

Multitask learning meets tensor factorization: task imputation via convex optimization

Authored by:

Masashi Sugiyama
Ryota Tomioka
Kishan Wimalawarne

Abstract

We study a multitask learning problem in which each task is parametrized by a weight vector and indexed by a pair of indices, which can be e.g. (consumer, time). The weight vectors can be collected into a tensor and the (multilinear-)rank of the tensor controls the amount of sharing of information among tasks. Two types of convex relaxations have recently been proposed for the tensor multilinear rank. However, we argue that both of them are not optimal in the context of multitask learning in which the dimensions or multilinear rank are typically heterogeneous. We propose a new norm, which we call the scaled latent trace norm and analyze the excess risk of all the three norms. The results apply to various settings including matrix and tensor completion, multitask learning, and multilinear multitask learning. Both the theory and experiments support the advantage of the new norm when the tensor is not equal-sized and we do not a priori know which mode is low rank.

1 Paper Body

We consider supervised multitask learning problems [1, 6, 7] in which the tasks are indexed by a pair of indices known as multilinear multitask learning (MLMTL) [17, 19]. For example, when we would like to predict the ratings of different aspects (e.g., quality of service, food, etc) of restaurants by different customers, the tasks would be indexed by aspects \times customers. When each task is parametrized by a weight vector over features, the goal would be to learn a features \times aspects \times customers tensor. Another possible task dimension would be time, since the ratings may change over time. This setting is interesting, because it would allow us to exploit the similarities across different customers as well as similarities across different aspects or time-points. Furthermore this would allow us to perform task imputation, that is to learn weights for tasks for which we have no training examples. On the other hand, the conventional matrix-based multitask learning (MTL) [2, 3, 13, 16] may fail to capture the

higher order structure if we consider learning a flat features \times tasks matrix and would require at least r samples, where r is the rank of the matrix to be learned, for each task. Recently several norms that induce low-rank tensors in the sense of Tucker decomposition or multilinear singular value decomposition [8, 9, 14, 25] have been proposed. The mean squared error for recovering a $n_1 \times \dots \times n_K$ tensor of multilinear rank (r_1, \dots, r_K) from its noisy version scale as $\sum_{k=1}^K \frac{1}{n_k} \text{rk}_k^2$ for the overlapped trace norm [23]. On the other hand, the $k=1$ error of the latent trace norm scales as $O(\min_k \text{rk}_k / \min_k n_k)$ in the same setting [21]. Thus while the latent trace norm has the better dependence in terms of the multilinear rank rk_k , it has the worse dependence in terms of the dimensions n_k . Tensors that arise in multitask learning typically have heterogeneous dimensions. For example, the number of aspects for a restaurant (quality of service, food, atmosphere, etc.) would be much 1

Table 1: Tensor denoising performance using different norms. The mean squared error $\|W - \hat{W}\|_F^2 / N$ is shown for the denoising algorithms (3) using different norms for tensors. Op

Overlapped trace norm	$\sum_{k=1}^K \text{rk}_k^2 / \sum_{k=1}^K n_k$
Latent trace norm	$\text{Op} \min_k \text{rk}_k / \min_k n_k$
Scaled latent trace norm	$\text{Op} \min_k (\text{rk}_k / n_k)$

smaller than the number of customers or the number of features. In addition, it is a priori unclear which mode (or dimension) would have the most redundancy or sharing that could be exploited by multitask learning. Some of the modes may have full ranks if there is no sharing of information along them. Therefore, both the latent trace norm and the overlapped trace norm would suffer either from the heterogeneous multilinear rank or the heterogeneous dimensions in this context. In this paper, we propose a modification to the latent trace norm whose mean squared error scales as $O(\min_k (\text{rk}_k / n_k))$ in the same setting, which is better than both the previously proposed extensions of trace norm for tensors. We study the excess risk of the three norms through their Rademacher complexities in various settings including matrix completion, multitask learning, and MLMTL. The new analysis allows us to also study the tensor completion setting, which was only empirically studied in [22, 23]. Our analysis consistently shows the advantage of the proposed scaled latent trace norm in various settings in which the dimensions or ranks are heterogeneous. Experiments on both synthetic and real data sets are also consistent with our theoretical findings.

