

From Common to Special: When Multi-Attribute Learning Meets Personalized Opinions

Zhiyong Yang,^{1,2} Qianqian Xu,¹ Xiaochun Cao,¹ Qingming Huang^{3,4 *}

¹SKLOIS, Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

²School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

³University of Chinese Academy of Sciences, Beijing, China

⁴Key Lab of Intell. Info. Process., Inst. of Comput. Tech., CAS, Beijing, China
{yangzhiyong, xuqianqian, caoxiaochun}@iie.ac.cn, qmhuang@ucas.ac.cn

Abstract

Visual attributes, which refer to human-labeled semantic annotations, have gained increasing popularity in a wide range of real world applications. Generally, the existing attribute learning methods fall into two categories: one focuses on learning user-specific labels separately for different attributes, while the other one focuses on learning crowd-sourced global labels jointly for multiple attributes. However, both categories ignore the joint effect of the two mentioned factors: the personal diversity with respect to the global consensus; and the intrinsic correlation among multiple attributes. To overcome this challenge, we propose a novel model to learn user-specific predictors across multiple attributes. In our proposed model, the diversity of personalized opinions and the intrinsic relationship among multiple attributes are unified in a common-to-special manner. To this end, we adopt a three-component decomposition. Specifically, our model integrates a common cognition factor, an attribute-specific bias factor and a user-specific bias factor. Meanwhile Lasso and group Lasso penalties are adopted to leverage efficient feature selection. Furthermore, theoretical analysis is conducted to show that our proposed method could reach reasonable performance. Eventually, the empirical study carried out in this paper demonstrates the effectiveness of our proposed method.

Introduction

Visual attributes, which describe human labeled properties (like *open*, *fashionable*) for a given image, have shown its great potential as a mid-level semantic cue to enhance a variety of applications including face verification (Song, Tan, and Chen 2014), person re-identification (Su et al. 2016; 2017), and zero-shot learning (Ji et al. 2017; Wang et al. 2017; Zhang and Saligrama 2016), *etc.* Generally speaking, there are two types of attributes: i) binary attributes express whether a property is absent or present in a given image (like *A is/is not open*); ii) relative attributes show the strength of an attribute conveyed in one image with respect to another image (like *A is more/similarly/less open than B*) (Parikh and Grauman 2011).

On one hand, the attribute predictors are often trained with the crowd-sourced global labels. The justification of such an approach is that there is only one unique ground truth

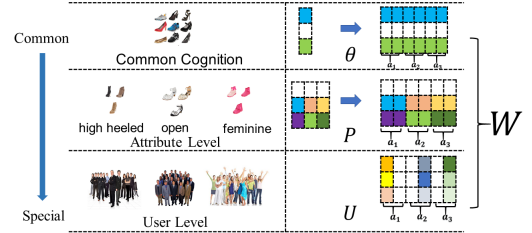


Figure 1: Illustration of the three-component decomposition. Here a_1 , a_2 and a_3 are the three mentioned attributes: high heeled, open and feminine. We assume that, for each attribute, there are two annotators who labeled the corresponding images. Note that we extend θ and P to match the size of U .

and the majority of annotators have to recognize this “*correct answer*” simultaneously. However, different annotators might very well have distinct preferences, such that participants of the crowdsourced experiments might vote under different criteria or conditions. It might be misleading to merely look at a global consensus while ignoring personal diversity. On the other hand, practical visual applications often involve simultaneously learning multiple attributes together. In such a case, different attributes may intrinsically share some common patterns. One reason is that these attributes convey similar semantic meaning. Another reason is that they use common subsets of low level image features. In that sense, training multiple attribute predictors independently might not be an appropriate protocol.

Based on the discussion above, our goal is to solve two problems simultaneously in this paper: 1) learning user-specific attributes, 2) learning multiple attributes together with their shared information.

For 1), (Kovashka and Grauman 2013) regard user-specific attribute learning as an adaption process. In this work, a generic model is first trained based on a large pool of crowd-sourced labels. Then a small user-specific dataset is employed to adapt the generic model to user-specific predictors. Meanwhile, (Kovashka and Grauman 2015) argue that one attribute may fit to different shades (interpretations)

*The corresponding author.

for different groups of persons. Correspondingly, the authors proposed an automatic shade discovery method to leverage group-wise user-specific attributes. Though these existing works are designed to deal with user-specific attributes, they neglect the mutual interactions between different attributes.

Multi-task learning framework is well known as a standard solution for 2). Recently, many efforts have been made to improve multi-task learning. (Ando and Zhang 2005) proposed an alternating structure optimization algorithm to decompose the predictive model of each task into two components: the task-specific feature mapping and task-shared feature mapping. For robust multi-task learning, (Chen, Zhou, and Ye 2011) proposed a corresponding method with a low-rank structure and a column-wise sparse structure. In (Gong, Ye, and Zhang 2012), the low rank structure proposed in (Chen, Zhou, and Ye 2011) was replaced by a row-wise sparse structure to leverage selection of a common subset of features. Since the tasks from the same group are closer to each other than those from a different group, (Zhou, Chen, and Ye 2011) proposed a clustering based multi-task learning framework. Motivated by the fact that the tasks should be related in terms of subsets of features, (Xu et al. 2015) proposed a novel multi-task learning method via task-feature co-clustering. As for applications, (Huang et al. 2014) proposed a robust dynamic multi-task method for trajectory regression.

There are also some existing works that focus on applying multi-task frameworks to attribute learning. One typical way to do this is to extend the existed multi-task algorithms to match attribute learning. For instance, in (Chen, Zhang, and Li 2014) the model proposed in (Gong, Ye, and Zhang 2012) was generalized to learn multiple relative attributes with their shared information. Meanwhile, some works employ deep learning methods to solve this problem by partially sharing the learned weights among different attributes. (Ehrlich et al. 2016) proposed a multi-task restricted boltzman machine so as to learn a shared feature representation for multiple facial attribute learning. (Hsieh, Hsu, and Chen 2017) incorporates identity and human attributes in learning discriminative face representations through a multi-task method. A deep multi-task learning approach was proposed in (Han et al. 2017) to jointly estimate multiple heterogeneous attributes from a single face image. (Hand and Chellappa 2017) also proposed a multi-task deep convolutional neural network with an auxiliary network at the top to capture attribute relationships. Though these works have successfully improved attribute learning with multi-task models, as was mentioned previously in this section, they all employ crowdsourced labels to train attribute predictors and ignore the disagreement among users.

Note that, except learning global labels for multiple attributes, multi-task frameworks are also suitable for learning user specific attributes where global patterns are necessary for capturing the public opinion, and task-specific patterns are indispensable as well for capturing user bias toward that public opinion. With such belief in mind, different with most of the previous works which partially met the requirement of our goal, we propose a hierarchical multi-task framework where task relationships are modeled on both the attribute

level and the user level.

The main contributions of this paper are two-fold:

- To match the hierarchical nature of the underlying problem, we propose a common-to-special decomposition of the model weights, which captures the general cognition pattern, attribute level bias and user specific bias, respectively. An optimization method is established based on the accelerated proximal gradient method.
- Theoretical analysis is performed in this paper. The corresponding results show that our proposed algorithm could attain reasonable performance.

Methodology

In this section, we'll present an attribute learning method to learn user specific labels across multiple attributes. We first introduce the notations used in this paper. Secondly, we propose our model formulation, which includes a common-to-special decomposition of the model weights and the corresponding objective function. Thirdly, we introduce our optimization method based on the accelerated proximal gradient method. Finally, the theoretical analysis is carried out to show the performance bound of our method.

Notations

In this paper, scalars, vectors, and matrices are denoted as lowercase letters (a), bold lower case letters (\mathbf{a}), and bold upper case letters (\mathbf{A}). \mathbf{X}_k denotes the k th row of \mathbf{X} . x_{ij} denotes the (i, j) entry of a matrix \mathbf{X} . $\mathbb{P}(\cdot)$ denotes a probability measure. $[a]$ denotes the set $\{1, 2, \dots, a\}$. Given an index set \mathcal{I} , $\mathbf{A}^{\mathcal{I}}$ denotes a matrix that contains all the corresponding rows of \mathbf{A} , while $\mathbf{a}^{\mathcal{I}}$ represents the vector that contains the corresponding elements of vector \mathbf{a} . $\|\cdot\|_p$ denotes the ℓ_p norm : $\|\mathbf{x}\|_p = (\sum_i x_i^p)^{(1/p)}$. $\|\cdot\|_{p,q}$

denotes the $\ell_{p,q}$ norm : $\left(\sum_i \left(\sum_j (x_{ij}^q)^{(1/q)}\right)^p\right)^{(1/p)}$. $\|\cdot\|$

denotes the ℓ_2 norm. $\langle \cdot, \cdot \rangle$ denotes the inner product for two matrices or two vectors. If $f(x) = o(g(x))$, we have

$\lim_{x \rightarrow +\infty} \frac{f(x)}{g(x)} = 0$. \oplus denotes the direct sum of two matrices.

