# Finite sample learning of moving targets

Nikolaus Vertovec [a], Kostas Margellos [b], Maria Prandini [c]

[a]*Department of Computer Science, University of Oxford, OX1 3PJ, UK*

[b]*Department of Engineering Science, University of Oxford, OX1 3PJ, UK*

[c]*Department of Electronics, Information and Bioengineering, Politecnico di Milano, Milano 20133, Italy*

**Abstract**

We consider a moving target that we seek to learn from samples. Our results extend randomized techniques developed in control and optimization for a constant target to the case where the target is changing. We derive a novel bound on the number of samples that are required to construct a probably approximately correct (PAC) estimate of the target. Furthermore, when the moving target is a convex polytope, we provide a constructive method of generating the PAC estimate using a mixed integer linear program (MILP). The proposed method is demonstrated on an application to autonomous emergency braking.

*Key words:* Statistical learning theory; Randomized methods; Probably approximately correct learning; Data-driven algorithms; Drifting target concept.

## 1 Introduction

The use of probabilistic and randomized methods to analyze and design systems affected by uncertainty has long been a key research area within the control community. Early attempts at dealing with uncertainty were focused on stochastic approaches with later research focusing on *worst-case* settings. Probabilistic approaches to robustness emerged to alleviate conservatism of worst-case considerations by resorting to probabilistic information. This rapprochement between more traditional stochastic and robust paradigms facilitates uncertainty quantification based on data. To this end, we consider algorithms based on uncertainty randomization known as *randomized algorithms* [27], which allow us to apply tools from statistical learning theory based on Vapnik-Chervonenkis (VC) theory to control [1, 27, 28]. In general, these developments can be cast as binary classification problems with the main focus being the provision of finite-sample complexity bounds. VC theoretic techniques require the so called VC dimension to be finite. The computation of the VC dimension is in general a difficult task for generic optimization problems. Under a convexity assumption, the so-called scenario approach

has offered a theoretically sound and efficient methodology to provide *a-priori* probabilistic feasibility guarantees for uncertain optimization programs, with uncertainty represented by means of scenarios and without resorting to VC theory [5–7, 10, 12]. These developments have been recently extended to the non-convex case, however, they typically involve *a posteriori* guarantees [9, 17]. Applications and sample complexity bounds of the aforementioned methodologies to control synthesis problems have been demonstrated in [13, 15, 27], while notable extensions involve trading feasibility to performance [8, 25, 26], applications in game theory [16], and sequential methods [27]. Connections between the scenario approach and statistical learning theory based on the notion of compression have been provided in [11, 23].

The aforementioned approaches can be considered in the context of learning an unknown labeling mechanism, whereby we independently draw $m$ samples from a domain $X \subseteq \mathbb{R}^n$, according to some possibly unknown probability distribution $\mathbb{P}$. Each sample is assigned a $\{0, 1\}$-valued label according to an unknown *target* labeling function, $f$. The learning problem involves characterizing sample complexity bounds for $m$, such that we can generate a hypothesis $h$ based on the labeled $m$-multisample that, with a prescribed confidence $1 - \delta$, provides the same labeling with the target function when it comes to a new sample $x$ up to a predefined accuracy

_____

*Email addresses:* `nikolaus.vertovec@cs.ox.ac.uk` (Nikolaus Vertovec), `kostas.margellos@eng.ox.ac.uk` (Kostas Margellos), `maria.prandini@polimi.it` (Maria Prandini).

level $\epsilon$, i.e.,

$$\mathbb{P}^m\big\{(x_1,\ldots,x_m) \in X^m :$$
$$\mathbb{P}\{x \in X : \ h(x) \neq f(x)\} \leq \epsilon\big\} \geq 1 - \delta, \qquad (1)$$

where $\mathbb{P}^m$ is the product probability measure. An algorithm that generates a hypothesis satisfying the above statement is said to be probably approximately correct (PAC) to accuracy $\epsilon$ if the left side of (1) approaches 1 as $m \to \infty$ [28, pg. 56]. We will refer to the labeling mechanism as being PAC learnable to accuracy $\epsilon$ if there exists an algorithm that is PAC to accuracy $\epsilon$.

In this paper, we will study a similar problem of finding a hypothesis satisfying (1), however, with the notable difference that we consider a *tracking problem* where the unknown labeling function is changing after each drawn sample. In light of this labeling mechanism changing in a structured manner as specified in the sequel, we will consider both the construction of the hypothesis as well as the minimum number of samples that are necessary, so as to, with a certain confidence, provide probabilistic bounds on the event of the hypothesis disagreeing with the subsequently received label. A similar tracking problem with an alternative structure of change imposed on the target is considered in [2, 3, 14, 21, 22]. Similar to the structure considered in this paper, [19] considers a setting that allows for variations in the change between samples. In [22] the distribution according to which samples are drawn is also considered to be changing, while recent work, such as [18], has considered adapting to a variable rate of change of the target concept.

We first provide a formal mathematical formulation of the tracking problem considered in this paper in Section 2. Our main contributions can be summarized as follows:

(1) In Section 3 we provide *a-priori* bounds on the minimum number of samples needed to generate a PAC to accuracy $\epsilon$ hypothesis. This analysis capitalizes on the aforementioned references, and in particular the work of [19]. However, we re-approach this formulation providing a PAC-type of result (that involves two layers of probability) rather than an expected value assessment. We also provide a remedy for a mathematical omission in the analysis of [19].

(2) In Section 4 we provide a constructive method of generating a hypothesis from a finite set of samples using a Mixed Integer Linear Program (MILP) when the class of targets is convex polytopes. Note that the analysis in all aforementioned references is of existential nature, and this constitutes the first constructive approach for a hypothesis that enjoys such tracking properties.

We demonstrate numerically our theoretical results in Section 5 on a case study that involves autonomous emergency braking and discuss practical improvements to excluding samples from consideration in the MILP. Finally, Section 6 provides some concluding remarks.

## 2 Learning moving targets

### 2.1 Problem statement

We consider the problem of learning a labeling mechanism that is changing in a structured manner (this structure will be specified in the sequel). To this end, we follow a sample-based approach, where each sample $x$ is generated independently from a domain $X \subseteq \mathbb{R}^n$, endowed with a $\sigma$-algebra $\mathcal{X}$. Let $\mathbb{P}$ denote the fixed (potentially unknown) probability measure over $\mathcal{X}$. We refer to $(x_1,\ldots,x_m) \in X^m$ as an $m$-multisample, where its elements $x_i \in X$ are independently and identically distributed (i.i.d.) according to $\mathbb{P}$. For each $i = 1,\ldots,m$, let $f_i(\cdot) : \ X \to \{0,1\}$ be a $\{0,1\}$-valued labeling function, referred to as a *target function*.

In our setting, each sample $x_i$ is labeled according to target function $f_i$, $i = 1,\ldots,m$, giving rise to the *labeled $m$-multisample* $\{(x_1, f_1(x_1)),\ldots,(x_m, f_m(x_m))\}$. Notice that each sample is labeled by means of a different target function. As we consider the target functions to be unknown, we only have access to the labels of specific samples, namely $\{f_i(x_i)\}_{i=1}^m$. A natural question that we seek to answer is whether we can construct a labeling mechanism $h_m \in \mathcal{H}$ that correctly (with a certain probability) predicts the label that would be assigned to the next sample $x$ by the unknown target function $f_{m+1}$. In other words, we seek to provide probabilistic guarantees that $h_m(x) = f_{m+1}(x)$, where $h_m$ is referred to as a *hypothesis* and constitutes an approximation/prediction of $f_{m+1}$. Notice that we introduce the subscript $m$ to our hypothesis to highlight that this is constructed on the basis of the labeled $m$-multisample. We refer to this problem, pictorially illustrated in Figure 1, as a tracking problem, as we seek to track a moving labeling mechanism.

While the target functions are considered to be unknown we will make the following assumption on the target and hypotheses function class.

**Assumption 1.** All target and hypotheses functions belong to the same class $\mathcal{H}$, i.e., $f_1,\ldots,f_m,f_{m+1},h_m \in \mathcal{H}$, and $\mathcal{H}$ is assumed to be known. We further assume that $\mathcal{H}$ has a finite $VC$ dimension.

**Remark 1.** In Section 4 we will consider $\mathcal{H}$ to be the class of non-empty convex polytopes with a certain maximum number of facets, but make no such restriction for the main results of Section 3.
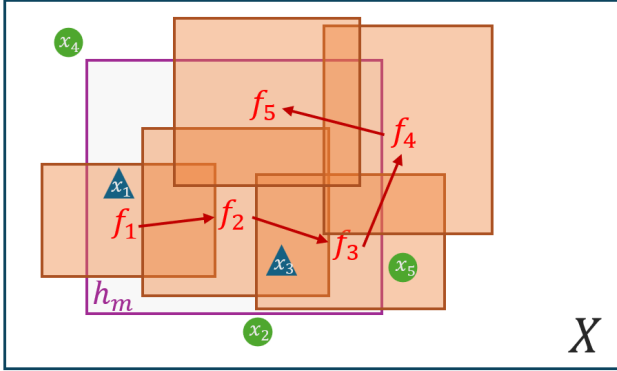
We formalize the tracking problem below.