2

Norms for tensors and their denoising performance

Let $W \in \mathbb{R}^{n_1 \times \dots \times n_K}$ be a K -way tensor. We denote the total number of entries by $N := \prod_{k=1}^K n_k$. A mode- k fiber of W is an n_k dimensional vector we obtain by fixing all but the k th index. The mode- k unfolding $W^{(k)}$ of W is the $n_k \times N/n_k$ matrix formed by concatenating all the N/n_k mode- k fibers along

columns. We say that W has multilinear rank (r_1, \dots, r_K) if $r_k = \text{rank}(W^{(k)})$. 2.1

Existing norms for tensors

First we review two norms proposed in literature in order to convexify tensor decomposition. The overlapped trace norm (see [12, 15, 18, 22]) is defined as the sum of the trace norms of the mode- k unfoldings as follows: $\|W\|_{\text{overlap}} = \sum_{k=1}^K \text{tr}(W^{(k)})$

where tr is the trace norm (also known as the nuclear norm) [10, 20], which is defined as the absolute sum of singular values. Romera-Paredes et al. [17] has used the overlapped trace norm in MLMTL. The latent trace norm [21, 22] is defined as the infimum over K tensors as follows: $\|W\|_{\text{latent}} = \inf_{W^{(1)} + \dots + W^{(K)} = W} \sum_{k=1}^K \text{tr}(W^{(k)})$

(2)

Table 1 summarizes the denoising performance in mean squared error analyzed in Tomioka and Suzuki [21] for the above two norms. The setting is as follows: we observe a noisy version Y of a tensor W with multilinear rank (r_1, \dots, r_K) and would like to recover W by solving $\min_W \|Y - W\|_F^2 + \lambda \|W\|$, (3) where $\|W\|$ is either the overlapped trace norm or the latent trace norm. We can see that while the latent trace norm has the better dependence in terms of the multilinear rank, it has the worse dependence in terms of the dimensions. Intuitively, the latent trace norm recognizes the mode with the lowest rank. However, it does not have a good control of the dimensions; in fact, the factor 2

$1/\min_k n_k$ comes from the fact that for a random tensor X with i.i.d. Gaussian entries, the expectation of the dual norm $\|X\|_{\text{latent}} = \max_k \text{tr}(X^{(k)})$ behaves like $\text{Op}(\max_k N/n_k)$, where Op is the operator norm. 2.2

A new norm

In order to correct the unfavorable behavior of the dual norm, we propose the scaled latent trace norm. It is defined similarly to the latent trace norm with weights $1/n_k$ as follows: $\|W\|_{\text{scaled}} = \sum_{k=1}^K \text{tr}(W^{(k)})/n_k$

Now the expectation of the dual norm $\|X\|_{\text{scaled}} = \max_k \text{tr}(X^{(k)})/n_k$ behaves like $\text{Op}(N)$ for X with random i.i.d. Gaussian entries and combined with the following relation $\|W\|_{\text{scaled}} \leq \min_k n_k \|W\|_{\text{latent}}$, (5) we obtain the scaling of the mean squared error in the last column of Table 1. We can see that the scaled latent norm recognizes the mode with the lowest rank relative to its dimension. $\|W\|_{\text{scaled}} =$

3

inf

$\sum_{k=1}^K \text{tr}(W^{(k)})/n_k$

Theory for multilinear multitask learning

We consider $T = P Q$ supervised learning tasks. Training samples $(x_{ipq}, y_{ipq})_{i=1}^S$ are provided for a relatively small fraction of the task index pairs $S \subseteq [P] \times [Q]$. Each task is parametrized by a weight vector $w_{pq} \in \mathbb{R}^d$, which can be collected into a 3-way tensor $W = (w_{pq}) \in \mathbb{R}^d \times P \times Q$ whose

(p, q) fiber is w_{pq} . We define the learning problem as follows: $\hat{W} = \arg\min_{W} L(W)$, W subject to $\|W\|_B \leq B_0$, (6)

$$\|W\|_B = \sum_{p,q} w_{pq}^2$$

where the norm $\|\cdot\|_B$ is either the overlapped trace norm, latent trace norm, or the scaled latent trace norm, and the empirical risk $L(W) = \frac{1}{m} \sum_{i=1}^m \sum_{p,q} (x_{ipq} - y_{ipq})^2$.

