

# Active Fourier Verifier: PAC Estimation of Model Properties with Influence Functions and Fourier Representations

Estimating all properties at once

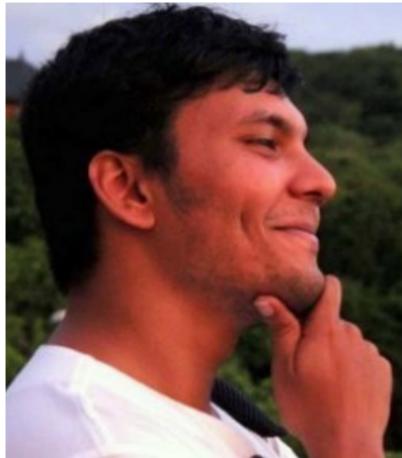
Ayoub Ajarra

23 de noviembre de 2023

*Inria* Joint work with



Bishwamittra Ghosh  
National University of Singapore



Debabrota Basu  
Équipe Scool, Univ. Lille, Inria

- Our problem: Given (restricted) access to a black box model  $h$ , we aim to estimate its properties: robustness and its discrimination toward sensitive subgroups and individuals.
- Our approach: In order to estimate those properties, we choose deterministic or random influence functions:
  - Random influence function for robustness and individual fairness
  - Deterministic influence function for group fairness

