\
&\quad + \frac{1}{n-1} \sum_{i=2}^n \left\{ \hat{E}^{S_i}(\hat{h}_i^*) + \lambda \|\hat{h}_i^*\|_{\sigma^*}^2 \right\}.
\end{aligned} \tag{18}$$

By the definition, we have that

$$\begin{aligned}
& \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(E^{\rho_i}(\bar{h}_i) + \lambda \|\bar{h}_i\|_{\hat{\sigma}}^2 \right) - \left(E^{\rho_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right) \right\} \leq 0, \\
& \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(\hat{E}^{S_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right) - \left(\hat{E}^{S_i}(\hat{h}_i^*) + \lambda \|\hat{h}_i^*\|_{\sigma^*}^2 \right) \right\} \\
& \geq 0, \\
& \frac{1}{n-1} \sum_{i=2}^n \left\{ \hat{E}^{S_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right\} \\
& \leq \frac{1}{n-1} \sum_{i=2}^n \left\{ \hat{E}^{S_i}(\hat{h}_i^*) + \lambda \|\hat{h}_i^*\|_{\sigma^*}^2 \right\}.
\end{aligned} \tag{19}$$

Therefore, we have

$$\begin{aligned}
& \left(E^{\rho_1}(\bar{h}_1) + \lambda \|\bar{h}_1\|_{\hat{\sigma}}^2 \right) - \left(E^{\rho_1}(\bar{h}_1^*) + \lambda \|\bar{h}_1^*\|_{\sigma^*}^2 \right) \\
& \leq \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(E^{\rho_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right) - \left(\hat{E}^{S_i}(\hat{h}_i) + \lambda \|\hat{h}_i\|_{\hat{\sigma}}^2 \right) \right\} \\
& \quad + \frac{1}{n-1} \sum_{i=2}^n \left\{ \left(E^{\rho_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right) \right. \\
& \quad \left. - \left(\hat{E}^{S_i}(\bar{h}_i^*) + \lambda \|\bar{h}_i^*\|_{\sigma^*}^2 \right) \right\}.
\end{aligned} \tag{20}$$

Then, the proposition follows. \square

In (16), there are 5 terms. The first term is called regularization error which depends on the approximation ability of hypothesis space to f_{ρ_1} . The estimation of this term has been discussed in [20] for reproducing kernel Hilbert space with Gaussian kernel with flexible variances.

The other four terms are called sample error. In the second and third terms, \hat{h}_1 and \bar{h}_1 are selected from $\mathcal{H}_{\hat{\sigma}}$ which is dependent on S_i , for $i = 2, \dots, n$, and is independent of S_1 . Therefore, when we take expectation of S_1 , $\mathcal{H}_{\hat{\sigma}}$ can be seen as a fixed function space. Consequently, these two terms can be estimated with the same method in the proof of Propositions 2.1 and 3.1 in [21]. In the fourth term, \bar{h}_i^* is a fixed function, for $i = 1, \dots, n$. Therefore, this term can also be estimated as in the proof of Proposition 2.1 in [21]. The last term is more difficult to deal with because \hat{h}_i , for $i = 2, \dots, n$, can not be considered in \mathcal{H}_{σ} , for any fixed $\sigma \in \Gamma$. Consequently, when sample number $m \rightarrow \infty$, the convergence rate of the sample error depends on that of the last term. Therefore, in the following section, we focus on the estimation for the bound of the last term in (16).

4. Error Analysis

In this section, we estimate the bound of the last term in (16). To bound this term, we have to estimate capacity of $\{\mathcal{H}_{\sigma} : \sigma \in \Gamma\}$. Here, the capacity is measured by the covering number.

Definition 7. For a subset \mathcal{F} of a metric space (X, d) and $\eta > 0$, the covering number $\mathcal{N}(\mathcal{F}, \eta, d)$ is defined to be the minimal integer $l \in \mathbb{N}$ such that there exist l disks with radius η covering \mathcal{F} .

