# An Optimal Control View of Adversarial Machine Learning

Xiaojin Zhu

Department of Computer Sciences, University of Wisconsin-Madison

**Abstract**

I describe an optimal control view of adversarial machine learning, where the dynamical system is the machine learner, the input are adversarial actions, and the control costs are defined by the adversary's goals to do harm and be hard to detect. This view encompasses many types of adversarial machine learning, including test-item attacks, training-data poisoning, and adversarial reward shaping. The view encourages adversarial machine learning researcher to utilize advances in control theory and reinforcement learning.

## 1 Adversarial Machine Learning is not Machine Learning

Machine learning has its mathematical foundation in concentration inequalities. This is a consequence of the independent and identically-distributed (i.i.d.) data assumption. In contrast, I suggest that adversarial machine learning may adopt optimal control as its mathematical foundation [3, 25]. There are telltale signs: adversarial attacks tend to be subtle and have peculiar non-i.i.d. structures – as control input might be.

## 2 Optimal Control

I will focus on deterministic discrete-time optimal control because it matches many existing adversarial attacks. Extensions to stochastic and continuous control are relevant to adversarial machine learning, too. The system to be controlled is called the plant, which is defined by the system dynamics:

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) \tag{1}$$

where $\mathbf{x}_t \in \mathbf{X}_t$ is the state of the system, $\mathbf{u}_t \in \mathbf{U}_t$ is the control input, and $\mathbf{U}_t$ is the control constraint set. The function $f$ defines the evolution of state under external control. The time index $t$ ranges from 0 to $T - 1$, and the time horizon $T$ can be finite or infinite. The quality of control is specified by the running cost:

$$g_t(\mathbf{x}_t, \mathbf{u}_t) \tag{2}$$

which defines the step-by-step control cost, and the terminal cost for finite horizon:

$$g_T(\mathbf{x}_T) \tag{3}$$

which defines the quality of the final state. The optimal control problem is to find control inputs $\mathbf{u}_0 \ldots \mathbf{u}_{T-1}$ in order to minimize the objective:

$$\min_{\mathbf{u}_0 \ldots \mathbf{u}_{T-1}} \quad g_T(\mathbf{x}_T) + \sum_{t=0}^{T-1} g_t(\mathbf{x}_t, \mathbf{u}_t) \tag{4}$$

$$\text{s.t.} \quad \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t), \ \mathbf{u}_t \in \mathbf{U}_t, \ \forall t$$

$$\mathbf{x}_0 \text{ given}$$

1

More generally, the controller aims to find control policies $\phi_t(\mathbf{x}_t) = \mathbf{u}_t$, namely functions that map observed states to inputs. In optimal control the dynamics $f$ is known to the controller. There are two styles of solutions: dynamic programming and Pontryagin minimum principle [2, 10, 17]. When $f$ is not fully known, the problem becomes either robust control where control is carried out in a minimax fashion to accommodate the worst case dynamics [28], or reinforcement learning where the controller probes the dynamics [23].

# 3 Adversarial Machine Learning as Control

Now let us translate adversarial machine learning into a control formulation. Adversarial machine learning studies vulnerability throughout the learning pipeline [4, 13, 20, 26]. As examples, I present training-data poisoning, test-time attacks, and adversarial reward shaping below. In all cases, the adversary attempts to control the machine learning system, and the control costs reflect the adversary's desire to do harm and be hard to detect.

Unfortunately, the notations from the control community and the machine learning community clash. For example, $\mathbf{x}$ denotes the state in control but the feature vector in machine learning. I will use the machine learning convention below.

## 3.1 Training-Data Poisoning

In training-data poisoning the adversary can modify the training data. The machine learner then trains a "wrong" model from the poisoned data. The adversary's goal is for the "wrong" model to be useful for some nefarious purpose. I use supervised learning for illustration.

### 3.1.1 Batch Learner

At this point, it becomes useful to distinguish batch learning and sequential (online) learning. If the machine learner performs batch learning, then the adversary has a degenerate *one-step* control problem. One-step control has not been the focus of the control community and there may not be ample algorithmic solutions to borrow from. Still, it is illustrative to pose batch training set poisoning as a control problem. I use Support Vector Machine (SVM) with a batch training set as an example below:

- The state is the learner's model $h : \mathbf{X} \mapsto \mathbf{Y}$. For instance, for SVM $h$ is the classifier parametrized by a weight vector $\mathbf{w}$. I will use $h$ and $\mathbf{w}$ interchangeably.

