Covert Learning: How to learn with an untrusted intermediary

Ran Canetti, Ari Karchmer



The perfect experiment design!

0



The perfect experiment design!

00



The perfect experiment design!

0







Covert Verifiable Learning Goals

- Learning: If Eve reports the experiment results truthfully, Alice learns a good model from her experiments
- Concept-hiding: No (maybe, little) information about the molecular relationship is leaked
- Hypothesis-hiding: No information about Alice's hypothesis or domain knowledge used to influence the hypothesis is leaked
- Verifiability: If Eve tampers with the results, she cannot deceive Alice into learning a faulty model (here let Alice have private access to some random ``ground truth" experiments)

PAC-verification [GRSY20]

- Learning: If Prover reports the experiment results truthfully, <u>Verifier</u> learns a good model
- **Concept-hiding**
- Hypothesis-hiding
- Verifiability: If Prover tampers with the results, she cannot deceive Verifier into learning a faulty model (here let <u>Verifier have private access to some</u> random ``ground truth" examples)

Cryptographic Sensing [KOS19]

- Learning: <u>Sensor</u> obtains an exact model/object from the experiments
- **Secrecy:** Object remains hidden to an adversary provided that it is drawn from a distribution of sufficient minimum entropy (entropic security)
- Hypothesis-hiding
- Verifiability



More Related Work Or: what this talk is not about

- Differentially private learning [KLNRS11]
 - Data privacy not *learner privacy*

- Verifiable computation [K92, M00]

 - Requires computational entity

Verify computation steps — not good learning outcomes

Example Application: Drug Discovery

- How does A prevent trade secrets being revealed by experimental design?
- How does A prevent the lab from "double-selling" data to competitor?
- How does A prevent the lab from returning faulty data (say, given access) to some random ground truth examples)?

Example Application: Drug Discovery

- Cannot execute protocol with the lab
- results at some point no way around this
- We must design an algorithm that interacts only with "Nature"
- Q+A only

• A must reveal some experiments at some point, while the lab must read the

More applications

- learning but can't execute a protocol
- E.g. Model extraction attacks Client tries to extract model by using active learning strategies



• Our work could be useful in any setting where we want to hide what we are



Contributions

- verifiability guarantees
 - (agnostic)
 - Covert Verifiable Learning even with **no privately accessed** examples
 - **verifiability** (perfect, statistical)

The **Covert Learning** model, which augments PAC learning with MQ with hiding guarantees

Covert Learning algorithms for parity functions and decision trees (agnostic)

The **Covert Verifiable Learning** model, which augments PAC learning with MQ with hiding and

Covert Verifiable Learning algorithms for parity functions and decision trees

Covert Verifiable Learning algorithms for O(log n)-juntas with strong privacy and

Contributions Today

- verifiability guarantees
 - (agnostic)
 - Covert Verifiable Learning even with **no** privately accessed examples
 - verifiability (perfect, statistical)

The **Covert Learning** model, which augments PAC learning with MQ with hiding guarantees

Covert Learning algorithms for parity functions and decision trees (agnostic)

The **Covert Verifiable Learning** model, which augments PAC learning with MQ with hiding and

Covert Verifiable Learning algorithms for parity functions and decision trees

Covert Verifiable Learning algorithms for O(log n)-juntas with strong privacy and

Itinerary

- Establish the Covert Learning model
- Covert Learning for noisy parities
- Covert Learning for heavy Fourier coefficients
- Establish the Covert Verifiable Learning model
- Transform Covert Learning into Covert Verifiable Learning
- Concluding questions and remarks

Alice's task? Learning with membership queries

- The structure-activity relationship is a boolean function $f: \{feature, \neg feature\}^n \rightarrow \{bind, \neg bind\}$
- class
- Experiments \iff synthesized membership queries

• Alice's initially hypothesized set of structure-activity models \iff A hypothesis



PAC Learning terminology

- For every $n \in \mathbb{N}$, A hypothesis class \mathscr{H}_n is a class of functions $\subseteq \{f: \{0,1\}^n \to \{0,1\}\}$
- For every $n \in \mathbb{N}$, \mathcal{D}_n is a **concept class**, where every $D_n \in \mathcal{D}_n$ is a distribution over $\{0,1\}^n \times \{0,1\}$
- For an example $(x, y) \sim D_n$, x is an input, and y is a label
- A loss function (w.r.t. a concept D_n) \mathscr{L}_{D_n} : $\mathscr{H}_n \to [0,1]$ quantifies how good a hypothesis $h \in \mathcal{H}$ is at approximating a concept. For instance misclassification error: $\mathcal{L}_D(h) = \Pr_{(x,y)\sim D} \left[h(x) \neq y \right]$