Model Formulation

Assume that we have n_a attributes to be evaluated, and that, for the i th attribute, we are given user-specific labels from n_{u_i} different workers. Then the training data could be represented as:

$$\mathcal{T} = \left\{ (\mathbf{X}^{(1,1)}, \mathbf{y}^{(1,1)}), \dots, (\mathbf{X}^{(n_a, n_{u_{n_a}})}, \mathbf{y}^{(n_a, n_{u_{n_a}})}) \right\}$$

where n_{ij} is the number of images the j th user for the i th attribute labeled. The input feature is preprocessed such that

$$\sum_{k=1}^{n_{ij}} \left(x_{k,l}^{(i,j)} \right)^2 = 1 \quad (1)$$

, where $x_{k,l}^{(i,j)}$ is the (k, l) th entry of $\mathbf{X}^{(i,j)}$.

For binary attributes, $\mathbf{X}^{(i,j)} \in \mathbb{R}^{n_{ij} \times d}$ is the corresponding feature matrix for the images that the j th user of the i th

attribute labeled. Each row of $\mathbf{X}^{(i,j)}$ represents the low-level feature for a corresponding image. $\mathbf{y}^{(i,j)} \in \{-1, 1\}^{n_{ij}}$ is the corresponding label vector¹. If $y_k^{(i,j)} = 1$, then the user thinks that the corresponding attribute is present in the k th image, otherwise it will be labeled as -1.

For relative attributes, we need to solve a ranking problem. The corresponding users are given a set of image pairs $\{(\mathbf{x}_{1,k}^{(i,j)}, \mathbf{x}_{2,k}^{(i,j)})\}_{k=1}^{n_{ij}}$. Since we adopt linear models in this paper, we define the k th row of $\mathbf{X}^{(i,j)}$ as $\mathbf{X}_k^{(i,j)} = \mathbf{x}_{1,k}^{(i,j)} - \mathbf{x}_{2,k}^{(i,j)}$. For the k th pair $y_k^{(i,j)} = 1$ if the user thinks that the corresponding attribute has a stronger expression in the former image (say 1 is more open than 2); $y_k^{(i,j)} = 0$ if the user thinks that both images show similar strength for the current attribute (say 1 is as open as 2); $y_k^{(i,j)} = -1$ if the user thinks that the corresponding attribute has a weaker expression in the former image (say 1 is less open than 2).

As mentioned in the introduction section, our goal is to learn a predictor $\mathbf{f}^{(i,j)}$ for each of the personalized label vectors $\mathbf{y}^{(i,j)}$. In this paper, we assume that $\mathbf{f}^{(i,j)}(\cdot)$ has a linear form :

$$\mathbf{f}^{(i,j)} = \mathbf{X}^{(i,j)} \mathbf{w}^{(i,j)} \quad (2)$$

where $\mathbf{w}^{(i,j)}$ is the corresponding model weight.

Now we are ready to introduce the modeling of $\mathbf{w}^{(i,j)}$. Note that our underlying problem could be comprehended in a common-to-special manner: there is a common pattern that captures the general cognition of a given object; an attribute-specific pattern is also necessary to express the attribute-level common pattern; finally, a user-specific factor is necessary to describe the personalized preference. As a result, we adopt a three-component additive decomposition of $\mathbf{w}^{(i,j)}$:

$$\mathbf{w}^{(i,j)} = \boldsymbol{\theta} + \mathbf{p}^{(i)} + \mathbf{u}^{(i,j)} \quad (3)$$

The practical meaning of these three components could be explained as follows:

- $\boldsymbol{\theta}$ (General Cognition Factor): $\boldsymbol{\theta} \in \mathbb{R}^{d \times 1}$ is the global factor that captures the overall cognition for the given class of object (say for shoes dataset, this factor captures the overall cognition about shoes). $\boldsymbol{\theta}$ is shared among all subtasks.
- $\mathbf{p}^{(i)}$ (Attribute Specific Bias Factor): An attribute-specific factor that captures the bias of the i th attribute with respect to the global cognition. For mathematical convenience, we denote $\mathbf{P} = [\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(n_a)}]$, and we have $\mathbf{P} \in \mathbb{R}^{d \times n_a}$.
- $\mathbf{u}^{(i,j)}$ (User Specific Bias Factor): A user specific factor that captures the personal bias for the j th user of the i th attribute. In order to simplify the mathematical expressions, we define $\mathbf{U}^{(t)} = [\mathbf{u}^{(t,1)}, \dots, \mathbf{u}^{(t,n_{u_t})}]$, $\mathbf{U} =$

¹If not explained, $A^{(i,j)}$ denotes the corresponding variable of A for the j th user of the i th attribute; $A^{(i)}$ denotes the corresponding variable for the i th attribute.

$[\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(n_a)}]$, thus we have $\mathbf{U} \in \mathbb{R}^{d \times n_u}$, where $n_u = \sum_{i=1}^{n_a} n_{u_i}$.

Since we adopt a linear form for $\mathbf{f}^{(i,j)}$, it is natural to assume that the real response $\mathbf{y}^{(i,j)}$ could be interpreted as the true predictor in our proposed model: $\mathbf{f}^{*(i,j)} = \mathbf{X}^{(i,j)} \mathbf{w}^{*(i,j)}$ plus a Gaussian noise $\boldsymbol{\delta}^{(i,j)} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$:

$$\mathbf{y}^{(i,j)} = \mathbf{f}^{*(i,j)} + \boldsymbol{\delta}^{(i,j)} \quad (4)$$

where $\mathbf{w}^{*(i,j)} = \boldsymbol{\theta}^* + \mathbf{p}^{*(i)} + \mathbf{u}^{*(i,j)}$.

As for the objective function, we adopt the least square loss as our empirical loss: $L(\mathbf{W})$ and a general regularizer $\Omega(\cdot)$. We could thus formulate our problem as (P1) as:

$$(P1) : \min_{\mathbf{W}} \underbrace{\sum_{i=1}^{n_a} \sum_{j=1}^{n_{u_i}} \|\mathbf{y}^{(i,j)} - \mathbf{X}^{(i,j)} \mathbf{w}^{(i,j)}\|^2}_{L(\mathbf{W})} + \Omega(\mathbf{W})$$

where $\mathbf{W} := \{\mathbf{w}^{(i,j)}\}_{(i,j)}$ is the set of all parameters.

For our problem, user annotated labels are often limited. Furthermore, the low level features for an image are located in a high dimensional space. Then it is necessary to encourage sparse models to reduce model complexity. To preserve the relationship among subtasks, we also need to leverage a shared feature subset. Consequently, we penalize $\boldsymbol{\theta}$ and \mathbf{P} with ℓ_1 norm and $\ell_{1,2}$ norm, respectively. Meanwhile, $\mathbf{u}^{(i,j)} \neq \mathbf{0}$ only when the corresponding user has a specific bias with respect to the popular opinion. We penalize \mathbf{U}^\top with $\ell_{1,2}$ norm to leverage column-wise sparsity. Above all, $\Omega(\mathbf{W})$ could be represented as follows :

$$\Omega(\mathbf{W}) = \lambda_1 \|\boldsymbol{\theta}\|_1 + \lambda_2 \|\mathbf{P}\|_{1,2} + \lambda_3 \|\mathbf{U}^\top\|_{1,2}$$

Putting all together, (P1) could be reformed as :

$$\min_{\mathbf{W}} \sum_{i=1}^{n_a} \sum_{j=1}^{n_{u_i}} \|\mathbf{y}^{(i,j)} - \mathbf{X}^{(i,j)} (\boldsymbol{\theta} + \mathbf{p}^{(i)} + \mathbf{u}^{(i,j)})\|^2 + \lambda_1 \|\boldsymbol{\theta}\|_1 + \lambda_2 \|\mathbf{P}\|_{1,2} + \lambda_3 \|\mathbf{U}^\top\|_{1,2} \quad (5)$$

Figure 1 illustrates the expected structure of the three components in the proposed model. It could be seen that both the attribute level and the user level task correlations are included in our proposed model.

To end this section, we introduce two important mathematical properties of (P1) as proposition 1 and proposition 2. Refer to our supplementary materials for a detailed proof of proposition 1 and proposition 2.