*Inria*

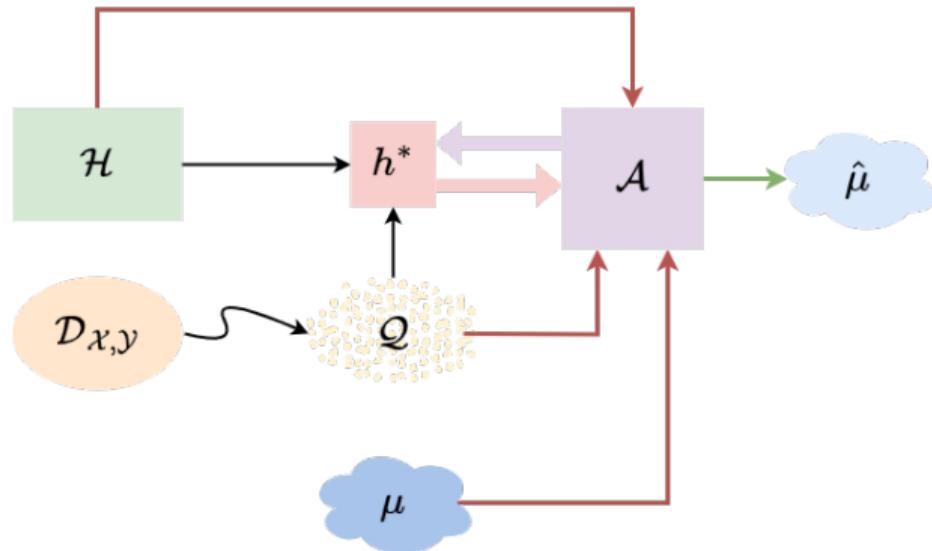
## Problem definition

Property auditing from ZKP



# Our problem setting

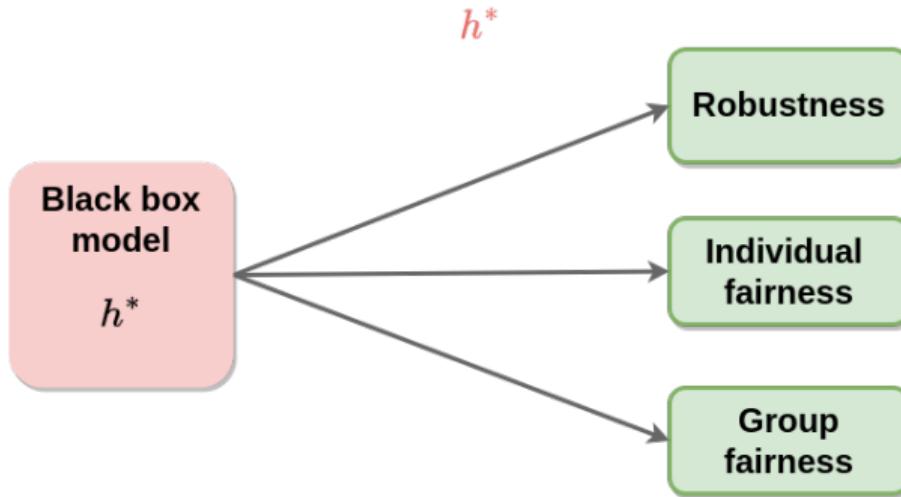
Property auditing: General setting



- Reconstruct then audit. [Tom Y. 2022]
- Sequential testing [Ben.C et al.2023]
- Ours: Estimation by embedding the property of interest in the space of Fourier expandable models.

# Zero-Knowledge Proof (ZKF) properties estimation: Unrealistic

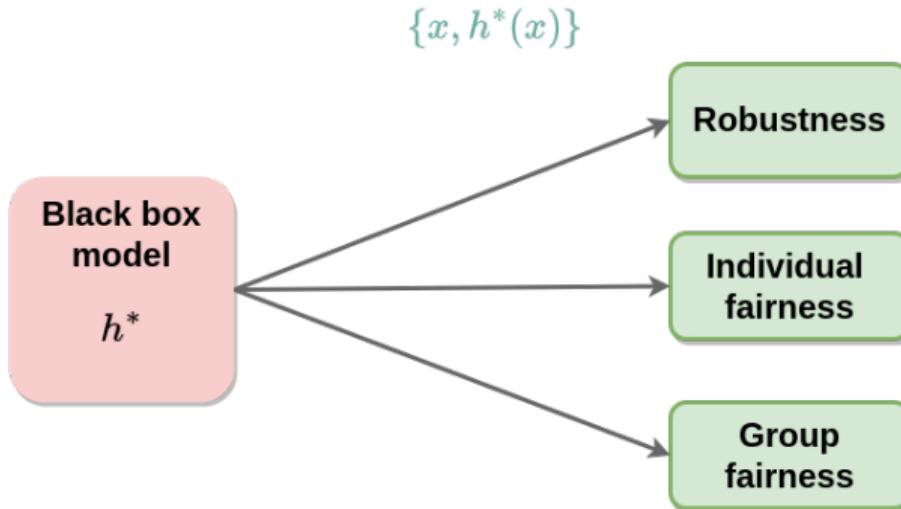
*The proof leaks no information about*



Estimating black-box model's properties in ZKP setting

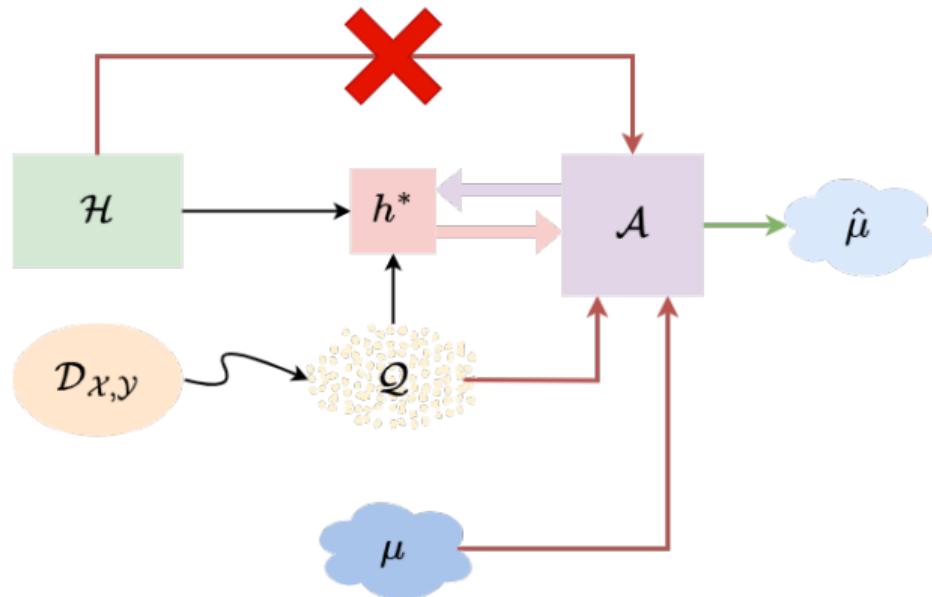
# Friendly Zero-Knowledge Proof (Friendly-ZKF) properties estimation: Our setting

