

---

# Parametric Task Learning

---

**Ichiro Takeuchi**

Nagoya Institute of Technology  
Nagoya, 466-8555, Japan  
takeuchi.ichiro@nitech.ac.jp

**Tatsuya Hongo**

Nagoya Institute of Technology  
Nagoya, 466-8555, Japan  
hongo.mllab.nit@gmail.com

**Masashi Sugiyama**

Tokyo Institute of Technology  
Tokyo, 152-8552, Japan  
sugi@cs.titech.ac.jp

**Shinichi Nakajima**

Nikon Corporation  
Tokyo, 140-8601, Japan  
nakajima.s@nikon.co.jp

## Abstract

We introduce an extended formulation of multi-task learning (MTL) called *parametric task learning (PTL)* that can systematically handle infinitely many tasks parameterized by a continuous parameter. Our key finding is that, for a certain class of PTL problems, the path of the optimal task-wise solutions can be represented as piecewise-linear functions of the continuous task parameter. Based on this fact, we employ a parametric programming technique to obtain the common shared representation across all the continuously parameterized tasks. We show that our PTL formulation is useful in various scenarios such as learning under non-stationarity, cost-sensitive learning, and quantile regression. We demonstrate the advantage of our approach in these scenarios.

## 1 Introduction

Multi-task learning (MTL) has been studied for learning multiple related tasks simultaneously. A key assumption behind MTL is that there exists a common shared representation across the tasks. Many MTL algorithms attempt to find such a common representation and at the same time to learn multiple tasks under that shared representation. For example, we can enforce all the tasks to share a common feature subspace or a common set of variables by using an algorithm introduced in [1, 2] that alternately optimizes the shared representation and the task-wise solutions.

Although the standard MTL formulation can handle only a finite number of tasks, it is sometimes more natural to consider infinitely many tasks parameterized by a continuous parameter, e.g., in *learning under non-stationarity* [3] where learning problems change over continuous time, *cost-sensitive learning* [4] where loss functions are asymmetric with continuous cost balance, and *quantile regression* [5] where the quantile is a continuous variable between zero and one. In order to handle these infinitely many parametrized tasks, we propose in this paper an extended formulation of MTL called *parametric-task learning (PTL)*.

The key contribution of this paper is to show that, for a certain class of PTL problems, the optimal common representation shared across infinitely many parameterized tasks can be obtainable. Specifically, we develop an alternating minimization algorithm à la [1, 2] for finding the entire continuum of solutions and the common feature subspace (or the common set of variables) among infinitely many parameterized tasks. Our algorithm exploits the fact that, for those classes of PTL problems, the path of task-wise solutions is piecewise-linear in the task parameter. We use the parametric programming technique [6, 7, 8, 9] for computing those piecewise linear solutions.

**Notations:** Let us denote by  $\mathbb{R}$ ,  $\mathbb{R}_+$ , and  $\mathbb{R}_{++}$  the set of real, nonnegative, and positive numbers, respectively, while we define  $\mathbb{N}_n := \{1, \dots, n\}$  for every natural number  $n$ . We denote by  $\mathcal{S}_{++}^d$  the set of  $d \times d$  positive definite matrices, and let  $I(\cdot)$  be the indicator function.

## 2 Review of Multi-Task Learning (MTL)

In this section, we review an MTL method developed in [1, 2]. Let  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathbb{N}_n}$  be the set of  $n$  training instances, where  $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^d$  is the input and  $y_i \in \mathcal{Y}$  is the output. We define  $w_i(t) \in [0, 1], t \in \mathbb{N}_T$  as the weight of the  $i^{\text{th}}$  instance for the  $t^{\text{th}}$  task, where  $T$  is the number of tasks. We consider an affine model  $f_t(\mathbf{x}) = \beta_{t,0} + \beta_t^\top \mathbf{x}$  for each task, where  $\beta_{t,0} \in \mathbb{R}$  and  $\beta_t \in \mathbb{R}^d$ . For notational simplicity, we define augmented vectors  $\tilde{\beta} := (\beta_0, \beta_1, \dots, \beta_d)^\top \in \mathbb{R}^{d+1}$  and  $\tilde{\mathbf{x}} := (1, x_1, \dots, x_d)^\top \in \mathbb{R}^{d+1}$ , and write the affine model as  $f_t(\mathbf{x}) = \tilde{\beta}_t^\top \tilde{\mathbf{x}}$ .

The multi-task feature learning method discussed in [1] is formulated as

$$\min_{\substack{\{\tilde{\beta}_t\}_{t \in \mathbb{N}_T} \\ D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1}} \sum_{t \in \mathbb{N}_T} \sum_{i \in \mathbb{N}_n} w_i(t) \ell_t(r(y_i, \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i)) + \frac{\gamma}{T} \sum_{t \in \mathbb{N}_T} \beta_t^\top D^{-1} \beta_t, \quad (1)$$

where  $\text{tr}(D)$  is the trace of  $D$ ,  $\ell_t : \mathbb{R} \rightarrow \mathbb{R}_+$  is the loss function for the  $t^{\text{th}}$  task incurred on the residual  $r(y_i, \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i)$ <sup>1</sup>, and  $\gamma > 0$  is the regularization parameter<sup>2</sup>. It was shown [1] that the problem (1) is equivalent to

$$\min_{\{\tilde{\beta}_t\}_{t \in \mathbb{N}_T}} \sum_{t \in \mathbb{N}_T} \sum_{i \in \mathbb{N}_n} w_i(t) \ell_t(r(y_i, \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i)) + \frac{\gamma}{T} \|B\|_{\text{tr}}^2,$$

where  $B$  is the  $d \times T$  matrix whose  $t^{\text{th}}$  column is given by the vector  $\beta_t$ , and  $\|B\|_{\text{tr}} := \text{tr}((BB^\top)^{1/2})$  is the *trace norm* of  $B$ . As shown in [10], the trace norm is the convex upper envelope of the rank of  $B$ , and (1) can be interpreted as the problem of finding a common feature subspace across  $T$  tasks. This problem is often referred to as *multi-task feature learning*. If the matrix  $D$  is restricted to be diagonal, the formulation (1) is reduced to *multi-task variable selection* [11, 12].