Proposition 1 (Global Optimality). *P1 is jointly convex with respect to $\boldsymbol{\theta}$, \mathbf{P} , \mathbf{U}*

Proposition 2 (Lipschitz Continuous Gradient). *Given two arbitrary feasible solutions \mathbf{W} and \mathbf{W}' , we have :*

$$\|\nabla L(\tilde{\mathbf{W}}) - \nabla L(\tilde{\mathbf{W}}')\| \leq \rho \|\tilde{\mathbf{W}} - \tilde{\mathbf{W}}'\|$$

where:

$$\rho = 6n_u \sqrt{(n_u + n_a + 1)} \max_{i,j} \left[\sigma_1 \left(\mathbf{X}^{(i,j)} \right) \right]^2$$

$$\tilde{W} = [\text{vec}(\boldsymbol{\theta})^\top, \text{vec}(\mathbf{P})^\top, \text{vec}(\mathbf{U})^\top]^\top$$

$$\tilde{W}' = [\text{vec}(\boldsymbol{\theta}')^\top, \text{vec}(\mathbf{P}')^\top, \text{vec}(\mathbf{U}')^\top]^\top, n_u = \sum_{i=1}^{n_a} n_{u_i}$$

Optimization Method

According to proposition 1 and proposition 2, $L(\mathbf{W})$ is a convex and smooth function while $\Omega(\mathbf{W})$ is a convex non-smooth function. Similar as the related literatures (Beck and Teboulle 2009; Chen, Zhou, and Ye 2011; Gong, Ye, and Zhang 2012), the accelerated proximal gradient method is employed to solve (P1).

According to proposition 1, $L(\cdot)$ is differentiable with Lipschitz continuous gradient. According to the basic mathematical properties of Lipschitz continuous functions, at iteration step k , for any reference point $\mathbf{W}^{ref_k} = (\boldsymbol{\theta}^{ref_k}, \mathbf{P}^{ref_k}, \mathbf{U}^{ref_k})$, $\exists \rho_k > 0$, such that :

$$\begin{aligned} L(\mathbf{W}) &\leq L(\mathbf{W}^{ref_k}) + \langle \nabla_P L(\mathbf{W}^{ref_k}), \Delta \mathbf{P} \rangle \\ &+ \langle \nabla_\theta L(\mathbf{W}^{ref_k}), \Delta \boldsymbol{\theta} \rangle + \langle \nabla_U L(\mathbf{W}^{ref_k}), \Delta \mathbf{U} \rangle \\ &+ \frac{\rho_k}{2} \|\Delta \mathbf{P}\|_F^2 + \frac{\rho_k}{2} \|\Delta \boldsymbol{\theta}\|_2^2 + \frac{\rho_k}{2} \|\Delta \mathbf{U}\|_F^2 \\ &\stackrel{def}{=} \hat{L}_{\mathbf{W}^{ref_k}, \rho}(\mathbf{W}) \end{aligned} \quad (6)$$

Following the Majorization-Minimization (MM) (Hunter 2004) scheme, at the k -th iteration, we could then solve (P2) instead of updating the weights based on original problem:

$$(P2) : (\boldsymbol{\theta}^k, \mathbf{P}^k, \mathbf{U}^k) := \underset{\boldsymbol{\theta}, \mathbf{P}, \mathbf{U}}{\operatorname{argmin}} \hat{L}_{\mathbf{W}^{ref_k}, \rho_k}(\mathbf{W}) + \Omega(\mathbf{W})$$

It is obvious that $\boldsymbol{\theta}, \mathbf{P}, \mathbf{U}$ are decoupled in $\hat{L}_{\mathbf{W}^{ref_k}, \rho_k}(\mathbf{W})$. Hence, solving (P2) is equivalent to solving the following three proximal subproblems simultaneously:

$$\boldsymbol{\theta}^k := \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}^k\|^2 + \frac{\lambda_1}{\rho_k} \|\boldsymbol{\theta}\|_1 \quad (7)$$

$$\mathbf{P}^k := \underset{\mathbf{P}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{P} - \tilde{\mathbf{P}}^k\|^2 + \frac{\lambda_2}{\rho_k} \|\mathbf{P}\|_{1,2} \quad (8)$$

$$\mathbf{U}^k := \underset{\mathbf{U}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{U} - \tilde{\mathbf{U}}^k\|^2 + \frac{\lambda_3}{\rho_k} \|\mathbf{U}^\top\|_{1,2} \quad (9)$$

where

$$\tilde{\boldsymbol{\theta}}^k = \boldsymbol{\theta}^{ref_k} - \frac{1}{\rho_k} \nabla_\theta L(\mathbf{W}^{ref_k})$$

$$\tilde{\mathbf{P}}^k = \mathbf{P}^{ref_k} - \frac{1}{\rho_k} \nabla_P L(\mathbf{W}^{ref_k})$$

$$\tilde{\mathbf{U}}^k = \mathbf{U}^{ref_k} - \frac{1}{\rho_k} \nabla_U L(\mathbf{W}^{ref_k})$$

All of these subproblems admit closed-form solutions :

$$\boldsymbol{\theta}_i^k := \operatorname{sign}(\tilde{\boldsymbol{\theta}}_i^k) \left(\left| \tilde{\boldsymbol{\theta}}_i^k \right| - \frac{\lambda_1}{\rho_k} \right)_+ \quad (10)$$

Algorithm 1: The accelerated proximal gradient method for solving (P1)

Input: $\mathcal{T}, \lambda_1, \lambda_2, \lambda_3, \rho_0 > 0, \eta > 1$

Output: $\boldsymbol{\theta}, \mathbf{P}, \mathbf{U}$

```

1 Initialize  $\boldsymbol{\theta}^0, \mathbf{P}^0, \mathbf{U}^0$ ;
2  $\boldsymbol{\theta}^{ref} := \boldsymbol{\theta}^0, \mathbf{P}^{ref} := \mathbf{P}^0, \mathbf{U}^{ref} := \mathbf{U}^0, t_1 := 1, k := 1$ ;
3 while Not Converged do
4   Solve  $\mathbf{W}^k = (\boldsymbol{\theta}^k, \mathbf{P}^k, \mathbf{U}^k)$  with Eq.(10)-Eq.(12);
5   Find the smallest  $i_k$  such that when  $\tilde{\rho} = \eta^{i_k} \tilde{\rho}_{k-1}$  :
      $L(\mathbf{W}^k) \leq \hat{L}_{\mathbf{W}^{ref_k}, \tilde{\rho}}(\mathbf{W}^k)$ ;
6    $\rho_k := \tilde{\rho}$ ;
7   update  $\mathbf{W}^k$  again;
8    $t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$ ;
9    $dt := \frac{t_k - 1}{t_{k+1}}$ ;
10   $\boldsymbol{\theta}^{ref} := \boldsymbol{\theta}^k + dt(\boldsymbol{\theta}^k - \boldsymbol{\theta}^{k-1})$ ;
11   $\mathbf{P}^{ref} := \mathbf{P}^k + dt(\mathbf{P}^k - \mathbf{P}^{k-1})$ ;
12   $\mathbf{U}^{ref} := \mathbf{U}^k + dt(\mathbf{U}^k - \mathbf{U}^{k-1})$ ;
13   $k := k + 1$ ;
14 end
```

$$\mathbf{P}_i^k := \left(1 - \frac{\lambda_2}{\rho_k \|\tilde{\mathbf{P}}_i^k\|} \right)_+ \tilde{\mathbf{P}}_i^k \quad (11)$$

$$\left(\mathbf{u}^{(i,j)} \right)^k := \left(1 - \frac{\lambda_3}{\rho_k \left\| \left(\tilde{\mathbf{u}}^{(i,j)} \right)^k \right\|} \right)_+ \left(\tilde{\mathbf{u}}^{(i,j)} \right)^k \quad (12)$$

Furthermore, the nestrov's acceleration strategy is employed to update the reference point \mathbf{W}^{ref} . Integrating all the results, an efficient algorithm to solve (P1) is introduced as Algorithm 1. According to the theoretical analyses proposed in (Beck and Teboulle 2009), \mathbf{W}^k could converge to a global optimal solution with rate $\mathcal{O}(\frac{1}{k^2})$, which is the provable optimal rate for first-order methods.

Theoretical Analysis

Following (Gong, Ye, and Zhang 2012), we will propose the performance bound of our algorithm based on assumption 1.

Here we define a set $\mathcal{N}(\mathbf{A})$ for a matrix (vector) \mathbf{A} as the indexes for zero rows (entries): $\mathcal{N}(\mathbf{A}) = \{i : \mathbf{A}_i = \mathbf{0}\}$, and $\mathcal{N}_\perp(\mathbf{A})$ as the complement of $\mathcal{N}(\mathbf{A})$: $\mathcal{N}_\perp(\mathbf{A}) = \{i : \mathbf{A}_i \neq \mathbf{0}\}$. $|\mathcal{N}(\mathbf{A})|$ is the number of zero rows (entries) of a matrix(vector) \mathbf{A} , and we have similar definition for $|\mathcal{N}_\perp(\mathbf{A})|$. Now we provide the basic assumption of the main result as assumption 1.

