

# Efficient Output Kernel Learning for Multiple Tasks

**Authored by:**

Matthias Hein  
Bernt Schiele  
Maksim Lapin  
Pratik Jawanpuria

## **Abstract**

The paradigm of multi-task learning is that one can achieve better generalization by learning tasks jointly and thus exploiting the similarity between the tasks rather than learning them independently of each other. While previously the relationship between tasks had to be user-defined in the form of an output kernel, recent approaches jointly learn the tasks and the output kernel. As the output kernel is a positive semidefinite matrix, the resulting optimization problems are not scalable in the number of tasks as an eigendecomposition is required in each step. Using the theory of positive semidefinite kernels we show in this paper that for a certain class of regularizers on the output kernel, the constraint of being positive semidefinite can be dropped as it is automatically satisfied for the relaxed problem. This leads to an unconstrained dual problem which can be solved efficiently. Experiments on several multi-task and multi-class data sets illustrate the efficacy of our approach in terms of computational efficiency as well as generalization performance.

## **1 Paper Body**

Multi-task learning (MTL) advocates sharing relevant information among several related tasks during the training stage. The advantage of MTL over learning tasks independently has been shown theoretically as well as empirically [1, 2, 3, 4, 5, 6, 7]. The focus of this paper is the question how the task relationships can be inferred from the data. It has been noted that naively grouping all the tasks together may be detrimental [8, 9, 10, 11]. In particular, outlier tasks may lead to worse performance. Hence, clustered multi-task learning algorithms [10, 12] aim to learn groups of closely related tasks. The information is then shared only within these clusters of tasks. This corresponds to learning the task covariance matrix, which we denote as the output kernel in this paper.

Most of these approaches lead to non-convex problems. In this work, we focus on the problem of directly learning the output kernel in the multi-task learning framework. The multi-task kernel on input and output is assumed to be decoupled as the product of a scalar kernel and the output kernel, which is a positive semidefinite matrix [1, 13, 14, 15]. In classical multi-task learning algorithms [1, 16], the degree of relatedness between distinct tasks is set to a constant and is optimized as a hyperparameter. However, constant similarity between tasks is a strong assumption and is unlikely to hold in practice. Thus recent approaches have tackled the problem of directly learning the output kernel. [17] solves a multi-task formulation in the framework of vector-valued reproducing kernel Hilbert spaces involving squared loss where they penalize the Frobenius norm of the output kernel as a regularizer. They formulate an invex optimization problem that they solve optimally. In comparison, [18] recently proposed an efficient barrier method to optimize a generic convex output kernel learning formulation. On the other hand, [9] proposes a convex formulation to learn low rank output kernel matrix by enforcing a trace constraint. The above approaches [9, 17, 18] solve the resulting optimization problem via alternate minimization between task parameters and the output kernel. Each step of the alternate minimization requires an eigen

value decomposition of a matrix having as size the number of tasks and a problem corresponding to learning all tasks independently. In this paper we study a similar formulation as [17]. However, we allow arbitrary convex loss functions and employ general  $p$ -norms for  $p \in (1, 2]$  (including the Frobenius norm) as regularizer for the output kernel. Our problem is jointly convex over the task parameters and the output kernel. Small  $p$  leads to sparse output kernels which allows for an easier interpretation of the learned task relationships in the output kernel. Under certain conditions on  $p$  we show that one can drop the constraint that the output kernel should be positive definite as it is automatically satisfied for the unconstrained problem. This significantly simplifies the optimization and our result could also be of interest in other areas where one optimizes over the cone of positive definite matrices. The resulting unconstrained dual problem is amenable to efficient optimization methods such as stochastic dual coordinate ascent [19], which scale well to large data sets. Overall we do not require any eigenvalue decomposition operation at any stage of our algorithm and no alternate minimization is necessary, leading to a highly efficient methodology. Furthermore, we show that this trick not only applies to  $p$ -norms but also applies to a large class of regularizers for which we provide a characterization. Our contributions are as follows: (a) we propose a generic  $p$ -norm regularized output kernel matrix learning formulation, which can be extended to a large class of regularizers; (b) we show that the constraint on the output kernel to be positive definite can be dropped as it is automatically satisfied, leading to an unconstrained dual problem; (c) we propose an efficient stochastic dual coordinate ascent based method for solving the dual formulation; (d) we empirically demonstrate the superiority of our approach in terms of generalization performance as well as significant reduction in training time compared to other methods learning the output kernel. The paper is organized

as follows. We introduce our formulation in Section 2. Our main technical result is discussed in Section 3. The proposed optimization algorithm is described in Section 4. In Section 5, we report the empirical results. All the proofs can be found in the supplementary material.

