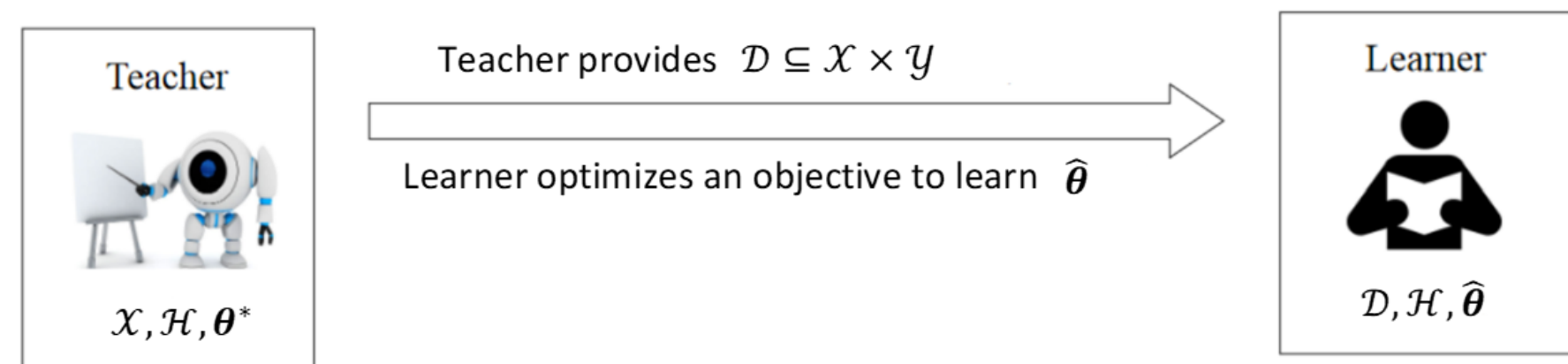


# THE TEACHING DIMENSION OF KERNEL PERCEPTRON

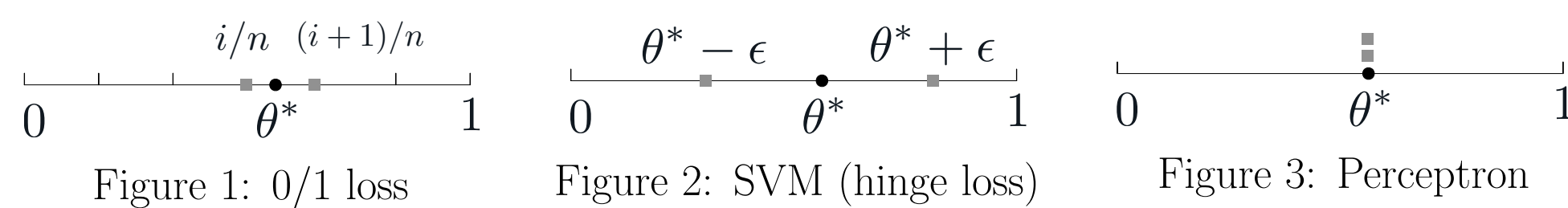
Akash Kumar\* Hanqi Zhang† Adish Singla\* Yuxin Chen†

\*Max Planck Institute for Software Systems (MPI-SWS) †University of Chicago

## Algorithmic Teaching: ERM Learner



### Teaching 1D threshold function to an ERM learner



|    | version-space                              | SVM  | perceptron                          |
|----|--|--|-------------------------------------|
| TD | 2 (discrete)                               | 2 (continuous)   | 2 (continuous)                      |
| TS | $\{(\frac{i}{n}, 0), (\frac{i+1}{n}, 1)\}$ | $\{(\theta^* - \epsilon, 0), (\theta^* + \epsilon, 1)\}$ | $\{(\theta^*, -1), (\theta^*, 1)\}$ |

### Constructive Setting

Teacher could provide *arbitrarily constructed* training set (teaching examples)  $\mathcal{D}$  in the support of the data distribution  $\mathcal{P}$ .

## Kernel Perceptron Learner

Fix the training set  $\mathcal{D} := \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where  $\mathbf{x}_i \in \mathbb{R}^d$  and hypothesis  $\theta \in \mathbb{R}^d$ .