Assumption 1. It is defined that : $\mathbf{X} \stackrel{def}{=} \oplus_{i,j} \mathbf{X}^{(i,j)}$,

$$\bar{\mathbf{W}} \stackrel{def}{=} [\mathbf{w}^{(1,1)\top}, \mathbf{w}^{(1,2)\top}, \dots, \mathbf{w}^{(n_a, n_{una})\top}]^\top$$

Let $0 \leq n_\theta \leq d$ be the upper bound of $|\mathcal{N}_\perp(\theta^*)|$, $0 \leq n_p \leq d$ be the upper bound of $|\mathcal{N}_\perp(\mathbf{P}^*)|$, $0 \leq n_{u,a} \leq n_u$ be the upper bound of $|\mathcal{N}_\perp(\mathbf{U}^{*\top})|$. We assume that there exists nonnegative $\kappa_\theta, \kappa_p, \kappa_{u,a}$:

$$\begin{aligned}\kappa_\theta &= \min_{\Gamma_\theta, \Gamma_P, \Gamma_U} \frac{\|X\Gamma_{\bar{W}}\|}{\sqrt{n_{\min}n_u}\|\Gamma_\theta^{\mathcal{N}_\perp(\theta)}\|_2} \\ \kappa_p &= \min_{\Gamma_\theta, \Gamma_P, \Gamma_U} \frac{\|X\Gamma_{\bar{W}}\|}{\sqrt{n_{\min}n_u}\|\Gamma_P^{\mathcal{N}_\perp(P)}\|_F} \\ \kappa_{u,a} &= \min_{\Gamma_\theta, \Gamma_P, \Gamma_U} \frac{\|X\Gamma_{\bar{W}}\|}{\sqrt{n_{\min}n_u}(\Gamma_U)^\top \mathcal{N}_\perp(U^\top) \|_F}\end{aligned}$$

where $\Gamma_\theta \in \mathbb{R}^d$ is a function of θ ; $\Gamma_P \in \mathbb{R}^{d \times n_a}$ is a function of \mathbf{P} ; $\Gamma_U \in \mathbb{R}^{d \times n_u}$ is a function of \mathbf{U} . Define $\Gamma_{w^{(i,j)}} = \Gamma_\theta + \Gamma_P^{(i)} + \Gamma_{u^{(i,j)}}$, then $\Gamma_{\bar{W}} = [\Gamma_{w^{(1,1)}}^\top, \Gamma_{w^{(1,2)}}^\top, \dots, \Gamma_{w^{(n_a, n_u n_a)}}^\top]^\top$. Furthermore, it is assumed that the following inequalities hold: $\|\Gamma_\theta^{\mathcal{N}(\theta)}\|_1 \leq \beta_\theta \|\Gamma_\theta^{\mathcal{N}_\perp(\theta)}\|_1$, $\|\Gamma_P^{\mathcal{N}(P)}\|_{1,2} \leq \beta_p \|\Gamma_P^{\mathcal{N}_\perp(P)}\|_{1,2}$, $\|\Gamma_U^{\mathcal{N}(U^\top)}\|_{1,2} \leq \beta_{u,a} \|\Gamma_U^{\mathcal{N}_\perp(U^\top)}\|_{1,2}$. where $\beta_\theta, \beta_p, \beta_{u,a}$ are positive scalars.

Note that the assumption on $\kappa_\theta, \kappa_p, \kappa_{u,a}$ is based on the restricted eigenvalue assumption (Bickel, Ritov, and Tsybakov 2008), which has been widely used in existing multi-task literatures (Chen, Zhou, and Ye 2011; Gong, Ye, and Zhang 2012).

According to the notations in assumption 1, the squared error between the predicted value $\mathbf{x}^\top \mathbf{w}$ and the real value \mathbf{f} could be formed as: $\|\mathbf{X}\bar{\mathbf{W}} - \mathbf{F}\|$, where \mathbf{F} is defined as:

$$\mathbf{F} = [\mathbf{f}^{*(1,1)\top}, \mathbf{f}^{*(1,2)\top}, \dots, \mathbf{f}^{*(n_a, n_u n_a)\top}]^\top$$

Let $\hat{\mathbf{W}} = (\hat{\theta}, \hat{\mathbf{P}}, \hat{\mathbf{U}})$ be an optimal solution of (P1). According to Eq.(4), we define $\mathbf{W}^* = (\theta^*, \mathbf{P}^*, \mathbf{U}^*)$. Our main result could be presented as Theorem 1.

Theorem 1 (Performance Bounds). Define $\alpha = 2\sigma\sqrt{dn_u} + t$, choose $\lambda_1, \lambda_2, \lambda_3$ as: $\lambda_1 \geq n_u\alpha, \lambda_2 \geq \tilde{n}\alpha, \lambda_3 \geq \alpha$, where $n_u = \sum_{i=1}^{n_a} n_{u_i}$ and $\tilde{n} = \sqrt{\sum_i n_{u_i}^2}$. Given

Assumption 1, let

$$\zeta = \frac{\lambda_1\sqrt{n_\theta}}{\kappa_\theta} + \frac{\lambda_2\sqrt{n_p}}{\kappa_p} + \frac{\lambda_3\sqrt{n_{u,a}}}{\kappa_{u,a}}$$

for $t > 0$, we have:

$$\mathbb{P}\left(\frac{1}{n_{\min}n_u}\|\mathbf{X}\bar{\mathbf{W}} - \mathbf{F}\|^2 \leq \left(\frac{2\zeta}{n_{\min}n_u}\right)^2\right) \geq \delta(t) \quad (13)$$

$$\mathbb{P}\left(\|\hat{\theta} - \theta^*\|_1 \leq \frac{2(\beta_\theta + 1)\sqrt{n_\theta}}{\kappa_\theta n_{\min}n_u}\zeta\right) \geq \delta(t) \quad (14)$$

$$\mathbb{P}\left(\|\hat{\mathbf{P}} - \mathbf{P}^*\|_{1,2} \leq \frac{2(\beta_p + 1)\sqrt{n_p}}{\kappa_p n_{\min}n_u}\zeta\right) \geq \delta(t) \quad (15)$$

$$\mathbb{P}\left(\|\hat{\mathbf{U}}^\top - \mathbf{U}^{*\top}\|_{1,2} \leq \frac{2(\beta_{u,a} + 1)\sqrt{n_{u,a}}}{\kappa_{u,a} n_{\min}n_u}\zeta\right) \geq \delta(t) \quad (16)$$

where $\delta(t) = 1 - \frac{1}{\sqrt{2\pi}Z_{dn_u}(t)}\exp\left(-\frac{Z_{dn_u}(t)}{2}\right)$; $n_{\min} = \min_{i,j} n_{ij}$ and $Z_{dn_u}(t) = t - dn_u \log(1 + \frac{t}{dn_u})$.

Remark 1. According to theorem 1, if

$$\zeta = o\left(n_{\min}n_u \min\left\{1, \frac{\kappa_\theta}{\sqrt{n_\theta}}, \frac{\kappa_p}{\sqrt{n_p}}, \frac{\kappa_{u,a}}{\sqrt{n_{u,a}}}\right\}\right)$$

We have: $\mathbb{E}(\mathbf{X}^{(i,j)}\hat{\mathbf{w}}^{(i,j)} - \mathbf{f}^{*(i,j)}) \rightarrow 0$, $\hat{\theta} \xrightarrow{\ell_1} \theta^*$, $\hat{\mathbf{P}} \xrightarrow{\ell_{1,2}} \mathbf{P}^*$, and $\hat{\mathbf{U}}^\top \xrightarrow{\ell_{1,2}} \mathbf{U}^{*\top}$ hold with high probability when $n_{\min} \rightarrow \infty$. Furthermore, though the proof of theorem 1 uses standard techniques developed in (Gong, Ye, and Zhang 2012), $\delta(t)$ in theorem 1 is a tighter probability bound than $\left(1 - \exp(-\frac{Z_{dn_u}(t)}{2})\right)$ proposed in (Gong, Ye, and Zhang 2012) and $\left(1 - n_u \exp(-\frac{Z_{dn_u}(t)}{2})\right)$ proposed in (Chen, Zhou, and Ye 2011), for sufficiently large t .

According to theorem 1 and remark 1, we see that our proposed method could both leverage good performance and estimate the parameters well with high probability.

Experiment

Now in this section, we show our experiment results on a simulated dataset, and two real world datasets respectively.