The true risk we are interested in minimizing is defined as follows: $L(W) = \mathbb{E}_{(x,y) \sim P} \sum_{p,q} (x_{ipq} - y_{ipq})^2$, $L(W) = \mathbb{E}_{(x,y) \sim P} \sum_{p,q} (x_{ipq} - y_{ipq})^2$

where P_{pq} is the distribution from which the samples $(x_{ipq}, y_{ipq})_{i=1}^m$ are drawn from.

$\mathbb{E} \|W\|_B$ with the expected dual norm $\mathbb{E} \|D\|_B$. The next lemma relates the excess risk $L(W) - L(\hat{W})$ through Rademacher complexity. Lemma 1. We assume that the output y_{ipq} is bounded as $|y_{ipq}| \leq b$, and the number of samples $m_{pq} \geq m_0 > 0$ for the observed tasks. We also assume that the loss function ℓ is Lipschitz continuous with the constant L , bounded in $[0, c]$ and $\ell(0) = 0$. Let W be any tensor such that $\|W\|_B \leq B_0$. Then with probability at least $1 - \delta$, any minimizer of (6) satisfies the following bound: $\mathbb{E} \|W\|_B \leq 2B \log(2/\delta) + \frac{1}{m} \sum_{p,q} \mathbb{E} \|D\|_B + \frac{1}{m} \sum_{p,q} \mathbb{E} \|D\|_B + \frac{1}{m} \sum_{p,q} \mathbb{E} \|D\|_B$ where $c = c + 1$, $\|D\|_B$ is the dual norm of $\|\cdot\|_B$, $\|D\|_B := \sum_{p,q} (p, q)^T S_m$. The tensor $D \in \mathbb{R}^{I \times J \times K}$ is defined as the sum $D = \sum_{p,q} (p, q)^T S_m$ where $Z \in \mathbb{R}^{I \times J \times K}$ is defined as $\{1\}_{ipq} x_{ipq}$, if $p = p^*$ and $q = q^*$, (p^*, q^*) th fiber of Z $= \sum_{i=1}^m x_{ipq}$, otherwise. Here $\{1\}_{ipq} \in \{0, 1\}$ are Rademacher random variables and the expectation in the above inequality is with respect to $\{1\}_{ipq}$, the random draw of tasks S , and the training samples $(x_{ipq}, y_{ipq})_{i=1}^m$. 3

Proof. The proof is a standard one following the line of [5] and it is presented in Appendix A. The next theorem computes the expected dual norm $\mathbb{E} \|D\|_B$ for the three norms for tensors (the proof can be found in Appendix B). Theorem 1. We assume that $C_{pq} := \mathbb{E}[x_{ipq} x_{ipq}] \geq d$ and there is a constant $R > 0$ such that $x_{ipq} \leq R$ almost surely. Let us define $D_1 := d + P Q$,

$$D_2 := P + dQ,$$

$$D_3 := Q + dP.$$

In order to simplify the presentation, we assume that $\max_k D_k \leq 3$ and $dP \leq \max(d_2, P^2, Q^2)$. For the overlapped trace norm, the latent trace norm, and the scaled latent trace norm, the expectation $\mathbb{E} \|D\|_B$ can be bounded as follows: (7) $\mathbb{E} \|D\|_B \leq C \min_k D_k \log D_k + \log D_k$, (8) $\mathbb{E} \|D\|_B \leq C \max_k D_k \log D_k + \log(\max_k D_k)$, (9) $\mathbb{E} \|D\|_B \leq C \max_k D_k \log D_k + \log(\max_k D_k)$ where C, C_1, C_2 are constants, $n_1 = d$, $n_2 = P$, and $n_3 = Q$. Furthermore, if $\max_k D_k \leq R^2 (\max_k D_k) \log(\max_k D_k)$, the $O(1/m)$ terms in the above inequalities can be dropped. Note