2

## The Output Kernel Learning Formulation

We first introduce the setting considered in this paper. We denote the number of tasks by  $T$ . We assume that all tasks have a common input space  $X$  and a common positive definite kernel function  $k : X \times X \rightarrow \mathbb{R}$ . We denote by  $\phi(\cdot)$  the feature map and by  $H_k$  the reproducing kernel Hilbert space (RKHS) [20] associated with  $k$ . The training data is  $(x_i, y_i, t_i)_{i=1}^n$ , where  $x_i \in X$ ,  $t_i$  is the task the  $i$ -th instance belongs to and  $y_i$  is the corresponding label. Moreover, we have a positive  $T \times T$  definite matrix  $\Gamma \in \mathbb{S}^+$  on the set of tasks  $\{1, \dots, T\}$ , where  $\mathbb{S}^+$  is the set of  $T \times T$  symmetric and positive semidefinite (p.s.d.) matrices. If one arranges the predictions of all tasks in a vector one can see multi-task learning as learning a vector-valued function in a RKHS [see 1, 13, 14, 15, 18, and references therein]. However, in this paper we use the one-to-one correspondence between real-valued and matrix-valued kernels, see [21], in order to limit the technical overhead. In this framework we define the joint kernel of input space and the set of tasks  $M : (X \times \{1, \dots, T\}) \times (X \times \{1, \dots, T\}) \rightarrow \mathbb{R}$  as

$M(x, s), (z, t) = k(x, z)\Gamma(s, t)$ , (1) We denote the corresponding RKHS of functions on  $X \times \{1, \dots, T\}$  as  $H_M$  and by  $\|\cdot\|_{H_M}$  the corresponding norm. We formulate the output kernel learning problem for multiple tasks as  $\min_{F, \gamma} C_L(y_i, F(x_i, t_i)) + \lambda \|F\|_{H_M}^2 + \gamma V(\gamma)$  (2)  $T, F \in H^2 \times \mathbb{S}^+, M_{i=1}^n$  where  $L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is the convex loss function (convex in the second argument),  $V(\gamma)$  is a convex regularizer penalizing the complexity of the output kernel  $\gamma$  and  $\lambda \in \mathbb{R}^+$  is the regularization parameter. Note that  $\|F\|_{H_M}^2$  implicitly depends also on  $\gamma$ . In the following we show that (2) can be reformulated into a jointly convex problem in the parameters of the prediction function and the output kernel  $\gamma$ . Using the standard representer theorem [20] (see the supplementary material) for fixed output kernel  $\gamma$ , one can show that the optimal solution  $F \in H_M$  of (2) can be written as  $F(x, t) =$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is} \phi(x) \phi(s)$$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is} M(x_i, s), (x, t) = \sum_{i=1}^n \sum_{s=1}^T \alpha_{is} k(x_i, x) \Gamma(s, t).$$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is}$$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is}$$

$$2$$

$$(3)$$

With the explicit form of the prediction function one can rewrite the main problem (2) as  $\min_{\alpha, \gamma} C_L(y_i, \sum_{i=1}^n \sum_{s=1}^T \alpha_{is} \phi(x_i) \phi(s)) + \lambda \|\alpha\|_{H_M}^2 + \gamma V(\gamma)$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is} \Gamma(s, s)$$

$$C$$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is}$$

$$L(y_i, \sum_{i=1}^n \sum_{s=1}^T \alpha_{is} \phi(x_i) \phi(s))$$

$$\sum_{i=1}^n \sum_{s=1}^T \alpha_{is}$$

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \end{aligned} \quad (4)$$

where  $\lambda_{st} = \lambda(r, s)$  and  $k_{ij} = k(x_i, x_j)$ . Unfortunately, problem (4) is not jointly convex in  $\lambda$  and  $\beta$  due to the product in the second term. A similar problem has been analyzed in [17]. They could show that for the squared loss and  $V(\beta) = \|\beta\|_F^2$  the corresponding optimization problem is invex and directly optimize it. For an invex function every stationary point is globally optimal [22]. We follow a different path which leads to a formulation similar to the one of [2] used for learning an input mapping (see also [9]). Our formulation for the output kernel learning problem is jointly convex in the task kernel  $\beta$  and the task parameters. We present a derivation for the general RKHS  $H_k$ , analogous to the linear case presented in [2, 9]. We use the following variable transformation,  $\beta_{it} =$

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \text{resp.} \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \end{aligned}$$

In the last expression  $\beta_{it}$  has to be understood as the pseudo-inverse if  $\beta$  is not invertible. Note that this causes no problems as in case  $\beta$  is not invertible, we can without loss of generality restrict  $\beta$  in (4) to the range of  $\beta$ . The transformation leads to our final problem formulation, where the squared prediction function  $F$  and its squared norm  $\|\beta\|_F^2$  can be written as  $F(x, t) =$