Fig. 1. At each iteration, we receive a single sample along with a $\{0, 1\}$-valued label. To illustrate this, consider the labeling mechanism as an indicator function over the orange set. The orange set will change between each drawn sample (we illustrate this by depicting the orange sets across multiple iterations). The green circles indicate a 0-label, while the blue triangles represent a 1-label. We seek to find a hypothesis on the basis of the labeling $\{(x_1, f_1(x_1)), \ldots, (x_m, f_m(x_m))\}$ that, with certain confidence, will agree with the subsequent (unknown) target function $f_{m+1}$ on a new sample. We depict an example of such a hypothesis with the purple rectangle.

**Problem 1** (Tracking Problem)**.** Let $\epsilon, \delta \in (0, 1)$ be any fixed accuracy and confidence level, respectively. Determine $m_0(\epsilon, \delta)$ such that for any number of labeled samples $m \geq m_0(\epsilon, \delta)$, namely, $\{(x_1, f_1(x_1)), \ldots, (x_m, f_m(x_m))\}$, we can construct a hypothesis $h_m \in \mathcal{H}$ such that

$$\mathbb{P}^m\big\{(x_1, \ldots, x_m) \in \mathrm{X}^m :$$
$$\mathbb{P}\{x \in \mathrm{X} : h_m(x) \neq f_{m+1}(x)\} \leq \epsilon_0 + \epsilon\big\} \geq 1 - \delta, \quad (2)$$

where $\epsilon_0 \in (0, 1)$.

In words, with confidence at least $1 - \delta$, the probability that the constructed hypothesis $h_m$ produces a label for a new sample $x$ that does not agree with the target function $f_{m+1}$ is at most $\epsilon_0 + \epsilon$. Notice that the statement we seek to provide is within the realm of PAC learning. Yet unlike more standard PAC statements, the accuracy is deteriorated by $\epsilon_0$; this is not user-chosen but rather depends on how the target function is moving. We specify this in the next section and show that its presence is the price to pay for providing such statements for moving targets, while $\epsilon_0 = 0$ for the specific case of a constant target.

### 2.2 Mathematical preliminaries and assumptions

To simplify notation, for any labeling functions $f, h$ we define their probabilistic and empirical disagreement, re-

spectively, as

$$\mathrm{er}(f, h) \coloneqq \mathbb{P}\{x \in \mathrm{X} : h(x) \neq f(x)\}, \quad (3)$$

$$\widehat{\mathrm{er}}_m(f, h) \coloneqq \frac{1}{m} \sum_{i=1}^{m} |f(x_i) - h(x_i)|, \quad (4)$$

where the empirical disagreement is computed on an $m$-multisample $\{(x_1, f_1(x_1)), \ldots, (x_m, f_m(x_m))\}$, hence we introduce the subscript $m$ in the definition of $\widehat{\mathrm{er}}_m(\cdot, \cdot)$ to emphasize this dependence. Notice that in (4), $|f(x_i) - h(x_i)| = 1$ if $f, h$ disagree on $x_i$, and zero otherwise. Under these definitions, for $\epsilon, \delta \in (0, 1)$, the statement of (2) can be equivalently written as $\mathbb{P}^m\{(x_1, \ldots, x_m) \in \mathrm{X}^m : \mathrm{er}(f_{m+1}, h_m) \leq \epsilon_0 + \epsilon\} \geq 1 - \delta$.

We first provide some preliminary results that will be invoked in the subsequent developments. Proposition 1 below is a direct consequence of Hoeffding's inequality (see e.g., [20], [27]).

**Proposition 1.** Let $p_1, \ldots, p_m \in [0, 1]$, and consider independent Bernoulli random variables $Y_1, \ldots, Y_m$ such that $\mathbb{P}\{Y_i = 1\} = p_i$ and $\mathbb{P}\{Y_i = 0\} = 1 - p_i$, for all $i = 1, \ldots, m$. For any $\tau > 0$ we then have that

$$\mathbb{P}^m\Big\{\sum_{i=1}^{m} Y_i - \sum_{i=1}^{m} p_i > \tau\Big\} \leq e^{-\frac{2\tau^2}{m}}. \quad (5)$$

The following result is a PAC-type bound that holds for any target function $f \in \mathcal{H}$. This is [1, Theorem 7] adapted to our notation.

**Theorem 1.** Fix $\epsilon, \delta \in (0, 1)$ and $\rho \in [0, 1)$. Fix any $f \in \mathcal{H}$, and denote by $d$ the VC dimension of $\mathcal{H}$. For any

$$m \geq \frac{5(\rho + \epsilon)}{\epsilon^2}\Big(\ln\frac{4}{\delta} + d\ln\frac{40(\rho + \epsilon)}{\epsilon^2}\Big) \quad (6)$$

we have that

$$\mathbb{P}^m\{(x_1, \ldots, x_m) \in X^m : \exists h \in \mathcal{H} \text{ such that}$$
$$\widehat{\mathrm{er}}_m(f, h) \leq \rho \text{ and } \mathrm{er}(f, h) > \rho + \epsilon\} \leq \delta. \quad (7)$$

In words, Theorem 1 states that the probability that there exists a hypothesis such that its empirical error $\widehat{\mathrm{er}}_m(f, h)$ with the target function is at most $\rho$ but the actual error $\mathrm{er}(f, h)$ is higher than $\rho + \epsilon$, is at most equal to $\delta$ (which is typically selected to be small). Note that unlike the tracking problems presented in [19], we consider two levels of probability rather than an expected value assessment.

For the subsequent developments we consider target functions that exhibit the following structure on the way the labeling is changing, i.e., the target is moving.

**Assumption 2.** Let $f_1, \ldots, f_m, f_{m+1} \in \mathcal{H}$, and consider $\underline{\mu}, \overline{\mu} \in (0, 1)$ with $\underline{\mu} \leq \overline{\mu}$. We assume that the average probability of disagreement of the previous labels with the label $f_{m+1}$, denoted by

$$\mu = \frac{1}{m} \sum_{i=1}^{m} \mathrm{er}(f_i, f_{m+1}), \tag{8}$$

is bounded such that $\underline{\mu} \leq \mu \leq \overline{\mu}$.

Assumption 2 implies that the target sets are changing but we impose a restriction (both upper and lower limits) on the probability that the labeling they produce changes. We refer to $\underline{\mu}, \overline{\mu}$ as the minimum and maximum, respectively, target change.

## 3 Finite sample probabilistic certificates

### 3.1 Main result

Problem 1 requires obtaining finite sample complexity bounds such that a hypothesis $h_m$ constructed on the basis of a labeled $m$-multisample tracks (probabilistically) the moving target function. In this section, we show that this is the case for hypotheses in minimal empirical disagreement on the $m$-multisample. We formalize the set of such hypotheses in the definition below.

**Definition 1** (minimal disagreement)**.** Consider a labeled $m$-multisample $\{(x_1, f_1(x_1)), \ldots, (x_m, f_m(x_m))\}$. We refer to the set

$$M_m := \operatorname*{arg\,min}_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^{m} |f_i(x_i) - h(x_i)| \tag{9}$$

as the set of hypotheses in $\mathcal{H}$ that minimize the empirical error with the labeled $m$-multisample. We then say that any $h \in M_m$ is in minimal disagreement with $f_1, \ldots, f_m$.

Let $h_m$ be a hypothesis in minimal disagreement with $f_1, \ldots, f_m$. We show that for this particular hypothesis choice, we can provide an answer to Problem 1, with $\epsilon_0 = 4\overline{\mu}$. We formalize this in the next theorem, which is the main result of this section.

**Theorem 2.** *Fix $\epsilon, \delta \in (0, 1)$. Denote by $d$ the VC dimension of $\mathcal{H}$, and consider Assumption 2 with $\overline{\mu} < \frac{1}{4}$. If we choose $m \geq m_0(\epsilon, \delta)$, where*

$$m_0(\epsilon, \delta) = \max \left\{ \frac{1}{2\underline{\mu}^2} \ln \frac{2}{\delta}, \right.$$
$$\left. \frac{5(4\overline{\mu} + \epsilon)}{\epsilon^2} \left( \ln \frac{8}{\delta} + d \ln \frac{40(4\overline{\mu} + \epsilon)}{\epsilon^2} \right) \right\}, \tag{10}$$

*we then have that for any $h_m \in M_m$,*

$$\mathbb{P}^m \{ (x_1, \ldots, x_m) \in \mathrm{X}^m :$$
$$\mathrm{er}(f_{m+1}, h_m) \leq 4\overline{\mu} + \epsilon \} \geq 1 - \delta. \tag{11}$$

*Proof.* Fix any $\epsilon, \delta \in (0, 1)$. We define the following events:

$$E = \{ (x_1, \ldots, x_m) \in \mathrm{X}^m : \mathrm{er}(f_{m+1}, h_m) > 4\overline{\mu} + \epsilon \},$$
$$A = \{ (x_1, \ldots, x_m) \in \mathrm{X}^m :$$
$$\frac{1}{m} \sum_{i=1}^{m} |f_i(x_i) - f_{m+1}(x_i)| > 2\mu \}. \tag{12}$$

$A$ is an approximation set as it includes the $m$-multisamples for which the empirical average disagreement $\frac{1}{m} \sum_{i=1}^{m} |f_i(x_i) - f_{m+1}(x_i)|$ is at least twice as big as the actual average disagreement $\mu$. $E$ plays the role of the error set, as by its definition, $\mathbb{P}^m\{E\} \leq \delta$ is the complementary statement to that of (11).

