Seminar on Machine Teaching

Saarland University – Winter Semester 2019

Adish Singla

4th Nov 2019





Machine Learning vs. Teaching



Why Machine Teaching?





edX

Adversarial settings aka training-set poisoning **Educational settings**

Explore Learn Record

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





VIRTAMED^O WE SIMULATE REALITY



🥐 🗆 🖸

Video credits: Virtamed – Zurich, Switzerland

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

Downy Woodpecker (Picoides pubescens) Research Grade



Description

KONICA MINOLTA DIGITAL CAMERA				Downy Woodpecker (Picoides pubescens)		
Activi	ity		Cumulative ID	Ds: 2 of 2	es pubescens/	
	sy25805 suggested an ID	Ø ID Withdrawn 2mo 🖌 🗸	0		2/3rds	; 2
	Woodpeckers and Allies Order Piciformes		✓ Agr	ree 7	Compare	1 About
۲	sy25805 suggested an ID	Timproving 2mo	Annotatio	ns		
Ĭ	Downy Woodpecker		Attribute	Value	Agree	Disagree
	Picoides pubescens	- compare - Agree	Sex	Select -		
			Life Stage	Solact -		

Community ID

Key challenge: Noise in the annotations

Image credits: CornellLab – Global Big Day

Follow -

What's this?

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



Machine Teaching: Key Components



Machine Teaching: Problem Space



Course Logistics: Different Parts

Part 1: Research Papers

- 1/3 of the score
- No formal feedback (each submission gets full points)
- Reports due by 10th Dec'19 (#1, 2, 3, 4) and by 10th Jan'20 (#5, 6, 7, 8)

Part 2: Slides Preparation

- 1/3 of the score
- Paper assignment on 11th Jan'20
- Final slides due by 15th Feb'20

Part 3: Final Talk

- 1/3 of the score
- Presentations to be scheduled between mid-Feb'20 to mid-Mar'20

Course Logistics: Research Papers

- Each report submitted as lastname_paper#.pdf
 - singla_1.pdf, singla_2.pdf, ...
- Typeset in latex using NeurIPS style files
 - <u>https://neurips.cc/Conferences/2019/PaperInformation/StyleFiles</u>
 - \usepackage[preprint]{neurips_2019} (Non-anonymous preprints)
- Structure the report as an extended review, e.g.,
 - Summarize the paper
 - Write down main strengths of this paper
 - Write down main weaknesses of this paper
 - Write down ways in which this paper could be improved
 - Write down ideas in which this work could be extended

Machine Teaching: Problem Space

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

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 $h: \mathcal{X} \to \{0,1\}$
- Target hypothesis $h^* \in \mathcal{H}$

Teaching Interaction

Start

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

At time t

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 ${\cal H}$
- Exact vs. approximate teaching

Optimization problem

• Find smallest sequence $\overrightarrow{S} = (x_1, x_2, ...)$ to achieve desired objective

$$\vec{S}^{opt} = \underset{\vec{S}}{\operatorname{argmin}} |\vec{S}|$$
 s.t. $h_t = h^*$
equivalent to

 $H_t = \{h^*\}$

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 $S = \{x_1, x_2, ...\}$ to cover $\mathcal{H} \setminus h^*$

Complexity of optimization

Theorem: Finding optimal teaching sequence \overrightarrow{S}^{opt} is NP-hard.

Teacher: Optimization Problem

Teaching problem is a Submodular Coverage problem

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

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

• Rewrite teaching problem as

 $S^{\text{opt}} = \underset{S}{\operatorname{argmin}} |S| \quad \text{s.t.} \quad F(S) \ge |\mathcal{H}| - 1$

Submodular Coverage problem

• F(.) satisfies submodularity: A notion of diminishing returns $F(\{a\} \cup S) - F(S) \ge F(\{a\} \cup \{b\} \cup S) - F(\{b\} \cup S)$

We can optimize using a greedy algorithm with provable guarantees

Teacher: Algorithm

Iterative greedy algorithm

- Input: \mathcal{H} , \mathcal{X} , h^*
- Initialize: set S ← Ø
- While $F(S) < |\mathcal{H}| 1$:
 - Select $x \leftarrow \operatorname{argmax}_{x' \in \mathcal{X}} F(x' \cup S) F(S)$
 - Provide x, $h^*(x)$ to learner
 - Update $S \leftarrow S \cup \{x\}$

Approximation guarantees

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