In order to solve the problem (1), the *alternating minimization algorithm* was suggested in [1] (see Algorithm 1). This algorithm alternately optimizes the task-wise solutions  $\{\tilde{\beta}_t\}_{t \in \mathbb{N}_T}$  and the common representation matrix  $D$ . It is worth noting that, when  $D$  is fixed, each  $\tilde{\beta}_t$  can be independently optimized (Step 1). On the other hand, when  $\{\tilde{\beta}_t\}_{t \in \mathbb{N}_T}$  are fixed, the optimization of the matrix  $D$  can be reduced to the minimization over  $d$  eigenvalues  $\lambda_1, \dots, \lambda_d$  of the matrix  $C := BB^\top$ , and the optimal  $D$  can be analytically computed (Step 2).

## 3 Parametric-Task Learning (PTL)

We consider the case where we have infinitely many tasks parametrized by a single continuous parameter. Let  $\theta \in [\theta_L, \theta_U]$  be a continuous task parameter. Instead of the set of weights  $w_i(t), t \in \mathbb{N}_T$ , we consider a weight function  $w_i : [\theta_L, \theta_U] \rightarrow [0, 1]$  for each instance  $i \in \mathbb{N}_n$ . In PTL, we learn a parameter vector  $\tilde{\beta}_\theta \in \mathbb{R}^{d+1}$  as a continuous function of the task parameter  $\theta$ :

$$\min_{\substack{\{\tilde{\beta}_\theta\}_{\theta \in [\theta_L, \theta_U]} \\ D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1}} \int_{\theta_L}^{\theta_U} \sum_{i \in \mathbb{N}_n} w_i(\theta) \ell_\theta(r(y_i, \tilde{\beta}_\theta^\top \tilde{\mathbf{x}}_i)) d\theta + \gamma \int_{\theta_L}^{\theta_U} \beta_\theta^\top D^{-1} \beta_\theta d\theta, \quad (2)$$

where, note that, the loss function  $\ell_\theta$  possibly depends on  $\theta$ .

As we will explain in the next section, the above PTL formulation is useful in various important machine learning scenarios including learning under non-stationarity, cost-sensitive learning, and

<sup>1</sup>For example,  $r(y_i, \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i) = (y_i - \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i)^2$  for regression problems with  $y_i \in \mathbb{R}$ , while  $r(y_i, \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i) = 1 - y_i \tilde{\beta}_t^\top \tilde{\mathbf{x}}_i$  for binary classification problems with  $y_i \in \{-1, 1\}$ .

<sup>2</sup>In [1],  $w_i(t)$  takes either 1 or 0. It takes 1 only if the  $i^{\text{th}}$  instance is used in the  $t^{\text{th}}$  task. We slightly generalize the setup so that each instance can be used in multiple tasks with different weights.

---

**Algorithm 1** ALTERNATING MINIMIZATION ALGORITHM FOR MTL [1]

---

- 1: **Input:** Data  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathbb{N}_n}$  and weights  $\{w_i(t)\}_{i \in \mathbb{N}_n, t \in \mathbb{N}_T}$ ;
- 2: **Initialize:**  $D \leftarrow I_d/d$  ( $I_d$  is  $d \times d$  identity matrix)
- 3: **while** convergence condition is not true **do**
- 4:   **Step 1:** For  $t = 1, \dots, T$  **do**

$$\tilde{\beta}_t \leftarrow \arg \min_{\tilde{\beta}} \sum_{i \in \mathbb{N}_n} w_i(t) \ell_t(r(y_i, \tilde{\beta}^\top \tilde{\mathbf{x}}_i)) + \frac{\gamma}{T} \beta^\top D^{-1} \beta$$

- 5:   **Step 2:**

$$D \leftarrow \frac{C^{1/2}}{\text{tr}(C)^{1/2}} = \arg \min_{D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1} \sum_{t \in \mathbb{N}_T} \beta_t^\top D^{-1} \beta_t,$$

where  $C := BB^\top$  whose  $(j, k)^{\text{th}}$  element is defined as  $C_{j,k} := \sum_{t \in \mathbb{N}_T} \beta_{tj} \beta_{tk}$ .

- 6: **end while**

- 7: **Output:**  $\{\tilde{\beta}_t\}_{t \in \mathbb{N}_T}$  and  $D$ ;
- 

quantile regression. However, at first glance, the PTL optimization problem (2) seems computationally intractable since we need to find infinitely many task-wise solutions as well as the common feature subspace (or the common set of variables if  $D$  is restricted to be diagonal) shared by infinitely many tasks.

Our key finding is that, for a certain class of PTL problems, when  $D$  is fixed, the optimal path of the task-wise solutions  $\tilde{\beta}_\theta$  is shown to be piecewise-linear in  $\theta$ . By exploiting this piecewise-linearity, we can efficiently handle infinitely many parameterized tasks, and the optimal solutions of those class of PTL problems can be exactly computed.

In the following theorem, we prove that the task-wise solutions  $\tilde{\beta}_\theta$  is piecewise-linear in  $\theta$  if the weight functions and the loss function satisfy certain conditions.

