

Machine Teaching

Machine teaching (MT) 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 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 where an iterative teacher would feed examples based on learners' status (current learnt models) **round by round**, such that the learner can converge to a target model within fewer rounds.

Motivation

Previous nonparametric teaching algorithms 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.

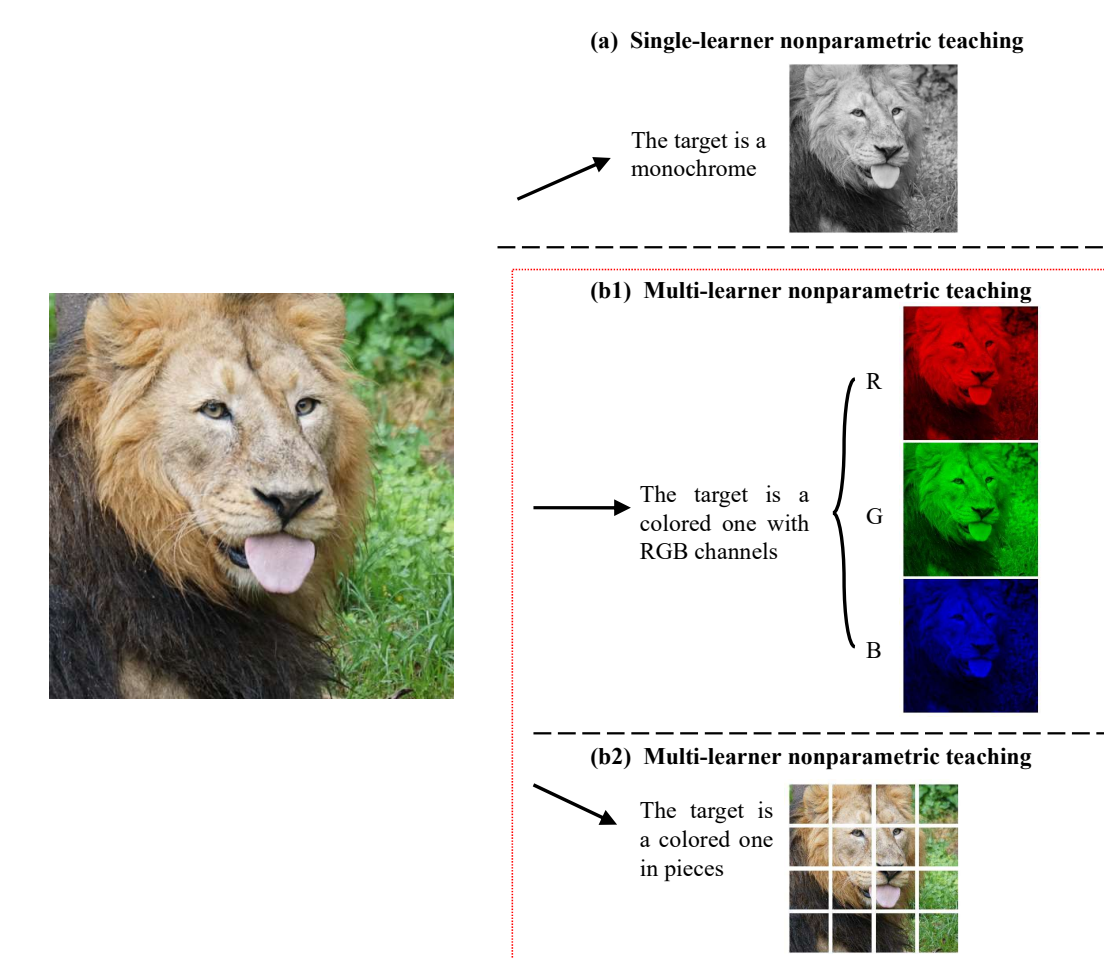


Figure: Comparison between the single-learner teaching and MINT.

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.
- ▶ By enabling the **communication** among multiple learners, learners can update themselves with a **linear combination** of current learnt functions of all learners. We study a communicated MINT where the teacher not only selects examples but also injects the **guidance** of communication.
- ▶ Under mild assumptions, we characterize the **efficiency** of our **multi-learner generalization** of nonparametric teaching. More importantly, we also **empirically** demonstrate its efficiency.