For $n \in \mathbb{N}$ and $\mathcal{H}_{\sigma, R} = \{f \in \mathcal{H}_{\sigma} : \|f\|_{\sigma} \leq R\}$, define

$$\mathcal{H}_{\sigma, R}^n = \{f = (f_1, \dots, f_n) : f_i \in \mathcal{H}_{\sigma, R}, i = 1, \dots, n\}. \tag{21}$$

For $f \in \mathcal{H}_{\sigma_1, R}^n$ and $g \in \mathcal{H}_{\sigma_2, R}^n$, define

$$\begin{aligned}
l_{\infty, n}(f, g) &= \frac{1}{n} \sum_{i=1}^n \|f_i - g_i\|_{\infty}, \\
d_{\infty}(f, \mathcal{H}_{\sigma_2, R}^n) &= \inf_{g \in \mathcal{H}_{\sigma_2, R}^n} l_{\infty, n}(f, g).
\end{aligned} \tag{22}$$

Then, define the distance from $\mathcal{H}_{\sigma_1, R}^n$ to $\mathcal{H}_{\sigma_2, R}^n$ as

$$\tilde{d}_{\infty}(\mathcal{H}_{\sigma_1, R}^n, \mathcal{H}_{\sigma_2, R}^n) = \sup_{f \in \mathcal{H}_{\sigma_1, R}^n} d_{\infty}(f, \mathcal{H}_{\sigma_2, R}^n). \quad (23)$$

Let $\mathcal{N}_R(\varepsilon, \Gamma, \tilde{d}_{\infty})$ denote the minimal integer $l \in \mathbb{N}$ such that there exist l parameters $\sigma_1, \dots, \sigma_l \in \Gamma$, such that

$$\sup_{\sigma \in \Gamma} \min_{i \in \{1, \dots, l\}} \tilde{d}_{\infty}(\mathcal{H}_{\sigma, R}^n, \mathcal{H}_{\sigma_i, R}^n) \leq \varepsilon. \quad (24)$$

Then, for $\mathcal{H}_{\Gamma, R}^n = \{f = (f_1, \dots, f_n) : f_i \in \mathcal{H}_{\sigma, R}, i = 1, \dots, n, \sigma \in \Gamma\}$, we have the following lemma.

Lemma 8. Consider the following:

$$\begin{aligned} & \mathcal{N}(2\varepsilon, \mathcal{H}_{\Gamma, R}^n, l_{\infty, n}) \\ & \leq \mathcal{N}_R(\varepsilon, \Gamma, \tilde{d}_{\infty}) \cdot \left(\sup_{\sigma \in \Gamma} \mathcal{N}(\varepsilon, \mathcal{H}_{\sigma, R}, l_{\infty, 1}) \right)^n. \end{aligned} \quad (25)$$

Proof. For any $f \in \mathcal{H}_{\Gamma, R}^n$, there is $\sigma \in \Gamma$ such that $f \in \mathcal{H}_{\sigma, R}^n$. By the definition of $\mathcal{N}_R(\varepsilon, \Gamma, \tilde{d}_{\infty})$, there is $\bar{\sigma} \in \{\sigma_1, \sigma_2, \dots, \sigma_{\mathcal{N}_R(\varepsilon, \Gamma, \tilde{d}_{\infty})}\}$ such that $\tilde{d}_{\infty}(\mathcal{H}_{\sigma, R}^n, \mathcal{H}_{\bar{\sigma}, R}^n) \leq \varepsilon$. Then, by the definition of \tilde{d}_{∞} , we have $d_{\infty}(f, \mathcal{H}_{\bar{\sigma}, R}^n) \leq \varepsilon$. And by the definition of d_{∞} , there is $\bar{f} \in \mathcal{H}_{\bar{\sigma}, R}^n$ satisfying $l_{\infty, n}(f, \bar{f}) \leq \varepsilon$.

By the definition of $\mathcal{N}(\varepsilon, \mathcal{H}_{\bar{\sigma}, R}^n, l_{\infty, n})$, there is $\tilde{f} \in \{f_{\bar{\sigma}}^1, \dots, f_{\bar{\sigma}}^{\mathcal{N}(\varepsilon, \mathcal{H}_{\bar{\sigma}, R}^n, l_{\infty, n})}\} \subseteq \mathcal{H}_{\bar{\sigma}, R}^n$ such that $l_{\infty, n}(\bar{f}, \tilde{f}) \leq \varepsilon$. Therefore, we can obtain

$$\mathcal{N}(2\varepsilon, \mathcal{H}_{\Gamma, R}^n, l_{\infty, n}) \leq \mathcal{N}_R(\varepsilon, \Gamma, \tilde{d}_{\infty}) \cdot \sup_{\sigma \in \Gamma} \mathcal{N}(\varepsilon, \mathcal{H}_{\sigma, R}^n, l_{\infty, n}). \quad (26)$$

