

# Bounds for Linear Multi-Task Learning

Andreas Maurer

Adalbertstr. 55  
D-80799 München

andreasmaurer@compuserve.com

**Abstract.** We give dimension-free and data-dependent bounds for linear multi-task learning where a common linear operator is chosen to preprocess data for a vector of task specific linear-thresholding classifiers. The complexity penalty of multi-task learning is bounded by a simple expression involving the margins of the task-specific classifiers, the Hilbert-Schmidt norm of the selected preprocessor and the Hilbert-Schmidt norm of the covariance operator for the total mixture of all task distributions, or, alternatively, the Frobenius norm of the total Gramian matrix for the data-dependent version. The results can be compared to state-of-the-art results on linear single-task learning.

## 1 Introduction

Simultaneous learning of different tasks under some common constraint, often called *multi-task learning*, has been tested in practice with good results under a variety of different circumstances (see [5], [8], [17], [18]). The technique has been analyzed theoretically and in some generality (see Baxter [6] and Zhang[18]). The purpose of this paper is to improve some of these theoretical results in a special case of practical importance, when input data is represented in a linear, potentially infinite dimensional space, and the common constraint is a linear preprocessor.

A simple way to understand multi-task learning and its potential advantages is perhaps agnostic learning with an input space  $\mathcal{X}$  and a finite set  $\mathcal{F}$  of hypotheses  $f : \mathcal{X} \rightarrow \{0, 1\}$ . For a hypothesis  $f \in \mathcal{F}$  let  $\text{er}(f)$  be the expected error and  $\hat{\text{er}}(f)$  the empirical error on a training sample  $S$  of size  $n$  (drawn iid from the underlying task distribution) respectively. Combining Hoeffding's inequality with a union bound one shows (see e.g. [1]), that with probability greater than  $1 - \delta$  we have for every  $f \in \mathcal{F}$  the error bound

$$\text{er}(f) \leq \hat{\text{er}}(f) + \frac{1}{\sqrt{2n}} \sqrt{\ln |\mathcal{F}| + \ln(1/\delta)}. \quad (1)$$

Suppose now that there are a set  $\mathcal{Y}$ , a rather large set  $\mathcal{G}$  of preprocessors  $g : \mathcal{X} \rightarrow \mathcal{Y}$ , and another set  $\mathcal{H}$  of classifiers  $h : \mathcal{Y} \rightarrow \{0, 1\}$  with  $|\mathcal{H}| \ll |\mathcal{F}|$ . For a cleverly chosen preprocessor  $g \in \mathcal{G}$  it will likely be the case that we find some  $h \in \mathcal{H}$  such that  $h \circ g$  has the same or even a smaller empirical error than the

best  $f \in \mathcal{F}$ . But this will lead to an improvement of the bound above (replacing  $|\mathcal{F}|$  by  $|\mathcal{H}|$ ) only if we choose  $g$  before seeing the data, otherwise we incur a large estimation penalty for the selection of  $g$  (replacing  $|\mathcal{F}|$  by  $|\mathcal{H} \circ \mathcal{G}|$ ).

The situation is improved if we have a set of  $m$  different learning tasks with corresponding task distributions and samples  $S_1, \dots, S_m$ , each of size  $n$  and drawn iid from the corresponding distributions. We now consider solutions  $h_1 \circ g, \dots, h_m \circ g$  for each of the  $m$  tasks where the preprocessing map  $g \in \mathcal{G}$  is *constrained to be the same for all tasks* and only the  $h_l \in \mathcal{H}$  specialize to each task  $l$  at hand. Again Hoeffding's inequality and a union bound imply that with probability greater  $1 - \delta$  we have for all  $(h_1, \dots, h_m) \in \mathcal{H}^m$  and every  $g \in \mathcal{G}$

$$\frac{1}{m} \sum_{l=1}^m \text{er}^l(h_l \circ g) \leq \frac{1}{m} \sum_{l=1}^m \text{er}^l(h_l \circ g) + \frac{1}{\sqrt{2n}} \sqrt{\ln |\mathcal{H}| + \frac{\ln |\mathcal{G}| + \ln(1/\delta)}{m}}. \quad (2)$$

Here  $\text{er}^l(f)$  and  $\text{er}^l(f)$  denote the expected error in task  $l$  and the empirical error on training sample  $S_l$  respectively. The left hand side above is an average of the expected errors, so that the guarantee implied by the bound is a little weaker than the usual PAC guarantees (but see Ben-David [7] for bounds on the individual errors). The first term on the right is the average empirical error, which a multi-task learning algorithm seeks to minimize. We can take it as an operational definition of task-relatedness relative to  $(\mathcal{H}, \mathcal{G})$  that we are able to obtain a very small value for this term. The remaining term, which bounds the estimation error, now exhibits the advantage of multi-task learning: Sharing the preprocessor implies sharing its cost of estimation, and the entropy contribution arising from the selection of  $g \in \mathcal{G}$  decreases with the number of learning tasks. Since by assumption  $|\mathcal{H}| \ll |\mathcal{F}|$ , the estimation error in the multi-task bound (2) can become much smaller than in the single task case (1) if the number  $m$  of tasks becomes large.

The choice of the preprocessor  $g \in \mathcal{G}$  can also be viewed as the selection of the hypothesis space  $\mathcal{H} \circ g$ . This leads to an alternative formulation of multi-task learning, where the common object is a hypothesis space chosen from a class of hypothesis spaces (in this case  $\{\mathcal{H} \circ g : g \in \mathcal{G}\}$ ), and the classifiers for the individual tasks are all chosen from the selected hypothesis space. Here we prefer the functional formulation of selecting a preprocessor instead of a hypothesis space, because it is more intuitive and sufficient in the situations which we consider.

The arguments leading to (2) can be refined and extended to certain infinite classes to give general bounds for multi-task learning ([6] and [18]). In this paper we concentrate on the case where the input space  $\mathcal{X}$  is a subset of the unit ball in a Hilbert space  $H$ , the class  $\mathcal{G}$  of preprocessors is a set  $\mathcal{A}$  of bounded symmetric linear operators on  $H$ , and the class  $\mathcal{H}$  is the set of classifiers  $h_v$  obtained by 0-thresholding linear functionals  $v$  in  $H$  with  $\|v\| \leq B$ , that is

$$h_v(x) = \text{sign}(\langle x, v \rangle) \text{ and } h_v \circ T(x) = \text{sign}(\langle Tx, v \rangle), x \in H, T \in \mathcal{G}, \|v\| \leq B.$$

The learner now searches for a multi-classifier  $\mathbf{h}_{\mathbf{v}} \circ T = (h_{v^1} \circ T, \dots, h_{v^m} \circ T)$  where the preprocessing operator  $T \in \mathcal{A}$  is the same for all tasks and only the vectors  $v^l$  specialize to each task  $l$  at hand. We desired multi-classifier  $\mathbf{h}_{\mathbf{v}} \circ T$  should have a small value of the average error

$$\text{er}(\mathbf{h}_{\mathbf{v}} \circ T) = \frac{1}{m} \sum_{l=1}^m \text{er}^l(h_{v^l} \circ T) = \frac{1}{m} \sum_{l=1}^m \Pr \{ \text{sign}(\langle TX^l, v^l \rangle) \neq Y^l \},$$

where  $X^l$  and  $Y^l$  are the random variables modelling input-values and labels for the  $l$ -th task. To guide this search we look for bounds on  $\text{er}(\mathbf{h}_{\mathbf{v}} \circ T)$  in terms of the total observed data for all tasks, valid uniformly for all  $\mathbf{v} = (v^1, \dots, v^m)$  with  $\|v^l\| \leq B$  and all  $T \in \mathcal{A}$ . We will prove the following :

**Theorem 1.** *Let  $\delta \in (0, 1)$ . With probability greater than  $1 - \delta$  it holds for all  $\mathbf{v} = (v^1, \dots, v^m) \in H$  with  $\|v^l\| \leq 1$  and all bounded symmetric operators  $T$  on  $H$  with  $\|T\|_{HS} \geq 1$ , and for all  $\gamma \in (0, 1)$  that*

$$\text{er}(\mathbf{h}_{\mathbf{v}} \circ T) \leq \text{er}_{\gamma}(\mathbf{v} \circ T) + \frac{8 \|T\|_{HS}}{\gamma \sqrt{n}} \sqrt{\|C\|_{HS} + \frac{1}{m}} + \sqrt{\frac{\ln \frac{4}{\delta \gamma}}{2nm}}.$$

Here  $\text{er}_{\gamma}(\mathbf{v} \circ T)$  is a margin-based empirical error estimate, bounded by the relative number of examples  $(X_i^l, Y_i^l)$  in the total training sample for all tasks  $l$ , where  $Y_i^l \langle TX_i^l, v^l \rangle < \gamma$  (see section 4).

The quantity  $\|T\|_{HS}$  is the *Hilbert-Schmidt norm* of  $T$ , defined for symmetric  $T$  by

$$\|T\|_{HS} = \left( \sum \lambda_i^2 \right)^{1/2},$$

where  $\lambda_i$  is the sequence of eigenvalues of  $T$  (counting multiplicities, see section 2).

$C$  is the *total covariance operator* corresponding to the mixture of all the task-input-distributions in  $H$ . Since data is constrained to the unit ball in  $H$  we always have  $\|C\|_{HS} \leq 1$  (see section 3).

The above theorem is the simplest, but not the tightest or most general form of our results. For example the factor 8 on the right hand side can be decreased to be arbitrarily close to 2, thereby incurring only a logarithmic penalty in the last term.

