

Template-Based Piecewise Affine Regression

Guillaume O. Berger † and Sriram Sankaranarayanan *

† UCLouvain, Louvain-la-Neuve, Belgium guillaume.berger@uclouvain.be

* University of Colorado Boulder, CO, USA srirams@colorado.edu

September 19, 2024

Abstract

We study the problem of fitting a piecewise affine (PWA) function to input–output data. Our algorithm divides the input domain into finitely many regions whose shapes are specified by a user-provided template and such that the input–output data in each region are fit by an affine function within a user-provided error tolerance. We first prove that this problem is NP-hard. Then, we present a top-down algorithmic approach for solving the problem. The algorithm considers subsets of the data points in a systematic manner, trying to fit an affine function for each subset using linear regression. If regression fails on a subset, the algorithm extracts a minimal set of points from the subset (an unsatisfiable core) that is responsible for the failure. The identified core is then used to split the current subset into smaller ones. By combining this top-down scheme with a set-covering algorithm, we derive an overall approach that provides optimal PWA models for a given error tolerance, where optimality refers to minimizing the number of pieces of the PWA model. We demonstrate our approach on three numerical examples that include PWA approximations of a widely used nonlinear insulin–glucose regulation model and a double inverted pendulum with soft contacts.

1 Introduction

The problem of identifying models from data is central to designing and verifying Cyber-Physical Systems (CPS). These models can predict the output of a subsystem for a given input or the next state of a dynamical system from the current state. Even if there is a physical basis for constructing a model of the system under investigation, it is often necessary to use data-driven modeling to augment these models to incorporate aspects of the system that cannot be easily modeled. CPS are often nonlinear and *multi-modal*, wherein different regions of the input/state space have different laws that govern the relationship between the inputs and outputs. In this paper, we study piecewise affine (PWA) regression. The goal of PWA regression is to fit a piecewise affine function to a given data set of input–output pairs $\{(x_k, y_k)\}_{k=1}^N$. The piecewise affine function splits the input domain into finitely many regions H_1, \dots, H_M and associates an affine function $f_i(x) = A_i x + b_i$ to each region H_i . This is illustrated in Fig. 1 below.

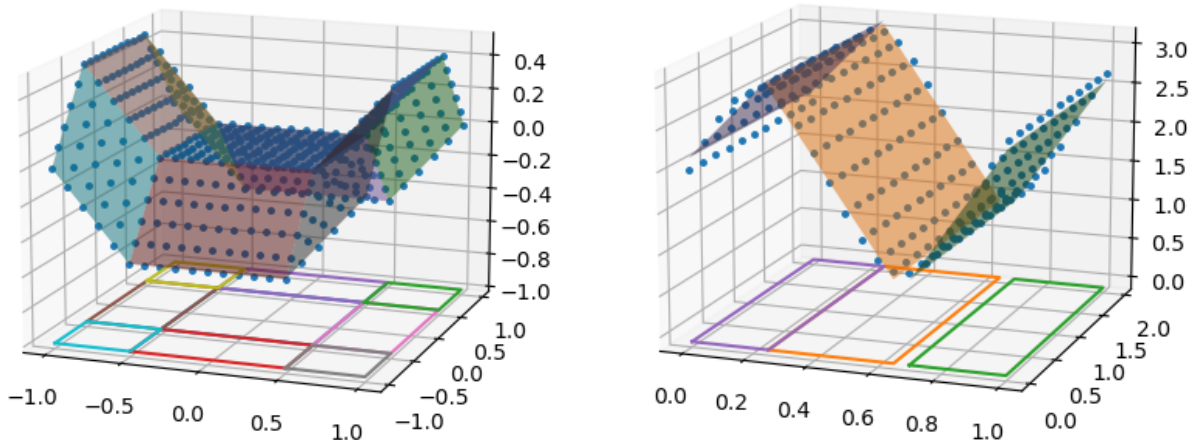


Figure 1: PWA regression of a set of input–output data points with rectangular template.

31 Further, we seek a PWA model that fits the given data while respecting a user-provided error
 32 bound ϵ and minimizing the number of regions (pieces). This problem has numerous applications
 33 including the identification of hybrid systems with state-based switching and simplifying nonlinear
 34 models using PWA approximations. Existing PWA regression approaches usually do not restrict how
 35 the input domain is split. For instance, an approach that simply specifies that the input domain is
 36 covered by polyhedral sets leads to high computational complexity for the regression algorithm [18].

37 In this paper, we restrict the possible shape of the polyhedral regions by requiring that each
 38 region H_i is described by a vector inequality $p(x) \leq c_i$, wherein p is a fixed, user-provided, vector-
 39 valued function, called the *template*, while the regions are obtained by varying the offset vector c_i .
 40 The resulting problem, called template-based PWA regression, allows us to split the input domain
 41 into pre-specified shapes using a suitable template. For instance, in Fig. 1, the input domain is
 42 divided into rectangular regions. Our contributions are as follows:

43 After introducing and formalizing the problem of template-based PWA regression (Sec. “Problem
 44 Statement”), we show that — similarly to the classical PWA regression problem [18] — the problem
 45 of template-based PWA regression is NP-hard in the dimension of the input space and the size
 46 of the template, but polynomial in the size of the data set (Sec. “Computational Complexity”).
 47 Next, we provide an algorithm for optimal bounded-error template-based PWA regression, i.e., with
 48 minimal number of regions while fitting the data within the given error tolerance (Sec. “Top-Down
 49 Algorithm”). Our algorithm is top-down because it breaks large sets of data into smaller ones until
 50 those can be fit by an affine function. A more detailed overview of the algorithm is provided later
 51 in this introduction, after the “Related Work”. Finally, we apply our algorithm on two practical
 52 problems (Sec. “Numerical Experiments”): the approximation of a nonlinear system, namely the
 53 insulin–glucose regulation process [12], with affine functions with rectangular domains, and the
 54 identification of a hybrid linear system consisting in an inverted double pendulum with soft contacts
 55 on the joints. For both applications, we show that template-based PWA regression is favorable
 56 compared to classical PWA regression both in terms of computation time and our ability to formulate
 57 models from the results. We also compare different templates for the identification of a hybrid linear
 58 system consisting of two carts with springs.

59 This paper is an extension of a preliminary version that appeared as part of the Learning for
 60 Decision and Control (L4DC) conference in May 2023 [7]. This paper extends our previous version
 61 by providing more detailed explanations of the algorithms and complete proofs of all the results.

62 We have also added a new numerical example involving the identification of a hybrid linear system
63 consisting of two carts with springs.

64 1.1 Related work

65 Piecewise affine systems and hybrid linear systems appear naturally in a wide range of applications
66 [16], or as approximations of more complex systems [10]. Therefore, the problems of switched affine
67 (SA) and piecewise affine (PWA) regression have received a lot of attention in the literature; see,
68 e.g., Paoletti et al. [27] and Lauer and Bloch [18] for surveys. Both problems are known to be
69 NP-hard [18]. The problem of SA regression can be formulated as a Mixed-Integer Program (MIP)
70 and solved using MIP solvers, but the complexity is exponential in the number of data points [27].
71 Vidal et al. [35] propose an efficient algebraic approach to solve the problem, but it is restricted to
72 noiseless data. Heuristics to solve the problem in the general case include greedy algorithms [5],
73 continuous relaxations of the MIP [21], block-coordinate descent (similar to k -mean regression)
74 algorithms [9, 17] and refinement of the algebraic approach using sum-of-squares relaxations [25];
75 however, these methods offer no guarantees of finding an (optimal) solution to the problem. As
76 for PWA regression, classical approaches include clustering-based methods [13], data classification
77 followed by geometric clustering [23] and block-coordinate descent algorithms [4]; however, these
78 methods are not guaranteed to find a (minimal) piecewise affine model. In this regard, our approach
79 considers a novel “top-down” approach that focuses on searching for subsets of the data that can be
80 part of the same affine model for the given error bounds.

81 Our approach contrasts with most other approaches in that it is top-down and focuses on refining
82 the domains (using the template assumption) of the pieces until affine fitting is possible. Indeed,
83 most approaches for PWA regression (e.g., [5], [13], [4]) use a two-step approach in which the
84 data are first clustered by solving a SA regression problem and then the clusters are separated
85 into polyhedral regions. Medhat et al. [19] and Yuan et al. [36] use a similar two-step approach
86 for learning hybrid linear automata. There are also approaches that learn the function and the
87 domains in one step: for instance, Breiman [10] for a special class of continuous PWA functions
88 called *Hinging Hyperplanes*, Sadraddini and Belta [31] using mixed-integer linear programming
89 (MILP), and Berger, Narasimhamurthy, and Sankaranarayanan [6] for a class of PWA systems,
90 called guarded linear systems. The class of Hinging Hyperplanes and guarded linear systems are
91 not comparable in general with those studied in this work.¹ Soto, Henzinger, and Schilling [33]
92 also learn the function and the domains simultaneously for hybrid linear systems; however, their
93 incremental approach is greedy, so that it does not come with guarantees of minimality. By contrast,
94 our approach guarantees to find a PWA function with minimal number of regions from the template.

95 Piecewise affine systems with constraints on the domain appear naturally in several applications
96 including biology [29] and mechanical systems with contact forces [1], or as approximations of
97 nonlinear systems [32]. Techniques for PWA regression with rectangular domains have been proposed
98 in Münz and Krebs [22] and Smarra et al. [32]; however, these approaches impose further restrictions
99 on the arrangement of the domains of the functions (e.g., forming a grid) and they are not guaranteed
100 to find a solution with a minimal number of pieces. In the one-dimensional case (time series),
101 optimal time series segmentation can be computed efficiently by using dynamic programming [2,
102 3, 26, 24], but the approach does not extend to higher dimension and solves a different problem
103 (min. # switches vs. min. # pieces). Our approach bears similarities with the “split-and-merge”
104 algorithm of Pavlidis and Horowitz [28], except that (i) we only split regions and never merge them
105 because by construction merging would lead to incompatible regions; (ii) our algorithm comes with

¹See for instance the example in Berger, Narasimhamurthy, and Sankaranarayanan [6, Lemma 2.5] for a system that can be described with the rectangular template but not as a guarded linear system.

106 guarantees of optimality; and (iii) we address the problem of PWA regression and not segmentation.
 107 As for the application involving mechanical systems with contact forces, a recent work by Jin et al.
 108 [15] proposes a heuristic based on minimizing a loss function to learn *linear complementary systems*.

109 1.2 Approach at a glance

110 Fig. 2 below shows the working of our algorithm on a simple data set with $N = 11$ points
 111 $(x_k, y_k) \in \mathbb{R} \times \mathbb{R}$ (see Plot I). We seek a PWA function that fits the data with error tolerance $\epsilon = 0.1$
 112 and with the smallest number of affine pieces (green lines in III, V and VI).