Note that, for any $\sigma \in \Gamma$, there holds $\mathcal{N}(\varepsilon, \mathcal{H}_{\sigma, R}^n, l_{\infty, n}) \leq (\mathcal{N}(\varepsilon, \mathcal{H}_{\sigma, R}, l_{\infty, 1}))^n$. Then the lemma follows. \square

Proposition 9. For $R = M/\sqrt{\lambda}$ and \hat{h}_i , for $i = 2, \dots, n$, defined in (10), there holds

$$\begin{aligned} & P \left\{ \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(\hat{h}_i) - \hat{E}^{S_i}(\hat{h}_i)\} \geq \varepsilon_1 \right\} \\ & \leq \mathcal{N} \left(\frac{\varepsilon_1}{16(\kappa R + M)}, \Gamma, d_{\infty} \right) \\ & \quad \cdot \left(\sup_{\sigma \in \Gamma} \mathcal{N} \left(\frac{\varepsilon_1}{16(\kappa R + M)}, \mathcal{H}_{\sigma, R}, l_{\infty, 1} \right) \right)^{n-1} \\ & \quad \cdot \exp \left\{ -\frac{(n-1)m\varepsilon_1^2}{8(\kappa R + M)^2} \right\}. \end{aligned} \quad (27)$$

Proof. By the definition of \hat{h}_i , we have $\|\hat{h}_i\|_{\hat{\sigma}} \leq R$, for $i = 1, 2, \dots, n$. Then, there holds

$$\begin{aligned} & P \left\{ \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(\hat{h}_i) - \hat{E}^{S_i}(\hat{h}_i)\} \geq \varepsilon_1 \right\} \\ & \leq P \left\{ \sup_{[f_2, \dots, f_n] \in \mathcal{H}_{\Gamma, R}^{n-1}} \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(f_i) - \hat{E}^{S_i}(f_i)\} \geq \varepsilon_1 \right\}. \end{aligned} \quad (28)$$

For $f = [f_2, \dots, f_n] \in \mathcal{H}_{\Gamma, R}^{n-1}$, denote $L_z(f) = (1/(n-1)) \sum_{i=2}^n \{E^{\rho_i}(f_i) - \hat{E}^{S_i}(f_i)\}$. Note that, for $f, g \in \mathcal{H}_{\Gamma, R}^{n-1}$, there holds

$$\begin{aligned} & |L_z(f) - L_z(g)| \\ & = \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(f_i) - \hat{E}^{S_i}(f_i)\} \\ & \quad - \frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(g_i) - \hat{E}^{S_i}(g_i)\} \\ & = \frac{1}{n-1} \sum_{i=2}^n \left\{ \int_Z (f_i(x) - y)^2 d\rho_i - \int_Z (g_i(x) - y)^2 d\rho_i \right\} \\ & \quad - \frac{1}{n-1} \sum_{i=2}^n \left\{ \frac{1}{m} \sum_{j=1}^m (f_i(x_j^i) - y_j^i)^2 \right. \\ & \quad \left. - \frac{1}{m} \sum_{j=1}^m (g_i(x_j^i) - y_j^i)^2 \right\} \\ & \leq 4(\kappa R + M) \cdot l_{\infty, n-1}(f, g). \end{aligned} \quad (29)$$

Therefore, we can find that $\mathcal{N} = \mathcal{N}(\varepsilon_1/8(\kappa R + M), \mathcal{H}_{\Gamma, R}^{n-1}, l_{\infty, n-1})$ balls $\{\mathcal{H}_{k, R}^{n-1}\}_{k=1}^{\mathcal{N}}$ such that the center \tilde{f}^k of each ball and point f^k in this ball satisfies $|L_z(\tilde{f}^k) - L_z(f^k)| \leq \varepsilon/2$. Therefore, the probability in (28) can be bounded by following expression:

$$\begin{aligned} & P \left\{ \sup_{[f_2, \dots, f_n] \in \bigcup_{k=1}^{\mathcal{N}(\varepsilon_1/8(\kappa R + M), \mathcal{H}_{\Gamma, R}^{n-1}, l_{\infty, n-1})} \mathcal{H}_{k, R}^{n-1}} L_z(f) \geq \varepsilon_1 \right\} \\ & \leq \mathcal{N} \left(\frac{\varepsilon_1}{8(\kappa R + M)}, \mathcal{H}_{\Gamma, R}^{n-1}, l_{\infty, n} \right) \\ & \quad \cdot \max_{k \in 1, \dots, \mathcal{N}} P \left\{ \sup_{[f_2, \dots, f_n] \in \mathcal{H}_{k, R}^{n-1}} L_z(f) \geq \varepsilon_1 \right\} \\ & \leq \mathcal{N} \left(\frac{\varepsilon_1}{8(\kappa R + M)}, \mathcal{H}_{\Gamma, R}^{n-1}, l_{\infty, n} \right) \\ & \quad \cdot \max_{k \in 1, \dots, \mathcal{N}} P \left\{ L_z(\tilde{f}^k) \geq \frac{\varepsilon_1}{2} \right\}. \end{aligned} \quad (30)$$

For the covering number \mathcal{N} , by Lemma 8, we have the estimate

$$\begin{aligned} & \mathcal{N}\left(\frac{\varepsilon_1}{8(\kappa R + M)}, \mathcal{H}_{\Gamma, R}^{n-1}, l_{\infty, n}\right) \\ & \leq \mathcal{N}\left(\frac{\varepsilon_1}{16(\kappa R + M)}, \Gamma, \tilde{d}_{\infty}\right) \\ & \cdot \left(\sup_{\sigma \in \Gamma} \mathcal{N}\left(\frac{\varepsilon_1}{16(\kappa R + M)}, \mathcal{H}_{\sigma, R}, l_{\infty, 1}\right)\right)^{n-1}. \end{aligned} \quad (31)$$

Using Hoeffding inequality, for any fixed \tilde{f}^k , we have

$$\begin{aligned} & P\left\{L_z(\tilde{f}^k) \geq \frac{\varepsilon_1}{2}\right\} \\ & = P\left\{\frac{1}{(n-1)m} \times \sum_{i=2}^n \sum_{j=1}^m \left\{ \int_Z (\tilde{f}_i^k(x) - y)^2 d\rho_i \right. \right. \\ & \quad \left. \left. - (\tilde{f}_i^k(x_j^i) - y_j^i)^2 \right\} \geq \frac{\varepsilon_1}{2}\right\} \\ & \leq \exp\left\{-\frac{(n-1)m\varepsilon_1^2}{8(\kappa R + M)^2}\right\}. \end{aligned} \quad (32)$$

Then, the proposition follows. \square

Finally, we can obtain a bound for the last term of (16) by Proposition 9.

Proposition 10. Let \hat{h}_i , for $i = 2, \dots, n$, be defined in (10). Assume, for all $\sigma \in \Gamma$ and for all $s > 0$, there holds

$$\ln \mathcal{N}(\eta, \mathcal{H}_{\sigma, 1}, l_{\infty, 1}) \leq C_0 \left(\frac{1}{\eta}\right)^s. \quad (33)$$

Then with confidence at least $1 - \delta$, there holds

$$\frac{1}{n-1} \sum_{i=2}^n \{E^{\rho_i}(\hat{h}_i) - \hat{E}^{\hat{S}_i}(\hat{h}_i)\} \leq \varepsilon_0, \quad (34)$$

where ε_0 is the solution of the following equation:

$$\begin{aligned} & \varepsilon_0^{2+s} - \frac{8(\ln(1/\delta))(\kappa R + M)^2}{(n-1)m} \varepsilon_0^s - \frac{8(\kappa R + M)^2}{m} \\ & \times C_0(16(\kappa R^2 + MR))^s \\ & - \frac{8(\kappa R + M)^2}{(n-1)m} \ln \mathcal{N}_R\left(\frac{\varepsilon_0}{16(\kappa R + M)}, \Gamma, d_{\infty}\right) = 0. \end{aligned} \quad (35)$$

Proof. Let

$$\begin{aligned} \delta & = \mathcal{N}_R\left(\frac{\varepsilon_1}{16(\kappa R + M)}, \Gamma, d_{\infty}\right) \\ & \cdot \left(\sup_{\sigma \in \Gamma} \mathcal{N}\left(\frac{\varepsilon_1}{16(\kappa R + M)}, \mathcal{H}_{\sigma, R}, l_{\infty}^1\right)\right)^{n-1} \\ & \cdot \exp\left\{-\frac{(n-1)m\varepsilon_1^2}{8(\kappa R + M)^2}\right\}. \end{aligned} \quad (36)$$

