# Machine Teaching

## Adish Singla

CMMRS, August 2019
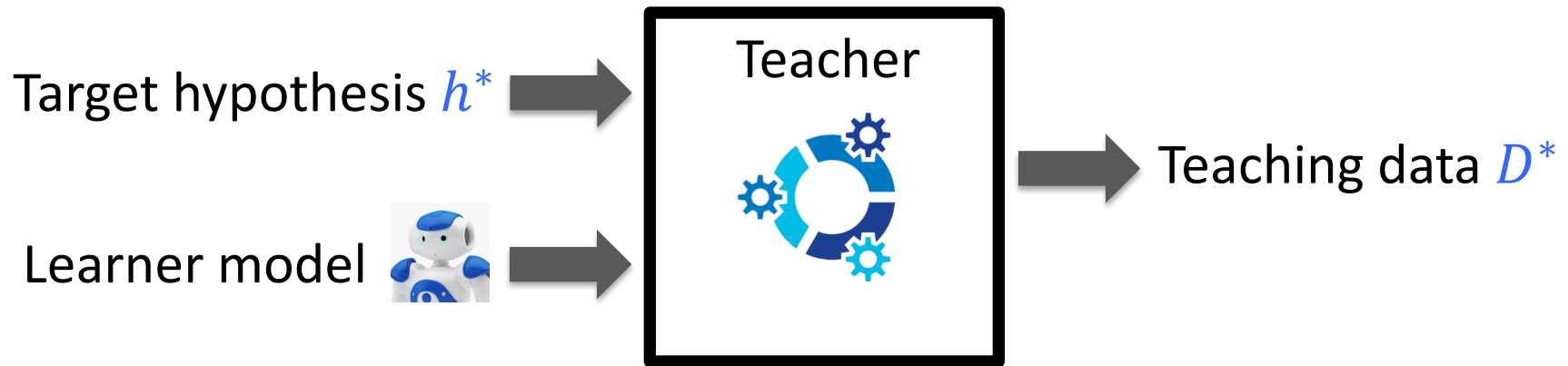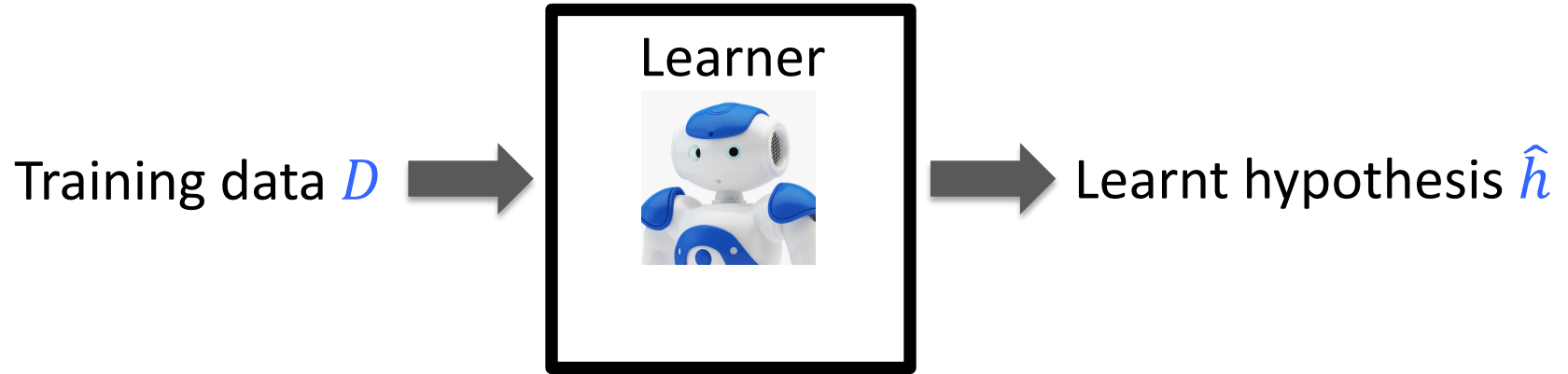
# Machine Learning vs. Teaching

