

# Online Machine Teaching under Learner Uncertainty: Gradient Descent Learners of a Quadratic Loss\*

Belen Martin-Urcelay<sup>†</sup>, Christopher J. Rozell<sup>†</sup>, and Matthieu R. Bloch<sup>†</sup>

---

**Abstract.** We revisit the framework of online machine teaching, a special case of active learning in which a teacher with full knowledge of a model attempts to train a learner by adaptively presenting examples. While online machine teaching example selection strategies are typically designed assuming omniscience, i.e., the teacher has absolute knowledge of the learner state, we show that efficient machine teaching is possible even when the teacher is uncertain about the learner initialization. Specifically, we consider the case of learners that perform gradient descent of a quadratic loss to learn a linear classifier, and propose an online machine teaching algorithm in which the teacher simultaneously learns the learner state while teaching the learner. We theoretically show that the learner’s mean square error decreases exponentially with the number of examples, thus achieving a performance similar to the omniscient case and outperforming two stage strategies that first attempt to make the teacher omniscient before teaching. We empirically illustrate our approach in the context of a cross-lingual sentiment analysis problem.

**Key words.** Machine teaching, active learning, online example selection, unknown initialization

**MSC codes.** 68W27, 68W40, 93C55

**1. Introduction.** The size of datasets used in modern machine learning has grown many-fold over the last decade, making the training of models on entire datasets frequently impractical [10], either because of the associated training time, training cost or incurred energy consumption and environmental cost. To circumvent these constraints, it is now common to only train models on a subset of examples. Using naive data selection strategies, such as randomly sampling a dataset, typically requires more examples than intentional strategies, such as active learning, by which the machine learning algorithm adaptively requests the labels of certain data points from a large pool of unlabeled examples [26]. Active learning has been successfully applied to a wide variety of settings, such as natural language processing [33, 4], data embedding [29, 7] or source localization [19, 21]. Machine Teaching (MT) considers a variation of the setup in which a knowledgeable expert knowing the ground truth model, the *teacher*, selects the examples fed to the machine learning algorithm, the *learner*. The aim of machine teaching is to exploit the teacher’s knowledge and identify the smallest set of examples to train the learner [34].

Machine teaching has proved useful in a variety of settings, ranging from an illustrative 1-Dimensional threshold classifier [35] to complex vocabulary learning platforms [30]. A crucial requirement in early machine teaching algorithms has been the need for consistent learners [8, 3], which directly discard all the hypotheses that do not agree with any training example.

---

\*Submitted to the editors May 27th 2024.

**Funding:** This work was funded by the Rafael del Pino Foundation, the Fulbright association and the International Peace Scholarship by P.E.O.

<sup>†</sup>Department of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA ([burcelay3@gatech.edu](mailto:burcelay3@gatech.edu), [crozell@gatech.edu](mailto:crozell@gatech.edu), [matthieu.bloch@ece.gatech.edu](mailto:matthieu.bloch@ece.gatech.edu) ).

Therefore, these algorithms do not perform well in the presence of noisy labels. The consistency requirement has been relaxed in recent literature [17, 16] by introducing the concept of *omniscient teaching*. An omniscient teacher possesses full knowledge of the learner, i.e., it is able to observe the state and dynamics of the learner during training. Under certain smoothness assumptions, the selection of examples reduces to a constrained convex optimization problem, for which a greedy machine teaching algorithm as in [17] achieves an exponential speed-up compared to random example selection. Nevertheless the omniscience requirement may pose practical implementation challenges [8].

First, we note that the initial state of an algorithm is often unknown. This is the case in adversarial attacks, such as training-state poisoning [12], in which attackers lack precise knowledge about the initial state of the targeted system, such as a spam filter. Unknown initial states also result from warm-starts [23], a technique by which pre-trained models are used to accelerate the learning process or transfer knowledge from related tasks. Second, we note that teacher and learner may operate in different feature spaces. For example, words may be embedded in different spaces for different languages and the mapping between language spaces may be unknown.

An existing approach to address the lack of omniscience is learning for omniscience [18, 16], which consists in introducing a preliminary probing phase during which the teacher queries the learner until enough feedback is gathered to accurately approximate the learner initial state. Unfortunately, this strategy requires many interactions between the teacher and learner during which the learner does not improve its model.

The present work aims to tackle the above limitations by developing an efficient machine teaching algorithm capable of boosting the convergence speed of learners even when the teacher is not fully omniscient. Our algorithm addresses the challenges related to unknown learner starting states and unknown orthogonal mappings between the learner and teacher feature spaces. Our main contribution is in realizing that jointly teaching the learner while estimating its parameters may offer significant and previously not identified gains. In particular:

1. We develop a non-omniscient machine teaching algorithm for gradient descent learners of a quadratic loss with unknown initializations. We prove that our algorithm achieves an exponential speed up compared to random example selection, without an explicit probing phase to estimate the learner initialization. Additionally, the exponential convergence guarantees hold under unknown orthonormal mappings between learner and teacher.
2. We draw connections between machine teaching and control theory. These connections allow us to leverage well-studied techniques, such as Kalman filters and Riccati recursions, to obtain theoretical guarantees on learning performance.
3. We empirically demonstrate the advantages of our framework over random sampling and probing based techniques, using the teaching of a binary sentiment classifier across languages as an example.

**2. Framework.** We now detail the framework of the machine teaching problem and introduce simplifying assumptions to make analytical progress in the non-omniscient setting. As illustrated in Figure 1, let the learner be a machine learning model parameterized by  $\theta$ . For instance,  $\hat{\theta}$  could represent an effective decision boundary. Machine teaching aims to guide

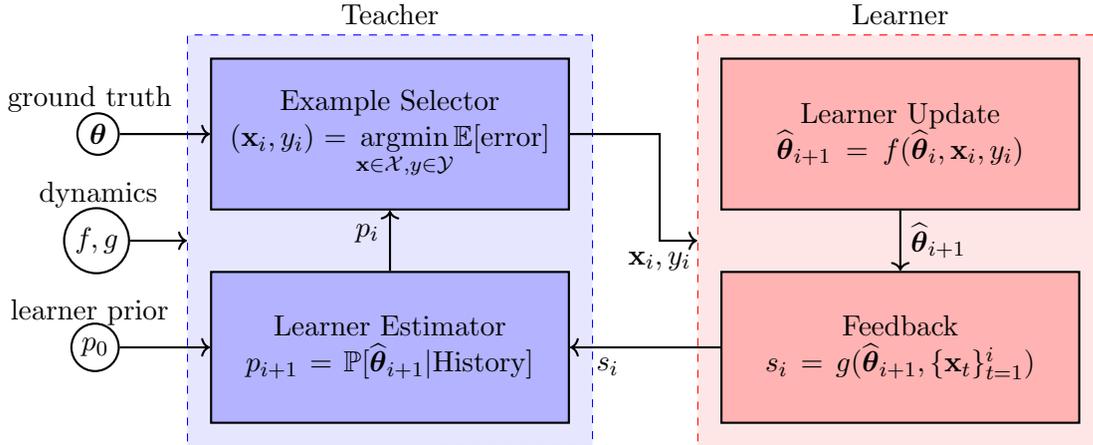


Figure 1: Block diagram of the online machine teaching framework. The goal of the teacher is to steer the learner towards the ground truth  $\theta$ , while simultaneously learning about the learner state  $\hat{\theta}_{i+1}$ .

80 the learner’s learned parameter,  $\hat{\theta}$ , towards the ground truth,  $\theta$ . Let the teacher be an entity  
 81 with knowledge of the ground truth  $\theta$  and selecting the examples presented to the learner.  
 82 At each time-step  $i$ , the teacher first presents an example and label pair  $(\mathbf{x}_i, y_i)$  from a pre-  
 83 determined pool  $(\mathcal{X}, \mathcal{Y})$  to the learner. The learner then uses the example to update its model  
 84  $\hat{\theta}_{i+1} = f(\hat{\theta}_i, \mathbf{x}_i, y_i)$  for some known function  $f$ . The learner may also provide some feedback  
 85 with information about its current state to the teacher  $s_i = g(\hat{\theta}_{i+1}, \{\mathbf{x}_t\}_{t=1}^i)$ , where  $g$  is some  
 86 known function. In the case of an omniscient teacher this feedback provides the exact learner  
 87 state.

88 Although we assume that the teacher knows the function  $f$  that the learner uses to update  
 89 its state, we emphasize that the teacher is not omniscient. Namely, the teacher does not  
 90 know the starting point of the learner,  $\hat{\theta}_0$ . Instead, we assume the teacher starts with a  
 91 prior Gaussian probability distribution  $\mathbf{p}_0$  for  $\hat{\theta}_0$ . We shall also consider the case in which the  
 92 teacher and learner do not share the same feature spaces: when the teacher selects an example  
 93  $\mathbf{x}$ , the learner observes  $\hat{\mathbf{x}} = \mathcal{G}(\mathbf{x})$ , where  $\mathcal{G}$  is an unknown orthonormal mapping between the  
 94 teacher and learner feature spaces.

95 For analytical tractability, we restrict our attention to a learner that performs gradient  
 96 descent to minimize the quadratic loss  $l(\hat{\theta}) := \frac{1}{2} \|\hat{\theta}^T \mathbf{x} - y\|_2^2$ . At each iteration the learner  
 97 updates its state according to

$$98 \quad (2.1) \quad \hat{\theta}_{i+1} = \hat{\theta}_i - \tau \left( \hat{\theta}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i,$$

99 where  $\tau \in \mathbb{R}^+$  is the learning rate, assumed known to the teacher. We denote the maximum  
 100 norm of the states by  $P$ , i.e.,  $\max_i \|\hat{\theta}_i\|_2^2 \leq P$ . We specifically look at teaching a linear binary  
 101 classifier  $\theta$ , s.t.  $\|\theta\|_2^2 \leq P$ . The classifier labels any example  $\mathbf{x} \in \mathcal{X}$  as  $y = \text{sign}(\theta^T \mathbf{x})$ . In  
 102 principle, one could attempt to extend the linear classifier to non-linear problems by mapping

103 the original non-linear space into a higher-dimensional feature space in which the data is  
 104 linearly separable, though this mapping is often hard to find in practice.

105 We consider synthesis based teaching [17] by which the teacher may provide any example  
 106 within a ball  $\mathcal{X} : \{\mathbf{x} = [1, x_1, \dots, x_{d-1}]^T \in \mathbb{R}^d; \|\mathbf{x}\|_2^2 \leq P_{\mathbf{x}}\}$ , together with any binary label in  
 107  $\mathcal{Y} : \{-1, 1\}$ . Following standard practice, the first coordinate of the examples is set to 1 to  
 108 allow for the parameter  $\boldsymbol{\theta}$  to account for both the direction and the offset of the hyperplane  
 109 characterizing the classifier. The freedom to synthetically generate examples may lead to non-  
 110 semantically-meaningful examples. To maintain interpretability, one can restrict the examples  
 111 space  $\mathcal{X}$  to data points that a teacher generates with a Variational AutoEncoder (VAE) trained  
 112 from a pre-defined dataset of meaningful examples. This restriction forces synthetic examples  
 113 to resemble the original training dataset, and thus be interpretable [24].