Agnostic PAC Learning

There is an algorithm A receiving enoug For any concept $D \in \mathcal{D}_n, \epsilon, \delta > 0$, A outputs $h \in H$ s.t. $\Pr\left[\mathscr{L}_D(h) \leq \mathscr{L}_D\right]^{\bullet}$ $\mathscr{L}_D(H) = \inf_{h \in H} \mathscr{L}_D(h)$

A hypothesis class H is agnostic PAC learnable (w.r.t. concept class \mathscr{D}_n) if (roughly)

- With membership queries: there is also an oracle that allows the learning algorithm to directly query the concept
- Agnostic: *H* need not contain a perfect hypothesis for the concept
 - *H* has a big impact on the resulting function which is output by $A \longrightarrow$ important choice, whatever H is



The Covert Learning Model

- $\mathscr{C}_n = \{H_n^i\}_{i \in [m]}$ a collection of hypothesis classes
- Learner input: $\epsilon, \delta < 1, H_n \in \mathscr{C}_n$
- Learner queries concept
- Learner outputs: for any $D_n \in \mathscr{D}_n$, a hypothesis $h \in H_n$ such that $\Pr_A \left[\mathscr{L}_D(h) \le \mathscr{L}_D(H) + \epsilon \right] \ge 1 - \delta$



The Covert Learning Model

- Adversary attempts to deduce information about either the concept $D_n \in \mathscr{D}_n$ or the hypothesis class $H_n \in \mathscr{C}_n$
- **Objection**: some concept classes we cannot hide, e.g. the constant function
- Adversary could apply Occam learning



The Covert Learning Model

There exists a p.p.t. simulator such that for every $H_n \in \mathscr{C}_n, D_n \in \mathscr{D}_n$ $Sim(S) \approx \{interaction transcript\}$ where $S \sim D_n$ (random examples of n



"Zero-knowledge-like" in the presence of public data set of random examples

Many learning problems thought to be computationally hard to learn from random examples (e.g. decision trees, small depth circuits, parities with noise)

Random examples are the minimum leakage, not just a result of the definition



Trivial example Covert Learning is easy when PAC learning is possible

- Every $H_n \in \mathscr{C}_n$ is efficiently PAC learnable
- Covert Learning algorithm:
 - Request a sufficiently large set of random examples from the oracle. •
 - Run the PAC learning algorithm.
 - The simulator returns the random examples given as input.
- Examples: Constant term DNFs, parities (in the absence of noise).

Itinerary

- Establish the Covert Learning model
- Covert Learning for noisy parities
- Covert Learning for heavy Fourier coefficients
- Establish the Covert Verifiable Learning model
- Transform Covert Learning into Covert Verifiable Learning
- Concluding questions and remarks



Parities with no noise are trivial, so what about the noisy case?

Focus on the concept hiding guarantee



- Search LPN: find s. \bullet
- subexponentially-hard

- **LPN distribution** [BFKL92]:
 - Sample $s \in \{0, 1\}^n$, $e \in \{0, 1\}$ from Bernoulli r.v. with mean
- Sample $a \in \{0, 1\}^n$ uniformly at random

Return $a, \langle a, s \rangle \oplus e$ (s is persistent over each example)

Decision LPN: distinguish $a, \langle a, s \rangle \oplus e$ from uniformly random. Both are thought to be



Covert Learning for noisy parities Low-noise LPN



- Search LPN: find s.
- **Decision LPN:** distinguish $a, \langle a, s \rangle \oplus e$ from subexponentially-hard

Is there a covert learning algorithm for learning s?

- Low-noise LPN distribution [Ale03]:
- Sample s $\in \{0, 1\}^n$, $e \in \{0, 1\}$ from Bernoulli r.v. with mean $p = 1/\sqrt{n}$
- Sample $a \in \{0, 1\}^n$ uniformly at random

Return a, $\langle a, s \rangle \oplus e$ (s is persistent over each example)

Decision LPN: distinguish $a, \langle a, s \rangle \oplus e$ from uniformly random. Both are thought to be



(actually, there exists a Covert Learning algorithm for noisy parities **unconditionally**. Recall previous example.)

Theorem: (informal) If low-noise LPN is hard, then there exists a Covert Learning algorithm for parities with respect to the low-noise LPN distribution









"Masked queries"









Random queries

query transcript

 \mathcal{S}

2n



Repetition and majority voting to decode bit-by-bit

Can we covertly learn more interesting concepts?