### Linear Perceptron

A homogeneous linear perceptron corresponding to hypothesis  $\theta$  is defined as:

$$f_\theta(\mathbf{x}) := \text{sign}(\theta \cdot \mathbf{x})$$

### Perceptron Loss Function

For a labelled point  $(\mathbf{x}, y)$ , hypothesis  $\theta$  we consider the perceptron loss  $\ell$ :

$$\ell(f_\theta(\mathbf{x}), y) := \max(-y \cdot f_\theta(\mathbf{x}), 0)$$

### Optimal Perceptron Algorithm (Learner)

For a given training set  $\mathcal{D}$ , we consider the optimal perceptron algorithm  $\mathcal{A}_{opt}$  which minimizes the perceptron loss as follows:

$$\mathcal{A}_{opt}(\mathcal{D}) := \arg \min_{\theta \in \mathbb{R}^d} \sum_{i=1}^n \ell(f_\theta(\mathbf{x}_i), y_i)$$

### Non-Linear Kernel Perceptron Learner

For a kernel operator  $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  inducing a *reproducing kernel Hilbert space* (RKHS)  $\mathcal{H}_{\mathcal{K}}$ , a non-linear kernel perceptron optimizes the following loss:

$$\mathcal{A}_{opt}(\mathcal{D}) := \arg \min_{\theta \in \mathcal{H}_{\mathcal{K}}} \sum_{i=1}^n \ell(f_\theta(\mathbf{x}_i), y_i)$$

where  $f_\theta(\cdot) = \sum_{i=1}^l \alpha_i \cdot \mathcal{K}(\mathbf{a}_i, \cdot)$  for some  $\{\mathbf{a}_i\}_{i=1}^l \subset \mathcal{X}$  and  $\alpha_i$  real. We also write  $f_\theta(\cdot) = \theta \cdot \Phi(\cdot)$  where  $\Phi : \mathcal{X} \rightarrow \mathcal{H}_{\mathcal{K}}$  is defined as feature map to  $\mathcal{K}$ . A reproducing kernel Hilbert space with  $\mathcal{K}$  could be decomposed as  $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle$  [SS01] for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ . Thus, we also identify  $f_\theta$  as  $\sum_{i=1}^l \alpha_i \cdot \Phi(\mathbf{a}_i)$ .

### Power Non-Linear Kernels

(polynomial)  $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle^k$  (Gaussian)  $\mathcal{K}(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}}$

## Our Contributions

We study the problem of teaching kernel perceptron (an ERM learner) under the *constructive* setting, where the teacher can construct arbitrary teaching examples in the support of the data distribution. This work extends the result of [LZ16] for teaching linear ERM learners.

- We formally define approximate teaching of kernel perceptron, and propose a novel measure of teaching complexity, namely the  $\epsilon$ -approximate teaching dimension ( $\epsilon$ -TD), which captures the complexity of teaching a “relaxed” target that is close to the target hypothesis in terms of the expected risk.
- We establish tight bounds on the teaching dimension of linear and polynomial perceptron. We exhibit optimal training sets that match these teaching dimensions.
- We show that for Gaussian kernelized perceptron, exact teaching is not possible with a finite set of examples, and then establish a  $d^{\mathcal{O}(\log^2 \frac{1}{\epsilon})}$  bound on the  $\epsilon$ -approximate teaching dimension.

Table 1: Main Results

|                            | linear      | polynomial                            | gaussian                                     |
|----------------------------|-------------|---------------------------------------|--|
| TD (exact)                 | $\Theta(d)$ | $\Theta\left(\binom{d+k-1}{k}\right)$ | $\infty$                                     |
| $\epsilon$ -approximate TD | -           | -                                     | $d^{\mathcal{O}(\log^2 \frac{1}{\epsilon})}$ |
| Assumption                 | -           | 1                                     | 2, 3   |

## Teaching Complexity

Fix a kernel perceptron learner, a kernel  $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  with the corresponding RKHS feature map  $\Phi(\cdot)$  and a target model  $\theta^* \in \mathcal{H}_{\mathcal{K}}$ .

### Teaching Set

A set of labelled points  $\mathcal{TS} \subseteq \mathcal{X} \times \mathcal{Y}$  provided by a helpful teacher for teaching a target hypothesis  $\theta^*$  to a kernel perceptron learner.