114 **3. Theoretical Guarantees.** Existing online teachers base their example selection criteria  
 115 on their knowledge of the learner state, which naturally prompts a number of questions: How  
 116 does a teacher handle learner uncertainty? Are there any convergence guarantees in that  
 117 case? We tackle these questions under two different settings: when the teacher receives no  
 118 information from the learner, and when the teacher receives some noisy feedback from the  
 119 learner at each iteration.

120 **3.1. Simultaneous Machine Teaching and Learning (SMTL) without Feedback.** As  
 121 a baseline, we first consider the situation in which the teacher receives no feedback from  
 122 the learner. At each iteration, the teacher only communicates with the learner via a single  
 123 example-label pair. We propose a greedy algorithm that chooses the example-label pair that  
 124 most reduces the expected error of the learned parameter from one iteration to the next. The  
 125 algorithm is motivated by the decomposition of the Mean-Square Error (MSE) of the learned  
 126 parameter as

$$\begin{aligned}
 127 \quad \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_{i+1} - \boldsymbol{\theta}\|_2^2 \mid H_i \right] &= \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \tau \left( \widehat{\boldsymbol{\theta}}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i - \boldsymbol{\theta}\|_2^2 \mid H_i \right] \\
 128 \quad &= \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \mid H_i \right] - \tau T(\mathbf{x}_i, y_i, \boldsymbol{\mu}_i, \mathbf{C}_i),
 \end{aligned}$$

129 where  $H_i := \{\mathbf{p}_0, (\mathbf{x}_t, y_t)_{t=1}^i\}$  refers to the history of past examples and labels, as well as  
 130 the prior distribution of  $\widehat{\boldsymbol{\theta}}_0$  known by the teacher. We set  $\boldsymbol{\mu}_i := \mathbb{E} \left[ \widehat{\boldsymbol{\theta}}_i \mid H_i \right]$  and  $\mathbf{C}_i :=$   
 131  $\mathbb{E} \left[ \widehat{\boldsymbol{\theta}}_i \widehat{\boldsymbol{\theta}}_i^T - \boldsymbol{\mu}_i \boldsymbol{\mu}_i^T \mid H_i \right]$  to represent the expectation and covariance matrix of the learner state,  
 132 respectively. We let  $T(\mathbf{x}_i, y_i, \boldsymbol{\mu}_i, \mathbf{C}_i) = \mathbb{E} \left[ 2(\widehat{\boldsymbol{\theta}}_i^T \mathbf{x}_i - y_i) \langle \widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}, \mathbf{x}_i \rangle - \tau(\widehat{\boldsymbol{\theta}}_i^T \mathbf{x}_i - y_i)^2 \|\mathbf{x}_i\|_2^2 \mid H_i \right]$   
 133 represent the expected improvement, i.e., how much the teacher expects the MSE to reduce  
 134 from time-step  $i$  to  $i + 1$ .

135 The proposed policy selects the example-label pair that most reduces the error from one  
 136 step to the next. Specifically, at time  $i$ , the teacher selects

$$137 \quad (3.1) \quad (\widehat{\mathbf{x}}_i, \widehat{y}_i) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}} T(\mathbf{x}, y, \boldsymbol{\mu}_i, \mathbf{C}_i).$$

138 **Lemma 3.1.** *The objective function in (3.1) is equivalent to*

$$139 \quad (3.2) \quad T(\mathbf{x}, y, \boldsymbol{\mu}_i, \mathbf{C}_i) = \underbrace{(2 - \tau \|\mathbf{x}\|_2^2) \mathbf{x}^T \mathbf{C}_i \mathbf{x}}_{\text{exploration}} + 2 \underbrace{(\boldsymbol{\theta}^T - \boldsymbol{\mu}_i^T) (y - \boldsymbol{\mu}_i^T \mathbf{x}) \mathbf{x}}_{\text{exploitation}} - \tau \underbrace{\|\mathbf{x}\|_2^2 (y - \boldsymbol{\mu}_i^T \mathbf{x})^2}_{\text{regularization}}.$$

140 **Lemma 3.1** follows from algebraic manipulations that are detailed in Section SM1.1. Note  
 141 that  $T$  is a fourth degree polynomial with  $d$  unknowns:  $y \in \{-1, 1\}$  and all but the first  
 142 coordinate of  $\mathbf{x}$ . The unconstrained absolute maximum of  $T$  may be calculated with standard  
 143 software such as Matlab’s *fmincon* function. Additionally, the teacher does not need to track  
 144 the probability distribution of the learner. The teacher only needs to track the first and second  
 145 order moments to compute equation (3.1) and select the appropriate example.

146 The maximization of (3.2) implicitly accounts for the trade-off between estimating the  
 147 learner state and teaching the ground truth to the learner. Under high uncertainty, corre-  
 148 sponding to large values in the covariance  $\mathbf{C}_i$ , the first term in (3.2) dominates. The first  
 149 term is an *exploration* component that promotes examples aligned with the direction of high-  
 150 est covariance, i.e., the examples that are most likely to decrease the teacher uncertainty  
 151 about the learner state. On the other hand, the second term promotes examples that steer  
 152 the estimated learner towards the ground truth, so the second term may be interpreted as an  
 153 *exploitation* component. As the distance between the estimated learner state and the ground  
 154 truth decreases, so does the relative weight of the exploitation term. The transition between  
 155 phases focused on exploitation and exploration is further analyzed in subsection SM3.1, which  
 156 examines the evolution of different sources of error. Lastly, the third term in (3.2) acts as a  
 157 regularizer that discourages the norm of the gradient from being too large. This *regularization*  
 158 term avoids abrupt and overly large updates in the learner state.

159 After sending the example and label pair to the learner, the teacher updates its estimation  
 160 of the learner state following the known dynamical model of the learner. The mean and  
 161 covariance are updated as

$$\begin{aligned} 162 \quad \boldsymbol{\mu}_{i+1} &:= \mathbb{E} \left[ \widehat{\boldsymbol{\theta}}_{i+1} \mid H_{i+1} \right] = \mathbb{E} \left[ \widehat{\boldsymbol{\theta}}_i - \tau \left( \widehat{\boldsymbol{\theta}}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i \mid H_i \right] = \boldsymbol{\mu}_i - \tau \left( \boldsymbol{\mu}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i. \\ 163 \quad \mathbf{C}_{i+1} &:= \mathbb{E} \left[ \widehat{\boldsymbol{\theta}}_{i+1} \widehat{\boldsymbol{\theta}}_{i+1}^T - \boldsymbol{\mu}_{i+1} \boldsymbol{\mu}_{i+1}^T \mid H_{i+1} \right] = \mathbb{E} \left[ \widehat{\boldsymbol{\theta}}_{i+1} \widehat{\boldsymbol{\theta}}_{i+1}^T - \boldsymbol{\mu}_{i+1} \boldsymbol{\mu}_{i+1}^T \mid H_i \right] \\ 164 \quad &= \mathbb{E} \left[ \left( \widehat{\boldsymbol{\theta}}_i - \tau \left( \widehat{\boldsymbol{\theta}}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i \right) \left( \widehat{\boldsymbol{\theta}}_i - \tau \left( \widehat{\boldsymbol{\theta}}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i \right)^T \mid H_i \right] \\ 165 \quad &\quad - \mathbb{E} \left[ \left( \boldsymbol{\mu}_i - \tau \left( \boldsymbol{\mu}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i \right) \left( \boldsymbol{\mu}_i - \tau \left( \boldsymbol{\mu}_i^T \mathbf{x}_i - y_i \right) \mathbf{x}_i \right)^T \mid H_i \right] \\ 166 \quad &= \mathbf{C}_i - \tau \mathbf{C}_i \mathbf{x}_i \mathbf{x}_i^T - \tau \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}_i + \tau^2 \mathbf{x}_i^T \mathbf{C}_i \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}_i. \end{aligned}$$

167 The first equality holds because, given the past history  $H_i$ , the teacher selects the next  
 168 example-label pair in a deterministic way: in the absence of feedback,  $H_{i+1}$  is completely  
 169 determined by  $H_i$ . We outline this approach, which we call Simultaneous Machine Teaching  
 170 and Learning (SMTL), in Algorithm 3.1.

171 Next, we characterize the convergence rate that SMTL provides. We recall the guarantees  
 172 for omniscient teaching as a baseline against the proposed algorithm.

**Algorithm 3.1** SMTL

---

```

1:  $\boldsymbol{\mu}_0, \mathbf{C}_0 \leftarrow p_0$ 
2: for  $i = 0, 1, 2, \dots$  do
3:   Select example:
       $(\mathbf{x}_i, y_i) \leftarrow \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}} T(\mathbf{x}, y, \boldsymbol{\mu}_i, \mathbf{C}_i)$ 
4:   Update estimations about learner:
       $\boldsymbol{\mu}_{i+1} \leftarrow \boldsymbol{\mu}_i - \tau (\boldsymbol{\mu}_i^T \mathbf{x}_i - y_i) \mathbf{x}_i$ 
       $\mathbf{C}_{i+1} \leftarrow \mathbf{C}_i - \tau \mathbf{C}_i \mathbf{x}_i \mathbf{x}_i^T - \tau \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}_i + \tau^2 \mathbf{x}_i^T \mathbf{C}_i \mathbf{x}_i \mathbf{x}_i^T$ 
5: end for

```

---

173 **Theorem 3.2.** [Adapted from [17, Theorem 4]] Consider a synthesis based omniscient  
174 teacher and a learner with updates given by (2.1). If  $\forall \widehat{\boldsymbol{\theta}}_i, \exists \gamma \in \mathbb{R}$  with  $|\gamma| \leq \frac{\sqrt{P}}{\|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2}$ ,  $\nu(\gamma) \in \mathbb{R}$   
175 and  $y' \in \{-1, 1\}$  s.t.  $0 < \tau (\widehat{\boldsymbol{\theta}}_i^T \mathbf{x}' - y') \mathbf{x}' \leq \nu(\gamma) < \frac{1}{\tau}$  for  $\mathbf{x}' = \gamma(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta})$ , then,

$$176 \quad \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \leq (1 - \tau\nu)^{2i} \|\widehat{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}\|_2^2.$$

177 **Theorem 3.2** applies to the specific case of our framework in which  $\forall i \boldsymbol{\mu}_i = \widehat{\boldsymbol{\theta}}_i$  and  $\mathbf{C}_i =$   
178  $\mathbf{0}$ . The theorem guarantees that an omniscient teacher teaches a classifier to a gradient  
179 descent learner exponentially fast with the number of examples, thereby offering a significant  
180 improvement compared to the linear convergence obtained when randomly selecting examples  
181 [22]. The auxiliary variables  $\gamma$  and  $\nu(\gamma)$  are related to the convergence speed. The guarantees  
182 for an omniscient teacher provide a baseline for the MSE of non-omniscient teachers. The  
183 following theorem offers a convergence rate guarantee in the non-omniscient scenario without  
184 feedback.

