

# Nonparametric Teaching for Multiple Learners

**Chen Zhang<sup>1</sup>, Xiaofeng Cao<sup>1</sup>, Weiyang Liu<sup>2,3</sup>, Ivor W. Tsang<sup>4</sup>, James T. Kwok<sup>5</sup>**

<sup>1</sup>Jilin University

<sup>2</sup>Max Planck Institute for Intelligent Systems

<sup>3</sup>University of Cambridge

<sup>4</sup>Agency for Science, Technology and Research

<sup>5</sup>Hong Kong University of Science and Technology

October 21, 2023

## 1. What is Machine Teaching?

## 2. Multi-learner nonparametric teaching

- 2.1 Teaching Settings
- 2.2 Vanilla Multi-learner Teaching
- 2.3 Communicated Multi-learner Teaching

## 3. Experiments and Results

# What is Machine Teaching?



Machine teaching (MT) [10, 11] considers the problem of how to design the most effective **teaching set**, typically with the **smallest amount** of (teaching) examples possible, to facilitate rapid learning of the **target models** by learners based on these examples.

It can be thought of as **an inverse of machine learning**, in the sense that the learner is to learn models on a given dataset, while the teacher is to seek a (minimal) dataset from a target model.

Depending on how teachers and learners **interact** with each other, MT can be carried out in either

- **batch** fashion [10, 7, 3, 8] which focuses on **single-round** interaction, that is, the most representative and effective teaching dataset are designed to be fed to the learner in one shot, or
- **iterative** fashion [4, 5, 6] where an iterative teacher would feed examples based on learners' status (current learnt models) **round by round**.

# Multi-learner nonparametric teaching

Previous nonparametric teaching algorithms [9] merely focus on the **single-learner setting** (*i.e.*, teaching a **scalar-valued** target model or function to a single learner). To empower them to fulfill the practical needs of complex tasks, we introduce a more comprehensive task called **Multi-learner Nonparametric Teaching** (MINT). In MINT, the teacher aims to instruct **multiple learners**, with each learner focusing on learning a **scalar-valued** target model.

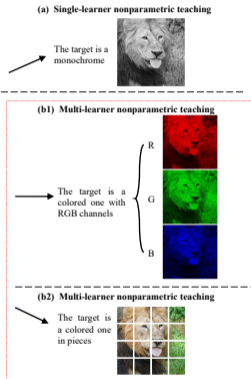
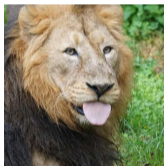


Figure: Comparison between the single-learner teaching and MINT.

## Main Contribution:

- By analyzing general **vector-valued RKHS**, we study the **multi-learner nonparametric teaching** (MINT), where the teacher selects examples based on a **vector-valued target function** (each component of it is a scalar-valued one for a single learner) such that **multiple** learners can learn its components simultaneously in a fast speed.
- Allowing the **communication** across multiple learners, that is, learners are allowed to carry out **linear combination** on current learnt functions of all learners, we investigate a communicated MINT where the teacher not only selects examples but also constructs a **matrix** as the guide of communication in each iteration.
- Under a mild assumption, we **theoretically** prove the efficiency of our **multi-learner generalization** of nonparametric teaching. We also **empirically** demonstrate its applicability and efficiency in extensive multi-learner experiments.

**Vector-valued Functional Optimization:** We define multi-learner nonparametric teaching as a **vector-valued functional minimization** over the collection of potential teaching sequences  $\mathcal{D}$  in the vector-valued reproducing kernel Hilbert space:

$$\mathcal{D}^* = \arg \min_{\mathcal{D} \in \mathbb{D}^d} \mathcal{M}(\hat{\mathbf{f}}^*, \mathbf{f}^*) + \lambda \cdot \text{len}(\mathcal{D}) \quad \text{s.t.} \quad \hat{\mathbf{f}}^* = \mathcal{A}(\mathcal{D}) \quad (1)$$

where  $\mathcal{M}$  denotes a discrepancy measure,  $\text{len}(\mathcal{D})$ , which is regularized by a constant  $\lambda$ , is the length of the teaching sequence  $\mathcal{D}$ , and  $\mathcal{A}$  represents the learning algorithm of learners. Specifically,  $\mathcal{A}$  is taken as  $\hat{\mathbf{f}}^* = \arg \min_{\mathbf{f} \in \mathcal{H}^d} \mathbb{E}_{(\mathbf{x}, \mathbf{y})} [\mathcal{L}(\mathbf{f}(\mathbf{x}), \mathbf{y})]$ ,

