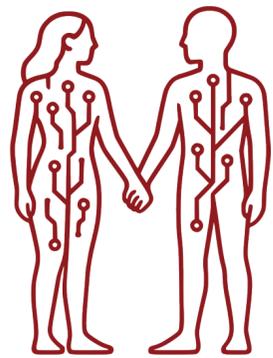


Learning in an Echo Chamber: Online Learning with Replay Adversary

Daniil Dmitriev*, Harald Eskelund Franck†, **Carolin Heinzler†**, Amartya Sanyal†



MACHINE LEARNING
UNIVERSITY OF COPENHAGEN

*University of Pennsylvania
†University of Copenhagen

Registration and travel support
for this presentation was provided by SIGACT.



Echo Chamber

**Everyone seems to agree with me,
so I must be right**

Sara Germain

Online Learning

Proper ($\hat{h}_t \in \mathcal{H}$) vs. Improper ($\hat{h}_t \in 2^{\mathcal{X}}$)

For $t = 1, \dots, T$:

Learner \mathcal{A} outputs a hypothesis $\hat{h}_t \in 2^{\mathcal{X}}$,

Realizable ($\exists f^* \in \mathcal{H}$)

Nature produces $x_t \in \mathcal{X}$ and reveals $(x_t, f^*(x_t))$ for some $f^* \in \mathcal{H}$,

Learner suffers loss $1 \{f^*(x_t) \neq \hat{h}_t(x_t)\}$

Adaptive adversary (adversary choosing $x_t \in \mathcal{X}$ every round t) vs. Stochastic adversary ($\exists \mathcal{D} : x_t \stackrel{i.i.d.}{\sim} \mathcal{D}$)

Online Learning with **Replay Adversary**

For $t = 1, \dots, T$:

Learner \mathcal{A} outputs a hypothesis $\hat{h}_t \in 2^{\mathcal{X}}$,

Nature produces $x_t \in \mathcal{X}$

Full loss NOT observed by learner

True label

Replay adversary reveals (x_t, y_t) with $y_t = \begin{cases} f^*(x_t) \text{ for some } f^* \in \mathcal{H} \\ \hat{h}_i(x_t) \text{ for some } i < t \end{cases}$

Learner suffers loss $1 \{y_t \neq \hat{h}_t(x_t) \text{ and } y_t = f^*(x_t)\}$

Replay label

Learning thresholds conservatively

\mathcal{H} class of thresholds on finite interval $[0,1]$, $|\mathcal{H}| = N$

The Learner

The Learner predicts with the smallest threshold (conservative and consistent).

$t = 1$

$$\hat{h}_1 = 0$$



$t = 1$

$$(x_1, y_1) = (1/4, 1)$$
$$\hat{h}_1(1/4) = 0 \neq y_1 \quad \times$$

A new positive point, which cannot be replay

→ Update classifier

Learning thresholds conservatively

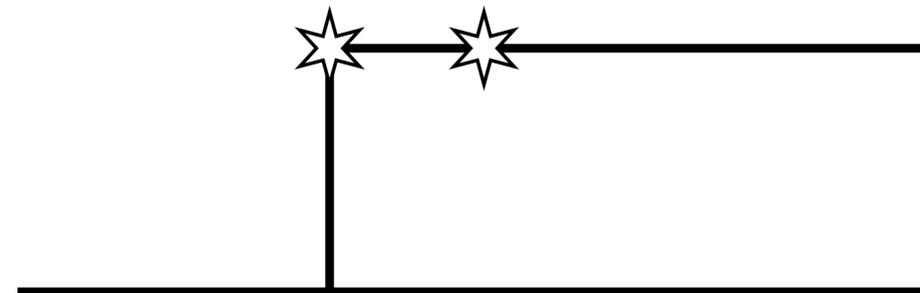
\mathcal{H} class of thresholds on finite interval $[0,1]$, $|\mathcal{H}| = N$

The Learner

The Learner predicts with the smallest threshold (conservative and consistent).

$t = 2$

$$\hat{h}_2 = [1/4, 1]$$



$t = 2$

$$(x_2, y_2) = (1/2, 1)$$

$$\hat{h}_2(1/8) = 1 = y_2 \quad \checkmark$$

Correct classification, no update

Notation: a function $h : \mathcal{X} \rightarrow \{0,1\} \leftrightarrow \text{supp}(h) = h^{-1}(1)$

Learning thresholds conservatively

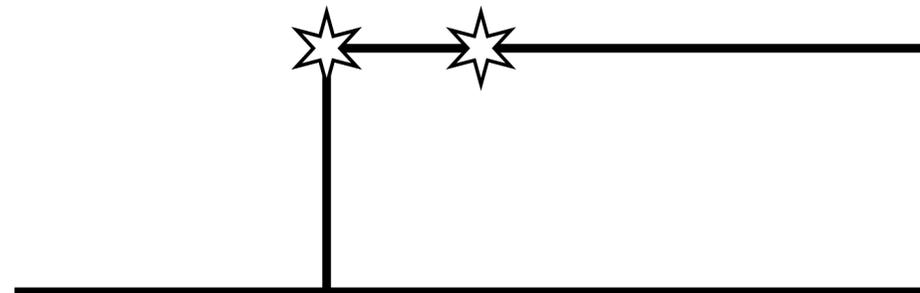
\mathcal{H} class of thresholds on finite interval $[0,1]$, $|\mathcal{H}| = N$

The Learner

The Learner predicts with the smallest threshold (conservative and consistent).