that the assumption that the norm of \mathbf{x}_{ipq} is bounded is natural because the target \mathbf{y}_{ipq} is also bounded. The parameter γ in the assumption $\mathbb{E} \|\mathbf{x}_{ipq}\|_2^2 \leq \gamma/d$ controls the amount of correlation in the data. Since $\text{Tr}(\mathbf{C}) = \mathbb{E} \|\mathbf{x}_{ipq}\|_2^2$, we have $\gamma = O(1)$ when the features are uncorrelated; on the other hand, we have $\gamma = O(d)$, if they lie in a one dimensional subspace. The number of γ samples $m - S = O(\max_k n_k)$ is enough to drop the $O(1/m - S)$ term even if $\gamma = O(1)$. Now we state the consequences of Theorem 1 for the three norms for tensors. The common assumptions are the same as in Lemma 1 and Theorem 1. We also assume $m - S = \gamma R^2 (\max_k n_k) \log(\max_k D_k) / \gamma$ to drop the $O(1/m - S)$ terms. Let \mathbf{W} be any $d \times P \times Q$ tensor with multilinear-rank (r_1, r_2, r_3) and bounded element-wise as $\|\mathbf{W}\|_{\max} \leq B$. Corollary 1 (Overlapped trace norm). With probability at least $1 - \delta$, any minimizer of (6) with $\|\cdot\|_{\text{W-overlap}}$ $B \gamma^{1/2} d P Q$ satisfies the following inequality: $\|\mathbf{W}\|_{\text{W-overlap}} \leq \log(2/\delta) \gamma^{1/2} L(\mathbf{W}) + c_1 B L(\mathbf{W}) \gamma^{1/2} \min(D_k \log D_k) + c_2 \gamma b + c_3$, $k = 1, \dots, m - S$ where $\gamma^{1/2} = (\sum_{k=1}^r D_k / 3)^{1/2}$ and c_1, c_2, c_3 are constants. Note that Tomioka et al. [23] obtained a bound that depends on $(\sum_{k=1}^r D_k / 3)^{1/2}$ instead of $\min(D_k \log D_k)$. Although the minimum may look better than the average, our bound has the worse constant $K = 3$ hidden in c_1 . The $\log D_k$ factor allows us to apply the above result to the setting of tensor completion as we show below. Corollary 2 (Latent trace norm). With probability at least $1 - \delta$, any minimizer of (6) with $\|\cdot\|_{\text{W-latent}}$ $B \min_k r_k d P Q$ satisfies the following inequality: $\|\mathbf{W}\|_{\text{W-latent}} \leq \log(2/\delta) \gamma^{1/2} L(\mathbf{W}) + c_1 B \min_k r_k \max(D_k \log D_k) + c_2 \gamma b + c_3$, $k = 1, \dots, m - S$ where c_1, c_2, c_3 are constants. Corollary 3 (Scaled γ latent trace norm). With probability at least $1 - \delta$, any minimizer of (6) with $\|\cdot\|_{\text{W-scaled}}$ $B \min_k (r_k / n_k) d P Q$ satisfies the following inequality: $\|\mathbf{W}\|_{\text{W-scaled}} \leq \log(2/\delta) \gamma^{1/2} L(\mathbf{W}) + c_1 B \min_k r_k \log(\max_k D_k) + c_2 \gamma b + c_3$, $k = 1, \dots, m - S$ where $n_1 = d, n_2 = P, n_3 = Q$, and c_1, c_2, c_3 are constants. 4

5

k, k, k, k

$\gamma^{1/2} \min(D_k \log D_k)$

k

$\gamma^{1/2} P Q \log(P Q)$

(r_1, r_2, r_3)

$(1, 1, m - S)$

$(\gamma, 1, m - S)$

$n_3!$

(r, P, r)

Tensor completion

$n_1! n_2!$

$Q!$

MLMTL [17] (heterogeneous case)

$d! P!$

MLMTL [17] (homogeneous case)