By condition on $\ln \mathcal{N}(\eta, \mathcal{H}_{\sigma, 1}, l_{\infty, 1}) \leq C_0(1/\eta)^s$, we have

$$\ln \mathcal{N}\left(\frac{\varepsilon_1}{16(\kappa R + M)}, \mathcal{H}_{\sigma, R}, l_{\infty, 1}\right) \leq C_0 \left(\frac{16(\kappa R^2 + MR)}{\varepsilon_1}\right)^s. \quad (37)$$

Then, ε_1 is not larger than ε_0 in the following equation:

$$\begin{aligned} \ln \delta & = \ln \mathcal{N}_R\left(\frac{\varepsilon_0}{16(\kappa R + M)}, \Gamma, d_{\infty}\right) \\ & + (n-1)C_0 \left(\frac{16(\kappa R^2 + MR)}{\varepsilon_0}\right)^s - \frac{(n-1)m}{8(\kappa R + M)^2} \varepsilon_0^2. \end{aligned} \quad (38)$$

Then by Proposition 9, the proposition follows. \square

Remark 11. Compare multitask learning with multiple Hilbert spaces to single task learning with multiple Hilbert space.

Recall that the least square regularized regression in $\{\mathcal{H}_{\sigma}, \sigma \in \Gamma\}$ for single task is defined as

$$(\hat{\sigma}_1, \hat{h}_1) = \arg \min_{h_1 \in \mathcal{H}_{\sigma}, \sigma \in \Gamma} \{\hat{E}^{\hat{S}_1}(h_1) + \lambda \|h_1\|_{\sigma}^2\}. \quad (39)$$

$E^{\rho_1}(\hat{h}_1) - E^{\rho_1}(f_{\rho_1})$ can be bounded by the sum of regularization error and sample error with similar method in Proposition 6. In the sample error, the term most difficult to estimated will be $E^{\hat{S}_1}(\hat{h}_1) - E^{\rho_1}(\hat{h}_1)$, because function \hat{h}_1 changed with the S_1 runs over function set $\{\mathcal{H}_{\sigma}, \sigma \in \Gamma\}$. By the same method in Proposition 10, with confidence $1 - \delta$, this term can be bounded by the solution $\tilde{\varepsilon}_0$ of the following equation:

$$\begin{aligned} & \tilde{\varepsilon}_0^{2+s} - \frac{8(\ln(1/\delta))(\kappa R + M)^2}{m} \tilde{\varepsilon}_0^s - \frac{8(\kappa R + M)^2}{m} \\ & \times C_0(16(\kappa R^2 + MR))^s \\ & - \frac{8(\kappa R + M)^2}{m} \ln \mathcal{N}_R\left(\frac{\tilde{\varepsilon}_0}{16(\kappa R + M)}, \Gamma, d_{\infty}\right) = 0. \end{aligned} \quad (40)$$

Obviously, we have $\varepsilon_0 < \tilde{\varepsilon}_0$. Therefore, multitask learning algorithm has potential faster learning rate.

Remark 12. Comparing multitask learning with multiple Hilbert spaces to single task learning with single Hilbert space.

In this paper, we set the hypothesis space as a set of Hilbert spaces. It is well known that hypothesis space with more

functions has stronger approximation ability and bigger complexity. Therefore, the regularization error may be smaller and sample error may be larger than that of algorithms with a single Hilbert space being the hypothesis space.

For least square regularized regression with a single Hilbert space \mathcal{H}_σ for single task, with confidence $1 - \delta$, the largest term in sample error can be bounded by the solution $\bar{\epsilon}_0$ of the following equation:

$$\begin{aligned} \bar{\epsilon}_0^{2+s} - \frac{8(\ln(1/\delta))(\kappa R + M)^2}{m} \bar{\epsilon}_0^s \\ - \frac{8(\kappa R + M)^2}{m} C_0 (16(\kappa R^2 + MR))^s = 0. \end{aligned} \quad (41)$$

If we assume $n = \mathcal{O}(m^\zeta)$ with some $\zeta > 0$ large enough, $(8(\kappa R + M)^2/(n-1)m) \ln \mathcal{N}_R(\epsilon_0/16(\kappa R + M), \Gamma, d_\infty)$ in Proposition 10 can converge to 0 fast. Then, we can obtain $\epsilon_0 \approx \bar{\epsilon}_0$, while the regularization error of multitask learning with multiple Hilbert spaces is smaller. Trading off the regularization error and sample error, we can obtain potential faster learning rate than that of single task learning.