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \end{aligned} \quad (5)$$

$$\begin{aligned} & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \\ & \sum_{i=1}^n \sum_{j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \left( \sum_{i,j=1}^n \sum_{s=1}^T \sum_{t=1}^T \lambda_{st} \right) \end{aligned}$$

$$\begin{aligned} & \mathbf{X}^T \mathbf{X} \\ (6) \end{aligned}$$

Before we analyze the convexity of this problem, we want to illustrate the connection to the formulations in [9, 17]. With the task weight vectors  $\mathbf{w}_t = \frac{1}{\sqrt{m}} \sum_{j=1}^m \mathbf{w}_j(x_j)$  we get predictions as  $F(\mathbf{x}, t) = \sum_{j=1}^m \mathbf{w}_j(x_j)$  and one can rewrite (2)

$$\begin{aligned} \mathbf{K} \mathbf{F} \mathbf{K}^T \mathbf{H} \mathbf{M} &= \\ \mathbf{T}^T \mathbf{X}^T \mathbf{X} \mathbf{T} & \\ \sum_{i,j=1}^n k(\mathbf{x}_i, \mathbf{x}_j) &= \sum_{r,s=1}^n \mathbf{w}_r^T \mathbf{w}_s \mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) \\ \mathbf{T}^T \mathbf{X} & \\ \sum_{i,j=1}^n & \\ \mathbf{w}_r^T & \\ \mathbf{H} \mathbf{M} &, \mathbf{w}_t \mathbf{I} . \\ \sum_{r,s=1}^n & \end{aligned}$$

This identity is known for vector-valued RKHS, see [15] and references therein. When  $\mathbf{T}$  is  $n \times n$  the identity matrix, then  $\mathbf{K} \mathbf{F} \mathbf{K}^T \mathbf{H} \mathbf{M} = \mathbf{T}^T \mathbf{X}^T \mathbf{X} \mathbf{T}$  and thus (2) is learning the tasks independently. As mentioned before the convexity of the expression of  $\mathbf{K} \mathbf{F} \mathbf{K}^T \mathbf{H} \mathbf{M}$  is crucial for the convexity of the full problem (6). The following result has been shown in [2] (see also [9]). **Lemma 1** Let  $R(\cdot)$  denote the range of  $\mathbf{T}^T \mathbf{X}^T \mathbf{X} \mathbf{T}$  and let  $\mathbf{T}^{\dagger}$  be the pseudoinverse. The extended  $\mathbf{T}^T \mathbf{X}^T \mathbf{X} \mathbf{T}$  function  $f: \mathbf{S}^+ \rightarrow \mathbf{R} \cup \{\infty\}$  defined as  $f(\mathbf{P}) = \mathbf{T}^T \mathbf{X}^T \mathbf{X} \mathbf{T} \mathbf{P}$  if  $\mathbf{P} \in R(\cdot)$ ,  $f(\mathbf{P}) = \infty$  else.

$\mathbf{P}$  is jointly convex. The formulation in (6) is similar to [9, 17, 18]. [9] uses the constraint  $\text{Trace}(\mathbf{V}) \leq 1$  instead of a regularizer  $V(\cdot)$  enforcing low rank of the output kernel. On the other hand, [17] employs squared Frobenius norm for  $V(\cdot)$  with squared loss function. [18] proposed an efficient algorithm for convex  $V(\cdot)$ . Instead we think that sparsity of  $\mathbf{P}$  is better to avoid the emergence of spurious relations between tasks and also leads to output kernels which are easier to interpret. Thus we propose to use the following regularization functional for the output kernel  $\mathbf{P}$ :  $V(\mathbf{P}) =$

$$\begin{aligned} & \mathbf{T}^T \mathbf{X}^T \mathbf{X} \mathbf{T} \mathbf{P} \\ & \sum_{t,t_0=1}^n \|\mathbf{w}_t - \mathbf{w}_{t_0}\|_p^p = \sum_{t,t_0=1}^n \|\mathbf{w}_t - \mathbf{w}_{t_0}\|_p^p , \\ & \sum_{t,t_0=1}^n \end{aligned}$$

for  $p \in [1, 2]$ . Several approaches [9, 17, 18] employ alternate minimization scheme, involving  $T$  costly eigendecompositions of  $\mathbf{T}^T \mathbf{T}$  matrix per iteration (as  $\mathbf{T}^T \mathbf{T} \in \mathbf{S}^+$ ). In the next section we show that for a certain set of values of  $p$  one can derive an unconstrained dual optimization problem which thus avoids the explicit minimization over the  $\mathbf{S}^+$  cone. The resulting unconstrained dual problem can then be easily optimized by stochastic coordinate ascent. Having explicit expressions of the primal variables  $\mathbf{P}$  and  $\mathbf{F}$  in terms of the dual variables allows us to get back to the original problem.

### Unconstrained Dual Problem Avoiding Optimization over $S+T$

The primal formulation (6) is a convex multi-task output kernel learning problem. The next lemma derives the Fenchel dual function of (6). This still involves the optimization over the primal variable  $T \in S+$ . A main contribution of this paper is to show that this optimization problem over the  $T \in S+$  cone can be solved with an analytical solution for a certain class of regularizers  $V(\cdot)$ . In the following we denote by  $\gamma_r := \{\gamma_i \mid t_i = r\}$  the dual variables corresponding to task  $r$  and by  $K_{rs}$  the kernel matrix  $(k(x_i, x_j) \mid t_i = r, t_j = s)$  corresponding to the dual variables of tasks  $r$  and  $s$ . Lemma 2 Let  $L_i$  be the conjugate function of the loss  $L_i : \mathbb{R} \rightarrow \mathbb{R}, u \mapsto L(y_i, u)$ , then  $q : \mathbb{R}^n \rightarrow \mathbb{R}, q(\gamma) = C$

