

## Lecture 10: Negative Results

Source: Lecture notes by  
Aaron Roth and Adam Smith

Lecturer: Uri Stemmer

In our course, we saw several methods that allow answering adaptive statistical queries. In particular, we saw:

- A computationally efficient algorithm that answers  $k$  adaptive queries using a sample of size

$$n \gtrsim \sqrt{k}$$

- A computationally inefficient algorithm that answers  $k$  adaptive queries using a sample of size

$$n \gtrsim (\log k)^3 \cdot (\log |X|)^2$$

Is there a computationally efficient algorithm that achieves significantly better results than the algorithm we saw? Today we will see that, in general, the answer is no.

**The idea:** We will present an attacker (analyst) who interacts with an efficient mechanism that answers queries. Using (approximately)  $n^2$  queries, our attacker will succeed in reconstructing almost the entire sample held by the mechanism. As we have already seen in our course, once the attacker reconstructs the sample, he can easily find a "bad" query.

For this, we will need to become familiar with the following cryptographic tool:

### **Fingerprinting codes (FPC)**

This is a cryptographic tool designed to try and prevent the illegal distribution of digital content. Suppose Netflix releases a new episode of a certain series to its subscribers every day. Netflix wants to prevent malicious users from downloading the episode and uploading it to the internet. To attempt to prevent this, they can embed invisible watermarks in the episodes they release, so that each subscriber receives a unique copy with different markings. Now, if Netflix finds the episode online, they can identify which subscriber distributed it.

But what happens if several subscribers decide to collaborate, compare their copies to try and remove the markings, and then upload a "processed" copy to the internet?

This problem is what FPC aims to solve. The idea relies on the following assumption:

- If two malicious users attempt to combine their copies to remove the markings, they will only be able to identify locations where their markings differ. However, if at a certain location they both have the exact same marking, they won't notice there is a marking there and will not remove it.

- Similarly, if a group of  $n$  malicious users tries to collaborate and remove the markings from their copies, they will not be able to identify markings that are common to all of them. That is, if at a certain point in the file they all have the exact same marking, that marking will remain in the "processed" copy they try to create.

Under these assumptions, FPC enables Netflix to identify the malicious users.

### **In more detail:**

- There are  $N$  users, out of which  $n$  users have decided to collaborate and upload episodes to the internet.
- Each day:
  - Netflix broadcasts a new episode with different markings for each user.
  - The  $n$  malicious users collaborate to try to remove or alter the markings in their copies, under the restriction that markings identical across all  $n$  malicious users cannot be removed. They then upload the processed copy they created to the internet.
    - **Formally, we will (seemingly) require even less from the malicious users: markings that are identical across all  $N$  users cannot be removed. (In particular, if markings identical across the  $n$  malicious users cannot be removed, then markings identical across all  $N$  users clearly cannot be removed either.)**
- Netflix observes the processed episode on the internet.
- After  $\ell$  such iterations, Netflix can identify (almost) all of the malicious users.

Formally, we consider the following game between an adversary  $\mathcal{P}$ , who controls  $n$  out of  $N$  users, and an algorithm  $\mathcal{F}$  (this is the FPC algorithm), which attempts to identify the  $n$  users controlled by the adversary.

#### **Algorithm IFPC( $\mathcal{P}, \mathcal{F}$ )**

- (1) Denote  $\text{ALIVE}^1 = [N]$
- (2) The Adversary  $\mathcal{P}$  selects a subset of users  $S^1 \subseteq [N]$  (unknown to  $\mathcal{F}$ )
- (3) For  $j = 1$  to  $\ell$  do
  - a. Algorithm  $\mathcal{F}$  outputs a vector  $c^j \in \{-1, 1\}^N$
  - b. Let  $c_{S^j}^j$  be the restriction of  $c^j$  to coordinates in  $S^j$ . This is given to  $\mathcal{P}$
  - c. The adversary  $\mathcal{P}$  outputs  $a^j \in \{-1, 1\}$ . This is given to  $\mathcal{F}$
  - d. Algorithm  $\mathcal{F}$  accuses a set of users  $I^j \subseteq [N]$ .
  - e. Let  $\text{ALIVE}^{j+1} = \text{ALIVE}^j \setminus I^j$ , and let  $S^{j+1} = S^j \setminus I^j$
- (4) The Adversary  $\mathcal{P}$  wins if it was consistent throughout the game. Formally,  $\mathcal{P}$  is consistent if for every  $j \in [\ell]$  there exists an index  $i \in \text{ALIVE}^j$  such that  $a^j = c_i^j$