- The control $\mathbf{u}_0$ is a whole training set, for instance $\mathbf{u}_0 = \{(\mathbf{x}_i, y_i)\}_{1:n}$.

- The control constraint set $\mathbf{U}_0$ consists of training sets available to the adversary; if the adversary can arbitrary modify a training set for supervised learning (including changing features and labels, inserting and deleting items), this could be $\mathbf{U}_0 = \cup_{n=0}^{\infty} (\mathbf{X} \times \mathbf{Y})^n$, namely all training sets of all sizes. This is a large control space.

- The system dynamics (1) is defined by the learner's learning algorithm. For the SVM learner, this would be empirical risk minimization with hinge loss $\ell()$ and a regularizer:

$$\mathbf{w}_1 = f(\mathbf{u}_0) \in \mathrm{argmin}_{\mathbf{w}} \sum_{i=1}^{n} \ell(\mathbf{w}, \mathbf{x}_i, y_i) + \lambda \|\mathbf{w}\|^2. \tag{5}$$

  The batch SVM does not need an initial weight $\mathbf{w}_0$. The adversary has full knowledge of the dynamics $f()$ if it knows the form (5), $\ell()$, and the value of $\lambda$.

- The time horizon $T = 1$.

- The adversary's running cost $g_0(\mathbf{u}_0)$ measures the poisoning effort in preparing the training set $\mathbf{u}_0$. This is typically defined with respect to a given "clean" data set $\tilde{\mathbf{u}}$ before poisoning in the form of

$$g_0(\mathbf{u}_0) = \text{distance}(\mathbf{u}_0, \tilde{\mathbf{u}}). \tag{6}$$

The running cost is domain dependent. For example, the distance function may count the number of modified training items; or sum up the Euclidean distance of changes in feature vectors.

- The adversary's terminal cost $g_1(\mathbf{w}_1)$ measures the lack of intended harm. The terminal cost is also domain dependent. For example:

   - If the adversary must force the learner into exactly arriving at some target model $\mathbf{w}^*$, then $g_1(\mathbf{w}_1) = \mathbb{I}_\infty[\mathbf{w}_1 \neq \mathbf{w}^*]$. Here $\mathbb{I}_y[z] = y$ if $z$ is true and 0 otherwise, which acts as a hard constraint.

   - If the adversary only needs the learner to get near $\mathbf{w}^*$ then $g_1(\mathbf{w}_1) = \|\mathbf{w}_1 - \mathbf{w}^*\|$ for some norm.

   - If the adversary wants to ensure that a specific future item $\mathbf{x}^*$ is classified $\epsilon$-confidently as positive, it can use $g_1(\mathbf{w}_1) = \mathbb{I}_\infty[\mathbf{w}_1 \notin \mathbf{W}^*]$ with the target set $\mathbf{W}^* = \{\mathbf{w} : \mathbf{w}^\top \mathbf{x}^* \geq \epsilon\}$. More generally, $\mathbf{W}^*$ can be a polytope defined by multiple future classification constraints.

With these definitions, the adversary's one-step control problem (4) specializes to

$$\min_{\mathbf{u}_0} \quad g_1(\mathbf{w}_1) + g_0(\mathbf{w}_0, \mathbf{u}_0) \tag{7}$$
$$\text{s.t.} \quad \mathbf{w}_1 = f(\mathbf{w}_0, \mathbf{u}_0)$$

Unsurprisingly, the adversary's one-step control problem is equivalent to a Stackelberg game and bi-level optimization (the lower level optimization is hidden in $f$), a well-known formulation for training-data poisoning [12, 21].

### 3.1.2 Sequential Learner

The adversary performs classic discrete-time control if the learner is sequential:

- The learner starts from an initial model $\mathbf{w}_0$, which is the initial state.