*The proof leaks only information of the form*



Estimating black-box model's properties in friendly ZKP setting

Can we reduce the amount of information communicated to  $\mathcal{A}$ ?



# Boolean functions & Fourier expansion in the context of learning theory

Motivation



### Theorem

Any bounded function  $h : \{-1, 1\}^n \rightarrow \mathcal{Y}$  can be uniquely written as:

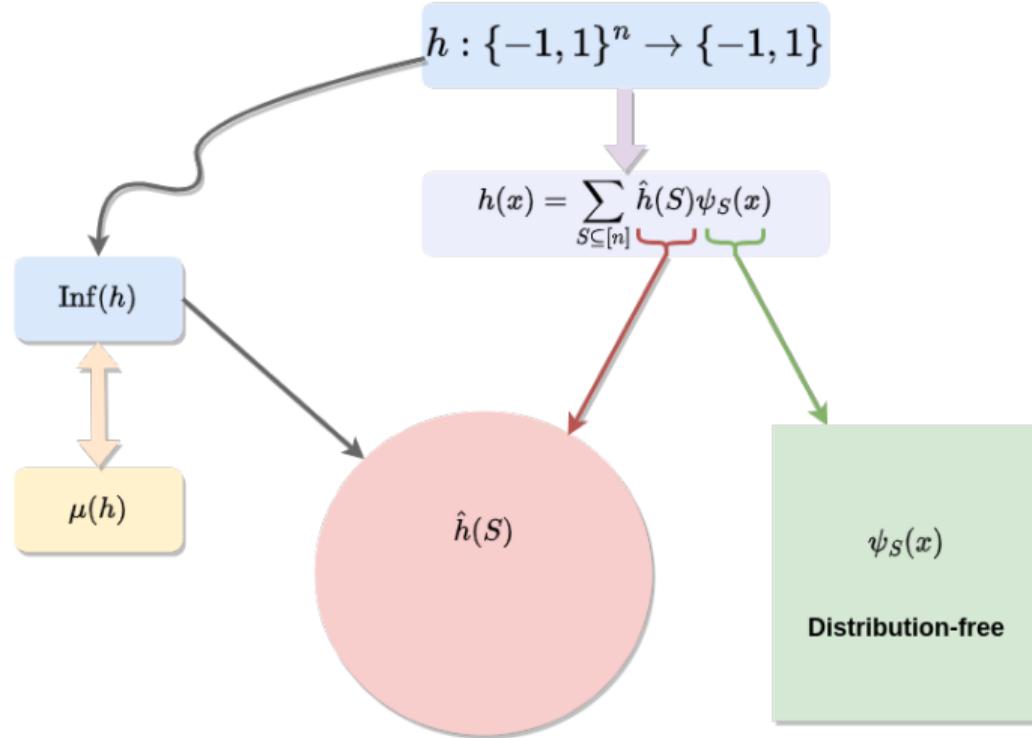
$$\forall (X_1, \dots, X_n) \in \{-1, 1\}^n : h(X_1, \dots, X_n) = \sum_{S \subseteq [n]} \hat{h}(S) \prod_{i \in S} X_i$$

- Assumption 1: We assume that  $h$  is bounded. ✓
- Assumption 2: The distribution over binary values is uniform. ✗
- Assumption 3: The feature space is a matrix of binary values. ✗