$\min_k r_k \max(D_k \log D_k)$

$dP \cdot Q \log(dQ)$
 k
 $d! \cdot d! \cdot d!$
 $\min(\text{nrkk})N \log(\max Dk)$
 $? \min(rP \cdot Q, dP \cdot r?) \log(dQ)$
 k
 $?(\min rk) d^2 \log(d^2)$
 $??r^{1/2} d^2 \log(d^2) \text{ } (?, 1, \text{---S---})$
 $?r^{1/2} \log(P \cdot Q) \text{ } ? (d, 1/d, P \cdot Q) (r, P, r?) \cdot PQ! \cdot d!$
 (heterogeneous MTL case)
 MTL [16] (homogeneous case)
 (r_1, r_2, r_3)
 $?r^{1/2} \log(d)^2$
 (r, d, d)
 $(1, r, r) \cdot d^2!$
 $d!$
 Matrix completion [11]
 $?(\min rk) d^2 \log(d^2)$
 $r \log(dQ) \cdot d \log(dQ)$
 2
 $r \log(d)$
 $P \cdot Q \log(P \cdot Q) \text{ } ?r^{1/2} (P + Q) \log(P + Q)$
 $(1, 1, \text{---S---}) \text{ } ? (d, 1/d, d^2)$
 Overlap $(?, B, \text{---S---}) (r_1, r_2, r_3)$
 $Q! \cdot P! \cdot 1!$
 (n_1, n_2, n_3)
 $r \log(d^2)$
 $r(P + Q) \log(P \cdot Q)$
 Scaled
 Sample complexities (per $1/?^2$) Latent

Table 2: Sample complexities of the overlapped trace norm, latent trace norm, and the scaled latent trace norm in various settings. The common factor $1/?^2$ is omitted from the sample complexities. The sample complexities are defined with respect to ---S--- for matrix completion, m for multitask learning, and $m\text{---S---}$ for tensor $?^3$ completion and multilinear multitask learning. In the heterogeneous cases, we assume $P \text{ } ? \text{ } r \text{ } ? \text{ } r?$. We define $?r^{1/2} = (\sum_{k=1}^r rk / K)^{1/2}$ and $N := n_1 \cdot n_2 \cdot n_3$.

We summarize the implications of the above corollaries for different settings in Table 2. We almost recover the settings for matrix completion [11] and multitask learning (MTL) [16]. Note that these simpler problems sometimes disguise themselves as the more general tensor completion or multilinear multitask learning problems. Therefore it is important that the new tensor based norms adapts to the simplicity of the problems in these cases. Matrix completion is when $d = ? = m = r_1 = 1$, and we assume that $r_2 = r_3 = r \text{ } ? \text{ } P, Q$. The sample complexities are the number of samples ---S--- that we need to make the leading term in Corollaries 1, 2, and 3 equal $?$. We can see that the overlapped

trace norm and the scaled latent trace norm recover the known result for matrix completion [11]. The plain latent trace norm requires $O(PQ)$ samples because it recognizes the first mode as the mode with the lowest rank 1. Although the rank r of the last two modes are low relative to their dimensions, the latent trace norm fails to recognize this. In multitask learning (MTL), only the first mode corresponding to features has a low rank r and the other two modes have full rank. Note that a tensor is a matrix when its multilinear rank is full except for one mode. We also assume that all the pairs (p, q) are observed ($|S| = PQ$) as in [16]. The sample complexities are defined the same way as above with respect to the number of samples m because $|S|$ is fixed. The homogeneous case is when $d = P = Q$. The heterogeneous case is when $P \neq Q \leq d$. Our bound for the overlapped trace norm is almost as good as the one in [16] but has an multiplicative $\log(PQ)$ factor (as oppose to their additive $\log(PQ)$ term) and $\frac{1}{2} \leq r \leq d$. Also note that the results in [16] can be applied when d is much larger than P and Q . Turning back to our bounds, both the latent trace norm and its scaled version can perform as well as knowing the mode with the lowest rank (the first mode) (see also [21]) when $d = P = Q$. However, when the dimensions are heterogeneous, similarly to the matrix completion case above, the plain latent trace norm fails to recognize the lowrank-ness of the first mode and requires $O(d)$ samples, because the second mode has the lowest rank P . In multilinear multitask learning (MLMTL) [17], any mode could possibly be low rank but it is a priori unknown. The sample complexities are defined the same way as above with respect to $m|S|$. The homogeneous case is when $d = P = Q$. The heterogeneous case is when the first mode or the third mode is low rank but $P \neq Q \leq d$. Similarly to the above two settings, the overlapped trace norm has a mild dependence on the dimensions but a higher dependence on the rank $\frac{1}{2} \leq r \leq d$. The latent trace norm performs as well as knowing the mode that has the lowest rank in the homogeneous case. However, it fails to recognize the mode with the lowest rank relative to its dimension. The scaled latent trace norm does this and although it has a higher logarithmic dependence, it is competitive in both cases. Finally, our bounds also hold for tensor completion. Although Tomioka et al. [22, 23] studied tensor completion algorithms, their analysis assumed that the inputs x_{ipq} are drawn from a Gaussian distribution, which does not hold for tensor completion. Note that in our setting x_{ipq} can be an indicator vector that has one in the j th position uniformly over $1, \dots, d$. In this case, $\frac{1}{2} = 1$. The sample complexities of different norms with respect to $m|S|$ is shown in the last row of Table 2. The sample complexity for the overlapped trace norm is the same as the one in [23] with a logarithmic factor. The sample complexities for the latent and scaled latent trace norms are new. Again we can see that while the latent trace norm recognize the mode with the lowest rank, the scaled latent trace norm is able to recognize the mode with the lowest rank relative to its dimension.

