

Hierarchical Gaussian Processes Model for Multi-task Learning

Ping Li, Songcan Chen*

*College of Computer Science and Technology, Nanjing University of Aeronautics and
Astronautics, Nanjing 210016, China
{ping.li.nj, s.chen}@nuaa.edu.cn*

Abstract

Multi-task learning (MTL) has been proved to improve performance of individual tasks by learning multiple related tasks together. Recently Nonparametric Bayesian Gaussian Process (GP) models have also been adapted to MTL and exhibit enough flexibility due to its non-parametric nature, thus can exempt from the assumption about the probability distributions of variables. To date, there have had two approaches proposed to implement GP-based MTL, i.e., cross-covariance-based and joint feature learning methods. Although successfully applied in scenarios such as face verification and collaborative filtering, these methods have their own drawbacks, for example, the cross-covariance-based method suffers from poor scalability because of the large covariance matrix involved; while the joint feature learning method can just implicitly incorporate relation between tasks, thus leading to a failure in explicitly exploiting the prior knowledge like correlation between tasks, which is crucial for further promoting MTLs. To address both issues, in this paper, we establish a two layer unified framework called Hierarchical Gaussian Process Multi-task Learning (HGPMT) method to jointly learn the latent shared features among tasks and a multi-task model. Furthermore, since the HGPMT does not need to involve the cross-covariance, its computational complexity is much lower. Finally, experimental results on both toy multi-task regression dataset and real datasets demonstrate

*Corresponding author
Email address: s.chen@nuaa.edu.cn (Songcan Chen)

its superiority in performance of multi-task learning to recently proposed approaches.

Keywords: GP-LVM, multi-task learning, feature learning, hierarchical model

1. Introduction

Multi-task learning (MTL) [1, 2, 3] is one of the main machine learning paradigms and can promote the improvement of model generalization performance by learning multiple related tasks together, especially for small number
5 of samples in each task [4, 5]. In MTL, sharing information among related tasks plays a crucial role and ensures that individual tasks can learn complementary knowledge to achieve the improvement of their own performance. To date, many methods for information sharing have been proposed and can roughly be divided into four categories: structure-sharing based [6, 7, 8], regularization-based
10 [9, 10, 11], task clustering based [4, 12, 13, 14] and multi-task feature learning based [3, 5, 15, 16, 17, 18]. Structure-sharing based methods assume that related tasks can be learned jointly by sharing certain hidden structures. These methods are well suited for multi-task learning since the structures learned from one task may be useful for other tasks [6]. Regularization-based methods, which
15 use certain regularizers to enforce the model parameters of each task to be close, can also be considered as sharing information among task-specific parameters. Task clustering based methods try to encode the task correlation knowledge of the task structure into the modeling or learn this knowledge directly from data [14]. Thus, regularization-based methods can be used to enforce the parameters
20 of dependent tasks to be as close to each other as possible. Multi-task feature learning (or selection) based method is also a widely used information sharing strategy which learns a set of shared features among related tasks. These features can be a subset of the original input features from feature selection [15, 16, 17] or newly formed features via feature transformation or learning [3, 5].

25 In their practical implementation, the above four methods can accommodate various base learners such as linear model [2], neural network [6], support

vector machine [9], Gaussian Process [14] and so on, according to the specific application scenarios. Among these base learners, *Gaussian Process* (GP) [19], as a flexible non-parametric method, has been shown to provide state-of-the-art performance in supervised learning task [20, 21] and widely used in many applications such as computer vision [22], recommendation system [23], geostatistics [24] and so on. Being a kernel based method, GP can effectively carry out non-linear learning with specific non-linear kernels. In addition to being highly accurate, it can also give an estimation of the prediction uncertainty which provides more information for the follow-up tasks [25]. Thus the MTL methods [14, 17, 26, 27] that adopt GPs as base learners both naturally inherit these properties and promote the performance of MTL. These GP-based MTL methods mainly involve the multi-task feature learning based [17, 26] and task clustering based methods (or cross-covariance-based methods) [14, 27]. However, there still exist some real-world problems that restrict their applicability. For example, in some application scenarios where some information (such as task-descriptor features, tasks similarities) about the explicit correlation among dependent tasks is available, the multi-task feature learning based methods can not explicitly exploit the information by just implicitly incorporating correlation among tasks, leading to a suboptimal solution. In other MTL application scenarios where there are a large number of dependent tasks, the clustering based GP methods will suffer from poor scalability because of the computation of large covariance matrix with the size dependent on the number of tasks and the number of the samples [28]. Although using a low-rank approximation for the task covariance matrix can improve scalability, the resulting expressive power may be impaired.

To address the problems described above, in this paper, we construct a two layer hierarchical gaussian process (HGPM) consisting of two gaussian (sub-) processes. HGPM jointly learns the latent shared features among tasks while explicitly utilizes the task relation information to improve the MTL performance. Specifically, the first layer i.e., a gaussian subprocess, of the HGPM is shared among all dependent tasks and used to capture the shared features. Us-

ing this structure, our HGPMT can fuse the information from dependent tasks and learn more representative features for the follow-up learning process. Then, the outputs of the first layer subprocess are transferred into the second layer one (corresponding to an output layer) to predict the target values of individual tasks. In this output layer, each node associates with an individual task with private parameters. Based on this hierarchical structure, we can embed explicit correlation prior among tasks by regularizing the parameters of the second layer. Furthermore, throughout the entire learning procedure, we need not construct and compute the cross-covariance matrix involved in cross-covariance-based methods. Thus, the computation complexity of HGPMT is much lower than the methods in [14, 27]. Finally we demonstrate the performance of HGPMT by comparing it with other related GP-based MTL methods on both synthetic and real-world datasets.

2. Related work

In this section, we give an introduction to the *Gaussian Process Regression* (GPR) [19] and *Gaussian Process Latent Variable Model* (GPLVM) [29] which are used in our MTL method. We also review the main existing MTL methods based on GP and discuss their disadvantages for specific applications.

2.1. GPR and GPLVM

GP models a finite set of random function variables $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^T$ as a joint gaussian distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix \mathbf{K} , if the function f has a GP prior, $f \sim \mathcal{GP}(\boldsymbol{\mu}, \mathbf{K})$, where in many cases we can specify a zero mean ($\boldsymbol{\mu} = \mathbf{0}$) and a kernel matrix \mathbf{K} (with hyper-parameter γ) as covariance matrix. Given a training dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ of N training samples, where $\mathbf{x}_n \in \mathbb{R}^D$ and $y_n \in \mathbb{R}$ denote n^{th} input variable and the corresponding continuous response variable respectively, our goal is to predict the response y^* of a new input \mathbf{x}^* based on this training dataset. In GPR, we model the response variable y_n as a noise-version of function value $f(\mathbf{x}_n)$,

$y_n \sim \mathcal{N}(f(\mathbf{x}_n), \sigma^2)$, where the distribution of noise ϵ_n is gaussian $\mathcal{N}(0, \sigma^2)$ with variance σ^2 . From the above definition, we can know that the joint probability of the response variables and latent function variables is $p(\mathbf{y}, \mathbf{f}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f})$. Thus, using bayesian theorem, we can obtain the distribution of the latent function value f^* is a gaussian distribution with mean $\mu(\mathbf{x}^*)$ and variance $var(\mathbf{x}^*)$:

$$\begin{aligned}\mu(\mathbf{x}^*) &= \mathbf{k}_{\mathbf{x}^* \mathbf{X}} (\mathbf{K}_{\mathbf{X} \mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \mathbf{y} \\ var(\mathbf{x}^*) &= k_{\mathbf{x}^* \mathbf{x}^*} - \mathbf{k}_{\mathbf{x}^* \mathbf{X}} (\mathbf{K}_{\mathbf{X} \mathbf{X}} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_{\mathbf{X} \mathbf{x}^*}\end{aligned}\tag{1}$$

where $\mathbf{k}_{\mathbf{x}^* \mathbf{X}}$ denotes the covariance between \mathbf{x}^* and the N training samples ($\mathbf{k}_{\mathbf{x}^* \mathbf{X}} = \mathbf{k}_{\mathbf{X} \mathbf{x}^*}^T$), $\mathbf{K}_{\mathbf{X} \mathbf{X}}$ denotes the covariance matrix of the N training samples.