### Teaching Dimension for Exact Parameters

We define the teaching dimension for *exact* parameter (upto decision boundary) of  $\theta^*$  corresponding to a kernel perceptron as  $TD(t\theta^*, \mathcal{A}_{opt})$ , which is the size of the smallest teaching set  $\mathcal{TS}$  such that  $\mathcal{A}_{opt}(\mathcal{TS}) = \{t\theta^*\}$  for some real  $t > 0$ , where

$$TD(\{t\theta^*\}, \mathcal{A}_{opt}) = \min_{\text{real } p > 0} TD(p\theta^*, \mathcal{A}_{opt}).$$

### Approximate Teaching

**Definition 1 ( $\epsilon$ -approximate TS).** For a given  $\epsilon > 0$ , we say  $\mathcal{TS} \subseteq \mathcal{X} \times \mathcal{Y}$  is an  $\epsilon$ -approximate teaching set wrt to  $\mathcal{P}$  if the kernel perceptron  $\hat{\theta} \in \mathcal{A}_{opt}(\mathcal{TS})$  satisfies

$$\left| \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} [\max(-y \cdot f^*(\mathbf{x}), 0)] - \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}} [\max(-y \cdot \hat{f}(\mathbf{x}), 0)] \right| \leq \epsilon \text{ where } f^*(\cdot) = \theta^* \cdot \Phi(\cdot), \hat{f}(\cdot) = \hat{\theta} \cdot \Phi(\cdot)$$

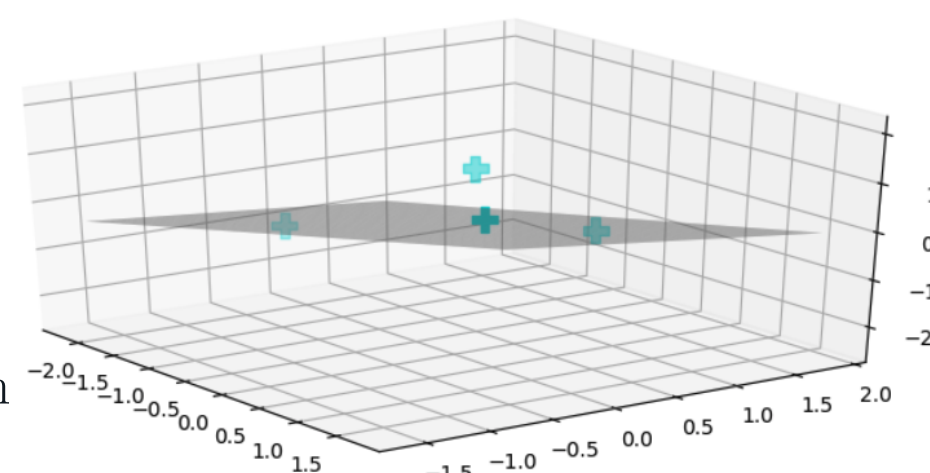
**Definition 2 ( $\epsilon$ -approximate TD).** For a given  $\epsilon > 0$ , we define  $\epsilon$ -TD( $\theta^*, \mathcal{A}_{opt}$ ) as the teaching dimension which is the size of the smallest teaching set for  $\epsilon$ -approximate teaching of  $\theta^*$  wrt  $\mathcal{P}$ .

## Teaching Homogeneous Linear Perceptron

We show  $TD(\{t\theta^*\}, \mathcal{A}_{opt}) = \Theta(d)$  via constructing a teaching set:

$$\begin{aligned} \mathbf{x}_i &= \mathbf{v}_i, & y_i &= 1 \quad \forall i \in [d-1]; \\ \mathbf{x}_d &= -\sum_{i=1}^{d-1} \mathbf{v}_i, & y_d &= 1; & \mathbf{x}_{d+1} &= \theta^*, & y_{d+1} &= 1 \end{aligned}$$

where  $\{\mathbf{v}_i\}_{i=1}^d$  is an orthogonal basis for  $\mathbb{R}^d$  which extends with  $\mathbf{v}_d = \theta^*$ .