Experiment Setting

For each subtask, we randomly split the corresponding samples into a training subset and test subset, with 40% and 80% of the samples selected as training set respectively. For each involved algorithm, the hyper-parameters are tuned based on a 3 fold cross validation on the training set, and the average performance of the test set on 5 different splits are recorded. It is important to note that, different with the setting of (Kovashka and Grauman 2013), we will not use any extra dataset for pre-training in this paper. Furthermore, the training data is preprocessed according to Eq.(1)

Simulated Dataset

Data Generation For simulated dataset, our goal is to verify that the proposed algorithm could reach reasonable performance based on our theoretical analysis. We here define a regression problem for this dataset. To this end, we generate simulated features and continuous user scores (but not discrete labels) for 5 attributes, and all n_{u_i} s are fixed as 10. Furthermore, we set the dimensionality d as 50. For the (i, j) th subtask, 300 samples are generated such that $\mathbf{X}^{(i,j)} \sim \mathcal{N}(0, 4\mathbf{I})$ and $\mathbf{y}^{(i,j)} = \mathbf{X}^{(i,j)}\mathbf{w}^{(i,j)} + \mathcal{N}(0, \mathbf{I})$. To leverage group sparsity of \mathbf{W} : θ is generated as $\theta \sim \mathcal{N}(1, 4 * \mathbf{1})$, and the first 15 elements are set as zero; \mathbf{P} is generated as $\mathbf{P} \sim \mathcal{N}(1, 5\mathbf{I})$ and the 20-35 th rows of \mathbf{P} are set as $\mathbf{0}$; \mathbf{U} is generated as $\mathbf{U} \sim \mathcal{N}(1, 10\mathbf{I})$ and the first 2 columns of each $\mathbf{U}^{(t)}$ are set as $\mathbf{0}$.

Competitors and Evaluation Metric To verify the effectiveness of our proposed algorithm, we compare our algorithm with four benchmark algorithms for regression: Support Vector Regression (SVR), Lasso regression (Lasso), Ridge regression (Ridge), and the Elastic Net. Meanwhile, all of these four benchmark algorithms are employed in a user-exclusive manner : the predictors for all the subtasks are trained separately as independent tasks. To evaluate the generalized performance of the algorithms, for a given attribute, we adopt \overline{NMSE} , the average value of normalized mean square error (NMSE) on all users for a given attribute, as the evaluation metric.

Table 1: Performance comparison for the simulated dataset with 40% samples selected as training data

SVR	Ridge	Lasso	Elastic Net	ours
1.000	0.830	2.99E-05	9.96E-05	2.70E-05
1.000	0.830	3.20E-05	9.92E-05	2.82E-05
1.000	0.830	3.15E-05	1.07E-04	2.78E-05
0.998	0.829	4.20E-05	1.14E-04	3.75E-05
1.005	0.834	2.88E-05	9.34E-05	2.72E-05
1.001	0.830	3.28E-05	1.03E-04	2.95E-05

Table 2: Performance comparison for the simulated dataset with 80% samples selected as training data

SVR	Ridge	Lasso	Elastic Net	ours
1.005	0.884	2.27E-05	5.19E-05	2.11E-05
1.009	0.888	2.37E-05	5.27E-05	2.17E-05
1.009	0.887	2.30E-05	5.14E-05	2.12E-05
1.010	0.889	2.79E-05	5.64E-05	2.60E-05
1.019	0.896	2.03E-05	4.79E-05	1.89E-05
1.010	0.889	2.35E-05	5.20E-05	2.18E-05

Performance Comparison According to Table 1 and Table 2, we could draw the following conclusions. The six rows in these two tables records the performance for attributes 1-5 and their average, respectively. On one hand, due to the inability to leverage sparse parameters, we see that performance of SVR and Ridge couldn't outperform the other three algorithms. On the other hand, our proposed algorithm reaches the best performance based on \overline{NMSE} , which verifies the effectiveness of our proposed algorithm.

Parameters Recovery Now we show the ability of our algorithm to recover the structured parameters. With the same simulated dataset, we select 80% the samples as training data. According to theorem 1, we set $t = 10$, $\lambda_1 = 2n_u\alpha$, $\lambda_2 = 2.5\tilde{n}\alpha$, $\lambda_3 = 32\alpha$. Figure 2 shows the resulting parameters of our algorithm. Note that we extend θ and P to $\mathbb{R}^{d \times n_u}$, so that they could match the size of U . As is shown in this figure, we conclude that all of these three sets of parameters could roughly reach their expected structure.

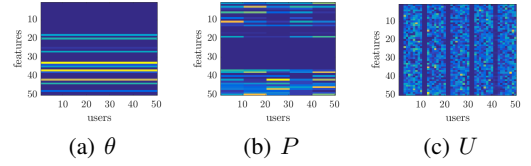


Figure 2: Figures for the magnitude of θ, P and U . For comparison convenience, both θ and P are extended to $\mathbb{R}^{d \times n_u}$ matrix

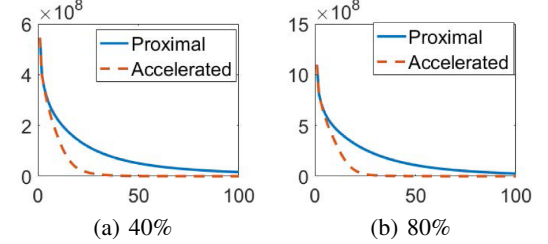


Figure 3: Effect of nestrov's Acceleration with (a) 40% random selected samples as training data and (b) 80% samples as training data. The y axis represents the average loss function of ($P1$) on 10 repetitions, and the x axis represents the iteration number.

Effect of Nestrov's Acceleration Strategy Next, we conduct empirical study the on the effect of Nestrov's Acceleration. To do this, we randomly select 40% and 80% samples as training data respectively. For each ratio, we run our algorithm 10 times for 100 iterations with different initial parameters. The average loss function per iteration both with and without the nestrov's acceleration strategy are presented in Figure 3. According to this figure, we conclude that, for both ratios, the accelerated algorithm starts to converge before the 50th iteration, which is much earlier than that of the ordinary proximal algorithm. We thus draw the conclusion that this acceleration strategy could successfully leverage faster convergence of our algorithm.

Shoes Dataset with Binary Attributes

Dataset Description For attribute learning, we use the shoes Dataset (Kovashka and Grauman 2013; Kovashka, Parikh, and Grauman 2015) which contains 14,658 online shopping images. Here we choose 6 user labeled attributes from the original dataset: bright, ornate, shiny, high, long, formal. For each of the attribute, user specific binary labels were collected on 60 images from 10 workers (Kovashka and Grauman 2013). In other words, we have $n_a = 6$, $n_{u_i} = 10$ and $n_{ij} = 60$. To form the feature of the images, we concatenate the GIST and color histograms provided by the original dataset.

Competitors and Evaluation Metric To show the effectiveness of our proposed algorithm on binary personal attribute learning, we compare our algorithm with three kinds of algorithms: global algorithms, user-specific algorithms and multi-task algorithms. For global algorithms, during

Table 3: Performance comparison for the binary attribute dataset

algorithm	accuracy	
	40%	80%
SVM	0.6278	0.6756
MLP	0.5057	0.4997
user exclusive	0.6549	0.6913
user adaptive	0.6771	0.6956
rMTFL-G	0.6421	0.6941
rMTFL-U	0.6659	0.7056
ours	0.6894	0.7121

the training phase, user specific labels are first processed to global labels via majority voting. For valid set and test set the user specific labels are used directly for performance evaluation. This kind of algorithms include: the Support Vector Machine(SVM), and Multi Layer Perceptron with single hidden layer (MLP). For user specific algorithms, we employ the user exclusive method where one SVM classifier is trained for each user independently and the user adaptive method proposed in (Kovashka and Grauman 2013). For multi-task algorithms, we employ rMTFL (Gong, Ye, and Zhang 2012), which shares similar model assumption as the proposed model, as the benchmark, and both the global version (rMTFL-G) and the user specific version (rMTFL-U) are considered. The setting of rMTFL-G is the same as that of the other global algorithms, except that the global classifiers for different attributes are trained in a multi-task manner. While for rMTFL-U, we regard all the user-specific classifiers as subtasks of rMTFL. To evaluate the generalized performance of different algorithms, the average value of the classification accuracy among all users is adopted as the performance metric.

Results Table 3 shows the average performance on 5 splits for all these algorithms when 40% and 80% of the samples are chosen as training data respectively. We could observe that our proposed algorithm reaches the best average accu-

racy.

Shoes Dataset with Relative Attributes

Dataset Description Here we use the same shoes dataset as the binary attribute experiment. The only difference is that the user specific labels are collected based on relative attribute between a pair of shoes images. For this task, we choose 6 attributes : pointy, bright, ornate, shiny, sporty and feminine. Similar as the previous subsection, we have $n_a = 6$, $n_{u_i} = 10$, $n_{ij} = 60$.