**Theorem 1** For any  $d \times d$  positive-definite matrix  $D \in \mathcal{S}_{++}^d$ , the optimal solution path of

$$\tilde{\beta}_\theta \leftarrow \arg \min_{\tilde{\beta}} \sum_{i \in \mathbb{N}_n} w_i(\theta) \ell_\theta(r(y_i, \tilde{\beta}^\top \tilde{\mathbf{x}}_i)) + \gamma \beta^\top D^{-1} \beta \quad (3)$$

for  $\theta \in [\theta_L, \theta_U]$  is written as a piecewise-linear function of  $\theta$  if the residual  $r(y, \tilde{\beta}^\top \tilde{\mathbf{x}})$  can be written as an affine function of  $\tilde{\beta}$ , and the weight functions  $w_i : [\theta_L, \theta_U] \rightarrow [0, 1]$ ,  $i \in \mathbb{N}_n$  and the loss function  $\ell : \mathbb{R} \rightarrow \mathbb{R}_+$  satisfy either of the following conditions **(a)** or **(b)**:

**(a)** All the weight functions are piecewise-linear functions, and the loss function is a convex piecewise-linear function which does not depend on  $\theta$ ;

**(b)** All the weight functions are piecewise-constant functions, and the loss function is a convex piecewise-linear function which depends on  $\theta$  in the following form:

$$\ell_\theta(r) = \sum_{h \in \mathbb{N}_H} \max\{(a_h + b_h r)(c_h + d_h \theta), 0\}, \quad (4)$$

where  $H$  is a positive integer, and  $a_h, b_h, c_h, d_h \in \mathbb{R}$  are constants such that  $c_h + d_h \theta \geq 0$  for all  $\theta \in [\theta_L, \theta_U]$ .

In the proof in Appendix A, we show that, if the weight functions and the loss function satisfy the conditions **(a)** or **(b)**, the problem (3) is reformulated as a *parametric quadratic program (parametric QP)*, where the parameter  $\theta$  only appears in the linear term of the objective function. As shown, for example, in [9], the optimal solution path of this class of parametric QP has a piecewise-linear form.

If  $\tilde{\beta}_\theta$  is piecewise-linear in  $\theta$ , we can exactly compute the entire solution path by using parametric programming. In machine learning literature, parametric programming is often used in the context

---

**Algorithm 2** ALTERNATING MINIMIZATION ALGORITHM FOR PTL
 

---

- 1: **Input:** Data  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathbb{N}_n}$  and weight functions  $w_i : [\theta_L, \theta_U] \rightarrow [0, 1]$  for all  $i \in \mathbb{N}_n$ ;
- 2: **Initialize:**  $D \leftarrow I_d/d$  ( $I_d$  is  $d \times d$  identity matrix)
- 3: **while** convergence condition is not true **do**
- 4:   **Step 1:** For all the continuum of  $\theta \in [\theta_L, \theta_U]$  **do**

$$\tilde{\beta}_\theta \leftarrow \arg \min_{\tilde{\beta}} \sum_{i \in \mathbb{N}_n} w_i(\theta) \ell_\theta(r(y_i, \tilde{\beta}^\top \tilde{\mathbf{x}}_i)) + \gamma \beta^\top D^{-1} \beta$$

by using parametric programming;

- 5:   **Step 2:**

$$D \leftarrow \frac{C^{1/2}}{\text{tr}(C)^{1/2}} = \arg \min_{D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1} \int_{\theta_L}^{\theta_U} \beta_\theta^\top D^{-1} \beta_\theta d\theta, \quad (5)$$

where  $(j, k)^{\text{th}}$  element of  $C \in \mathbb{R}^{d \times d}$  is defined as  $C_{j,k} := \int_{\theta_L}^{\theta_U} \beta_{\theta,j} \beta_{\theta,k} d\theta$ ;

- 6: **end while**

- 7: **Output:**  $\{\tilde{\beta}_\theta\}$  for  $\theta \in [\theta_L, \theta_U]$  and  $D$ ;
- 

of *regularization path-following* [13, 14, 15]<sup>3</sup>. We start from the solution at  $\theta = \theta_L$ , and follow the path of the optimal solutions while  $\theta$  is continuously increased. This is efficiently conducted by exploiting the piecewise-linearity.

Our proposed algorithm for solving the PTL problem (2) is described in Algorithm 2, which is essentially a continuous version of the MTL algorithm shown in Algorithm 1. Note that, by exploiting the piecewise linearity of  $\beta_\theta$ , we can compute the integral at Step 2 (Eq. (5)) in Algorithm 2.

Algorithm 2 can be changed to parametric-task variable selection if Step 2 is replaced with

$$D \leftarrow \text{diag}(\lambda_1, \dots, \lambda_d) \text{ where } \lambda_j = \frac{\sqrt{\int_{\theta_L}^{\theta_U} \beta_{\theta,j}^2 d\theta}}{\sum_{j' \in \mathbb{N}_d} \sqrt{\int_{\theta_L}^{\theta_U} \beta_{\theta,j'}^2 d\theta}} \text{ for all } j \in \mathbb{N}_d,$$

which can also be computed efficiently by exploiting the piecewise linearity of  $\beta_\theta$ .

## 4 Examples of PTL Problems

In this section, we present three examples where our PTL formulation (2) is useful.

**Binary Classification Under Non-Stationarity** Suppose that we observe  $n$  training instances sequentially, and denote them as  $\{(\mathbf{x}_i, y_i, \tau_i)\}_{i \in \mathbb{N}_n}$ , where  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$ , and  $\tau_i$  is the time when the  $i^{\text{th}}$  instance is observed. Without loss of generality, we assume that  $\tau_1 < \dots < \tau_n$ . Under non-stationarity, if we are requested to learn a classifier to predict the output for a test input  $\mathbf{x}$  observed at time  $\tau$ , the training instances observed around time  $\tau$  should have more influence on the classifier than others.

Let  $w_i(\tau)$  denote the weight of the  $i^{\text{th}}$  instance when training a classifier for a test point at time  $\tau$ . We can for example use the following triangular weight function (see Figure 1):

$$w_i(\tau) = \begin{cases} 1 + s^{-1}(\tau_i - \tau) & \text{if } \tau - s \leq \tau_i < \tau, \\ 1 - s^{-1}(\tau_i - \tau) & \text{if } \tau \leq \tau_i < \tau + s, \\ 0 & \text{otherwise,} \end{cases} \quad (6)$$

where  $s > 0$  determines the width of the triangular time windows. The problem of training a classifier for time  $\tau$  is then formulated as

$$\min_{\tilde{\beta}} \sum_{i \in \mathbb{N}_n} w_i(\tau) \max(0, 1 - y_i \tilde{\beta}^\top \tilde{\mathbf{x}}_i) + \gamma \|\tilde{\beta}\|_2^2,$$

where we used the hinge loss.