**Definition 1:** Algorithm  $\mathcal{F}$  is an FPC with parameters  $n, N, \ell, \varepsilon$  if for any adversary  $\mathcal{P}$ , the following conditions hold:

$$(1) \quad \Pr \left[ \begin{array}{c} \text{A good user is} \\ \text{accused as malicious} \end{array} \right] \leq \varepsilon$$

$$(2) \quad \Pr[\text{the adversary } \mathcal{P} \text{ wins}] = \Pr \left[ \begin{array}{c} \mathcal{P} \text{ is consistent throughout} \\ \text{all of the execution} \end{array} \right] \leq \varepsilon$$

where the probabilities are over the execution of the game *IFPC*

### What's happening here?

Throughout the execution, we accuse users of being malicious, and they are "removed from the game" (in particular, they are removed from the set  $\mathcal{S}^j$ ). We require that, with high probability, no honest user is accused of being malicious.

Additionally, we require that at some point during the execution, because malicious users are progressively removed, the adversary  $\mathcal{P}$  will no longer be able to satisfy the consistency requirement. Since we assume that attackers cannot produce "processed episodes" without meeting the consistency requirement, this means the attackers will no longer be able to produce processed episodes.

**Theorem 2:** For every  $1 \leq n \leq N$ , there exists an FPC with parameters  $n, N, \ell, \varepsilon$ , where  $\varepsilon = \frac{1}{N}$  and  $\ell = O(n^2 \cdot \log N)$

### Partial Proof of Theorem 2:

We provide a partial proof of the theorem, using the parameter  $\ell = \mathbf{poly}(n)$  instead of  $\ell \approx n^2$ .

Our current goal is to construct an algorithm  $\mathcal{F}$  for the aforementioned IFPC game. Recall that, intuitively, during the game, algorithm  $\mathcal{F}$  attempts to "accuse" more and more malicious users, until (hopefully) it accuses all of the malicious users.

We will begin by designing an algorithm (denoted  $\mathcal{F}_1$ ) whose goal is to accuse a single malicious user from the set of malicious users. Then, we will iteratively use this algorithm to identify additional malicious users.

**Algorithm  $\mathcal{F}_1$**

(1) Let  $Z$  be the following matrix:

$$\begin{pmatrix} \mathbf{0} \text{ block} & \mathbf{1st block} & \mathbf{2dn block} & \mathbf{nth block} \\ 000 \dots 000 & 111 \dots 111 & 111 \dots 111 & 111 \dots 111 \\ 000 \dots 000 & 000 \dots 000 & 111 \dots 111 & 111 \dots 111 \\ 000 \dots 000 & 000 \dots 000 & 000 \dots 000 & \dots 111 \dots 111 \\ \vdots & \vdots & \vdots & 111 \dots 111 \\ 000 \dots 000 & 000 \dots 000 & 000 \dots 000 & 111 \dots 111 \end{pmatrix}$$

In this matrix, there are  $(n + 1)$  types of columns, with  $\tilde{O}(n^2)$  copies of each type of column. Therefore, the total number of columns in this matrix is  $\tilde{O}(n^3)$ .