Besides been widely used in regression, GP can also be applied in the non-linear dimensionality reduction models which are called *Gaussian Process Latent Variable Model* (GPLVM) [29, 21, 30]. It assumes that we have a observed data matrix $\mathbf{Y} \in \mathbb{R}^{N \times D}$ where N and D are the number and dimensionality of samples, respectively. As a generative model, GP-LVM assumes that the observed data is generated from a set of low dimensional variables $\mathbf{X} \in \mathbb{R}^{N \times Q}$ ($Q \ll D$). The generation process of the y_{nd} is $y_{nd} = f_d(\mathbf{x}_n) + \epsilon_{nd}$, where y_{nd} is the d^{th} dimension of the n^{th} sample. ϵ_{nd} is the noise with gaussian prior $\epsilon_{nd} \sim \mathcal{N}(0, \sigma^2)$. f_d is a nonlinear function with GP prior $f_d \sim \mathcal{GP}(\mathbf{0}, \mathbf{K})$. Thus, we can get the marginal likelihood $p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})$ by using bayesian theorem and integrating out f_d ,

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \prod_{d=1}^D \frac{1}{(2\pi)^{\frac{1}{2}} |\mathbf{K} + \sigma^2 \mathbf{I}|^{\frac{1}{2}}} e^{-\frac{1}{2} \mathbf{y}_{:,d}^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}_{:,d}}\tag{2}$$

where $\boldsymbol{\theta}$ denotes the hyper-parameters of kernel function and noise. $\mathbf{y}_{:,d}$ denotes the d^{th} column of matrix \mathbf{Y} . The goal of GP-LVM is to learn a low dimensional representation \mathbf{X} of high dimensional data \mathbf{Y} . Thus, we can maximize the marginal likelihood with respect to \mathbf{X} and the hyper-parameter $\boldsymbol{\theta}$ to find the optimal value $\hat{\mathbf{X}}$ and $\hat{\boldsymbol{\theta}}$, as shown in eq.(3).

$$\{\hat{\mathbf{X}}, \hat{\boldsymbol{\theta}}\} = \arg \max_{\mathbf{X}, \boldsymbol{\theta}} p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})\tag{3}$$

110 *2.2. Gaussian Process for Multi-task Learning*

GP provides a flexible method for the bayesian nonparametric inference of nonlinear latent functions and has been shown to generalize well when the training data is limited [31]. It also has been used in building multi-task learning models and verified its superiority in application such as physiological time-series analysis [32], facial expression recognition [33], emotion recognition [34] 115 and so on. A typical implementation of GP-based MTL is by forcing the hyper-parameters of multiple GPs to be the same [35] or to have the same prior [36], in this way, achieving the purpose of multi-task joint learning. However, this kind of method [35, 36] often fails to capture the relation among dependent tasks due 120 to its relative inflexibility [37].

Another important class of GP-based MTL method is joint feature learning based approach [17, 26] which assumes that correlation among tasks can be constructed by learning a set of shared features. For example, in [17], a GP-based MTL method is proposed by using an ARD kernel function [38] for 125 multi-task joint feature learning and taking into account the issue of shared hyper-parameters in [35]. It allows each task to have its own set of ARD hyper-parameters and share the sparsity patterns over the task specific hyper-parameters. In [26], a multi-kernel GP method is proposed and can also be considered as a joint feature learning method with task-specific weights. In 130 addition, this method has a strong relationship with generalized linear models, sparse factor analysis and matrix factorization, thus can be applied in dictionary learning and collaborative filtering. While these joint feature learning methods have obtained convincing results, they can just implicitly incorporate relation between tasks, thus leading to a failure in explicitly exploiting the prior 135 knowledge like correlation between tasks, which is crucial for further promoting MTLs.

In addition to the above GP-based MTL methods, researchers also explored the cross-covariance-based MTL methods [14, 27], in which the covariance matrix is constructed across all samples and task pairs. Specifically, in [14], the 140 covariance matrix of multi-task GP is modeled as $\mathbf{K} = \mathbf{K}^f \otimes \mathbf{K}^x$, where \otimes

denotes the *Kronecker* product, \mathbf{K}^f is a positive semi-definite (PSD) matrix that specifies the inter-task similarities, \mathbf{K}^x is a covariance function over inputs, $K_{lk}^f = k^f(\mathbf{t}_l, \mathbf{t}_k)$ the measurement of correlations between the l^{th} and the k^{th} task. Specifically, if task-descriptor feature \mathbf{t} is known, we can utilize this prior information to calculate the inter-task similarity matrix directly. Furthermore, we can also infer \mathbf{K}^f from data by defining it as a “free form” parameter being optimised. By incorporating the asymmetric dependency structures, Leen et al. [27] proposed a focused Gaussian Process model which introduces an “explaining away” model for each of the additional tasks to model their unrelated variation. Although these cross-covariance-based methods can utilize prior information of task similarities, they have been limited in applications by the computation burden [28]. We can highlight this drawback by the following prediction for a new sample \mathbf{x}^* ,

$$f_i(\mathbf{x}^*) = (\mathbf{k}_i^f \otimes \mathbf{k}_*^x)^T \boldsymbol{\Sigma}^{-1} \mathbf{y}, \quad \boldsymbol{\Sigma} = \mathbf{K}^f \otimes \mathbf{K}^x + \mathbf{D} \otimes \mathbf{I} \quad (4)$$

where \mathbf{k}_i^f is the l^{th} column of \mathbf{K}^f , \mathbf{k}_*^x is the vector covariances between the test point \mathbf{x}^* and the training points, \mathbf{D} is a $M \times M$ diagonal matrix whose $(l, l)^{th}$ element is σ_l^2 . As we can see from Eq.(4), inverting the $MN \times MN$ matrix in cross-covariance-based methods is quite time-consuming. Although some scalable optimization methods have been adopted, their loss in prediction accuracy is unavoidable [28].

In next section, we attempt to address the above problems by introducing a two layer hierarchical GP MTL method that combines the capabilities of non-linear feature learning and GP-based MTL method. This method enjoys two advantages: firstly, instead of learning MTL model on the original features, it uses the first layer gaussian subprocess to fuse the multiple task information and obtain more representative features. Secondly, it uses the second layer gaussian subprocess to construct the MTL model with a set of task-specific parameters. Thus, explicit correlation prior of dependent tasks can be utilized by using regularization based methods on these parameters.

3. The Proposed Method

3.1. Model Definition

In the multi-task learning scenarios, we assume that there are N training samples $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^T$, where $\mathbf{x}_n \in \mathbb{R}^D$ denotes the n^{th} input variable. The outputs of the T tasks are denoted by matrix $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_N]^T$, where $\mathbf{y}_n \in \mathbb{R}^T$ denotes the T tasks outputs that corresponded to the n^{th} input. In many cases, matrix \mathbf{Y} can be incomplete, meaning that for an input sample, we do not observe all the target values of the dependent tasks. The goal of MTL is to predict the target values of a new input sample based on the training dataset.

In this paper, we combine the feature learning and multi-task learning into a unified framework. The detailed generation processes of observed variables \mathbf{X} and \mathbf{Y} are shown in Eqs.(5) and (6). In this process, we construct the connection between feature learning and multi-task learning by using GPLVM. Specifically, we assume both \mathbf{X} and \mathbf{Y} are generated by a set of shared latent variables $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]$, where $\mathbf{z}_n \in \mathbb{R}^L$ is the latent variable corresponding to both the n^{th} input \mathbf{x}_n and output \mathbf{y}_n . The complete generation process mainly contains two parts: the generation of \mathbf{X} and the generation of \mathbf{Y} .

The generation of \mathbf{X} . In this paper, we use GPLVM to model the generation process of input variables \mathbf{X} . Following the definition of GPLVM, we assume the d^{th} dimension of the input \mathbf{x}_n is obtained by summing a non-linear function $\psi_d(\mathbf{z}_n)$ and a gaussian noise τ_{nd} , where $\tau_{nd} \sim \mathcal{N}(0, \epsilon^2)$. Similar to the GPLVM, we assume that each non-linear function ψ_d has an identical GP prior $\psi_d \sim \mathcal{GP}(\boldsymbol{\mu}_0, \mathbf{K}_0)$.

$$\psi_d \sim \mathcal{GP}(\boldsymbol{\mu}, \mathbf{K}_0), \quad d = 1, \dots, D \quad (5a)$$

$$x_{nd} \sim \mathcal{N}(\psi_d(\mathbf{z}_n), \epsilon^2), \quad n = 1, \dots, N \quad (5b)$$

With this approach, our method projects the data from their original spaces into a shared feature space, which can learn more representative features for the information fusion of dependent tasks. This kind of technique has been widely

used in many learning methods [8, 39, 40], which learns a set of feature vectors in
 190 a shared feature space across tasks or views for the further exploration. However,
 to the best of our knowledge, this is the first work that uses GPLVM for the
 multi-task feature learning. The excellent non-linear feature learning capability
 of GPLVM enables our method to learn better features for multi-tasks prediction
 as confirmed in the experimental section below.