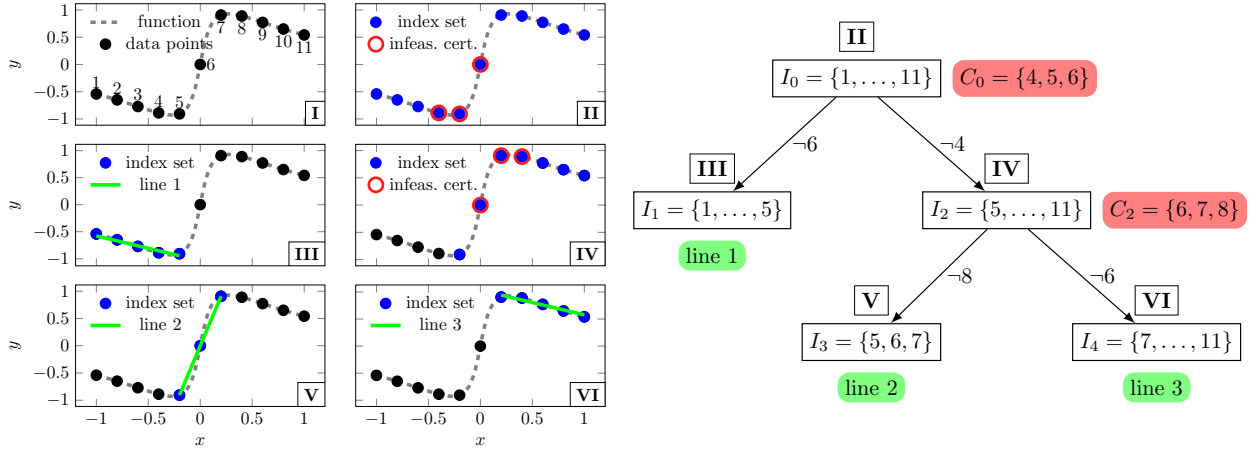


Figure 2: *Left.* Illustration of our algorithm on a simple data set with 11 data points $(x_k, y_k) \in \mathbb{R} \times \mathbb{R}$.
Right. the index sets explored by our algorithm.

113 The algorithm works as follows. At the very first step, the approach tries to fit a single straight
 114 line through all the 11 points. This corresponds to the *index set* $I_0 = \{1, \dots, 11\}$ (see II) where
 115 the points are indexed as in I. However, no such line can fit the points for the given ϵ . Hence, our
 116 approach generates an *infeasibility certificate* that identifies the indices $C_0 = \{4, 5, 6\}$ as a cause
 117 of infeasibility (see II). In other words, we cannot have all three points in C_0 be part of the same
 118 piece of the PWA function we seek. As explained in the paper, infeasibility certificates can be
 119 computed efficiently using Linear Programming. Moreover, in case of infeasibility, we can always
 120 find an infeasibility certificate with cardinality at most $d + 2$, where d is the dimension of the input
 121 (here, $d = 1$).² Our approach then splits I_0 into two subsets $I_1 = \{1, \dots, 5\}$ and $I_2 = \{5, \dots, 11\}$.
 122 These two sets are maximal intervals with respect to set inclusion and do not contain C_0 . The
 123 set I_1 can be fit by a single straight line with tolerance ϵ (see III). However, considering I_2 , we
 124 notice once again that a single straight line cannot be fit (see IV). We identify the set $C_2 = \{6, 7, 8\}$
 125 as an infeasibility certificate and our algorithm splits I_2 into maximal subsets $I_3 = \{5, 6, 7\}$ and
 126 $I_4 = \{7, \dots, 11\}$. Each of these subsets can be fit by a straight line (see V and VI). Thus, our
 127 approach finishes by discovering three affine pieces that cover all the points $\{1, \dots, 11\}$. Note that
 128 although the data point indexed by 5 is part of two pieces, we can resolve this “tie” in an arbitrary
 129 manner by assigning 5 to the first piece and removing it from the second; the same holds for the
 130 data point indexed by 7.

²This is a consequence of Farkas’ Lemma and Carathéodory’s Theorem; see Lemma 2.

131 **Notation** Given two vectors or matrices u and v , their horizontal (resp. vertical) concatenation is
 132 denoted by $[u, v]$ (resp. $[u; v]$). For positive integers d and e and a scalar α , we denote by $[\alpha]_d$ (resp.
 133 $[\alpha]_{e,d}$) the vector in \mathbb{R}^d (resp. matrix in $\mathbb{R}^{e \times d}$) whose components are all equal to α . Finally, given
 134 an natural number n , we let $[n] = \{1, \dots, n\}$.

135 2 Problem Statement

136 Given a data set of $N \in \mathbb{N}_{>0}$ input–output pairs $\{(x_k, y_k)\}_{k=1}^N \subseteq \mathbb{R}^d \times \mathbb{R}^e$, the problem of piecewise
 137 affine (PWA) regression aims at finding a piecewise affine (PWA) function that fits the data within
 138 some given error tolerance $\epsilon \geq 0$. Formally, a PWA function f over a domain $D \subseteq \mathbb{R}^d$ is defined by
 139 covering D with M regions H_1, \dots, H_M and associating an affine function $f_i(x) = A_i x + b_i$ with
 140 each H_i :

$$f(x) = \begin{cases} A_1 x + b_1 & \text{if } x \in H_1, \\ A_2 x + b_2 & \text{if } x \in H_2, \\ \vdots & \\ A_M x + b_M & \text{if } x \in H_M. \end{cases}$$

141 Note that if $H_i \cap H_j \neq \emptyset$ for some $i \neq j$, then f is no longer a function. However, in such a case, we
 142 may “break the tie” by defining $f(x) = f_i(x)$ wherein $i = \min \{j : x \in H_j\}$.

143 **Problem 1** (PWA regression). Given a data set $\{(x_k, y_k)\}_{k=1}^N$ and an error bound $\epsilon \geq 0$, find M
 144 regions $H_i \subseteq \mathbb{R}^d$ and affine functions $f_i(x) = A_i x + b_i$ such that

$$\forall k, \exists i : x_k \in H_i \quad \text{and} \quad \forall k, \forall i : x_k \in H_i \Rightarrow \|y_k - f_i(x_k)\|_\infty \leq \epsilon. \quad (1)$$

145 Furthermore, we restrict the domain H_i of each affine piece by specifying a *template*, which can
 146 be any function $p : \mathbb{R}^d \rightarrow \mathbb{R}^h$. Given a template p and a vector $c \in \mathbb{R}^h$, we define the set $H(c)$ as

$$H(c) = \{x \in \mathbb{R}^d : p(x) \leq c\}, \quad (2)$$

147 wherein \leq is elementwise and $c \in \mathbb{R}^h$ parameterizes the set $H(c)$. We let $\mathcal{H} = \{H(c) : c \in \mathbb{R}^h\}$
 148 denote the set of all regions in \mathbb{R}^d described by the template p .

149 Fixing a template allows to control the complexity of the domains, and thus of the overall PWA
 150 function. This allows among others to mitigate overfitting. The *rectangular* template $p(x) = [x; -x]$
 151 defines regions $H(c)$ that form boxes in \mathbb{R}^d . Indeed, for two vectors $c_1 \leq c_2$, $H([c_2; -c_1])$ defines the
 152 box $\{x \in \mathbb{R}^d : c_1 \leq x \leq c_2\}$. Similarly, allowing pairwise differences between individual variables as
 153 components of p yields the *octagon domain* [20]. Fig. 1 illustrates PWA functions with rectangular
 154 domains. Thus, we define the template-based piecewise affine (TPWA) regression problem:

155 **Problem 2** (TPWA regression). Given a data set $\{(x_k, y_k)\}_{k=1}^N$, a template $p : \mathbb{R}^d \rightarrow \mathbb{R}^h$ and an
 156 error bound $\epsilon > 0$, find M regions $H_i \in \mathcal{H}$ and affine functions $f_i(x) = A_i x + b_i$ such that (1) is
 157 satisfied.

158 Prob. 2 can be posed as a *decision problem*: Given a bound \hat{M} , is there a TPWA function with
 159 $M \leq \hat{M}$ pieces? Alternatively, we can pose it as an optimization problem: Find a TPWA function
 160 with minimum number M of pieces.

161 Although a solution to the decision problem can be used repeatedly to solve the optimization
 162 problem, we will focus on directly solving the optimization problem in this paper. Prob. 2 is closely
 163 related to the well-known problem of switched affine (SA) regression, in which one aims to explain
 164 the data with a finite number of affine functions, but there is no assumption on which function

165 may explain each data point (x_k, y_k) . In other words, SA aims to identify a sufficient number of
 166 modes and dynamics corresponding to each mode without necessarily explaining how the modes are
 167 assigned to each point in the domain.

168 **Problem 3** (SA regression). Given a data set $\{(x_k, y_k)\}_{k=1}^N$ and an error bound $\epsilon \geq 0$, find M
 169 affine functions $f_i(x) = A_i x + b_i$ such that $\forall k, \exists i: \|y_k - f_i(x)\|_\infty \leq \epsilon$.

170 3 Computational Complexity

171 The problem of SA regression (Prob. 3) is known to be NP-hard, even for $M = 2$ [18, Sec. 5.2.4]. In
 172 this section, we show that the same holds for the decision version of Prob. 2. We study the problem
 173 in the RAM model, wherein the problem input size is $N(d + e) + h$, where $p: \mathbb{R}^d \rightarrow \mathbb{R}^h$.

174 **Theorem 1** (NP-hardness). *The decision version of Prob. 2 is NP-hard, even for $M = 2$ and*
 175 *rectangular templates $(p(x) = [x; -x])$.*

176 The proof reduces Prob. 3 which is known to be NP-hard to Prob. 2.

177 *Proof.* Without loss of generality, we restrict to piecewise *linear* models since piecewise affine models
 178 can be obtained from linear ones by augmenting each data point x_k with a component equal to 1.

179 We reduce Prob. 3 to Prob. 2 as follows. Consider an instance of Prob. 3 consisting of a data set
 180 $\mathcal{D} = \{(x_k, y_k)\}_{k=1}^N \subseteq \mathbb{R}^d \times \mathbb{R}^e$ and tolerance ϵ . From \mathcal{D} , we build another data set $\mathcal{D}' \subseteq \mathbb{R}^{d+N} \times \mathbb{R}^e$
 181 with $|\mathcal{D}'| = 4N$ as follows. For each $k \in [N]$, we let $\chi_k \in \mathbb{R}^N$ be the indicator vector of the k^{th}
 182 component. We define

$$\mathcal{D}' = \bigcup_{\sigma \in \{-1, 1\}} \bigcup_{k=1}^N [\{([\sigma x_k; \chi_k], \sigma y_k), ([0]_d; \chi_k), [\sigma \epsilon]_e\}].$$

183 In other words, for each data point (x_k, y_k) in the original dataset \mathcal{D} , we add four data points of the
 184 form $([x_k; \chi_k], y_k)$, $([-x_k; \chi_k], -y_k)$, $([[0]_d; \chi_k], [\epsilon]_e)$, $([[0]_d; \chi_k], [-\epsilon]_e)$ wherein χ_k is a vector of size
 185 N with a 1 entry in the k^{th} position and a 0 entry everywhere else and $[0]_d$ is a vector of d zeros
 186 wherein $[\epsilon]_e$ is a vector with e entries each of which is ϵ .