$$\begin{aligned} & \max_{T \in S+} \sum_{i=1}^n X_i L_i(\gamma_{rs} h_r, K_{rs} \gamma_s) \quad (7) \end{aligned}$$

is the dual function of (6), where  $\gamma \in \mathbb{R}^n$  are the dual variables. The primal variable  $T \in \mathbb{R}^n \times T$  in (6) and  $P$  the prediction function  $F$  can be expressed in terms of  $\gamma$  and  $\gamma$  as  $\gamma_i s = \gamma_i \gamma_{st_i}$  and  $n F(x, s) = \sum_{j=1}^n \gamma_j \gamma_{st_j} k(x_j, x)$  respectively, where  $t_j$  is the task of the  $j$ -th training example. We now focus on the remaining maximization problem in the dual function in (7)

$$\begin{aligned} & \max_{T \in S+} \sum_{i=1}^n X_i \gamma_{rs} h_r, K_{rs} \gamma_s \in V(\gamma). \quad (8) \end{aligned}$$

This is a semidefinite program which is computationally expensive to solve and thus prohibits to scale the output kernel learning problem to a large number of tasks. However, we show in the following that this problem has an analytical solution for a subset of the regularizers  $V(\gamma) = \sum_{r,s=1}^p \gamma_{rs} \gamma_{rs}$  for  $p \geq 1$ . For better readability we defer a more general result towards the end of the section. The basic idea is to relax the constraint on  $T \in S+T$  in (8) so that it is equivalent to the computation of the conjugate  $V^*$  of  $V$ . If the maximizer of the relaxed problem is positive semi-definite, one has found the solution of the original problem. Theorem 3 Let  $k \in \mathbb{N}$  and  $p = \max$

$$\begin{aligned} & \max_{T \in S+} \sum_{i=1}^n X_i \gamma_{rs} \gamma_{rs} \\ & \text{then with } \gamma_{rs} = \frac{1}{2} \sum_{r,s=1}^p \gamma_{rs} \gamma_{rs} \\ & \text{we have} \\ & \sum_{r,s=1}^p \gamma_{rs} \gamma_{rs} \geq \frac{1}{2} \sum_{r,s=1}^p \gamma_{rs} \gamma_{rs} \quad (9) \end{aligned}$$

(10)

Plugging the result of the previous theorem into the dual function of Lemma 2 we get for  $k \in \mathbb{N}$  and  $p = 2k+1$  with  $V(\lambda) = k!k^p$  the following unconstrained dual of our main problem (6):  $\max_{\lambda} \Phi_C$

$\lambda \in \mathbb{R}$

$\lambda \in \mathbb{X}$

$T = 2k+1$   $2k+1$   $X = 2k+1$   $L(\lambda) = \sum_{i=1}^k h_i$ ,  $K_{rs} = \sum_{i=1}^k C_{4k+2-2k} r, s=1$   $i=1$

(11)

Note that by doing the variable transformation  $\lambda_i := \Phi_C$  we effectively have only one hyperparameter in (11). This allows us to cross-validate more efficiently. The range of admissible values for  $p$  in Theorem 3 lies in the interval  $(1, 2]$ , where we get for  $k = 1$  the value  $p = 2$  and as  $k \rightarrow \infty$

Table 1: Examples of regularizers  $V(\lambda)$  together with their generating function  $\Phi$  and the explicit form of  $\Phi$  in terms of the dual variables,  $\lambda_{rs} = 2\lambda h_{rs}$ ,  $K_{rs} = \sum_{i=1}^k C_{4k+2-2k} r, s=1$   $i=1$   $\Phi = V(\lambda) = \sum_{r,s=1}^k \lambda(\lambda_{rs})$ .  $\Phi = T$

$\Phi(z)$

$\lambda_{rs}$

$V(\lambda)$

$z = 2k+1$ ,

$k \in \mathbb{N}$

$e^z =$

$P =$

$2k+1$   $2k$

$k$

$z = 0$   $k!$

$\cosh(z) = 1 =$

$z = 2k+1$   $(2k)!$

$P =$

$T = P$

$2k$

$2k+1$   $rs$

$-\lambda_{rs} - 2k+1$

$r, s=1$

$\Phi = T = P = \log(\lambda) = \lambda$  if  $\lambda_{rs} \geq 0$ ,  $r, s = 1, \dots, k$  else  $\lambda = 0$ .

$T = P = \lambda_{rs} \operatorname{arcsinh}(\lambda_{rs}) = 1 + \lambda_{rs} + T = 2$

$e^{\lambda_{rs} \operatorname{arcsinh}(\lambda_{rs})}$

$r, s=1$

we have  $p \geq 1$ . The regularizer for  $p = 2$  together with the squared loss has been considered in the primal in [17, 18]. Our analytical expression of the dual is novel and allows us to employ stochastic dual coordinate ascent to solve the involved primal optimization problem. Please also note that by optimizing the dual, we have access to the duality gap and thus a well-defined stopping criterion. This is in contrast to the alternating scheme of [17, 18] for the primal problem which involves costly matrix operations. Our runtime experiments show that our solver for (11) outperforms the solvers of [17, 18].