The generation of \mathbf{Y} . From the above definition, we have constructed a
 shared feature learning structure. However, in supervised learning scenarios, a
 further challenge is how to improve the discriminative capabilities of the learned
 features for specific tasks. Accordingly, we pay more attention to the learning
 and selection of more discriminative features for multi-task fusion. Specifically,
 as show in Eq.(6), we define each output y_{nt} as the sum of a latent function
 $f_t(\mathbf{z}_n)$ and a gaussian noise v_{nt} , where f_t denotes the latent function associated
 with the t^{th} task and $v_{nt} \sim \mathcal{N}(0, \sigma_t^2)$. In the conventional GPLVM, the latent
 functions f_t ($t = 1, \dots, T$) are integrated out to take the full bayesian inference
 of the latent variables \mathbf{X} . However, in order to achieve supervised learning and
 multi-task information fusion, we assume that there exist M intermediate latent
 functions (or variables) $\{\phi_m\}_{m=1}^M$ and accordingly, define the t^{th} ($t = 1, \dots, T$)
 latent function f_t to be the weighted sum of the M nonlinear latent functions
 $\{\phi_m\}_{m=1}^M$ with the weight vector $\mathbf{w}_t \in \mathbb{R}^M$. Thus, we can jointly learn the
 T latent functions $\{f_t\}_{t=1}^T$ of dependent tasks by make they share a set of
 intermediate latent function $\{\phi_m\}_{m=1}^M$. Next following the idea of GPLVM,
 we impose a GP prior on the latent function $\phi_m \sim \mathcal{GP}(\boldsymbol{\mu}_m, \mathbf{K}_m)$ to derive a
 bayesian inference of the method.

$$\phi_m \sim \mathcal{GP}(\boldsymbol{\mu}_m, \mathbf{K}_m), \quad m = 1, \dots, M \quad (6a)$$

$$f_t = \sum_{m=1}^M w_{tm} \phi_m, \quad t = 1, \dots, T \quad (6b)$$

$$y_{nt} \sim \mathcal{N}(y_{nt} | f_t(\mathbf{x}_n), \sigma_t^2), \quad n = 1, \dots, N \quad (6c)$$

195 Through such a set of weight parameter $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_T]^T$ and interme-

diate latent functions $\{\phi_m\}_{m=1}^M$, we can effectively carry out MTL by learning these parameters and the latent features simultaneously. Furthermore, although we do not seek a deeper exploration on the dependence structure of the dependent tasks in this paper, some special regularization terms based on the priori information can also be imposed on \mathbf{W} . Objectively speaking, the generation process of \mathbf{Y} in our method is similar to multi-task multi-kernel learning method (MTMKL) in [26]. However, compared with the MTMKL which does not contain the learning of representative new features, HGPMT can jointly learn discriminative features and MTL prediction model, thus resulting in better prediction as shown in our experimental results.

3.2. Model Learning

In the model learning procedure, we wish to learn the parameters of the model and the latent variable \mathbf{Z} . In general, this can be achieved by maximizing the marginal likelihood $p(\mathbf{X}, \mathbf{Y} | \mathbf{Z}, \boldsymbol{\theta})$ where we use $\boldsymbol{\theta}$ to denote all the parameters involved in the model. However, since the generation of \mathbf{Y} is significantly different from the conventional GPLVM, the derivation of marginal likelihood is difficult. In this subsection, we give a detailed formulation of the marginal likelihood and the corresponding optimization algorithm.

From the generation process, we can obtain that \mathbf{X} and \mathbf{Y} are conditional independent given the latent variable \mathbf{Z} and hyper-parameters $\boldsymbol{\theta}$. Thus, we can write the joint marginal likelihood of \mathbf{X} and \mathbf{Y} as

$$\mathcal{L} = p(\mathbf{X}, \mathbf{Y} | \mathbf{Z}, \boldsymbol{\theta}) = p(\mathbf{X} | \mathbf{Z}, \boldsymbol{\theta}) p(\mathbf{Y} | \mathbf{Z}, \boldsymbol{\theta}) \quad (7)$$

By taking log of the Eq.(7), we get the log marginal likelihood,

$$\ln p(\mathbf{X}, \mathbf{Y} | \mathbf{Z}, \boldsymbol{\theta}) = \sum_{d=1}^D \ln p(\mathbf{x}_{:,d} | \mathbf{Z}, \boldsymbol{\theta}) + \sum_{t=1}^T \ln p(\mathbf{y}_{:,t} | \mathbf{Z}, \boldsymbol{\theta}) \quad (8)$$

where $\mathbf{x}_{:,d}$ and $\mathbf{y}_{:,t}$ denote the d^{th} column of \mathbf{X} and the t^{th} column of \mathbf{Y} respectively. Obviously, the computation of \mathcal{L} mainly contains two terms, $p(\mathbf{x}_{:,d} | \mathbf{Z}, \boldsymbol{\theta})$ and $p(\mathbf{y}_{:,t} | \mathbf{Z}, \boldsymbol{\theta})$. Since the first term can be computed analytically as the conven-

tional GPLVM, the primary challenge is the computation of $p(\mathbf{y}_{:,t}|\mathbf{Z},\boldsymbol{\theta})$ which involves a linear transformation of the latent functions $\{\phi_m\}_{m=1}^M$.

We follow the approach in GPR [19] and use the following formulation to calculate the distribution $p(\mathbf{y}_{:,t}|\mathbf{Z},\boldsymbol{\theta})$.

$$p(\mathbf{y}_{:,t}|\mathbf{Z},\boldsymbol{\theta}) = \int_{\mathbf{f}_t} p(\mathbf{y}_{:,t}|\mathbf{f}_t,\boldsymbol{\theta})p(\mathbf{f}_t|\mathbf{Z},\boldsymbol{\theta})d\mathbf{f}_t \quad (9)$$

225 where $\mathbf{f}_t = [f_{t1}, \dots, f_{tN}]$ denotes the latent variables associated with the t^{th} task.

We define matrices $\boldsymbol{\Phi} = [\phi_1, \dots, \phi_M]$ and $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_T]$. Since $\phi_m \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_0)$, we can conclude that the distribution of $\mathbf{f}_t = \boldsymbol{\Phi}\mathbf{w}_t$ is also a gaussian (The linear combination of gaussian distributed variables is itself gaussian). Thus, we can specify $p(\mathbf{f}_t)$ by finding its mean and covariance. As shown in
230 Eq.(10), the mean is $\mathbf{0}$, since $E[\phi_m] = \mathbf{0}$ ($m = 1, \dots, M$).

$$\boldsymbol{\mu}_{\mathbf{f}_t} = \mathbb{E}[\mathbf{f}_t] = \mathbb{E}[\boldsymbol{\Phi}]\mathbf{w}_t = \mathbf{0} \quad (10)$$

The derivation of the covariance is difficult because of the linear transformation. In this paper, we use the following equation to approximate the covariance $\mathbb{E}[\mathbf{f}_t\mathbf{f}_t^T]$.

$$\mathbf{K}_{\mathbf{f}_t} = \mathbb{E}[\mathbf{f}_t\mathbf{f}_t^T] = \mathbb{E}[\boldsymbol{\Phi}\mathbf{w}_t\mathbf{w}_t^T\boldsymbol{\Phi}^T] \approx \mathbb{E}[\boldsymbol{\Phi}\mathbf{D}_t\boldsymbol{\Phi}^T] \quad (11)$$

In Eq.(11), we assume that matrix $\mathbf{w}_t\mathbf{w}_t^T$ can be approximated by a diagonal
235 matrix \mathbf{D}_t where the m^{th} diagonal element $d_{mm} = w_{tm}^2$. This formulation is reasonable as it can be considered that we have ignored the correlation among the columns of matrix $\boldsymbol{\Phi}$. As a result, we can obtain \mathbf{K}_t by the following equation.