- The control input at time $t$ is $\mathbf{u}_t = (\mathbf{x}_t, y_t)$, namely the $t^{th}$ training item for $t = 0, 1, \ldots$

- The dynamics is the sequential update algorithm of the learner. For example, the learner may perform one step of gradient descent:

$$\mathbf{w}_{t+1} = f(\mathbf{w}_t, \mathbf{u}_t) = \mathbf{w}_t - \eta_t \nabla \ell(\mathbf{w}_t, \mathbf{x}_t, y_t). \tag{8}$$

- The adversary's running cost $g_t(\mathbf{w}_t, \mathbf{u}_t)$ typically measures the effort of preparing $\mathbf{u}_t$. For example, it could measure the magnitude of change $\|\mathbf{u}_t - \tilde{\mathbf{u}}_t\|$ with respect to a "clean" reference training sequence $\tilde{\mathbf{u}}$. Or it could be the constant 1 which reflects the desire to have a short control sequence.

- The adversary's terminal cost $g_T(\mathbf{w}_T)$ is the same as in the batch case.

The problem (4) then produces the optimal training sequence poisoning. Earlier attempts on sequential teaching can be found in [1, 18, 19].

3

## 3.2 Test-Time Attack

Test-time attack differs from training-data poisoning in that a machine learning model $h : \mathbf{X} \mapsto \mathbf{Y}$ is already-trained and given. Also given is a "test item" $\mathbf{x}$. There are several variants of test-time attacks, I use the following one for illustration: The adversary seeks to minimally perturb $\mathbf{x}$ into $\mathbf{x}'$ such that the machine learning model classifies $\mathbf{x}$ and $\mathbf{x}'$ differently. That is,

$$\min_{\mathbf{x}'} \quad \text{distance}(\mathbf{x}, \mathbf{x}') \tag{9}$$
$$\text{s.t.} \quad h(\mathbf{x}) \neq h(\mathbf{y}).$$

The distance function is domain-dependent, though in practice the adversary often uses a mathematically convenient surrogate such as some $p$-norm $\|\mathbf{x} - \mathbf{x}'\|_p$.

One way to formulate test-time attack as optimal control is to treat the test-item itself as the state, and the adversarial actions as control input. Let us first look at the popular example of test-time attack against image classification:

- Let the initial state $\mathbf{x}_0 = \mathbf{x}$ be the clean image.

- The adversary's control input $\mathbf{u}_0$ is the vector of pixel value changes.

- The control constraint set is $\mathbf{U}_0 = \{\mathbf{u} : \mathbf{x}_0 + \mathbf{u} \in [0, 1]^d\}$ to ensure that the modified image has valid pixel values (assumed to be normalized in $[0, 1]$).

- The dynamical system is trivially vector addition: $\mathbf{x}_1 = f(\mathbf{x}_0, \mathbf{u}_0) = \mathbf{x}_0 + \mathbf{u}_0$.

- The adversary's running cost is $g_0(\mathbf{x}_0, \mathbf{u}_0) = \text{distance}(\mathbf{x}_0, \mathbf{x}_1)$.

- The adversary's terminal cost is $g_1(\mathbf{x}_1) = \mathbb{I}_\infty[h(\mathbf{x}_1) = h(\mathbf{x}_0)]$. Note the machine learning model $h$ is only used to define the hard constraint terminal cost; $h$ itself is not modified.

With these definitions this is a one-step control problem (4) that is equivalent to the test-time attack problem (9).

This control view on test-time attack is more interesting when the adversary's actions are sequential $\mathbf{U}_0, \mathbf{U}_1, \ldots$, and the system dynamics render the action sequence non-commutative. The adversary's running cost $g_t$ then measures the effort in performing the action at step $t$. One limitation of the optimal control view is that the action cost is assumed to be additive over the steps.

## 3.3 Defense Against Test-Time Attack by Adversarial Training

Some defense strategies can be viewed as optimal control, too. One defense against test-time attack is to require the learned model $h$ to have the large-margin property with respect to a training set. Let $(\mathbf{x}, y)$ be any training item, and $\epsilon$ a margin parameter. Then the large-margin property states that the decision boundary induced by $h$ should not pass $\epsilon$-close to $(\mathbf{x}, y)$:

$$\forall \mathbf{x}' : (\|\mathbf{x}' - \mathbf{x}\|_p \leq \epsilon) \Rightarrow h(\mathbf{x}') = y. \tag{10}$$

This is an uncountable number of constraints. It is relatively easy to enforce for linear learners such as SVMs, but impractical otherwise.