187 Also, we let $p: \mathbb{R}^{d+N} \rightarrow \mathbb{R}^{2(d+N)}$ be the rectangular template in \mathbb{R}^{d+N} .

188 *Main step:* We show that Prob. 3 with \mathcal{D} , ϵ and $M = 2$ has a solution if and only if Prob. 2
 189 with \mathcal{D}' , p , ϵ and $M = 2$ has a solution.

190 *Proof of “if direction”.* Assume that Prob. 2 has a solution given by $H_1, H_2 \subseteq \mathbb{R}^{d+N}$ and
 191 $A_1, A_2 \in \mathbb{R}^{e \times (d+N)}$, and for each i , decompose $A_i = [B_i, C_i]$, wherein $B_i \in \mathbb{R}^{e \times d}$ and $C_i \in \mathbb{R}^{e \times N}$.
 192 We will show that B_1, B_2 provide a solution to Prob. 3.

193 Therefore, fix $k \in [N]$. Using the pigeon-hole principle, let $i \in \{1, 2\}$ be such that at least two
 194 points in $\{[x_k; \chi_k], [-x_k; \chi_k], [[0]_d; \chi_k]\}$ belong to H_i . Then, by the convexity of H_i , it holds that
 195 $[[0]_d; \chi_k] \in H_i$. For definiteness, assume that $[x_k; \chi_k] \in H_i$. Since H_1, H_2 and A_1, A_2 provide a
 196 solution to Prob. 2, it follows that

$$\|y_k - B_i x_k - C_i \chi_k\|_\infty \leq \epsilon, \quad \|[\epsilon]_e - C_i \chi_k\|_\infty \leq \epsilon, \quad \|[\epsilon]_e + C_i \chi_k\|_\infty \leq \epsilon.$$

197 The last two conditions imply that $C_i \chi_k = 0$, so that $\|y_k - B_i x_k\|_\infty \leq \epsilon$. Since k was arbitrary, this
 198 shows that B_1, B_2 provide a solution to Prob. 3; thereby proving the “if direction”.

199 *Proof of “only if direction”.* Assume that Prob. 3 has a solution given by $A_1, A_2 \in \mathbb{R}^{e \times d}$. Then,
 200 for each $k \in [N]$, define the intervals $I_{1,k}, I_{2,k} \subseteq \mathbb{R}$ as follows: $I_{i,k} = [0, 1]$ if $\|y_k - A_i x_k\|_\infty \leq \epsilon$,

201 and $I_{i,k} = \{0\}$ otherwise. Now, define the rectangular regions $H_1, H_2 \subseteq \mathbb{R}^{d+K}$ as follows: $H_i =$
 202 $\mathbb{R}^d \times I_{i,1} \times \cdots \times I_{i,N}$. Also define the matrices $B_1, B_2 \in \mathbb{R}^{e \times (d+N)}$ as follows: $B_i = [A_i, [0]_{e,K}]$. We
 203 will show that H_1, H_2 and B_1, B_2 provide a solution to Prob. 2.

204 Therefore, fix $k \in [N]$ and $i \in \{1, 2\}$. First, assume $\|y_k - A_i x_k\|_\infty \leq \epsilon$. We show that (a) $[x_k; \chi_k]$,
 205 $[-x_k; \chi_k]$ and $[[0]_d; \chi_k]$ belong to H_i , and (b)

$$\|y_k - B_i[x_k; \chi_k]\|_\infty \leq \epsilon, \quad \|-y_k - B_i[-x_k; \chi_k]\|_\infty \leq \epsilon, \quad \|[\pm \epsilon]_e - B_i[[0]_d; \chi_k]\|_\infty \leq \epsilon.$$

206 This is direct (a) by the definition of $I_{i,k}$, and (b) by the definition of B_i . Now, assume that
 207 $\|y_k - A_i x_k\|_\infty \leq \epsilon$ does not hold. We show that $[x_k; \chi_k], [-x_k; \chi_k]$ do not belong to H_i . This is
 208 direct since $1 \notin I_{i,k}$. Thus, we have shown that H_1, H_2 and B_1, B_2 provide a solution to Prob. 2;
 209 thereby proving the ‘‘only if direction’’.

210 Hence, we have built a polynomial-time reduction from Prob. 3 to Prob. 2. Since Prob. 3 is
 211 NP-hard [18, Sec. 5.2.4], this shows that Prob. 2 is NP-hard as well. \square

212 *Remark 1.* Note that the reduction from Prob. 3 to Prob. 2 in the above proof relies on the fact that
 213 $M = 2$. Two comments are due here. First, the fact that Prob. 2 is NP-hard with $M = 2$ implies
 214 that Prob. 2 is NP-hard with any $M \geq 2$. Indeed, if Prob. 2 can be solved in polynomial time for
 215 some $M = \hat{M} > 2$, then one can add spurious data points (e.g., at a far distance of the original
 216 data points) to enforce the value of $\hat{M} - 2$ affine pieces of the PWA function. The satisfiability
 217 of Prob. 2 with $M = \hat{M}$ and the augmented data set is then equivalent to the satisfiability of
 218 Prob. 2 with $M = 2$ and the original data set. Second, given $\hat{M} \geq 2$ and any template p , a
 219 construction similar to the one used in the above proof can be used to reduce Prob. 3 to Prob. 2
 220 at the cost of introducing a small gap in the reduction. Indeed, fix $\lambda > 0$ and consider the data
 221 set $\mathcal{D}' = \bigcup_{t=1}^{M+1} \{([x_k; t\lambda\chi_k], y_k)\}_{k=1}^N$. Then, one can show that if Prob. 2 with \mathcal{D}' , p , $\epsilon = \hat{\epsilon}(1 - \frac{2}{\lambda})$
 222 and $M = \hat{M}$ has a solution, then Prob. 3 with \mathcal{D} , $\epsilon = \hat{\epsilon}$ and $M = \hat{M}$ has a solution. The gap
 223 corresponds to the factor $1 - \frac{2}{\lambda}$, which can be made arbitrarily close to one. \triangleleft

224 So, we showed that Prob. 2 is NP-hard, thereby implying no known algorithm with complexity
 225 polynomial in the problem input size $N(d + e) + h$. Nevertheless, one can show that for fixed
 226 dimension d , template size h and number of pieces M , the complexity is polynomial in the data size
 227 N .

228 **Theorem 2** (Polynomial complexity in N). *For fixed dimension d , template $p : \mathbb{R}^d \rightarrow \mathbb{R}^h$ and*
 229 *number of pieces M , the complexity of Prob. 2 is bounded by $O(N^{Mh})$.*

230 The following notation will be useful in the proof of Theorem 2 below and also later in the paper.
 231 For every $c \in \mathbb{R}^h$, let $I(c) = \{k \in [N] : x_k \in H(c)\}$ be the set of all indices k such that $x_k \in H(c)$.
 232 Also, let $\mathcal{I} = \{I(c) : c \in \mathbb{R}^h\}$ be the set of all such index sets.

233 *Proof.* The crux of the proof is to realize that $|\mathcal{I}| \leq N^h + 1$.

234 For every $c \in \mathbb{R}^h$, define $P(c) = \{p(x_k) : k \in [N], p(x_k) \leq c\}$ and let $\mathcal{P} = \{P(c) : c \in \mathbb{R}^h\}$. First,
 235 we prove that $|\mathcal{P}| \leq N^h + 1$. For convenience, we write $p(x) = [p^1(x), \dots, p^h(x)]$. Note that for any
 236 $c \in \mathbb{R}^h$ such that $P(c) \neq \emptyset$, it holds that $P(c) = P(\hat{c})$ where $\hat{c} = \max(\{p(x_k) : k \in [N], p(x_k) \leq c\})$
 237 and the ‘‘max’’ is elementwise. Therefore, we can identify at most h elements x_{k_1}, \dots, x_{k_h} wherein
 238 $k_1, \dots, k_h \in [N]$ such that $\hat{c} \in \{[p^1(x_{k_1}), \dots, p^h(x_{k_h})] : k_1, \dots, k_h \in [N]\}$. Each element can be seen
 239 as ‘‘fixing’’ the maximum value along some dimension of $p(x)$. Hence, there are at most N^h distinct
 240 such \hat{c} . This implies that there are at most N^h distinct nonempty sets $P(c)$, concluding the proof
 241 that $|\mathcal{P}| \leq N^h + 1$.

242 Next, observe that there is a one-to-one correspondence between \mathcal{P} and \mathcal{I} given by: $P(c) \mapsto I(c)$.
 243 Indeed, it is clear that if $I(c_1) = I(c_2)$, then $P(c_1) = P(c_2)$. On the other hand, if $I(c_1) \not\subseteq I(c_2)$,

244 then there is at least one k such that $p(x_k) \leq c_1$ but $p(x_k) \not\leq c_2$. This implies that $P(c_1) \not\subseteq P(c_2)$.
 245 Therefore, $|\mathcal{P}| = |\mathcal{I}| \in O(N^h)$.

246 Now, Prob. 2 can be solved by enumerating the $L = N^h$ nonempty index sets I_1, \dots, I_L in \mathcal{I} ,
 247 and keeping only those I_ℓ for which we can fit an affine function over the data $\{(x_k, y_k)\}_{k \in I_\ell}$ with
 248 error bound ϵ . Next, we enumerate all combinations of M such index sets that cover the indices
 249 $[N]$. There are at most L^M such combinations. This concludes the proof of the theorem. \square

250 *Remark 2.* Note that a similar result holds for Prob. 3 [18, Theorem 5.4]. The proof of Theorem 2 is
 251 however simpler than that in [18] because in our case we can use the template to build the different
 252 regions. \triangleleft

253 The algorithm presented in the proof of Theorem 2, although polynomial in the size of the data
 254 set, can be quite expensive in practice. For instance, in dimension $d = 2$, with rectangular regions
 255 (i.e., $h = 4$) and $N = 100$ data points, one would need to solve $N^h = 10^8$ regression problems,³ each
 256 of which is a linear program.

257 In the next section, we present an algorithm for TPWA regression that is generally several orders
 258 of magnitude faster by using a *top-down* approach that avoids having to systematically enumerate
 259 all the elements in \mathcal{I} .

260 4 Top-Down Algorithm

261 We first define the concept of compatible and maximal compatible index sets. We will assume an
 262 instance of Prob. 2 with data $\{(x_k, y_k)\}_{k=1}^N$, template p , and error bound ϵ .

263 **Definition 1** (Maximal compatible index set). An index set $I \subseteq [N]$ is *compatible* if (a) $I \in \mathcal{I}$ and
 264 (b) there is an affine function $f(x) = Ax + b$ such that $\forall k \in I, \|y_k - f(x_k)\|_\infty \leq \epsilon$. A compatible
 265 index set I is *maximal* if there is no compatible index set I' such that $I \subsetneq I'$.

266 The key idea is that we can restrict ourselves to searching for *maximal* compatible index sets in
 267 order to find a solution to Prob. 2.