$$\mathbf{K}_{\mathbf{f}_t} \approx \sum_{m=1}^M w_{tm}^2 \mathbb{E}[\phi_m\phi_m^T] = \sum_{m=1}^M w_{tm}^2 \mathbf{K}_m \quad (12)$$

where $\mathbb{E}[\phi_m\phi_m^T]$ is the covariance of gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{K}_m)$. From the
240 above formulation, $P(\mathbf{f}_t)$ is a gaussian distribution with mean $\mathbf{0}$ and covariance $\sum_{m=1}^M w_{tm}^2 \mathbf{K}_m$ as shown in Eq.(13),

$$p(\mathbf{y}_{:,t}|\mathbf{Z},\boldsymbol{\theta}) = \int_{\boldsymbol{\Phi}} p(\mathbf{y}_{:,t}|\boldsymbol{\Phi},\boldsymbol{\theta})p(\boldsymbol{\Phi}|\mathbf{Z},\boldsymbol{\theta})d\boldsymbol{\Phi} = \mathcal{N}(\mathbf{0}, \mathbf{C}_t) \quad (13)$$

where $\mathbf{C}_t = \sum_{m=1}^M w_{tm}^2 \mathbf{K}_m + \sigma_t^2 \mathbf{I}$. Similar to the derivation of GPLVM in Section 2.1, we can get $\ln p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\theta})$

$$\ln p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\theta}) = \sum_{d=1}^D \ln p(\mathbf{X}_{:,d}|\mathbf{Z}, \boldsymbol{\theta}) = \sum_{d=1}^D \ln \mathcal{N}(0, \mathbf{K}_0 + \epsilon^2 \mathbf{I}) \quad (14)$$

Substituting Eq.(13) and Eq.(14) into Eq.(8), we can obtain the log marginal
 245 likelihood \mathcal{L} directly.

With the above derivation, the HGPMT realize the simultaneous learning for multiple tasks by allowing these tasks to share a set of kernel functions with task-specific parameters. In fact, this method can be considered as a multi-kernel learning which has been proven perform better in learning non-stationary
 250 function [41]. It is worthy to point out that although a further exploration for the correlation information of dependent tasks is not made in this paper, we can likewise also use regularization method to improve the performance of HGPMT by imposing prior information of tasks relation on the parameters..

As the optimization process in GPR and GPLVM, the maximization of the
 255 log marginal likelihood can be done using efficient gradient-based optimization algorithms such as conjugate gradients [42]. By using the chain rule, the gradient $\partial \ln p(\mathbf{X}, \mathbf{Y}|\mathbf{Z}, \boldsymbol{\theta})/\partial \boldsymbol{\theta}$ can be computed analytically as

$$\frac{\partial \ln p(\mathbf{X}, \mathbf{Y}|\mathbf{Z}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{d=1}^D \frac{\partial \ln p(\mathbf{x}_{:,d}|\mathbf{Z}, \boldsymbol{\theta})}{\partial \mathbf{C}_0} \cdot \frac{\partial \mathbf{C}_0}{\partial \boldsymbol{\theta}} + \sum_{t=1}^T \frac{\partial \ln p(\mathbf{y}_{:,t}|\mathbf{Z}, \boldsymbol{\theta})}{\partial \mathbf{C}_t} \cdot \frac{\partial \mathbf{C}_t}{\partial \boldsymbol{\theta}} \quad (15)$$

where $\mathbf{C}_0 = \mathbf{K}_0 + \epsilon^2 \mathbf{I}$. From Eq.(8), there are two kinds of parameters to be learned, i.e., the hyper-parameters $\boldsymbol{\theta}$ (including the weight matrix \mathbf{W} , the kernel
 260 parameters $\{\mathcal{A}_m\}_{m=0}^M$ corresponding to the kernel matrices $\{\mathbf{K}_m\}_{m=0}^M$ and the noises parameters $\{\sigma_t\}_{t=1}^T$ and ϵ) and the latent variables \mathbf{Z} . In order to speed up the convergence, we first train a conventional GPLVM on the input variables \mathbf{X} and then use the learned latent variables to initial \mathbf{Z} . The optimization procedure is detailed in Algorithm 1 which involves two main steps: In the
 265 first step, maximizing \mathcal{L} with respect to \mathbf{W} , $\{\sigma_t\}_{t=1}^T$, ϵ and $\{\mathcal{A}_m\}_{m=0}^M$ for fixed \mathbf{Z} . Alternatively, in the second step, maximizing \mathcal{L} with respect to \mathbf{Z} for fixed parameters which have obtained from the first step. Such an alternating

procedure is terminated until that the variation of \mathcal{L} between two iterations is less than some predefined threshold.

Algorithm 1 Hyper-parameters and latent variables learning

Inputs: The training samples $\{\mathbf{X}, \mathbf{Y}\}$, the number of kernels M and the dimensionality of the latent variables L .

Outputs: The weight matrix \mathbf{W} , the noise parameters $\{\sigma_t\}_{t=1}^T$ and ϵ , the kernel parameters $\{\mathcal{A}_m\}_{m=0}^M$ and the latent variables \mathbf{Z} .

1: Initial parameters \mathbf{W} , $\{\sigma_t\}_{t=1}^T$, ϵ and $\{\mathcal{A}_m\}_{m=0}^M$ randomly.

2: Initial \mathbf{Z} by training a GPLVM on \mathbf{X} .

3: **repeat**

4: Maximize \mathcal{L} with respect to $\{\sigma_t\}_{t=1}^T$, ϵ and $\{\mathcal{A}_m\}_{m=0}^M$ for fixed \mathbf{Z}

$$\arg \max_{\{\sigma_t\}_{t=1}^T, \epsilon, \{\mathcal{A}_m\}_{m=0}^M} \ln p(\mathbf{X}, \mathbf{Y} | \mathbf{Z}, \{\sigma_t\}_{t=1}^T, \epsilon, \{\mathcal{A}_m\}_{m=0}^M)$$

5: Maximize \mathcal{L} with respect to \mathbf{Z} for fixed $\{\sigma_t\}_{t=1}^T$, ϵ and $\{\mathcal{A}_m\}_{m=0}^M$

$$\arg \max_{\mathbf{Z}} \ln p(\mathbf{X}, \mathbf{Y} | \mathbf{Z}, \{\sigma_t\}_{t=1}^T, \epsilon, \{\mathcal{A}_m\}_{m=0}^M)$$

6: **until** the variation of \mathcal{L} between two iterations is smaller than some predefined threshold.

7: **return** \mathbf{W} , $\{\sigma_t\}_{t=1}^T$, ϵ , $\{\mathcal{A}_m\}_{m=0}^M$ and \mathbf{Z}

270 In each iteration step above, the main computation cost of HGPMT is inverting the $N \times N$ total samples covariance matrix and its total computational complexity is $\mathcal{O}(N^3)$. However, the cross-covariance-based methods (Kronecker GP [14] and focus MT-GP [27]) require to invert the $T \times N$ covariance matrix, giving the complexity of $\mathcal{O}(T^3 N^3)$. Thus, the HGPMT has much lower computational complexity than the cross-covariance-based methods. For a comparison, 275 we also list the computational complexity of GPMTFS [17] and MTMKL [26] in Table 1, where N_{max} denotes the largest sample size in tasks $t \in 1, \dots, T$. Although GPMTFS and MTMKL have lower or same complexity, they fail to utilize the correlation information among dependent tasks and have a lower 280 accuracy than our HGPMT in the experimental section.

3.3. Prediction of New Samples

The goal of supervised learning is to predict the target value of a new input sample \mathbf{x}_* . For the HGPMT, its prediction process mainly contains two steps:

Table 1: Time complexity of the GP-based methods.

Methods	HGPMT	GPMTFS	Kronecker GP	MTMKL	focused MT-GP
Complexity	$\mathcal{O}(N^3)$	$\mathcal{O}(N_{max}^3)$	$\mathcal{O}(T^3N^3)$	$\mathcal{O}(N^3)$	$\mathcal{O}(T^3N^3)$

feature learning and prediction.

285 In the feature learning step, the algorithm uses the first layer of HGPMT to learn the latent representation \mathbf{z}_* of \mathbf{x}_* . Since we have learned the parameters of HGPMT, \mathbf{z}_* can be learned in the same way as in [29]. This is reasonable, since the generation process of the input \mathbf{x} is modeled as a complete GPLVM which can be extrapolated to new test samples. Thus, given a testing data \mathbf{x}_* ,
 290 we can derive the marginal likelihood of $p(\mathbf{x}_*|\mathbf{Z}, \mathbf{z}_*, \mathbf{X})$ from Eq.(1), then, \mathbf{z}_* can be optimized by maximizing $p(\mathbf{x}_*|\mathbf{Z}, \mathbf{z}_*, \mathbf{X})$ with gradients-based method.