## Teaching Non-Linear Kernel Perceptron

### Main Ideas and Results for Polynomial Kernel Perceptron

- $\mathcal{H}_{\mathcal{K}}$  is isomorphic to space of homogeneous polynomials.
- Assumption 1:** For the target model  $\theta^* \in \mathcal{H}_{\mathcal{K}}$ , we assume that there exist  $(r-1)$  linearly independent polynomials on the orthogonal subspace of  $\theta^*$  in  $\mathcal{H}_{\mathcal{K}}$  of the form  $\{\Phi(\mathbf{z}_i)\}_{i=1}^{r-1}$  where  $\forall i \mathbf{z}_i \in \mathcal{X}$ .
- Without **Assumption 1**,  $\theta^*$  can't be taught *exactly* and infeasibility of *approximate* teaching for some pathological cases.
- Under **Assumption 1**, we achieve a tight bound on teaching dimension (Table 1)

### Gaussian Kernel Perceptron

- Key Ideas
  - Truncate the Taylor features of the Gaussian Kernel to obtain a finite dimensional kernel.
  - Inspired from polynomial setting, we make a necessary assumption on projection of  $\theta^*$  in truncated space.
  - Analyze the solutions of perceptron algorithm and pick bounded solutions. (**Assumption 3**)
- Truncated Taylor Features

$$\underbrace{\Phi_{k,\lambda}(\mathbf{x}) = e^{-\frac{\|\mathbf{x}\|^2}{2\sigma^2}} \cdot \frac{\sqrt{C_\lambda^k}}{\sqrt{k!}\sigma^k} \cdot \mathbf{x}^\lambda}_{\text{coordinated Gaussian feature map: } \Phi(\cdot)} \longrightarrow \underbrace{\tilde{\Phi}_{k,\lambda}(\mathbf{x}) = \Phi_{k,\lambda}(\mathbf{x})}_{\text{Taylor features } \forall k \leq s} \longrightarrow \underbrace{\tilde{\mathcal{K}}(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x}\|^2 + \|\mathbf{x}'\|^2}{2\sigma^2}} \sum_{k=0}^s \frac{1}{k!} \left(\frac{\langle \mathbf{x}, \mathbf{x}' \rangle}{\sigma^2}\right)^k}_{\text{truncated kernel with projection map}} \longrightarrow \mathbb{P} \text{ (projection map)}$$

### Assumptions and Reformulation of Perceptron Algorithm

**Assumption 2 (Existence of orthogonal classifiers):** For some  $\epsilon > 0$ ,  $\exists r = r(\theta^*, \epsilon)$  such that  $\mathbb{P}\theta^*$  has  $r-1$  linear independent projections on the orthogonal subspace of  $\mathbb{P}\theta^*$  in  $\mathcal{H}_{\mathcal{K}}$  of the form  $\{\Phi(\mathbf{z}_i)\}_{i=1}^{r-1}$  such that  $\forall i \mathbf{z}_i \in \mathcal{X}$ .

Fix  $\mathcal{TS}_{\theta^*} := \{(\mathbf{z}_i, 1), (\mathbf{z}_i, -1)\}_{i=1}^{r-1} \cup \{(\mathbf{a}, 1)\}$  where  $\mathbb{P}\theta^* \cdot \mathbb{P}\Phi(\mathbf{a}) > 0$  and  $\Phi(\mathbf{a}) \cdot \Phi(\mathbf{z}_i) \leq Q\epsilon$  (constant  $Q$ ).

$$\mathcal{A}_{opt}(\mathcal{TS}_{\theta^*}) := \arg \min_{\beta_0 \in \mathbb{R}, \gamma \in \mathbb{R}^{r-1}} \max_{i=1}^{2r-1} (\ell(\beta_0, \gamma, \mathbf{x}_i, y_i), 0) \quad (1)$$