## Distribution free setting

Assumption 2: Generalizing to Agnostic Fourier expansion to the distribution

- Method 1: Instead of expanding the model on the basis of parity functions, we can use Gram-Schmidt-type process.
- Method 2: Instead of expanding the model  $h$  in the basis of parity functions, we expand the new model  $\frac{h(x)}{\mathcal{D}_{\mathcal{X}}(x)}$



## Computing target model properties

Target model's properties in terms of Fourier coefficients



$x \sim \mathcal{D}_{\mathcal{X}}$ ,  $\rho \in [-1, 1]$ ,  $l \in [n]$

- $x' \sim \mathcal{N}_\rho(x) \iff \forall i \in [n] : x'_i = \begin{cases} x_i & \text{with probability } \frac{1+\rho}{2} \\ -x_i & \text{with probability } \frac{1-\rho}{2} \end{cases}$
- $x' \sim \mathcal{N}_{\rho,l}(x) \iff \forall i \in [n]^l : x'_i = \begin{cases} x_i & \text{with probability } \frac{1+\rho}{2} \\ -x_i & \text{with probability } \frac{1-\rho}{2} \end{cases}$

## Random influence functions

$\text{Inf}_\rho(h) \triangleq \mathbb{P}_{\substack{x \sim \mathcal{D}_{\mathcal{X}} \\ x' \sim \mathcal{N}_\rho}} [h(x) \neq h(x')] \text{ measures stability.}$

$\text{Inf}_{\rho,l}(h) \triangleq \mathbb{P}_{\substack{x \sim \mathcal{D}_{\mathcal{X}} \\ x' \sim \mathcal{N}_{\rho,l}}} [h(x) \neq h(x')] \text{ measure individual fairness.}$

## Deterministic influence functions

Let A be a sensitive attribute

$$\text{Inf}_A(h) \triangleq \mathbb{P}_{x \sim \mathcal{D}_{\mathcal{X}}} [h(x) \neq h(x^{\sim A})]$$

$$\mu_{GF_A}(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) = y | x \in A^+] - \mathbb{P}_{x \sim \mathcal{D}} [h(x) = y | x \in A^-]$$