268 **Lemma 1.** *Let M be given. Prob. 2 has a solution if and only if it has a solution wherein the*
 269 *regions correspond to maximal compatible index sets.*

270 *Proof.* The “if direction” is clear. We prove the “only if direction”. Consider a solution of Prob. 2
 271 with regions H_1, \dots, H_M . For each $i \in [M]$, there is a maximal compatible index set $I_i = I(c_i)$ such
 272 that $H_i \cap \{x_k\}_{k=1}^N \subseteq H(c_i)$. Since $\{x_k\}_{k=1}^N \subseteq \bigcup_{i=1}^M H_i$, it holds that $\{x_k\}_{k=1}^N \subseteq \bigcup_{i=1}^M H(c_i)$. Hence,
 273 $H(c_1), \dots, H(c_M)$, along with affine functions $f_i(x) = A_i x + b_i$ satisfying (b) in Def. 1, provide a
 274 solution to Prob. 2, concluding the proof. \square

275 The main result of this section is that maximal compatible index sets can be computed by
 276 using a recursive *top-down* approach as follows (implemented in Algo. 1). Consider the lattice
 277 $\langle \mathcal{I}, \subseteq \rangle$ consisting of elements of \mathcal{I} ordered by set inclusion. Our algorithm starts at the very top
 278 of this lattice with a set of points $I = [N]$ and descends until we find maximal compatible index
 279 sets. At each step, we consider a current set $I \in \mathcal{I}$ (initially, $I = [N]$) that is a candidate for being
 280 compatible and check it for compatibility. If I is not compatible, we find subsets $I_1, \dots, I_S \subsetneq I$
 281 using the FINDSUBSETS procedure, which is required to be *consistent*, as defined below:

³In theory, by using Sauer–Shelah’s lemma (see, e.g., [14, Lemma 6.2.2]), this number can be reduced to $\sum_{i=1}^h \binom{K}{i} \approx 4 \times 10^6$. This is because the *VC dimension* of rectangular regions in \mathbb{R}^d is $2d$.

Algorithm 1: Top-down algorithm to compute maximal compatible index sets.

Data: Data set $\{(x_k, y_k)\}_{k=1}^N$, template p
Result: Collection \mathcal{S} of all maximal compatible index sets

```

1  $\mathcal{S} \leftarrow \emptyset$  // compatible sets so far
2  $\mathcal{U} \leftarrow \{[N]\}$  // sets to be processed
3  $\mathcal{V} \leftarrow \emptyset$  // already explored
4 while  $\mathcal{U} \neq \emptyset$  do
5   Pop an index set  $I$  from  $\mathcal{U}$  // removes  $I$  from  $\mathcal{U}$ 
6   if  $I$  is a subset of a set of  $\mathcal{S}$  then
7     // do nothing
8   else if  $I$  is compatible then
9     Add  $I$  to  $\mathcal{S}$ 
10    Remove strict subsets of  $I$  from  $\mathcal{S}$  // Removed sets are not maximal
11  else
12     $(I_1, \dots, I_S) \leftarrow \text{FINDSUBSETS}(I)$  // must be consistent
13    Add to  $\mathcal{U}$  all  $I_1, \dots, I_S$  that are not in  $\mathcal{V}$ 
14  end
15  Add  $I$  to  $\mathcal{V}$ 
16 end
17 return  $\mathcal{S}$ 

```

282 **Definition 2** (Consistency). Given a non-compatible index set $I \in \mathcal{I}$, a collection of index sets
283 $I_1, \dots, I_S \in \mathcal{I}$ is said to be *consistent* w.r.t. I if (a) for each s , $I_s \subsetneq I$ and (b) for every compatible
284 index set $J \subseteq I$, there is s such that $J \subseteq I_s$.

285 We will assume that the procedure FINDSUBSETS is implemented such that for any non-
286 compatible index set I , the collection of sets I_1, \dots, I_S output by FINDSUBSETS(I) is consistent
287 w.r.t. I .

288 **Theorem 3** (Correctness of Algo. 1). *If FINDSUBSETS satisfies that for every non-compatible index*
289 *set $I \in \mathcal{I}$, the output of FINDSUBSETS(I) is consistent w.r.t. I , then Algo. 1 always terminates and*
290 *the output \mathcal{S} contains all the maximal compatible index sets.*

291 *Proof.* Termination follows from the fact that each index set $I \in \mathcal{I}$ is picked at most once, because
292 when some $I \in \mathcal{I}$ is picked, it is then added to the collection \mathcal{V} of visited index sets, so that it
293 cannot be added a second time to \mathcal{U} (line 12). Since \mathcal{I} is finite, this implies that the algorithm
294 terminates in a finite number of steps.

295 Now, we prove that, upon termination, any maximal compatible index set is in the output \mathcal{S} of
296 the algorithm. Therefore, let J be a maximal compatible index set. Then, among all sets I picked
297 during the execution of the algorithm and satisfying $J \subseteq I$, let I^* have minimal cardinality. Such
298 an index set exists since $J \subseteq [N]$. We will show that:

299 *Main result.* $I^* = J$.

300 *Proof of main result.* For a proof by contradiction, assume that $I^* \neq J$. Since J is maximal and
301 $J \subsetneq I^*$, I^* is not compatible. Hence, the index sets $(I_1, \dots, I_S) = \text{FINDSUBSETS}(I^*)$ were added
302 to \mathcal{U} (line 11). Using the assumption on FINDSUBSETS, let s be such that $J \subseteq I_s \subsetneq I^*$. Since I_s
303 must have been picked during the execution of the algorithm, this contradicts the minimality of the
304 cardinality of I^* , concluding the proof of the main result.

Algorithm 2: An implementation of FINDSUBSETS using infeasibility certificates

Data: Data set $\{(x_k, y_k)\}_{k=1}^N$, template $p = [p^1, \dots, p^h]$, non-compatible index set $I = I(c)$ where $c = [c^1, \dots, c^h]$, infeasibility certificate $C \subseteq I$

Result: A collection of index sets I_1, \dots, I_S consistent w.r.t. I

```

1 foreach  $s = 1, \dots, h$  do
2    $\hat{c}^s \leftarrow \max \{p^s(x_k) : k \in I, p^s(x_k) < \max_{\ell \in C} p^s(x_\ell)\}$ 
3   Define  $I_s = I([c^1, \dots, c^{s-1}, \hat{c}^s, c^{s+1}, \dots, c^h])$ 
4 end
5 return all nonempty index sets  $I_1, \dots, I_h$ 

```

305 Thus, J was picked during the execution of the algorithm. Since it is compatible, it was added
 306 to \mathcal{S} at the iteration at which it was picked (line 8), and since it is maximal, it is not removed at
 307 later iterations. Hence, upon termination, $J \in \mathcal{S}$. Since J was arbitrary, we conclude that upon
 308 termination, \mathcal{S} contains all maximal compatible index sets.

309 Finally, we show that upon termination, \mathcal{S} contains only maximal compatible index sets. This
 310 follows from the fact that, at each iteration of the algorithm, for any distinct $I_1, I_2 \in \mathcal{S}$, it holds that
 311 $I_1 \not\subseteq I_2$ and $I_2 \not\subseteq I_1$. Indeed, when I_1 is added to \mathcal{S} , all subsets of I_1 are removed from \mathcal{S} (line 9)
 312 and are ignored at later iterations (line 6). The same holds for I_2 . This concludes the proof of the
 313 theorem. \square

314 4.1 Implementation of FindSubsets using infeasibility certificates

315 We now explain how to implement FINDSUBSETS so that it is consistent. For that, we use infeasibility
 316 certificates, which are index sets that are not compatible.

317 **Definition 3** (Infeasibility certificate). An index set $C \subseteq [N]$ is an *infeasibility certificate* if there
 318 is no affine function $f(x) = Ax + b$ such that $\forall k \in C, \|y_k - f(x_k)\|_\infty \leq \epsilon$.

319 Note that any incompatible index set I contains an infeasibility certificate $C \subseteq I$ (e.g., $C = I$).
 320 However, it is quite useful to extract an infeasibility certificate C that is as small as possible. We
 321 explain below how to compute small infeasibility certificates. Thereafter, from an infeasibility
 322 certificate $C \subseteq I$, one can compute a consistent collection of index subsets of I by tightening each
 323 component of the template *independently*, in order to exclude a minimal nonzero number of indices
 324 from the infeasibility certificate, while keeping the other components unchanged. This results in
 325 an implementation of FINDSUBSETS that satisfies the consistency property, described in Algo. 2.
 326 Fig. 3 below shows an illustration for rectangular regions.

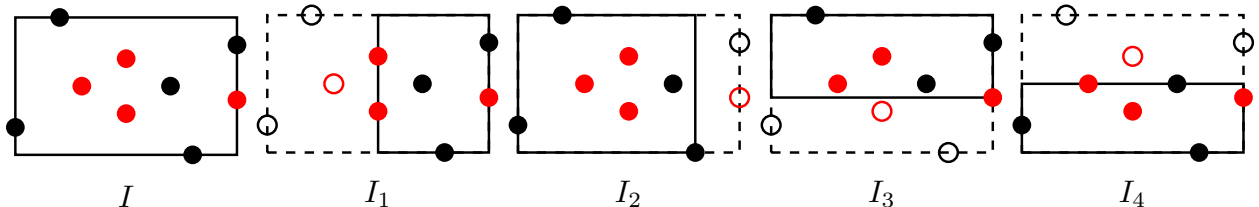


Figure 3: FINDSUBSETS implemented by Algo. 2 with rectangular regions. The red dots represent the infeasibility certificate C . Each I_s excludes at least one point from C by moving one face of the box but keeping the others unchanged.

327 **Theorem 4** (Correctness of Algo. 2). *For every non-compatible index set $I \in \mathcal{I}$, the output I_1, \dots, I_S*
 328 *of Algo. 2 is consistent w.r.t. I .*

329 *Proof.* Observe that if C is an infeasibility certificate, then every $I \subseteq [N]$ satisfying $C \subseteq I$ is not
 330 compatible. Now, let $J \subseteq I$ be compatible. Using that $C \not\subseteq J$, let s be a component such that
 331 $\max_{k \in J} p^s(x_k) < \max_{k \in C} p^s(x_k)$. It holds that $J \subseteq I_s$. Since J was arbitrary, this concludes the
 332 proof. \square

333 4.2 Finding infeasibility certificates

334 We outline the process of finding an infeasibility certificate $C \subseteq I$ for a given non-compatible index
 335 set $I \subseteq [N]$. Recall that the data is of the form $\{(x_k, y_k)\}_{k=1}^N$, wherein for each $k \in [N]$, $x_k \in \mathbb{R}^d$
 336 and $y_k \in \mathbb{R}^e$. For simplicity, assume that the output is scalar, i.e. $e = 1$, or equivalently, $y_k \in \mathbb{R}$ for
 337 all $k \in [N]$. We will subsequently show how technique for the scalar case will extend to the case of
 338 $e > 1$.