A special case results from restricting the set of candidate preprocessors to  $\mathcal{P}_d$ , the set of orthogonal projections in  $H$  with  $d$ -dimensional range. In this case learning amounts to the selection of a  $d$ -dimensional subspace  $M$  of  $H$  and of an  $m$ -tuple of vectors  $v^l$  in  $M$  (components of  $v^l$  orthogonal to  $M$  are irrelevant to the projected data). All operators  $T \in \mathcal{P}_d$  satisfy  $\|T\|_{HS} = \sqrt{d}$ , which can then be substituted in the above bound. This covers the case considered by Ando and Zhang ([18]), where a practical algorithm for this type of multi-task learning is presented.

The bound in the above theorem is dimension free, it does not require the data distribution in  $H$  to be confined to a finite dimensional subspace. Almost to the contrary: Suppose that the input data is distributed uniformly on  $M \cap S_1$  where  $M$  is a  $k$ -dimensional subspace in  $H$  and  $S_1$  is the sphere consisting of vectors with unit norm in  $H$ . Then  $C$  has the  $k$ -fold eigenvalue  $1/k$ , the remaining eigenvalues being zero. Therefore  $\|C\|_{HS} = 1/\sqrt{k}$ , so part of the bound above decreases to zero as the dimensionality of the data-distribution increases. The fact that our bounds are dimension free (in contrast to those in [18], for example) allows their general use for multi-task learning in kernel-induced Hilbert spaces (see [9]).

If we compare the second term on the right hand side to the estimation error bound in (2), we can recognize a certain similarity: Loosely speaking we can identify  $\|T\|_{HS}^2/m$  with the cost of estimating the operator  $T$ , and  $\|T\|_{HS}^2 \|C\|_{HS}$  with the cost of finding the linear classifiers  $v_1, \dots, v_m$ . The order of dependence on the number of tasks  $m$  is the same in Theorem 1 as in (2).

In the limit  $m \rightarrow \infty$  it is preferable to use a different bound (see Theorems 6 and 7), at the expense of slower convergence in  $m$ . The main inequality of the theorem then becomes

$$\text{er}(\mathbf{h}_{\mathbf{v}} \circ T) \leq \text{er}_{\gamma}(\mathbf{v} \circ T) + \frac{2 \|T^2\|_{HS}^{1/2}}{(1-\epsilon)^2 \gamma \sqrt{n}} \left( \|C\|_{HS}^2 + \frac{3}{m} \right)^{1/4} + \sqrt{\frac{\ln \frac{1}{\delta \gamma \epsilon^2}}{2nm}}. \quad (3)$$

for some very small  $\epsilon > 0$  to be fixed in advance. If  $T$  is an orthogonal projection with  $d$ -dimensional range then  $\|T^2\|_{HS}^{1/2} = d^{1/4}$ , so for a large number of tasks  $m$  the bound on the estimation error becomes approximately

$$\frac{2d^{1/4} \|C\|_{HS}^{1/2}}{\gamma \sqrt{n}}.$$

One of the best dimension-free bounds for linear single-task learning (see e.g. Bartlett and Mendelson [2] or Lemma 4 below) would give  $2/(\gamma \sqrt{n})$  for this term, if all data is constrained to unit vectors. We therefore expect superior estimation for multi-task learning of  $d$ -dimensional projections with large  $m$ , whenever  $d^{1/4} \|C\|_{HS}^{1/2} \ll 1$ . If we again assume the data-distribution to be uniform on  $M \cap S_1$  with  $M$  a  $k$ -dimensional subspace, this is the case whenever  $d \ll k$ , that is, qualitatively speaking, whenever the dimension of the utilizable part of the data is considerably smaller than the dimension of the total data distribution.

The results stated above give some insights, but they have the practical disadvantage of being unobservable, because they depend on the properties of the covariance operator  $C$ , which in turn depends on an unknown data distribution. One way to solve this problem is using the fact that the finite-sample approximations to the covariance operator have good concentration properties (see Theorem 3 below). The corresponding result is:

**Theorem 2.** *With probability greater than  $1 - \delta$  in the sample  $\mathbf{X}$  it holds for all  $v_1, \dots, v_m \in H$  with  $\|v_l\| \leq 1$  and all bounded symmetric operators  $T$  on  $H$  with  $\|T\|_{HS} \geq 1$ , and for all  $\gamma \in (0, 1)$  that*

$$\frac{1}{m} \sum_{l=1}^m er(h_{v_l} \circ T) \leq \frac{1}{m} \sum_{l=1}^m er_{\hat{\gamma}}(v_l \circ T) + \frac{8 \|T\|_{HS}}{\gamma \sqrt{n}} \sqrt{\frac{1}{mn} \|\hat{C}(\mathbf{X})\|_{Fr} + \frac{1}{m}} + \sqrt{\frac{9 \ln \frac{8}{\delta \gamma}}{2nm}}.$$

where the  $\|\hat{C}(\mathbf{X})\|_{Fr}$  is the Frobenius norm of the gramian.

By definition

$$\|\hat{C}(\mathbf{X})\|_{Fr} = \left( \sum_{l,r,i,j} \langle X_i^l, X_j^r \rangle^2 \right)^{1/2}.$$

Here  $X_i^l$  is the random variable describing the  $i$ -th data point in the sample corresponding to the  $l$ -th task. The corresponding label  $Y_i^l$  enters only in the empirical margin error. The quantity  $(mn)^{-1} \|\hat{C}(\mathbf{X})\|_{Fr}$  can be regarded as an approximation to  $\|C\|_{HS}$ , valid with high probability, so that Theorem 2 is a sample based version of Theorem 1.

In section 2 we introduce the necessary terminology and results on Hilbert-Schmidt operators and in section 3 the covariance operator of random elements in a Hilbert space. Section 4 gives a formal definition of multi-task systems and a general PAC bound in terms of Rademacher complexities. For the readers benefit a proof of this bound is given in an appendix, where we follow the path prepared by Kolchinskii and Panchenko ([11]) and Bartlett and Mendelsson ([2]). In section 5 we study the Rademacher complexities of linear multi-task systems. In section 6 we give bounds for non-interacting systems, which are essentially equivalent to single-task learning, and derive bounds for proper, interacting multi-task learning, including the above mentioned results. We conclude with the construction of an example to demonstrate the advantages of multi-task learning.

## 2 Hilbert-Schmidt operators

For a fixed real, separable Hilbert space  $H$ , with inner product  $\langle \cdot, \cdot \rangle$  and norm  $\|\cdot\|$ , we define a second real, separable Hilbert space consisting of *Hilbert-Schmidt operators*. With  $HS$  we denote the real vector space of operators  $T$  on  $H$  satisfying  $\sum_{i=1}^{\infty} \|Te_i\|^2 \leq \infty$  for every orthonormal basis  $(e_i)_{i=1}^{\infty}$  of  $H$ . Every  $T \in HS$  is bounded. For  $S, T \in HS$  and an orthonormal basis  $(e_i)$  the series  $\sum_i \langle Se_i, Te_i \rangle$  is absolutely summable and independent of the chosen basis. The number  $\langle S, T \rangle_{HS} = \sum \langle Se_i, Te_i \rangle$  defines an inner product on  $HS$ , making it into

a Hilbert space. We denote the corresponding norm with  $\|\cdot\|_{HS}$  in contrast to the usual operator norm  $\|\cdot\|_\infty$  (see Reed and Simon [16] for background on functional analysis). We use  $HS^*$  to denote the set of symmetric Hilbert-Schmidt operators. For every member of  $HS^*$  there is a complete orthonormal basis of eigenvectors, and for  $T \in HS^*$  the norm  $\|T\|_{HS}$  is just the  $\ell_2$ -norm of its sequence of eigenvalues. With  $HS^+$  we denote the members of  $HS^*$  with only nonnegative eigenvalues.

We use two simple maps from  $H$  or  $H^2$  to  $HS$  to relate the geometries of objects in  $H$  to the geometry in  $HS$ .

**Definition 1.** Let  $x, y \in H$ . We define two operators  $Q_x$  and  $G_{x,y}$  on  $H$  by

$$\begin{aligned} Q_x z &= \langle z, x \rangle x, \quad \forall z \in H \\ G_{x,y} z &= \langle x, z \rangle y, \quad \forall z \in H. \end{aligned}$$

We will frequently use parts of the following lemma, the proof of which is very easy.

**Lemma 1.** Let  $x, y, x', y' \in H$  and  $T \in HS$ . Then

- (i)  $Q_x \in HS^+$  and  $\|Q_x\|_{HS} = \|x\|^2$ .
- (ii)  $\langle Q_x, Q_y \rangle_{HS} = \langle x, y \rangle^2$ .
- (iii)  $\langle T, Q_x \rangle_{HS} = \langle Tx, x \rangle$ .
- (iv)  $\langle T^*T, Q_v \rangle_{HS} = \|Tv\|^2$ .
- (v)  $Q_y Q_x = \langle x, y \rangle G_{x,y}$ .
- (vi)  $G_{x,y} \in HS$  and  $\|G_{x,y}\|_{HS} = \|x\| \|y\|$ .
- (vii)  $\langle G_{x,y}, G_{x',y'} \rangle_{HS} = \langle x, x' \rangle \langle y, y' \rangle$
- (viii)  $\langle T, G_{x,y} \rangle_{HS} = \langle Tx, y \rangle$ .