Finally, note that even for suboptimal dual variables  $\lambda$ , the corresponding  $\lambda$  matrix in (10) is positive semidefinite. Thus we always get a feasible set of primal variables. Characterizing the set of convex regularizers  $V$  which allow an analytic expression for the dual function The previous theorem raises the question for which class of convex, separable regularizers we can get an analytical expression of the dual function by explicitly solving the optimization problem (8) over the positive semidefinite cone. A key element in the proof of the previous theorem is the characterization of functions  $f : \mathbb{R}^T \rightarrow \mathbb{R}$  which when applied elementwise  $T f(A) = (f(a_{ij}))_{i,j=1}^T$  to a positive semidefinite matrix  $A \in \mathbb{S}^+_T$  result in a p.s.d. matrix, that is  $T f(A) \in \mathbb{S}^+_T$ . This set of functions has been characterized by Hiai [23].  $T$ . We denote by  $f(A) = (f(a_{ij}))_{i,j=1}^T$  the elementwise application of  $f$  to  $A$ . Theorem 4 ([23]) Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  and  $A \in \mathbb{S}^+_T$  if and only if  $f$  is analytic  $\Rightarrow f(A) \in \mathbb{S}^+_T$  wise application of  $f$  to  $A$ . It holds  $\forall T \geq 2, A \in \mathbb{S}^+_T, P \in \mathbb{R}^T, k \geq 0$  and  $f(x) = \sum_{k=0}^{\infty} a_k x^k$  with  $a_k \geq 0$  for all  $k \geq 0$ .

Note that in the previous theorem the condition on  $f$  is only necessary when we require the implication to hold for all  $T$ . If  $T$  is fixed, the set of functions is larger and includes even (large) fractional powers, see [24]. We use the stronger formulation as we want that the result holds without any restriction on the number of tasks  $T$ . Theorem 4 is the key element used in our following characterization of separable regularizers of  $\lambda$  which allow an analytical expression of the dual function. Theorem 5 Let  $\lambda : \mathbb{R} \rightarrow \mathbb{R}$  be analytic on  $\mathbb{R}$  and given as  $\lambda(z) = \sum_{k=0}^{\infty} a_{k+1} z^k$  where  $a_k \in \mathbb{R}, \forall k \geq 0$ . If  $\lambda$  is convex, then,  $V(\lambda) := \sum_{r,s=1}^T \lambda(a_{rs})$ , is a convex function  $V : \mathbb{R}^T \rightarrow \mathbb{R}$  and  $\max_{\lambda \in \mathbb{R}^T} h(\lambda, \lambda) = V(\lambda) = V^*(\lambda) =$

$$\sum_{r=1}^T \lambda(a_{rr})$$

$$T \times$$

$$\sum_{r,s=1}^T a_{rs},$$

$$(12)$$

$$r,s=1$$

$$T \times T \text{ where the global maximizer fulfills } \lambda \in \mathbb{S}^+_T \text{ if } \lambda \in \mathbb{S}^+_T \text{ and } \sum_{r,s=1}^T a_{rs} =$$

$$P \in \mathbb{R}^T$$

$$k=0$$

$$a_k \geq 0 \text{ for all } k \geq 0.$$

Table 1 summarizes e.g. of functions  $\lambda$ , the corresponding  $V(\lambda)$  and the maximizer  $\lambda^*$  in (12).

4

Optimization Algorithm

The dual problem (11) can be efficiently solved via decomposition based methods like stochastic dual coordinate ascent algorithm (SDCA) [19]. SDCA enjoys low computational complexity per iteration and has been shown to scale effortlessly to large scale optimization problems. 5

Algorithm 1 Fast MTL-SDCA Input: Gram matrix  $K$ , label vector  $y$ , regularization parameter and relative duality gap parameter Output:  $\lambda^*$  ( $\lambda^*$  is computed from  $\lambda^*$  using our result in 10) Initialize  $\lambda = \lambda(0)$  repeat Randomly choose a dual variable  $\lambda_i$  Solve for  $\lambda$  in (13) corresponding to  $\lambda_i$   $\lambda_i \leftarrow \lambda_i + \eta$  until Relative duality gap is below



Our algorithm for learning the output kernel matrix and task parameters is summarized in Algorithm 1 (refer to the supplementary material for more details). At each step of the iteration we optimize the dual objective over a randomly chosen  $\tau_i$  variable. Let  $t_i = r$  be the task corresponding to  $\tau_i$ . We apply the update  $\tau_i \leftarrow \tau_i + \eta$ . The optimization problem of solving (11) with respect to  $\eta$  is as follows:

$$\min_{\eta} L(\tau_i + \eta)/C + \eta (a\tau_i^2 + 2b\tau_i + c) + 2(b\tau_i + c)\tau_i + c\tau_i^2, \quad (13)$$

$$s_6 = r$$

$$s, z_6 = r$$

where  $a = k_{ii}$ ,  $b = \sum_j k_{ij} \tau_j$ ,  $c = h^T s$ ,  $Ks = z$  and  $\eta = C(4k\tau_i^2 + 2k\tau_i + 2k)$ . This one-dimensional convex optimization problem is solved efficiently via Newton method. The complexity of the proposed algorithm is  $O(T)$  per iteration. The proposed algorithm can also be employed for learning output kernels regularized by generic  $V(\cdot)$ , discussed in the previous section. Special case  $p = 2$  ( $k = 1$ ): For certain loss functions such as the hinge loss, the squared loss, etc.,  $L(\tau_i + \eta)$  yields a linear or a quadratic expression in  $\eta$ . In such cases problem (13) reduces to finding the roots of a cubic equation, which has a closed form expression. Hence, our algorithm is highly efficient with the above loss functions when  $\tau$  is regularized by the squared Frobenius norm.