In the prediction step, after learning the latent variable \mathbf{z}_* , we can utilize the prediction method used as in GPR to get the target values of each task. This process is similar to GPR prediction except that the ordinary kernel is replaced
 295 by a weighted combination of multiple common kernels. Given the estimated weight matrix \mathbf{W} , kernel parameters and and noise parameter, it is straightforward to find out the combined kernel and then predict the corresponding target values for latent variable \mathbf{z}_* . The whole prediction algorithm is detailed in Algorithm 2.

Algorithm 2 Prediction of new samples

Inputs: The training samples $\{\mathbf{X}, \mathbf{Y}\}$ and the latent variables \mathbf{Z} , the weight matrix \mathbf{W} , the noise parameters $\{\sigma_t\}_{t=1}^T$ and ϵ , the kernel parameters $\{\mathcal{A}_m\}_{m=0}^M$.

Outputs: the prediction mean $\mu(\mathbf{x}^*)_t$ and the variance $var(\mathbf{x}^*)_t$ of \mathbf{x}^* for $t = 1, \dots, T$.

1: Optimize for \mathbf{z}^* :

$$\arg \max_{\mathbf{z}^*} \ln p(\mathbf{x}_*|\mathbf{Z}, \mathbf{z}_*, \mathbf{X}, \{\sigma_t\}_{t=1}^T, \{\mathcal{A}_m\}_{m=1}^M)$$

2: Calculate the mean $\mu(\mathbf{x}^*)_t$ and the variance $var(\mathbf{x}^*)_t$ for $t = 1, \dots, T$

$$\begin{aligned} \mu(\mathbf{x}^*)_t &= \mathbf{k}_{f_t(\mathbf{z}^* \mathbf{Z})} (\sigma_t^2 \mathbf{I} + \mathbf{K}_{f_t(\mathbf{Z} \mathbf{Z})})^{-1} \mathbf{y}_t \\ var(\mathbf{x}^*)_t &= k_{f_t(\mathbf{z}^* \mathbf{z}^*)} - \mathbf{k}_{f_t(\mathbf{z}^* \mathbf{Z})} (\sigma_t^2 \mathbf{I} + \mathbf{K}_{f_t(\mathbf{Z} \mathbf{Z})})^{-1} \mathbf{k}_{f_t(\mathbf{Z} \mathbf{z}^*)} \end{aligned} \tag{16}$$

3: **return** $\mu(\mathbf{x}^*)_t$ and $var(\mathbf{x}^*)_t$ for $t = 1, \dots, T$.

300 4. Experiments and Analysis

In this section, we demonstrate the effectiveness of the proposed HGPMT method by comparing it with recently proposed related GP-based methods (Kronecker GP [14], focused MT-GP [27], GPMTFS [17] and MTMKL [26]) on both synthetic and real landmarks datasets. In the synthetic data test, we evaluate
305 the accuracy of HGPMT in learning multiple tasks together and its improvement by using task relation information. This can be demonstrated by comparing HGPMT with independent GPs which are trained separately on each task. We also conduct several experiments on three benchmark datasets and compare the result with the related GP models. Throughout our experiment, we use Squared
310 Exponential (SE) kernel in HGPMT, Kronecker GP, focused MT-GP, MTMKL methods and ARD kernel in GPMTFS method. During the experiment, since focused MTGP is a transfer machine learning method (where we should specify a primary task and a set of secondary tasks), we randomly choose a task from the training dataset as the primary task and the rest tasks as the secondary
315 tasks. Then, we use this setting to test the focus MTGP throughout the experiment. The performance measurement used in the experiment are Mean Square Error (MSE) and explained variance. Specifically, for the synthetic datasets, we use MSE to intuitively demonstrate the difference between the predicted value and observed value. For the real landmarks datasets, we measure the performance of the MTL methods using explained variance which reflects not only the
320 quality of the regression, but also the distribution of the independent variables.

4.1. Synthetic Data Test

To demonstrate the effectiveness of the HGPMT in multi-task learning, we construct an artificial dataset which contains 12 tasks. The correlations are built
325 into the data by jointly drawing samples of all tasks from the same gaussian process $\mathcal{GP}(\mathbf{0}, \mathbf{K}^f \otimes \mathbf{K}^x)$, where \mathbf{K}^x is a non-stationary kernel as shown in Eq.(17), \mathbf{K}^f a PSD matrix that specifies the inter-task similarities. We draw 81 samples for each task and then add gaussian noise to each sample. In order

to test the model on unobserved samples, we randomly draw 17 samples of a
 330 continuous region in each task as the test set. Since having already obtained
 the inter-task similarity matrix \mathbf{K}^f , we use this matrix to compute a Laplacian
 matrix and construct a graph regularization term on \mathbf{W} as mentioned in [43].
 We call such kind HGPMT as regularized-HGPMT method. Our intention
 here is to explore whether HGPMT and regularized-HGPMT can utilize the
 335 correlation among tasks and give better performance in those unobserved regions
 than GPs trained independently.

$$\mathcal{K}(x, x') = \exp\left(\frac{-x^2 - x'^2}{20}\right) (4\cos(0.5(x - x')) + \cos(2(x - x'))) \quad (17)$$

During the experimental process, we test the HGPMT and the regularized-
 HGPMT with different task numbers ($T \in \{2, 4, 6, 8, 10, 12\}$). For compari-
 son, we also train T independent Gaussian process regressors respectively cor-
 340 responding to T tasks. The HGPMT, the regularized-HGPMT and the inde-
 pendent GPs all adopt the SE kernel as their covariance function. Specifical-
 ly, 12 SE kernels are used in both HGPMT and regularized-HGPMT. Figure
 1 shows the synthetic samples (* and \circ denote the observed and unobserved
 samples respectively), the predicted means of HGPMT (the black solid lines),
 345 regularized-HGPMT(the green solid lines) and independent GPs (the dotted
 lines) for $T = 6$. As seen from the figure, in most cases, HGPMT, regularized-
 HGPMT and independent GPs can fit the data well to dense data region. How-
 ever, in sparse data region, HGPMT significantly outperforms independent G-
 PPs while regularized-HGPMT obtains the best result, indicating the HGPMT
 350 and the regularized-HGPMT can more sufficiently exploit the correlation a-
 mong tasks and make more accurate prediction on unobserved samples. More-
 over, regularized-HGPMT can further improve the performance by utilizing the
 known task relation information.

To further explore the performance, we compare the HGPMT, the regularized-
 355 HGPMT with the independent GPs by gradually increasing the number of tasks
 (from 2 to 12). The MSE of the each task is shown in Figure 2. As we can see, the
 MSEs of independent GPs remain unchanged with the increase of task number.

However, the MSEs of HGPMT and regularized-HGPMT decrease obviously with the increasing of the related task number. Moreover for varying number
 360 of related tasks, regularized-HGPMT always has a lower MSE than HGPMT, which demonstrates the effectiveness of using task relation information.

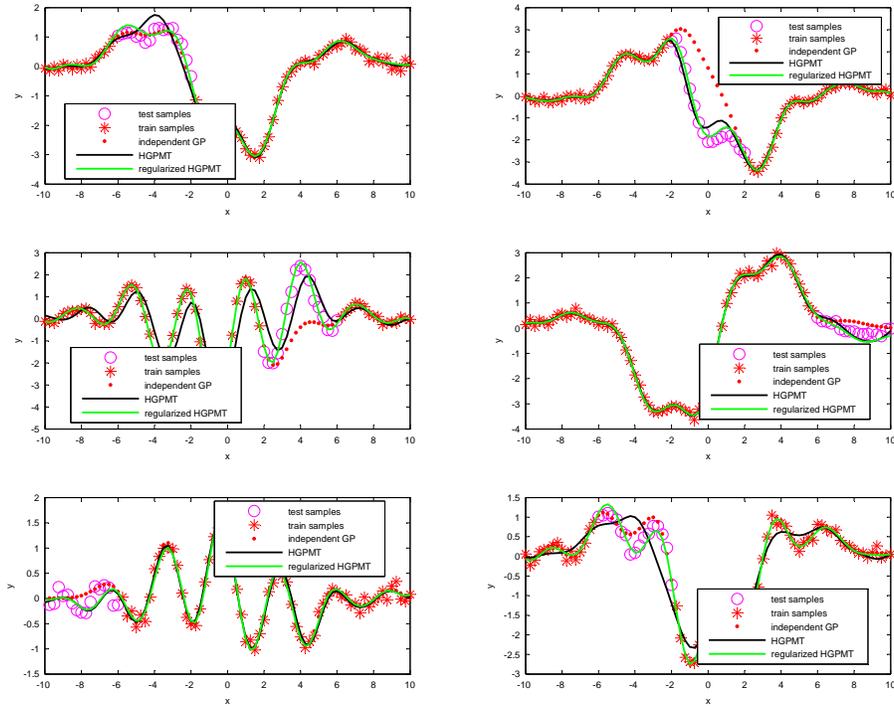


Figure 1: The predicted mean curves of HGPMT and independent GPs for $T = 6$. * and o denote observed and unobserved samples respectively. The black solid, green solid and dotted lines denote the predicted means of HGPMT, regularized-HGPMT and independent GPs, respectively.