- (2) Let  $Y$  be a matrix obtained by applying a random permutation to the columns of matrix  $Z$ .
- (3) For  $\ell = \tilde{O}(n^3)$  rounds, algorithm  $\mathcal{F}_1$  outputs (one after the other) the columns of the matrix  $Y$  as the watermark vectors it is supposed to output in the IFPC game
- (4) Let  $\vec{a} \in \{0,1\}^\ell$  denote the vector of responses that algorithm  $\mathcal{F}_1$  received from the adversary during step (3).
- (5) For each column type/block  $i$ , let  $\mathbf{count}_i$  denote number of columns of type  $i$  for which the adversary returned the response  $\mathbf{1}$ . Note that this is a number between  $\mathbf{0}$  and  $\tilde{O}(n^2)$
- (6) Return an index  $i \in [n]$  (i.e., the accused user) for which the following holds:

$$\mathbf{count}_i > \mathbf{count}_{i-1} + \tilde{O}(n)$$

First, we will show that step (6) is well-defined, meaning that an index  $i$  as described indeed exists. Recall that we assume the adversary must be consistent (otherwise, we win, and we don't care). Thus:

- For all columns of type  $\mathbf{0}$ , the adversary must return the response  $\mathbf{0}$ .
- For all columns of type  $\mathbf{n}$ , the adversary must return the response  $\mathbf{1}$ .

This implies that:

$$\begin{aligned} \mathbf{count}_0 &= \mathbf{0} \\ \mathbf{count}_n &= \tilde{O}(n^2) \end{aligned}$$

Therefore, since there are  $n + 1$  types of columns, there must be an index  $i$  that satisfies the condition from step (6) in the algorithm.

Next, we will show that, with high probability, an index  $i$  that satisfies this condition must correspond to a "malicious" user. To this end, let  $i$  be an index of an "honest" user. We will show that the probability of  $i$  satisfying the condition is negligible.

To see this, note that since user  $i$  is "honest," then the adversary cannot distinguish between columns from block  $(i - 1)$  and columns from block  $i$ , because the only row that differentiates them is row  $i$ , and the adversary never sees the entries of this row.

Thus, intuitively, the adversary's responses must behave similarly across these two blocks, and there cannot be a "jump" of  $\approx n$  in the counts.

More formally: Fix the adversary's responses and the entire permutation, except for the partition of the columns (after the permutation) between those belonging to block  $(i - 1)$  and those belonging to block  $i$ . Since the adversary does not see row  $i$ , this partition between the columns of  $(i - 1)$  and  $i$  remains completely random, even given all the adversary's responses.

Let  $\vec{a}_i$  denote the set of the adversary's responses corresponding to the columns from blocks  $(i - 1)$  and  $i$ . Now the random variables  $\text{count}_{i-1}$  and  $\text{count}_i$  have the same marginal distribution: each of them is the sum of a randomly chosen half of the bits in  $\vec{a}_i$ . Therefore, these two variables have exactly the same expectation. Furthermore, by the Hoeffding bound, with overwhelming probability, the sums will not deviate by more than  $\tilde{O}(n)$  from their expectation (since each sum involves  $\tilde{O}(n^2)$  bits). Hence, the probability of a large "jump" occurring for any "honest" user  $i$  is negligible.

### Where do we stand in the proof of Theorem 2?

We have shown that using  $\tilde{O}(n^3)$  rounds, we can accuse one "malicious" user. We then run this algorithm  $n$  times. Each run is conducted only on the set of users who are still "alive", and in each iteration we remove one user from the set of "alive" users. In total, after  $\tilde{O}(n^4)$  rounds, we will identify all the malicious users (or the adversary will "lose" somewhere along the way, meaning they will fail to maintain consistency, in which case we also win).

Q.E.D. (Theorem 2)

**Now, let us see how FPC can be used to derive negative results for the problem of answering adaptively chosen statistical queries. Specifically, we will show that efficient mechanisms cannot answer more than  $n^2$  queries using a sample of size  $n$ .**

Simplifying assumption: The query-answering mechanism satisfies the following property:

**Definition:** A query-answering mechanism is called "natural" if, when it holds a sample  $S \in X^n$  and receives a query  $q: X \rightarrow \{\pm 1\}$ , its output does not depend on the values of the query outside the sample. That is, for any two queries  $q$  and  $q'$  that agree on all the elements of the sample  $S$ , the answer returned by the algorithm is identical (or distributed identically).

**Theorem:**

Let  $M$  be a natural mechanism that answers  $k$  adaptive queries over a domain of size  $|X| \geq 2000n$  using a sample of size  $n$ . Then,  $n = \Omega(\sqrt{k})$ .