5

### Empirical Results

In this section, we present our results on benchmark data sets comparing our algorithm with existing approaches in terms of generalization accuracy as well as computational efficiency. Please refer to the supplementary material for additional results and details.

#### Multi-Task Data Sets

We begin with the generalization results in multi-task setups. The data sets are as follows: a) Sarcos: a regression data set, aim is to predict 7 degrees of freedom of a robotic arm, b) Parkinson: a regression data set, aim is to predict the Parkinson's disease symptom score for 42 patients, c) Yale: a face recognition data with 28 binary classification tasks, d) Landmine: a data set containing binary classifications from 19 different landmines, e) MHC-I: a bioinformatics data set having 10 binary classification tasks, f) Letter: a handwritten letters data set with 9 binary classification tasks. We compare the following algorithms: Single task learning (STL), multi-task methods learning the output kernel matrix (MTL [16], CMTL [12], MTRL [9]) and approaches that learn both input and output kernel matrices (MTFL [11], GMTL [10]). Our proposed formulation (11) is denoted by FMTLp. We consider three different values for the  $p$ -norm:  $p = 2$  ( $k = 1$ ),  $p = 4/3$  ( $k = 2$ ) and  $p = 8/7$  ( $k = 4$ ). Hinge and -SVR loss functions were employed for classification and regression problems respectively. We follow the experimental protocol described in [11]. Table 2 reports the performance of the algorithms averaged over ten random train-test splits. The proposed FMTLp attains the best generalization accuracy in general. It outperforms the baseline MTL as well as MTRL and CMTL, which solely learns the output kernel matrix. Moreover, it achieves an overall better

performance than GMTL and MTFL. The FMTL<sub>p=4/3,8/7</sub> give comparable generalization to  $p = 2$  case, with the additional benefit of learning sparser and more interpretable output kernel matrix (see Figure 1). 1

The performance of STL, MTL, CMTL and MTFL are reported from [11]. 6

Table 2: Mean generalization performance and the standard deviation over ten train-test splits. Data set

STL
MTL
CMTL
MTFL
GMTL
MTRL
Regression data sets: Explained Variance (%) Sarcos 40.5?7.6 34.5?10.2 33.0?13.4 49.9?6.3 45.8?10.6 41.6?7.1 2.7?3.6 16.8?10.8 33.6?9.4 12.0?6.8 Parkinson 2.8?7.5 4.9?20.0 Classification data sets: AUC (%) Yale 93.4?2.3 96.4?1.6 Landmine 74.6?1.6 76.4?0.8 MHC-I 69.3?2.1 72.3?1.9 Letter 61.2?0.8 61.0?1.6 95.2?2.1 75.9?0.7 72.6?1.4 60.5?1.1 97.0?1.6 76.4?1.0 71.7?2.2 60.5?1.8 91.9?3.2 76.7?1.2 72.5?2.7 61.2?0.9 46.7?6.9 27.0?4.4
2
2
4
4
4
6
6
6
8
8
8
10
10
10
12
12
12
14
14
14
16
16
16
18
18 4

6  
8  
10  
12  
14  
16  
18  
 $p = 8/7$   
50.3?5.8 48.4?5.8 27.0?4.4 27.0?4.4  
96.1?2.1 97.0?1.2 97.0?1.4 76.1?1.0 76.8?0.8 76.7?1.0 71.5?1.7 71.7?1.9 70.8?2.1  
60.3?1.4 61.4?0.7 61.5?1.0  
2  
2  
FMTLP  $p = 4/3$   
 $p=2$   
96.8?1.4 76.4?0.9 70.7?1.9 61.4?1.0  
18 2  
4  
( $p = 2$ )  
6  
8  
10  
12  
14  
16  
18  
2  
4  
6  
( $p = 4/3$ )  
8  
10  
12  
14  
16  
18  
( $p = 8/7$ )

Figure 1: Plots of  $\text{---}\text{?}\text{---}$  matrices (rescaled to  $[0,1]$  and averaged over ten splits) computed by our solver FMTLP for the Landmine data set for different  $p$ -norms, with cross-validated hyper-parameter values. The darker regions indicate higher value. Tasks (landmines) numbered 1-10 correspond to highly foliated regions and those numbered 11-19 correspond to bare earth or desert regions. Hence, we expect two groups of tasks (indicated by the red squares). We can observe that the learned  $\text{?}$  matrix at  $p = 2$  depicts much more spurious task relationships than the ones at  $p = 4/3$  and  $p = 8/7$ . Thus, our sparsifying regularizer improves interpretability.