---

<sup>3</sup>In regularization path-following, one computes the optimal solution path w.r.t. the regularization parameter, whereas we compute the optimal solution path w.r.t. the task parameter  $\theta$ .

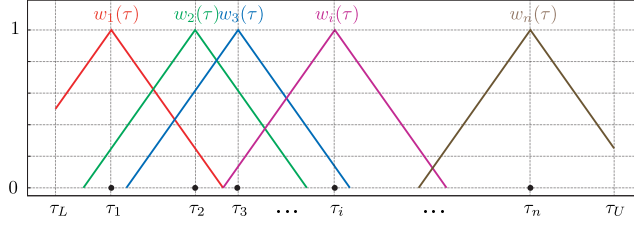


Figure 1: Examples of weight functions  $\{w_i(\tau)\}_{i \in \mathbb{N}_n}$  in non-stationary time-series learning. Given a training instances  $(\mathbf{x}_i, y_i)$  at time  $\tau_i$  for  $i = 1, \dots, n$  under non-stationary condition, it is reasonable to use the weights  $\{w_i(\tau)\}_{i \in \mathbb{N}_n}$  as shown here when we learn a classifier to predict the output of a test input at time  $\tau$ .

If we have the belief that a set of classifiers for different time should have some common structure, we can apply our PTL approach to this problem. If we consider a time interval  $\tau \in [\tau_L, \tau_U]$ , the parametric-task feature learning problem is formulated as

$$\min_{\substack{\{\tilde{\beta}(\tau)\}_{\tau \in [\tau_L, \tau_U]} \\ D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1}} \int_{\tau_L}^{\tau_U} \sum_{i \in \mathbb{N}_n} w_i(\tau) \max(0, 1 - y_i \tilde{\beta}_\tau^\top \tilde{\mathbf{x}}_i) d\tau + \gamma \int_{\tau_L}^{\tau_U} \beta_\tau^\top D^{-1} \beta_\tau d\tau. \quad (7)$$

Note that the problem (7) satisfies the condition **(a)** in Theorem 1.

**Joint Cost-Sensitive Learning** Next, let us consider cost-sensitive binary classification. When the costs of false positives and false negatives are unequal, or when the numbers of positive and negative training instances are highly imbalanced, it is effective to use the *cost-sensitive learning* approach [16]. Suppose that we are given a set of training instances  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathbb{N}_n}$  with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$ . If we know that the ratio of the false positive and false negative costs is approximately  $\theta : (1 - \theta)$ , it is reasonable to solve the following cost-sensitive SVM [17]:

$$\min_{\tilde{\beta}} \sum_{i \in \mathbb{N}_n} w_i(\theta) \max(0, 1 - y_i \tilde{\beta}^\top \tilde{\mathbf{x}}_i) + \gamma \|\beta\|_2^2,$$

where the weight  $w_i(\theta)$  is defined as

$$w_i(\theta) = \begin{cases} \theta & \text{if } y_i = -1, \\ 1 - \theta & \text{if } y_i = +1. \end{cases}$$

When the exact false positive and false negative costs in the test scenario are unknown [4], it is often desirable to train several cost-sensitive SVMs with different values of  $\theta$ . If we have the belief that a set of classifiers for different cost ratios should have some common structure, we can apply our PTL approach to this problem. If we consider an interval  $\theta \in [\theta_L, \theta_U]$ ,  $0 < \theta_L < \theta_U < 1$ , the parametric-task feature learning problem is formulated as

$$\min_{\substack{\{\tilde{\beta}_\theta\}_{\theta \in [\theta_L, \theta_U]} \\ D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1}} \int_{\theta_L}^{\theta_U} \sum_{i \in \mathbb{N}_n} w_i(\theta) \max(0, 1 - y_i \tilde{\beta}_\theta^\top \tilde{\mathbf{x}}_i) d\theta + \gamma \int_{\theta_L}^{\theta_U} \beta_\theta^\top D^{-1} \beta_\theta d\theta. \quad (8)$$

The problem (8) also satisfies the condition **(a)** in Theorem 1. Figure 2 shows an example of joint cost-sensitive learning applied to a toy 2D binary classification problem.

**Joint Quantile Regression** Given a set of training instances  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathbb{N}_n}$  with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$  drawn from a joint distribution  $P(\mathbf{X}, Y)$ , *quantile regression* [19] is used to estimate the conditional  $\tau^{\text{th}}$  quantile  $F_{Y|X=\mathbf{x}}^{-1}(\tau)$  as a function of  $\mathbf{x}$ , where  $\tau \in (0, 1)$  and  $F_{Y|X=\mathbf{x}}$  is the cumulative distribution function of the conditional distribution  $P(Y|X=\mathbf{x})$ . Jointly estimating multiple conditional quantile functions is often useful for exploring the stochastic relationship between  $\mathbf{X}$  and  $Y$  (see Section 5 for an example of joint quantile regression problems). Linear quantile regression along with  $L_2$  regularization [20] at order  $\tau \in (0, 1)$  is formulated as

$$\min_{\tilde{\beta}} \sum_{i \in \mathbb{N}_n} \rho_\tau(y_i - \tilde{\beta}^\top \tilde{\mathbf{x}}_i) + \gamma \|\beta\|_2^2, \quad \rho_\tau(r) := \begin{cases} (1 - \tau)|r| & \text{if } r \leq 0, \\ \tau|r| & \text{if } r > 0. \end{cases}$$