We can bound $\mathbb{P}^m\{E\}$ as

$$\mathbb{P}^m\{E\} = \mathbb{P}^m\{E \cap A\} + \mathbb{P}^m\{E \cap \overline{A}\}$$
$$\leq \mathbb{P}^m\{A\} + \mathbb{P}^m\{E \cap \overline{A}\}, \tag{13}$$

where $\overline{A}$ denotes the complement of $A$. The inequality is since $\mathbb{P}^m\{E \cap A\} \leq \mathbb{P}^m\{A\}$. To show (11), we can equivalently establish that $\mathbb{P}^m\{E\} \leq \delta$. To achieve this, it suffices to show that $\mathbb{P}^m\{A\} \leq \delta/2$ and $\mathbb{P}^m\{E \cap \overline{A}\} \leq \delta/2$ [1].

*Case $\mathbb{P}^m\{A\} \leq \delta/2$:* For each $i = 1, \ldots, m$, set $Y_i = |f_i(x_i) - f_{m+1}(x_i)|$ and $p_i = \mathrm{er}(f_i, f_{m+1})$ so that $\mathbb{P}\{Y_i = 1\} = p_i$ and $\mathbb{P}\{Y_i = 0\} = 1 - p_i$. Notice that $Y_1, \ldots, Y_m$ are independent Bernoulli random variables, and by (8), $\sum_{i=1}^{m} p_i = m\mu$. Under this variables assignment, and selecting $\tau = m\mu$, $\mathbb{P}^m\{A\}$ coincides with the left-hand side of (5). We then have that

$$\mathbb{P}^m\{A\} \leq e^{-2m\mu^2} \leq e^{-2m\underline{\mu}^2}, \tag{14}$$

where the first inequality is due to Proposition 1, and the second one is since $\mu \geq \underline{\mu}$ by Assumption 2.

By inspection of (14), to ensure that $\mathbb{P}^m\{A\} \leq \delta/2$, it suffices to show that $e^{-2m\underline{\mu}^2} \leq \delta/2$. By taking the

---

[1] Splitting the confidence equally between these two terms is not necessary; further optimizing the split would have minor effect on the final sample size bound as the confidence appears inside the logarithm. As such we do not pursue this here to simplify the analysis.

logarithm making $m$ the argument, we conclude that if

$$m \geq \frac{1}{2\underline{\mu}^2} \ln \frac{2}{\delta} \implies \mathbb{P}^m\{A\} \leq \frac{\delta}{2}. \qquad (15)$$

*Case* $\mathbb{P}^m\{E \cap \overline{A}\} \leq \delta/2$: We have that

$$\mathbb{P}^m\{E \cap \overline{A}\}$$
$$\leq \mathbb{P}^m\{(x_1, \ldots, x_m) \in \mathrm{X}^m : \mathrm{er}(f_{m+1}, h_m) > 4\overline{\mu} + \epsilon$$
$$\text{and } \frac{1}{m} \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)| \leq 2\overline{\mu}\}, \quad (16)$$

where the first statement in the right-hand side of (16) is the event $E$, and the second one encompasses $\overline{A}$. To see the latter, notice that $\overline{A}$ requires $\frac{1}{m} \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)| \leq 2\mu$, and $\mu \leq \overline{\mu}$ due to Assumption 2.

We have assumed that for any $m$-multisample, $h_m$ is chosen from $M_m$. By (9), since $h_m \in M_m$, $\sum_{i=1}^m |f_i(x_i) - h_m(x_i)| \leq \sum_{i=1}^m |f_i(x_i) - h(x_i)|$ for any $h \in \mathcal{H}$. However, since we also have that $f_{m+1} \in \mathcal{H}$, we have that for any $m$-multisample,

$$\sum_{i=1}^m |f_i(x_i) - h_m(x_i)| \leq \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)|. \quad (17)$$

Moreover, we have that

$$\widehat{\mathrm{er}}_m(f_{m+1}, h_m) = \frac{1}{m} \sum_{i=1}^m |f_{m+1}(x_i) - h_m(x_i)|$$
$$\leq \frac{1}{m} \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)| + \frac{1}{m} \sum_{i=1}^m |f_i(x_i) - h_m(x_i)|$$
$$\leq \frac{2}{m} \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)|, \qquad (18)$$

where the equality is due to (4), and the first inequality is by adding and subtracting $f_i(x_i)$ in each term in the summation and applying the triangle inequality. The last inequality is due to (17).

Since (18) holds for any $m$-multisample, we have that

$$\{(x_1, \ldots, x_m) \in \mathrm{X}^m : \frac{1}{m} \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)| \leq 2\overline{\mu}\}$$

$$\subseteq \{(x_1, \ldots, x_m) \in \mathrm{X}^m : \widehat{\mathrm{er}}_m(f_{m+1}, h_m) \leq 4\overline{\mu}\}. \quad (19)$$

As a result, by (16) and (19) we obtain

$$\mathbb{P}^m\{E \cap \overline{A}\}$$
$$\leq \mathbb{P}^m\{(x_1, \ldots, x_m) \in \mathrm{X}^m : \mathrm{er}(f_{m+1}, h_m) > 4\overline{\mu} + \epsilon$$
$$\text{and } \widehat{\mathrm{er}}_m(f_{m+1}, h_m) \leq 4\overline{\mu}\}. \qquad (20)$$

Notice that (20) takes the form of (7), with $f_{m+1}$, $h_m$ and $4\overline{\mu}$ in place of $f$, $h$ and $\rho$, respectively. Theorem 1 with $\delta/2$ in place of $\delta$ implies that

$$m \geq \frac{5(4\overline{\mu} + \epsilon)}{\epsilon^2} \left( \ln \frac{8}{\delta} + d \ln \frac{40(4\overline{\mu} + \epsilon)}{\epsilon^2} \right)$$
$$\implies \mathbb{P}^m\{E \cap \overline{A}\} \leq \frac{\delta}{2}. \qquad (21)$$

By (15) and (21), we obtain that if $m \geq m_0(\epsilon, \delta)$, where $m_0(\epsilon, \delta)$ is as in (10), we have that $\mathbb{P}^m\{E\} \leq \delta$, thus concluding the proof. □

The proof of Theorem 2 is inspired by [19, Theorem 1]. However, the result therein does not involve two layers of probability and effectively provides a bound on the expectation of the probability of incorrectly tracking the target. Moreover, only an upper bound on the target change is considered in [19]. This is due to the fact that a term similar to $e^{-2m\mu^2}$ was bounded by $e^{-2m\overline{\mu}^2}$, which is, however, not valid as $\mu \leq \overline{\mu}$. Here we correct this issue by introducing a lower bound on the target change, resulting in equation (14).

The sample size bound in (10) depends polynomially on $1/\epsilon$ and logarithmically on $\delta$. This implies that we could make the confidence $1 - \delta$ high without an unaffordable increase on the number of samples required. Figure 2 illustrates the number of samples as a function of $\epsilon$. The color code corresponds to different values of $\underline{\mu}, \overline{\mu}$. It should be noted that the overall accuracy level for the prediction properties of our hypothesis is $4\overline{\mu} + \epsilon$. Even though $\epsilon$ is user-chose, $\overline{\mu}$ is a property of the target, and as such the labeling mechanism is considered to be PAC learnable to accuracy $4\overline{\mu}$. Therefore, insightful accuracy levels can be achieved if $\overline{\mu}$ is relatively low, i.e., for moderately changing target functions.

**Remark 2** (Effect of $\underline{\mu}, \overline{\mu}$). As evident from (10), the minimum number of samples that need to be generated is the maximum of two terms: the first one depends only on the minimum target change $\underline{\mu}$, while the second one depends on the maximum target change $\overline{\mu}$. These two sample size bounds that comprise $m_0(\epsilon, \delta)$ emanate from bounding

(1) the event $A$, that the empirical average disagreement $\frac{1}{m} \sum_{i=1}^m |f_i(x_i) - f_{m+1}(x_i)|$ is at least twice as big as the actual average disagreement $\mu$. Bounding this term is responsible for the first sample size bound in (10).