4.2. School Data

School dataset [43] is a real-world multi-task dataset from Inner London
 365 Education Authority. This dataset contains 15362 examination records of students in 139 secondary schools from the year 1985 to 1987. It mainly consists two kinds of features: the student-related features (including year of the exam, gender, VR band and ethnic group) and the school-related features (percentage

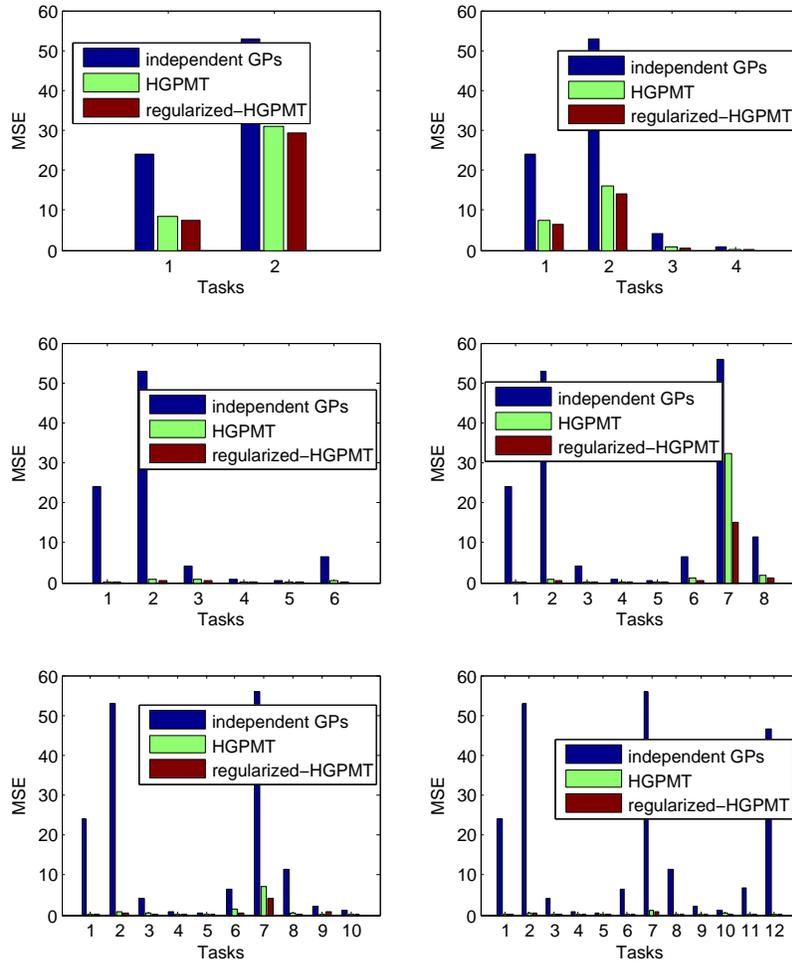


Figure 2: The MSEs of HGPMT, regularized-HGPMT and independent GPs for $T = \{2, 4, 6, 8, 10, 12\}$.

of students eligible for free school meals in the school, percentage of students in
 VR band one in the school, gender of the school (i.e. male, female, mixed) and
 school denomination). Our goal is to predict the exam score of a student that
 370 belongs to a specific school by considering each school as a task. It should be
 noted that for the dataset, we have access to the school-related features, thus
 can use these features to compute its Laplacian matrix which is in turn used to
 construct a graph regularization term as mentioned in [43]. With this method,
 375 we can utilize the tasks relation information to improve the performance of
 MTL. In the procedure of data pretreatment, we use the method described in
 [43] to create dummy variables for each discrete student-related feature. At last,
 we obtain 27 input features and then use them to train the models. In order
 to fully explore HGPMT, we compare it with closely related GP-based MTL
 380 methods with varying number of tasks and samples.

Firstly, we compare the HGPMT with the related models with task number-
 s ($T \in \{2, 50, 100, 139\}$). Specifically, we randomly draw T tasks and make 10
 random splits of the corresponding samples into training (75%) and test (25%).
 Then, we use these datasets to obtain the optimal hyper-parameters and eval-
 385 uate the performances of the models. The performances are measured by the
 explained variance (a widely used performance measure for multi-task learning).
 The explained variance is defined to be the total variance of the data minus the
 sum squared error on the test set as a percentage of the total data variance,
 which is a percentage version of the standard R^2 error measure for regression
 390 for the test data. Specifically, it is defined as

$$\begin{aligned} \textit{Explained variance} &= \frac{\textit{total variance} - \textit{sum of square errors}}{\textit{total variance}} \\ &= 1 - \textit{MSE}/\textit{variance} \end{aligned} \quad (18)$$

where *total variance* denotes the sum of the individual variances, *variance* =
 $\frac{\textit{total variance}}{N}$, $\textit{MSE} = \frac{\textit{sum of square errors}}{N}$ and N denotes the number of sam-
 ples. Obviously, the *explained variance* (18) is consistent with *MSE* in MTL
 learning methods. A lower *MSE* corresponds to a higher *explained variance*.
 395 Thus, a high value of explained variance is preferred over a low value. The re-

sult is shown in Figure 3. From Figure 3, we can conclude that the joint feature learning based methods (HGPMT and GPMTFS) have better MTL data fusion capability and can obtain more accurate results. As a transfer learning method, the focus MT-GP has a lower explained variance than other methods, since the correlation information among the secondary tasks is not fully utilized. We also
400 can see that the accuracies of all the models become higher with the increasing number of tasks. However, the HGPMT obtains higher explained variance than other models, implying that the HGPMT outperforms other compared methods. While the regularized-HGPMT shows the best result, demonstrating that
405 our HGPMT can make sufficient use of the correlation information of dependent tasks to improve the performance of MTL.

Secondly, we also compare the HGPMT with the related methods in performance by using varying size of training samples. Specifically, we randomly draw $h\%$ ($h \in \{20, 40, 60, 80, 100\}$) samples of the training set in each of the 10
410 random splits. Then, we take this subset and the corresponding test set as final train set and test set, respectively. Figure 4 shows how the training size affects the performance of the models. As we can see, when the number of samples in each task is relatively low, the explained variances of all the models show poor performances. As the training size increases, the performances of all the models
415 are accordingly improved. However, the HGPMT and the regularized-HGPMT are more improved in performance than other models, proving the effectiveness of our methods.

4.3. Personal Computer Data

Personal Computer Data is also a real word dataset including the ratings
420 of personal computers by students. Specifically, there are 190 students and 20 personal computers. Each computer contains 14 binary features (such as price, RAM, Hard Disk, Color of unit and so on). Students were asked to rate their likelihood of each computer. The likelihoods are denoted by integer number ranged from 0 to 10. Thus, we can treat students and computers as tasks
425 and samples, respectively. Although there is no task correlation information

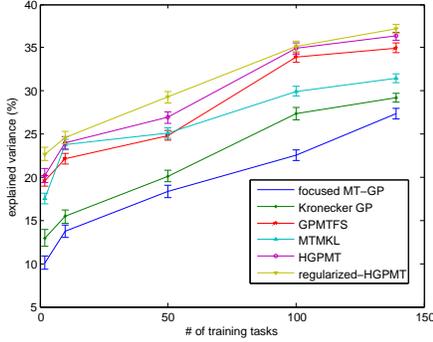


Figure 3: Explained variance vs varying number of training tasks on School Data.

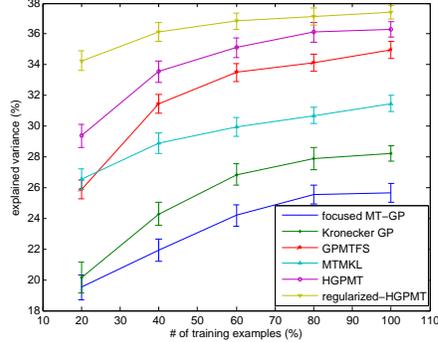


Figure 4: Explained variance vs varying number of training samples on School Data.

available for this dataset, our HGPMT method can still give more accurate results.