185 **Theorem 3.3.** Consider a synthesis based teacher following SMTL and a learner with up-  
186 dates given by (2.1). If  $\forall \widehat{\boldsymbol{\theta}}_i, \exists \gamma \in \mathbb{R}$  with  $|\gamma| \leq \frac{\sqrt{P}}{\|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2}$ ,  $\nu(\gamma) \in \mathbb{R}$  and  $y' \in \{-1, 1\}$  s.t.  
187

$$188 \quad (3.3) \quad 0 < \left( \widehat{\boldsymbol{\theta}}_i^T \gamma(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' - \frac{1}{\tau\gamma} \right)^2 \leq \nu^2 < \frac{1}{\tau^2\gamma^2},$$

189 then,

$$190 \quad \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \mid H_{i-1} \right] \leq (\tau\gamma\nu)^{2i} \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}\|_2^2 \mid H_0 \right].$$

191 *Proof.* We base our proof on [17, Theorem 4]. The expected evolution of the MSE from  
192 iteration  $i$  to iteration  $i + 1$  is described by

$$193 \quad (3.4) \quad \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_{i+1} - \boldsymbol{\theta}\|_2^2 \mid H_i \right] = \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \mid H_i \right] - \tau T(\widehat{\mathbf{x}}_i, \widehat{y}_i, \boldsymbol{\mu}_i, \mathbf{C}_i)$$

194 where

$$195 \quad T(\widehat{\mathbf{x}}_i, \widehat{y}_i, \boldsymbol{\mu}_i, \mathbf{C}_i) = \mathbb{E} \left[ 2 \left( \widehat{\boldsymbol{\theta}}_i^T \widehat{\mathbf{x}}_i - \widehat{y}_i \right) \langle \widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}, \widehat{\mathbf{x}}_i \rangle - \tau \left( \widehat{\boldsymbol{\theta}}_i^T \widehat{\mathbf{x}}_i - \widehat{y}_i \right)^2 \|\widehat{\mathbf{x}}_i\|_2^2 \mid H_i \right]$$

196 represents the expected MSE improvement at the  $i$ -th iteration when selecting the example-  
 197 label pair  $(\widehat{\mathbf{x}}_i, \widehat{y}_i)$ . We analyze the objective function  $T$  at  $(\mathbf{x}' = \gamma(\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}), y')$ , for some aux-  
 198 iliary parameter  $\gamma \in \mathbb{R}$ , to obtain the following lower bound:

(3.5)

$$\begin{aligned}
 199 \quad & T(\mathbf{x}', y', \boldsymbol{\mu}_i, \mathbf{C}_i) \\
 200 \quad &= \mathbb{E} \left[ 2 \left( \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' \right) \langle \widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}, \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) \rangle - \tau \left( \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' \right)^2 \|\gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta})\|_2^2 \middle| H_i \right] \\
 201 \quad &= \gamma \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \left( 2 \left( \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' \right) - \tau \gamma \left( \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' \right)^2 \right) \middle| H_i \right] \\
 202 \quad &= \tau \gamma^2 \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \left( \frac{1}{\tau^2 \gamma^2} - \left( \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' - \frac{1}{\tau \gamma} \right)^2 \right) \middle| H_i \right] \\
 203 \quad &= \frac{1}{\tau} \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \middle| H_i \right] - \tau \gamma^2 \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \left( \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - y' - \frac{1}{\tau \gamma} \right)^2 \middle| H_i \right] \\
 204 \quad &\geq \frac{1}{\tau} \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \middle| H_i \right] (1 - (\tau \gamma \nu)^2),
 \end{aligned}$$

205 where the last inequality holds because of Assumption (3.3).

206 The teacher selects the example-label pair that maximizes the expected improvement in  
 207 MSE. By definition of argmax in (3.1),  $T(\widehat{\mathbf{x}}_i, \widehat{y}_i, \boldsymbol{\mu}_i, \mathbf{C}_i) \geq T(\mathbf{x}', y', \boldsymbol{\mu}_i, \mathbf{C}_i), \forall \mathbf{x}' \in \mathcal{X}, \forall y' \in \mathcal{Y}$ .  
 208 Combining this inequality with (3.5) and (3.4) we obtain

$$\begin{aligned}
 209 \quad & \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_{i+1} - \boldsymbol{\theta}\|_2^2 \middle| H_i \right] \leq \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \middle| H_{i-1} \right] - \tau T(\mathbf{x}', y', \boldsymbol{\mu}_i, \mathbf{C}_i) \\
 210 \quad & \leq \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \middle| H_{i-1} \right] - \tau \frac{1}{\tau} \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \middle| H_{i-1} \right] (1 - (\tau \gamma \nu)^2) \\
 211 \quad & \leq (\tau \gamma \nu)^2 \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2^2 \middle| H_{i-1} \right] \\
 212 \quad & \leq (\tau \gamma \nu)^{2(i+1)} \mathbb{E} \left[ \|\widehat{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}\|_2^2 \middle| H_0 \right],
 \end{aligned}$$

213 where the shifts in the history index hold because, without feedback, the example selection  
 214 criteria is deterministic given the prior distribution of the learner  $\mathbf{p}_0 = H_0$ . ■

215 **Corollary 3.4.** *Let the learning rate be  $0 < \tau < \frac{2P}{3}$ . Any learner with updates given by (2.1)*  
 216 *converges exponentially with the number of examples when taught by a synthesis based teacher*  
 217 *following the SMTL algorithm.*

218 *Proof.* To guarantee exponential convergence, it is sufficient to show that Theorem 3.3 is  
 219 applicable for a learning rate  $\tau \in (0, \frac{2P}{3})$ , i.e., that Assumption (3.3) holds. The following

220 three inequalities are sufficient conditions for Assumption (3.3) to hold:

$$221 \quad (3.6) \quad \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) > y,$$

$$222 \quad (3.7) \quad -\widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) + \frac{2}{\tau \gamma} > -y,$$

$$223 \quad (3.8) \quad \gamma^2 \widehat{\boldsymbol{\theta}}_i^T (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) - \gamma y - \frac{1}{\tau} \neq 0.$$

224 Recall that  $\max\{\|\boldsymbol{\theta}\|_2^2, \max_i \|\widehat{\boldsymbol{\theta}}_i\|_2^2\} \leq P$ . Selecting  $y' = -1$  and  $0 < \gamma < \min\left\{\frac{1}{P}, \frac{\sqrt{P}}{\|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2}\right\}$   
 225 we show that all requirements (3.6-3.8) hold.

226 We fulfill (3.6) because

$$227 \quad \widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) = \gamma \|\widehat{\boldsymbol{\theta}}_i\|_2 \left( \|\widehat{\boldsymbol{\theta}}_i\|_2 - \|\boldsymbol{\theta}\|_2 \cos(\angle \widehat{\boldsymbol{\theta}}_i, \boldsymbol{\theta}) \right) > -\gamma P > -1 = y,$$

228 where the operator  $\angle \cdot, \cdot$  refers to the angle between two vectors. Next, we note that

$$229 \quad (3.9) \quad -\widehat{\boldsymbol{\theta}}_i^T \gamma (\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}) + \frac{2}{\tau \gamma} > -2\gamma P + \frac{2}{\tau \gamma} > -2 + \frac{2P}{\tau}.$$

230 As we restrict the step-size,  $\tau \in (0, \frac{2P}{3})$ , we may further lower bound (3.9) as

$$231 \quad -2 + \frac{2P}{\tau} > -2 + 3 = 1 = -y,$$

232 so (3.7) is also fulfilled.

233 The left hand side in (3.8) is a non-degenerate quadratic equation with respect to  $\gamma$ , with  
 234 at most two roots. As the interval  $(0, \frac{1}{P})$  is continuous, it must contain non-root values, so  
 235 there must exist a  $\gamma \in \left(0, \min\left\{\frac{1}{P}, \frac{\sqrt{P}}{\|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2}\right\}\right)$  for which (3.8) also holds. Since (3.3) holds, we  
 236 may directly apply Theorem 3.3 to conclude the proof. ■

237 Theorem 3.3 shows that SMTL achieves an exponential behavior similar to omniscient  
 238 teaching. To guarantee the desired exponential convergence of the learner to the ground truth  
 239 with respect to the number of examples, we require  $\mathbb{E}[\|\widehat{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}\|_2^2 \mid H_0] < \infty$ . This requirement  
 240 is a characteristic of most machine learning models, as in general, the starting point of learning  
 241 algorithms is bounded. A sufficient condition for this to hold in our system is  $P < \infty$ . In  
 242 addition, Corollary 3.4 asserts that the assumptions of Theorem 3.3 are fulfilled as long as  
 243 the learning rate is not too large.

244 Following SMTL, a gradient descent learner described by (2.1) needs  $\mathcal{O}(\log \frac{1}{\epsilon} \mathbb{E}[\|\widehat{\boldsymbol{\theta}}_0 - \boldsymbol{\theta}\|_2^2])$   
 245 example-label pairs to learn an  $\epsilon$ -approximation of the ground truth model. This convergence  
 246 rate is of the same order as the one achieved by omniscient teaching, while relaxing the  
 247 assumption about knowledge of the exact learner initialization.

248 The performance guarantees also hold in the case of rescalable pool based teaching with  
 249 a rich enough example set. The approach and its analysis are detailed in Section SM2.  
 250 Additionally, Lemma 3.5 below extends the problem to settings in which the example space  
 251 of the learner suffers an unknown orthonormal transformation with respect to the example  
 252 space of the teacher.

253 **Lemma 3.5.** Let  $\mathcal{G}$  be an unknown orthonormal transformation describing the mapping  
 254 from the feature space of the teacher to the learner. For every example  $\tilde{\mathbf{x}}$  selected by the  
 255 teacher according to SMTL, the learner observes  $\hat{\mathbf{x}} = \mathcal{G}(\tilde{\mathbf{x}})$  and updates its state accord-  
 256 ing to (2.1). If  $\forall \tilde{\boldsymbol{\theta}}_i, \exists \gamma \in \mathbb{R}$  with  $|\gamma| \leq \frac{\sqrt{P}}{\|\tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}\|_2}$ ,  $\nu(\gamma) \in \mathbb{R}$  and  $y' \in \{-1, 1\}$  s.t.  $0 <$   
 257  $\left(\tilde{\boldsymbol{\theta}}_i^T \gamma(\tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}) - y' - \frac{1}{\tau\gamma}\right)^2 \leq \nu^2 < \frac{1}{\tau^2\gamma^2}$  then,

$$258 \quad \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_{i-1} \right] \leq (\tau\gamma\nu)^{2(i+1)} \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_0 - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_0 \right],$$

259 where  $\tilde{\boldsymbol{\theta}}_i = \mathcal{G}^T(\hat{\boldsymbol{\theta}}_i)$  and  $\tilde{\boldsymbol{\theta}}$  represent the learner state and ground truth respectively, in the  
 260 teacher feature space.