**Proof Sketch:**

Consider the following attack. The adversary defines a target distribution  $\mathfrak{D}$ , from which the mechanism  $M$  receives a sample (the adversary does not see the sample). The adversary then asks (approximately)  $n^2$  adaptive queries. For at least one of these queries, the mechanism will fail to answer accurately.

| <b>Attack against a natural mechanism <math>M</math></b>   |
|--|
| <p><i>Setup:</i></p> <ul style="list-style-type: none"> <li>• Given parameter <math>n</math>, let <math>N = 2000n</math> and let <math>\mathfrak{D}</math> be the uniform distribution over <math>\{1, \dots, N\}</math></li> <li>• Let <math>S</math> be a sample containing <math>n</math> iid elements from <math>\mathfrak{D}</math>. The sample <math>S</math> is given to the mechanism <math>M</math> (but not to the adversary)</li> </ul> <p><i>The Attack:</i></p> <ol style="list-style-type: none"> <li>(1) Initialize an <math>(n, N, \ell, \varepsilon)</math>-FPC algorithm <math>\mathcal{F}</math>, where <math>\ell = \tilde{O}(n^2)</math> and <math>\varepsilon \leq \frac{1}{N}</math></li> <li>(2) Let <math>T^0 = \emptyset</math></li> <li>(3) For <math>j = 1</math> to <math>\ell</math> do             <ol style="list-style-type: none"> <li>a. Let <math>c^j \in \{\pm 1\}^N</math> be the vector chosen by <math>\mathcal{F}</math></li> <li>b. Define the query <math>\hat{q}^j</math> such that <math>\hat{q}^j(i) = \begin{cases} c_i^j &amp; , \text{ if } i \notin T^{j-1} \\ 0 &amp; , \text{ else} \end{cases}</math></li> <li>c. Ask <math>\hat{q}^j</math> to the mechanism <math>M</math> and obtain answer <math>a^j</math>. Round <math>a^j</math> to <math>\bar{a}^j \in \{\pm 1\}</math></li> <li>d. Give <math>\bar{a}^j</math> to <math>\mathcal{F}</math> and let <math>I^j \subseteq [N]</math> be the set of accused users. Let <math>T^j \leftarrow T^{j-1} \cup I^j</math></li> </ol> </li> </ol> |

First, note that by the properties of the FPC, with probability at least  $1 - \frac{1}{N}$ , in every round  $j$ , it holds that  $I^j \subseteq S$ , since we assume that false accusations occur with probability at most  $\frac{1}{N}$ .

In this case, for every  $j$ , it holds that  $|T^j| \leq n$ . Let us assume this is indeed the case

Additionally, according to the guarantees of FPC, with probability at least  $1 - \frac{1}{N}$ , there exists a round  $j$  where the response  $\bar{a}^j$  is inconsistent. In particular, this means that at time  $j$ , the vector  $c^j$

must be either identically  $-1$  or identically  $1$  (since if this vector contains both  $-1$ 's and  $1$ 's, then by definition,  $\bar{\mathbf{a}}^j$  cannot be inconsistent.)

This means that at time  $j$ , one of two cases occurs:

Case 1:  $\bar{\mathbf{a}}^j = -1$  but  $\mathbf{c}^j \equiv 1$ .

Case 2:  $\bar{\mathbf{a}}^j = 1$  but  $\mathbf{c}^j \equiv -1$ .

In either case, this implies that the response returned by the mechanism is very far from the expected value of the query in that round. Specifically, in Case 1,

$$\hat{q}^j(\mathfrak{D}) \geq 1 - \frac{|T^j|}{N} \geq 1 - \frac{n}{N} \geq 1 - \frac{1}{2000}$$

And in Case 2,

$$\hat{q}^j(\mathfrak{D}) \leq -1 + \frac{|T^j|}{N} \leq -1 + \frac{n}{N} \leq -1 + \frac{1}{2000}$$

This shows that the rounded answer  $\bar{\mathbf{a}}^j$  is very far from the "true answer"  $\hat{q}^j(\mathfrak{D})$ , to the extent that the unrounded answer (the answer provided by the mechanism  $\mathbf{M}$ ) is also very inaccurate.

q.e.d.