During the experimental process, we compare HGPMT with related models by using different numbers of tasks and training samples. The performance is also measured by using explained variance. Firstly, as described in Section 4.2, we randomly draw samples of $T \in \{2, 50, 100, 150, 190\}$ tasks and then make 10 random splits of the these samples into training dataset (80%) and test dataset (20%). Since each task has 20 samples, there are 16 and 4 samples in the training and test dataset for each task, respectively. The explained variances of HGPMT and related models are shown in Figure 5. It is worth noting that when the number of tasks is relatively low, all of the models have a poor performance, due to the small sample size problem. Since all the tasks in personal computer data are dependent, increasing the number of tasks can cause the improvement of all the models. Secondly, we also evaluate the performance of HGPMT with respect to the number of training samples. Specifically, we draw $h\%$ ($h = 20, 40, 60, 80, 100$) samples of the training samples and all the test samples as training set and test set. The explained variances of all models are shown in Figure 6. As we can see, HGPMT also gives the better result than other methods. It is worth to note that there are no task-related features in this dataset, thus we do not compare the regularized-HGPMT with the other

methods.

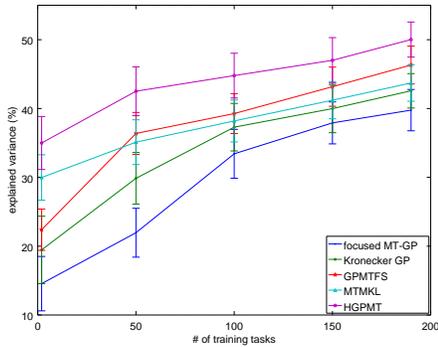


Figure 5: Explained variance vs varying number of training tasks on Personal Computer Data.

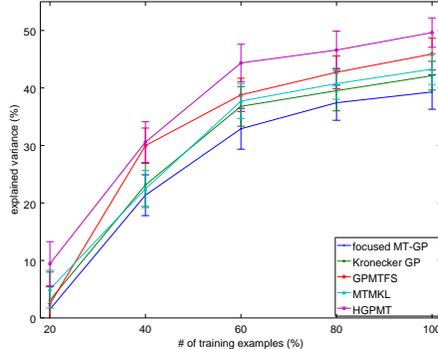


Figure 6: Explained variance vs varying number of training samples on Personal Computer Data.

4.4. SARCOS Data

The SARCOS data relates to an inverse dynamics problem for a seven degrees-of-freedom SARCOS anthropomorphic robot arm. Our goal is to map from a 21-dimensional input space (7 joint positions, 7 joint velocities, 7 joint accelerations) to the corresponding 7 joint torques. It is obvious that there are 7 tasks and 21 input features. In this experiment, we compare the HGPMT with independent GPs and MTMK method which is also a multi-kernel MTL method. Specifically, the experiment mainly includes two aspects: firstly, we explore the influence of kernel function number $M \in \{4, 9, 15\}$ and latent variable dimensionality $L \in \{10, 30\}$ on the performance of HGPMT; secondly, we validate whether HGPMT can learn better combination of multiple kernels and give better performance than MTMK. We randomly sampled 100, 200, 300, 400, 500 samples for training set and 2000 samples for test set. Then we perform 10 random splits of the data and run independent GPs, MTMK, HGPMT on these datasets. The mean and the standard deviation of the explained variances in the 10 trials are shown in Table 2. As we can see, the performances of both HGPMT and MTMK are affected by the number of kernel functions. The performance of HGPMT is also affected by the latent variable dimensionality. In

465 many cases, HGPMT can outperform the MTMK method. Furthermore, with the same number of kernels, our HGPMT has higher accuracy than MTMK, indicating that HGPMT can learn better combination of multiple kernels for multi-kernel learning.

Table 2: Comparison of MTL methods on the Sarcos dataset in terms of explained variance where $MTMK_j$ denotes MTMK method with $M = i$ and $HGPMT_{i,j}$ denotes HGPMT method with $M = i, L = j$.

	#100	#200	#300	#400	#500
<i>independent GPs</i>	0.5583±0.0053	0.5745±0.0015	0.5838±0.0017	0.5965±0.0013	0.6012±0.0012
$MTMK_4$	0.5990±0.0026	0.6746±0.0014	0.7218±0.0017	0.7300±0.0012	0.7636±0.0011
$MTMK_9$	0.6713±0.0039	0.7493±0.0016	0.7993±0.0018	0.8143±0.0012	0.8365±0.0011
$MTMK_{15}$	0.6767±0.0042	0.7548±0.0014	0.8064±0.0019	0.8217±0.0014	0.8428±0.0013
$HGPMT_{4,10}$	0.6923±0.0024	0.7536±0.0015	0.7941±0.0016	0.8354±0.0012	0.8436±0.0011
$HGPMT_{4,30}$	0.6836±0.0025	0.7651±0.0011	0.8197±0.0012	0.8567±0.0011	0.8655±0.0010
$HGPMT_{9,10}$	0.6836±0.0039	0.7754±0.0017	0.8382±0.0014	0.8491±0.0012	0.8656±0.0011
$HGPMT_{9,30}$	0.7092±0.0042	0.7956±0.0011	0.8598±0.0014	0.8650±0.0015	0.8750±0.0013
$HGPMT_{15,10}$	0.7127±0.0038	0.7927±0.0011	0.8597±0.0016	0.8622±0.0013	0.8750±0.0011
$HGPMT_{15,30}$	0.7073±0.0026	0.8100±0.0013	0.8649±0.0015	0.8700±0.0014	0.8922±0.0012

4.5. Comparison of time complexity

470 In Section 3.2, we have given the theoretical analysis of the computation complexity of all methods. In order to further evaluate the efficiency of HGPMT, we also conduct an empirical comparison among Kronecker GP [14], GPMTFS [17], MTMKL [26], focused MT-GP [27] and our HGPMT. Their training times on the School and Personal Computer data sets are respectively shown in Figures 7 and 8 from which we can witness that both Kronecker GP and focused
475 MT-GP are more time-consuming, especially for large scale problems (where the numbers of tasks and training samples are both large), while both HGPMT and regularized-HGPMT have comparable time complexity to GPMTFS and MTMKL. Furthermore, both HGPMT and regularized-HGPMT give higher ac-
480 curacy as shown in Sections 4.1 and 4.2. All these results verified our theoretical analysis in Section 3.2.

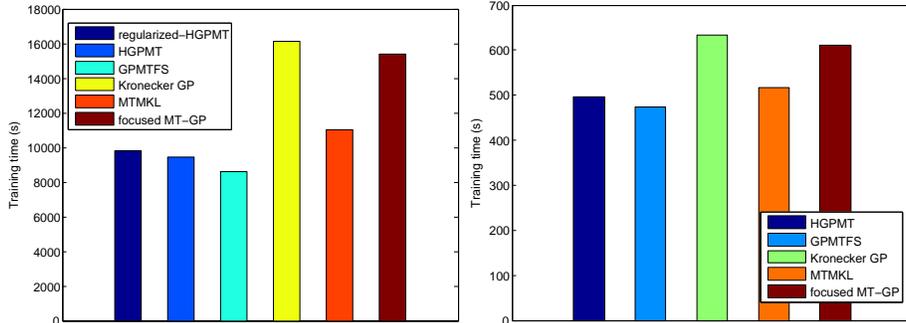


Figure 7: Training time comparison on School Data. Figure 8: Training time comparison on Personal Computer Data.

5. conclusion

In this paper, we establish a two-layer unified framework for jointly learning the latent shared features among tasks and the multi-task prediction model. As a multi-task learning method, this framework can embed task relation prior information into the task-specific parameter and effectively address the problem involving in the GP-based multi-task feature learning methods which are unable to handle such information. Furthermore, our method has much lower computational complexity comparing with the cross-covariance-based methods which can utilize the task relation information by constructing a large cross-covariance matrix. Experimental results show that this method has better multi-task learning performance and can also be considered as multi-kernel learning method which can efficiently learn the non-stationary functions. In addition, although we just use the graph-regularization-based method to embed task relation information into our modeling, there are also many other applicable methods, e.g., multi-task lasso [15], task clustering method [12]. And further exploring them to promote our HGPMT also deserves to be attempted.

References

- [1] M. L. Seltzer, J. Droppo, Multi-task learning in deep neural networks for improved phoneme recognition, in: Proceedings of the IEEE International

Conference on Acoustics, Speech and Signal Processing, 2013, pp. 6965–6969.