Teaching Settings

Vector-valued Functional Optimization: We define multi-learner nonparametric teaching as a **vector-valued functional minimization** over the collection of potential teaching sequences \mathbb{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 λ , 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 $\mathcal{L}(\mathbf{f}(\mathbf{x}), \mathbf{y}) = \sum_{i=1}^d \mathcal{L}_i(f_i(x_{i,j_i}), y_{i,j_i}) = E_{\mathbf{x}} [\sum_{i=1}^d \mathcal{L}_i(f_i, y_{i,j_i})]$, where \mathcal{L}_i is the **convex** loss for i -th learner.

Vanilla Multi-learner Teaching

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

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

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

Theorem 2 (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 $\min_{i \in \mathbb{N}_d} (m_{i,t}(\mu_i) + m_{i,t}''(\mu_i) \sigma_i^2 / 2)$ w.r.t. t :

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

where $0 < \tilde{\eta} = \min_{i \in \mathbb{N}_d} \eta_i^t \leq 1/(2L_{\mathcal{L}} M_K)$, and given a small constant $\epsilon > 0$ it would take approximately $\mathcal{O}(2 \mathbb{E}_{\mathbf{x} \sim [\mathbb{P}_i(x_i)]^d} [\mathcal{L}(\mathbf{f}^0)] / (d\tilde{\eta}\epsilon))$ iterations to reach a **stationary point**.

Lemma 3 (Sufficient Descent for multi-learner GFT). Under the same assumption, if $\eta_i^t \leq \frac{1}{2L_{\mathcal{L}} M_K}$ for all $i \in \mathbb{N}_d$, the GFT teachers can achieve a **greater** reduction in the multi-learner loss \mathcal{L} :

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

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

Theorem 4 (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 :

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

where $\tilde{\eta}$ has the same definition as before.

Communicated Multi-learner Teaching

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

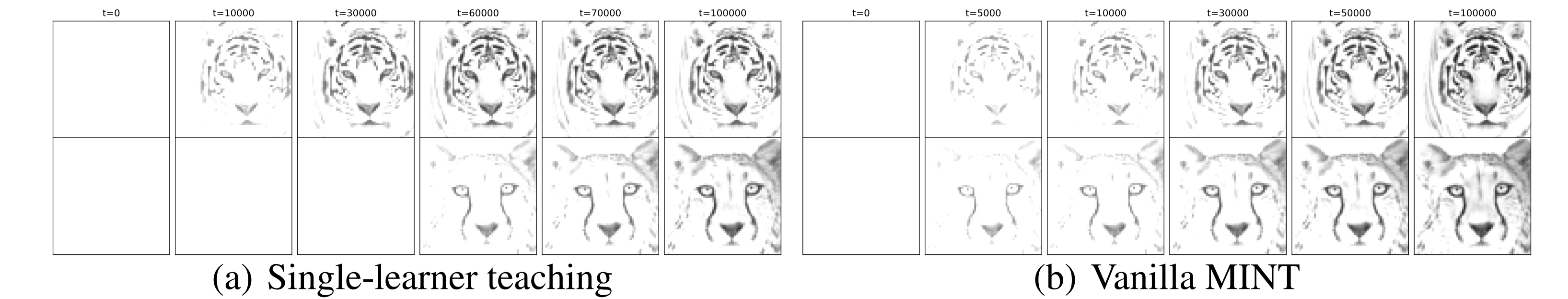
Theorem 7 Suppose the **communication** in the t -th iteration of multiple learners is denoted by the **matrix** A^t and returns $\mathbf{f}_{A^t}^{t+1} \in \mathcal{H}^d$, for both RFT and GFT we have:

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

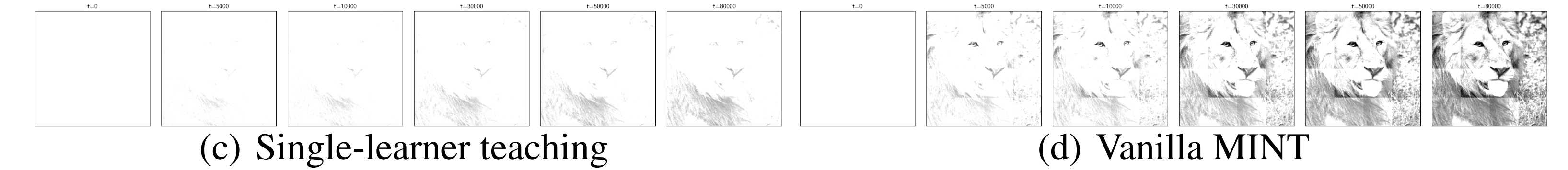
Experiments and Results

▶ MINT in gray scale.

Simultaneous teaching of a tiger and a cheetah.



Teaching of a lion by partition.



▶ MINT in three (RGB) channels.