4

Experiments

We conducted several experiments to evaluate performances of tensor based multitask learning setting we have discussed in Section 3. In Section 4.1, we

discuss simulation we conducted using synthetic data sets. In Sections 4.2 and 4.3, we discuss experiments on two real world data sets, namely the Restaurant data set [26] and School Effectiveness data set [3, 4]. Both of our real world data sets have heterogeneous dimensions (see Figure 2) and it is a priori unclear across which mode has the most amount of information sharing. 4.1

Synthetic data sets

The true $d \times P \times Q$ tensor \mathcal{W} was generated by first sampling a $r_1 \times r_2 \times r_3$ core tensor and then multiplying random orthonormal matrix to each of its modes. For each task $(p, q) \in [P] \times [Q]$, we generated training set of m vectors $(x_{ipq}, y_{ipq})_{i=1}^m$ by first sampling x_{ipq} from the standard normal distribution and then computing $y_{ipq} = x_{ipq} w_{pq} + \epsilon_i$, where ϵ_i was drawn from a zero-mean normal distribution with variance 0.1. We used the penalty formulation of (6) with the squared loss and selected the regularization parameter λ using two-fold cross validation on the training set from the range 0.01 to 10 with the interval 0.1. In addition to the three norms for tensors we discussed in the previous section, we evaluated the matrix-based multitask learning approaches that penalizes the trace norm of the unfolding of \mathcal{W} at specific modes. The conventional convex multitask learning [2, 3, 16] corresponds to one of these approaches that penalizes the trace norm of the first unfolding $\text{tr}(\mathcal{W}^{(1)})$. The convex MLMTL in [17] corresponds to the overlapped trace norm. In the first experiment, we chose $d = P = Q = 10$ and $r_1 = r_2 = r_3 = 3$. Therefore, both the dimensions and the multilinear rank are homogeneous. The result is shown in Figure 1(a). The overlapped trace norm performed the best, the matrix-based approaches performed next, and the latent trace norm and the scaled latent trace norm were the worst. The scaling of the latent trace norm had no effect because the dimensions were homogeneous. Since the sample complexities for all the methods were the same in this setting (see Table 2), the difference in the performances could be explained by a constant factor $K(= 3)$ that is not shown in the sample complexities. In the second experiment, we chose the dimensions to be homogeneous as $d = P = Q = 10$, but $(r_1, r_2, r_3) = (3, 6, 8)$. The result is shown in Figure 1(b). In this setting, the (scaled) latent trace norm and the mode-1 regularization performed the best. The lower the rank of the corresponding mode, the lower were the error of the matrix-based MTL approaches. The overlapped trace norm was somewhat in the middle of the three matrix-based approaches. 6

0.016

0.022

Overlapped trace norm Latent trace norm Scaled Latent trace norm Mode
1 Mode 2 Mode 3

0.015

0.014

0.024

Overlapped trace norm Latent trace norm Scaled Latent trace norm Mode
1 Mode 2 Mode 3

0.02

0.018

	Overlapped trace norm	Latent trace norm	Scaled Latent trace norm	Mode
1	Mode 2	Mode 3		
	0.022			
	0.02			
	0.013			
	MSE			
	MSE			
	MSE			
	0.018	0.016		
	0.016	0.012		
	0.014			
	0.011			
	0.012			
	0.014			
	0.01	10		
	20			
	30			
	40			
	50			
	60			
	70			
	80			
	90			
	100			
	110			
	Sample size			
	0.01	10		
	0.012			
	20			
	30			
	40			
	50			
	60			
	70			
	80			
	90			
	100			
	110			
	0.01	10		
	Sample size			

(a) Synthetic experiment for the case when both the dimensions and the ranks are homogeneous. The true tensor is $10 \times 10 \times 10$ with multilinear rank $(3, 3, 3)$.