Training data $D$ → Learner → Learnt hypothesis $\hat{h}$

Target hypothesis $h^*$ → Teacher

Learner model → Teacher → Teaching data $D^*$

# Why Machine Teaching?



Adversarial settings
aka training-set poisoning
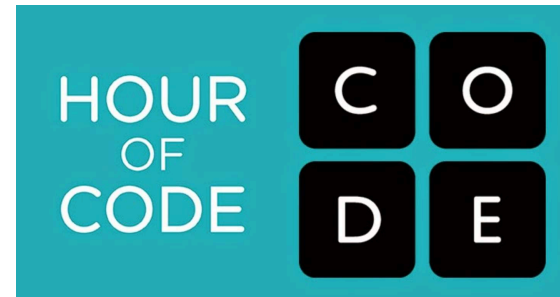
Educational settings

# Applications: Online Education via MOOCs

- Astronomical growth with over 100 million students
- Over 10,000 courses offered online

**Key challenge**: Dropout rate of over 95%

# Applications: Skill Assessment and Practice





- Over 10 million problems solved per year on ASSISTments
- Over 0.8 billion hours of code by 100 million students

**Key limitation**: No automated or personalized curriculum of problems

# Applications: Training Simulators

# Applications: Language Learning



- Over 300+ million students
- Based on **spaced repetition** of flash cards

Can we compute **optimal personalized schedule** of repetition?

# Applications: Biodiversity Monitoring



**Key challenge**: Noise in the annotations

# Applications: Biodiversity Monitoring



- Teaching helps increase awareness and engagement
- Labeled data is crucial for training machine learning systems

Can we **teach** participants to label more accurately?

# Machine Teaching: Applications

**Online education via MOOCs**

**Skill assessment and practice**

**Educational settings**

**Language learning**

**Biodiversity monitoring**

**Training simulators**

# Machine Teaching: Key Components



Application

Learner's model

Teacher's algorithm

- Type and complexity of task

- Type and model of learning agent

- Teacher's knowledge and observability

# Course Outline

**Part 1**: Different viewpoints of the problem space

- **Information-theoretic** models of teaching
- **Cognitive** models of teaching

**Part 2**: Designing algorithms for teaching people

- **Classification** rules for biodiversity monitoring
- **Vocabulary** for language learning
- **Policies** for performing sequential tasks

- Type and complexity of task



- Type and model of learning agent



- Teacher's knowledge and observability

# An Example: 1-D Threshold Function

- Task: Classify animal image as  **Weevil** ⬤  or  **Vespula** ⭕

- $\mathcal{X}$: Set of images, each $x \in \mathcal{X}$ is associated with a contrast level

- $\mathcal{H}$: Set of hypotheses, each $h \in \mathcal{H}$ is a binary threshold classifier

- $h^*$: True classifier

W W W W W W W V V V V

# An Example: 1-D Threshold Function

- Learning setting (Passive): avg. size of $D$ is $\Theta(n)$



- Learning setting (Active): size of $D$ is $\Theta(\log n)$



- Teaching setting: size of $D$ is $2$

# Teaching Binary Functions

- Set of unlabeled examples $\mathcal{X}$

- Hypotheses class $\mathcal{H}$ as a set of binary functions $h : \mathcal{X} \rightarrow \{0,1\}$

- Target hypothesis $h^* \in \mathcal{H}$

$\mathcal{X}$

|  | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ |
|---|---|---|---|---|---|
| $h_1$ | 1 | 1 | 1 | 1 | 1 |
| $h_2$ | 0 | 1 | 1 | 1 | 1 |
| $h_3$ | 0 | 0 | 1 | 1 | 1 |
| $h^* = h_4$ | 0 | 0 | 0 | 1 | 1 |
| $h_5$ | 0 | 0 | 0 | 0 | 1 |

$\mathcal{H}$

$\mathcal{X}$

|  | $x_1$ | $x_2$ | $x_3$ | $x_4$ |
|---|---|---|---|---|
| $h_1$ | 1 | 0 | 0 | 0 |
| $h_2$ | 0 | 1 | 0 | 0 |
| $h_3$ | 0 | 0 | 1 | 0 |
| $h_4$ | 0 | 0 | 0 | 1 |
| $h^* = h_5$ | 0 | 0 | 0 | 0 |

$\mathcal{H}$

# Teaching Interaction

**Start**

- Learner starts at $h_0 \in \mathcal{H}$

**At time $t$**



Teacher

$\mathcal{X}, \mathcal{H}, h^*$

Teacher receives an estimate of $h_{t-1}$

Teacher selects $x_t$, provides $x_t, h^*(x_t)$

Learner updates to $h_t$

Learner

$\mathcal{X}, \mathcal{H}, h_{t-1}$

**Stop**

- When $h_t = h^*$

# Learner Model: Version space learning

## Notion of version space

- Maintain a set of eligible hypotheses
  - Start with $H_0 = \mathcal{H}$
- At time $t$, remove hypothesis inconsistent with $x_t, h^*(x_t)$
  - $H_t = H_{t-1} \setminus \{h \in \mathcal{H} \mid h(x_t) \neq h^*(x_t)\}$