261 *Proof.* Let  $\mathcal{G}$  be an orthonormal transformation from the teacher feature space, whose  
 262 elements are identified by  $\tilde{\cdot}$ , to the learner feature space, whose elements are identified by  $\hat{\cdot}$ .  
 263 Let  $\mathcal{G}^T$  denote the inverse mapping from the learner to the teacher feature space. By definition  
 264 of an orthonormal transformation,  $\mathcal{G}$  preserves the inner product, i.e.,  $\langle \hat{\boldsymbol{\theta}}_i, \hat{\mathbf{x}} \rangle = \langle \tilde{\boldsymbol{\theta}}_i, \tilde{\mathbf{x}} \rangle$ . Thus,  
 265 we write the learner updates from iteration  $i$  to  $i + 1$  as

$$266 \quad \hat{\boldsymbol{\theta}}_{i+1} = \hat{\boldsymbol{\theta}}_i - \tau \left( \hat{\boldsymbol{\theta}}_i^T \hat{\mathbf{x}}_i - y_i \right) \hat{\mathbf{x}}_i = \tilde{\boldsymbol{\theta}}_i - \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right) \mathcal{G}(\tilde{\mathbf{x}}_i).$$

267 The error metric is given by the expected squared distance between the ground truth  $\tilde{\boldsymbol{\theta}}$   
 268 and the teacher's estimation about the learner state  $\tilde{\boldsymbol{\theta}}_{i+1}$  in the teacher feature space. As the  
 269 mapping is invertible  $\mathcal{G}^T(\mathcal{G}(\mathbf{x})) = \mathbf{x}$ , we may decompose the MSE as

$$\begin{aligned} 270 \quad \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_{i+1} - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_i \right] &= \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_i - \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right) \mathcal{G}^T \mathcal{G}(\tilde{\mathbf{x}}_i) - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_i \right] \\ 271 &= \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}} - \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right) \tilde{\mathbf{x}}_i \right\|_2^2 \middle| H_i \right] \\ 272 &= \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_i \right] \\ 273 &\quad + \mathbb{E} \left[ -2 \left\langle \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}, \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right) \tilde{\mathbf{x}}_i \right\rangle + \left\| \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right) \tilde{\mathbf{x}}_i \right\|_2^2 \middle| H_i \right] \\ 274 &= \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_i \right] \\ 275 \quad (3.10) &\quad - \tau \mathbb{E} \left[ 2 \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right) \left\langle \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}, \tilde{\mathbf{x}}_i \right\rangle - \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}}_i - y_i \right)^2 \|\tilde{\mathbf{x}}_i\|_2^2 \middle| H_i \right]. \end{aligned}$$

276 The SMTL algorithm selects the example-label pair in the teacher feature space as

$$277 \quad (\tilde{\mathbf{x}}_i, y_i) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}} \mathbb{E} \left[ 2 \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}} - y \right) \left\langle \tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}, \tilde{\mathbf{x}} \right\rangle - \tau \left( \tilde{\boldsymbol{\theta}}_i^T \tilde{\mathbf{x}} - y \right)^2 \|\tilde{\mathbf{x}}\|_2^2 \middle| H_i \right],$$

278 such that the MSE is greedily minimized. This is equivalent to the teacher's behavior when  
 279 the teacher and the learner share the same feature space. Therefore, we apply the inequality  
 280 (3.5) to upper-bound (3.10) as

$$\mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_{i+1} - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_i \right] \leq (\tau\gamma\nu)^{2(i+1)} \mathbb{E} \left[ \left\| \tilde{\boldsymbol{\theta}}_0 - \tilde{\boldsymbol{\theta}} \right\|_2^2 \middle| H_0 \right]. \quad \blacksquare$$

Lemma 3.5 shows that SMTL is invariant to rotations and reflections. Simultaneously teaching and learning provides an exponential speed up even when the learner and teacher do not share a representation space, but there exists an unknown orthonormal transformation between the teacher and learner feature spaces. This result extends the applicability of SMTL to various real-world problems, such as the cross-lingual sentiment analysis discussed in Section 4.2.

### 3.2. Simultaneous Machine Teaching and Learning with noisy Feedback (SMTL-F).

We now analyze the situation in which the teacher receives some feedback from the learner. Without knowledge of the exact learner state, previous approaches [18] propose a dedicated probing phase in which the teacher exploits the feedback to obtain an accurate estimation of the learner state, then allowing the teacher to proceed as if it were omniscient. We show that the teacher may instead simultaneously learn the learner state and teach the ground truth to the learner, thereby, avoiding an explicit probing phase that improves the learner’s estimate without teaching.

For analytical tractability, we consider the case in which the feedback from the learner is given by

$$(3.11) \quad s_i = \hat{\boldsymbol{\theta}}_{i+1}^T \mathbf{x}_i + w_i,$$

where  $w_i \sim \mathcal{N}(0, \sigma^2)$  represents some random noise that accounts for imperfections in the communication channel between learner and teacher. The feedback is a noisy measurement of the learner certainty regarding the latest example classification. Specifically, the learner returns a noisy function of the distance and direction from the latest example to its current classifier. A large positive value of  $s_i$  suggests that the learner probably classifies the latest example  $\mathbf{x}_i$  as class 1. Similarly, a large negative value of  $s_i$  suggests that a classification of  $\mathbf{x}_i$  in class -1 is more probable. On the other hand, a value of  $s_i$  around 0 suggests that the example lies close to the learner classification boundary. Note that recovering the high dimensional true parameter  $\hat{\boldsymbol{\theta}}_{i+1} \in \mathbb{R}^d$  from this noisy scalar  $s_i \in \mathbb{R}$  is not straightforward.

At each time step, the teacher has access to two sources of information about the learner state. First, the teacher directly observes the noisy feedback. Second, the teacher knows the dynamical model of the learner and may predict its future state based on its current estimate. Kalman filtering is a well-known approach to optimally leverage these two sources of information.

The proposed Simultaneous Machine Teaching and Learning algorithm with noisy Feedback (SMTL-F) is summarized in Algorithm 3.2. The teacher interleaves the greedy example selection strategy given by (3.1), with a Kalman filter to achieve optimal tracking. Lines 4, 5 and 6 of Algorithm 3.2 outline the computations required to track mean and covariance of the learner state.

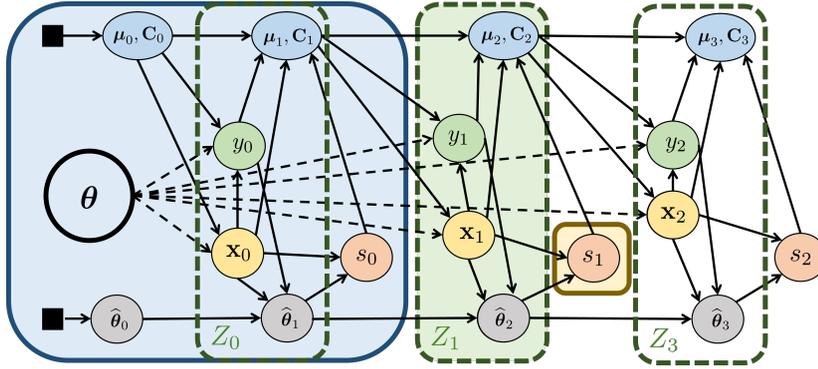


Figure 2: Functional dependence graph showing causal relationships between the teacher estimators about the learner  $\mu, \mathbf{C}$ , the true learner state  $\hat{\theta}$ , the ground truth  $\theta$ , the example-label pairs  $\{\mathbf{x}, y\}$  and the feedback  $s$ . We observe that the system state  $Z_i = \{\mathbf{x}_i, y_i, \mu_{i+1}, \mathbf{C}_{i+1}, \hat{\theta}_{i+1}\}$  is Markovian and that the feedback is conditionally independent of the past given the current state.

318 **Theorem 3.6.** Consider a learner that updates according to (2.1) and provides some feed-  
 319 back according to (3.11). The estimator in SMTL-F then is the **optimal estimator**. Ad-  
 320 ditionally, when  $\tau \leq \frac{2}{P_{\mathbf{x}}}$ , the covariance of the teacher estimation about the learner state is  
 321 monotonically non increasing

$$322 \quad \|\mathbf{C}_{i+1}\|_{\infty} \leq \|\mathbf{C}_i\|_{\infty},$$

323 where  $\|\mathbf{C}\|_{\infty} = \lim_{k \rightarrow \infty} \|\mathbf{C}^k\|^{1/k}$ .

324 *Proof.* To prove Theorem 3.6, we must prove that the learner state estimator in SMTL-F  
 325 is both optimal in the Bayesian sense and that it exhibits stable behavior with monotonically  
 326 non increasing covariance.

### 327 Optimality of the Estimator in SMTL-F

328 We define the system state  $Z_i = \{\mathbf{x}_i, y_i, \mu_{i+1}, \mathbf{C}_{i+1}, \hat{\theta}_{i+1}\}$ . The functional dependence  
 329 graph in Figure 2 shows that the state  $Z_i$  d-separates [5, Definition 2.14] the latest feedback  
 330  $s_i$  from the ground truth, past learner states and past feedback. Therefore, the current  
 331 feedback is conditionally independent of the history given the system state

$$332 \quad \mathbb{P}\left[s_i, \theta, \hat{\theta}_0, \mu_0, \mathbf{C}_0, \{Z_t\}_{t=0}^{i-1} \mid Z_i\right] = \mathbb{P}\left[s_i \mid Z_i\right] \mathbb{P}\left[\theta, \hat{\theta}_0, \mu_0, \mathbf{C}_0, \{Z_t\}_{t=0}^{i-1} \mid Z_i\right]$$

$$333 \quad = \mathbb{P}\left[s_i \mid Z_i\right] \mathbb{P}\left[H_{i-1} \mid Z_i\right].$$

334 Figure 2 also shows that  $Z_i$  d-separates  $Z_{i-1}$  from  $Z_{i+1}$ , therefore, the state is Markovian  
 335  $Z_{i-1} \rightarrow Z_i \rightarrow Z_{i+1}$ . We also observe that any state is independent of past feedback given the  
 336 previous state, so that

$$337 \quad \mathbb{P}\left[Z_{i+1} \mid \{Z_t, s_t\}_{t=0}^i\right] = \mathbb{P}\left[Z_{i+1} \mid Z_i\right].$$

338 Combining the conditional independence with the fact that both learner state and feed-  
 339 back are Gaussian random variables shows that the system follows a Gauss-Markov model.  
 340 Consequently, the Kalman Filter is the Bayesian optimal filter [6]. Moreover, the distribu-  
 341 tions are jointly Gaussian, so we only need to keep track of the mean and covariance matrices  
 342 to obtain the optimal estimator of the learner state. SMTL-F implements the known closed  
 343 form solution of the Kalman Filter for Gauss-Markov models [15]. Hence, SMTL-F obtains  
 344 the optimal posterior probability density function of the learner state in a tractable way.