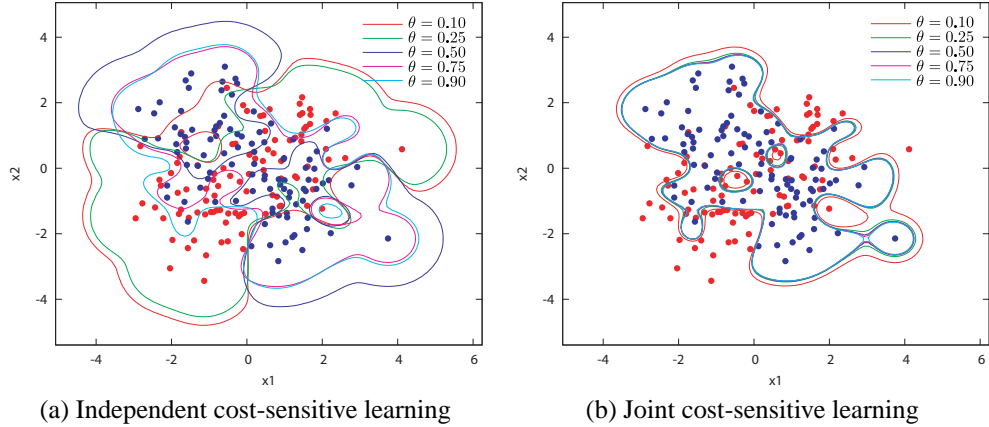


Figure 2: An example of joint cost-sensitive learning on 2D toy dataset (2D input  $\mathbf{x}$  is expanded to  $n$ -dimension by radial basis functions centered on each  $\mathbf{x}_i$ ). In each plot, the decision boundaries of five cost-sensitive SVMs ( $\theta = 0.1, 0.25, 0.5, 0.75, 0.9$ ) are shown. (a) Left plot is the results obtained by independently training each cost-sensitive SVMs. (b) Right plot is the results obtained by jointly training infinitely many cost-sensitive SVMs for all the continuum of  $\theta \in [0.05, 0.95]$  using the methodology we present in this paper (both are trained with the same regularization parameter  $\gamma$ ). When independently trained, the inter-relationship among different cost-sensitive SVMs looks inconsistent (c.f., [18]).

If we have the belief that a family of quantile regressions at various  $\tau \in (0, 1)$  have some common structure, we can apply our PTL framework to joint estimation of the family of quantile regressions. This PTL problem satisfies the condition (b) in Theorem 1, and is written as

$$\min_{\{\beta_\tau\}_{\tau \in (0,1)}, D \in \mathcal{S}_{++}^d, \text{tr}(D) \leq 1} \int_0^1 \sum_{i \in \mathbb{N}_n} \rho_\tau(y_i - \beta_\tau^\top \mathbf{x}_i) d\tau + \gamma \int_0^1 \beta_\tau^\top D^{-1} \beta_\tau d\tau,$$

where we do not need any weighting and omit  $w_i(\tau) = 1$  for all  $i \in \mathbb{N}_n$  and  $\tau \in [0, 1]$ .

## 5 Numerical Illustrations

In this section, we illustrate various aspects of PTL with the three examples discussed in the previous section.

**Artificial Example for Learning under Non-stationarity** We first consider a simple artificial problem with non-stationarity, where the data generating mechanism gradually changes. We assume that our data generating mechanism produces the training set  $\{(\mathbf{x}_i, y_i, \tau_i)\}_{i \in \mathbb{N}_n}$  with  $n = 100$  as follows. For each  $\tau_i \in \{0, 1, \frac{2\pi}{n}, 2\frac{2\pi}{n}, \dots, (n-1)\frac{2\pi}{n}\}$ , the output  $y_i$  is first determined as  $y_i = 1$  if  $i$  is odd, while  $y_i = -1$  if  $i$  is even. Then,  $\mathbf{x}_i \in \mathbb{R}^d$  is generated as

$$x_{i1} \sim N(y_i \cos \tau_i, 1^2), x_{i2} \sim N(y_i \sin \tau_i, 1^2), x_{ij} \sim N(0, 1^2), \forall j \in \{3, \dots, d\}, \quad (9)$$

where  $N(\mu, \sigma^2)$  is the normal distribution with mean  $\mu$  and variance  $\sigma^2$ . Namely, only the first two dimensions of  $\mathbf{x}$  differ in two classes, and the remaining  $d - 2$  dimensions are considered as noise. In addition, according to the value of  $\tau_i$ , the means of the class-wise distributions in the first two dimensions gradually change. The data distributions of the first two dimensions for  $\tau = 0, 0.5\pi, \pi, 1.5\pi$  are illustrated in Figure 3. Here, we applied our PT feature learning approach with triangular time windows in (6) with  $s = 0.25\pi$ . Figure 4 shows the mis-classification rate of PT feature learning (PTFL) and ordinary independent learning (IND) on a similarly generated test sample with size 1000. When the input dimension  $d = 2$ , there is no advantage for learning common features since these two input dimensions are important for classification. On the other hand, as  $d$  increases, PT feature learning becomes more and more advantageous. Especially when the regularization parameter  $\gamma$  is large, the independent learning approach is completely deteriorated as  $d$  increases, while PTFL works reasonably well in all the setups.

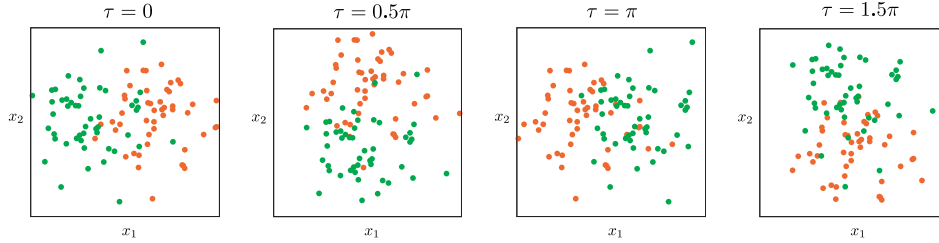


Figure 3: The first 2 input dimensions of artificial example at  $\tau = 0, 0.5\pi, \pi, 1.5\pi$ . The class-wise distributions in these two dimensions gradually change with  $\tau \in [0, 2\pi]$ .