**Question:** Where did we use the simplifying assumption that the mechanism is "natural"?

**Answer:** The guarantees of FPC hold only under the rules of the IFPC game. An important point is that in the IFPC game, the adversary  $\mathcal{P}$  only observes the values  $\mathbf{c}_{S^j}^j$ , i.e., only the coordinates of the vector  $\mathbf{c}^j$  corresponding to the elements in  $S$  that have not yet been accused as malicious.

In contrast, in the last attack we presented, we give the mechanism  $\mathbf{M}$  the entire vector  $\mathbf{c}^j$  (this vector defines the statistical query we present to the mechanism). In general, this falls outside the rules of the IFPC game and breaks its guarantees. However, under our simplifying assumption—that the mechanism is "natural"—we know that it "ignores" the values of this query outside the sample, i.e., it "does not look" at the coordinates of  $\mathbf{c}^j$  it is not supposed to see. Hence, the attack works.

**Note:** We did not make any assumptions here about the computational power of the mechanism. In other words, what we have shown is that a query-answering mechanism, if it is "natural," cannot answer more than  $n^2$  queries—even if it is not computationally limited.

**Note:** To eliminate this simplifying assumption, we can use encryption to hide from the mechanism the content of the query in the coordinates it is not supposed to see. This effectively forces the mechanism to behave like a "natural" mechanism, and the attack still works. (However, in this case, we would need to assume that the mechanism is computationally limited, as otherwise, it could break the encryption and access the hidden coordinates.)

# Negative result for adaptive streaming

Let us see how to use the above impossibility results to also derive a negative result for **adaptive streaming**. The difference from the impossibility result we previously saw for adaptive streaming is that now we will obtain a negative result for a streaming problem defined by a real-valued function (it is going to be a negative result for an “estimation” problem rather than a “search” problem).

## The Plan:

We will present a streaming problem that can be easily solved in the **non-adaptive** model but cannot be solved with small space in the **adaptive** model. Specifically, we will show that if this problem can be solved effectively in the adaptive model, it would imply the existence of an efficient mechanism that answers adaptive queries too effectively, which we know is impossible.

## The SADA Problem:

- Each update  $s_i$  in the stream is either a point  $s_i = p_i \in X$  or a function  $s_i = h_i: X \rightarrow \{0, 1\}$ .
- **Goal:** At each step  $i$ , after receiving the update  $s_i$ , return an approximation of the average of the most recent function received over the multiset containing all the points received in the stream so far.

## Formally:

- Let  $S = (s_1, s_2, \dots, s_i)$  be a prefix of the stream.
- Let  $f_S$  denote the most recent function appearing in  $S$ .
- Let  $P_S = \{s_j : j \in [i] \ \& \ s_j \in X\}$  denote the multiset of all points appearing in  $S$ .
- Define the target function as:

$$g(S) = \frac{1}{|P_S|} \sum_{x \in P_S} f_S(x)$$

Intuitively, in the **non-adaptive world**, the SADA problem is easy to solve because at each step, we can estimate the average of  $f_S$  using a small random sample from  $P_S$ . That is, we do not need to store the entire multiset  $P_S$  (which could be very large) to approximate the target function.

But, in the **adaptive world**, we know that a small sample would not suffice in order to answer many adaptive queries. So, if we want to approximate average of  $f_S$  using a sample from  $P_S$ , then intuitively our sample size must be large which requires large space.

**Theorem (Informal):** There exists a computationally efficient algorithm for the SADA problem in the non-adaptive model that uses  $O(\text{polylog}(m))$  memory, where  $m$  is the length of the stream.

### **Proof Sketch:**

All we need to do is understand how a streaming algorithm can maintain a random sample of the points that have appeared in the stream. Once we know how to do this, we can answer the queries using the sample, and achieve accuracy (in the non-adaptive model) using Chernoff bound.

Moreover, it suffices to show a streaming algorithm that maintains **a single random sample** throughout the stream, because if we can do this, we can simply run  $y$  independent instances of the algorithm in parallel to obtain a sample of size  $y$ .