345 Stability of the Estimator in SMTL-F

346 Next, we show that the estimation of the learner state derived by SMTL-F is stable, in  
 347 the sense that the uncertainty about the learner state is monotonically non increasing. The  
 348 detailed proofs of all auxiliary lemmas are in [Section SM1](#) of the supplemental material. We  
 349 start by deriving the discrete-time algebraic Riccati recursion of the system

350 **Lemma 3.7.** *The dynamic Riccati equation describing the evolution of the teacher’s covari-*  
 351 *ance about the learner state is given by*

$$352 \quad (3.12) \quad \mathbf{C}_{i+1} = \mathbf{F}_i \mathbf{C}_i \mathbf{F}_i^T \mathbf{T}_i,$$

353 where  $\mathbf{F}_i = \mathbf{I} - \mathbf{x}_i \mathbf{x}_i^T$  is the Hermitian state transition matrix at the  $i$ -th iteration and  $\mathbf{T}_i =$   
 354  $\mathbf{I} - (\mathbf{x}_i^T \mathbf{F}_i \mathbf{C}_i \mathbf{F}_i \mathbf{x}_i + \sigma^2)^{-1} \mathbf{x}_i \mathbf{x}_i^T \mathbf{F}_i \mathbf{C}_i \mathbf{F}_i$ , where  $\mathbf{I}$  represents the identity matrix.

355 As a stepping stone towards proving the stability of SMTL-F, we analyze the spectral  
 356 radius of the factors in the Riccati equation (3.12).

357 **Lemma 3.8.** *The spectral radius of  $\mathbf{F}_i$  is 1 for  $\tau \leq \frac{2}{P_x}$ .*

358 **Lemma 3.9.** *The spectral radius of  $\mathbf{T}_i$  is 1.*

359 Lastly, we take the submultiplicative matrix norm  $\|\cdot\|_\infty := \lim_{k \rightarrow \infty} \|\cdot\|^k\|^{1/k}$  on both sides  
 360 of the Riccati recursion (3.12),

$$361 \quad (3.13) \quad \|\mathbf{C}_{i+1}\|_\infty \leq \|\mathbf{C}_i\|_\infty \|\mathbf{F}_i\|_\infty^2 \|\mathbf{T}_i\|_\infty.$$

362 Gelfand’s formula guarantees that  $\rho(\mathbf{A}) = \|\mathbf{A}\|_\infty$  [11], where the operator  $\rho(\cdot)$  represent  
 363 the spectral radius of a matrix. Applying this result together with [Lemma 3.8](#) and [Lemma 3.9](#)  
 364 to (3.13) we obtain

$$365 \quad \|\mathbf{C}_{i+1}\|_\infty \leq \|\mathbf{C}_i\|_\infty \rho(\mathbf{F}_i)^2 \rho(\mathbf{T}_i) \leq \|\mathbf{C}_i\|_\infty,$$

366 which proves that  $\|\mathbf{C}_i\|_\infty$  is monotonically non increasing. ■

367 In the presence of feedback, the estimation of the learner state derived by SMTL-F is both  
 368 optimal (it achieves the smallest expected error) and stable (the uncertainty about the learner  
 369 state is monotonically non increasing).

370 **4. Empirical Performance.** We now analyze the empirical performance of the algorithms  
 371 in a synthetic 2D binary classification problem as well as in a real cross-lingual sentiment  
 372 analysis problem. The code with the algorithms to replicate the experiments is available  
 373 online<sup>1</sup>.

---

<sup>1</sup><https://github.com/BelenMU/SMTL/tree/main>

**Algorithm 3.2** SMTL-F

- 
- 1:  $\boldsymbol{\mu}_0, \mathbf{C}_0 \leftarrow p_0$ .
  - 2: **for**  $i = 0, 1, 2, \dots$  **do**
  - 3:   Select example:  
 $(\mathbf{x}_i, y_i) \leftarrow \underset{\mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}}{\operatorname{argmax}} T(\mathbf{x}, y, \boldsymbol{\mu}_i, \mathbf{C}_i)$
  - 4:   Estimator - Predict:  
 $\boldsymbol{\mu}_{i+1|i} \leftarrow \boldsymbol{\mu}_i - \tau (\boldsymbol{\mu}_i^T \mathbf{x}_i - y_i) \mathbf{x}_i$   
 $\mathbf{C}_{i+1|i} \leftarrow \mathbf{C}_i - \tau \mathbf{C}_i \mathbf{x}_i \mathbf{x}_i^T - \tau \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}_i + \tau^2 \mathbf{x}_i^T \mathbf{C}_i \mathbf{x}_i \mathbf{x}_i^T \mathbf{C}_i$
  - 5:   Estimator<sub>T</sub> - Observe feedback:  
 $s_i \leftarrow \widehat{\boldsymbol{\theta}}_{i+1}^T \mathbf{x}_i + w_i$
  - 6:   Estimator - Update estimation:  
 $\mathbf{K}_{i+1} \leftarrow \mathbf{C}_{i+1|i} \mathbf{x}_i (\mathbf{x}_i^T \mathbf{C}_{i+1|i} \mathbf{x}_i + \sigma^2)^{-1}$   
 $\boldsymbol{\mu}_{i+1} \leftarrow \boldsymbol{\mu}_{i+1|i} + \mathbf{K}_{i+1} (s_i - \boldsymbol{\mu}_{i+1|i}^T \mathbf{x}_i)$   
 $\mathbf{C}_{i+1} \leftarrow (\mathbf{I} - \mathbf{K}_{i+1} \mathbf{x}_i^T) \mathbf{C}_{i+1|i} (\mathbf{I} - \mathbf{K}_{i+1} \mathbf{x}_i^T)^T + \sigma^2 \mathbf{K}_{i+1} \mathbf{K}_{i+1}^T$
  - 7: **end for**
- 

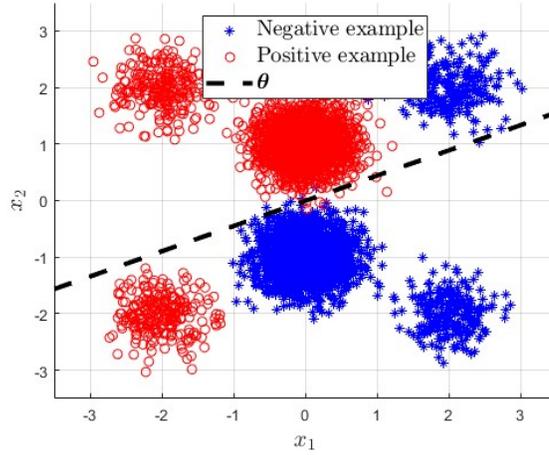


Figure 3: Synthetic dataset synth2 [32].

374 **4.1. Synthetic Dataset.** We first compare the performance of the SMTL and SMTL-F  
 375 algorithms against the state of the art online machine teaching methods with a synthetic  
 376 dataset. We generate a standard 2D binary dataset, shown in Figure 3, following the procedure  
 377 outlined in [32].

378 We validate the proposed online algorithms against the baseline omniscient teaching algo-  
 379 rithm. Figure 4a shows the evolution of the learner error  $\|\widehat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}\|_2$  as more examples are  
 380 presented. We observe that the error decreases exponentially fast for both online algorithms  
 381 as well as for the omniscient teacher, hence offering a significant improvement compared to the  
 382 rate of traditional Stochastic Gradient Descent (SGD) in which examples are chosen randomly.

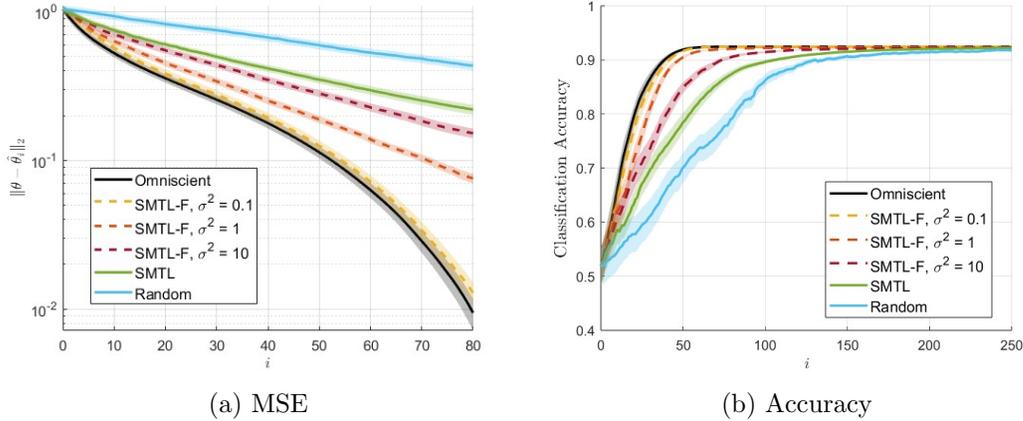


Figure 4: Performance comparison between algorithms on the synthetic dataset. All online MT algorithms achieve an exponential speed-up w.r.t. randomly selecting examples. Within the exponential convergence of the MSE, the lower the noise level of the feedback the faster the MSE decreases and the classification accuracy increases.

383 However, within the exponential rates, the omniscient teacher performs the best because it  
 384 has the most information about the learner state.

385 In the presence of feedback, tracking the learner is a good strategy to bridge the gap  
 386 in performance between the omniscient teacher and the no-feedback case. The MSE of the  
 387 SMTL-F is lower bounded by the MSE of omniscient teaching and upper bounded by the MSE  
 388 of SMTL. As the feedback noise level decreases, SMTL-F approaches the omniscient teacher  
 389 performance. In fact, as Figure 4a shows, under feedback with very low noise levels, SMTL-F  
 390 rapidly achieves a precise estimation of the learner state, becoming a *de facto* omniscient  
 391 teacher.

392 Although we use the squared distance between the learner and the ground truth as a per-  
 393 formance metric, the ultimate objective is to achieve a good classification accuracy. Figure 5  
 394 shows how these two metrics are intertwined: a learner close to the ground truth, i.e., a low  
 395  $\|\hat{\theta}_i - \theta\|_2^2$ , implies a good classification accuracy. The same relationship holds for different  
 396 datasets, as analyzed in Section SM3.2. This justifies a posteriori why the proposed online  
 397 algorithms focus on non increasing  $\|\hat{\theta}_i - \theta\|_2^2$ , as this is a good heuristic for classification accu-  
 398 racy improvement. The relationship between both metrics is highly non-linear, meaning that  
 399 an improvement on the learner state can strongly improve the classification accuracy when the  
 400 state is far from the ground truth. Once the learner is sufficiently close to the ground truth,  
 401 fine-tuning the learner’s state yields a much less significant change in classification accuracy.  
 402 This behavior highlights the benefits of SMTL-F: for sufficiently low noise levels on the feed-  
 403 back, teachers following SMTL-F are able to keep up with the omniscient teacher until a high  
 404 enough accuracy is reached, at which point fine-tuning of the learner state no longer has a  
 405 significant impact on classification accuracy.