$t = 3$

$$\hat{h}_3 = [1/4, 1]$$



In general

Easy to see:

$$\mathcal{M}_T = \Theta(\min\{N, T\})$$



Learning thresholds NOT conservatively

\mathcal{H} class of thresholds on finite interval $[0,1]$, $|\mathcal{H}| = N$

Other Learner

Predict with threshold left of smallest threshold (not conservative, but consistent).

$t = 1$

$$\hat{h}_1 = [1/2, 1]$$



$t = 1$

$$(x_1, y_1) = (1/4, 1)$$
$$\hat{h}_1(1/4) = 0 \neq y_1 \quad \times$$

**A new positive point, which cannot be replay
→ Update classifier**

Learning thresholds NOT conservatively

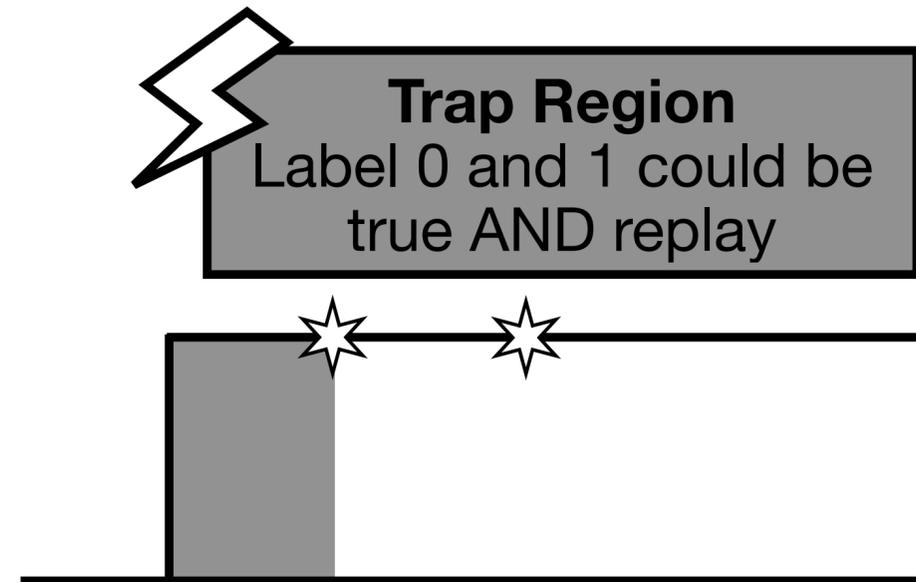
\mathcal{H} class of thresholds on finite interval $[0,1]$, $|\mathcal{H}| = N$

Other Learner

Predict with threshold left of smallest threshold (not conservative, but consistent).

$t = 2$

$$\hat{h}_2 = [1/8, 1]$$



$t = 2$

$$(x_2, y_2) = (1/2, 1)$$

$$\hat{h}_2(1/2) = 1 = y_2 \quad \checkmark$$

Correct classification, no update

Learning thresholds NOT conservatively

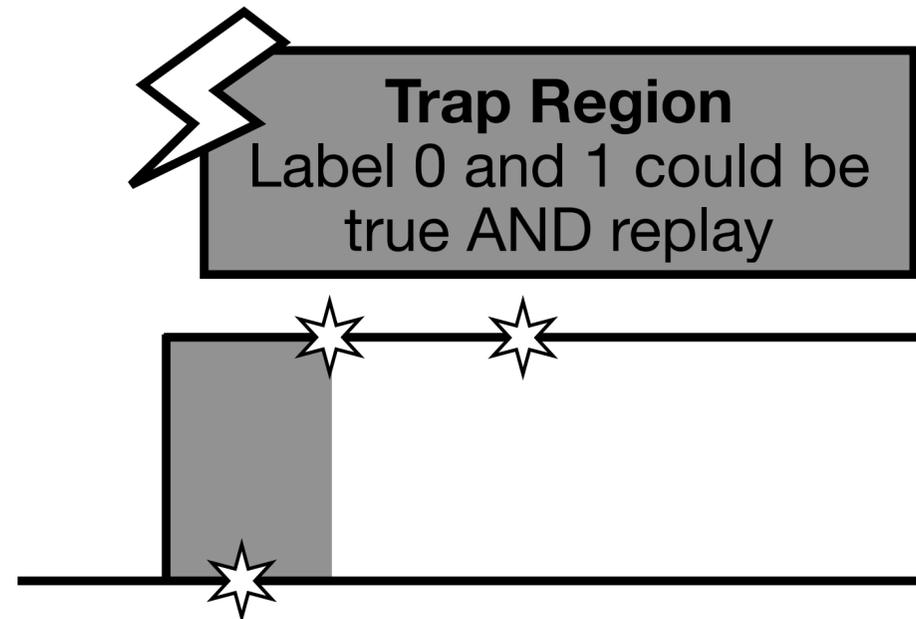
\mathcal{H} class of thresholds on finite interval $[0,1]$, $|\mathcal{H}| = N$

Other Learner

Predict with threshold left of smallest threshold (not conservative, but consistent).