## Version space learner

- Learner starts at $h_0 \in \mathcal{H}$, $H_0 = \mathcal{H}$
- At time $t$:
  - Learner receives $x_t, h^*(x_t)$ and updates $H_t$
  - Learner selects a new hypothesis $h_t \in H_t$ at random

# Teacher: Optimization Problem

## Analysis setting

- **Worst-case** vs. average case
- **Finite** vs. infinite/continuous $\mathcal{H}$
- **Exact** vs. approximate teaching

## Optimization problem

- Find smallest sequence $\overleftrightarrow{S} = (x_1, x_2, \ldots)$ to achieve desired objective

$$\overleftrightarrow{S}^{\text{opt}} = \operatorname*{argmin}_{\overleftrightarrow{S}} |\overleftrightarrow{S}| \qquad \text{s.t.} \qquad \boxed{\begin{array}{c} h_t = h^* \\ \text{equivalent to} \\ H_t = \{h^*\} \end{array}}$$

# Teacher: Optimization Problem

**Teaching problem is equivalent to Set Cover problem**

- $\mathcal{H} \setminus \{h^*\}$ is the set of elements to remove or "cover"

- Each $x$ covers a subset $\mathcal{H}(x) = \{h \in \mathcal{H} \mid h(x) \neq h^*(x)\}$

- Find smallest set $S = \{x_1, x_2, \ldots\}$ to cover $\mathcal{H} \backslash h^*$

**Complexity of optimization**

**Theorem**: Finding optimal teaching sequence $\overleftrightarrow{S}^{\mathrm{opt}}$ is NP-hard.

# Teacher: Optimization Problem

## Teaching problem is a Submodular Coverage problem

- Define set function $F: 2^{\mathcal{X}} \rightarrow \mathbb{R}_+$ as

$$F(S) = |\bigcup_{x \in S} \mathcal{H}(x)| \text{ where } S \subseteq \mathcal{X}$$

- Rewrite teaching problem as

$$S^{\text{opt}} = \underset{S}{\arg\min} |S| \qquad \text{s.t.} \qquad F(S) \geq |\mathcal{H}| - 1$$

## Submodular Coverage problem

- $F(.)$ satisfies submodularity: A notion of diminishing returns

$$F(\{a\} \cup S) - F(S) \geq F(\{a\} \cup \{b\} \cup S) - F(\{b\} \cup S)$$

We can optimize using a greedy algorithm with provable guarantees

# Teacher: Algorithm

## Iterative greedy algorithm

- **Input:** $\mathcal{H}, \mathcal{X}, h^*$

- **Initialize:** set $S \leftarrow \emptyset$

- While $F(S) < |\mathcal{H}| - 1$:

  - Select $x \leftarrow \text{argmax}_{x' \in \mathcal{X}} F(x' \cup S) - F(S)$

  - Provide $x, h^*(x)$ to learner

  - Update $S \leftarrow S \cup \{x\}$

## Approximation guarantees

**Theorem**: Let $S^{\text{gr}}$ be the set provided by the algorithm and $\overleftrightarrow{S}^{\text{opt}}$ denote the optimal teaching sequence. Then, $|S^{\text{gr}}| \leq |\overleftrightarrow{S}^{\text{opt}}| \cdot \log(|\mathcal{H}|)$.

# Complexity Measures: TD

## Notion of teaching complexity: Teaching dimension TD

- Introduced by [Goldman, Kearns '95]
- Analysis setting
  - randomized version space learner
  - worst-case analysis
  - finite size hypothesis class
  - exact teaching

## Formal definition of TD

- Length of optimal teaching sequence for $h^*$ is $|TS(h^*; \mathcal{H}, \mathcal{X})|$
- Teaching dimension is defined as

$$TD(\mathcal{H}, \mathcal{X}) := \max_{h^* \in \mathcal{H}} |TS(h^*; \mathcal{H}, \mathcal{X})|$$