*Proof.* For  $x = 0$  (iii) is obvious. For  $x \neq 0$  chose an orthonormal basis  $(e_i)_{i=1}^\infty$ , so that  $e_1 = x / \|x\|$ . Then  $e_1$  is the only nonzero eigenvector of  $Q_x$  with eigenvalue  $\|x\| > 0$ . Also

$$\langle T, Q_x \rangle_{HS} = \sum_i \langle T e_i, Q_x e_i \rangle = \langle Tx, Q_x x \rangle / \|x\|^2 = \langle Tx, x \rangle,$$

which gives (iii). (ii), (i) and (iv) follow from substitution of  $Q_y$ ,  $Q_x$  and  $T^*T$  respectively for  $T$ . (v) follows directly from the definition when applied to any  $z \in H$ . Let  $(e_k)_{k=1}^\infty$  be any orthonormal basis. Then  $x = \sum_k \langle x, e_k \rangle e_k$ , so by boundedness of  $T$

$$\begin{aligned} \langle Tx, y \rangle &= \left\langle T \sum_k \langle x, e_k \rangle e_k, y \right\rangle = \sum_k \langle T e_k, \langle x, e_k \rangle y \rangle = \sum_k \langle T e_k, G_{x,y} e_k \rangle \\ &= \langle T, G_{x,y} \rangle_{HS}, \end{aligned}$$

which is (viii). Similarly

$$\begin{aligned} \langle G_{x,y}, G_{x',y'} \rangle_{HS} &= \sum_k \langle \langle x, e_k \rangle y, \langle x', e_k \rangle y' \rangle = \langle y, y' \rangle \sum_k \langle x, e_k \rangle \langle x', e_k \rangle \\ &= \langle x, x' \rangle \langle y, y' \rangle, \end{aligned}$$

which gives (vii) and (vi).  $\square$

The following application of Lemma 1 is the key to our main results.

**Lemma 2.** *Let  $T \in HS$  and  $w_1, \dots, w_m$  and  $v_1, \dots, v_m$  vectors in  $H$  with  $\|v_i\| \leq B$ . Then*

$$\sum_{l=1}^m \langle Tw_l, v_l \rangle \leq B \|T\|_{HS} \left( \sum_{l,r} |\langle w_l, w_r \rangle| \right)^{1/2}$$

and

$$\sum_{l=1}^m \langle Tw_l, v_l \rangle \leq Bm^{1/2} \|T^*T\|_{HS}^{1/2} \left( \sum_{l,r} \langle w_l, w_r \rangle^2 \right)^{1/4}$$

*Proof.* Without loss of generality assume  $B = 1$ . Using Lemma 1 (viii), Schwartz' inequality in  $HS$  and 1 (vii) we have

$$\begin{aligned} \sum_{l=1}^m \langle Tw_l, v_l \rangle &= \left\langle T, \sum_{l=1}^m G_{w_l, v_l} \right\rangle_{HS} \leq \|T\|_{HS} \left\| \sum_{l=1}^m G_{w_l, v_l} \right\|_{HS} \\ &= \|T\|_{HS} \left( \sum_{l,r} \langle w_l, w_r \rangle \langle v_l, v_r \rangle \right)^{1/2} \\ &\leq \|T\|_{HS} \left( \sum_{l,r} |\langle w_l, w_r \rangle| \right)^{1/2}. \end{aligned}$$

This proves the first inequality. Also, using Schwartz' inequality in  $H$  and  $\mathbb{R}^m$ , Lemma 1 (iv) and Schwartz' inequality in  $HS$

$$\begin{aligned} \sum_{l=1}^m \langle Tw_l, v_l \rangle &\leq \left( \sum_{l=1}^m \|v_l\|^2 \right)^{1/2} \left( \sum_{l=1}^m \|Tw_l\|^2 \right)^{1/2} \leq \sqrt{m} \left\langle T^*T, \sum_{l=1}^m Q_{w_l} \right\rangle_{HS}^{1/2} \\ &\leq \sqrt{m} \|T^*T\|_{HS}^{1/2} \left\| \sum_{l=1}^m Q_{w_l} \right\|_{HS}^{1/2} = \sqrt{m} \|T^*T\|_2^{1/2} \left( \sum_{l,r} \langle w_l, w_r \rangle^2 \right)^{1/4} \quad \square \end{aligned}$$

The set of  $d$ -dimensional, orthogonal projections in  $H$  is denoted with  $\mathcal{P}_d$ . We have  $\mathcal{P}_d \subset HS^*$  and if  $P \in \mathcal{P}_d$  then  $\|P\|_{HS} = \sqrt{d}$  and  $P^2 = P$ .

An operator  $T$  is called *trace-class* if  $\sum_{i=1}^{\infty} \langle Te_i, e_i \rangle$  is an absolutely convergent series for every orthonormal basis  $(e_i)_{i=1}^{\infty}$  of  $H$ . In this case the number  $tr(T) = \sum_{i=1}^{\infty} \langle Te_i, e_i \rangle$  is called the *trace* of  $T$  and it is independent of the chosen basis.

If  $\mathcal{A} \subset HS^*$  is a set of symmetric and bounded operators in  $H$  we use the notation

$$\|\mathcal{A}\|_{HS} = \sup \{ \|T\|_{HS} : T \in \mathcal{A} \} \text{ and } \mathcal{A}^2 = \{ T^2 : T \in \mathcal{A} \}.$$

### 3 Vector- and operator-valued random variables

Let  $(\Omega, \Sigma, \mu)$  be a probability space with expectation  $E[F] = \int_{\Omega} F d\mu$  for a random variable  $F : \Omega \rightarrow \mathbb{R}$ . Let  $X$  be a random variable with values in  $H$ , such that  $E[\|X\|] \leq \infty$ . The linear functional  $v \in H \mapsto E[\langle X, v \rangle]$  is bounded by  $E[\|X\|]$  and thus defines (by the Riesz Lemma) a unique vector  $E[X] \in H$  such that  $E[\langle X, v \rangle] = \langle E[X], v \rangle, \forall v \in H$ , with  $\|E[X]\| \leq E[\|X\|]$ .

We now look at the random variable  $Q_X$ , with values in  $HS$ . Suppose that  $E[\|X\|^2] \leq \infty$ . Passing to the space  $HS$  of Hilbert-Schmidt operators the above construction can be carried out again: By Lemma 1 (i)  $E[\|Q_X\|_{HS}] = E[\|X\|^2] \leq \infty$ , so there is a unique operator  $E[Q_X] \in HS$  such that  $E[\langle Q_X, T \rangle_{HS}] = \langle E[Q_X], T \rangle_{HS}, \forall T \in HS$ .

**Definition 2.** *The operator  $E[Q_X]$  is called the covariance operator of  $X$ .*

**Lemma 3.** *The covariance operator  $E[Q_X] \in HS^+$  has the following properties.*

- (i)  $\|E[Q_X]\|_{HS} \leq E[\|Q_X\|_{HS}]$ .
- (ii)  $\langle E[Q_X]y, z \rangle = E[\langle y, X \rangle \langle z, X \rangle], \forall y, z \in H$ .
- (iii)  $\text{tr}(E[Q_X]) = E[\|X\|^2]$ .
- (iv) For  $H$ -valued independent  $X_1$  and  $X_2$  with  $E[\|X_i\|^2] \leq \infty$  we have

$$\langle E[Q_{X_1}], E[Q_{X_2}] \rangle_{HS} = E[\langle X_1, X_2 \rangle^2].$$

*Proof.* (i) follows directly from the construction, (iv) from the identity  $\langle E[Q_{X_1}], E[Q_{X_2}] \rangle_{HS} = E[\langle Q_{X_1}, Q_{X_2} \rangle_{HS}]$ . Let  $y, z \in H$ . Then using 1 (viii) we get

$$\begin{aligned} \langle E[Q_X]y, z \rangle &= \langle E[Q_X], G_{y,z} \rangle_{HS} = E[\langle Q_X, G_{y,z} \rangle_{HS}] = E[\langle Q_X y, z \rangle] \\ &= E[\langle y, X \rangle \langle z, X \rangle] \end{aligned}$$

and hence (ii). We have with orthonormal basis  $(e_k)_{k=1}^{\infty}$  and using (ii)

$$\text{tr}(E[Q_X]) = \sum_k \langle E[Q_X]e_k, e_k \rangle = \sum_k E[\langle e_k, X \rangle \langle e_k, X \rangle] = E[\|X\|^2],$$

which gives (iii). Substitution of an eigenvector  $v$  for both  $y$  and  $z$  in (ii) shows that the corresponding eigenvalue must be nonnegative, so  $E[Q_X] \in HS^+$ .  $\square$

Property (ii) above is sometimes taken as the defining property of the covariance operator (see [4]).

If  $X$  is distributed uniformly on  $M \cap S_1$ , where  $M$  is a  $k$ -dimensional subspace and  $S_1$  the unit sphere in  $H$ , then  $E[\langle X, y \rangle^2] = \langle E[Q_X]y, y \rangle$  is zero if and only



if  $y \in M^\perp$ , so the range of  $E[Q_X]$  is  $M$ , so there are exactly  $k$ -eigenvectors corresponding to non-zero eigenvalues of  $E[Q_X]$ . By symmetry these eigenvalues must all be equal, and by (iii) above they sum up to 1, so  $E[Q_X]$  has the  $k$ -fold eigenvalue  $1/k$ , with zero as the only other eigenvalue. It follows that  $\|E[Q_X]\|_{HS} = 1/\sqrt{k}$ . We have given this derivation to illustrate the tendency of the Hilbert-Schmidt norm of the covariance operator of a distribution concentrated on unit vectors to decrease with the effective dimensionality of the distribution. This idea is relevant to the interpretation of our results.

The fact that  $HS$  is a separable Hilbertspace just as  $H$  allows to define an operator  $E[T]$  whenever  $T$  is a random variable with values in  $HS$  and  $E[\|T\|_{HS}] < \infty$ . Also any result valid in  $H$  has a corresponding analogue valid in  $HS$ . We quote a corresponding operator-version of a Theorem of Christianini and Shawe-Taylor [15] on the concentration of independent vector-valued random variables.