Adversarial training can be viewed as a heuristic to approximate the uncountable constraint (10) with a finite number of active constraints: one performs test-time attack against the current $h$ from $\mathbf{x}$ to find an adversarial item $\mathbf{x}^{(1)}$, such that $\|\mathbf{x}^{(1)} - \mathbf{x}\|_p \leq \epsilon$ but $h(\mathbf{x}^{(1)}) \neq y$. Instead of adding a single constraint $h(\mathbf{x}^{(1)}) = y$, an additional training item $(\mathbf{x}^{(1)}, y)$ is then added to the training set. The machine learning algorithm learns a different $h$, with the hope (but not constraining) that $h(\mathbf{x}^{(1)}) = y$. This process repeats for $k$ iteration, resulting in $k$ additional training items $(\mathbf{x}^{(i)}, y)$ for $i = 1 \ldots k$.

It should be clear that such defense is similar to training-data poisoning, in that the defender uses data to modify the learned model. This is especially interesting when the learner performs sequential updates. One way to formulate adversarial training defense as control is the following:

- The state is the model $h_t$. Initially $h_0$ can be the model trained on the original training data.

- The control input $\mathbf{u}_t = (\mathbf{x}_t, y_t)$ is an additional training item with the trivial constraint set $\mathbf{U}_t = \mathbf{X} \times \mathbf{y}$.

- The dynamics $h_{t+1} = f(h_t, \mathbf{u}_t)$ is one-step update of the model, e.g. by back-propagation.

- The defender's running cost $g_t(h_t, \mathbf{u}_t)$ can simply be 1 to reflect the desire for less effort (the running cost sums up to $k$).

- The defender's terminal cost $g_T(h_T)$ penalizes small margin of the final model $h_T$ with respect to the original training data.

Of course, the resulting control problem (4) does not directly utilize adversarial examples. One way to incorporate them is to restrict $\mathbf{U}_t$ to a set of adversarial examples found by invoking test-time attackers on $h_t$, similar to the heuristic in [7]. These adversarial examples do not even need to be successful attacks.

## 3.4   Adversarial Reward Shaping

When adversarial attacks are applied to sequential decision makers such as multi-armed bandits or reinforcement learning agents, a typical attack goal is to force the latter to learn a wrong policy useful to the adversary. The adversary may do so by manipulating the rewards and the states experienced by the learner [11, 14].

To simplify the exposition, I focus on adversarial reward shaping against stochastic multi-armed bandit, because this does not involve deception through perceived states. To review, in stochastic multi-armed bandit the learner at iteration $t$ chooses one of $k$ arms, denoted by $I_t \in [k]$, to pull according to some strategy [6]. For example, the $(\alpha, \psi)$-Upper Confidence Bound (UCB) strategy chooses the arm

$$I_t \in \text{argmax}_{i \in [k]} \hat{\mu}_{i, T_i(t-1)} + \psi^{*-1} \left( \frac{\alpha \log t}{T_i(t-1)} \right) \tag{11}$$

where $T_i(t-1)$ is the number of times arm $i$ has been pulled up to time $t-1$, $\hat{\mu}_{i, T_i(t-1)}$ is the empirical mean of arm $i$ so far, and $\psi^*$ is the dual of a convex function $\psi$. The environment generates a stochastic reward $r_{I_t} \sim \nu_{I_t}$. The learner updates its estimate of the pulled arm:

$$\hat{\mu}_{I_t, T_{I_t}(t)} = \frac{\hat{\mu}_{I_t, T_{I_t}(t-1)} T_{I_t}(t-1) + r_{I_t}}{T_{I_t}(t-1) + 1} \tag{12}$$

which in turn affects which arm it will pull in the next iteration. The learner's goal is to minimize the pseudo-regret $T\mu^{\max} - \mathbb{E} \sum_{t=1}^{T} \mu_{I_t}$ where $\mu_i = \mathbb{E}\nu_i$ and $\mu^{\max} = \max_{i \in [k]} \mu_i$. Stochastic multi-armed bandit strategies offer upper bounds on the pseudo-regret.