Competitors and Evaluation Metric For relative attribute learning, we also adopt the aforementioned three kinds of algorithms as benchmarks. The global models include: ranksvm models for relative attribute learning (Parikh and Grauman 2011) (rel_attr), ranknet (Burges et al. 2005), and rankboost (Freund et al. 2003). For the user exclusive model we train one rel_attr model for each user independently. And user adaptive model is the same as (Kovashka and Grauman 2013) except that we do not use any extra data for pre-training. rMTFL-G and rMTFL-U are the same as that used in the previous subsection except that the input feature for an image pair is processed as what mentioned in the “Model Formulation” section to fit these algorithms to ranking problems. We use the average ranking accuracy among all users for all tasks as our evaluation metric.

Results Table 4 show the performance comparison when 40% and 80% of the samples are chosen as training data respectively. It is concluded that our proposed algorithm could reach reasonable improvements with respect to the benchmarks, which demonstrate its effectiveness.

Conclusion

In this paper, we propose a hierarchical multi-task model for user specific attribute learning across multiple attributes with a common-to-special decomposition of the model weights. Specifically, our model weights include a common cognition factor, an attribute-specific factor and a user specific factor. The well-known accelerated proximal gradient method is employed to solve this model. Based on assumption 1, we prove theoretically that the proposed algorithm could both leverage good performance and estimate the true parameters well with high probability. The experiment results further verify the effectiveness of our proposed model.

Acknowledgments

The research of Zhiyong Yang and Qingming Huang was supported in part by National Natural Science Foundation of China: 61332016, U1636214, 61650202 and 61620106009, in part by Key Research Program of Frontier Sciences, CAS: QYZDJ-SSW-SYS013. The research of Qianqian Xu was supported in part by National Natural Science Foundation of China (No. 61672514, 61390514, 61572042), CCF-Tencent Open Research Fun. The research of Xiaochun Cao was Supported by National Key Research and Development Plan (No.2016YFB0800603), National Natural Science Foundation of China (No.U1636214, 61422213).

Table 4: Performance comparison for the relative attribute dataset

algorithm	accuracy	
	40%	80%
rel_attr	0.4797	0.5195
RankNet	0.4791	0.4721
RankBoost	0.4669	0.5251
user exclusive	0.4753	0.5303
user adaptive	0.4777	0.5336
rMTFL-G	0.4807	0.5074
rMTFL-U	0.4838	0.5433
ours	0.5119	0.5546

References

- Ando, R. K., and Zhang, T. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research* 6:1817–1853.
- Beck, A., and Teboulle, M. 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences* 2(1):183–202.
- Bickel, P. J.; Ritov, Y.; and Tsybakov, A. B. 2008. Simultaneous analysis of lasso and dantzig selector. *Annals of Statistics* 37(4):1705–1732.
- Boyd, S., and Vandenberghe, L. 2004. *Convex Optimization*. Cambridge University Press.
- Burges, C. J. C.; Shaked, T.; Renshaw, E.; Lazier, A.; Deeds, M.; Hamilton, N.; and Hullender, G. N. 2005. Learning to rank using gradient descent. In *proceedings of the the Twenty-Second International Conference on Machine Learning*, 89–96.
- Chen, L.; Zhang, Q.; and Li, B. 2014. Predicting multiple attributes via relative multi-task learning. In *proceedings of the Conference on Computer Vision and Pattern Recognition, 2014*, 1027–1034.
- Chen, J.; Zhou, J.; and Ye, J. 2011. Integrating low-rank and group-sparse structures for robust multi-task learning. In *proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 42–50.
- Ehrlich, M.; Shields, T. J.; Almaev, T.; and Amer, M. R. 2016. Facial attributes classification using multi-task representation learning. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2016*, 752–760.
- Freund, Y.; Iyer, R. D.; Schapire, R. E.; and Singer, Y. 2003. An efficient boosting algorithm for combining preferences. *Journal of Machine Learning Research* 4:933–969.
- Gong, P.; Ye, J.; and Zhang, C. 2012. Robust multi-task feature learning. In *proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '12*, 895–903.
- Han, H.; Jain, A. K.; Shan, S.; and Chen, X. 2017. Heterogeneous face attribute estimation: A deep multi-task learning approach. *arxiv 1706.00906*.
- Hand, E. M., and Chellappa, R. 2017. Attributes for improved attributes: A multi-task network utilizing implicit and explicit relationships for facial attribute classification. In *proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 4068–4074.
- Hsieh, H.; Hsu, W. H.; and Chen, Y. 2017. Multi-task learning for face identification and attribute estimation. In *proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing, 2017*, 2981–2985.
- Huang, A.; Xu, L.; Li, Y.; and Chen, E. 2014. Robust dynamic trajectory regression on road networks: A multi-task learning framework. In *proceedings of the International Conference on Data Mining, 2014*, 857–862.
- Hunter, D. R. 2004. Mm algorithms for generalized bradley-terry models. *Annals of Statistics* 32(1):384–406.
- Ji, Z.; Sun, Y.; Yu, Y.; Guo, J.; and Pang, Y. 2017. Semantic softmax loss for zero-shot learning. *arxiv 1705.07692*.
- Kovashka, A., and Grauman, K. 2013. Attribute adaptation for personalized image search. In *proceedings of the International Conference on Computer Vision 2013*, 3432–3439.
- Kovashka, A., and Grauman, K. 2015. Discovering attribute shades of meaning with the crowd. *International Journal of Computer Vision* 114(1):56–73.
- Kovashka, A.; Parikh, D.; and Grauman, K. 2015. Whittle-search: Interactive image search with relative attribute feedback. *International Journal of Computer Vision* 115(2):185–210.
- Parikh, D., and Grauman, K. 2011. Relative attributes. In *proceedings of the IEEE International Conference on Computer Vision, 2011*, 503–510.
- Song, F.; Tan, X.; and Chen, S. 2014. Exploiting relationship between attributes for improved face verification. *Computer Vision and Image Understanding* 122:143–154.
- Su, C.; Zhang, S.; Xing, J.; Gao, W.; and Tian, Q. 2016. Deep attributes driven multi-camera person re-identification. In *proceedings of the 14th European Conference on Computer Vision, 2016*, 475–491.
- Su, C.; Zhang, S.; Yang, F.; Zhang, G.; Tian, Q.; Gao, W.; and Davis, L. S. 2017. Attributes driven tracklet-to-tracklet person re-identification using latent prototypes space mapping. *Pattern Recognition* 66:4–15.
- Wallace, D. L. 1959. Bounds on normal approximations to student’s and the chi-square distributions. *Annals of Mathematical Statistics* 30(4):1121–1130.
- Wang, Y.; Kwok, J. T.; Yao, Q.; and Ni, L. M. 2017. Zero-shot learning with a partial set of observed attributes. In *proceedings of the International Joint Conference on Neural Networks*, 3777–3784.
- Xu, L.; Huang, A.; Chen, J.; and Chen, E. 2015. Exploiting task-feature co-clusters in multi-task learning. In *proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 1931–1937.
- Zhang, Z., and Saligrama, V. 2016. Zero-shot learning via a joint latent similarity embedding. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 6034–6042.
- Zhou, J.; Chen, J.; and Ye, J. 2011. Clustered multi-task learning via alternating structure optimization. In *proceedings of the 25th Annual Conference on Neural Information Processing Systems*, 702–710.

Appendix

Proof of Propositions

Proof of Proposition 1. *The proof directly follows the operations that preserves the convexity of a function (Boyd and Vandenberghe 2004)* \square

Proof of Proposition 2. Denote Δ_{ij} as :

$$\Delta_{ij} = 2(\mathbf{X}^{(i,j)\top} \mathbf{X}^{(i,j)} \mathbf{w}^{(i,j)} - \mathbf{X}^{(i,j)\top} \mathbf{y}^{(i,j)}) \quad (17)$$

We could reach :

$$\nabla_{\boldsymbol{\theta}} L(\mathbf{W}) = \sum_{i=1}^{n_a} \sum_{j=1}^{n_{u_i}} \Delta_{i,j}$$

$$\nabla_{\mathbf{p}^{(i)}} L(\mathbf{W}) = \sum_{j=1}^{n_{u_i}} \Delta_{i,j}$$

$$\nabla_{\mathbf{u}^{(i,j)}} L(\mathbf{W}) = \Delta_{i,j}$$

We further denote

$$\begin{aligned} dL_{i,j} &= 2\mathbf{X}^{(i,j)\top} \mathbf{X}^{(i,j)} (\mathbf{w}^{(i,j)} - \mathbf{w}'^{(i,j)}) \\ d\boldsymbol{\theta} &= \nabla_{\boldsymbol{\theta}} L(\mathbf{W}) - \nabla_{\boldsymbol{\theta}'} L(\mathbf{W}') \\ d\mathbf{p}^{(i)} &= \nabla_{\mathbf{p}^{(i)}} L(\mathbf{W}) - \nabla_{\mathbf{p}'^{(i)}} L(\mathbf{W}') \\ d\mathbf{u}^{(i,j)} &= \nabla_{\mathbf{u}^{(i,j)}} L(\mathbf{W}) - \nabla_{\mathbf{u}'^{(i,j)}} L(\mathbf{W}') \end{aligned}$$