**Theorem 3.** *Suppose that  $T_1, \dots, T_m$  are independent random variables in  $H$  with  $\|T_i\| \leq 1$ . Then for all  $\delta > 0$  with probability greater than  $\delta$  we have*

$$\left\| \frac{1}{m} \sum_{i=1}^m E[T_i] - \frac{1}{m} \sum_{i=1}^m T_i \right\|_{HS} \leq \frac{2}{\sqrt{m}} \left( 1 + \sqrt{\frac{\ln(1/\delta)}{2}} \right).$$

Apply this with  $T_i = Q_{X_i}$  where the  $X_i$  are iid  $H$ -valued with  $\|X_i\| \leq 1$ . The theorem then shows that the covariance operator  $E[Q_X]$  can be approximated in  $HS$ -norm with high probability by the empirical estimates  $(1/m) \sum_i Q_{X_i}$ . The quantity

$$\left\| \sum_i Q_{X_i} \right\|_{HS} = \left( \sum_{i,j} \langle X_i, X_j \rangle^2 \right)^{1/2}$$

is the Frobenius norm of the Gramian (or kernel-) matrix  $\hat{C}(\mathbf{X})_{ij} = \langle X_i, X_j \rangle$ , denoted  $\|\hat{C}(\mathbf{X})\|_{Fr}$ . An immediate corollary to the above is, that  $(1/m) \|\hat{C}(\mathbf{X})\|_{Fr}$  is with high probability a good approximation of  $\|E[Q_X]\|_{HS}$ . In the proof of Theorem 2 we will not need this fact however.

## 4 Multi-task systems and general bounds

For our discussion of multi-task learning we concentrate on binary labeled data. Let  $(\Omega, \Sigma, \mu)$  be a probability space. We assume that there are  $m$  independent random variables  $Z^l = (X^l, Y^l) : \Omega \rightarrow \mathcal{X} \times \{-1, 1\}$ , where

- $l \in \{1, \dots, m\}$  identifies one of the  $m$  learning tasks,
- $X^l$  models the input data of the  $l$ -th task, distributed in a set  $\mathcal{X}$ , called the *input space*.
- $Y^l \in \{-1, 1\}$  models the output-, or label-data of the  $l$ -th task.
- For each  $l \in \{1, \dots, m\}$  there is an  $n$ -tuple of independent random variables  $(Z_i^l)_{i=1}^n = (X_i^l, Y_i^l)_{i=1}^n$ , where each  $Z_i^l$  is identically distributed to  $Z^l$ .

The random variable  $\mathbf{Z} = (Z_i^l)_{(i,l)=(1,1)}^{(n,m)}$  is called the *training sample* or training data. We also write  $\mathbf{X} = (X_i^l)_{(i,l)=(1,1)}^{(n,m)}$ . We use the superscript  $l$  to identify learning tasks running from 1 to  $m$ , the subscript  $i$  to identify data points in the sample, running from 1 to  $n$ . We will use the notations  $\mathbf{x} = (x_i^l)_{(i,l)=(1,1)}^{(n,m)}$  for generic members of  $(\mathcal{X}^n)^m$  and  $\mathbf{z} = (z_i^l)_{(i,l)=(1,1)}^{(n,m)} = (\mathbf{x}, \mathbf{y}) = (x_i^l, y_i^l)_{(i,l)=(1,1)}^{(n,m)}$  for generic members of  $((\mathcal{X} \times \{-1, 1\})^n)^m$ .

A *multiclassifier* is a map  $\mathbf{h} : \mathcal{X} \rightarrow \{-1, 1\}^m$ . We write  $\mathbf{h} = (h^1, \dots, h^m)$  and interpret  $h^l(x)$  as the label assigned to the vector  $x$  when the task is known to be  $l$ . The average error of a multiclassifier  $\mathbf{h}$  is the quantity

$$\text{er}(\mathbf{h}) = \frac{1}{m} \sum_{l=1}^m \Pr \{h^l(X^l) \neq Y^l\},$$

which is just the misclassification probability averaged over all tasks. Typically a classifier is chosen from some candidate set minimizing some error estimate based on the training data  $\mathbf{Z}$ . Here we consider zero-threshold classifiers  $\mathbf{h}_{\mathbf{f}}$  which arise as follows:

Suppose that  $\mathcal{F}$  is a class of vector valued functions  $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^m$  with  $\mathbf{f} = (f^1, \dots, f^m)$ . A function  $\mathbf{f} \in \mathcal{F}$  defines a multi-classifier  $\mathbf{h}_{\mathbf{f}} = (h_{\mathbf{f}}^1, \dots, h_{\mathbf{f}}^m)$  through  $h_{\mathbf{f}}^l(x) = \text{sign}(f^l(x))$ . To give uniform error bounds for such classifiers in terms of empirical estimates, we define for  $\gamma > 0$  the margin functions

$$\phi_{\gamma}(t) = \begin{cases} 1 & \text{if } t \leq 0 \\ 1 - t/\gamma & \text{if } 0 < t < \gamma \\ 0 & \text{if } \gamma \leq t \end{cases},$$

and for  $\mathbf{f} \in \mathcal{F}$  the random variable

$$\hat{\text{er}}_{\gamma}(\mathbf{f}) = \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n \phi_{\gamma}(Y_i^l f^l(X_i^l)),$$

called the *empirical  $\gamma$ -margin error* of  $\mathbf{f}$ . The following Theorem (taken from [2] with minor modifications to adapt to the multi-task situation and combined with the model-selection lemma 15.5 in [1], Lemma 5 in this paper) gives a bound on  $\text{er}(\mathbf{h}_{\mathbf{f}})$  in terms of  $\hat{\text{er}}_{\gamma}(\mathbf{f})$ , valid with high probability uniformly in  $\mathbf{f} \in \mathcal{F}$  and  $\gamma$ .

**Theorem 4.** *Let  $\epsilon, \delta \in (0, 1)$*

(i) *With probability greater than  $1 - \delta$  it holds for all  $\mathbf{f} \in \mathcal{F}$  and all  $\gamma \in (0, 1)$  that*

$$\text{er}(\mathbf{h}_{\mathbf{f}}) \leq \hat{\text{er}}_{\gamma}(\mathbf{f}) + \frac{1}{\gamma(1-\epsilon)} E \left[ \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) \right] + \sqrt{\frac{\ln(1/(\delta\gamma\epsilon))}{2nm}}.$$

(ii) *With probability greater than  $1 - \delta$  it holds for all  $\mathbf{f} \in \mathcal{F}$  and all  $\gamma \in (0, 1)$  that*

$$\text{er}(\mathbf{h}_{\mathbf{f}}) \leq \hat{\text{er}}_{\gamma}(\mathbf{f}) + \frac{1}{\gamma(1-\epsilon)} \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) + \sqrt{\frac{9 \ln(2/(\delta\gamma\epsilon))}{2nm}}.$$

Here  $\hat{\mathcal{R}}_n^m(\mathcal{F})$  is the *empirical Rademacher complexity* in the sense of the following

**Definition 3.** Let  $\{\sigma_i^l : l \in \{1, \dots, m\}, i \in \{1, \dots, n\}\}$  be a collection of independent random variables, distributed uniformly in  $\{-1, 1\}$ . The empirical Rademacher complexity of a class  $\mathcal{F}$  of functions  $\mathbf{f} : \mathcal{X} \rightarrow \mathbb{R}^m$  is the function  $\hat{\mathcal{R}}_n^m(\mathcal{F})$  defined on  $(\mathcal{X}^n)^m$  by

$$\hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{x}) = E_\sigma \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{2}{mn} \sum_{l=1}^m \sum_{i=1}^n \sigma_i^l f^l(x_i^l) \right].$$

For the readers convenience we give a proof of Theorem 4 in the appendix.

The bounds in the Theorem each involve three terms. The last one expresses the dependence of the estimation error on the confidence parameter  $\delta$  and a model-selection penalty  $\ln(1/(\gamma\epsilon))$  for the choice of the margin  $\gamma$ . Note that it generally decreases as  $1/\sqrt{nm}$ . This is not an a priori advantage of multi-task learning, but a trivial consequence of the fact that we estimate an average of  $m$  probabilities (in contrast to Ben David [7] where bounds are valid for each individual task - of course under more restrictive assumptions). The  $1/\sqrt{nm}$  decay however implies that even for moderate values of  $m$  and  $n$  the parameter  $\epsilon$  in Theorem 4 can be chosen very small, so that the factor  $1/(1-\epsilon)$  in the second term on the right of the two bounds is very close to unity.

The second term involves the complexity of the function class  $\mathcal{F}$ , either as measured in terms of the distribution of the random variable  $\mathbf{X}$  or in terms of the observed sample. Since the distribution of  $\mathbf{X}$  is unobservable in practice, the bound (i) is primarily of theoretical importance, while (ii) can be used to drive an algorithm which selects the multi-classifier  $\mathbf{h}_{\mathbf{f}^*}$ , where  $(\mathbf{f}^*, \gamma) \in \mathcal{F} \times (0, 1)$  are chosen to minimize the right side of the bound with given  $\delta, \epsilon$ . It is questionable if minimizing upper bounds is a good strategy, but it can serve as a motivating guideline.

Of key importance in the analysis of these algorithms is the empirical Rademacher complexity  $\hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X})$ , as observed on the sample  $\mathbf{X}$ , and its expectation, measuring respectively the sample- and distribution-dependent complexities of the function class  $\mathcal{F}$ . Bounds on these quantities can be substituted in Theorem 4 to give corresponding error bounds.

## 5 The Rademacher complexity of linear multi-task learning

We will now concentrate on multi-task learning in the linear case, when the data lives in a real, separable Hilbert space  $H$ , by means of some kernel induced-embedding (see [9]), the details of which will not concern us at this point. We therefore take  $H$  as input space  $\mathcal{X}$ , so that the random variables  $X^l$  take values in  $H$  for all  $l \in \{1, \dots, m\}$ , and we generally require  $\|X^l\| \leq 1$ . The case  $\|X^l\| = 1$  where the data is constrained to the unit sphere in  $H$  is of particular interest, corresponding to a class of radial basis function kernels.