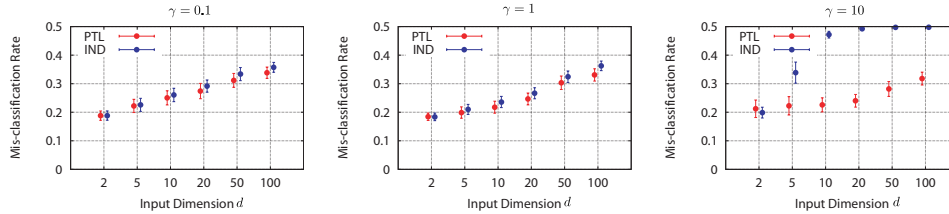


Figure 4: Experimental results on artificial example under non-stationarity. Mis-classification rate on test sample with size 1000 for various setups  $d \in \{2, 5, 10, 20, 50, 100\}$  and  $\gamma \in \{0.1, 1, 10\}$  are shown. The red symbols indicate the results of our PT feature learning (PTFL) whereas the blue symbols indicate ordinary independent learning (IND). The plotted are average (and standard deviation) over 100 replications with different random seeds. All the differences except  $d = 2$  are statistically significant ( $p < 0.01$ ).

**Joint Cost-Sensitive SVM Learning on Benchmark Datasets** Here, we report the experimental results on joint cost-sensitive SVM learning discussed in Section 4. Although our main contribution is not just claiming favorable generalization properties of parametric task learning solutions, we compared, as an illustration, the generalization performances of PT feature learning (PTFL) and PT variable selection (PTVS) with the ordinary independent learning approach (IND). In PTFL and PTVS, we learned common feature subspaces and common sets of variables shared across the continuum of cost-sensitive SVM for  $\theta \in [0.05, 0.95]$  for 10 benchmark datasets (see Table 1). In each data set, we divided the entire sample into training, validation, and test sets with almost equal size. The average test errors (and the standard deviation) of 10 different data splits are reported in Table 1. The total test errors for cost-sensitive SVMs with  $\theta = 0.1, 0.2, \dots, 0.9$  are defined as  $\sum_{\theta \in \{0.1, \dots, 0.9\}} (\theta \sum_{i: y_i = -1} I(f_\theta(\mathbf{x}_i) > 0) + (1 - \theta) \sum_{i: y_i = 1} I(f_\theta(\mathbf{x}_i) \leq 0))$ , where  $f_\theta$  is the trained SVM with the cost ratio  $\theta$ . Model selection was conducted by using the same criterion on validation sets. We see that, in most cases, PTFL or PTVS had better generalization performance than IND.

**Joint Quantile Regression** Finally, we applied PT feature learning to joint quantile regression problems. Here, we took a slightly different approach from what was described in the previous section. Given a training set  $\{(\mathbf{x}_i, y_i)\}_{i \in \mathbb{N}_n}$ , we first estimated conditional mean function  $E[Y|\mathbf{X} = \mathbf{x}]$  by least-square regression, and computed the residual  $r_i := y_i - \hat{E}[Y|\mathbf{X} = \mathbf{x}_i]$ , where  $\hat{E}$  is the estimated conditional mean function. Then, we applied PT feature learning to  $\{(\mathbf{x}_i, r_i)\}_{i \in \mathbb{N}_n}$ , and estimated the conditional  $\tau^{\text{th}}$  quantile function as  $\hat{F}_{Y|\mathbf{X}=\mathbf{x}}^{-1}(\tau) := \hat{E}[Y|\mathbf{X} = \mathbf{x}_i] + \hat{f}_{\text{res}}(\mathbf{x}|\tau)$ , where  $\hat{f}_{\text{res}}(\cdot|\tau)$  is the estimated  $\tau^{\text{th}}$  quantile regression fitted to the residuals.

When multiple quantile regressions with different  $\tau$ s are independently learned, we often encounter a notorious problem known as *quantile crossing* (see Section 2.5 in [5]). For example, in Figure 5(a), some of the estimated conditional quantile functions *cross* each other (which never happens in the true conditional quantile functions). One possible approach to mitigate this problem is to assume a model on the heteroscedastic structure. In the simplest case, if we assume that the data is *homoscedastic* (i.e., the conditional distribution  $P(Y|\mathbf{x})$  does not depend on  $\mathbf{x}$  except its location),

Table 1: Average (and standard deviation) of test errors obtained by joint cost-sensitive SVMs on benchmark datasets.  $n$  is the sample size,  $d$  is the input dimension, Ind indicates the results when each cost-sensitive SVM was trained independently, while PTFL and PTVS indicate the results from PT feature learning and PT feature selection, respectively. The bold numbers in the table indicate the best performance among three methods.

Data Name	$n$	$d$	Ind	PTFL	PTVS
<i>Parkinson</i>	195	20	32.30 (10.60)	<b>30.21 (9.09)</b>	30.25 (8.53)
<i>Breast Cancer Diagnostic</i>	569	30	20.36 (7.77)	<b>18.49 (6.15)</b>	19.46 (5.89)
<i>Breast Cancer Prognostic</i>	194	33	48.97 (12.92)	49.28 (9.83)	<b>48.68 (5.89)</b>
<i>Australian</i>	690	14	117.97 (22.97)	<b>106.25 (12.66)</b>	111.22 (15.95)
<i>Diabetes</i>	768	8	185.90 (21.13)	179.89 (16.31)	<b>175.95 (16.26)</b>
<i>Fourclass</i>	862	2	181.69 (22.13)	179.30 (14.25)	<b>178.67 (19.24)</b>
<i>German</i>	1000	24	242.21 (18.35)	<b>219.66 (16.22)</b>	237.20 (15.78)
<i>Splice</i>	1000	60	179.80 (24.22)	<b>151.69 (18.02)</b>	183.54 (21.27)
<i>SVM Guide</i>	300	10	175.70 (15.55)	<b>170.16 (9.99)</b>	179.76 (14.76)
<i>DVowel</i>	528	10	<b>175.16 (13.78)</b>	175.74 (9.37)	175.50 (7.38)

quantile regressions at different  $\tau$ s can be obtained by just vertically shifting other quantile regression function (see Figure 5(f)).