## Robustness

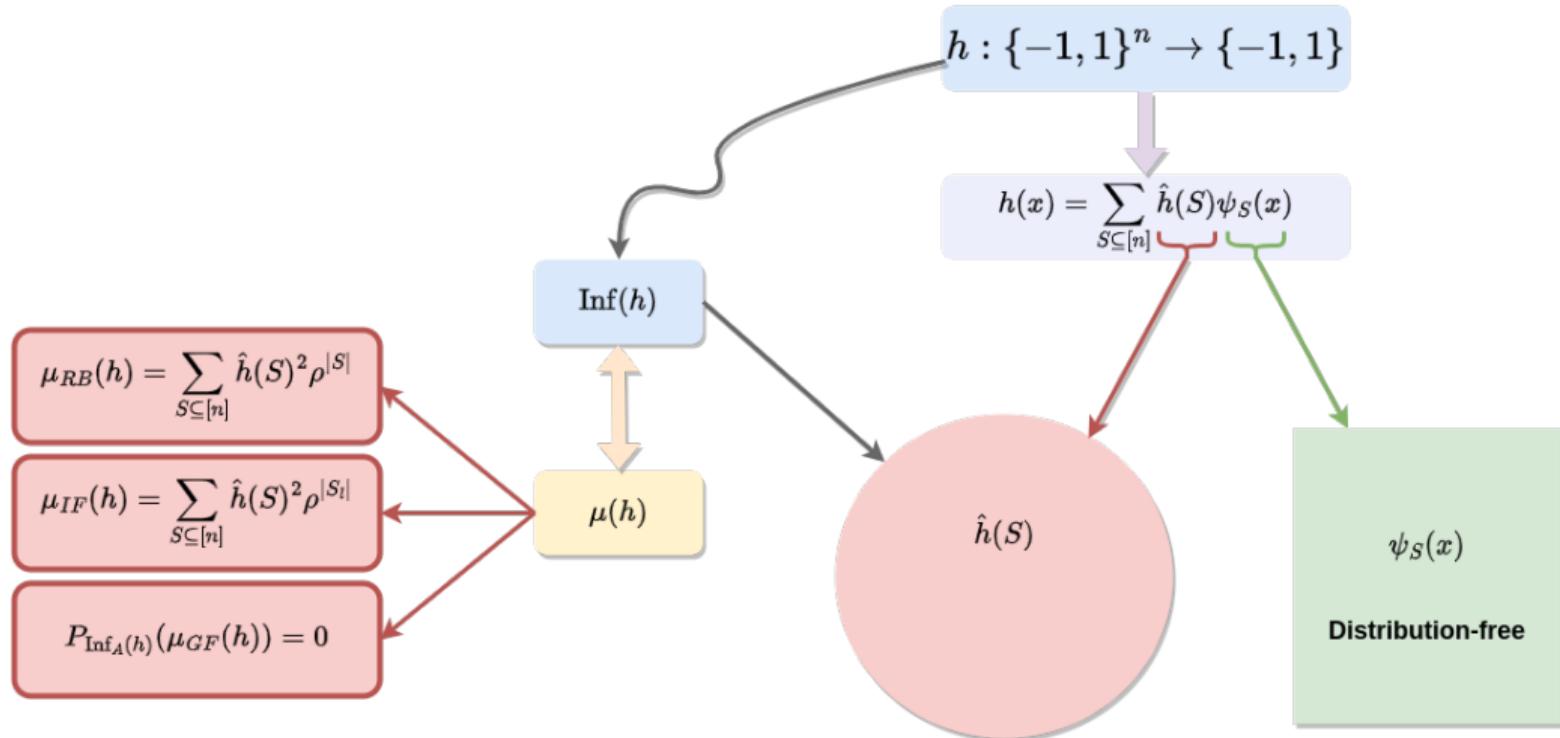
$$\mu_{RB}(h^*) = \sum_{S \subseteq [n]} \hat{h}^*(S)^2 \rho^{|S|}$$

## Individual Fairness

$$\mu_{IF}(h^*) = \sum_{S \subseteq [n]} \hat{h}^*(S)^2 \rho^{|S_1|}$$

## Group Fairness (Informal)

Given the assumption that the marginal distribution is invariant under the flip membership action,  $\text{Inf}_A(h^*)$  is polynomial in  $\mu_{GF}(h^*)$ .



# Hardness of computing target model properties & Universal lower bounds



## Theorem

Given a threshold  $\tau \in \mathbb{R}$ , the problem of testing significant Fourier coefficients with respect to the threshold  $\tau$  is NP-complete.

The degree of a boolean function is the degree of the polynomial of its Fourier representation.

### Example

we can split every function  $h : \{-1, 1\}^n \rightarrow \{-1, 1\}$  to its different spectral levels,  
$$h = \sum_{i=1}^n h^{(i)} \text{ where } h^{(i)} = \sum_{\substack{S \subseteq [n] \\ |S|=i}} \hat{h}(S) \psi_S$$

### Theorem

Let  $d \in \mathbb{N}$ ,  $h \in \mathcal{H}_d$ , where  $\mathcal{H}_d$  is the concept class of Boolean functions of degree at most  $d$ . In other words,  $h$  is  $\epsilon$ -concentrated in some subset of size at most  $d$ .

Algorithm  $\mathcal{A}_\mu$  outputs an  $(\epsilon, \delta)$ -PAC estimate of  $\mu(h)$  with

$$\Omega\left((1 - \epsilon)2^{d-2} \log_2 n - (d + 1)2^{d-2} \log_2(1 - \delta)\right)$$

queries.