(2) the event $E \cap \overline{A}$, that the empirical average disagreement is less than twice as the actual average disagreement $\mu$, yet that the true probability of disagreement between the hypothesis, $h_m$, and the

5

subsequent label, $f_{m+1}$, is more than $4\overline{\mu}+\epsilon$. Bounding this term is responsible for the second sample size bound in (10).

With reference to Figure 2, for high values of $\epsilon$ the first sample size bound in (10) dominates, which is independent of $\epsilon$, hence that part of each curve is constant. On the contrary, for lower values of $\epsilon$ the second sample size bound in (10) becomes the dominant one.

If $\underline{\mu}$ is sufficiently low (the target could move slowly), then the first sample size bound in (10) dominates. Intuitively, this implies that if the target could move slowly, then learning the actual probability of change from the empirical one (this is encoded in the definition of the event $A$) requires more samples, as with few samples we might get misleading results due to observing a faster target change than the true average change in the target. With reference to Figure 2, the minimum number of samples required increases as $\underline{\mu}$ decreases (compare the constant part of the curves).

If we now allow for a large change of the target, encoded by a large $\overline{\mu}$, then the second sample size bound in (10) dominates. This implies that we need a sufficiently high number of samples to, with high confidence, bound the event that the true change with respect to the subsequent label, $f_{m+1}$, is not considerably lower than the observed, empirical change (encoded by event $E \cap \overline{A}$). Intuitively, if the target is moving fast, then incorrectly predicting the label of a new sample if the empirical error is low (event $E \cap \overline{A}$) requires more samples. This is since with fewer samples we may get into a situation with a low empirical error, however, due to the target changing fast the error when it comes into predicting the label of yet another sample may be significantly higher. With reference to Figure 2, for any fixed $\epsilon$, the minimum number of samples required increases as $\overline{\mu}$ increases (compare the non-constant part of the curves).

To account for both cases and make sure that the probability of both events $A$ and $E \cap \overline{A}$ is sufficiently low, we take the maximum of the associated sample size bounds.

**Remark 3** (Constant target). The case of a constant target can be obtained as a direct byproduct of the proof of Theorem 2. To see this, notice that a constant target implies that $\underline{\mu} = \overline{\mu} = 0$, i.e., if all target functions are the same, their mutual error is zero. As such, $f_i(x_i)$ and $f_{m+1}(x_i)$ will always be in agreement. We present the proof of Theorem 2 under a constant target assumption in the Appendix. As a result, the sample size bound is identical to that of Theorem 1 with $\rho = 0$. This implies, that when it comes to providing guarantees for the minimal disagreement hypothesis and for the case where the target is constant, Theorem 2 specializes to the result of Theorem 1 with $\rho = 0$. With reference to Problem 1, notice also that in this case, $\epsilon_0 = 4\overline{\mu} = 0$.
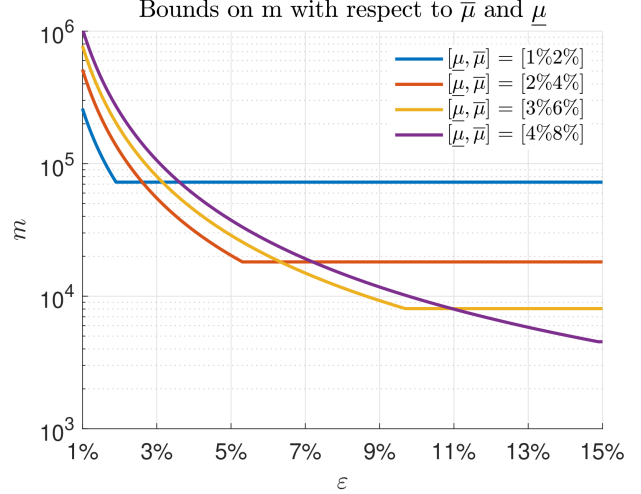


Bounds on m with respect to $\overline{\mu}$ and $\underline{\mu}$

Fig. 2. Number of samples required according to (10) for different accuracy levels $\epsilon$ and $\delta = 10^{-6}$ with VC dimension 4. The color code corresponds to different values of $\underline{\mu}, \overline{\mu}$. Notice that the term dependent on $\underline{\mu}$ in (10) does not depend on $\epsilon$ and thus constitutes the constant dominant at higher levels of $\epsilon$.

## 4 Hypothesis computation

We now consider the construction of the hypothesis $h_m$ that minimizes the empirical error with respect to the labeled $m$-multisample, i.e., $h_m \in M_m$. For the remainder of the paper, we will assume that the domain X is compact. Furthermore, we consider the labelling functions $f_i$, $i = 1, \ldots, m$, to be defined as

$$f_i(x) = \mathbb{1}_{B_i}(x) = \begin{cases} 1 & \text{if } x \in B_i \\ 0 & \text{otherwise} \end{cases}, \qquad (22)$$

with the sets $B_i$, $i = 1, \ldots, m$, being non-empty convex polytopes in $\mathbb{R}^n$, each of them having at most $n_f$ facets. As the hypothesis belongs to the same class with the target functions, we seek to find a convex polytope, denoted by $B_{h_m}$, such that the hypothesis $h_m$ defined as

$$h_m(x) = \mathbb{1}_{B_{h_m}}(x) = \begin{cases} 1 & \text{if } x \in B_{h_m} \\ 0 & \text{otherwise} \end{cases}, \qquad (23)$$

is in minimal disagreement with the observed labels. Since $B_{h_m}$ is a convex polytope with at most $n_f$ facets we represent it by means of $n_f$ linear inequality constraints as $Ax + b \leq 0$, where $A \in \mathbb{R}^{n_f \times n}$ and $b \in \mathbb{R}^{n_f}$. Denote each row-vector of $A$ (respectively, $b$) by $a_j$ ($b_j$), $j = 1, \ldots, n_f$. For each $j = 1, \ldots, n_f$, $a_j x + b_j$ denotes then a facet of $B_{h_m}$. We make this parameterization explicit by denoting the convex polytope as $B_{h_m}(A, b)$. Moreover, we assume that for each $j = 1, \ldots, n_f$, $(a_j^\top, b_j) \in C_j \subset \mathbb{R}^{n_f+1}$, where $C_j$ is some arbitrarily large compact set that contains the origin in its interior. The purpose of the set will become clear in the sequel.

We show how to construct a Mixed Integer Linear Program (MILP) that will in turn return the parameterization of $B_{h_m}$, namely $A$ and $b$, that results in a hypothesis $h_m \in M_m$. To this end, let $I_1$ and $I_0$ be the set of sample indices for which the label is 1 and 0, respectively, i.e.,

$$I_1 = \{i \in \{1, \ldots, m\} \text{ such that } f_i(x_i) = 1\}, \quad (24)$$
$$I_0 = \{i \in \{1, \ldots, m\} \text{ such that } f_i(x_i) = 0\}. \quad (25)$$

We instantiate the MILP that returns the minimal disagreement hypothesis in the following main steps:

*1. Disagreement with the sample indices in $I_1$.* Fix any $i \in I_1$, and let $x_i$ be the associated sample. Fix also a parameterization $A, b$ of $B_{h_m}$. If $h_m(x_i) = f_i(x_i) = 1$, i.e., the label that a hypothesis, constructed on the basis of $B_{h_m}(A, b)$, provides on $x_i$ agrees with that of $f_i$, then $x_i \in B_{h_m}(A, b)$ since $i \in I_1$. We thus have that

$$\begin{aligned} x_i \in B_{h_m}&(A, b) \\ &\iff a_j x_i + b_j \leq 0, \ \forall j = 1, \ldots, n_f. \end{aligned} \quad (26)$$

However, we are seeking a hypothesis that is in minimal disagreement with the samples, rather than in zero disagreement. As such, we want to allow for a certain number of incorrect labels, or equivalently, we want to allow violating the right-hand side of (26). Therefore, we introduce the slack variables $s_{ij} \geq 0$, $j = 1, \ldots, n_f$, $i \in I_1$. As such, for each $i \in I_1$, we consider the relaxed constraints

$$a_j x_i + b_j \leq s_{ij}, \ \forall j = 1, \ldots, n_f. \quad (27)$$

By means of (26) and the definition of $h_m$, enforcing (27), implies that

$$\begin{cases} h_m(x_i) \neq f_i(x_i) & \text{if } \sum_{j=1}^{n_f} s_{ij} > 0, \\ h_m(x_i) = f_i(x_i) & \text{otherwise.} \end{cases} \quad (28)$$

In words, if $\sum_{j=1}^{n_f} s_{ij} > 0$ (which is satisfied if at least one $s_{ij}$, $j = 1, \ldots, n_f$, is positive as the slack variables are non-negative) implies that the hypothesis $h_m$ disagrees with the target function $f_i$ on the sample $x_i$. If all slack variables are zero, then $h_m$ agrees with $f_i$ on $x_i$, $i \in I_1$.