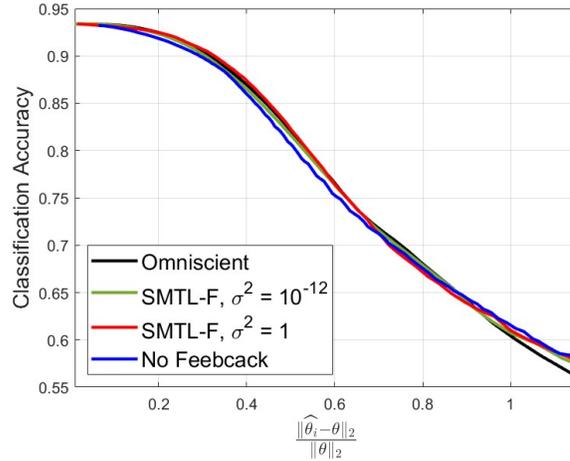


Figure 5: Correspondence between classification accuracy and the learner’s distance to the ground truth for different algorithms.

406 The graphs in [Figure 5](#) and [Figure SM3](#) show that all algorithms exhibit similar relationships  
 407 between MSE and classification accuracy. This behavior suggests that there are implicit  
 408 *trajectories* that all the online machine teaching algorithms approximately follow, and that the  
 409 speed at which learners travel along the *trajectories*, measured in terms of number examples,  
 410 strongly depends on the teacher’s knowledge about the learner. Said differently, the feedback  
 411 provided by the learner does not seem to provide advantages in terms of trajectory, it only  
 412 seems to affect how fast the learner reaches a low  $\|\hat{\theta}_i - \theta\|_2^2$  value.

413 [Figure 4b](#) summarizes the performance of the online algorithms, as measured by the clas-  
 414 sification accuracy. Machine teaching outperforms random example selection. With more  
 415 information about the learner, the classification accuracy of the learner improves faster with  
 416 respect to the number of examples.

417 We explore how learner initializations impact algorithm performance. We randomly ini-  
 418 tialize 50 learners and compare the resulting variation in performance. The shaded regions in  
 419 [Figure 4](#) represent the standard error between initializations. Notably, online machine teach-  
 420 ing not only outperforms random example selection but also enhances robustness as SMTL  
 421 and SMTL-F exhibit significantly lower variance. This finding suggests that the proposed  
 422 algorithms offer more consistent and stable results under different starting conditions, making  
 423 them a favorable choice for various applications. Online machine teaching mitigates the im-  
 424 pact of learner initializations on performance, and this effect is further diminished as feedback  
 425 noise decreases.

426 We further validate SMTL-F against the Learning for Omniscience (LfO) algorithm [[18](#),  
 427 [16](#)]. There are two distinct phases of the LfO algorithm corresponding to the probing and  
 428 teaching phases. At first, the teacher focuses solely on decreasing its uncertainty about the  
 429 learner state, i.e.,  $(\hat{\mathbf{x}}, \hat{y}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}} \|\mathbf{C}_i\|_2$ . As [Figure 6a](#) shows, at first teachers fol-  
 430 lowing LfO reduce their uncertainty about the learner much faster than SMTL and SMTL-F.

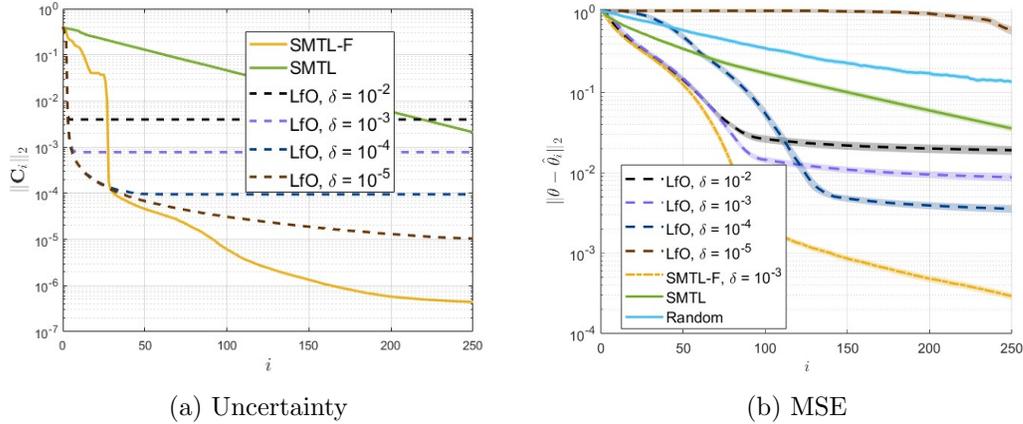


Figure 6: Performance of the proposed online machine teaching algorithms against the state of the art. SMTL-F outperforms LfO with  $\sigma^2 = 10^{-3}$  by continuously updating its estimation about the learner while teaching, avoiding an explicit probing phase.

431 However, getting an accurate estimation of the learner is done at the expense of teaching the  
 432 ground truth. As Figure 6b shows, the MSE of the learner remains constant during the first  
 433 iterations as the learner does not update its state immediately [18]. The teaching phase of  
 434 the LfO algorithm starts once the uncertainty about the learner state is sufficiently low, i.e.  
 435  $\|C_i\|_2 < \delta$ , for a given threshold  $\delta \in \mathbb{R}^+$ . Then, the teacher proceeds as if it were omniscient  
 436 using its latest estimation.

437 Figure 6b shows the performance of SMTL-F against LfO when  $\sigma^2 = 10^{-3}$ . We observe  
 438 that having separate learning and teaching phases negatively impacts the overall performance  
 439 of the algorithm. If the probing phase is too short, the teacher does not have an accurate  
 440 estimation of the learner, so it is not able to teach it efficiently and the error decreases much  
 441 slower than with SMTL-F, which continuously improves its estimation of the learner. On the  
 442 other hand, a longer probing phase leads to an accurate estimation of the learner state but  
 443 requires many iterations without teaching in which the error does not decrease. In practice,  
 444 LfO with a long probing phase, i.e., low  $\delta$ , is unable to catch up with the online algorithm that  
 445 has been teaching all along. The proposed algorithm with noisy feedback avoids the costly  
 446 probing phase, while still obtaining an accurate and ever-improving estimation of the learner  
 447 state.

448 These experiments confirm that jointly teaching the learner while estimating its param-  
 449 eters offers significant gains.

450 **4.2. Cross-lingual Sentiment Analysis.** Language can be harnessed to understand the  
 451 attitude of individuals [25]. Towards this goal, binary sentiment word classification aims to  
 452 accurately label words according to their connotation as positive (e.g., love) or negative (e.g.,  
 453 death). Traditionally, research on lingual sentiment analysis has focused on a few languages  
 454 that have a large amount of annotated data [9]. To tackle this resource imbalance, cross-lingual

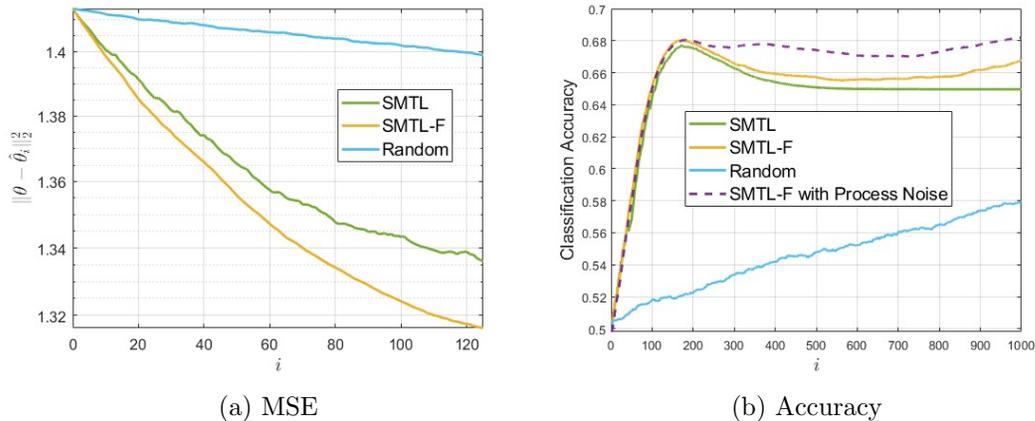


Figure 7: Performance on the cross-lingual sentiment analysis problem. Online machine teaching algorithms speed up the teaching. Adding process noise in the Kalman update reduces the drop in performance caused by non-orthogonalities in the mapping between Spanish and Italian words.

455 adaptation [1, 14, 28] aims to transfer the knowledge of languages with plentiful resources to  
 456 languages with few resources. In this section, we apply SMTL and SMTL-F to tackle the  
 457 cross-lingual sentiment analysis problem. We assume that the teacher has access to a linear  
 458 sentiment classifier in the word-space created from a Spanish dictionary. The teacher aims to  
 459 teach a learner working on the word-space created from an Italian dictionary to accurately  
 460 classify Italian words.

461 We use existing monolingual word embeddings<sup>2</sup> [2] and normalize each word vector. Pre-  
 462 vious work [31] empirically shows that the mapping of normalized word vectors between lan-  
 463 guages is accurately described by an orthonormal transformation. Hence following Lemma 3.5,  
 464 SMTL is suitable for cross-lingual knowledge transfer, even if the explicit mapping between  
 465 the Spanish and Italian word embeddings is unknown.

466 The teacher works in the Spanish word-embedding. At each iteration, the teacher selects  
 467 the example-label pair according to (3.1) where  $\mathcal{X}$  is the set of embedded Spanish words.  
 468 We limit the examples to a finite dataset by selecting the 10000 most common words. This  
 469 extension of synthesis-based teaching to a pool-based setting is detailed in Section SM2.1.  
 470 We use Google Translate<sup>3</sup> to translate each example from Spanish to Italian. The learner  
 471 only sees the embedding corresponding to the translated word in the Italian vector space.  
 472 Figure 7a shows the evolution of the MSE when a teacher working in the Spanish word space  
 473 teaches a word sentiment classifier to a learner in the Italian word space. Machine teaching  
 474 decreases the error significantly faster than random example selection.