Our PT feature learning approach, when applied to the joint quantile regression problem, allows us to *interpolate* these two extreme cases. Figure 5 shows a joint QR example on the bone mineral density (BMD) data [21]. We applied our approach after expanding univariate input  $x$  to a  $d = 5$  dimensional vector by using evenly allocated RBFs. When (a)  $\gamma \rightarrow 0$ , our approach is identical with independently estimating each quantile regression, while it coincides with homoscedastic case when (f)  $\gamma \rightarrow \infty$ . In our experience, the best solution is usually found somewhere between these two extremes: in this example, (d)  $\gamma = 5$  was chosen as the best model by 10-fold cross-validation.

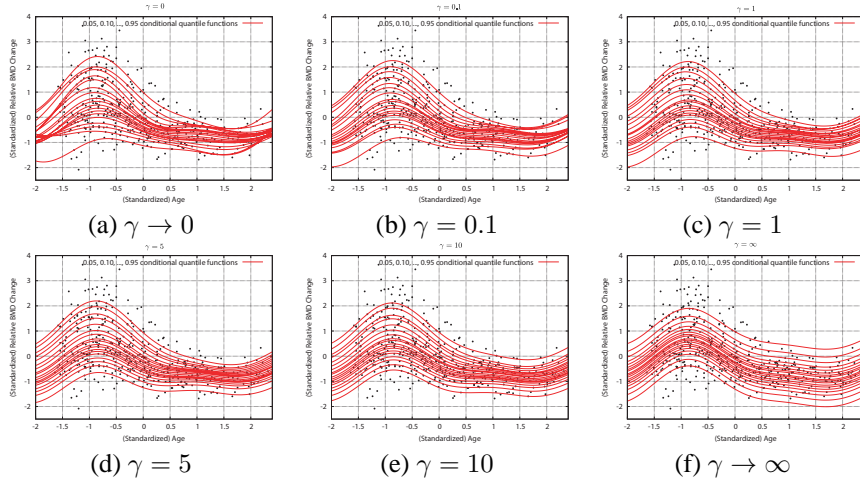


Figure 5: Joint quantile regression examples on BMD data [21] for six different  $\gamma$ s.

## 6 Conclusions

In this paper, we introduced parametric-task learning (PTL) approach that can systematically handle infinitely many tasks parameterized by a continuous parameter. We illustrated the usefulness of this approach by providing three examples that can be naturally formulated as PTL. We believe that there are many other practical problems that falls into this PTL framework.

## Acknowledgments

The authors thank the reviewers for fruitful comments. IT, MS, and SN thank the support from MEXT Kakenhi 23700165, JST CREST Program, MEXT Kakenhi 23120004, respectively.



## References

- [1] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *Advances in Neural Information Processing Systems*, volume 19, pages 41–48. 2007.
- [2] A. Argyriou, C. A. Micchelli, M. Pontil, and Y. Ying. A spectral regularization framework for multi-task structure learning. In *Advances in Neural Information Processing Systems*, volume 20, pages 25–32. 2008.
- [3] L. Cao and F. Tay. Support vector machine with adaptive parameters in financial time series forecasting. *IEEE Transactions on Neural Networks*, 14(6):1506–1518, 2003.
- [4] F. R. Bach, D. Heckerman, and E. Horvits. Considering cost asymmetry in learning classifiers. *Journal of Machine Learning Research*, 7:1713–41, 2006.
- [5] R. Koenker. *Quantile Regression*. Cambridge University Press, 2005.
- [6] K. Ritter. On parametric linear and quadratic programming problems. *mathematical Programming: Proceedings of the International Congress on Mathematical Programming*, pages 307–335, 1984.
- [7] E. L. Allgower and K. George. Continuation and path following. *Acta Numerica*, 2:1–63, 1993.
- [8] T. Gal. *Postoptimal Analysis, Parametric Programming, and Related Topics*. Walter de Gruyter, 1995.
- [9] M. J. Best. An algorithm for the solution of the parametric quadratic programming problem. *Applied Mathematics and Parallel Computing*, pages 57–76, 1996.
- [10] M. Fazel, H. Hindi, and S. P. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *Proceedings of the American Control Conference*, volume 6, pages 4734–4739, 2001.
- [11] B. A. Turlach, W. N. Venables, and S. J. Wright. Simultaneous variable selection. *Technometrics*, 47:349–363, 2005.
- [12] G. Obozinski, B. Taskar, and M. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, 20(2):231–252, 2010.
- [13] M. R. Osborne, B. Presnell, and B. A. Turlach. A new approach to variable selection in least squares problems. *IMA Journal of Numerical Analysis*, 20(20):389–404, 2000.
- [14] B. Efron and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32(2):407–499, 2004.
- [15] T. Hastie, S. Rosset, R. Tibshirani, and J. Zhu. The entire regularization path for the support vector machine. *Journal of Machine Learning Research*, 5:1391–415, 2004.
- [16] Y. Lin, Y. Lee, and G. Wahba. Support vector machines for classification in nonstandard situations. *Machine Learning*, 46:191–202, 2002.
- [17] M. A. Davenport, R. G. Baraniuk, and C. D. Scott. Tuning support vector machine for min-max and Neyman-Pearson classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2010.
- [18] G. Lee and C. Scott. Nested support vector machines. *IEEE Transactions on Signal Processing*, 58(3):1648–1660, 2010.
- [19] R. Koenker. *Quantile Regression*. Cambridge University Press, 2005.
- [20] I. Takeuchi, Q. V. Le, T. Sears, and A. J. Smola. Nonparametric quantile estimation. *Journal of Machine Learning Research*, 7:1231–1264, 2006.
- [21] L. K. Bachrach, T. Hastie, M. C. Wang, B. Narasimhan, and R. Marcus. Acquisition in healthy Asian, hispanic, black and caucasian youth. a longitudinal study. *The Journal of Clinical Endocrinology and Metabolism*, 84:4702–4712, 1999.
- [22] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.

## Appendix A: Proof of Theorem 1

First, we prove the claim for the condition **(a)**. Let us divide the interval  $[\theta_L, \theta_U]$  into finite number of segments so that, within each segment, the weight vector  $\mathbf{w}(\theta) := (w_1(\theta), \dots, w_n(\theta))^\top \in [0, 1]^n$  changes linearly with  $\theta$ , and denote the breakpoints of those segments as  $\theta_L = \theta_0 < \theta_1 < \dots < \theta_s < \dots < \theta_S = \theta_U$ , where  $S$  is the number of those segments.