(b) Synthetic experiment for the case when the dimensions are homogeneous but the ranks are heterogeneous. The true tensor is $10 \times 10 \times 10$ with multilinear rank $(3, 6, 8)$.

20
30
40
50
60
70
80
90
100
110
Sample size

(c) Synthetic experiment for the case when both the dimensions and the ranks are heterogeneous. The true tensor is $10 \times 3 \times 10$ with multilinear rank $(3, 3, 8)$.

Figure 1: Results for the synthetic data sets. In the last experiment, we chose both the dimensions and the multilinear rank to be heterogeneous as $(d, P, Q) = (10, 3, 10)$ and $(r_1, r_2, r_3) = (3, 3, 8)$. The result is shown in Figure 1(c). Clearly the first mode had the lowest rank relative to its dimension. However, the latent trace norm recognizes the second mode as the mode with the lowest rank and performed similarly to the mode-2 regularization. The overlapped trace norm performed better but it was worse than the mode-1 regularization. The scaled latent trace norm performed comparably to the mode-1 regularization. 4.2

Restaurant data set

The Restaurant data set contains data for a recommendation system for restaurants where different customers have given ratings to different aspects of each restaurant. Following the same approach as in [17] we modelled the problem as a MLMTL problem with $d = 45$ features, $P = 3$ aspects, and $Q = 138$ customers. The total number of instances for all the tasks were 3483 and we randomly selected training set of sizes 400, 800, 1200, 1600, 2000, 2400, and 2800. When the size was small many tasks contained no training example. We also selected 250 instances as the validation set and the rest was used as the test set. The regularization parameter for each norm was selected by minimizing the mean squared error on the validation set from the candidate values in the interval $[50, 1000]$ for the overlapped, $[0.5, 40]$ for the latent, $[6000, 20000]$ for the scaled latent norms, respectively. We also evaluated matrix-based MTL approaches on different modes and ridge regression (Frobenius norm regularization; abbreviated as RR) as baselines. The convex MLMTL in [17] corresponds to the overlapped trace norm. The result is shown in Figure 2(a). We found the multilinear rank of the solution obtained by the overlapped trace norm to be typically $(1, 3, 3)$. This was consistent with the fact that the performances of the mode-1 regularization and the ridge regression were equal. In other words, the effective dimension of the first mode (features) was one instead of 45. The latent trace norm recognized the first mode as the mode with the lowest rank and it failed to take advantage of the low-rank-ness of the second and the third modes. The scaled latent trace norm was able to perform the best matching

than the rank of the other two modes. Clearly the scaled latent trace norm performed the best matching with the performance of the mode-2 regularization; probably the second mode had the most redundancy. The performance of the overlapped trace norm was comparable or slightly better than the mode-1 regularization. The percentage of the explained variance of the latent trace norm exceeds 30 % around sample size 4000 (around 30 samples per school), which is higher than the Hierarchical Bayes [4] (around 29.5 %) and matrix-based MTL [3] (around 26.7 %) that used around 80 samples per school.

5

Discussion

Using tensors for modeling multitask learning [17, 19] is a promising direction that allows us to take advantage of similarity of tasks in multiple dimensions and even make prediction for a task with no training example. However, having multiple modes, we would have to face with more hyperparameters to choose in the conventional nonconvex tensor decomposition framework. Convex relaxation of tensor multilinear rank allows us to side-step this issue. In fact, we have shown that the sample complexity of the latent trace norm is as good as knowing the mode with the lowest rank. This is consistent with the analysis of [21] in the tensor denoising setting (see Table 1). In the setting of tensor-based MTL, however, the notion of mode with the lowest rank may be vacuous because some modes may have very low dimension. In fact, the sample complexity of the latent trace norm can be as bad as not using any low-rank-ness at all if there is a mode with dimension lower than the rank of the other modes. The scaled latent trace norm we proposed in this paper recognizes the mode with the lowest rank relative to its dimension and lead to the competitive sample complexities in various settings we have shown in Table 2. Acknowledgment: MS acknowledges support from the JST CREST program.