Acknowledgments

This work was supported in part by the Fundamental Research Funds for the Central Universities under Grant ZY1118, National Technology Support Program under Grant 2012BAH05B01, and the National Science Foundation of China under Grants 11101024 and 11171014.

References

- [1] R. Caruana, "Multitask learning," *Machine Learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [2] Y. Zhang and D.-Y. Yeung, "A convex formulation for learning task relationships in multi-task learning," in *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI '10)*, pp. 733–742, Catalina Island, Calif, USA, July 2010.
- [3] T. Evgeniou and M. Pontil, "Regularized multi-task learning," in *Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 109–117, New York, NY, USA, August 2004.
- [4] A. Argyriou, T. Evgeniou, and M. Pontil, "Convex multi-task feature learning," *Machine Learning*, vol. 73, no. 3, pp. 243–272, 2008.
- [5] J. Liu, S. Ji, and J. Ye, "Multi-task feature learning via efficient L2,1-norm minimization," in *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence (UAI '09)*, pp. 339–348, Montreal, Canada, June 2009.
- [6] B. Bakker and T. Heskes, "Task clustering and gating for bayesian multitask learning," *Journal of Machine Learning Research*, vol. 4, no. 1, pp. 83–99, 2004.
- [7] A. Argyriou, C. A. Micchelli, M. Pontil, and Y. M. Ying, "A spectral regularization framework for multi-task structure learning," in *Proceedings of the 21st Annual Conference on Advances in Neural Information Processing Systems (NIPS '07)*, December 2007.
- [8] J. Guinney, Q. Wu, and S. Mukherjee, "Estimating variable structure and dependence in multitask learning via gradients," *Machine Learning*, vol. 83, no. 3, pp. 265–287, 2011.
- [9] Y. Xue, X. Liao, L. Carin, and B. Krishnapuram, "Multi-task learning for classification with Dirichlet process priors," *Journal of Machine Learning Research*, vol. 8, pp. 35–63, 2007.
- [10] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2010.
- [11] J. Baxter, "Learning internal representations," in *Proceedings of the Workshop on Computational Learning Theory (COLT '95)*, Morgan Kaufmann, San Mateo, Calif, USA, 1995.
- [12] N. Intrator and S. Edelman, "Making a low-dimensional representation suitable for diverse tasks," *Connection Science*, vol. 8, no. 2, pp. 205–224, 1996.
- [13] S. Thrun, "Is learning the n-th thing any easier than learning the first?" in *Proceedings of the Advances in Neural Information Processing Systems (NIPS '96)*, D. Touretzky and M. Mozer, Eds., 1996.
- [14] T. Heskes, "Solving a huge number of similar tasks: a combination of multi-task learning and a hierarchical Bayesian approach," in *Proceedings of the International Conference on Machine Learning (ICML '98)*, 1998.
- [15] B. Romera-Paredes, A. Argyriou, N. Berthouze, and M. Pontil, "Exploiting unrelated tasks in multi-task learning," in *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS '12)*, pp. 951–959, La Palma, Spain, 2012.
- [16] J. Baxter, "A model of inductive bias learning," *Journal of Artificial Intelligence Research*, vol. 12, pp. 149–198, 2000.
- [17] R. K. Ando and T. Zhang, "A framework for learning predictive structures from multiple tasks and unlabeled data," *Journal of Machine Learning Research*, vol. 6, pp. 1817–1853, 2005.
- [18] S. Ben-David and R. S. Borbely, "A notion of task relatedness yielding provable multiple-task learning guarantees," *Machine Learning*, vol. 73, no. 3, pp. 273–287, 2008.
- [19] M. Solnon, S. Arlot, and F. Bach, "Multi-task regression using minimal penalties," *Journal of Machine Learning Research*, vol. 13, pp. 2773–2812, 2012.
- [20] Y. M. Ying and D.-X. Zhou, "Learnability of Gaussians with flexible variances," *Journal of Machine Learning Research*, vol. 8, pp. 249–276, 2007.
- [21] Q. Wu, Y. Ying, and D.-X. Zhou, "Learning rates of least-square regularized regression," *Foundations of Computational Mathematics*, vol. 6, no. 2, pp. 171–192, 2006.