### **Reservoir Sampling [Vitter, 1985]:**

1. Init  $s = \perp$
2. When the next item  $u_i$  is read, toss a biased coin and with probability  $1/i$  let  $s = u_i$  (note that we need to maintain both the element  $s$  and a counter for  $i$ )
3. Output  $s$

**Theorem:** Assume that the elements in the stream come from a domain of size  $n$ , specifically  $u_i \in [n]$ , and that the stream has length  $m$ . Then, the above algorithm uses  $O(\log(n) + \log(m))$  memory, and its output is an element uniformly sampled from the stream, i.e., for every  $i \in [m]$  we have  $\Pr[s = u_i] = \frac{1}{m}$ .

**Proof:** By induction on  $i \in [m]$

Assume the theorem holds after the first  $j$  steps. That is, after  $j$  steps, for every  $i \in [j]$ , it holds that

$$\Pr[s = u_i] = \frac{1}{j}$$

We will show that the property holds after the next step:

- By construction, for  $u_{j+1}$ , it holds that  $\Pr[s = u_{j+1}] = \frac{1}{j+1}$
- And for  $i \in [j]$  we have

$$\Pr[s = u_i] = \Pr \left[ \begin{array}{c} s = u_i \\ \text{after } j \text{ steps} \end{array} \right] \cdot \Pr[u_i \text{ survives}] = \frac{1}{j} \cdot \left(1 - \frac{1}{j+1}\right) = \frac{1}{j} \cdot \frac{j}{j+1} = \frac{1}{j+1}$$

q.e.d.

**In contrast, as we will now see, the problem becomes difficult in the adaptive model:**

**Theorem (Informal):** Any computationally efficient algorithm for the SADA problem in the adaptive model must use memory  $\Omega(\text{poly}(m))$ .

The difficulty arises from the inability to efficiently answer many adaptive queries using a small sample. As we have seen, if we want to answer many adaptive queries, our sample size must be large, and therefore, intuitively, the memory required to solve the SADA problem must also be large.

Formally we will prove this via reduction: If there exists an adaptive streaming algorithm that solves the SADA problem "too well", we could construct from it a mechanism that answers adaptive queries "too well", which is impossible.

So let us assume that we have an algorithm  $\mathcal{A}$  that solves the SADA problem in the adaptive world. To simplify the proof, assume that  $\mathcal{A}$  is deterministic. We will construct from  $\mathcal{A}$  a mechanism  $\mathcal{M}$  that answers queries as follows:

| <b>Algorithm <math>\mathcal{M}</math></b>   |
|---|
| <p><b>Input:</b> Dataset <math>P = (p_1, \dots, p_n) \in X^n</math> containing <math>n</math> elements from the domain <math>X</math></p> <p><b>Algorithm used:</b> An adversarially robust streaming algorithm <math>\mathcal{A}</math> for the SADA problem</p> <ol style="list-style-type: none"><li>1. For <math>p \in P</math> feed the update <math>p</math> to <math>\mathcal{A}</math></li><li>2. For <math>i = 1, 2, \dots, k</math>:<ol style="list-style-type: none"><li>a) Obtain next function <math>f_i: X \rightarrow \{0, 1\}</math></li><li>b) Feed the update <math>f_i</math> to <math>\mathcal{A}</math> and obtain an answer <math>z_i</math></li><li>c) Output <math>z_i</math></li></ol></li></ol> |

**Let us try and analyze this algorithm:**

- Fix a distribution  $\mathfrak{D}$  over  $X$  and consider the execution of  $\mathcal{M}$  on a sample  $P$  containing i.i.d. samples from  $\mathfrak{D}$ .
- By construction, if algorithm  $\mathcal{A}$  succeeds in solving the SADA problem, then at every step  $i$  we have that

$$z_i \approx \frac{1}{|P|} \sum_{p \in P} f_i(p)$$

- This is good, but it's not exactly what we want... We want to ensure that  $\mathcal{M}$ 's answers are accurate with respect to the distribution  $\mathfrak{D}$ , not just the empirical average over  $P$ !

**Discussion:** How do we complete the proof?