Table 3: Mean accuracy and the standard deviation over five train-test splits.

Data set	
STL	
MTL-SDCA	
MNIST USPS	
	84.1 $\pm$ 0.3 90.5 $\pm$ 0.3
	86.0 $\pm$ 0.2 90.6 $\pm$ 0.2
5.2	
GMTL	
MTRL	
p=2	
FMTLp -H p = 4/3	
p = 8/7	
p=2	
FMTLp -S p = 4/3	
p = 8/7	
	84.8 $\pm$ 0.3 85.6 $\pm$ 0.4 86.1 $\pm$ 0.4 85.8 $\pm$ 0.4 86.2 $\pm$ 0.4 82.2 $\pm$ 0.6 82.5 $\pm$ 0.4 82.4 $\pm$ 0.3 91.6 $\pm$ 0.3
	92.4 $\pm$ 0.2 92.4 $\pm$ 0.2 92.6 $\pm$ 0.2 92.6 $\pm$ 0.1 87.2 $\pm$ 0.4 87.7 $\pm$ 0.3 87.5 $\pm$ 0.3

#### Multi-Class Data Sets

The multi-class setup is cast as T one-vs-all binary classification tasks, corresponding to T classes. In this section we experimented with two loss functions: a) FMTLp -H ? the hinge loss employed in SVMs, and b) FMTLp -S ? the squared loss employed in OKL [17]. In these experiments, we also compare our results with MTL-SDCA, a state-of-the-art multi-task feature learning method [25]. USPS & MNIST Experiments: We followed the experimental protocol detailed in [10]. Results are tabulated in Table 3. Our approach FMTLp -H obtains better accuracy than GMTL, MTRL and MTL-SDCA [25] on both data sets. MIT Indoor67 Experiments: We report results on the MIT Indoor67 benchmark [26] which covers 67 indoor scene categories. We use the train/test split (80/20 images per class) provided by the authors. FMTLp -S achieved the accuracy of 73.3% with p = 8/7. Note that this is better than the ones reported in [27] (70.1%) and [26] (68.24%). SUN397 Experiments: SUN397 [28] is a challenging scene classification benchmark [26] with 397 classes. We use m = 5, 50 images per class for training, 50 images per class for testing and report the average accuracy over the 10 standard splits. We employed the CNN features extracted with the 7

Table 4: Mean accuracy and the standard deviation over ten train-test splits on SUN397. m

STL	
MTL	
MTL-SDCA	
5 50	
	40.5 $\pm$ 0.9 55.0 $\pm$ 0.4
	42.0 $\pm$ 1.4 57.0 $\pm$ 0.2
	41.2 $\pm$ 1.3 54.8 $\pm$ 0.3
p=2	

FMTLp -H  $p = 4/3$   
 $p = 8/7$   
 $p=2$   
 FMTLp -S  $p = 4/3$   
 $p = 8/7$   
 41.5?1.1 55.1?0.2  
 41.6?1.3 55.6?0.3  
 41.6?1.2 55.1?0.3  
 44.1?1.3 58.6?0.1  
 44.1?1.1 58.5?0.1  
 44.0?1.2 58.6?0.2  
 (Time by baseline) / (Time by FMTL2?S)  
 3  
 10  
 FMTL ?S 2  
 2  
 ConvexOKL  
 Time (log10 scale), s  
 10  
 OKL 1  
 10  
 0  
 10  
 ?1  
 10  
 ?2  
 10  
 50  
 100  
 150  
 200 250 300 Number of Tasks  
 350  
 400  
 (a)  
 20 18 16  
 MIT Indoor67, OKL SUN397, OKL MIT Indoor67, ConvexOKL . SUN397,  
 ConvexOKL  
 14 12 10 8 6 4 2 0 3  
 3.5  
 4  
 4.5  
 5 5.5 Log10(?)  
 6  
 6.5  
 7  
 (b)

Figure 2: (a) Plot compares the runtime of various algorithms with varying number of tasks on SUN397. Our approach FMTL2 -S is 7 times faster than OKL [17] and 4.3 times faster than ConvexOKL [18] when the number of tasks is maximum. (b) Plot showing the factor by which FMTL2 S outperforms OKL and ConvexOKL over the hyper-parameter range on various data sets. On SUN397, we outperform OKL and ConvexOKL by factors of 5.2 and 7 respectively. On MIT Indoor67, we are better than OKL and ConvexOKL by factors of 8.4 and 2.4 respectively. convolutional neural network (CNN) [26] using Places 205 database. The results are tabulated in Table 4. The ? matrices computed by FMTLp -S are discussed in the supplementary material. 5.3

#### Scaling Experiment

We compare the runtime of our solver for FMTL2 -S with the OKL solver of [17] and the ConvexOKL solver of [18] on several data sets. All the three methods solve the same optimization problem. Figure 2a shows the result of the scaling experiment where we vary the number of tasks (classes). The parameters employed are the ones obtained via cross-validation. Note that both OKL and ConvexOKL algorithms do not have a well defined stopping criterion whereas our approach can easily compute the relative duality gap (set as  $10^{-3}$ ). We terminate them when they reach the primal objective value achieved by FMTL2 -S. Our optimization approach is 7 times and 4.3 times faster than the alternate minimization based OKL and ConvexOKL, respectively, when the number of tasks is maximal. The generic FMTLp=4/3,8/7 are also considerably faster than OKL and ConvexOKL. Figure 2b compares the average runtime of our FMTLp -S with OKL and ConvexOKL on the crossvalidated range of hyper-parameter values. FMTLp -S outperform them on both MIT Indoor67 and SUN397 data sets. On MNIST and USPS data sets, FMTLp -S is more than 25 times faster than OKL, and more than 6 times faster than ConvexOKL. Additional details of the above experiments are discussed in the supplementary material.