339 For the scalar case, the goal is to find an affine function $f(x) = a^\top x + b$ wherein $a \in \mathbb{R}^d$ and
 340 $b \in \mathbb{R}$ so that for all $k \in I$, $|y_k - f(x_k)| \leq \epsilon$. However, since I is non-compatible, no such function
 341 exists by definition. Therefore, the system of linear inequalities involving unknowns $(a, b) \in \mathbb{R}^d \times \mathbb{R}$
 342 is infeasible:

$$\begin{aligned} x_k^\top a + b - y_k &\leq \epsilon \quad \forall k \in I, \\ -x_k^\top a - b + y_k &\leq \epsilon \quad \forall k \in I. \end{aligned}$$

343 Note that each constraint of the form $|\alpha| \leq \beta$ is expanded as two constraints $\alpha \leq \beta$ and $-\alpha \leq \beta$.
 344 By applying Farkas' Lemma or a theorem of the alternative [30, Theorem 21.3], and simplifying the
 345 result, we conclude that the system of inequalities above is infeasible (i.e. I is non-compatible) if
 346 and only if there exists a multiplier $\lambda_k \in \mathbb{R}$ for each $k \in I$ such that the following system of linear
 347 constraints is feasible:

$$\begin{aligned} \sum_{k \in I} \lambda_k x_k &= [0]_d, \\ \sum_{k \in I} \lambda_k &= 0, \\ \sum_{k \in I} \lambda_k y_k + \sum_{k \in I} |\lambda_k| \epsilon &\leq -1. \end{aligned} \tag{3}$$

348 Thus, in order to check whether a given index set I is non-compatible, we simply formulate
 349 the system of inequalities (3) involving unknowns $\lambda_k \in \mathbb{R}$ for each $k \in I$ and attempt to find
 350 a feasible solution using a Linear Programming (LP) solver. If feasible, we conclude that I is
 351 non-compatible. In fact, given any solution $(\lambda_k)_{k \in I}$ to the system of inequalities above, we can
 352 extract a corresponding infeasibility certificate as $C = \{k \in I : \lambda_k \neq 0\}$. It is easy to see why: any
 353 $k \in I$ with $\lambda_k = 0$ indicates that $I \setminus \{k\}$ remains a non-compatible set since such a variable λ_k can
 354 be removed from the system of inequalities while retaining feasibility.

355 **Lemma 2.** *If I is non-compatible, then there exists an infeasibility certificate $C \subseteq I$ such that*
 356 *$|C| \leq d + 2$.*

357 *Proof.* The system of constraints (3) has $|I|$ unknowns and $d + 2$ constraints of which $d + 1$ are
 358 equality constraints. We may write $\lambda_k = \alpha_k - \beta_k$ for variables $\alpha_k, \beta_k \geq 0$ and therefore $|\lambda_k| = \alpha_k + \beta_k$.
 359 Under this transformation, the system of constraints can be rewritten as

$$\begin{aligned} \sum_{k \in I} (\alpha_k - \beta_k) x_k &= [0]_d, \\ \sum_{k \in I} (\alpha_k - \beta_k) &= 0, \\ \sum_{k \in I} (\alpha_k - \beta_k) y_k + \sum_{k \in I} (\alpha_k + \beta_k) \epsilon &\leq -1, \\ \alpha_k, \beta_k &\geq 0, \quad \forall k \in I. \end{aligned}$$

360 With this transformation, the system above is a standard-form Linear Program with $d+2$ constraints
361 and decision variables $\alpha_k, \beta_k \geq 0$ for each $k \in I$. We treat the objective function for this LP as the
362 constant 0. Any basic feasible solution will have at most $d+2$ non-zero variables. Furthermore,
363 from the theory of Linear Programming, we know that if a system is feasible, then it has a basic
364 feasible solution [34]. Translating this back to the original system (3), we get that there is a solution
365 involving at most $d+2$ non-zero values for λ_k . That is, there exists an infeasibility certificate C
366 such that $|C| \leq d+2$. \square

367 So far, we have assumed that the outputs y_k are scalar, i.e., $e = 1$. However, if $e > 1$, we can
368 simply use the previous analysis by focusing separately on each component of the output vectors y_k .
369 This is possible because if a given index set I is non-compatible, then there must exist a component
370 $s \in [e]$ such that I is non-compatible even if the data set is restricted to $\{(x_k, y_k^s)\}_{k=1}^N$, wherein
371 y_k^s denotes the s^{th} component of y_k . Indeed, if it were not the case, then for each component s ,
372 we could find an affine function $f_s(x) = a_s^\top x + b_s$ such that for all $k \in I$, $|f_s(x_k) - y_k^s| \leq \epsilon$. By
373 letting $A = [a_1^\top; \dots; a_e^\top]$ and $b = [b_1, \dots, b_e]^\top$, we see that $f(x) = Ax + b$ satisfies that for all
374 $k \in I$, $\|f(x_k) - y_k\|_\infty \leq \epsilon$. This contradicts the assumption that I is non-compatible. Thus, the
375 infeasibility certificate for $e > 1$ is extracted by simply considering each output component in turn,
376 thus reducing the problem to the scalar case considered above.

377 We conclude that checking whether a set I is non-compatible and if so, finding an infeasibility
378 certificate $C \subseteq I$, can be solved by posing the system of constraints (3) and solving it using an
379 algorithm such as the simplex algorithm.

380 4.2.1 Good infeasibility certificates for the top-down approach

381 The implementation of FINDSUBSETS boils down to finding infeasibility certificates, which can be
382 done as explained above. However, not all certificates will be as good in terms of overall complexity
383 of the top-down approach. To exclude non-compatible index sets more rapidly, it is desirable that
384 the points in the certificate are “spatially concentrated” in the input domain. This means that the
385 points $\{x_k\}_{k \in C}$ are close to each other w.r.t. some distance metric.

386 We illustrate the benefit of spatially concentrated certificates with the example used to illustrate
387 the top-down approach in the introduction.

388 *Example 1.* Consider the TPWA regression problem in Fig. 2, introduced in the section “Approach at
389 a glance”. In Plot II, we obtained the infeasibility certificate $C_0 = \{4, 5, 6\}$. Note that $\hat{C}_0 = \{1, 5, 11\}$
390 is another infeasibility certificate that could have been obtained as a result of solving the system (3).
391 However, C_0 is more “spatially concentrated” in the sense that the points in C_0 are closer to each
392 other than those in \hat{C}_0 .

393 Recall that using C_0 as the infeasibility certificate allowed to find the compatible subset
394 $I_1 = \{1, \dots, 5\}$ directly and the compatible subsets $I_3 = \{5, 6, 7\}$ and $I_4 = \{7, \dots, 11\}$ at the
395 subsequent steps.

396 However, if the certificate $\hat{C}_0 = \{1, 5, 11\}$ were used, then we would have obtained $\hat{I}_1 = \{2, \dots, 11\}$
397 and $\hat{I}_2 = \{1, \dots, 10\}$. Note that both \hat{I}_1, \hat{I}_2 are non-compatible because they contain C_0 and thus
398 further steps of our procedure are needed until we find maximally compatible sets. \triangleleft

399 We now refine the above argument with a volume-contraction argument to discuss what would
400 be the complexity of the overall top-down algorithm if all certificates are spatially concentrated in
401 the input domain.

402 *Example 2.* Consider a PWA function f with M pieces whose domains H_1, \dots, H_M are rectangles,
403 as illustrated in Fig. 4. Let $N \in \mathbb{N}_{>0}$ and consider a data set $\mathcal{D} = \{(x_k, y_k)\}_{k=1}^N$. We aim to solve

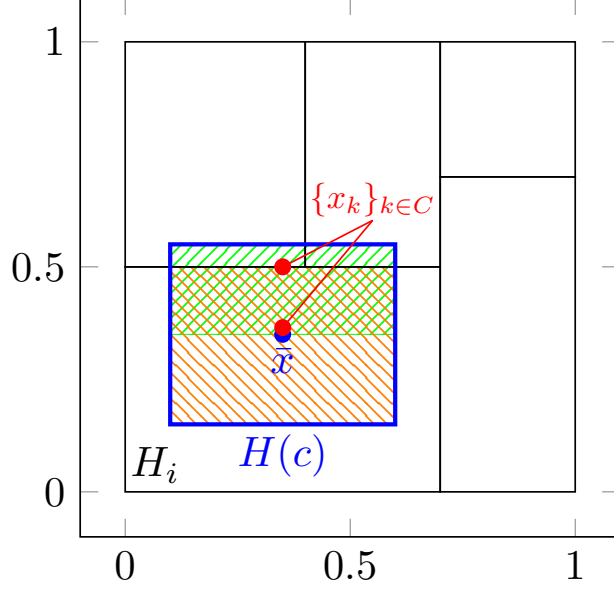


Figure 4: Illustration of FINDSUBSETS with a spatially concentrated certificate. The green and orange hatched rectangles illustrate two possible cases for $H(c_s)$ output by FINDSUBSETS.

404 Prob. 2 with \mathcal{D} , $\epsilon = 0$ and the rectangular template. We will discuss the complexity of the top-down
 405 approach presented in Algo. 1 if all certificates are spatially concentrated. This means that $\{x_k\}_{k \in C}$
 406 consists in $d + 2$ points concentrated around the center \bar{x} of $H(c)$, where C is the certificate for $I(c)$
 407 (see Fig. 4).

408 This will imply that the rectangles $H(c_1), \dots, H(c_h)$ computed by FINDSUBSETS satisfy that
 409 for all $s \in [h]$, either the volume of $H(c_s)$ is half of that of $H(c)$ since one face is tight at \bar{x} (see the
 410 green rectangle in Fig. 4) or $H(c_s)$ has one more face near the boundary of H_i compared to $H(c)$
 411 (see the orange rectangle in Fig. 4). By adding the natural assumption that all regions H_j have a
 412 volume of at least $\nu \in (0, 1]$ and discarding regions with volume smaller than ν , we get that the
 413 algorithm cannot divide the volume of a region more than $-\log_2(\nu)$ times. Hence, the depth of the
 414 tree underlying the algorithm is upper bounded by $h - \log_2(\nu)$. Since, each node of the tree has at
 415 most h children (the subsets given by FINDSUBSETS), the number of rectangles encountered during
 416 the algorithm is upper bounded by $h^{h - \log_2(\nu)}$. Note that this upper bound on the complexity of the
 417 algorithm is independent of the data size N . This concludes the example. \triangleleft

418 To conclude this section on good certificates, we explain briefly how spatially concentrated
 419 certificates can be computed by adding a cost function to the Linear Program (3). For simplicity, we
 420 will assume that the output dimension $e = 1$. For the case when $e > 1$, we can apply our approach
 421 to each component of the output in turn.

422 Given a center point $\bar{x} = \frac{1}{N} \sum_{i \in I} x_i$ for a non-compatible index set $I \subseteq [N]$, we consider the
 423 following Linear Program with variables $\lambda_k \in \mathbb{R}$ for each $k \in I$,

$$\begin{aligned}
 & \text{minimize} && \sum_{k \in I} |\lambda_k| \|x_k - \bar{x}\|^2 \\
 & \text{s.t.} && \sum_{k \in I} \lambda_k x_k = [0]_d, \\
 & && \sum_{k \in I} \lambda_k = 0, \\
 & && \sum_{k \in I} \lambda_k y_k + \sum_{k \in I} |\lambda_k| \epsilon \leq -1.
 \end{aligned} \tag{4}$$

424 The objective function of (4) tends to put zero value to λ_k whenever $\|x_k - \bar{x}\|$ is large. This promotes

Algorithm 3: Extra step at the beginning of each iteration of Algo. 1

Data: \mathcal{S} and \mathcal{U} at the iteration, N

Result: BREAK if we can extract from \mathcal{S} an optimal cover of $[N]$ with compatible index sets; otherwise, CONTINUE

- 1 Let α be the size of an optimal cover of $[N]$ by index sets in \mathcal{S} // **Note:** $\alpha = \infty$ if $[N]$ cannot be covered by \mathcal{S} .
 - 2 Let β be the size of an optimal cover of $[N]$ by index sets in $\mathcal{S} \cup \mathcal{U}$
 - 3 **if** $\alpha \leq \beta$ **then return** BREAK **else return** CONTINUE
-

Algorithm 4: Top-down algorithm for Prob. 2.