# Complexity Measures: TD

**Examples for computing TD**

| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | $|TS(h^*)|$ |
|---|---|---|---|---|---|---|
| $h_1$ | 1 | 1 | 1 | 1 | 1 | 1 |
| $h_2$ | 0 | 1 | 1 | 1 | 1 | 2 |
| $h_3$ | 0 | 0 | 1 | 1 | 1 | 2 |
| $h_4$ | 0 | 0 | 0 | 1 | 1 | 2 |
| $h_5$ | 0 | 0 | 0 | 0 | 1 | 2 |

$$TD(\mathcal{H}, \mathcal{X}) = 2$$

| | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $|TS(h^*)|$ |
|---|---|---|---|---|---|
| $h_1$ | 1 | 0 | 0 | 0 | 1 |
| $h_2$ | 0 | 1 | 0 | 0 | 1 |
| $h_3$ | 0 | 0 | 1 | 0 | 1 |
| $h_4$ | 0 | 0 | 0 | 1 | 1 |
| $h_5$ | 0 | 0 | 0 | 0 | 4 |

$$TD(\mathcal{H}, \mathcal{X}) = 4$$

# Complexity Measures: TD vs. VCD

## Notion of learning complexity: VCD

- Introduced by [Vapnik, Chervonenkis '71]

- Sample complexity bounds for learning grow as $\Theta\big(VCD(\mathcal{H}, \mathcal{X})\big)$

## A fundamental question: TD vs. VCD?

- $TD(\mathcal{H}, \mathcal{X})$ is $O\big(VCD(\mathcal{H}, \mathcal{X})\big)$?
- There exists problems with
  - $TD(\mathcal{H}, \mathcal{X}) \ll O\big(VCD(\mathcal{H}, \mathcal{X})\big)$
  - $TD(\mathcal{H}, \mathcal{X}) \gg O\big(VCD(\mathcal{H}, \mathcal{X})\big)$

|       | $x_1$ | $x_2$ | $x_3$ | $x_4$ |
|-------|-------|-------|-------|-------|
| $h_1$ | 1     | 0     | 0     | 0     |
| $h_2$ | 0     | 1     | 0     | 0     |
| $h_3$ | 0     | 0     | 1     | 0     |
| $h_4$ | 0     | 0     | 0     | 1     |
| $h_5$ | 0     | 0     | 0     | 0     |

$$TD(\mathcal{H}, \mathcal{X}) = 4$$

$$VCD(\mathcal{H}, \mathcal{X}) = 1$$

# Improved Notions of TD: RTD

## Teaching an "adversarial" learner: Classic TD

- Simple classes can be difficult to teach

## Teaching a "cooperative" learner: Recursive TD (RTD)

- Introduced by [Zilles et al. @ COLT'08]
- $RTD(\mathcal{H}, \mathcal{X})$ is $O\big(VCD(\mathcal{H}, \mathcal{X})\big)$? [Simon, Zilles @ COLT'15]
- An active area of research
  - $O\big(d\,2^d \log\log |\mathcal{H}|\big)$ [Moran et al. @ FOCS'15]
  - $O\big(d\,2^d\big)$ [Chen et al. @ NIPS' 16]
  - $O(d^2)$ [Hu et al. @ COLT' 17]
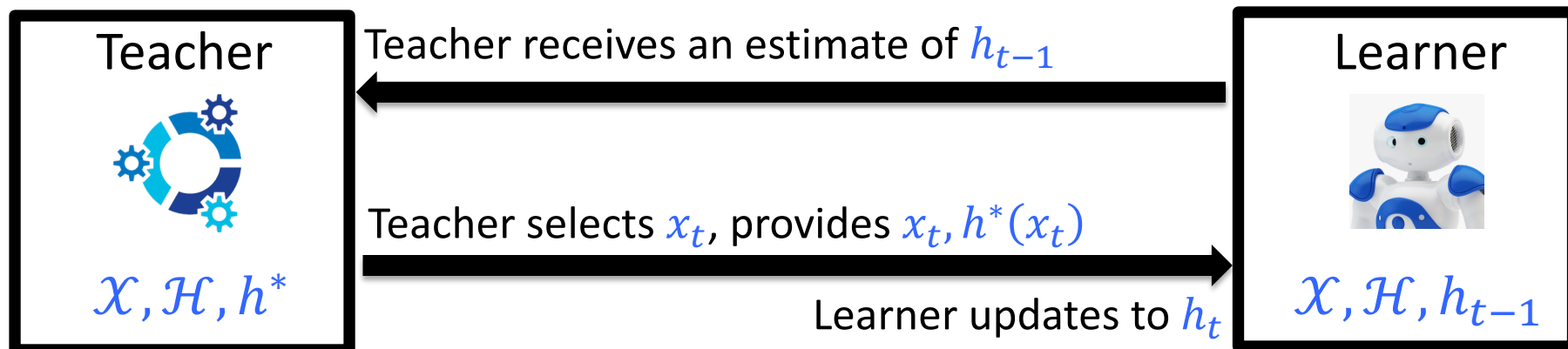
  where $d$ denotes $VCD(\mathcal{H}, \mathcal{X})$