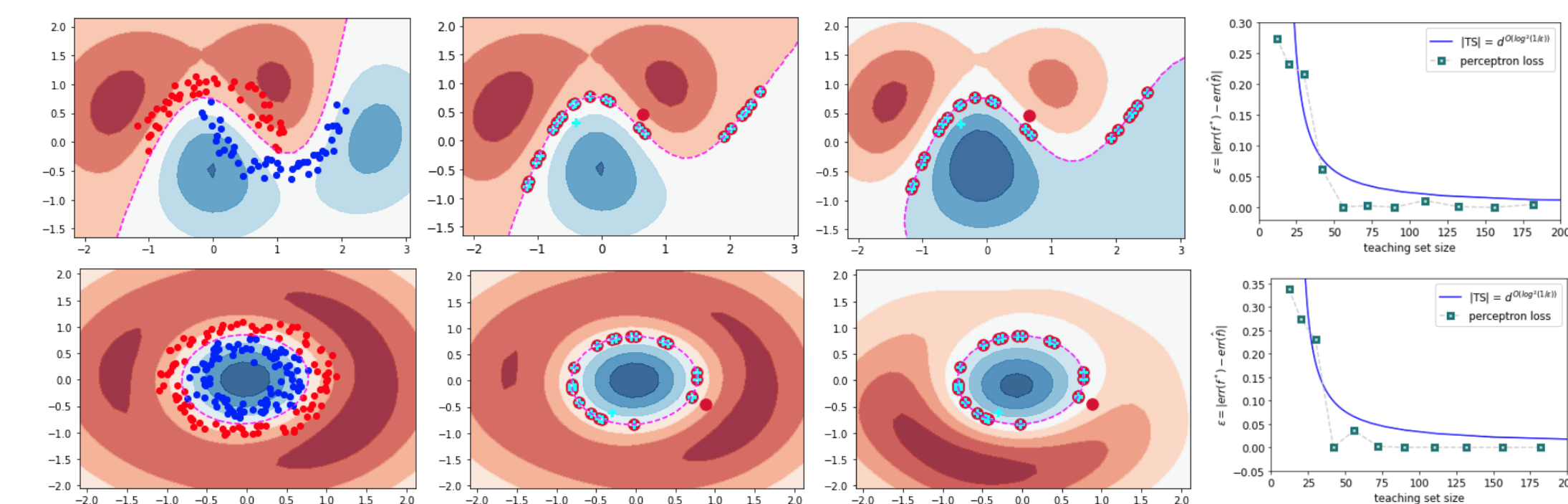
where for any  $i \in [2r-1]$   $\ell(\beta_0, \gamma, \mathbf{x}_i, y_i) = -y_i \cdot (\beta_0 \cdot \mathcal{K}(\mathbf{a}, \mathbf{x}_i) + \sum_{j=1}^{r-1} \gamma_j \cdot \mathcal{K}(\mathbf{z}_j, \mathbf{x}_i))$ .

Note,  $\hat{\theta} \in \mathcal{A}_{opt}(\mathcal{TS}_{\theta^*})$  has the form  $\beta_0 \cdot \mathcal{K}(\mathbf{a}, \cdot) + \sum_{j=1}^{r-1} \gamma_j \cdot \mathcal{K}(\mathbf{z}_j, \cdot)$  and target model has the form  $\theta^* = \sum_{i=1}^l \alpha_i \cdot \mathcal{K}(\mathbf{a}_i, \cdot)$ .

**Assumption 3 (Bounded Cone):** The learner optimizes to a solution  $\hat{\theta}$  for Eq. (1) with bounded coefficients, i.e.,  $\sum_{i=1}^l |\alpha_i|$  and  $|\beta_0| + \sum_{j=1}^{r-1} |\gamma_j|$  are bounded where  $\hat{\theta} \in \mathcal{H}_{\mathcal{K}}$  has the form  $\hat{\theta} = \beta_0 \cdot \mathcal{K}(\mathbf{a}, \cdot) + \sum_{j=1}^{r-1} \gamma_j \cdot \mathcal{K}(\mathbf{z}_j, \cdot)$ .

### Main Results

- Under **Assumptions 2-3** and  $\epsilon > 0$ , for any  $\hat{f} \in \mathcal{A}_{opt}(\mathcal{TS}_{\theta^*})$   $|f^*(\mathbf{x}) - \hat{f}(\mathbf{x})| \leq \epsilon$ .
- Under **Assumptions 2-3** and  $\epsilon > 0$ , the teaching set  $\mathcal{TS}_{\theta^*}$  constructed for Eq. (1) is an  $\epsilon$ -approximate teaching set. (Table 1)



## References

- [SS01] Bernhard Scholkopf and Alexander J. Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Cambridge, MA, USA: MIT Press, 2001. ISBN: 0262194759.
- [LZ16] Ji Liu and Xiaojin Zhu. “The Teaching Dimension of Linear Learners”. In: *Journal of Machine Learning Research* 17:162 (2016), pp. 1–25.