It thus follows that

$$\begin{aligned} & \|\nabla L(\tilde{\mathbf{W}}) - \nabla L(\tilde{\mathbf{W}}')\| \\ &= \sqrt{\|d\boldsymbol{\theta}\|^2 + \sum_i \|d\mathbf{p}^{(i)}\|^2 + \sum_i \sum_j \|d\mathbf{u}^{(i,j)}\|^2} \\ &\stackrel{(I)}{\leq} \|d\boldsymbol{\theta}\| + \sum_i \|d\mathbf{p}^{(i)}\| + \sum_i \sum_j \|d\mathbf{u}^{(i,j)}\| \\ &= \left\| \sum_{i=1}^{n_a} \sum_{j=1}^{n_{u_i}} dL_{i,j} \right\| + \sum_{i=1}^{n_a} \left(\left\| \sum_{j=1}^{n_{u_i}} dL_{i,j} \right\| \right) + \sum_{i=1}^{n_a} \sum_{j=1}^{n_{u_i}} \|dL_{i,j}\| \\ &\leq 6C \sum_{i=1}^{n_a} \sum_{j=1}^{n_{u_i}} \|(\mathbf{w}^{(i,j)} - \mathbf{w}'^{(i,j)})\| \\ &\leq 6C \left(\sum_i n_{u_i} \|\mathbf{p}^{(i)} - \mathbf{p}'^{(i)}\| + n_u \|\boldsymbol{\theta} - \boldsymbol{\theta}'\| \right. \\ &\quad \left. + \sum_{i,j} \|\mathbf{u}^{(i,j)} - \mathbf{u}'^{(i,j)}\| \right) \\ &\stackrel{(II)}{\leq} 6C n_u \sqrt{n_a + 1 + n_u} \|\tilde{\mathbf{W}} - \tilde{\mathbf{W}}'\| \end{aligned} \tag{18}$$

where $C = \max_{i,j} [\sigma_1(\mathbf{X}^{(i,j)})]^2$; (I) follow that $\sqrt{\sum_i \mathbf{a}_i} \leq \sum_i \sqrt{\mathbf{a}_i}$, $\forall \mathbf{a} \in \mathbb{R}^n$; (II) is due to the fact that $\|\mathbf{a}\|_1 \leq \sqrt{n} \|\mathbf{a}\|_2$, $\forall \mathbf{a} \in \mathbb{R}^n$ \square

Proof of the main result

We first proof three fundamental lemmas which are the precursors of our main result.

Lemma 1. $\forall \mathbf{A}, \hat{\mathbf{A}} \in \mathbb{R}^{d \times m}$, $m \geq 1$, we have :

$$\|\hat{\mathbf{A}} - \mathbf{A}\|_\ell + \|\mathbf{A}\|_\ell - \|\hat{\mathbf{A}}\|_\ell \leq 2\|(\hat{\mathbf{A}} - \mathbf{A})^{\mathcal{N}_\perp(\mathbf{A})}\|_\ell$$

holds for all norm $\|\cdot\|_\ell$ such that $\forall \mathcal{C} \subset [d]$ and $\mathcal{C}_\perp = [d] \setminus \mathcal{C}$:

$$\|\mathbf{A}\|_\ell = \|\mathbf{A}^\mathcal{C}\|_\ell + \|\mathbf{A}^{\mathcal{C}_\perp}\|_\ell \tag{19}$$

Proof. Since $\mathbf{A}^{\mathcal{N}(\mathbf{A})} = 0$, we have : $\|(\hat{\mathbf{A}} - \mathbf{A})^{\mathcal{N}(\mathbf{A})}\|_\ell = \|(\hat{\mathbf{A}})^{\mathcal{N}(\mathbf{A})}\|_\ell$

Hence :

$$\begin{aligned} & \|(\hat{\mathbf{A}} - \mathbf{A})^{\mathcal{N}(\mathbf{A})}\|_\ell + \|\mathbf{A}\|_\ell - \|\hat{\mathbf{A}}\|_\ell \\ &\stackrel{(*)}{=} \|\mathbf{A}\|_\ell - \|\hat{\mathbf{A}}^{\mathcal{N}_\perp(\mathbf{A})}\|_\ell \stackrel{(**)}{\leq} \|(\hat{\mathbf{A}} - \mathbf{A})^{\mathcal{N}_\perp(\mathbf{A})}\|_\ell \end{aligned} \tag{20}$$

where $(*)$ follows Eq.(19) and $(**)$ follows the fact that $\|\mathbf{A}\| - \|\mathbf{B}\| \leq \|\mathbf{A} - \mathbf{B}\|$.

With Eq.(19) and Eq.(20), the correctness of lemma 1 is straightforward \square

Lemma 2. Let $\chi^2(d)$ be a chi-square random variable with degree of freedom d , then the following inequality holds :

$$\mathbb{P}(\chi^2(d) \geq d + t) \leq \frac{1}{\sqrt{2\pi}Z_d(t)} \exp\left(-\frac{Z_d(t)}{2}\right), \forall t > 0$$

where $Z_d(t) = t - d \log(1 + \frac{t}{d})$,

Proof. According to the Wallace inequality (Wallace 1959), we have :

$$\mathbb{P}(\chi^2(d) \geq d + t) \leq \mathbb{P}(\mathcal{N}_{0,1} \geq \sqrt{Z_d(t)}) \tag{21}$$

where $\mathcal{N}_{0,1}$ is a random variable subject to $\mathcal{N}(0, 1)$.

Moreover, let $u = \sqrt{Z_d(t)}$:

$$\begin{aligned} \mathbb{P}(\mathcal{N}_{0,1} \geq u) &= \int_u^{+\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx \\ &\leq \int_u^{+\infty} \frac{x}{u} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) dx \\ &= \frac{1}{\sqrt{2\pi}u} \exp\left(-\frac{u^2}{2}\right) \end{aligned} \tag{22}$$

It is obvious that lemma 2 directly follows (21) and (22). \square

Lemma 3. Let $\alpha = 2\sigma\sqrt{dn_u + t}$, choose $\lambda_1, \lambda_2, \lambda_3$ as : $\lambda_1 \geq n_u \alpha$, $\lambda_2 \geq \tilde{n} \alpha$, $\lambda_3 \geq \alpha$, where $n_u = \sum_{i=1}^{n_a} n_{u_i}$ and

$\tilde{n} = \sqrt{\sum_i n_{u_i}^2}$. let $\hat{\mathbf{W}} = (\hat{\boldsymbol{\theta}}, \hat{\mathbf{P}}, \hat{\mathbf{U}})$ be an optimal solution of

(P1), and $\mathbf{W} = (\boldsymbol{\theta}, \mathbf{P}, \mathbf{U})$ be an arbitrary feasible solution such that $\mathbf{W} \neq \hat{\mathbf{W}}$, then

$$\begin{aligned} & \sum_{i,j} \left\| \mathbf{x}^{(i,j)} \hat{\mathbf{w}}^{(i,j)} - \mathbf{f}^{*(i,j)} \right\|^2 \leq \sum_{i,j} \left\| \mathbf{x}^{(i,j)} \mathbf{w}^{(i,j)} - \mathbf{f}^{*(i,j)} \right\|^2 \\ & + 2\lambda_1 \left\| (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^{\mathcal{N}_\perp(\boldsymbol{\theta})} \right\|_1 + 2\lambda_2 \left\| (\hat{\mathbf{P}} - \mathbf{P})^{\mathcal{N}_\perp(\mathbf{P})} \right\|_{1,2} \\ & + 2\lambda_3 \left\| (\hat{\mathbf{U}}^\top - \mathbf{U}^\top)^{\mathcal{N}_\perp(\mathbf{U}^\top)} \right\|_{1,2} \end{aligned} \tag{23}$$

holds with probability at least:

$$1 - \frac{1}{\sqrt{2\pi}Z_{dn_u}(t)} \exp\left(-\frac{Z_{dn_u}(t)}{2}\right), \forall t > 0 \text{ where } n_u = \sum_{i=1}^{n_a} n_{u_i}$$