475 The performance further improves when the learner provides feedback about its state to

<sup>2</sup><http://ixa2.si.ehu.es/martetxe/vecmap/es.emb.txt.gz>

<sup>3</sup><https://translate.google.com>

476 the teacher. As orthonormal transformations preserve inner-products, we follow the framework  
 477 described in Section 3.2. The feedback from the learner to the teacher is described as

$$478 \quad s_i = \widehat{\boldsymbol{\theta}}_{i+1}^T \mathcal{G}(\mathbf{x}_i) = \mathcal{G}^{-1}(\widehat{\boldsymbol{\theta}}_{i+1})^T \mathbf{x}_i + w_i,$$

479 where  $\mathcal{G}$  is the unknown orthonormal mapping of word embeddings from the teacher to the  
 480 learner language space. As the real mapping is not exactly an orthonormal transformation,  
 481 we introduce  $w_i \sim \mathcal{N}(0, \sigma^2)$  to account for the deviations from the perfect orthonormality  
 482 assumption.

483 We estimate the noise level  $\sigma^2$  from the information exchanged between teacher and  
 484 learner. The teacher samples  $N$  random pairs of words  $(\tilde{\mathbf{x}}_a, \tilde{\mathbf{x}}_b)$ , the learner observes the  
 485 corresponding word pairs in the learner word space  $(\widehat{\mathbf{x}}_a, \widehat{\mathbf{x}}_b)$ , computes each pair’s inner prod-  
 486 uct and transmits the resulting products to the teacher. The teacher then calculates the  
 487 differences in inner-products between the pairs of words in the teacher language and the  
 488 learner language. The variance among these differences becomes the estimator for  $\sigma^2$ ,

$$489 \quad \sigma^2 \approx \frac{1}{N} \sum_{n=1}^N (\tilde{\mathbf{x}}_{a,n}^T \tilde{\mathbf{x}}_{b,n} - \widehat{\mathbf{x}}_{a,n}^T \widehat{\mathbf{x}}_{b,n})^2,$$

490 where  $(\tilde{\mathbf{x}}_{a,n}, \tilde{\mathbf{x}}_{b,n})$  is the  $n$ -th pair of words sampled by the teacher. As Figure 7a shows,  
 491 incorporating learner feedback with this estimator further improves the rate at which the  
 492 MSE decreases.

493 We also test the learner accuracy for classifying a preexisting sentiment lexicon in Italian<sup>4</sup>.  
 494 The results are shown in Figure 7b. Online machine teaching algorithms are superior to  
 495 random selection of examples. In fact, 50 examples selected by SMTL or SMTL-F achieve the  
 496 same classification accuracy as 1000 randomly selected examples.

497 **4.2.1. Deviations from Orthogonal Mappings.** As the mapping between languages is  
 498 not perfectly orthonormal, the teacher model of the learner dynamical system is slightly  
 499 inaccurate. This could lead to instances in which the teacher is certain of its learner state  
 500 estimation, but this estimation is inaccurate. This would explain the dip in accuracy observed  
 501 in Figure 7b. In this section, we further analyze this conjecture; i.e., we investigate how  
 502 deviations from the orthonormality assumption in Lemma 3.5 affect the performance of SMTL-  
 503 F. We also propose an extension of the algorithm to account for the deviations, and diminish  
 504 the performance dips they cause.

505 To empirically understand how SMTL-F performs under non-orthogonal transformations,  
 506 we modify the synthetic experiments in Section 4.1. We create a new learner example space  
 507 by rotating each example

$$508 \quad \widehat{\mathbf{x}}_i = \text{Rotate}(\tilde{\mathbf{x}}_i, \phi + z_i),$$

509 where the degrees of rotation  $\phi + z_i$  are composed of a deterministic amount, unknown to the  
 510 teacher, along with an additional random rotation. The deterministic rotation, denoted by  $\phi$ ,  
 511 is sampled from a uniform distribution  $\phi \sim \mathcal{U}(0, 2\pi)$  and remains constant for all examples.

---

<sup>4</sup><https://www.kaggle.com/datasets/rtatman/sentiment-lexicons-for-81-languages>

512 On the other hand, the random rotation  $z_i$  is sampled independently for each example from  
 513 a Gaussian distribution  $z_i \sim \mathcal{N}(0, z^2)$  which adds an extra random degree of rotation to each  
 514 instance.

515 **Figure 8** shows that as the examples deviate further from the perfect orthonormal trans-  
 516 formation, a dip in accuracy appears. This behavior gives credence to our conjecture that the  
 517 drop in performance in **Figure 7b** is caused by deviations from the assumption of orthogonal  
 518 mapping between languages.

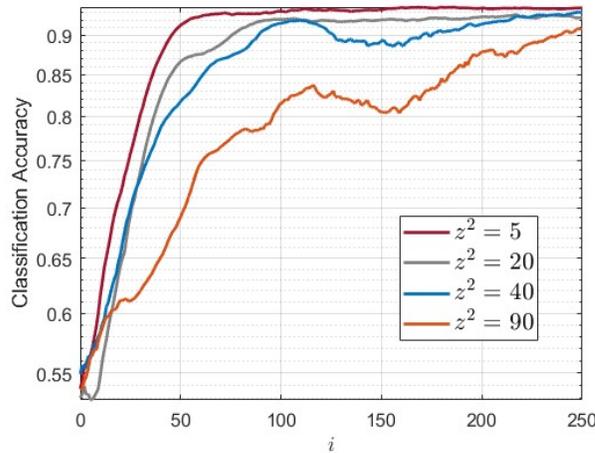


Figure 8: Performance of SMTL-F for examples deviated by  $\mathcal{N}(0, z^2)$  from perfect orthonormal mapping between teacher and learner feature spaces. Deviations lead to a performance dip.

519 The SMTL-F algorithm assumes a perfect knowledge of the dynamical system of the  
 520 learner. However, the examples from the teacher to the learner space do not always experience  
 521 the same rotation so, in practice, the teacher may not be able to exactly determine the  
 522 evolution of the learner state. The teacher overcomes the estimation error when observing  
 523 more feedback from the learner, which is consistent with previous works [20, 13] showing that  
 524 interactivity mitigates the impact of imperfect knowledge and mismatches.

525 Another approach is to account for the mapping imperfections by introducing process  
 526 noise in the dynamical model of the learner. Let  $\mathbf{r}_i$  denote the difference between the teacher’s  
 527 example mapped in a perfectly orthogonal way  $\mathcal{G}(\tilde{\mathbf{x}})$  and the corresponding example in the  
 528 learner space  $\hat{\mathbf{x}}_i$ ; i.e.,  $\mathbf{r}_i = \hat{\mathbf{x}}_i - \mathcal{G}(\tilde{\mathbf{x}}_i)$ . Then, the evolution of the learner state from iteration  $i$

529 to  $i + 1$  is given by

$$\begin{aligned}
530 \quad \hat{\boldsymbol{\theta}}_{i+1} &= \hat{\boldsymbol{\theta}}_i - \tau \left( \hat{\boldsymbol{\theta}}_i^T \hat{\mathbf{x}}_i - y_i \right) \hat{\mathbf{x}}_i = (\mathbf{I} - \tau \hat{\mathbf{x}}_i \hat{\mathbf{x}}_i^T) \hat{\boldsymbol{\theta}}_i + \tau y \hat{\mathbf{x}}_i \\
531 \quad &= (\mathbf{I} - \tau (\mathcal{G}(\tilde{\mathbf{x}}_i) + \mathbf{r}_i) (\mathcal{G}(\tilde{\mathbf{x}}_i) + \mathbf{r}_i)^T) \hat{\boldsymbol{\theta}}_i + \tau y_i (\mathcal{G}(\tilde{\mathbf{x}}_i) + \mathbf{r}_i) \\
532 \quad &= (\mathbf{I} - \tau \mathcal{G}(\tilde{\mathbf{x}}_i) \mathcal{G}(\tilde{\mathbf{x}}_i)^T - \tau \mathcal{G}(\tilde{\mathbf{x}}_i) \mathbf{r}_i^T - \tau \mathbf{r}_i \mathcal{G}(\tilde{\mathbf{x}}_i)^T - \tau \mathbf{r}_i \mathbf{r}_i^T) \hat{\boldsymbol{\theta}}_i + \tau y_i (\mathcal{G}(\tilde{\mathbf{x}}_i) + \mathbf{r}_i) \\
533 \quad &= (\mathbf{I} - \tau \mathcal{G}(\tilde{\mathbf{x}}_i) \mathcal{G}(\tilde{\mathbf{x}}_i)^T) \hat{\boldsymbol{\theta}}_i + \tau y_i \mathcal{G}(\tilde{\mathbf{x}}_i) - \tau \underbrace{(\mathcal{G}(\tilde{\mathbf{x}}_i) \mathbf{r}_i^T + \mathbf{r}_i \mathcal{G}(\tilde{\mathbf{x}}_i)^T + \mathbf{r}_i \mathbf{r}_i^T)}_{\mathbf{v}_i} \hat{\boldsymbol{\theta}}_i + \tau y_i \mathbf{r}_i \\
534 \quad &= \hat{\boldsymbol{\theta}}_i - \tau \left( \hat{\boldsymbol{\theta}}_i^T \mathcal{G}(\tilde{\mathbf{x}}_i) - y_i \right) \mathcal{G}(\tilde{\mathbf{x}}_i) - \tau \underbrace{(\mathcal{G}(\tilde{\mathbf{x}}_i) \mathbf{r}_i^T + \mathbf{r}_i \mathcal{G}(\tilde{\mathbf{x}}_i)^T + \mathbf{r}_i \mathbf{r}_i^T)}_{\mathbf{v}_i} \hat{\boldsymbol{\theta}}_i + \tau y_i \mathbf{r}_i.
\end{aligned}$$

535 From a control perspective, the deviations from perfect orthogonal mappings create unknowns  
536 in the dynamical system, these unknowns  $\mathbf{v}_i$  are random variables referred as process noise.

537 Dealing with process noise is a known and well investigated problem in control theory  
538 [27, Chapter 7]. We leave the best modeling of this process noise for future work. For now,  
539 we model the deviations from orthogonality in a naive way by assuming that the noise is  
540 independent and identically distributed (i.i.d.) Gaussian, namely,  $\mathbf{v}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{v}})$ . Under  
541 this assumption, the covariance extrapolation for the Kalman update becomes

$$542 \quad \mathbf{C}_{i+1|i} = (\mathbf{I} - \tau \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T) \mathbf{C}_{i|i} (\mathbf{I} - \tau \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T)^T + \boldsymbol{\Sigma}_{\mathbf{v}}.$$

543 Despite the simplicity of the process noise model, we observe a significant improvement in  
544 performance. The dashed purple line in Figure 7b shows that assuming Gaussian process  
545 noise smooths the performance curve. We diminish the drop in performance in the cross-  
546 lingual experiment by accounting for the deviations from the orthogonal mapping between  
547 Italian and Spanish words with Gaussian process noise.