We write  $C^l$  for the covariance operator  $E[Q_{X^l}]$  and  $C$  for the total covariance operator  $C = (1/m) \sum_l C^l$ , corresponding to a uniform mixture of distributions. By Lemma 3 we have  $\|C^l\|_{HS} \leq \text{tr}(C^l) = E[\|X^l\|^2] \leq 1$ .

Let  $B > 0$ , let  $T$  be a fixed symmetric, bounded linear operator on  $H$  with  $\|T\|_\infty \leq 1$ , and let  $\mathcal{A}$  be a set of symmetric, bounded linear operators  $T$  on  $H$ , all satisfying  $\|T\|_\infty \leq 1$ . We will consider the vector-valued function classes

$$\begin{aligned} \mathcal{F}_B &= \{x \in H \mapsto (v_1, \dots, v_m)(x) := (\langle x, v_1 \rangle, \dots, \langle x, v_m \rangle) : \|v_i\| \leq B\} \\ \mathcal{F}_B \circ T &= \{x \in H \mapsto (v_1, \dots, v_m) \circ T(x) := (\langle Tx, v_1 \rangle, \dots, \langle Tx, v_m \rangle) : \|v_i\| \leq B\} \\ \mathcal{F}_B \circ \mathcal{A} &= \{x \in H \mapsto (v_1, \dots, v_m) \circ T(x) : \|v_i\| \leq B, T \in \mathcal{A}\}. \end{aligned}$$

The algorithms which chose from  $\mathcal{F}_B$  and  $\mathcal{F}_B \circ T$  are essentially trivial extensions of linear single-task learning, where the tasks do not interact in the selection of the individual classifiers  $v_i$ , which are chosen independently. In the case of  $\mathcal{F}_B \circ T$  the preprocessing operator  $T$  is chosen before seeing the training data. Since  $\|T\|_\infty \leq 1$  we have  $\mathcal{F}_B \circ T \subseteq \mathcal{F}_B$ , so that we can expect a reduced complexity for  $\mathcal{F}_B \circ T$  and the key question becomes if the choice of  $T$  (possibly based on experience with other data) was lucky enough to allow for a sufficiently low empirical error.

The non-interacting classes  $\mathcal{F}_B$  and  $\mathcal{F}_B \circ T$  are important for comparison to  $\mathcal{F}_B \circ \mathcal{A}$  which represents proper multi-task learning. Here the preprocessing operator  $T$  is selected from  $\mathcal{A}$  in response to the data. The constraint that  $T$  be the same for all tasks forces an interaction of tasks in the choice of  $T$  and  $(v_1, \dots, v_m)$ , deliberately aiming for a low empirical error. At the same time we also have  $\mathcal{F}_B \circ \mathcal{A} \subseteq \mathcal{F}_B$ , so that again a reduced complexity is to be expected, giving a smaller contribution to the estimation error. The promise of multi-task learning is based on the combination of these two ideas: Aiming for a low empirical error, using a function class of reduced complexity.

We first look at the complexity of the function class  $\mathcal{F}_B$ . The proof of the following lemma is essentially the same as the proof of Lemma 22 in [2].

**Lemma 4.** *We have*

$$\begin{aligned} \hat{\mathcal{R}}_n^m(\mathcal{F}_B)(\mathbf{x}) &\leq \frac{2B}{nm} \sum_{l=1}^m \left( \sum_{i=1}^n \|x_i^l\|^2 \right)^{1/2} \\ E \left[ \hat{\mathcal{R}}_n^m(\mathcal{F}_B)(\mathbf{X}) \right] &\leq \frac{2B}{\sqrt{n}} \frac{1}{m} \sum_{l=1}^m \left( E \left[ \|X^l\|^2 \right] \right)^{1/2} = \frac{2B}{\sqrt{n}} \frac{1}{m} \sum_{l=1}^m \text{tr}(C^l)^{1/2} \end{aligned}$$

*Proof.* Using Schwartz' and Jensen's inequality and the independence of the  $\sigma_i^l$  we get

$$\begin{aligned} \hat{\mathcal{R}}_n^m(\mathcal{F}_B)(\mathbf{x}) &= E_\sigma \left[ \sup_{v_1, \dots, v_m, \|v_l\| \leq B} \frac{2}{nm} \sum_{l=1}^m \left\langle \sum_{i=1}^n \sigma_i^l x_i^l, v_l \right\rangle \right] \\ &\leq BE_\sigma \left[ \frac{2}{nm} \sum_{l=1}^m \left\| \sum_{i=1}^n \sigma_i^l x_i^l \right\| \right] \\ &\leq \frac{2B}{nm} \sum_{l=1}^m \left( E_\sigma \left[ \left\| \sum_{i=1}^n \sigma_i^l x_i^l \right\|^2 \right] \right)^{1/2} \\ &= \frac{2B}{nm} \sum_{l=1}^m \left( \sum_{i=1}^n \|x_i^l\|^2 \right)^{1/2}. \end{aligned}$$

Jensen's inequality then gives the second conclusion  $\square$

The first bound in the lemma is just the average of the bounds given by Bartlett and Mendelson in [2] on the empirical complexities for the various task-components of the sample. For inputs constrained to the unit sphere in  $H$ , when  $\|X^l\| = 1$ , both bounds become  $2B/\sqrt{n}$ , which sets the mark for comparison with the interacting case  $\mathcal{F}_B \circ \mathcal{A}$ . For motivation we next look at the case  $\mathcal{F}_B \circ T$ , working with a fixed linear preprocessor  $T$  of operator norm bounded by 1. Using the above bound we obtain

$$\hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ T)(\mathbf{x}) = \hat{\mathcal{R}}_n^m(\mathcal{F}_B)(T\mathbf{x}) \leq \frac{2B}{nm} \sum_{l=1}^m \left( \sum_{i=1}^n \|Tx_i^l\|^2 \right)^{1/2}, \quad (4)$$

which is always bounded by  $B/\sqrt{n}$ , because  $\|Tx\| \leq \|x\|, \forall x$ . Using Lemma 1 (iv) we can rewrite the right side above as

$$\frac{2B}{\sqrt{n}} \frac{1}{m} \sum_{l=1}^m \left\langle T^2, \frac{1}{n} \sum_{i=1}^n Q_{x_i^l} \right\rangle_{HS}^{1/2}.$$

Taking the expectation and using the concavity of the root function gives, with two applications of Jensen's inequality and an application of Schwartz' inequality (in  $HS$ ),

$$E \left[ \hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ T)(\mathbf{X}) \right] \leq \frac{2B}{\sqrt{n}} \|T^2\|_{HS}^{1/2} \|C\|_{HS}^{1/2},$$

which can be significantly smaller than  $B/\sqrt{n}$ , for example if  $T$  is a  $d$ -dimensional projection, and the data-distribution is spread well over a much more than  $d$ -dimensional submanifold of the unit ball in  $H$ , as explained in the introduction and section 3. If we substitute the bound above in Theorem 4 we obtain an inequality which looks like (3) in the limit  $m \rightarrow \infty$ .

We now consider the case where  $T$  is chosen from some set  $\mathcal{A}$  (symmetric, bounded) candidate operators on the basis of the same sample  $\mathbf{X}$ , simultaneous to the determination of the classification vectors  $v_1, \dots, v_l$ . We give two bounds each for the Rademacher complexity and its expectation, one which is somewhat similar to other bounds for multi-task learning (e.g. (2)) and another one which is tighter in the limit when the number of tasks  $m$  goes to infinity.

**Theorem 5.** *The following inequalities hold*

$$\hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ \mathcal{A})(\mathbf{x}) \leq \frac{2B \|\mathcal{A}\|_{HS}}{\sqrt{n}} \sqrt{\frac{1}{mn} \|\hat{C}(\mathbf{x})\|_{Fr} + \frac{1}{m}} \quad (5)$$

$$\hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ \mathcal{A})(\mathbf{x}) \leq \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{\sqrt{n}} \left( \left( \frac{1}{mn} \|\hat{C}(\mathbf{x})\|_{Fr} \right)^2 + \frac{2}{m} \right)^{1/4} \quad (6)$$

$$E \left[ \hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ \mathcal{A})(\mathbf{X}) \right] \leq \frac{2B \|\mathcal{A}\|_{HS}}{\sqrt{n}} \sqrt{\|C\|_{HS} + \frac{1}{m}} \quad (7)$$

$$E \left[ \hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ \mathcal{A})(\mathbf{X}) \right] \leq \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{\sqrt{n}} \left( \|C\|_{HS}^2 + \frac{3}{m} \right)^{1/4}. \quad (8)$$

*Proof.* Fix  $\mathbf{x}$  and define vectors  $w_l = w_l(\sigma) = \sum_{i=1}^n \sigma_i^l x_i^l$  depending on the Rademacher variables  $\sigma_i^l$ . Then by Lemma 2 and Jensen's inequality

$$\begin{aligned} \hat{\mathcal{R}}(\mathcal{F}_B \circ \mathcal{A})(\mathbf{x}) &= E_\sigma \left[ \sup_{T \in \mathcal{A}} \sup_{v_1, \dots, v_m, \|v_l\| \leq B} \frac{2}{nm} \sum_{l=1}^m \langle T w_l, v_l \rangle \right] \quad (9) \\ &\leq \frac{2B}{nm} \|\mathcal{A}\|_{HS} E_\sigma \left[ \left( \sum_{l,r} |\langle w_l, w_r \rangle| \right)^{1/2} \right] \\ &\leq \frac{2B}{nm} \|\mathcal{A}\|_{HS} \left( \sum_{l,r} E_\sigma [|\langle w_l, w_r \rangle|] \right)^{1/2}. \end{aligned}$$