#### 6

#### Conclusion

We proposed a novel formulation for learning the positive semi-definite output kernel matrix for multiple tasks. Our main technical contribution is our analysis of a certain class of regularizers on the output kernel matrix where one may drop the positive semi-definite constraint from the optimization problem, but still solve the problem optimally. This leads to a dual formulation that can be efficiently solved using stochastic dual coordinate ascent algorithm. Results on benchmark multi-task and multi-class data sets demonstrates the effectiveness of the proposed multi-task algorithm in terms of runtime as well as generalization accuracy. Acknowledgments. P.J. and M.H. acknowledge the support by the Cluster of Excellence (MMCI). 8

## 2 References

- [1] T. Evgeniou, C. A. Micchelli, and M. Pontil. Learning multiple tasks with kernel methods. *JMLR*, 6:615?637, 2005. [2] A. Argyriou, T. Evgeniou, and

M. Pontil. Convex multi-task feature learning. *ML*, 73:243?272, 2008. [3] K. Lounici, M. Pontil, A. B. Tsybakov, and S. van de Geer. Taking advantage of sparsity in multi-task learning. In *COLT*, 2009. [4] A. Jalali, P. Ravikumar, S. Sanghavi, and C. Ruan. A dirty model for multi-task learning. In *NIPS*, 2010. [5] P. Jawanpuria and J. S. Nath. Multi-task multiple kernel learning. In *SDM*, 2011. [6] A. Maurer, M. Pontil, and B. Romera-paredes. Sparse coding for multitask and transfer learning. In *ICML*, 2013. [7] P. Jawanpuria, J. S. Nath, and G. Ramakrishnan. Generalized hierarchical kernel learning. *JMLR*, 16:617? 652, 2015. [8] R. Caruana. Multitask learning. *ML*, 28:41?75, 1997. [9] Y. Zhang and D. Y. Yeung. A convex formulation for learning task relationships in multi-task learning. In *UAI*, 2010. [10] Z. Kang, K. Grauman, and F. Sha. Learning with whom to share in multi-task feature learning. In *ICML*, 2011. [11] P. Jawanpuria and J. S. Nath. A convex feature learning formulation for latent task structure discovery. In *ICML*, 2012. [12] L. Jacob, F. Bach, and J. P. Vert. Clustered multi-task learning: A convex formulation. In *NIPS*, 2008. [13] C. A. Micchelli and M. Pontil. Kernels for multitask learning. In *NIPS*, 2005. [14] A. Caponnetto, C. A. Micchelli, M. Pontil, and Y. Ying. Universal multi-task kernels. *JMLR*, 9:1615? 1646, 2008. ? [15] M. A. Alvarez, L. Rosasco, and N. D. Lawrence. Kernels for vector-valued functions: a review. *Foundations and Trends in Machine Learning*, 4:195?266, 2012. [16] T. Evgeniou and M. Pontil. Regularized multi-task learning. In *KDD*, 2004. [17] F. Dinuzzo, C. S. Ong, P. Gehler, and G. Pillonetto. Learning output kernels with block coordinate descent. In *ICML*, 2011. [18] C. Ciliberto, Y. Mroueh, T. Poggio, and L. Rosasco. Convex learning of multiple tasks and their structure. In *ICML*, 2015. [19] S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss. *JMLR*, 14(1):567?599, 2013. [20] B. Schölkopf and A. Smola. *Learning with Kernels*. MIT Press, 2002. [21] M. Hein and O. Bousquet. Kernels, associated structures and generalizations. Technical Report TR-127, Max Planck Institute for Biological Cybernetics, 2004. [22] A. Ben-Israel and B. Mond. What is invexity ? *J. Austral. Math. Soc. Ser. B*, 28:1?9, 1986. [23] F. Hiai. Monotonicity for entrywise functions of matrices. 431(8):1125 ? 1146, 2009.

Linear Algebra and its Applications,

[24] R. A. Horn. The theory of infinitely divisible matrices and kernels. *Trans. Amer. Math. Soc.*, 136:269?286, 1969. [25] M. Lapin, B. Schiele, and M. Hein. Scalable multitask representation learning for scene classification. In *CVPR*, 2014. [26] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *NIPS*, 2014. [27] M. Koskela and J. Laaksonen. Convolutional network features for scene recognition. In *Proceedings of the ACM International Conference on Multimedia*, 2014. [28] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.