548

## REFERENCES

- 549 [1] M. ABDALLA AND G. HIRST, *Cross-Lingual Sentiment Analysis Without (Good) Translation*, in Proc. of  
550 International Joint Conference on Natural Language Processing., 7 2017.
- 551 [2] M. ARTETXE, G. LABAKA, AND E. AGIRRE, *Learning principled bilingual mappings of word embeddings*  
552 *while preserving monolingual invariance*, in Proc. of Conference on Empirical Methods in Natural  
553 Language Processing, Association for Computational Linguistics, 2016, pp. 2289–2294, [https://doi.](https://doi.org/10.18653/v1/d16-1250)  
554 [org/10.18653/v1/d16-1250](https://doi.org/10.18653/v1/d16-1250).
- 555 [3] F. J. BALBACH AND T. ZEUGMANN, *Recent Developments in Algorithmic Teaching*, in Proc. of Interna-  
556 tional Conference on Language and Automata Theory and Applications, vol. 5457, Heidelberg, 2009,  
557 Springer, pp. 1–18, [https://doi.org/10.1007/978-3-642-00982-2\\_{-}1](https://doi.org/10.1007/978-3-642-00982-2_{-}1).
- 558 [4] M. BAYER AND C. REUTER, *ActiveLLM: Large Language Model-based Active Learning for Textual Few-*  
559 *Shot Scenarios*, (2024), <http://arxiv.org/abs/2405.10808>.
- 560 [5] M. BLOCH AND J. BARROS, *Physical-Layer Security: From Information Theory to Security Engineering*,  
561 Cambridge University Press, 2011.
- 562 [6] J. B. M. BRIAN D. O. ANDERSON, *Optimal Filtering*, 1979.
- 563 [7] G. CANAL, S. FENU, AND C. ROZELL, *Active ordinal querying for tuplewise similarity learning*, AAAI  
564 2020 - 34th AAAI Conference on Artificial Intelligence, (2020), pp. 3332–3340, [https://doi.org/10.](https://doi.org/10.1609/aaai.v34i04.5734)  
565 [1609/aaai.v34i04.5734](https://doi.org/10.1609/aaai.v34i04.5734).
- 566 [8] S. DASGUPTA, D. HSU, S. POULIS, AND X. ZHU, *Teaching a black-box learner*, in Proc. of International  
567 Conference on Machine Learning, PMLR, 2019, pp. 1547–1555.
- 568 [9] K. DASHTIPOUR, S. PORIA, A. HUSSAIN, E. CAMBRIA, A. Y. HAWALAH, A. GELBUKH, AND Q. ZHOU,  
569 *Multilingual Sentiment Analysis: State of the Art and Independent Comparison of Techniques*, Cog-  
570 nitive Computation, 8 (2016), pp. 757–771, <https://doi.org/10.1007/S12559-016-9415-7/TABLES/2>,  
571 <https://link.springer.com/article/10.1007/s12559-016-9415-7>.
- 572 [10] S. FREITAS, E. LABER, P. LAZERA, AND M. MOLINARO, *Time-constrained learning*, Pattern Recognition,  
573 142 (2023), p. 109672, <https://doi.org/10.1016/j.patcog.2023.109672>.
- 574 [11] I. GELFAND, *Normierte Ringe*, Matematischeskii Sbornik, 9 (1941), pp. 3–24.
- 575 [12] I. J. GOODFELLOW, J. SHLENS, AND C. SZEGEDY, *Explaining and harnessing adversarial examples*, in  
576 Proc. of International Conference on Learning Representations, 2015, pp. 1–11.
- 577 [13] C. GUERRA, F. S. MELO, AND M. LOPES, *FIT: Using Feature Importance to Teach Classification Tasks to*  
578 *Unknown Learners*, in Proc. of Progress in Artificial Intelligence, Lisbon, 2022, Springer, pp. 440–451,  
579 [https://doi.org/10.1007/978-3-031-16474-3\\_{-}36](https://doi.org/10.1007/978-3-031-16474-3_{-}36).
- 580 [14] X. HE, H. ZHANG, W. CHAO, AND D. WANG, *Semi-supervised learning on cross-lingual sentiment*  
581 *analysis with space transfer*, in Proc. of International Conference on Big Data Computing Ser-  
582 vice and Applications, Institute of Electrical and Electronics Engineers, 8 2015, pp. 371–377,  
583 <https://doi.org/10.1109/BigDataService.2015.57>.
- 584 [15] R. E. KALMAN, *A New Approach to Linear Filtering and Prediction Problems*, Journal of Basic Engi-  
585 neering, 82 (1960), pp. 35–45, <https://doi.org/10.1115/1.3662552>.
- 586 [16] P. KAMALARUBAN, R. DEVIDZE, V. CEVHER, AND A. SINGLA, *Interactive Teaching Algorithms for*  
587 *Inverse Reinforcement Learning*, in Proc. of International Joint Conference on Artificial Intelligence,  
588 Macao, China, 5 2019, AAAI Press, pp. 2692–2700.
- 589 [17] W. LIU, B. DAI, A. HUMAYUN, C. TAY, C. YU, L. B. SMITH, J. M. REHG, AND L. SONG, *Iterative*  
590 *machine teaching*, in Proc. of International Conference on Machine Learning, vol. 5, JMLR, 2017,  
591 pp. 3390–3412.
- 592 [18] W. LIU, B. DAI, X. LI, Z. LIU, J. M. REHG, AND L. SONG, *Towards black-box iterative machine teaching*,  
593 in Proc. of International Conference on Machine Learning, J. Dy and A. Krause, eds., PMLR, 2018,  
594 pp. 3141–3149.
- 595 [19] Y. LU, T. JAVIDI, AND S. LAZEBNIK, *Adaptive Object Detection Using Adjacency and Zoom Prediction*, in  
596 Proc. of Conference on Computer Vision and Pattern Recognition, vol. 2016-Decem, 2016, pp. 2351–  
597 2359, <https://doi.org/10.1109/CVPR.2016.258>.
- 598 [20] F. S. MELO, C. GUERRA, AND M. LOPES, *Interactive Optimal Teaching with Unknown Learners*, in  
599 Proc. of International Joint Conference on Artificial Intelligence, California, 7 2018, pp. 2567–2573,  
600 <https://doi.org/10.24963/ijcai.2018/356>.

- 601 [21] P. NARAYANAMURTHY AND U. MITRA, *Uncertainty-Based Non-Parametric Active Peak Detection*, in  
602 Proc. of International Symposium on Information Theory, 2022, pp. 1–10, [http://arxiv.org/abs/  
603 2205.02376](http://arxiv.org/abs/2205.02376).
- 604 [22] A. NEMIROVSKI, A. JUDITSKY, G. LAN, AND A. SHAPIRO, *Robust Stochastic Approximation Approach  
605 to Stochastic Programming*, SIAM Journal on Optimization, 19 (2009), pp. 1574–1609, [https://doi.  
606 org/10.1137/070704277](https://doi.org/10.1137/070704277).
- 607 [23] S. J. PAN AND Q. YANG, *A survey on transfer learning*, Transactions on Knowledge and Data Engineering,  
608 22 (2010), pp. 1345–1359, <https://doi.org/10.1109/TKDE.2009.191>.
- 609 [24] Z. QIU, W. LIU, T. Z. XIAO, Z. LIU, U. BHATT, Y. LUO, A. WELLER, AND B. SCHÖLKOPF, *Iterative  
610 Teaching by Data Hallucination*, in Proc. of International Conference on Artificial Intelligence and  
611 Statistics, F. Ruiz, J. Dy, and J.-W. van de Meent, eds., PMLR, 4 2023, pp. 9892–9913.
- 612 [25] M. D. ROCKLAGE AND R. H. FAZIO, *The Evaluative Lexicon: Adjective use as a means of assessing and  
613 distinguishing attitude valence, extremity, and emotionality*, Journal of Experimental Social Psychol-  
614 ogy, 56 (2015), pp. 214–227, <https://doi.org/10.1016/j.jesp.2014.10.005>.
- 615 [26] B. SETTLES, *Active Learning*, Synthesis Lectures on Artificial Intelligence and Machine Learning, 6 (2012),  
616 pp. 1–114, <https://doi.org/10.2200/S00429ED1V01Y201207AIM018>.
- 617 [27] D. SIMON, *Kalman filter generalizations*, in Optimal State Estimation: Kalman, H-Infinity, and Nonlinear  
618 Approaches, John Wiley & Sons, Inc., Hoboken, NJ, USA, 5 2006, ch. 7, pp. 183–227, [https://doi.  
619 org/10.1002/0470045345](https://doi.org/10.1002/0470045345).
- 620 [28] Q. SUN, M. S. AMIN, B. YAN, C. MARTELL, V. MARKMAN, A. BHASIN, AND J. YE, *Transfer learning  
621 for bilingual content classification*, in Proc. of Conference on Knowledge Discovery and Data Mining,  
622 vol. 2015-Augus, New York, NY, USA, 2015, ACM, pp. 2147–2156, [https://doi.org/10.1145/2783258.  
623 2788575](https://doi.org/10.1145/2783258.2788575).
- 624 [29] O. TAMUZ, C. LIU, S. BELONGIE, O. SHAMIR, AND A. T. KALAI, *Adaptively learning the crowd kernel*,  
625 Proc. of International Conference on Machine Learning, (2011), pp. 673–680.
- 626 [30] J. WHITEHILL AND J. MOVELLAN, *Approximately Optimal Teaching of Approximately Optimal Learners*,  
627 IEEE Transactions on Learning Technologies, 11 (2018), pp. 152–164, [https://doi.org/10.1109/TLT.  
628 2017.2692761](https://doi.org/10.1109/TLT.2017.2692761).
- 629 [31] C. XING, D. WANG, C. LIU, AND Y. LIN, *Normalized word embedding and orthogonal transform for  
630 bilingual word translation*, in Proc. of Conference of Human Language Technologies, 2015, pp. 1006–  
631 1011, <https://doi.org/10.3115/v1/n15-1104>.
- 632 [32] Y. YANG AND M. LOOG, *A benchmark and comparison of active learning for logistic regression*, Pattern  
633 Recognition, 83 (2018), pp. 401–415, <https://doi.org/10.1016/J.PATCOG.2018.06.004>.
- 634 [33] Z. ZHANG, E. STRUBELL, AND E. HOVY, *A Survey of Active Learning for Natural Language Processing*,  
635 Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP  
636 2022, (2022), pp. 6166–6190.
- 637 [34] X. ZHU, *Machine Teaching: An Inverse Problem to Machine Learning and an Approach Toward Optimal  
638 Education*, Proc. of International Joint Conference on Artificial Intelligence, 29 (2015), [https://doi.  
639 org/10.1609/aaai.v29i1.9761](https://doi.org/10.1609/aaai.v29i1.9761).
- 640 [35] X. ZHU, A. SINGLA, S. ZILLES, AND A. N. RAFFERTY, *An Overview of Machine Teaching*, arXiv, (2018),  
641 pp. 1–18, <http://arxiv.org/abs/1801.05927>.