2 References

[1] R. K. Ando and T. Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *J. Mach. Learn. Res.*, 6:1817-1853, 2005. [2] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In B. Schölkopf, J. Platt, and T. Hoffman, editors, *Adv. Neural. Inf. Process. Syst.* 19, pages 41-48. MIT Press, Cambridge, MA, 2007.

8

[3] A. Argyriou, M. Pontil, Y. Ying, and C. A. Micchelli. A spectral regularization framework for multi-task structure learning. In J. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Adv. Neural. Inf. Process. Syst.* 20, pages 25-32. Curran Associates, Inc., 2008. [4] B. Bakker and T. Heskes. Task clustering and gating for bayesian multitask learning. *J. Mach. Learn. Res.*, 4:83-99, 2003. [5] P. L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *J. Mach. Learn. Res.*, 3:463-482, 2002. [6] J. Baxter. A model of inductive bias learning. *J. Artif. Intell. Res.*, 12:149-198, 2000. [7] R. Caruana. Multitask learning. *Machine*

learning, 28(1):41?75, 1997. [8] L. De Lathauwer, B. De Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.*, 21(4):1253?1278, 2000. [9] L. De Lathauwer, B. De Moor, and J. Vandewalle. On the best rank-1 and rank-(R_1, R_2, \dots, R_N) approximation of higher-order tensors. *SIAM J. Matrix Anal. Appl.*, 21(4):1324?1342, 2000. [10] M. Fazel, H. Hindi, and S. P. Boyd. A Rank Minimization Heuristic with Application to Minimum Order System Approximation. In *Proc. of the American Control Conference*, 2001. [11] R. Foygel and N. Srebro. Concentration-based guarantees for low-rank matrix reconstruction. *Arxiv preprint arXiv:1102.3923*, 2011. [12] S. Gandy, B. Recht, and I. Yamada. Tensor completion and low-rank tensor recovery via convex optimization. *Inverse Problems*, 27:025010, 2011. [13] S. M. Kakade, S. Shalev-Shwartz, and A. Tewari. Regularization techniques for learning with matrices. *J. Mach. Learn. Res.*, 13(1):1865?1890, 2012. [14] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455?500, 2009. [15] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. In *Proc. ICCV*, 2009. [16] A. Maurer and M. Pontil. Excess risk bounds for multitask learning with trace norm regularization. Technical report, *arXiv:1212.1496*, 2012. [17] B. Romera-Paredes, H. Aung, N. Bianchi-Berthouze, and M. Pontil. Multilinear multitask learning. In *Proceedings of the 30th International Conference on Machine Learning*, pages 1444?1452, 2013. [18] M. Signoretto, L. De Lathauwer, and J. Suykens. Nuclear norms for tensors and their use for convex multilinear estimation. Technical Report 10-186, ESAT-SISTA, K.U.Leuven, 2010. [19] M. Signoretto, L. De Lathauwer, and J. A. K. Suykens. Learning tensors in reproducing kernel hilbert spaces with multilinear spectral penalties. Technical report, *arXiv:1310.4977*, 2013. [20] N. Srebro, J. D. M. Rennie, and T. S. Jaakkola. Maximum-margin matrix factorization. In L. K. Saul, Y. Weiss, and L. Bottou, editors, *Adv. Neural. Inf. Process. Syst.* 17, pages 1329?1336. MIT Press, Cambridge, MA, 2005. [21] R. Tomioka and T. Suzuki. Convex tensor decomposition via structured Schatten norm regularization. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Adv. Neural. Inf. Process. Syst.* 26, pages 1331?1339. 2013. [22] R. Tomioka, K. Hayashi, and H. Kashima. Estimation of low-rank tensors via convex optimization. Technical report, *arXiv:1010.0789*, 2011. [23] R. Tomioka, T. Suzuki, K. Hayashi, and H. Kashima. Statistical performance of convex tensor decomposition. In *Adv. Neural. Inf. Process. Syst.* 24, pages 972?980. 2011. [24] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.*, 12(4): 389?434, 2012. [25] L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279?311, 1966. [26] B. Vargas-Govea, G. Gonzalez-Serna, and R. Ponce-Medell?n. Effects of relevant contextual features in the performance of a restaurant recommender system. In *Proceedings of 3rd Workshop on Context-Aware Recommender Systems*. 2011.