With adversarial reward shaping, an adversary fully observes the bandit. The adversary intercepts the environmental reward $r_{I_t}$ in each iteration, and may choose to modify ("shape") the reward into

$$r_{I_t} + u_t$$

with some $u_t \in \mathbb{R}$ before sending the modified reward to the learner. The adversary's goal is to use minimal reward shaping to force the learner into performing specific wrong actions. For example, the adversary may want the learner to frequently pull a particular target arm $i^* \in [k]$. It should be noted that the adversary's goal may not be the exact opposite of the learner's goal: the target arm $i^*$ is not necessarily the one with the worst mean reward, and the adversary may not seek pseudo-regret maximization.

Adversarial reward shaping can be formulated as stochastic optimal control:

- The state $s_t$, now called control state to avoid confusion with the Markov Decision Process states experienced by an reinforcement learning agent, consists of the sufficient statistic tuple at time $t$:

$$s_t = (T_1(t-1), \hat{\mu}_{1, T_1(t-1)}, \dots, T_k(t-1), \hat{\mu}_{k, T_k(t-1)}, I_t).$$

- The control input is $u_t \in \mathbf{U}_t$ with $\mathbf{U}_t = \mathbb{R}$ in the unconstrained shaping case, or the appropriate $\mathbf{U}_t$ if the rewards must be binary, for example.

- The dynamics $s_{t+1} = f(s_t, u_t)$ is straightforward via empirical mean update (12), $T_{I_t}$ increment, and new arm choice (11).

- The adversary's running cost $g_t(s_t, u_t)$ reflects shaping effort and target arm achievement in iteration $t$. For instance,
$$g_t(s_t, u_t) = u_t^2 + \mathbb{I}_\lambda[I_t \neq i^*]. \tag{13}$$
where $\lambda > 0$ is a trade off parameter.

- There is not necessarily a time horizon $T$ or a terminal cost $g_T(s_T)$.

The control state is stochastic due to the stochastic reward $r_{I_t}$ entering through (12).

# 4 Advantages of the Optimal Control View

There are a number of potential benefits in taking the optimal control view:

- It offers a unified conceptual framework for adversarial machine learning;

- The optimal control literature provides efficient solutions when the dynamics $f$ is known and one can take the continuous limit to solve the differential equations [15];

- Reinforcement learning, either model-based with coarse system identification or model-free policy iteration, allows approximate optimal control when $f$ is unknown, as long as the adversary can probe the dynamics [8, 9];

- A generic defense strategy may be to limit the controllability the adversary has over the learner.

- I mention in passing that the optimal control view applies equally to machine teaching [27, 29], and thus extends to the application of personalized education [22, 24].

I need to point out some limitations:

- Having a unified optimal control view does not automatically produce efficient solutions to the control problem (4). For adversarial machine learning applications the dynamics $f$ is usually highly nonlinear and complex. Furthermore, in graybox and blackbox attack settings $f$ is not fully known to the attacker. They affect the complexity in finding an optimal control.

- The adversarial learning setting is largely non-game theoretic, though there are exceptions [5, 16].

These problems call for future research from both machine learning and control communities.

# References

[1] Scott Alfeld, Xiaojin Zhu, and Paul Barford. Data poisoning attacks against autoregressive models. In *The Thirtieth AAAI Conference on Artificial Intelligence (AAAI-16)*, 2016.

[2] Michael Athans and Peter L Falb. *Optimal control: An introduction to the theory and its applications.* Courier Corporation, 2013.

[3] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control.* Athena Scientific, 4th edition, 2017.

[4] Battista Biggio and Fabio Roli. Wild patterns: Ten years after the rise of adversarial machine learning. *CoRR*, abs/1712.03141, 2017.

[5] Michael Brückner and Tobias Scheffer. Stackelberg games for adversarial prediction problems. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 547–555. ACM, 2011.

[6] Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.

[7] Qi-Zhi Cai, Min Du, Chang Liu, and Dawn Song. Curriculum adversarial training. In *The 27th International Joint Conference on Artificial Intelligence (IJCAI)*, 2018.