Itinerary

- Establish the Covert Learning model
- Covert Learning for noisy parities
- Covert Learning for heavy Fourier coefficients
- Establish the Covert Verifiable Learning model
- Transform Covert Learning into Covert Verifiable Learning
- Concluding questions and remarks



Boolean Fourier Analysis Review

- Let $f: \mathbb{F}_2^n \to \mathbb{R}$. In particular, we care about functions $f: \mathbb{F}_2^n \to \{-1, 1\}$
- For $S \subseteq [n]$, we define $\chi_S : \mathbb{F}_2^n \to \mathbb{R}$ by χ_S

The Fourier expansion of $f: \mathbb{F}_2^n \to \mathbb{R}$ is f(x)

coefficient on S

- Every function can be uniquely represented by its Fourier expansion
- We say that $\hat{f}(S)$ is "heavy" if $\hat{f}(S) \ge 1/poly(n)$, and the degree of $\hat{f}(S)$ is |S|

$$g(x)_{\widehat{f(S)}} = (-1)^{\sum_{i \in S} x_i}$$

$$f(x) = \sum_{S \subseteq [n]} \widehat{f}(S) \chi_S(x) \text{ where } \widehat{f}(S) \text{ is the Fourier}$$

Why Learn Fourier Coefficients?

- By definition, $\hat{f}(S) = \mathbb{E}_x[f(x)\chi_S(x)]$
- If we know S such that $\hat{f}(S) \ge \tau$, then we have a hypothesis that agrees with f on $1/2 + \tau/2$ fraction of inputs (just take χ_S)
- Good for learning w.r.t. uniform distribution
- Finding all heavy Fourier coefficients of degree O(log n) is strong enough to efficiently learn poly(n)-size decision trees. Not known in plain PAC learning model
- Our goal: find (covertly) the O(log n) Fourier coefficients of a function



 $f: \{0,1\}^n \to \{-1,1\}$



- Membership query q bypasses distribution over inputs and returns f(q)
- **Goal:** find Fourier coefficients of f while satisfying privacy guarantees

Covert Learning of low degree Fourier coefficients Highlighting hypothesis-hiding

- Learner has a hypothesis (due to expert domain knowledge) that all heavy Fourier coefficients of a function are contained in a subset $T \subseteq [n]$.
- Valuable: this knowledge could improve the efficiency of a normal learning algorithm (e.g. KM) algorithm).
- **Leaky:** repeated restrictions of KM queries clearly reveals information on T
- How to hide this expert domain knowledge?

hypothesis class in indexed by $T \subseteq [n]$

Covert Learning then hides all information about T

Theorem: (informal) Under the subexponential LPN assumption, the collection of all low-degree Fourier subsets is covertly learnable.

Answer: do Covert Learning for collection of hypothesis classes, where each



Covert Learning of low degree Fourier coefficients Squared-log entropy assumption

- Extend the technique from the previous result: "masked queries"

hardness with *n*-bit secrets of just $\Theta(log^2(n))$ entropy (due to [YZ16])

Need decisional LPN variant (sqlogLPN) that provides super-polynomial

• Variant implied by subexponential LPN: $2^{\sqrt{n}}$ -hard with $2^{\sqrt{n}}$ samples [YZ16]







- Let otp ~ sqlogLPNⁿ
- For query $q \in \{0,1\}^n$, let $\hat{q} = q \oplus \text{otp}$ ("masked query" technique)
- What can we say about $f(\hat{q})$?
- $\mathbb{E}_{\phi,q,\hat{q}}\left[\phi(f(\hat{q})) \cdot \chi_k(q)\right] \geq \Omega(\tau n^{-c})$ In particular, if $\hat{f}(S) \ge 1/poly(n)$, then $\mathbb{E}[\phi(f(\hat{q}))\chi_k(q)] \ge 1/poly(n)$.

• Let ϕ be a randomized mapping defined $\phi(f(\hat{q})) := f(\hat{q}) \cdot \chi_s(r)$, where s is the sale of PN secret used on the other for a and r is a random *n*-bit string **Lemma:** If $k \subseteq [n]$ s.t. $|k| = O(\log n)$ and $\hat{f}(k) \ge \tau$, then there exists a small constant c such that

The lemma provides a "Goldreich-Levin type" environment

$$\mathbb{E}_{\phi,q,\hat{q}} \left[\phi(f(\hat{q})) \cdot \chi_k(q) \right] \geq \\ \text{Noisy Predictor}$$

 $\geq \Omega(\tau n^{-c})$

• The Goldriech-Levin reduction/algorithm gives a method for extracting a k s.t. $\hat{f}(k) \ge \tau$ in time $poly(n, \frac{1}{\tau})$ (w.h.p.)

- Run the Goldreich-Levin algorithm in the "masked query regime"
- Can run Goldreich-Levin on some subset of interesting indices
 - This encodes secret hypothesis ${\cal T}$
- Pseudorandomness hides ${\cal T}$
- At the end we obtained $S \in T : |S| \le O(\log n) , \hat{f}(S) \ge 1/poly(n)$





Covert Learning of poly(n) size decision trees

Sufficient to produce an approximation to f that is competitive with the optimal poly(n) -size decision tree approximating f.