Proof. Since $\hat{\mathbf{W}}$ is an optimal solution of (P1), for any feasible solution $\mathbf{W} \neq \hat{\mathbf{W}}$, we have :

$$\begin{aligned} \sum_{i,j} \left\| \mathbf{x}^{(i,j)} \hat{\mathbf{w}}^{(i,j)} - \mathbf{y}^{(i,j)} \right\|^2 &\leq \sum_{i,j} \left\| \mathbf{x}^{(i,j)} \mathbf{w}^{(i,j)} - \mathbf{y}^{(i,j)} \right\|^2 \\ &+ \lambda_1 (\|\boldsymbol{\theta}\|_1 - \|\hat{\boldsymbol{\theta}}\|_1) + \lambda_2 (\|\mathbf{P}\|_{1,2} - \|\hat{\mathbf{P}}\|_{1,2}) \\ &+ \lambda_3 (\|\mathbf{U}^\top\|_{1,2} - \|\hat{\mathbf{U}}^\top\|_{1,2}) \end{aligned} \quad (24)$$

$$\text{Let } \Delta \mathbf{w}_{ij} = \underbrace{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})}_{\Delta \boldsymbol{\theta}} + \underbrace{(\hat{\mathbf{p}}^{(i)} - \mathbf{p}^{(i)})}_{\Delta \mathbf{p}^{(i)}} + \underbrace{(\hat{\mathbf{u}}^{(i,j)} - \mathbf{u}^{(i,j)})}_{\Delta \mathbf{u}^{(i,j)}},$$

according to Eq.(4), we have :

$$\begin{aligned} \sum_{i,j} \left\| \mathbf{x}^{(i,j)} \hat{\mathbf{w}}^{(i,j)} - \mathbf{f}^{*(i,j)} \right\|^2 &\leq \sum_{i,j} \left\| \mathbf{x}^{(i,j)} \mathbf{w}^{(i,j)} - \mathbf{f}^{*(i,j)} \right\|^2 \\ &+ \lambda_1 (\|\boldsymbol{\theta}\|_1 - \|\hat{\boldsymbol{\theta}}\|_1) + \lambda_2 (\|\mathbf{P}\|_{1,2} - \|\hat{\mathbf{P}}\|_{1,2}) \\ &+ \lambda_3 (\|\mathbf{U}^\top\|_{1,2} - \|\hat{\mathbf{U}}^\top\|_{1,2}) + 2 \underbrace{\sum_{i,j} \langle Z_{ij}, \Delta \mathbf{w}_{ij} \rangle}_{\Delta L} \end{aligned} \quad (25)$$

where $Z_{ij} = \mathbf{X}^{(i,j)\top} \delta^{(i,j)} \in \mathbb{R}^d$

Let $\mathbf{Z} = [Z_{11}, Z_{1,2}, \dots, Z_{n_a, n_{u_{n_a}}}]$, we could bound ΔL further as:

$$\begin{aligned} \Delta L &= 2 \sum_{i,j} \langle Z_{ij}, \Delta \boldsymbol{\theta} \rangle + 2 \sum_{i,j} \langle Z_{ij}, \Delta \mathbf{p}^{(i)} \rangle \\ &+ 2 \sum_{i,j} \langle Z_{ij}, \Delta \mathbf{u}^{(i,j)} \rangle \\ &\stackrel{(a)}{\leq} 2 \|\mathbf{Z}\|_F \left(n_u \|\Delta \boldsymbol{\theta}\|_1 + \tilde{n} \|\Delta \mathbf{P}\|_{1,2} + \|\Delta \mathbf{U}^\top\|_{1,2} \right) \end{aligned} \quad (26)$$

where (a) is the result of the fact that $\|x\|_2 \leq \|x\|_1, \forall x \in \mathbb{R}^n$ and the cauchy-schwarz inequality.

For $Z_{ij,k}$: the k th element of Z_{ij} , we have :

$$Z_{ij,k} = \sum_{l=1}^{n_{ij}} x_{lk}^{(i,j)} \delta_l^{(i,j)}$$

Based on Eq.(1), it is easy to show that $\frac{Z_{ij,k}^2}{\sigma^2} \sim \chi^2(1)$.

Hence, we get :

$$\frac{\|\mathbf{Z}\|_F^2}{\sigma^2} \sim \chi^2(dn_u)$$

where $n_u = \sum_{i=1}^{n_a} n_{u_i}$.

It follows lemma 2 that

$$\mathbb{P}[2\|\mathbf{Z}\|_F \leq \alpha] \geq 1 - \frac{1}{\sqrt{2\pi Z_{dn_u}(t)}} \exp\left(-\frac{Z_{dn_u}(t)}{2}\right), \forall t > 0$$

Suppose that $2\|\mathbf{Z}\|_F \leq \alpha$ holds, based on the requirement of this lemma and (26), we attain :

$$\Delta L \leq \lambda_1 \|\Delta \boldsymbol{\theta}\|_1 + \lambda_2 \|\Delta \mathbf{P}\|_{1,2} + \lambda_3 \|\Delta \mathbf{U}^\top\|_{1,2} \quad (27)$$

Now we could end the proof, since (23) follows (25), (27), and lemma 1. \square

Now we're ready to proof the main result based on assumption 1 and lemma 3.

Proof of Theorem 1. Let $\mathbf{w}^{(i,j)} = \mathbf{w}^{*(i,j)}$, with lemma3, if (23) holds, then we have :

$$\begin{aligned} \|\mathbf{X}\bar{\mathbf{W}} - \mathbf{F}\|^2 &\leq 2\lambda_1 \|(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^{\mathbb{N}_\perp(\boldsymbol{\theta}^*)}\|_1 \\ &+ 2\lambda_2 \|(\hat{\mathbf{P}} - \mathbf{P}^*)^{\mathbb{N}_\perp(\mathbf{P}^*)}\|_{1,2} \\ &+ 2\lambda_3 \|(\hat{\mathbf{U}}^\top - \mathbf{U}^{*\top})^{\mathbb{N}_\perp(\mathbf{U}^{*\top})}\|_{1,2} \end{aligned}$$

Take $\Gamma_{\theta^*} = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*$, $\Gamma_{P^*} = \hat{\mathbf{P}} - \mathbf{P}^*$, $\Gamma_{U^*} = \hat{\mathbf{U}} - \mathbf{U}^*$, $\Gamma_{\bar{\mathbf{W}}} = \bar{\mathbf{W}} - \bar{\mathbf{W}}^*$ according to assumption 1, we attain :

$$\begin{aligned} \|(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^{\mathbb{N}_\perp(\boldsymbol{\theta}^*)}\|_1 &\leq \sqrt{n_\theta} \|(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^{\mathbb{N}_\perp(\boldsymbol{\theta}^*)}\|_2 \\ &\leq \frac{\sqrt{n_\theta}}{\kappa_\theta \sqrt{n_{min} n_u}} \|\mathbf{X}\bar{\mathbf{W}} - \mathbf{F}\| \end{aligned} \quad (28)$$

Similarly, we could reach :

$$\|(\hat{\mathbf{P}} - \mathbf{P}^*)^{\mathbb{N}_\perp(\mathbf{P}^*)}\|_{1,2} \leq \frac{\sqrt{n_p}}{\kappa_p \sqrt{n_{min} n_u}} \|\mathbf{X}\bar{\mathbf{W}} - \mathbf{F}\| \quad (29)$$

$$\begin{aligned} \|(\hat{\mathbf{U}}^\top - \mathbf{U}^{*\top})^{\mathbb{N}_\perp(\mathbf{U}^{*\top})}\|_{1,2} \\ \leq \frac{\sqrt{n_{u,a}}}{\kappa_{u,a} \sqrt{n_{min} n_u}} \|\mathbf{X}\bar{\mathbf{W}} - \mathbf{F}\| \end{aligned} \quad (30)$$

Combining lemma3, (28)-(30), we could reach that (12) holds.

Following assumption 1, we get :

$$\|(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)\|_1 \leq (\beta_\theta + 1) \|(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^{\mathbb{N}_\perp(\boldsymbol{\theta}^*)}\|_1 \quad (31)$$

$$\|(\hat{\mathbf{P}} - \mathbf{P}^*)\|_{1,2} \leq (\beta_p + 1) \|(\hat{\mathbf{P}} - \mathbf{P}^*)^{\mathbb{N}_\perp(\mathbf{P}^*)}\|_{1,2} \quad (32)$$

$$\begin{aligned} \|(\hat{\mathbf{U}}^\top - \mathbf{U}^{*\top})\|_{1,2} \\ \leq (\beta_{u,a} + 1) \|(\hat{\mathbf{U}}^\top - \mathbf{U}^{*\top})^{\mathbb{N}_\perp(\mathbf{U}^{*\top})}\|_{1,2} \end{aligned} \quad (33)$$

Then (13)-(15) directly follows lemma 3 and (31)-(33). \square