# Improved Notions of TD: $TD_\sigma$

## Teaching models for classic TD or RTD

- Order of examples and learner's feedback does not matter

## Teaching a "state-dependent" learner: $TD_\sigma$

- Introduced in our recent work [NeurIPS'18, arXiv'19]
- Generalizes existing notions of teaching dimensions
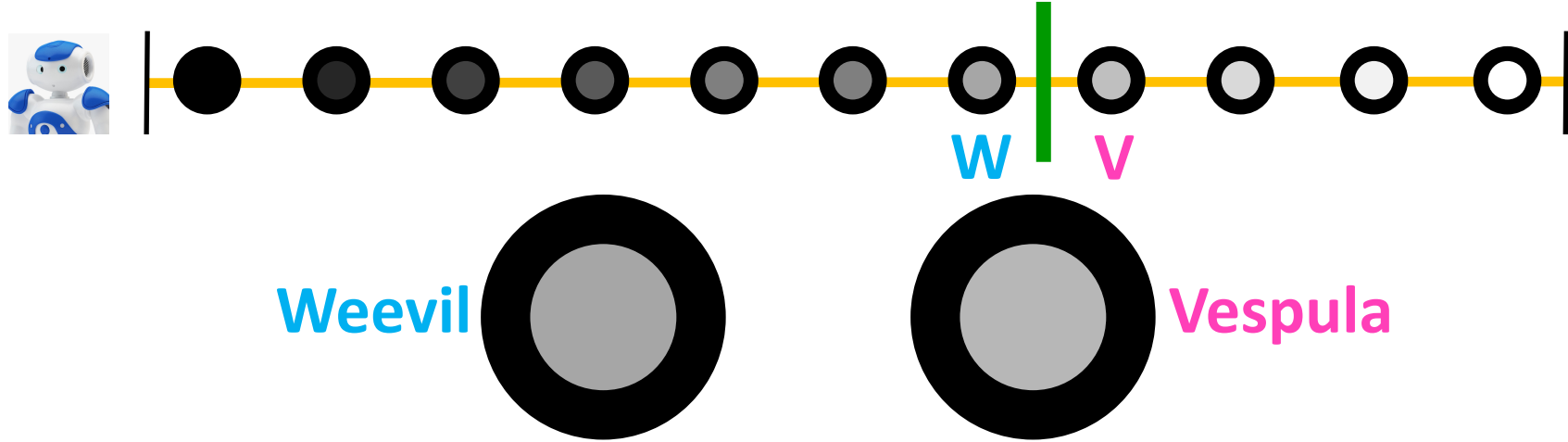- Provides necessary conditions when feedback matters



Teacher

$\mathcal{X}, \mathcal{H}, h^*$

Teacher receives an estimate of $h_{t-1}$

Teacher selects $x_t$, provides $x_t, h^*(x_t)$

Learner updates to $h_t$

Learner

$\mathcal{X}, \mathcal{H}, h_{t-1}$

# Teaching Binary Functions

- Understanding TD vs. VCD relation
  - see work by Sandra Zilles: http://www2.cs.uregina.ca/~zilles/

- Teaching complexity for ML models (e.g., SVM)
  - see work by Jerry Zhu: http://pages.cs.wisc.edu/~jerryzhu/

# Teaching Binary Functions to People

- Teaching setting: size of $D$ is $2$



Weevil

Vespula

W    V

# Teaching Binary Functions to People

- Teaching setting: size of $D$ is $2$



- Limited inference power and noise

- Mismatch in representation for $\mathcal{X}$, $\mathcal{H}$

- Limited memory

- Engagement

- Interpretability (e.g., teaching via labels vs. features)

- Safety (e.g., when teaching physical tasks)

- Fairness (e.g., when teaching a class)

More suitable for poisoning attacks, less for educational settings

- Type and complexity of task



- Type and model of learning agent



- Teacher's knowledge and observability