Then, consider a segment defined on  $\theta \in [\theta_s, \theta_{s+1}]$ ,  $s \in \{0, \dots, S-1\}$ , and denote the weight vectors at  $\theta_s$  and  $\theta_{s+1}$  as  $\mathbf{w}(\theta_s)$  and  $\mathbf{w}(\theta_{s+1})$ , respectively. The problem of computing the solution path within this segment is written as the following parametric optimization problem

$$\tilde{\boldsymbol{\beta}}_\mu \leftarrow \arg \min_{\boldsymbol{\beta}} \sum_{i \in \mathbb{N}_n} ((1-\mu)w_i(\theta_s) + \mu w_i(\theta_{s+1})) \ell_{(1-\mu)\theta_s + \mu\theta_{s+1}}(r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i)) + \gamma \boldsymbol{\beta}^\top D^{-1} \boldsymbol{\beta} \quad (10)$$

for  $\mu \in [0, 1]$ .

Since the loss function  $\ell_\theta$  does not depend on  $\theta$  and is convex piecewise-linear in  $r$ , we can write  $\ell_\theta$  as

$$\ell_\theta(r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i)) = \sum_{h \in \mathbb{N}_H} \max\{\phi_{ih} + \psi_{ih} \cdot r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i)\},$$

where  $\phi_{ih}, \psi_{ih} \in \mathbb{R}$ ,  $(i, h) \in \mathbb{N}_n \times \mathbb{N}_H$  are constants, and  $H$  is the number of pieces of the piecewise-linear loss function  $\ell_\theta$  (see, for example, section 4.3.1 of [22]).

Using slack variables  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n$ , the parametric programming problem in (10) is rewritten as

$$\begin{aligned} \{\tilde{\boldsymbol{\beta}}_\mu, \boldsymbol{\xi}_\mu\} &\leftarrow \arg \min_{\tilde{\boldsymbol{\beta}}, \boldsymbol{\xi}} ((1-\mu)\mathbf{w}(\theta_s) + \mu\mathbf{w}(\theta_{s+1}))^\top \boldsymbol{\xi} + \gamma \boldsymbol{\beta}^\top D^{-1} \boldsymbol{\beta} \\ \text{s.t.} \quad &\xi_i \geq \phi_{ih} + \psi_{ih} \cdot r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i) \text{ for all } (i, h) \in \mathbb{N}_n \times \mathbb{N}_H \end{aligned} \quad (11)$$

with respect to  $\mu \in [0, 1]$ . The problem (11) belongs to the class of *parametric QP* (note that, when  $\mu$  is fixed, the problem (11) is quadratic program with respect to  $\tilde{\boldsymbol{\beta}}$  and  $\boldsymbol{\xi}$ , which has a quadratic objective function and a set of linear constraints.). As shown, for example, in [6, 9], a parametric quadratic program which contains the parameter ( $\mu$ ) in the linear term of the quadratic objective function are shown to have a solution path in piecewise-linear form.

Similarly for the condition **(b)**, we consider a segment defined on  $\theta \in [\theta_t, \theta_{s+1}]$ ,  $s \in \{0, \dots, S-1\}$ , in which the weight vector  $\mathbf{w}(\theta)$  is constant (and thus omitted hereafter). Using slack variables  $\xi_{ih}$  for  $i \in \mathbb{N}_n$  and  $h \in \mathbb{N}_H$

$$\begin{aligned} &\min_{\tilde{\boldsymbol{\beta}}} \sum_{i \in \mathbb{N}_n} \sum_{h \in \mathbb{N}_H} \max\{(a_h + b_h \cdot r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i))(c_h + d_h \theta), 0\} + \gamma \boldsymbol{\beta}^\top D^{-1} \boldsymbol{\beta} \\ \Leftrightarrow &\min_{\tilde{\boldsymbol{\beta}}} \sum_{h \in \mathbb{N}_H} (c_h + d_h \theta) \sum_{i \in \mathbb{N}_n} \max\{(a_h + b_h \cdot r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i)), 0\} + \gamma \boldsymbol{\beta}^\top D^{-1} \boldsymbol{\beta} \\ \Leftrightarrow &\min_{\tilde{\boldsymbol{\beta}}, \boldsymbol{\xi}} \sum_{h \in \mathbb{N}_H} (c_h + d_h \theta) \sum_{i \in \mathbb{N}_n} \xi_{ih} + \gamma \boldsymbol{\beta}^\top D^{-1} \boldsymbol{\beta} \\ \text{s.t.} \quad &\xi_{ih} \geq a_h + b_h \cdot r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i), \xi_{ih} \geq 0 \forall (i, h) \in \mathbb{N}_n \times \mathbb{N}_H. \end{aligned}$$

The parametric programming problem in Theorem 1 **(b)** is thus written as

$$\begin{aligned} \{\tilde{\boldsymbol{\beta}}_\theta, \boldsymbol{\xi}_\theta\} &\leftarrow \min_{\tilde{\boldsymbol{\beta}}, \boldsymbol{\xi}} \sum_{h \in \mathbb{N}_H} (c_h + d_h \theta) \sum_{i \in \mathbb{N}_n} \xi_{ih} + \gamma \boldsymbol{\beta}^\top D^{-1} \boldsymbol{\beta} \\ \text{s.t.} \quad &\xi_{ih} \geq a_h + b_h \cdot r(y_i, \tilde{\boldsymbol{\beta}}^\top \tilde{\mathbf{x}}_i), \xi_{ih} \geq 0 \forall (i, h) \in \mathbb{N}_n \times \mathbb{N}_H \end{aligned}$$

for  $\theta \in [\theta_s, \theta_{s+1}]$ , and it also belongs to parametric QP, meaning that the optimal solution path is shown to be piecewise linear in  $\theta$ .  $\square$