// Lines 1--3 from Algo. 1

1 **while** *true* **do**

2 | **if** *Algo. 3 outputs* BREAK **then return** an optimal cover of $[N]$ using index sets from \mathcal{S}
| // Lines 5--14 from Algo. 1

3 **end**

425 proximity of the point x_k to \bar{x} when $\lambda_k \neq 0$.⁴

426 4.3 Early stopping using set cover algorithms

427 Finally, Algo. 1 can be made more efficient by enabling early termination if $[N]$ is optimally covered
428 by the compatible index sets computed so far. For that, we add an extra step at the beginning of
429 each iteration, that consists in (i) computing a lower bound β on the size of an optimal cover of $[N]$
430 with compatible index sets; and (ii) checking whether we can extract from \mathcal{S} a collection of β index
431 sets that form a cover of $[N]$. The extra step returns BREAK if (ii) is successful. An implementation
432 of the extra step is provided in Algo. 3.

433 The soundness of Algo. 3 follows from the following lemma.

434 **Lemma 3.** *Let β be as in Algo. 3. Then, any cover of $[N]$ with compatible index sets has size at*
435 *least β .*

436 *Proof.* The crux of the proof relies on the observation from the proof of Theorem 3 that for any
437 compatible index set $I \in \mathcal{I}$, there is $J \in \mathcal{S} \cup \mathcal{U}$ such that $I \subseteq J$. It follows that for any cover of $[N]$
438 with M compatible index sets, there is a cover of $[N]$ with M index sets in $\mathcal{S} \cup \mathcal{U}$. Since β is the
439 smallest size of such a cover, this concludes the proof. \square

440 The implementation of the extra step in Algo. 1 is provided in Algo. 4. The correctness of the
441 algorithm follows from that of Algo. 1 (Theorem 3) and Algo. 3 (Lemma 3). In conclusion, we
442 provided an algorithm for optimal TPWA regression.

443 **Theorem 5** (Optimal TPWA regression). *Algo. 4 solves Prob. 2 with minimal M .*

444 *Proof.* Let I_1, \dots, I_M be the output of Algo. 4. For each $i \in [M]$, let $H_i = H(c_i)$ where $I_i = I(c_i)$
445 and let $f_i(x) = A_i x + b_i$ be as in (b) of Def. 1. The fact that H_1, \dots, H_M and f_1, \dots, f_M is a solution
446 to Prob. 2 follows from the fact that I_1, \dots, I_M is a cover of $[N]$ and the definition of f_1, \dots, f_M .

⁴Note that L^1 regularization costs are often used in machine learning to induce sparsity of the optimal solution [8, p. 304]. Here, we use a *weighted* L^1 regularization cost to induce a sparsity pattern dictated by the geometry of the problem.

447 The fact that it is a solution with minimal M follows from the optimality of I_1, \dots, I_M among all
 448 covers of $[N]$ with compatible index sets. \square

449 *Remark 3.* To solve the optimal set cover problems (known to be NP-hard in general) in Algo. 3, we
 450 use MILP formulations. The complexity of solving these MILPs grows as $2^{|\mathcal{S}|}$ and $2^{|\mathcal{S} \cup \mathcal{U}|}$ respectively.
 451 However, in our numerical experiments (see next section), we observed that the gain of stopping the
 452 algorithm early if an optimal cover is found systematically outbalanced the computational cost of
 453 solving the set cover problems. Furthermore, if one is satisfied with a sub-optimal solution, they
 454 can use an approximation algorithm, such as the greedy algorithm, which outputs a cover whose
 455 size is within some factor $t(N) \geq 1$ of the optimal set cover size [11]. In this case, Algo. 3 outputs
 456 BREAK if $\alpha \leq t(N)\beta$. \triangleleft

457 5 Numerical Experiments

458 In this section, we demonstrate the applicability of our algorithm on three numerical examples.⁵
 459 We also compare it with the MILP and PARC [4] approaches to solve SA and PWA regression, and
 460 we discuss the impact of different templates in terms of simplicity of the model and efficiency of the
 461 algorithm.

462 5.1 PWA approximation of insulin–glucose regulation model

463 Dalla Man, Rizza, and Cobelli [12] present a nonlinear model of insulin–glucose regulation that has
 464 been widely used to test artificial pancreas devices for treatment of type-1 diabetes. The model is
 465 nonlinear and involves 10 state variables. However, the nonlinearity arises mainly from the term
 466 U_{id} (insulin-dependent glucose utilization) involving two state variables, say x_1 and x_2 (namely, the
 467 level of insulin in the interstitial fluid, and the glucose mass in rapidly equilibrating tissue):

$$U_{\text{id}}(x_1, x_2) = \frac{(3.2667 + 0.0313x_1)x_2}{253.52 + x_2}.$$

468 We consider the problem of approximating U_{id} with a PWA model, thus converting the entire model
 469 into a PWA model. Therefore, we simulated trajectories and collected $N = 100$ values of x_1 , x_2 and
 470 $U_{\text{id}}(x_1, x_2)$; see Fig. 5(a). For three different values of the error tolerance, $\epsilon \in \{0.2, 0.1, 0.05\}$, we
 471 used Algo. 4 to compute a PWA regression of the data with rectangular domains. The results of the
 472 computations are shown in Fig. 5(b,c,d). The computation times are respectively 1, 22 and 112 secs.
 473 Finally, we evaluate the accuracy of the PWA regression for the modeling of the glucose-insulin
 474 evolution by simulating the system with U_{id} replaced by the PWA models. The results are shown in
 475 Fig. 5(e,f). We see that the PWA model with $\epsilon = 0.05$ induces a prediction error less than 2% over
 476 the whole simulation interval, which is a significant improvement compared to the PWA models
 477 with only 1 affine piece ($\epsilon = 0.2$) or 2 affine pieces ($\epsilon = 0.1$).

478 Finally, we compare with SA regression and classical PWA regression. To find a SA model,
 479 we solved Prob. 3 with $\epsilon = 0.05$ and $M = 3$ using a MILP approach. The computation is very
 480 fast (< 0.5 secs); however, the computed clusters of data points (see Fig. 6) do not allow to learn
 481 a (simple) PWA model, thereby hindering the derivation of a model for U_{id} that can be used for
 482 simulation and analysis.

⁵The implementation is made in Julia, with Gurobi 11.0, under academic license, as LP and MILP solver (including for the optimal set cover problems). Our approach uses the standard set data structures available in Julia for manipulating index sets in order to implement the key steps of Algorithm 1. All computations were made on a laptop with processor Intel Core i7-7600u and 16GB RAM running Windows.

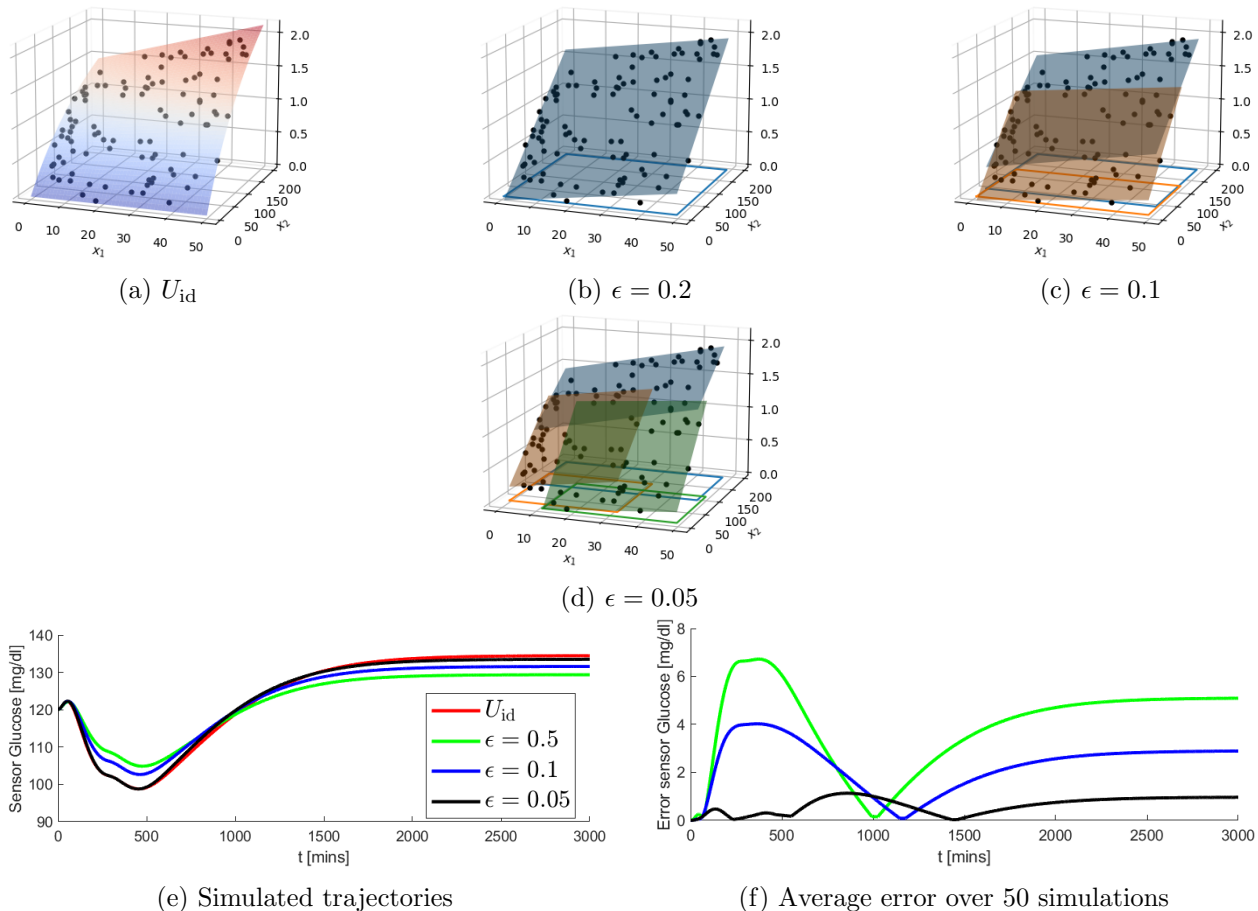


Figure 5: Glucose–insulin system. (a) 100 sampled points (black dots) on the graph of U_{id} (surface). (b,c,d) Optimal TPWA regression for various error tolerances ϵ . (e) Simulations using the nonlinear model versus the PWA approximations. (f) Error between nonlinear and PWA models averaged over 50 simulations with different initial conditions.