Now we have

$$E_\sigma \left[ \|w_l\|^2 \right] = \sum_{i=1}^n \sum_{j=1}^n E_\sigma [\sigma_i^l \sigma_j^l] \langle x_i^l, x_j^l \rangle = \sum_{i=1}^n \|x_i^l\|^2. \quad (10)$$

Also, for  $l \neq r$ , we get, using Jensen's inequality and independence of the Rademacher variables,

$$\begin{aligned} (E_\sigma [|\langle w_l, w_r \rangle|])^2 &\leq E_\sigma \left[ \langle w_l, w_r \rangle^2 \right] \quad (11) \\ &= \sum_{i=1}^n \sum_{j=1}^n \sum_{i'=1}^n \sum_{j'=1}^n E_\sigma [\sigma_i^l \sigma_j^r \sigma_{i'}^l \sigma_{j'}^r] \langle x_i^l, x_j^r \rangle \langle x_{i'}^l, x_{j'}^r \rangle \\ &= \sum_{i,j=1}^n \langle x_i^l, x_j^r \rangle^2. \end{aligned}$$

Taking the square-root and inserting it together with (10) in (9) we obtain the following intermediate bound

$$\hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ \mathcal{A})(\mathbf{x}) \leq \frac{2B \|\mathcal{A}\|_{HS}}{nm} \left( \sum_{l=1}^m \sum_{i=1}^n \|x_i^l\|^2 + \sum_{l \neq r} \left( \sum_{i,j=1}^n \langle x_i^l, x_j^r \rangle^2 \right)^{1/2} \right)^{1/2} \quad (12)$$

By Jensen's inequality we have

$$\frac{1}{m^2} \sum_{l \neq r} \left( \sum_{i,j=1}^n \langle x_i^l, x_j^r \rangle^2 \right)^{1/2} \leq \left( \frac{1}{m^2} \sum_{l,r=1}^m \sum_{i,j=1}^n \langle x_i^l, x_j^r \rangle^2 \right)^{1/2} = \frac{1}{m} \|\hat{C}(\mathbf{x})\|_{Fr},$$

which together with (12) and  $\|x_i^l\| \leq 1$  implies (5).

To prove (6) first use the second part of Lemma 2 and Jensen's inequality to get

$$\hat{\mathcal{R}}(\mathcal{F}_B \circ \mathcal{A})(\mathbf{x}) \leq \frac{2B}{n\sqrt{m}} \|\mathcal{A}^2\|_{HS}^{1/2} \left( \sum_{l,r} E_\sigma [\langle w_l, w_r \rangle^2] \right)^{1/4}. \quad (13)$$

Now we have  $E_\sigma [\sigma_i^l \sigma_j^l \sigma_{i'}^{l'} \sigma_{j'}^{l'}] \leq \delta_{ij} \delta_{i'j'} + \delta_{ii'} \delta_{jj'} + \delta_{ij'} \delta_{ji'}$  so

$$\begin{aligned} E_\sigma [\langle w_l, w_l \rangle^2] &\leq \sum_{i,j=1}^n \left( \|X_i^l\|^2 \|X_j^l\|^2 + 2 \langle X_i^l, X_j^l \rangle^2 \right) \\ &\leq 2 \left( \sum_{i=1}^n \|X_i^l\|^2 \right)^2 + \sum_{i,j=1}^n \langle X_i^l, X_j^l \rangle^2 \end{aligned}$$

Inserting this together with (11) in (13) gives

$$\hat{\mathcal{R}}_n^m(\mathcal{F}_B \circ \mathcal{A})(\mathbf{x}) \leq \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{n\sqrt{m}} \left( 2 \sum_{l=1}^m \left( \sum_{i=1}^n \|x_i^l\|^2 \right)^2 + \sum_{l,r=1}^m \sum_{i,j=1}^n \langle x_i^l, x_j^r \rangle^2 \right)^{1/4}, \quad (14)$$

which together with  $\|x_i^l\| \leq 1$  gives (6).

Taking the expectation of (12), using Jensen's inequality,  $\|X^l\| \leq 1$  and independence of  $X^l$  and  $X^r$  for  $l \neq r$ , and Jensen's inequality again, we get

$$\begin{aligned}
& E \left[ \hat{\mathcal{R}}_n^m (\mathcal{F}_B \circ \mathcal{A}) (\mathbf{X}) \right] \\
& \leq \frac{2B \|\mathcal{A}\|_{HS}}{nm} \left( nm + \sum_{l \neq r} \left( E \left[ \sum_{i,j=1}^n \langle X_i^l, X_j^r \rangle^2 \right] \right)^{1/2} \right)^{1/2} \\
& = \frac{2B \|\mathcal{A}\|_{HS}}{nm} \left( \sum_{l \neq r} \left( E \left[ \left\langle \frac{1}{n} \sum_{i=1}^n Q_{X_i^l}, \frac{1}{n} \sum_{j=1}^n Q_{X_j^r} \right\rangle_{HS} \right] \right)^{1/2} + nm \right)^{1/2} \\
& = \frac{2B \|\mathcal{A}\|_{HS}}{\sqrt{n}} \left( \frac{1}{m^2} \sum_{l \neq r} \langle E[Q_{X^l}], E[Q_{X^r}] \rangle_{HS}^{1/2} + \frac{1}{m} \right)^{1/2} \\
& \leq \frac{2B \|\mathcal{A}\|_{HS}}{\sqrt{n}} \left( \left\langle \frac{1}{m} \sum_{l=1}^m E[Q_{X^l}], \frac{1}{m} \sum_{r=1}^m E[Q_{X^r}] \right\rangle_{HS}^{1/2} + \frac{1}{m} \right)^{1/2},
\end{aligned}$$

which gives (7). In a similar way we obtain from (14)

$$\begin{aligned}
& E \left[ \hat{\mathcal{R}}_n^m (\mathcal{F}_B \circ \mathcal{A}) (\mathbf{X}) \right] \\
& \leq \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{n\sqrt{m}} \left( 2mn^2 + mn + \sum_{l,r} \sum_{\substack{i,j=1 \\ (l,i) \neq (r,j)}}^n \langle E[Q_{X_i^l}], E[Q_{X_j^r}] \rangle^2 \right)^{1/4} \\
& \leq \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{n\sqrt{m}} \left( 3mn^2 + m^2n^2 \left\| E \left[ \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n Q_{X_i^l} \right] \right\|_{HS}^2 \right)^{1/4},
\end{aligned}$$

which gives (8)  $\square$

## 6 Bounds for multi-task learning

Inserting the bounds of Theorem 5 in Theorem 4 immediately gives

**Theorem 6.** *Let  $\mathcal{A}$  be a set of bounded, symmetric operators in  $H$  and  $\epsilon, \delta \in (0, 1)$*

*(i) With probability greater than  $1 - \delta$  it holds for all  $\mathbf{f} = (v_1, \dots, v_m) \circ T \in \mathcal{F}_B \circ \mathcal{A}$  and all  $\gamma \in (0, 1)$  that*

$$er(\mathbf{h}_{\mathbf{f}}) \leq e\hat{r}_{\gamma}(\mathbf{f}) + \frac{1}{\gamma(1-\epsilon)}A + \sqrt{\frac{\ln(1/(\delta\gamma\epsilon))}{2nm}},$$



where  $A$  is either

$$A = \frac{2B \|\mathcal{A}\|_{HS}}{\sqrt{n}} \sqrt{\|C\|_{HS} + \frac{1}{m}} \quad (15)$$

or

$$A = \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{\sqrt{n}} \left( \|C\|_{HS}^2 + \frac{3}{m} \right)^{1/4}. \quad (16)$$

(ii) With probability greater than  $1 - \delta$  it holds for all  $\mathbf{f} = (v_1, \dots, v_m) \circ T \in \mathcal{F}_B \circ \mathcal{A}$  and for all  $\gamma \in (0, 1)$  that

$$er(\mathbf{h}_{\mathbf{f}}) \leq e\hat{r}_{\gamma}(\mathbf{f}) + \frac{1}{\gamma(1-\epsilon)} A(\mathbf{X}) + \sqrt{\frac{9 \ln(2/(\delta\gamma\epsilon))}{2nm}},$$

where the random variable  $A(\mathbf{X})$  is either

$$A(\mathbf{X}) = \frac{2B \|\mathcal{A}\|_{HS}}{\sqrt{n}} \sqrt{\frac{1}{mn} \|\hat{C}(\mathbf{x})\|_{Fr} + \frac{1}{m}}$$

or

$$A(\mathbf{X}) = \frac{2B \|\mathcal{A}^2\|_{HS}^{1/2}}{\sqrt{n}} \left( \left( \frac{1}{mn} \|\hat{C}(\mathbf{x})\|_{Fr} \right)^2 + \frac{2}{m} \right)^{1/4}.$$

We can invoke again Lemma 5 to stratify over different operator norms. This will give a very similar looking result, which we state in abbreviated fashion.

**Theorem 7.** *Theorem 6 holds with the following modifications:*

- The class  $\mathcal{F}_B \circ \mathcal{A}$  is replaced by all of  $\mathcal{F}_B \circ HS^*$ .
- $\|\mathcal{A}\|_{HS}$  and  $\|\mathcal{A}^2\|_{HS}$  are replaced by  $\|T\|_{HS} \vee 1$  and  $\|T^2\|_{HS} \vee 1$  respectively.
- $(1 - \epsilon)$  and  $(\delta\gamma\epsilon)$  are replaced by  $(1 - \epsilon)^2$  and  $(\delta\gamma\epsilon^2)$  respectively.

The passage from  $\|T\|_{HS}$  to  $\|T\|_{HS} \vee 1$  is an artifact introduced by the stratification. We could also require  $\|T\|_{HS} \geq 1$ . Setting  $\epsilon = 1/2$  gives Theorem 1 and Theorem 2.