where  $(\mathbf{x}, \mathbf{y}) \in \mathcal{X}^d \times \mathcal{Y}^d$  and  $(\mathbf{x}, \mathbf{y}) \sim [\mathbb{Q}_i(x_i, y_i)]^d$ . Evaluated at an example vector  $(\mathbf{x}, \mathbf{y}) = [(x_{i,j_i}, y_{i,j_i})]^d$  with the example index  $j_i \in \mathbb{N}_k$ , the **multi-learner convex** loss  $\mathcal{L}$  therein is  $\mathcal{L}(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{i=1}^d \mathcal{L}_i(f_i(x_{i,j_i}), y_{i,j_i}) = E_{\mathbf{x}} [[\mathcal{L}_i(f_i, y_{i,j_i})]^d]$ , where  $\mathcal{L}_i$  is the **convex** loss for  $i$ -th learner.

We investigate MINT in the **gray-box setting**, which is equivalent to the one considered in [9]. To facilitate the theoretical analysis, we adopt some moderate **assumptions regarding  $\mathcal{L}_i$  and kernels**, which align with those made in [9].

## Assumption 1

Each loss  $\mathcal{L}_i(f_i), i \in \mathbb{N}_d$  is  $L_{\mathcal{L}_i}$ -Lipschitz smooth, *i.e.*,  $\forall f_i, f'_i \in \mathcal{H}, x_i \in \mathcal{X}$  and  $i \in \mathbb{N}_d$

$$|E_{x_i} [\nabla_f \mathcal{L}_i(f_i)] - E_{x_i} [\nabla_f \mathcal{L}_i(f'_i)]| \leq L_{\mathcal{L}_i} |E_{x_i} [f_i] - E_{x_i} [f'_i]|,$$

where  $L_{\mathcal{L}_i} \geq 0$  is a constant. To simplify the notation, we assume that  $L_{\mathcal{L}_i} = L_{\mathcal{L}}$  for all  $i \in \mathbb{N}_d$ .

## Assumption 2

Each kernel  $K(x, x') \in \mathcal{H}$  is bounded, *i.e.*,  $\forall x, x' \in \mathcal{X}, K(x, x') \leq M_K$ , where  $M_K \geq 0$  is a constant.

# Vanilla Multi-learner Teaching

In tackling MINT, we begin by examining a basic scenario in which **multiple learners concurrently** learns corresponding components of a vector-valued target function **without communication** between them [2, 1].

## Lemma 3 (Sufficient Descent for multi-learner RFT)

Suppose there are  $d$  learners, and the example **mean** for each learner is

$\mu_i = \mathbb{E}_{x_i \sim \mathbb{P}_i(x_i)}(x_i) < \infty$ , and the **variance**  $\sigma_i^2 = \mathbb{E}_{x_i \sim \mathbb{P}_i(x_i)}(x_i - \mu_i)^2 < \infty$ ,  $i \in \mathbb{N}_d$ .

Under the **Lipschitz smooth and bounded kernel assumptions**, if  $\eta_i^t \leq \frac{1}{2L_{\mathcal{L}} \cdot M_K}$  for all  $i \in \mathbb{N}_d$ , then RFT teachers can, **on average**, reduce the multi-learner loss  $\mathcal{L}(\mathbf{f})$  by:

$$\mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^{t+1}) - \mathcal{L}(\mathbf{f}^t)] \leq -\frac{\tilde{\eta}^t}{2} \sum_{i=1}^d (m_{i,t}(\mu_i) + \frac{m''_{i,t}(\mu_i)}{2} \sigma_i^2) \leq -\frac{\tilde{\eta}^t d}{2} \cdot \min_{i \in \mathbb{N}_d} \left( m_{i,t}(\mu_i) + \frac{m''_{i,t}(\mu_i)}{2} \sigma_i^2 \right), \quad (2)$$

where  $\tilde{\eta}^t = \min_{i \in \mathbb{N}_d} \eta_i^t$  and  $m_{i,t}(\dot{x}) := E_{\dot{x}}[(\nabla_f \mathcal{L}_i(\mathbf{f})|_{\mathbf{f}=\mathbf{f}_i^t})^2]$ .