[8] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack on graph structured data. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1115–1124, Stockholmsmssan, Stockholm Sweden, 10–15 Jul 2018. PMLR.

[9] Yang Fan, Fei Tian, Tao Qin, and Tie-Yan Liu. Learning to teach. In *ICLR*, 2018.

[10] Terry L Friesz. *Dynamic optimization and differential games*, volume 135. Springer Science & Business Media, 2010.

[11] Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *arXiv*, 2017.

[12] Matthew Jagielski, Alina Oprea, Battista Biggio, Chang Liu, Cristina Nita-Rotaru, and Bo Li. Manipulating machine learning: Poisoning attacks and countermeasures for regression learning. *The 39th IEEE Symposium on Security and Privacy*, 2018.

[13] Anthony D. Joseph, Blaine Nelson, Benjamin I. P. Rubinstein, and J. D. Tygar. *Adversarial Machine Learning*. Cambridge University Press, 2018. in press.

[14] Kwang-Sung Jun, Lihong Li, Yuzhe Ma, and Xiaojin Zhu. Adversarial attacks on stochastic bandits. In *Advances in Neural Information Processing Systems (NIPS)*, 2018.

[15] L. Lessard, X. Zhang, and X. Zhu. An Optimal Control Approach to Sequential Machine Teaching. *ArXiv e-prints*, October 2018.

[16] Bo Li and Yevgeniy Vorobeychik. Scalable Optimization of Randomized Operational Decisions in Adversarial Classification Settings. In Guy Lebanon and S. V. N. Vishwanathan, editors, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pages 599–607, San Diego, California, USA, 09–12 May 2015. PMLR.

[17] Daniel Liberzon. *Calculus of variations and optimal control theory: A concise introduction*. Princeton University Press, 2011.

[18] Weiyang Liu, Bo Dai, Ahmad Humayun, Charlene Tay, Chen Yu, Linda B Smith, James M Rehg, and Le Song. Iterative machine teaching. In *International Conference on Machine Learning*, pages 2149–2158, 2017.

[19] Weiyang Liu, Bo Dai, Xingguo Li, Zhen Liu, James M. Rehg, and Le Song. Towards black-box iterative machine teaching. In *ICML*, volume 80 of *JMLR Workshop and Conference Proceedings*, pages 3147–3155. JMLR.org, 2018.

[20] Daniel Lowd and Christopher Meek. Adversarial learning. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 641–647. ACM, 2005.

[21] Shike Mei and Xiaojin Zhu. Using machine teaching to identify optimal training-set attacks on machine learners. In *The Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.

[22] Kaustubh Patil, Xiaojin Zhu, Lukasz Kopec, and Bradley Love. Optimal teaching for limited-capacity human learners. In *Advances in Neural Information Processing Systems (NIPS)*, 2014.

[23] B. Recht. A Tour of Reinforcement Learning: The View from Continuous Control. *ArXiv e-prints*, June 2018.

[24] Ayon Sen, Purav Patel, Martina A. Rau, Blake Mason, Robert Nowak, Timothy T. Rogers, and Xiaojin Zhu. Machine beats human at sequencing visuals for perceptual-fluency practice. In *Educational Data Mining*, 2018.

[25] Emanuel Todorov. Optimal control theory. *Bayesian brain: probabilistic approaches to neural coding*, pages 269–298, 2006.

[26] Yevgeniy Vorobeychik and Murat Kantarcioglu. Adversarial machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 12(3):1–169, 2018.

[27] Xiaojin Zhu. Machine teaching: an inverse problem to machine learning and an approach toward optimal education. In *The Twenty-Ninth AAAI Conference on Artificial Intelligence (AAAI "Blue Sky" Senior Member Presentation Track)*, 2015.

[28] Xiaojin Zhu, Ji Liu, and Manuel Lopes. No learner left behind: On the complexity of teaching multiple learners simultaneously. In *The 26th International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.

[29] Xiaojin Zhu, Adish Singla, Sandra Zilles, and Anna N. Rafferty. An Overview of Machine Teaching. *ArXiv e-prints*, January 2018. https://arxiv.org/abs/1801.05927.