*2. Disagreement with the sample indices in $I_0$.* Fix any $i \in I_0$, and let $x_i$ be the associated sample. Fix also a parameterization $A, b$ of $B_{h_m}$. If $h_m(x_i) = f_i(x_i) = 0$, i.e., the hypothesis and the target function $f_i$ agree on $x_i$, then $x_i \notin B_{h_m}(A, b)$. This exclusion can imply that the sample $x_i$ would violate the half-space constraint encoding the facets of $B_{h_m}(A, b)$ for at least one facet. This can be written as a logical constraint; employing the developments of [4, 24], we equivalently reformulate it to mixed-integer inequalities by introducing the binary variables $z_{ij} \in \{0, 1\}$, $j = 1, \ldots, n_f$, $i = 1, \ldots, m$. Let $M_j = \sup_{x \in X, (a_j^\top, b_j) \in C_j} a_j x + b_j$,

$m_j = \inf_{x \in X, (a_j, b_j) \in C_j} a_j^\top x + b_j$, $j = 1, \ldots, n_f$. Note that these exist and are finite, as X and $C_j$, $j = 1, \ldots, n_f$, are assumed to be compact. We then have that

$$\begin{aligned} & x_i \notin B_{h_m}(A, b) \\ & \iff \begin{cases} a_j x_i + b_j \leq M_j(1 - z_{ij}), \ \forall j = 1, \ldots, n_f, \\ a_j x_i + b_j > m_j z_{ij}, \ \forall j = 1, \ldots, n_f, \\ \sum_{j=1}^{n_f} z_{ij} \leq n_f - 1. \end{cases} \end{aligned}$$
$$(29)$$

Notice that if $z_{ij} = 0$, then the first inequality in (29) becomes $a_j x_i + b_j \leq M_j$ (trivially satisfied by the definition of $M_j$), while the second one reduces to $a_j x_i + b_j > 0$. The latter implies then that $x_i \notin B_{h_m}(A, b)$ as it violates the constraint of its $j$-th facet. On the contrary, if $z_{ij} = 1$, then the first inequality in (29) implies that $x_i$ is within the half-space defined by the $j$-th facet of $B_{h_m}(A, b)$ [2]. For $x_i$ to be inside $B_{h_m}(A, b)$, i.e., $x_i \in B_{h_m}(A, b)$, this has to be the case for all $j = 1, \ldots, n_f$, or equivalently $\sum_{j=1}^{n_f} z_{ij} = n_f$. This justifies the last constraint in (29).

Since we only seek a hypothesis in minimal (rather than in zero) disagreement with the target functions, we relax these constraints by introducing slack variables $s_{ij} \geq 0$, $j = 1, \ldots, n_f$, $i \in I_0$. As such, for each $i \in I_0$, the associated relaxed constraints are given by

$$\begin{cases} a_j x_i + b_j \leq M_j(1 - z_{ij}), \ \forall j = 1, \ldots, n_f, \\ a_j x_i + b_j > m_j z_{ij} - s_{ij}, \ \forall j = 1, \ldots, n_f, \\ \sum_{j=1}^{n_f} z_{ij} \leq n_f - 1. \end{cases} \quad (30)$$

Notice that we do not need to introduce a slack variable in the first inequality in (30), as this becomes non-redundant only if $z_{ij} = 1$. In this case, however, satisfying the resulting inequality would already mean disagreeing with the target, so we do not need to relax that condition. By means of (29) and the definition of $h_m$, enforcing (30) leads to the same disagreement implications as in (28).

*3. Minimizing disagreements.* In view of constructing the hypothesis that is in minimal disagreement with the target functions, we need to be able to count the number of disagreements. However, if $i \in I_1$ we have a disagreement if $x_i \notin B_{h_m}(A, b)$, while if $i \in I_0$ we have a disagreement if $x_i \in B_{h_m}(A, b)$. By (28) and the discussion below (30), disagreement happens if $\sum_{j=1}^{n_f} s_{ij} > 0$. If we

---

[2] Note that if $a_j x_i + b_j = m_j$, for $z_{ij} = 1$, the second inequality in (29) would not be satisfied. This limiting case where $a_j x_i + b_j$ admits its lowest value is not an issue in the numerical implementation (see Remark 4) as a tolerance parameter is introduced to "implement" strict inequalities numerically. Alternatively, we could choose any finite $m_j < \inf_{x \in X, (a_j, b_j) \in C_j} a_j^\top x + b_j$, $j = 1, \ldots, n_f$, rather than choosing $m_j$ exactly equal to its lowest admissible value.

introduce the binary variable $v_i \in \{0,1\}$, $i = 1, \ldots, m$, defined as

$$v_i = \begin{cases} 1 & \text{if } \sum_{j=1}^{n_f} s_{ij} > 0, \\ 0 & \text{otherwise}, \end{cases} \qquad (31)$$

then, the total number of disagreements that we seek to minimize is given by $\sum_{i=1}^{m} v_i$.

For each $j = 1, \ldots, n_f$, we have assumed that $(a_j^\top, b_j) \in C_j$, where $C_j$ is compact and contains the origin in its interior. As such, $M_j > 0$ and $m_j < 0$ for all $j = 1, \ldots, n_f$. Therefore, by (27) and the definition of $M_j$, $s_{ij} \leq M_j$, for all $i \in I_1$. Similarly, by (30) and the definition of $m_j$, $s_{ij} < -m_j$, for all $i \in I_0$. Notice that this follows from requiring the right-hand side in the second inequality of (30) to be greater than or equal to the worst-case lower bound of $a_j x_i + b_j$, namely, $m_j$, for the case where $z_{ij} = 0$ that this constraint becomes nontrivial. Summing the across $j = 1, \ldots, n_f$, we obtain

$$\begin{cases} \sum_{j=1}^{n_f} s_{ij} \in [0, \sum_{j=1}^{n_f} M_j], & \text{if } i \in I_1 \\ \sum_{j=1}^{n_f} s_{ij} \in [0, -\sum_{j=1}^{n_f} m_j), & \text{if } i \in I_0. \end{cases} \qquad (32)$$

The logical implication in (31) is then reformulated as

$$\begin{cases} \sum_{j=1}^{n_f} s_{ij} - v_i \sum_{j=1}^{n_f} M_j \leq 0, & \text{if } i \in I_1, \\ \sum_{j=1}^{n_f} s_{ij} + v_i \sum_{j=1}^{n_f} m_j < 0, & \text{if } i \in I_0. \end{cases} \qquad (33)$$

To see the equivalence between (33) and (31), consider the former inequality in (33). Notice that if $\sum_{j=1}^{n_f} s_{ij} > 0$ then this implies that we must have $v_i \sum_{j=1}^{n_f} M_j > 0$ which, since $M_j > 0$ for all $j = 1, \ldots, n_f$, implies that $v_i = 1$. On the other hand, if $\sum_{j=1}^{n_f} s_{ij} = 0$, then (33) implies that $v_i \sum_{j=1}^{n_f} M_j \geq 0$. However, since we are seeking the minimal disagreement hypothesis and hence we will be minimizing $\sum_{i=1}^{m} v_i$, the minimum value of $v_i$ for which the previous inequality is satisfied is $v_i = 0$. A similar reasoning applies also to the equivalence between the second inequality in (33) and (31).

*4. Minimal disagreement MILP.* The MILP that results in a hypothesis that is in minimal disagreement with respect to the target functions on the $m$-multisample,

i.e., $h_m \in M_m$, is given by:

$$\underset{A, b, \left\{ \{z_{ij}, s_{ij}\}_{j=1}^{n_f} \right\}_{i=1}^{m}, \{v_i\}_{i=1}^{m}}{\text{minimize}} \sum_{i=1}^{m} v_i \qquad (34)$$

subject to

$$\forall i \in I_1 : \begin{cases} a_j x_i + b_j \leq s_{ij}, \ \forall j = 1, \ldots, n_f, \\ \sum_{j=1}^{n_f} s_{ij} - v_i \sum_{j=1}^{n_f} M_j \leq 0, \end{cases} \qquad (35)$$

$$\forall i \in I_0 : \begin{cases} a_j x_i + b_j \leq M_j(1 - z_{ij}), \ \forall j = 1, \ldots, n_f, \\ a_j x_i + b_j > m_j z_{ij} - s_{ij}, \ \forall j = 1, \ldots, n_f, \\ \sum_{j=1}^{n_f} z_{ij} \leq n_f - 1, \\ \sum_{j=1}^{n_f} s_{ij} + v_i \sum_{j=1}^{n_f} m_j < 0. \end{cases} \qquad (36)$$

The constraints in (35) correspond to (27) and the first inequality in (33), encoding (relaxed) agreement on the sample with $i \in I_1$, and determining disagreements for this case, respectively. Similarly, the constraints in (36) correspond to (30) and the second inequality in (33), and admit a similar interpretation.