## 7 An Example

We conclude with an idealized example of a multi-task system to which our bounds can be applied.

Fix numbers  $m, k \in \mathbb{N}$  and a 'radius'  $\rho > 0$ . For each  $l \in \{1, \dots, m\}$  let  $\alpha^l$  be a real random variable distributed uniformly on the set

$$\left[ \frac{(l-1)\pi}{m}, \frac{l\pi}{m} \right) \cup \left[ \pi + \frac{(l-1)\pi}{m}, \pi + \frac{l\pi}{m} \right),$$

such that  $\alpha^l$  and  $\alpha^r$  are independent for  $l \neq r$ . For  $l = 1, \dots, m$  let  $V^l$  be independent random variables distributed uniformly on a sphere of radius  $\sqrt{1 - \rho^2}$

in  $\mathbb{R}^k$ . All the  $\alpha^l$  and all the  $V^l$  are now mutually independent. We now define a random variable  $X^l$  with values in  $\mathbb{R}^{k+2} \subset \ell_2$  by

$$X^l = (\rho \cos \alpha^l, \rho \sin \alpha^l, V_1^l, \dots, V_k^l)$$

and a labeling variable  $Y^l$  with values in  $\{-1, 1\}$  by

$$Y^l = \begin{cases} 1 & \text{if } \alpha^l \in \left[ \frac{(l-1)\pi}{m}, \frac{l\pi}{m} \right) \\ -1 & \text{if } \alpha^l \in \left[ \pi + \frac{(l-1)\pi}{m}, \pi + \frac{l\pi}{m} \right) \end{cases}.$$

We have defined a multi-task system with the following properties:

- The data is distributed on the unit-sphere.
- All the relevant data is contained in the first two coordinates.
- The remaining  $k$  coordinates are filled up with noise of an amplitude  $\sqrt{1 - \rho^2}$ .
- The optimal unit vector to classify the  $l$ -th task is given by  $(\cos \beta^l, \sin \beta^l, 0, \dots, 0)$  where  $\beta^l = (l - 1/2)/m$ . It has a margin of  $\rho_m = \rho \cos(\pi/m) \rightarrow \rho$  as  $m \rightarrow \infty$ .

To appreciate the difficulty of learning imagine the variable  $X^l$  to be 'shuffled' by an unknown unitary transformation  $U$  on  $\ell_2$ , that is  $X^l \leftarrow U \circ X^l$ .

Suppose we have a margin error of zero at margin  $\rho_m$ , for unit vectors  $v^1, \dots, v^m$ . How certain can we be, that our classifiers are any good? We will study the behaviour of our bounds for non-interacting and interacting learning under the following conditions:

- The dimensionality  $k$  of the noise is large.
- The amplitude  $\rho$  of the relevant coordinate values is very small, so that they become buried in irrelevant information.
- The number of learning tasks  $m$  is large.

Since we have a margin error of zero we only need to consider the bounds on the estimation error. Moreover, because  $m$  is large we will neglect the last term in our bounds, which depends on the confidence parameter. According to Theorem 4 we are left with the term

$$\frac{E \left[ \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) \right]}{\rho_m}$$

as an error bound, depending on the function class  $\mathcal{F}$  from which the multi-classifier was chosen. For non-interacting learning, which also corresponds to the single task case in (Bartlett Mendelson) we immediately obtain the bound

$$\frac{2}{\rho_m \sqrt{n}} \tag{17}$$

or a sample complexity (required sample size)  $M(\epsilon, \rho, k, m)$ , depending on a maximal allowed error and the other parameters  $\rho$ ,  $k$  and  $m$ , of

$$M(\epsilon, \rho, k) = \frac{4}{\rho_m^2 \epsilon^2}.$$

To compute the bound for interacting multi-task learning we first need to find  $\|C\|_{HS}$ . A lengthy calculation gives

$$\|C\|_{HS}^2 = \frac{\rho^2}{2} + \frac{1 - \rho^2}{k}.$$

If we substitute this into (16) of Theorem 6 and work with  $d$ -dimensional projections we obtain the bound

$$\frac{2d^{1/4}}{\rho_m \sqrt{n}} \left( \frac{\rho^2}{2} + \frac{1 - \rho^2}{k} + \frac{3}{m} \right)^{1/4}. \quad (18)$$

This is superior to the non-interacting bound (17) if

$$d \left( \frac{\rho^2}{2} + \frac{1 - \rho^2}{k} + \frac{3}{m} \right) < 1,$$

which will be the case for sufficiently large  $m$  and  $k$  if  $d\rho^2/2 < 1$ . Since we work with small values of  $\rho$ , the latter will happen even for unnecessarily large values of  $d$ , but of course the best choice for interacting learning is with 2-dimensional projections.

With 2-dimensional projections, in the limit of a large number of tasks  $m$  and a large dimension  $k$  of the noise distribution, the interacting bound (18) becomes

$$\frac{2}{\sqrt{\rho n}},$$

which in comparison with the non-interacting bound (17) shows that an effective margin has been improved from  $\rho$  to  $\sqrt{\rho}$ . There is a corresponding reduction in sample complexity, which for interacting learning now is

$$M(\epsilon, \rho, k) = \frac{4}{\rho \epsilon^2}.$$

For  $\rho \approx 10^{-1}$  this already amounts to a factor of 10.

We conclude with the remark that previous bounds on multi-task learning as in [6] and [18] suffer from linear scaling with the dimension  $k$  and behave poorly on on this example.

## Appendix

In this section we give a proof of Theorem 4 for the readers convenience. Most of this material is combined from [1], [2], [3] and [18], and we make no claim to originality for any of it. A preliminary result is

**Theorem 8.** Let  $\mathcal{F}$  be a  $[0, 1]^m$ -valued function class on a space  $\mathcal{X}$ , and  $\mathbf{X} = (X_i^l)_{(l,i)=(1,1)}^{(m,n)}$  a vector of  $\mathcal{X}$ -valued independent random variables where for fixed  $l$  and varying  $i$  all the  $X_i^l$  are identically distributed. Fix  $\delta > 0$ . Then with probability greater than  $1 - \delta$  we have for all  $\mathbf{f} = (f^1, \dots, f^m) \in \mathcal{F}$

$$\frac{1}{m} \sum_{l=1}^m E[f^l(X_1^l)] \leq \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n f^l(X_i^l) + \mathcal{R}_n^m(\mathcal{F}) + \sqrt{\frac{\ln(1/\delta)}{2mn}}.$$

We also have with probability greater than  $1 - \delta$  for all  $\mathbf{f} = (f^1, \dots, f^m) \in \mathcal{F}$ , that

$$\frac{1}{m} \sum_{l=1}^m E[f^l(X_1^l)] \leq \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n f^l(X_i^l) + \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) + \sqrt{\frac{9 \ln(2/\delta)}{2mn}}.$$

*Proof.* Let  $\Psi$  be the function on  $\mathcal{X}^{mn}$  given by

$$\Psi(\mathbf{x}) = \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{m} \sum_{l=1}^m \left( E[f^l(X_1^l)] - \frac{1}{n} \sum_{i=1}^n f^l(X_i^l) \right)$$

and let  $\mathbf{X}'$  be an iid copy of the  $\mathcal{X}^{mn}$ -valued random variable  $\mathbf{X}$ . Then

$$\begin{aligned} E[\Psi(\mathbf{X})] &= E_{\mathbf{X}} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{mn} E_{\mathbf{X}'} \left[ \sum_{l=1}^m \sum_{i=1}^n \left( f^l((X_i^l)') - f^l(X_i^l) \right) \right] \right] \\ &\leq E_{\mathbf{X}\mathbf{X}'} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n \left( f^l((X_i^l)') - f^l(X_i^l) \right) \right] \\ &= E_{\mathbf{X}\mathbf{X}'} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n \sigma_i^l \left( f^l((X_i^l)') - f^l(X_i^l) \right) \right], \end{aligned}$$

for any realization  $\boldsymbol{\sigma} = (\sigma_i^l)$  of the Rademacher variables, because the expectation  $E_{\mathbf{X}\mathbf{X}'}$  is symmetric under the exchange  $(X_i^l)' \longleftrightarrow X_i^l$ . Hence

$$E[\Psi(\mathbf{X})] \leq E_{\mathbf{X}} E_{\boldsymbol{\sigma}} \left[ \sup_{\mathbf{f} \in \mathcal{F}} \frac{2}{mn} \sum_{l=1}^m \sum_{i=1}^n \sigma_i^l f^l(X_i^l) \right] = \mathcal{R}_n^m(\mathcal{F}).$$

Now fix  $\mathbf{x} \in \mathcal{X}^{mn}$  and let  $\mathbf{x}' \in \mathcal{X}^{mn}$  be as  $\mathbf{x}$ , except for one modified coordinate  $(x_i^l)'$ . Since each  $f^l$  has values in  $[0, 1]$  we have  $|\Psi(\mathbf{x}) - \Psi(\mathbf{x}')| \leq 1/(mn)$ . So by the one-sided version of the bounded difference inequality (see McDiarmid [12])

$$\Pr \left\{ \Psi(\mathbf{X}) > E_{\mathbf{X}'} [\Psi(\mathbf{X}')] + \sqrt{\frac{\ln(1/\delta)}{2mn}} \right\} \leq \delta.$$

Together with the above bound on  $E[\Psi(\mathbf{X})]$  and the definition of  $\Psi$  this gives the first conclusion.