$t = 2$

$$\hat{h}_2 = [1/8, 1]$$



$t = 2$

$$(x_2, y_2) = (3/16, 0)$$
$$\hat{h}_2(3/16) = 1 \neq y_2 \quad ?$$

Sample in Trap Region

→ No matter what the learner predicts $\exists f^*$ s.t. the learner errs

The Learner: Closure-Based Algorithm

Definition

Define the \mathcal{H} -closure as $\text{clos}_{\mathcal{H}}(Y) := \bigcap_{h \in \mathcal{H}: Y \subseteq h} h$ for any $Y \subseteq \mathcal{X}$

\mathcal{H} is intersection-closed* if $\forall S \subseteq \mathcal{H} : \bigcap_{h \in S} h \in \mathcal{H}$.

* (over arbitrary intersections)

The Learner

The Learner predicts $\hat{h}_t = \text{clos}_{\mathcal{H}}(\{x_i \mid y_i = 1, i \in \{1, \dots, t-1\}\})$

Measure of Complexity

Definition

Define Threshold dimension $\text{ThD}(\mathcal{H})$ largest k , such that
 $\exists x_1, \dots, x_k \in \mathcal{X}$ and $h_0, h_1, \dots, h_k \in \mathcal{H}$ with $h_i(x_j) = 1 \{j \leq i\}$

$$\text{ExThD}(\mathcal{H}) := \min_{f \subseteq \mathcal{X}} \text{ThD}(\overline{\mathcal{H}^f})$$

Complexity measure for Learning with Replay Adversary

Where $\overline{\mathcal{H}} := \{ \bigcap_{h \in S} h \mid S \subseteq \mathcal{H} \}$

and $\mathcal{H}^f := \{ h^f \mid h \in \mathcal{H} \}$ and $h^f := \{ x \in \mathcal{X} \mid h(x) \neq f(x) \}$

For \mathcal{H} thresholds on finite intervals:
 $\text{ExThD}(\mathcal{H}) = C \cdot \text{ThD}(\mathcal{H}) = C \cdot N$

For \mathcal{H} intersection-closed:
 $\text{ExThD}(\mathcal{H}) = C \cdot \text{ThD}(\mathcal{H})$

For \mathcal{H} general:
 $\text{ExThD}(\mathcal{H})$

Results I

Number of mistakes: $\mathcal{M}_T(\mathcal{A}) = \sum_{t=1}^T 1\{y_t \neq \hat{h}_t(x_t) \text{ and } y_t = f^*(x_t)\}$

$$\text{ExThD}(\mathcal{H}) = O(N)$$

Hypothesis Class	Adaptive Adv. \mathcal{M}_T	Stochastic Adv. $\mathbb{E}[\mathcal{M}_T]$
Thresholds on $[N]$	$\Theta(\min\{N, T\})$	$\Theta(\min\{N, \log T\})$

Results I

Number of mistakes: $\mathcal{M}_T(\mathcal{A}) = \sum_{t=1}^T 1\{y_t \neq \hat{h}_t(x_t) \text{ and } y_t = f^*(x_t)\}$

$\text{ExThD}(\mathcal{H}) = O(N)$

Hypothesis Class	Adaptive Adv. \mathcal{M}_T	Stochastic Adv. $\mathbb{E}[\mathcal{M}_T]$
Thresholds on $[N]$	$\Theta(\min\{N, T\})$	$\Theta(\min\{N, \log T\})$
Intersection-Closed \mathcal{H}	$\Theta(\text{ThD}(\mathcal{H}))$	$O(\min\{\text{ThD}(\mathcal{H}), d_{\text{VC}}(\mathcal{H}) \log T\})$ $\Omega(\min\{\text{ThD}(\mathcal{H}), \log T\})$

$\text{ExThD}(\mathcal{H}) = O(\text{ThD}(\mathcal{H}))$

Same upper bound would depend on $d_{\text{VC}}(\overline{\mathcal{H}^f})$

Gap for dependence on T

Results II

A finite \mathcal{H} is **properly learnable** in the replay setting,



there exists a f -representation of the class \mathcal{H} such that
 $\mathcal{H}^f := \{h^f \mid h \in \mathcal{H}\}$ is **intersection-closed**.

Open Questions/Extensions

- Open Question: Close gap for Stochastic Adversary
- Open Question: Existence of f -representation for \mathcal{H}^f intersection-closed?
- Extension: Limit the power of the Replay Adversary (e.g. Replay only after k times)

**Everyone seems to agree with me,
so I must be right**

Sara Germain