The objective function $\sum_{i=1}^{m} v_i$ involves minimizing the total number of disagreements. We use the volume of the convex polytope parameterized by $A, b$, namely, $\text{vol}(A, b)$, as a tie-break rule to single out a unique solution in case of multiple minimizers. Once the optimal $A, b$ is determined, we can construct $B_{h_m}(A, b)$, and hence $h_m$ by means of (23).

**Remark 4.** Note that for the samples indexed by $i \in I_0$, the second inequality in the disagreement constraints in (29) (and hence also (30)) are strict. From a numerical point of view, to implement these constraints we can turn them into non-strict inequalities, where following [4] we introduce a tolerance parameter $\varrho \geq 0$, fixed to the numerical solver precision.

We can then replace the second inequality in (29) by

$$a_j x_i + b_j \geq \varrho + (m_j - \varrho) z_{ij}, \ \forall j = 1, \ldots, n_f.$$

Similarly, the second inequality in (30) should be replaced by $a_j x_i + b_j \geq \varrho + (m_j - \varrho) z_{ij} - s_{ij}, \ \forall j = 1, \ldots, n_f$. As a result, the second and fourth inequalities in (36) should, respectively, become

$$a_j x_i + b_j \geq \varrho + (m_j - \varrho) z_{ij} - s_{ij}, \ \forall j = 1, \ldots, n_f$$
$$\sum_{j=1}^{n_f} s_{ij} + v_i \sum_{j=1}^{n_f} (m_j - \varrho) \leq 0.$$

Once such a $\rho$ parameter is introduced, for samples indexed by $i \in I_0$, the condition $\sum_{j=1}^{n_f} s_{ij} > 0$ is necessary but not sufficient for the hypothesis to disagree with the target. To see this, notice that if $z_{ij} = 0$ then the second

inequality in (30) would become non-redundant, and result in $a_j x_i + b_j \geq \varrho - s_{ij}$. Due to the presence of $\varrho > 0$, if $s_{ij} > 0$ but $\varrho - s_{ij} > 0$, then $x_i$ may still be outside of a facet of the convex polytope thus agreeing with the target (recall that label agreement here means being outside $B_{h_m}(A, b)$) despite the fact that the associated slack is non-zero. As a result, the MILP in (34)-(36) minimizes an upper bound on the total number of disagreements.

## 5 Numerical example

### 5.1 Problem set-up

We demonstrate numerically our theoretical developments on a case study that involves Autonomous Emergency Braking (AEB) systems. Furthermore, we consider the computational feasibility of the MILP and introduce an approach to discard redundant samples, thus reducing the constraints of the MILP.

Let us consider a car driving along a road while receiving measurements of the distance $l$ to any vehicle or obstacle ahead, as well as its velocity $v$. If the braking distance at the current velocity exceeds the available distance to the car or obstacle ahead, we want the AEB system to engage the brakes autonomously. The necessary braking distance in case of an emergency stop can be calculated by setting the braking force times the distance equal to the kinetic energy of the vehicle. Thus if

$$\frac{1}{2} v^2 \frac{m}{F} \leq l, \tag{37}$$

where $m$ is the vehicle mass and $F$ is the braking force, then there is a sufficient distance to the vehicle or obstacle ahead, hence the corresponding measurement is classified as safe. In view of (37) depending on $v^2$, hereafter we consider $x = (l, v^2)$ as the measurement vector.

The braking force will depend on the friction coefficient of the brakes and will deteriorate over time. Similarly, the vehicle mass will depend on the fuel, passengers, and cargo, which will also change over time. In line with our theoretical developments, we consider $x_i$, $i = 1, \ldots, m$, to be independent measurements, with the index $i$ acting as a time-stamp. Let $F_i$ denote the corresponding braking force, which depends on $i$ to reflect the change of the friction coefficient, and let $m_i$ denote the vehicle mass, which will also depend on $i$ to reflect changes to the vehicle mass. This dependence of $F_i$ and $m_i$ on $i$ induces a different labeling function $f_i$. In particular, we label a sample $x = (l, v^2)$ by means of

$$f_i(x) = \begin{cases} 1 & \text{if } \frac{1}{2} v^2 \frac{m_i}{F_i} \leq l \\ 0 & \text{otherwise.} \end{cases} \tag{38}$$

For the construction of the hypothesis, we collect a measurement $x_i$ after each engagement of the vehicle's

brakes. In addition to obtaining $x_i$, we assume to obtain a brake performance measurement $\frac{m_i}{F_i}$ from which to construct the label $f_i$. Furthermore, we assume that we have knowledge of the expected minimum and maximum degradation of the braking performance, allowing us to obtain values for $\underline{\mu}$ and $\overline{\mu}$, respectively.

Using numerical values for the braking parameters and vehicle mass as defined in the sequel, the evolution of the braking performance is shown in Figure 3. For visual clarity, we only depict a random subset of the samples.

For an effective AEB system, we want to classify a new sample $x$ as safe or unsafe without having to first engage the brakes to receive a measurement of the new braking performance, which depends on the unknown braking force $F_{m+1}$ and mass $m_{m+1}$. In light of this, we will utilize the results from Section 3 to construct a hypothesis on the basis of a labeled $m$-multisample, allowing us to a-priori classify a sample $x$ as safe or unsafe.

Aligned with the theoretical developments of Section 4 we consider our hypothesis to be a convex polytope in the 2D plane as illustrated in Figure 4. We assume that matrix $A$ parameterizing the convex polytope is fixed, and is given by

$$A = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \\ -\cos\theta & -\sin\theta \\ \sin\theta & -\cos\theta \end{bmatrix}, \tag{39}$$

where $\theta := \tan^{-1} \frac{m}{2F}$ denotes the rotation of the convex polytope. Since the safety label is defined by a single half-plane, only one of the facets of the convex polytope becomes relevant, namely $a_3 = [-\cos\theta \ -\sin\theta]$. This observation reduces the VC-dimension (employed in the sample complexity bounds) to $d = 1$. Considering the evolution of the braking performance as shown in Figure 3, the rotation of the convex polytope is minimal. Since the inclusion of the variable rotation introduces a nonlinearity into the MILP, for the sake of clarity, we will consider the angle $\theta$ to be fixed in the subsequent computation of the hypothesis (however leave the true safety label unchanged).

To reduce the size of the MILP (34 - 36), we will consider how samples can be discarded prior to computing the hypothesis. Recall that the MILP minimizes the total number of disagreements between the hypothesis and the sample labels. Thus, at best, we can obtain zero disagreement between the hypothesis and the samples in $I_1$. For the AEB example, this implies that the halfplane constructed by the hypothesis will need to lie to the left of all samples in $I_1$. This is illustrated in Figure 5 by the cyan dotted line. However, any halfplane that lies further
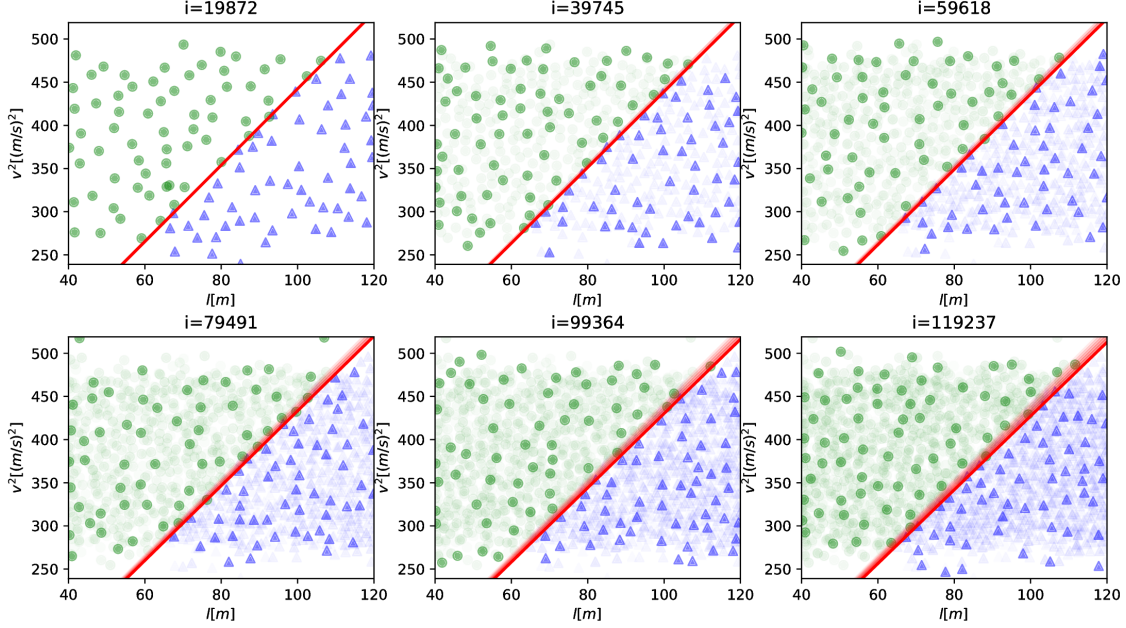
9