With  $\mathbf{x}$  and  $\mathbf{x}'$  as above we have  $\left| \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{x}) - \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{x}') \right| \leq 2/(mn)$ , so by the other tail of the bounded difference inequality

$$\Pr \left\{ \mathcal{R}_n^m(\mathcal{F}) < \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) + \sqrt{\frac{4 \ln(1/\delta)}{2mn}} \right\} \leq \delta,$$

which, combined with the first conclusion in a union bound, gives the second conclusion  $\square$

We quote the following folklore theorem (see for example [3]) bounding the Rademacher complexity of a function class composed with a fixed Lipschitz function.

**Theorem 9.** *Let  $\mathcal{F}$  be an  $\mathbb{R}^m$ -valued function class on a space  $\mathcal{X}$  and suppose that  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  has Lipschitz constant  $L$ . Let*

$$\phi \circ \mathcal{F} = \{ (\phi \circ f^1, \dots, \phi \circ f^m) : (f^1, \dots, f^m) \in \mathcal{F} \}.$$

Then

$$\hat{\mathcal{R}}_n^m(\phi \circ \mathcal{F}) \leq L \hat{\mathcal{R}}_n^m(\mathcal{F}).$$

Suppose now that  $\mathcal{F}$  is an  $\mathbb{R}^m$ -valued function class on  $\mathcal{X}$ . For  $\mathbf{f} = (f^1, \dots, f^m)$  define functions  $\mathbf{f}' = (f'^1, \dots, f'^m)$  and  $\mathbf{f}'' = (f''^1, \dots, f''^m)$ , from  $\mathcal{X} \times \{-1, 1\}$  to  $\mathbb{R}^m$  or  $[0, 1]^m$  respectively, by

$$f'^l(x, y) = y f^l(x) \text{ and } f''^l(x, y) = \phi_\gamma \circ f^l(x, y) = \phi_\gamma(y f^l(x))$$

and function classes  $\mathcal{F}' = \{\mathbf{f}' : \mathbf{f} \in \mathcal{F}\}$  and  $\mathcal{F}'' = \{\mathbf{f}'' : \mathbf{f} \in \mathcal{F}\}$ . It follows from the definition of  $\hat{\mathcal{R}}$  that  $\hat{\mathcal{R}}_n^m(\mathcal{F}')(\mathbf{x}, \mathbf{y}) = \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{x})$  for all  $(\mathbf{x}, \mathbf{y}) \in (\mathcal{X} \times \{-1, 1\})^{nm}$ . Since  $\phi_\gamma$  is Lipschitz with constant  $\gamma^{-1}$ , the previous theorem implies that

$$\hat{\mathcal{R}}_n^m(\mathcal{F}'')(\mathbf{X}, \mathbf{Y}) \leq \gamma^{-1} \hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) \text{ and } \mathcal{R}_n^m(\mathcal{F}'') \leq \gamma^{-1} \mathcal{R}_n^m(\mathcal{F}). \quad (19)$$

On the other hand, for every  $\mathbf{f} = (f^1, \dots, f^m) \in \mathcal{F}$  we have

$$\begin{aligned} \text{er}(\mathbf{h}_{\mathbf{f}}) &= \frac{1}{m} \sum E [1_{(-\infty, 0]}(Y_1^l f^l(X_1^l))] \\ &\leq \frac{1}{m} \sum E [\phi_\gamma \circ (f^l)^l(X_1^l, Y_1^l)] \\ &= \frac{1}{m} \sum E [(f'')^l(X_1^l, Y_1^l)] \end{aligned} \quad (20)$$

and

$$\frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n f''^l(X_i^l, Y_i^l) = \frac{1}{mn} \sum_{l=1}^m \sum_{i=1}^n \phi_\gamma(Y_i^l f^l(X_i^l)) = \text{er}_\gamma(\mathbf{f}). \quad (21)$$

Applying Theorem to the class  $\mathcal{F}''$  and substitution of (20), (21) and (19) yield

**Theorem 10.** Let  $\mathcal{F}$  be a  $\mathbb{R}^m$ -valued function class on a space  $\mathcal{X}$ ,  $\gamma \in (0, 1)$  and  $(\mathbf{X}, \mathbf{Y}) = (X_i^l, Y_i^l)_{(l,i)=(1,1)}^{(m,n)}$  a vector of  $\mathcal{X} \times \{-1, 1\}$ -valued independent random variables where for fixed  $l$  and varying  $i$  all the  $(X_i^l, Y_i^l)$  are identically distributed. Fix  $\delta > 0$ . Then with probability greater than  $1 - \delta$  we have for all  $\mathbf{f} \in \mathcal{F}$

$$er(\mathbf{h}_{\mathbf{f}}) \leq e\hat{r}_{\gamma}(\mathbf{f}) + \gamma^{-1}\mathcal{R}_n^m(\mathcal{F}) + \sqrt{\frac{\ln(1/\delta)}{2mn}}.$$

We also have with probability greater than  $1 - \delta$  for all  $\mathbf{f} \in \mathcal{F}$ , that

$$er(\mathbf{h}_{\mathbf{f}}) \leq e\hat{r}_{\gamma}(\mathbf{f}) + \gamma^{-1}\hat{\mathcal{R}}_n^m(\mathcal{F})(\mathbf{X}) + \sqrt{\frac{9\ln(2/\delta)}{2mn}}.$$

To arrive at Theorem 4 we still need to convert this into a statement valid with high probability for all margins  $\gamma \in (0, 1)$ . This is done following the techniques described in [1], using the following lemma (a copy of Lemma 15.5 from [1]):

**Lemma 5.** Suppose

$$\{F(\alpha_1, \alpha_2, \delta) : 0 < \alpha_1, \alpha_2, \delta \leq 1\}$$

is a set of events such that:

(i) For all  $0 < \alpha \leq 1$  and  $0 < \delta \leq 1$ ,

$$\Pr\{F(\alpha, \alpha, \delta)\} \leq \delta.$$

(ii) For all  $0 < \alpha_1 \leq \alpha \leq \alpha_2 \leq 1$  and  $0 < \delta_1 \leq \delta \leq 1$ ,

$$F(\alpha_1, \alpha_2, \delta_1) \subseteq F(\alpha, \alpha, \delta).$$

Then for  $0 < a, \delta < 1$ ,

$$\Pr\left(\bigcup_{\alpha \in (0,1]} F(\alpha a, \alpha, \delta \alpha (1-a))\right) \leq \delta.$$

Applications of this lemma follow a standard pattern. We give only one example, where we apply it to the event

$$F(\alpha_1, \alpha_2, \delta) = \left\{ \exists \mathbf{f} \in \mathcal{F} \text{ s.t. } er(\mathbf{h}_{\mathbf{f}}) > e\hat{r}_{\alpha_2}(\mathbf{f}) + \alpha_1^{-1}\mathcal{R}_n^m(\mathcal{F}) + \sqrt{\frac{\ln(1/\delta)}{2mn}} \right\}.$$

Condition (i) follows from the previous theorem, condition (ii) from the fact that the right side in the inequality increases if we decrease  $\delta$  and  $\alpha_1$  and increase  $\alpha_2$ . If we replace  $a$  by  $1 - \epsilon$  and  $\alpha$  by  $\gamma$ , then the conclusion of the lemma becomes the first conclusion of Theorem 4. The second conclusion of Theorem 4 and the application in Theorem 7 are handled similarly.

## References

1. M. Anthony and P. L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, Cambridge, UK, 1999.
2. P. L. Bartlett and S. Mendelson. Rademacher and Gaussian Complexities: Risk Bounds and Structural Results. *Journal of Machine Learning Research*, 2002.
3. P. Bartlett, O. Bousquet and S. Mendelson. Local Rademacher complexities. Available online: <http://www.stat.berkeley.edu/~bartlett/papers/bbm-lrc-02b.pdf>.
4. P. Baxendale. Gaussian measures on function spaces. *Amer. J. Math.*, 98:891-952, 1976.
5. J. Baxter, Theoretical Models of Learning to Learn, in *Learning to Learn*, S. Thrun, L. Pratt Eds. Springer 1998
6. J. Baxter, A Model of Inductive Bias Learning, *Journal of Artificial Intelligence Research* 12: 149-198, 2000
7. S. Ben-David and R. Schuller. Exploiting task relatedness for multiple task learning. In *COLT 03*, 2003.
8. R. Caruana, Multitask Learning, in *Learning to Learn*, S. Thrun, L. Pratt Eds. Springer 1998.
9. Nello Cristianini and John Shawe-Taylor, Support Vector Machines, *Cambridge University Press*, 2000.
10. T. Evgeniou and M. Pontil, Regularized multi-task learning. *Proc. Conference on Knowledge Discovery and Data Mining*, 2004.
11. V. Koltchinskii and D. Panchenko, Empirical margin distributions and bounding the generalization error of combined classifiers, *The Annals of Statistics*, Vol. 30, No 1, 1-50.
12. Colin McDiarmid, Concentration, in *Probabilistic Methods of Algorithmic Discrete Mathematics*, p. 195-248. Springer, Berlin, 1998.
13. C. A. Miccheli and M. Pontil, Kernels for multi-task learning. Available online, 2005.
14. S. Mika, B. Schölkopf, A. Smola, K.-R. Müller, M. Scholz and G. Rätsch. Kernel PCA and De-noising in Feature Spaces, in *Advances in Neural Information Processing Systems* 11, 1998.
15. J. Shawe-Taylor, N. Cristianini, Estimating the moments of a random vector, *Proceedings of GRETSI 2003 Conference*, I: 47-52, 2003.
16. Michael Reed and Barry Simon. *Functional Analysis*, part I of *Methods of Mathematical Physics*, Academic Press, 1980.
17. S. Thrun, Lifelong Learning Algorithms, in *Learning to Learn*, S. Thrun, L. Pratt Eds. Springer 1998
18. R. K. Ando, T. Zhang, A framework for learning predictive structures from multiple tasks and unlabeled data, Report RC23462, IBM T.J. Watson Research Center, 2004.