- [2] Y. Zhang, D.-Y. Yeung, A convex formulation for learning task relationships in multi-task learning, in: Proceedings of the 26th Conference on
505 Uncertainty in Artificial Intelligence, 2010, pp. 733–742.
- [3] P. Gong, J. Ye, C. Zhang, Robust multi-task feature learning, in: Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining, 2012, pp. 895–903.
- [4] C. Ciliberto, Y. Mroueh, T. Poggio, L. Rosasco, Convex learning of multiple tasks and their structure, in: Proceedings of the 32nd International
510 Conference on Machine Learning, 2015, pp. 1548–1557.
- [5] H. Fei, J. Huan, Structured feature selection and task relationship inference for multi-task learning, Knowledge and Information Systems 35 (2013) 171–180.
- [6] X. Chu, W. Ouyang, W. Yang, X. Wang, Multi-task recurrent neural network for immediacy prediction, in: Proceedings of the IEEE International
515 Conference on Computer Vision, 2015, pp. 3352–3360.
- [7] Y. Huang, W. Wang, L. Wang, T. Tan, Multi-task deep neural network for multi-label learning, in: Proceedings of the IEEE International Conference
520 on Image Processing, 2013, pp. 2897–2900.
- [8] Z. Zhang, P. Luo, C. C. Loy, X. Tang, Facial landmark detection by deep multi-task learning, in: Proceedings of the European Conference on Computer Vision, 2014, pp. 94–108.
- [9] T. Evgeniou, M. Pontil, Regularized multi-task learning, in: Proceedings of
525 the 10th ACM SIGKDD international conference on Knowledge discovery and data mining, 2004, pp. 109–117.

- [10] T. Kato, H. Kashima, M. Sugiyama, K. Asai, Multi-task learning via conic programming, in: Proceedings of the 20th International Conference on Neural Information Processing Systems, 2007, pp. 737–744.
- 530 [11] Y. Li, X. Tian, M. Song, D. Tao, Multi-task proximal support vector machine, Pattern Recognition 48 (2015) 3249–3257.
- [12] B. Bakker, T. Heskes, Task clustering and gating for bayesian multitask learning, Journal of Machine Learning Research 4 (2004) 83–99.
- [13] Y. Xue, X. Liao, L. Carin, B. Krishnapuram, Multi-task learning for classification with dirichlet process priors, Journal of Machine Learning Research 8 (2007) 35–63.
- 535 [14] E. V. Bonilla, K. M. A. Chai, C. K. Williams, Multi-task gaussian process prediction, in: Proceedings of the 20th International Conference on Neural Information Processing Systems, 2007, pp. 153–160.
- 540 [15] H. Liu, M. Palatucci, J. Zhang, Blockwise coordinate descent procedures for the multi-task lasso, with applications to neural semantic basis discovery, in: Proceedings of the 26th International Conference on Machine Learning, 2009, pp. 649–656.
- [16] Y. Zhang, D.-Y. Yeung, Q. Xu, Probabilistic multi-task feature selection, in: Proceedings of the 23rd International Conference on Neural Information Processing Systems, 2010, pp. 2559–2567.
- 545 [17] P. Srijiith, S. Shevade, Gaussian process multi-task learning using joint feature selection, in: Proceedings of the Joint European Conference on Machine Learning and Knowledge Discovery in Databases, 2014, pp. 98–113.
- 550 [18] J. He, J.-F. Hu, X. Lu, W.-S. Zheng, Multi-task mid-level feature learning for micro-expression recognition, Pattern Recognition 66 (2017) 44–52.

- [19] C. E. Rasmussen, C. K. Williams, Gaussian processes for machine learning, MIT press Cambridge, 2006.
- 555 [20] T. X. Lu C, Surpassing human-level face verification performance on lfw with gaussianface, in: Proceedings of the 29th AAAI Conference on Artificial Intelligence, 2015, pp. 3811–3819.
- [21] G. Song, S. Wang, Q. Huang, Q. Tian, Similarity gaussian process latent variable model for multi-modal data analysis, in: Proceedings of the IEEE
560 International Conference on Computer Vision, 2015, pp. 4050–4058.
- [22] C. Long, G. Hua, A. Kapoor, A joint gaussian process model for active visual recognition with expertise estimation in crowdsourcing, International Journal of Computer Vision 116 (2016) 1–25.
- [23] F. Bohnert, I. Zukerman, D. F. Schmidt, Using gaussian spatial processes to
565 model and predict interests in museum exhibits, in: Proceedings of the 7th International Conference on Intelligent Techniques for Web Personalization & Recommender Systems, 2009, pp. 13–19.
- [24] D. Nychka, S. Bandyopadhyay, D. Hammerling, F. Lindgren, S. Sain, A multiresolution gaussian process model for the analysis of large spatial
570 datasets, Journal of Computational & Graphical Statistics (2014) 579–599.
- [25] C. X. Li, B. Marlin, A scalable end-to-end gaussian process adapter for irregularly sampled time series classification, in: Proceedings of the 29th International Conference on Neural Information Processing Systems, 2016, pp. 1804–1812.
- 575 [26] M. K. Titsias, M. Lázaro-Gredilla, Spike and slab variational inference for multi-task and multiple kernel learning, in: Proceedings of the 24th International Conference on Neural Information Processing Systems, 2011, pp. 2339–2347.
- [27] G. Leen, J. Peltonen, S. Kaski, Focused multi-task learning in a gaussian
580 process framework, Machine learning 89 (2012) 157–182.

- [28] Y. Zhang, D.-Y. Yeung, Multi-task learning using generalized t process, in: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, 2010, pp. 964–971.
- [29] N. D. Lawrence, Gaussian process latent variable models for visualisation of high dimensional data, in: Proceedings of the 16th International Conference on Neural Information Processing Systems, 2003, pp. 329–336.
- [30] Z. Dai, J. Hensman, N. Lawrence, Spike and slab gaussian process latent variable models, arXiv preprint arXiv:1505.02434 (2015).
- [31] R. Urtasun, D. J. Fleet, A. Hertzmann, P. Fua, Priors for people tracking from small training sets, in: Proceedings of the 10th IEEE International Conference on Computer Vision, 2005, pp. 403–410.
- [32] R. Dürichen, M. A. F. Pimentel, L. Clifton, A. Schweikard, D. A. Clifton, Multitask gaussian processes for multivariate physiological time-series analysis, IEEE Transactions on Biomedical Engineering 62 (2015) 314–322.
- [33] O. Rudovic, M. Pantic, I. Patras, Coupled gaussian processes for pose-invariant facial expression recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence 35 (2013) 1357–1369.
- [34] S. Gómez-González, M. A. Álvarez, H. F. García, J. I. Ríos, A. A. Orozco, Discriminative training for convolved multiple-output gaussian processes, in: Proceedings of the Iberoamerican Congress on Pattern Recognition, 2015, pp. 595–602.
- [35] N. D. Lawrence, J. C. Platt, Learning to learn with the informative vector machine, in: Proceedings of the 21th International Conference on Machine Learning, 2004, pp. 65–73.
- [36] K. Yu, V. Tresp, A. Schwaighofer, Learning gaussian processes from multiple tasks, in: Proceedings of the 22nd international conference on Machine learning, 2005, pp. 1012–1019.

- [37] G. Steidl, Bayesian multitask classification with gaussian process priors, *IEEE Transactions on Neural Networks* 22 (2011) 2011–2021.
- 610 [38] J. Jacobs, Bayesian support vector regression with automatic relevance determination kernel for modeling of antenna input characteristics, *IEEE Transactions on Antennas and Propagation* 60 (2012) 2114–2118.
- [39] C. H. Ek, P. H. Torr, N. D. Lawrence, Gaussian process latent variable models for human pose estimation, in: *Proceedings of the 4th international conference on Machine learning for multimodal interaction, 2007*, pp. 132–143.
- 615 [40] S. Eleftheriadis, O. Rudovic, M. Pantic, Shared gaussian process latent variable model for multi-view facial expression recognition, in: *International Symposium on Visual Computing, 2013*, pp. 527–538.
- [41] S. Vasudevan, F. Ramos, E. Nettleton, H. Durrant-Whyte, Non-stationary dependent gaussian processes for data fusion in large-scale terrain modeling, in: *Proceedings of the IEEE International Conference on Robotics and Automation*, 2011, pp. 1875–1882.
- 620 [42] Fletcher, R., Reeves, M. C., Function minimization by conjugate gradients, *Computer Journal* 7 (1964) 149–154.
- 625 [43] T. Evgeniou, C. A. Micchelli, M. Pontil, Learning multiple tasks with kernel methods, *Journal of Machine Learning Research* 6 (2005) 615–637.