Fig. 3. The evolution of the braking performance over time. Green circles indicate samples with label 0, while blue triangles show samples with label 1. The bold red halfplane represents the true safety label at the given iteration, while the opaque halfplanes show the safety boundary at previous iterations.
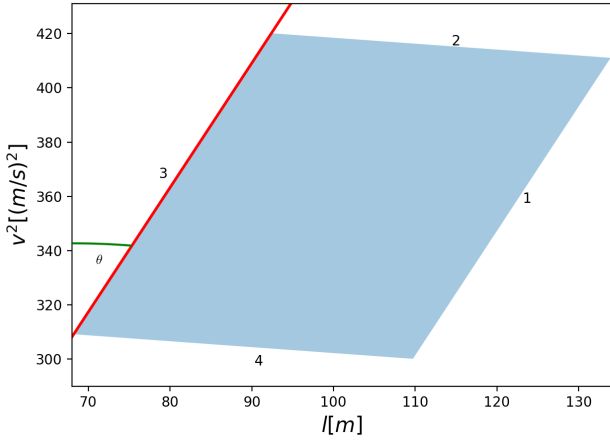


Fig. 4. Illustration of the facets of the convex polytope. Since the safety label relies in this case on a single halfplane (drawn in red), we only need to consider facet 3.

to the left of this line, would unnecessarily label samples from $I_0$ with label *one*, increasing the total number of label disagreements. However, the constructed MILP aims at minimizing the total number of label disagreements; as such the cyan dotted halfplane would always be preferable, rendering any samples in the cyan-colored area redundant. If we have non-zero disagreement with respect to the samples from $I_1$, the halfplane will lie fur-

ther to the right of the cyan dotted one. Similarly, to obtain zero disagreement between the hypothesis and the samples in $I_0$, the halfplane constructed by the hypothesis will need to lie to the right of all samples in $I_0$. Following a similar argumentation as before, any sample in the magenta-colored area will not change the solution of the MILP.

Since we know that the samples in both the blue and the magenta regions will not affect the solution of the MILP, we can discard these samples prior to computing the hypothesis, resulting in only the red samples in Figure 5 being considered. Following similar arguments, it can be possible to discard redundant samples also in the setting of higher-order convex polytopes. However, generalizing the proposed methodology to achieve this is case-dependent and is not pursued further here.

### 5.2   Simulation results

While no knowledge of the distribution of the samples $(l, v^2)$ needs to be known for generating the hypothesis, for simulation purposes we draw $l$ from a uniform distribution over the interval $[40m, 120m]$ and draw $v^2$ from a normal distribution with mean $\overline{v^2} = (70km/h)^2$ and standard deviation $\sigma_{v^2} = (20km/h)^2$. The performance of the brakes at each time step will deteriorate by a factor of $\omega_F$, i.e. $F_{i+1} = \omega_F F_i$, where $\omega_F$ is a random variable drawn from a normal distribution with mean

$\mu = (1 - 3 \cdot 10^{-7})$ and standard deviation $\sigma = 10^{-6}$. The initial car mass is $m = 900kg$ and will randomly change by a factor of $\omega_m$, where $\omega_m$ is a random variable drawn from a normal distribution with mean $\mu = 1$ and standard deviation $\sigma = 10^{-3}$.

For the construction of the hypothesis, the confidence level is chosen as $\delta = 10^{-6}$ with an accuracy of $\epsilon = 1\%$. For the satisfaction of Assumption 2, we choose $\frac{1}{m} \sum_{i=1}^{m} \mathrm{er}(f_i, f_{m+1})$ to be bounded by $\overline{\mu} \leq 2\%$ and $\underline{\mu} \geq 0.78\%$. By Theorem 2 it then follows that we need at least $119, 237$ samples to accurately predict the safety label of the subsequent timesteps.

Using the aforementioned discarding approach, we can discard 95% of the samples prior to instantiating the MILP constraints. The discarding approach is illustrated in Figure 5 where, for the purpose of visualization, we omit samples close to one another to prevent the image from being cluttered. The hypothesis in minimal dis-
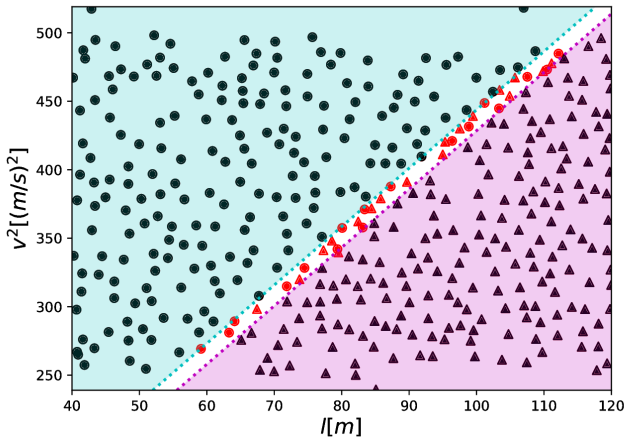


Fig. 5. All samples in black are discarded, while the red samples are kept for the computation of the hypothesis. This results in 95% of the samples being omitted, greatly improving the computational feasibility.

agreement with the labeled samples, computed by means of the MILP (34 - 36), is shown in Figure 6. Solving the MILP took 561 seconds, making the deployment of the approach computationally feasible. The number of violations, $v$, is 1335. We have made all code for generating and reproducing our results available online [3].

We empirically validate our risk level by means of Monte Carlo simulations. For each run, we generate a new labeling mechanism $f_{m+1}$, corresponding to the random deterioration of the braking force and change to the vehicle mass. We then draw 5000 samples for which we evaluate the corresponding label (by means of $f_{m+1}$) and

[3] www.vertovec.info/code/learning-moving-targets

compare this with the label assigned by means of the hypothesis constructed by our methodology, thus calculating $\widehat{\mathrm{er}}_m(f_{m+1}, h_m)$. We repeat this for 500 runs, each time generating a new label $f_{m+1}$. In Figure 7 the frequency of certain $\widehat{\mathrm{er}}_m(f_{m+1}, h_m)$ values is shown.

Recall that $\mu$ was upper bounded by 2%, such that for the chosen $\epsilon = 1\%$, Theorem 2 implies that $\mathrm{er}(f_{m+1}, h_m) \leq 4\overline{\mu} + \epsilon = 9\%$ with high confidence. The Monte Carlo simulation supports this, with the average empirical disagreement $\widehat{\mathrm{er}}_m(f_{m+1}, h_m)$ being approximately 2.4% (see Figure 7), well below the theoretically predicted 9%.
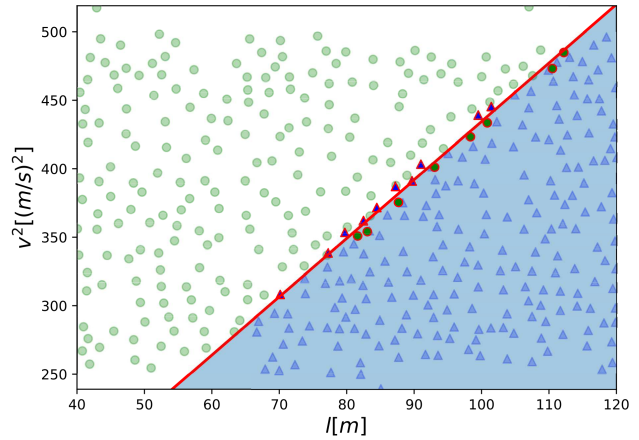


Fig. 6. Generated hypothesis; we only show the halfplane responsible for the labeling, illustrated by red. For visual ease, we randomly omit samples close to one another. Red samples are violations as defined in (31).
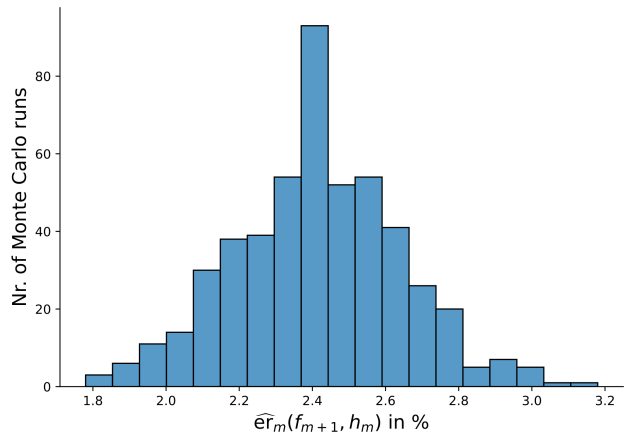


Fig. 7. Empirical distribution of the disagreement $\widehat{\mathrm{er}}_m(f_{m+1}, h_m)$, constructed by means of 500 Monte Carlo runs.