Theorem: (informal) Under the subexponential LPN assumption, the collection of all subsets of decision trees is covertly learnable.



Itinerary

- Establish the Covert Learning model
- Covert Learning for noisy parities
- Covert Learning for heavy Fourier coefficients
- Establish the Covert Verifiable Learning model
- Transform Covert Learning into Covert Verifiable Learning
- Concluding questions and remarks



Covert Verifiable Learning A malicious adversary

- Interactive between the learner and adversarial intermediary (AI)
- AI monitors access to the membership oracle
- The learner request queries from oracle, but responses my be corrupted by AI



Covert Verifiable Learning The verifiability guarantee

except with small probability (similar to soundness)

- Allow the learner to abort
- Access to set of uniformly random "ground truth" examples

Inputs: *E*

We say that verifiability is computational if I is p.p.t..

- Goal: If AI chooses to corrupt, the learner should not output a faulty hypothesis
 - Verifiability for learning algorithm A for collection \mathscr{C}_n

$$H_n \in \mathscr{C}_n, \epsilon, \delta, \mathcal{S}$$

For any $D_n \in \mathcal{D}_n$, any $H_n \in \mathcal{C}_n$, $\mathcal{S} \sim D_n^m$, any Al I which corrupts oracle responses,

$$\Pr\left[\mathscr{L}(A) > \mathscr{L}(H_n) + \epsilon \mid A \neq \text{abort}\right] < \delta$$



Covert Verifiable Learning How to extend privacy?

Real:

Adversary (a distinguisher) chooses $H_n \in \mathscr{C}_n, D_n \in \mathscr{D}_n, \varepsilon, \delta$

 \mathcal{S} is drawn from D_n



Learner gets $\mathcal{S}, H_n, \epsilon, \delta$, Al gets ϵ, δ

Learner tries to learn H_n with access to oracle. Al sees the learner's queries and responses and is given the chance to modify the responses. At the end of interaction, Al outputs a string denoted by real \mathcal{D}_A^n

Output
$$\left(H_n, \epsilon, \delta, \mathcal{S}, \operatorname{real}_{A, I}^{\mathscr{D}_n}\right)$$

Ideal:

Adversary (a distinguisher) chooses $H_n \in \mathscr{C}_n, D_n \in \mathscr{D}_n, \varepsilon, \delta$

 \mathcal{S}' is drawn from D_n



Sim gets access to $\mathcal{S}, \mathcal{S}', H_n, \epsilon, \delta$, where \mathcal{S} is the set of examples given to the real learner. Al gets ϵ, δ

Sim "interacts" with the oracle. Al "views" the queries and responses and may change the responses. The simulator outputs a string, $ideal_I^{Sim}$

Output
$$\left(H_{n}, \epsilon, \delta, \mathcal{S}, \text{ideal}_{I}^{Sim}\right)$$



Itinerary

- Establish the Covert Learning model
- Covert Learning for noisy parities
- Covert Learning for heavy Fourier coefficients
- Establish the Covert Verifiable Learning model
- Transform Covert Learning into Covert Verifiable Learning
- Concluding questions and remarks



Augmenting Covert Learning with Verifiability

Claim Suppose that on one iteration, w.p. $1/2 + \epsilon$, the AI can corrupt at least 1 oracle response without causing the learner to abort. Then the squared-log LPN problem is efficiently solved with advantage ϵ .

- Follows from the fact that AI is always caught lying in case c = 1

- Idea: Alternate "testing" and "learning"
- Eventually AI is caught or one round is an uncorrupted "learning" phase
- Verifiability follows from learning guarantee of basic Covert Learning algorithm

```
Input: S = random examples
Repeat r times:
   Flip coin c
   If c = 0
        covert_learn()
   If c = 1
        Query section of S
        abort if results
        inconsistent
```

Theorem: (informal) Under the subexponential LPN assumption, the collection of all subsets of low-degree Fourier subsets is verifiably covertly learnable.

Theorem: (informal) Under the subexponential LPN assumption, the collection of all subsets of decision trees is verifiably covertly learnable.

Next Questions

Little questions

- numbers)?
- relax the fully agnostic setting)

Big question

lacksquare(Verifiable) Learning algorithms?

Can we obtain more Covert Learning algorithms (e.g. new learning problems, large fields, real

In the verifiability setting, can we remove the assumption of privately accessed dataset? (Yes, if we

When is there a general *compiler* for turning plain PAC learning with MQ algorithms into Covert