483 5.2 Hybrid system identification: double pendulum with soft contacts

484 We consider a hybrid linear system consisting in an inverted double pendulum with soft contacts
 485 at the joints, as depicted in Fig. 7(a). This system has nine linear modes, depending on whether
 486 the contact force of each joint is inactive, active on the left or active on the right [1]. Our goal is
 487 to learn these linear modes as well as their domain of validity, from data. For that, we simulated
 488 trajectories and collected $N = 250$ sampled values of θ_1 , θ_2 and the force applied on the lower joint.
 489 We used Algo. 4 to compute a PWA regression of the data with rectangular domains and with error
 490 tolerance $\epsilon = 0.01$. The result is shown in Fig. 7(b). The number of iterations of the algorithm was
 491 about 23000 for a total time of 800 secs.

492 We see that the affine pieces roughly divide the state space into a grid of 3×3 regions. This is
 493 consistent with our ground-truth model, in which the contact force at each joint has three linear
 494 modes depending only on the angle made at the joint. The PWA regression provided by Algo. 4
 495 allows us to learn this feature of the system from data, without assuming anything about the system
 496 except that the domains of the affine pieces are rectangular.

497 We compare with SA regression and classical PWA regression. The MILP approach to solve the

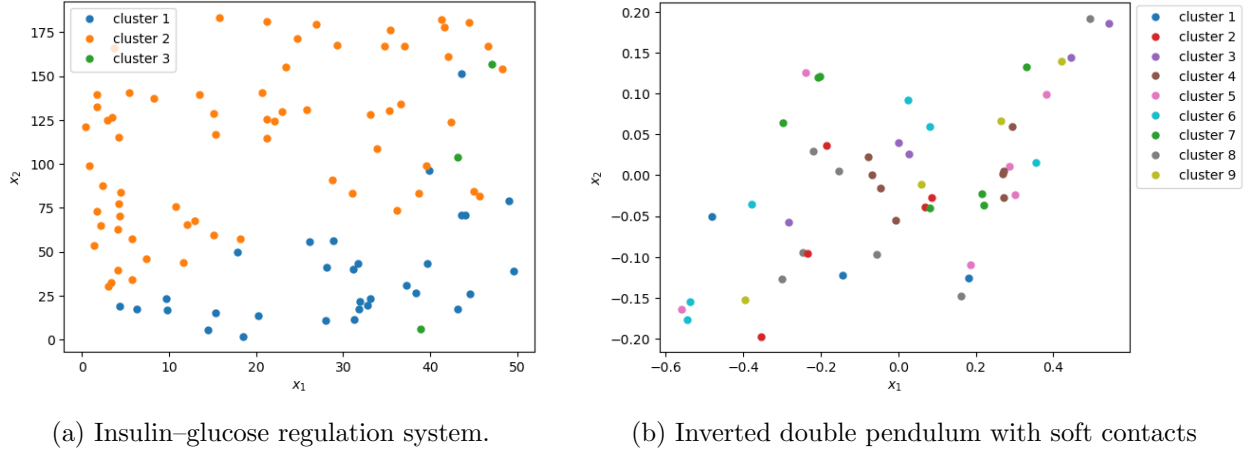


Figure 6: Clusters of data points from SA regression.

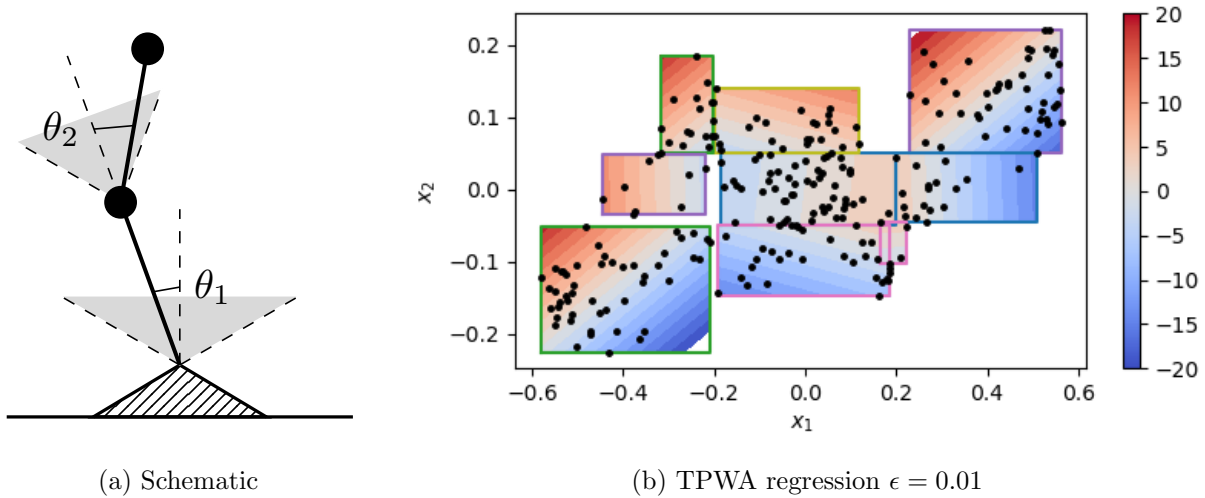
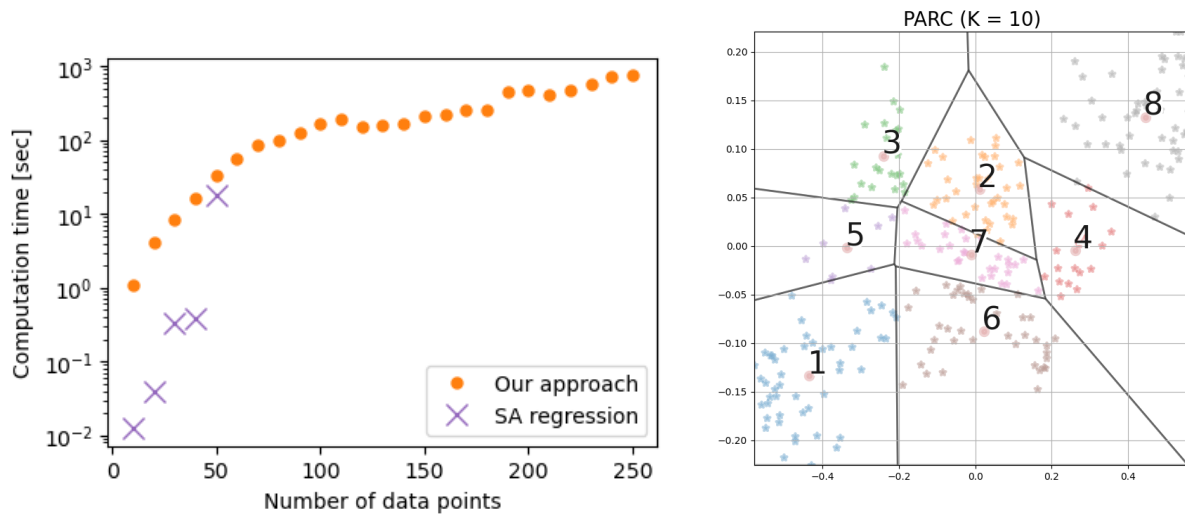


Figure 7: Inverted double pendulum with soft contacts. (a): Elastic contact forces apply when θ is outside gray region, (b): Optimal TPWA regression of the data with rectangular domains.

498 SA regression (Prob. 3) with $\epsilon = 0.01$ and $M = 9$ could not handle more than 51 data points within
 499 reasonable time (1000 secs); see Fig. 8(a). Furthermore, the computed clusters of data points (see
 500 Fig. 6) do not allow to learn a (simple) PWA model, thereby hindering to learn important features
 501 of the system. Last but not least, we compare with the recent tool PARC [4].⁶ The fitting accuracy
 502 on training data is high ($R^2 = 0.995$). The resulting partition of the input space is depicted in
 503 Fig. 8(b). As we can see, PARC finds a PWA function with 8 modes, although an upper bound
 504 (K) of 10 was given. However, the regions do not align with the axis (as this is not enforced by the
 505 algorithm). Consequently, regions with a small number of samples (e.g., lower-right) are missing,
 506 while regions with many samples (e.g., central) are overly divided.

⁶We used the default parameters proposed on the webpage <https://github.com/bemporad/PyPARC>.



(a) Computation times: our approach vs MILP

(b) PWA regression using PARC

Figure 8: (a): Comparison with MILP approach for SA regression. Time limit is set to 1000 secs. (b) Partition of the input space by the PWA function computed using the state-of-the-art tool PARC [4].

5.3 Hybrid system identification: carts with springs

We consider a hybrid linear system consisting in two carts with springs, as depicted in Fig. 9(a). The force applied on the left cart has four linear modes, depending on the values of x_1 and x_2 . Our goal is to learn these linear modes as well as their domain of validity, from data. For that, we used $N = 400$ data obtained by gridding the input domain $[0, L]^2$ uniformly. We used Algo. 4 to compute a PWA regression with error tolerance $\epsilon = 0.01$. We considered two templates: first the “rectangular–octagonal” template wherein each region can have up to eight faces consisting in horizontal, vertical and oblique lines (see Fig. 9(b) for an example); then, we compared with the rectangular template.

The results are shown in Fig. 9(c,d). The running time of the algorithm was about 950 secs for the rectangular–octagonal template, and 120 secs for the rectangular template. It is natural that the rectangular–octagonal template takes more time because we allow for degrees of freedom in the shape of the regions. However, we observe in Fig. 9(c,d) that the most expressive template gives better result in terms of simplicity of the PWA function.

6 Conclusion

To conclude, we have presented an algorithm for fitting piecewise affine models wherein each piece ranges over a region whose shape is dictated by a user-provided template. The complexity of the problem has been analyzed in terms of the number of data points, the dimension of the input domain and the template, and the desired number of pieces of the model. We have presented a top-down algorithm that explores subsets of the data guided by the concept of infeasibility certificates. Finally, our implementation provides some interesting applications of this approach to cyber-physical systems. Despite these contributions, the problem of identifying hybrid systems from data remains a computationally hard problem and the computational challenges of providing precise