## 6 Conclusion

We considered learning a moving target from a finite set of samples and showed that, when the labeling mechanism changes in a structured manner, it remains PAC learnable, meeting certain accuracy-confidence levels. Furthermore, for the class of convex polytopes, we presented a constructive method to generate the hypothesis based on a Mixed Integer Linear Program (MILP). We illustrated the applicability of our theoretical developments to a case study involving an Autonomous Emergency Braking (AEB) system. Future work aims at considering the distribution according to which samples are drawn to also be changing, similarly to [22].

## A Recovering the constant target case

Following Remark 3, we show in the next result how Theorem 2 specializes to obtain probabilistic guarantees for a minimal disagreement hypothesis $h_m \in M_m$, for the case where the target is constant.

**Theorem 3.** *Fix $\epsilon, \delta \in (0,1)$. Denote by $d$ the VC dimension of $\mathcal{H}$, and consider $\underline{\mu} = \overline{\mu} = 0$. If we choose $m \geq m_0(\epsilon, \delta)$, where*

$$m_0(\epsilon, \delta) = \frac{5}{\epsilon}\left( \ln\frac{4}{\delta} + d\ln\frac{40}{\epsilon} \right), \qquad (A.1)$$

*we then have that for any $h_m \in M_m$,*

$$\mathbb{P}^m\{(x_1, \ldots, x_m) \in \mathrm{X}^m : \\ \mathrm{er}(f_{m+1}, h_m) \leq \epsilon\} \geq 1 - \delta. \quad (A.2)$$

*Proof.* Fix any $\epsilon, \delta \in (0,1)$. We follow the same proofline with Theorem 2, but since $\underline{\mu} = \overline{\mu} = 0$, all target functions are identical. To this end, let $f_i = f$, for all $i = 1, \ldots, m, m+1$. We define the following event:

$$E = \{(x_1, \ldots, x_m) \in \mathrm{X}^m : \mathrm{er}(f, h_m) > \epsilon\},$$
$$\widehat{E} = \{(x_1, \ldots, x_m) \in \mathrm{X}^m : \widehat{\mathrm{er}}_m(f, h_m) = 0\}. \quad (A.3)$$

$\widehat{E}$ is the set of $m$-multisamples for which the empirical average disagreement between the (constant) target and the hypothesis, namely, $\frac{1}{m}\sum_{i=1}^m |f(x_i) - h_m(x_i)|$, is equal to zero. Notice that since $h_m \in M_m$ (a minimal disagreement hypothesis), for any $m$-multisample, $\sum_{i=1}^m |f(x_i) - h_m(x_i)| \leq \sum_{i=1}^m |f(x_i) - h(x_i)|$ for any $h \in \mathcal{H}$. Since the target function $f$ itself is an element of $\mathcal{H}$, taking $h = f$ in the aforementioned statement directly leads to $\sum_{i=1}^m |f(x_i) - h_m(x_i)| \leq 0$, and hence $\mathbb{P}^m\{\widehat{E}\} = 1$.

To establish (A.2) it suffices to show that $\mathbb{P}^m\{E\} \leq \delta$. To this end, we have that

$$\mathbb{P}^m\{E\} \\ = \mathbb{P}^m\{E \cap \widehat{E}\} \\ = \mathbb{P}^m\{(x_1, \ldots, x_m) \in \mathrm{X}^m : \mathrm{er}(f, h_m) > \epsilon \\ \text{and } \widehat{\mathrm{er}}_m(f, h_m) = 0\}. \quad (A.4)$$

where the first equality is since $\mathbb{P}^m\{\widehat{E}\} = 1$, and the second one follows from the definition of $E$ and $\widehat{E}$.

Notice that (A.4) takes the form of (7), with $f_{m+1}$, $h_m$ and 0 in place of $f$, $h$ and $\rho$, respectively. Theorem 1 implies then that

$$m \geq \frac{5}{\epsilon}\left( \ln\frac{4}{\delta} + d\ln\frac{40}{\epsilon} \right) \implies \mathbb{P}^m\{E \cap \widehat{E}\} \leq \delta. \quad (A.5)$$

Therefore, by (A.4) and (A.5) we have that $\mathbb{P}^m\{E\} \leq \delta$, thus concluding the proof. $\qquad \square$

## References

[1] T. Alamo, R. Tempo, and E. F. Camacho. Randomized strategies for probabilistic solutions of uncertain feasibility and optimization problems. *IEEE Transactions on Automatic Control*, 54(11):2545–2559, 2009.

[2] P. L. Bartlett, S. Ben-David, and S. R. Kulkarni. Learning changing concepts by exploiting the structure of change. *Machine Learning*, 41(2):153–174, 2000.

[3] R. D. Barve and P. M. Long. On the complexity of learning from drifting distributions. *Information and Computation*, 138(2):170–193, Nov. 1997.

[4] A. Bemporad and M. Morari. Control of systems integrating logic, dynamics, and constraints. *Automatica*, 35(3):407–427, 1999.

[5] G. Calafiore and M. Campi. The scenario approach to robust control design. *IEEE Transactions on Automatic Control*, 51(5):742–753, May 2006.

[6] G. C. Calafiore. Random convex programs. *SIAM Journal on Optimization*, 20(6):3427–3464, 2010.

[7] M. C. Campi and S. Garatti. The exact feasibility of randomized solutions of uncertain convex programs. *SIAM Journal on Optimization*, 19(3):1211–1230, 2008.

[8] M. C. Campi and S. Garatti. A sampling-and-discarding approach to chance-constrained optimization: Feasibility and optimality. *Journal of Optimization Theory and Applications*, 148(2):257–280, Oct. 2010.

[9] M. C. Campi and S. Garatti. Wait-and-judge scenario optimization. *Mathematical Programming*, 167(1):155–189, July 2016.

[10] M. C. Campi and S. Garatti. *Introduction to the Scenario Approach*. Society for Industrial and Applied Mathematics, Philadelphia, PA, Nov. 2018.

[11] M. C. Campi and S. Garatti. Compression, generalization and learning, 2023.

[12] M. C. Campi, S. Garatti, and M. Prandini. The scenario approach for systems and control design. *Annual Reviews in Control*, 33(2):149–157, 2009.

[13] M. C. Campi, S. Garatti, and F. A. Ramponi. A general scenario theory for nonconvex optimization and decision making. *IEEE Transactions on Automatic Control*, 63(12):4067–4078, 2018.

[14] K. Crammer, E. Even-Dar, Y. Mansour, and J. W. Vaughan. Regret minimization with concept drift. In *Proceedings of the 23rd Annual Conference on Learning Theory (COLT)*, 2010.

[15] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4):633–679, Aug. 2019.

[16] F. Fele and K. Margellos. Probably approximately correct nash equilibrium learning. *IEEE Transactions on Automatic Control*, 66(9):4238–4245, 2021.

[17] S. Garatti and M. C. Campi. Risk and complexity in scenario optimization. *Mathematical Programming*, 191(1):243–279, Jan. 2022.

[18] S. Hanneke, V. Kanade, and L. Yang. *Learning with a Drifting Target Concept*, page 149–164. Springer International Publishing, 2015.

[19] D. P. Helmbold and P. M. Long. Tracking drifting concepts by minimizing disagreements. *Machine Learning*, 14(1):27–45, Jan. 1994.

[20] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.

[21] A. Kuh, T. Petsche, and R. Rivest. Learning time-varying concepts. In R. Lippmann, J. Moody, and D. Touretzky, editors, *Advances in Neural Information Processing Systems*, volume 3. Morgan-Kaufmann, 1990.

[22] P. M. Long. The complexity of learning according to two models of a drifting environment. *Machine Learning*, 37(3):337–354, 1999.

[23] K. Margellos, M. Prandini, and J. Lygeros. On the connection between compression learning and scenario based single-stage and cascading optimization problems. *IEEE Transactions on Automatic Control*, 60(10):2716–2721, 2015.

[24] M. Morari. Hybrid system analysis and control via mixed integer optimization. *IFAC Proceedings Volumes*, 34(25):1–12, 2001. 6th IFAC Symposium on Dynamics and Control of Process Systems 2001, Jejudo Island, Korea, 4-6 June 2001.

[25] L. Romao, K. Margellos, and A. Papachristodoulou. Probabilistic feasibility guarantees for convex scenario programs with an arbitrary number of discarded constraints. *Automatica*, 149:1–9, 2023.

[26] L. Romao, A. Papachristodoulou, and K. Margellos. On the exact feasibility of convex scenario programs with discarded constraints. *IEEE Transactions on Automatic Control*, 68(4):1986–2001, 2023.

[27] R. Tempo, G. Calafiore, and F. Dabbene. *Randomized algorithms for analysis and control of uncertain systems*. Communications and control engineering series. Springer, London, 2005.

[28] M. Vidyasagar. *Learning and Generalisation*. Springer London, 2003.