### Theorem 4 (Convergence for multi-learner RFT)

Suppose the **vector-valued** model for multiple learners is initialized with  $\mathbf{f}^0 \in \mathcal{H}^d$  and returns  $\mathbf{f}^t \in \mathcal{H}^d$  after  $t$  iterations, we have the **upper bound** of  $\min_{i \in \mathbb{N}_d} \left( m_{i,t}(\mu_i) + m''_{i,t}(\mu_i) \sigma_i^2 / 2 \right)$  w.r.t.  $t$ :

$$\min_{i \in \mathbb{N}_d} \left( m_{i,t-1}(\mu_i) + m''_{i,t-1}(\mu_i) \sigma_i^2 / 2 \right) \leq 2 \mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^0)] / (d\eta t), \quad (3)$$

where  $0 < \eta = \min_{l \in \{0\} \cup \mathbb{N}_{t-1}} \tilde{\eta}^l \leq 1 / (2L_{\mathcal{L}} \cdot M_K)$ , and given a small constant  $\epsilon > 0$  it would take approximately

$\mathcal{O} \left( 2(\mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^0)] - \epsilon) / (d\eta \min_{i \in \mathbb{N}_d} \left( m_{i,t-1}(\mu_i) + m''_{i,t-1}(\mu_i) \sigma_i^2 / 2 \right)) \right)$  iterations to reduce the **multi-learner** loss  $\mathcal{L}$  to a **sufficiently small** value and to reach a **stationary point** in terms of  $\mathcal{L}$ .

## Lemma 5 (Sufficient Descent for multi-learner GFT)

Under the same assumption, if  $\eta_i^t \leq \frac{1}{2L_{\mathcal{L}} \cdot M_K}$  for all  $i \in \mathbb{N}_d$ , the GFT teachers can achieve a **greater** reduction in the multi-learner loss  $\mathcal{L}$ :

$$\mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^{t+1}) - \mathcal{L}(\mathbf{f}^t)] \leq -\frac{\tilde{\eta}^t}{2} \sum_{i=1}^d m_{i,t}(x_i^{t*}) \leq -\frac{\tilde{\eta}^t d}{2} \cdot \min_{i \in \mathbb{N}_d} m_{i,t}(x_i^{t*}), \quad (4)$$

where  $\tilde{\eta}^t$  and  $m_{i,t}(\cdot)$  retain their previous meaning.

## Theorem 6 (Convergence for multi-learner GFT)

Suppose the **vector-valued** model for multiple learners is initialized with  $\mathbf{f}^0 \in \mathcal{H}^d$  and returns  $\mathbf{f}^t \in \mathcal{H}^d$  after  $t$  iterations, we have the **upper bound** of  $\min_{i \in \mathbb{N}_d} m_{i,t}(x_i^{t*})$  w.r.t.  $t$ :

$$\min_{i \in \mathbb{N}_d} m_{i,t-1}(x_i^{t-1*}) \leq \frac{2}{d\eta t} \mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^0)] + \frac{1}{d} \sum_{l=0}^{t-1} \sum_{i=1}^d \left( \|x_i^{l*} - \mu_i\|_2 \right), \quad (5)$$

where  $\eta$  has the same definition as before, and given a small constant  $\epsilon > 0$  it would need around  $\mathcal{O} \left( 2(\mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^0)] - \epsilon) / (d\eta \min_{i \in \mathbb{N}_d} m_{i,t-1}(x_i^{t-1*})) \right)$  iterations to decrease the **multi-learner** loss  $\mathcal{L}$  to a **sufficiently small** value and to reach a **stationary point** in terms of  $\mathcal{L}$ .

# Communicated Multi-learner Teaching



An infant would **integrate** previously learnt knowledge to grasp a new target concept, such as comprehending what a zebra is by combining the learnt ideas of horses and black-and-white stripes. Such an efficient paradigm motivates us to explore the **communicated MINT**, which enables the **communication** between learners.

## Proposition 5

If the proximity between  $\mathbf{f}^t$  and  $\mathbf{f}^*$  is **sufficiently close**, meaning that  $\|\mathbf{f}^t - \mathbf{f}^*\|_{\mathcal{H}^d} \leq \epsilon$  where  $\epsilon$  is a tiny positive constant, then  $A^t$  equals the **identity matrix**  $I_d$ .

## Lemma 6

Under **Lipschitz smooth** assumption, the **communication** across learners will result in a **reduction** of the **multi-learner convex** loss  $\mathcal{L}$  by

$$0 \leq \mathcal{L}(\mathbf{f}^t) - \mathcal{L}(A^t \mathbf{f}^t) \leq 2L_{\mathcal{L}} \|\mathbf{f}^t - \mathbf{f}^*\|_{\mathcal{H}^d}.$$

## Theorem 7

Suppose the **communication** in the  $t$ -th iteration of multiple learners is denoted by the **matrix**  $A^t$  and returns  $\mathbf{f}_{A^t}^{t+1} \in \mathcal{H}^d$ , for both RFT and GFT we have:

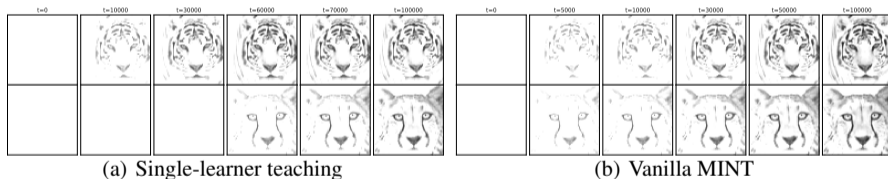
$$\mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}_{A^t}^{t+1}) - \mathcal{L}(\mathbf{f}^t)] \leq \mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}_{A^t}^{t+1}) - \mathcal{L}(A^t \mathbf{f}^t)] \leq 0.$$

# Experiments and Results

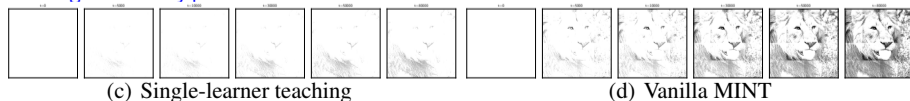
Testing the teaching of a **multi-learner (vector-valued) target model**, MINT presents more satisfactory performance than repeatedly carrying out the single-learner teaching, which is **consistent with our theoretical findings**.

- **MINT in gray scale.**

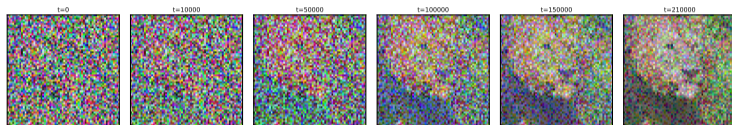
Simultaneous teaching of a tiger and a cheetah.



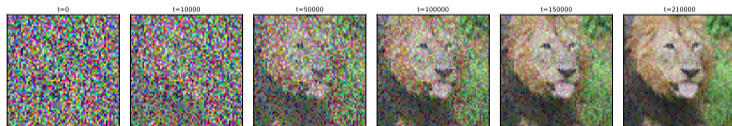
Teaching of a lion by partition.



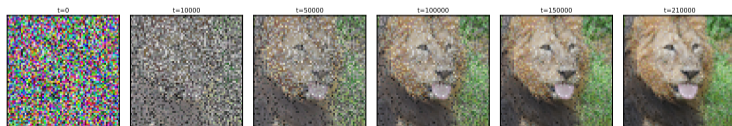
- MINT in three (RGB) channels.



(a) Single-learner teaching.



(b) Vanilla MINT.



(c) Communicated MINT.

**Thank you for listening!**



- [1] Nicolò Cesa-Bianchi, Pierre Laforgue, Andrea Paudice, et al. Multitask online mirror descent. Transactions of Machine Learning Research.
- [2] Sham M Kakade, Shai Shalev-Shwartz, and Ambuj Tewari. Regularization techniques for learning with matrices. The Journal of Machine Learning Research, 13(1):1865–1890, 2012.
- [3] Akash Kumar, Hanqi Zhang, Adish Singla, and Yuxin Chen. The teaching dimension of kernel perceptron. In AISTATS, 2021.
- [4] Weiyang Liu, Bo Dai, Ahmad Humayun, Charlene Tay, Chen Yu, Linda B Smith, James M Rehg, and Le Song. Iterative machine teaching. In ICML, 2017.
- [5] Weiyang Liu, Bo Dai, Xingguo Li, Zhen Liu, James Rehg, and Le Song. Towards black-box iterative machine teaching. In ICML, 2018.
- [6] Weiyang Liu, Zhen Liu, Hanchen Wang, Liam Paull, Bernhard Schölkopf, and Adrian Weller. Iterative teaching by label synthesis. In NeurIPS, 2021.
- [7] Farnam Mansouri, Yuxin Chen, Ara Vartanian, Jerry Zhu, and Adish Singla. Preference-based batch and sequential teaching: Towards a unified view of models. In NeurIPS, 2019.

- [8] Hong Qian, Xu-Hui Liu, Chen-Xi Su, Aimin Zhou, and Yang Yu. The teaching dimension of regularized kernel learners. In ICML, 2022.
- [9] Chen Zhang, Xiaofeng Cao, Weiyang Liu, Ivor Tsang, and James Kwok. Nonparametric iterative machine teaching. In ICML, 2023.
- [10] Xiaojin Zhu. Machine teaching: An inverse problem to machine learning and an approach toward optimal education. In AAAI, 2015.
- [11] Xiaojin Zhu, Adish Singla, Sandra Zilles, and Anna N Rafferty. An overview of machine teaching. arXiv preprint arXiv:1801.05927, 2018.