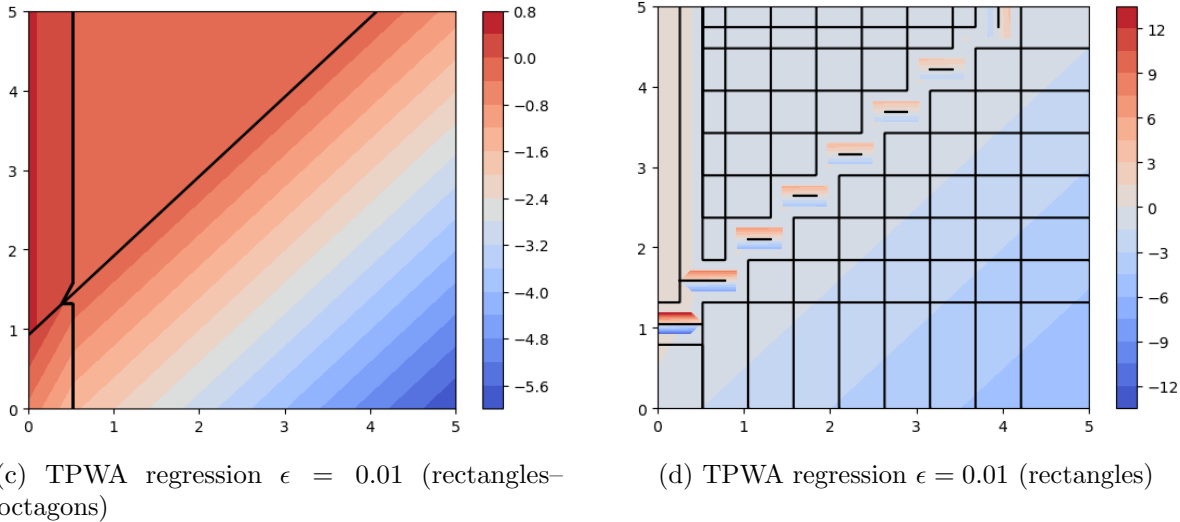
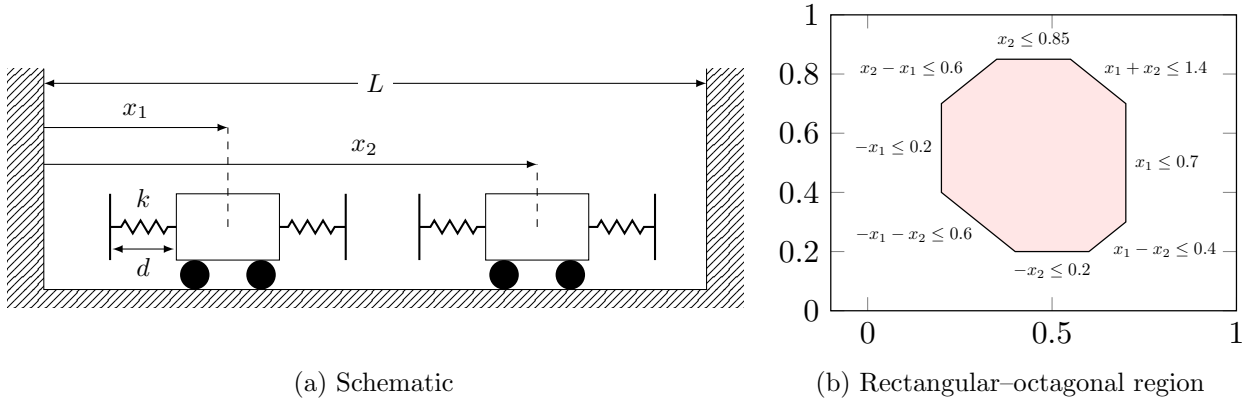


Figure 9: Carts with springs. (a): Elastic contact forces apply when the springs are compressed. (b): Example of rectangular–octagonal region. (c): Optimal TPWA regression of the data with rectangular–octagonal domains. (d): Optimal TPWA regression of the data with rectangular domains.

530 solutions with mathematical guarantees remain formidable. Our future work will investigate the use
 531 of better data structures to help scale our algorithms to larger and higher dimensional data sets. We
 532 are also investigating other approaches to PWA identification involving regions that are separated
 533 by arbitrary hyperplanes rather than fixed templates. Finally, we are interested in connections
 534 between the approach presented here and ideas from computational geometry. In particular, the
 535 link between the VC dimension of the shapes used to specify the regions in our PWA model and the
 536 complexity of the regression procedure offers interesting venues for future work.

537 **Author Contribution** GB and SS conceived the method and the experimental setup, produced
 538 the technical results and wrote the paper together. GB implemented the code and produced the
 539 experimental results.

540 **Funding Statement** This research was supported by grants from the Federation Wallonie–
 541 Bruxelles (WBI) and the US National Science Foundation (NSF) under award numbers 1836900
 542 and 1932189.

543 **Competing Interests** None.

544 **Ethics Statements** No experiments involving humans or animals were conducted in the context
545 of this research.

546 **Data Availability Statement** The code and data used in this paper are available at https://github.com/guberger/L4DC2023_Tool_PWA_Regression_Top-Down.
547

548 **Connections References** Paoletti N, Woodcock J (2023). How to ensure safety of learning-
549 enabled cyber-physical systems? *Research Directions: Cyber-Physical Systems*. **1**, e2, 1–2. <https://doi.org/10.1017/cbp.2023.2>
550

551 **References**

- 552 [1] Alp Aydinoglu, Victor M Preciado, and Michael Posa. “Contact-aware controller design
553 for complementarity systems”. In: *2020 IEEE International Conference on Robotics and
554 Automation (ICRA)*. IEEE, 2020, pp. 1525–1531. DOI: 10.1109/ICRA40945.2020.9197568.
- 555 [2] Richard Bellman. “On the approximation of curves by line segments using dynamic program-
556 ming”. In: *Communications of the ACM* **4.6** (1961), p. 284. DOI: 10.1145/366573.366611.
- 557 [3] Richard Bellman and Robert Roth. “Curve fitting by segmented straight lines”. In: *Journal of
558 the American Statistical Association* **64.327** (1969), pp. 1079–1084. DOI: 10.1080/01621459.
559 1969.10501038.
- 560 [4] Alberto Bemporad. “A piecewise linear regression and classification algorithm with application
561 to learning and model predictive control of hybrid systems”. In: *IEEE Transactions on
562 Automatic Control* (2022). DOI: 10.1109/TAC.2022.3183036.
- 563 [5] Alberto Bemporad et al. “A bounded-error approach to piecewise affine system identification”.
564 In: *IEEE Transactions on Automatic Control* **50.10** (2005), pp. 1567–1580. DOI: 10.1109/
565 TAC.2005.856667.
- 566 [6] Guillaume Berger, Monal Narasimhamurthy, and Sriram Sankaranarayanan. “Algorithms for
567 identifying flagged and guarded linear systems”. In: *Proceedings of the 27th ACM International
568 Conference on Hybrid Systems: Computation and Control*. ACM. 2024, pp. 1–13. DOI: 10.
569 1145/3641513.3650140.
- 570 [7] Guillaume O Berger and Sriram Sankaranarayanan. “Template-based piecewise affine regres-
571 sion”. In: *Proceedings of The 5th Annual Learning for Dynamics and Control Conference*.
572 Ed. by Nikolai Matni, Manfred Morari, and George J Pappas. Vol. 211. Proceedings of Machine
573 Learning Research. 2023, pp. 509–520.
- 574 [8] Stephen Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge, UK: Cambridge
575 University Press, 2004. DOI: 10.1017/CB09780511804441.
- 576 [9] Paul S Bradley and Olvi L Mangasarian. “K-plane clustering”. In: *Journal of Global optimiza-
577 tion* **16.1** (2000), pp. 23–32. DOI: 10.1023/A:1008324625522.
- 578 [10] Leo Breiman. “Hinging hyperplanes for regression, classification, and function approximation”.
579 In: *IEEE Transactions on Information Theory* **39.3** (1993), pp. 999–1013. DOI: 10.1109/18.
580 256506.

- 581 [11] Vasek Chvatal. “A greedy heuristic for the set-covering problem”. In: *Mathematics of Operations*
582 *Research* 4.3 (1979), pp. 233–235. DOI: 10.1287/moor.4.3.233.
- 583 [12] Chiara Dalla Man, Robert A Rizza, and Claudio Cobelli. “Meal simulation model of the glucose-
584 insulin system”. In: *IEEE Transactions on biomedical engineering* 54.10 (2007), pp. 1740–1749.
585 DOI: 10.1109/TBME.2007.893506.
- 586 [13] Giancarlo Ferrari-Trecate et al. “A clustering technique for the identification of piecewise affine
587 systems”. In: *Automatica* 39.2 (2003), pp. 205–217. DOI: 10.1016/S0005-1098(02)00224-8.
- 588 [14] Sarel Har-Peled. *Geometric approximation algorithms*. Vol. 173. Mathematical Surveys and
589 Monographs. Providence, RI: American Mathematical Society, 2011.
- 590 [15] Wanxin Jin et al. “Learning linear complementarity systems”. In: *Proceedings of The 4th*
591 *Annual Learning for Dynamics and Control Conference*. Ed. by Roya Firoozi et al. Vol. 168.
592 Proceedings of Machine Learning Research. 2022, pp. 1137–1149.
- 593 [16] Raphaël M Jungers. *The joint spectral radius: theory and applications*. Berlin: Springer, 2009.
594 DOI: 10.1007/978-3-540-95980-9.
- 595 [17] Fabien Lauer. “Estimating the probability of success of a simple algorithm for switched linear
596 regression”. In: *Nonlinear Analysis: Hybrid Systems* 8 (2013), pp. 31–47. DOI: 10.1016/j.
597 nahs.2012.10.001.
- 598 [18] Fabien Lauer and Gérard Bloch. *Hybrid system identification: theory and algorithms for*
599 *learning switching models*. Cham: Springer, 2019. DOI: 10.1007/978-3-030-00193-3.
- 600 [19] Ramy Medhat et al. “A framework for mining hybrid automata from input/output traces”. In:
601 *2015 International Conference on Embedded Software (EMSOFT)*. IEEE. 2015, pp. 177–186.
602 DOI: 10.1109/EMSOFT.2015.7318273.
- 603 [20] Antoine Miné. “The octagon abstract domain”. In: *Higher-Order and Symbolic Computation*
604 19.1 (2006), pp. 31–100. DOI: 10.1007/s10990-006-8609-1.
- 605 [21] Eberhard Münz and Volker Krebs. “Continuous optimization approaches to the identification
606 of piecewise affine systems”. In: *IFAC Proceedings Volumes* 38.1 (2005), pp. 349–354. DOI:
607 10.3182/20050703-6-CZ-1902.00342.
- 608 [22] Eberhard Münz and Volker Krebs. “Identification of hybrid systems using a priori knowledge”.
609 In: *IFAC Proceedings Volumes* 35.1 (2002), pp. 451–456. DOI: 10.3182/20020721-6-ES-
610 1901.00563.
- 611 [23] Hayato Nakada, Kiyotsugu Takaba, and Tohru Katayama. “Identification of piecewise affine
612 systems based on statistical clustering technique”. In: *Automatica* 41.5 (2005), pp. 905–913.
613 DOI: 10.1016/j.automatica.2004.12.005.
- 614 [24] Necmiye Ozay. “An exact and efficient algorithm for segmentation of ARX models”. In:
615 *2016 American Control Conference (ACC)*. IEEE. 2016, pp. 38–41. DOI: 10.1109/ACC.2016.
616 7524888.
- 617 [25] Necmiye Ozay, Constantino Lagoa, and Mario Sznaier. “Robust identification of switched
618 affine systems via moments-based convex optimization”. In: *Proceedings of the 48th IEEE*
619 *Conference on Decision and Control (CDC) held jointly with 2009 28th Chinese Control*
620 *Conference*. IEEE, 2009, pp. 4686–4691. DOI: 10.1109/CDC.2009.5399962.
- 621 [26] Necmiye Ozay et al. “A sparsification approach to set membership identification of switched
622 affine systems”. In: *IEEE Transactions on Automatic Control* 57.3 (2012), pp. 634–648. DOI:
623 10.1109/TAC.2011.2166295.