*Inria*

## Upper bounds

PAC Fourier Auditor





---

## Algorithm Property estimation in the Boolean domain

---

Input: Confidence parameter  $\delta$ , target error  $\epsilon$ , Sensitive attribute A, perturbation parameters:  $\rho, l, q$  queries

$$\{x_k, h(x_k)\}_{k \in [q]} \leftarrow MQ(h, q)$$

$$L_h \leftarrow CFS(\{x_k, h(x_k)\}_{k \in [q]}, \tau)$$

for  $S \in L_h$  do

$$\hat{h}(S) \leftarrow \frac{1}{q} \sum_{k=1}^q h(x_k) \psi_S(x_k)$$

end for

Output:  $\{\hat{\mu}_{Rob_\rho}, \hat{\mu}_{IF}, \hat{\mu}_{GFA}\}$

---

Algorithm CFS

---

Input:  $\tau, \delta \in (0, 1)$

$$\varepsilon \leftarrow \tau^2/4$$

Instantiate set L

for  $k = 0, \dots, n$  do

    for  $S \subseteq [k]$  do

$\tilde{W}^{S,k} \leftarrow W^{S,k}$  (estimate sum up to accuracy  $\varepsilon$  w.p.  $\geq 1 - \delta$ )

        if  $\tilde{W}^{S,k} \leq \frac{\tau^2}{2}$  then

            Discard  $W^{S,k}$

    end if

end for

Output: L (A list of single elements)

---

## Upper bounds for robustness estimation

$$\mathcal{O}\left(\left(\frac{(C+1)\rho^{\min(S_{\mathcal{F}}, \bar{S}_{\mathcal{F}})}|\mathcal{F}|}{\epsilon}\right)^2 \log \frac{1}{\delta}\right)$$

$$|\mathcal{F}| = \mathcal{O}\left(\frac{n}{\tau^2}\right)$$

$$S_{\mathcal{F}} = \min\{|S| : S \in \mathcal{F}\}$$

$$\bar{S}_{\mathcal{F}} = \min\{|S| : S \notin \mathcal{F}\}$$

## Upper bounds for Individual Fairness estimation

$$\mathcal{O}\left(\left(\frac{(C+1)\rho^{\min(S_{l,\mathcal{F}}, S_{l,\overline{\mathcal{F}}})}|\mathcal{F}|}{\epsilon}\right)^2 \log \frac{1}{\delta}\right)$$

$$|\mathcal{F}| = \mathcal{O}\left(\frac{n}{\tau^2}\right)$$

$$S_{l,\mathcal{F}} = \min\{|S_l| : S_l \in \mathcal{F}\}$$

$$S_{l,\overline{\mathcal{F}}} = \min\{|S_l| : S_l \notin \mathcal{F}\}$$

## Upper bounds for group fairness

$$\mathcal{O}\left(\left(\frac{(C+1)|\mathcal{F}| + \frac{1}{4}}{\epsilon^2}\right)^2 \log \frac{1}{\delta}\right)$$

$$|\mathcal{F}| = \mathcal{O}\left(\frac{n}{\tau^2}\right)$$

## Adversarial reconstruction

Our algorithm guarantees an adversarial reconstruction of the target model only on  $2\epsilon$ -concentrated region over the Fourier space.

## Extension to categorical & continuous domains

Group characters & Lováz extension



# Extension to general domain

Function extension

Categorical domain:

$$\cdot h(x) = \sum_{\zeta} \hat{h}(\zeta) \omega_p(\langle \zeta, x \rangle)$$

Continuous domain:

$$\cdot h(x) = \sum_{S \subseteq [n]} \hat{h}(S) m_i x_i$$

Categorical domain:

- Robustness:

$$\left(\frac{1}{1-p}\right)^n \sum_{\zeta: \zeta_j \neq 0} |\hat{f}(\zeta)|^2 + \sum_{\zeta: \zeta_j = 0} |\hat{f}(\zeta)|^2$$

Continuous domain:

- Robustness

$$\frac{1-\rho}{2} \sum_{S \subseteq [n]} \hat{h}(S)^2 \left(1 + |S| \zeta \frac{1+\rho}{1-\rho}\right)$$

Categorical domain:

- Robustness:

$$\left(\frac{1}{1-p}\right)^n \sum_{\zeta: \zeta_j \neq 0} |\hat{f}(\zeta)|^2 + \sum_{\zeta: \zeta_j = 0} |\hat{f}(\zeta)|^2$$

- Individual fairness:

$$\frac{1}{p} \sum_{\zeta} |\hat{f}(\zeta)|^2 \cos\left(\frac{2\pi}{p} \sum_{T \in \mathbb{F}_p^n} \langle \zeta, T \rangle\right)$$

Continuous domain:

- Robustness

$$\frac{1-\rho}{2} \sum_{S \subseteq [n]} \hat{h}(S)^2 \left(1 + |S| \zeta \frac{1+\rho}{1-\rho}\right)$$

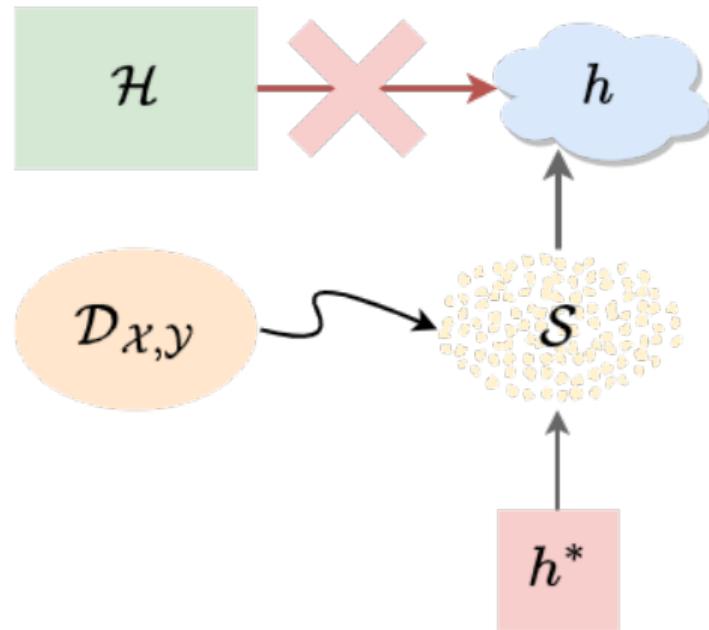
- Individual fairness

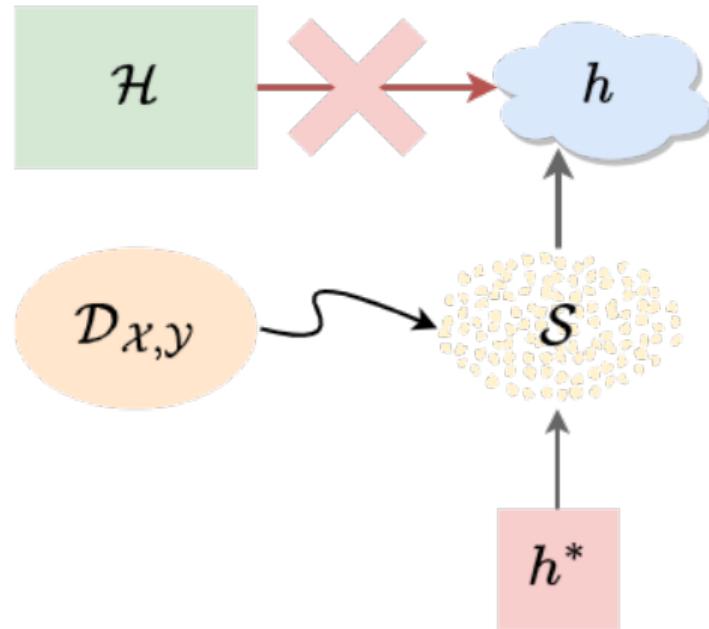
$$\frac{1-\rho}{2} \sum_{S \subseteq [n]} \hat{h}(S)^2 \left(1 + |S_1| \zeta \frac{1+\rho}{1-\rho}\right)$$

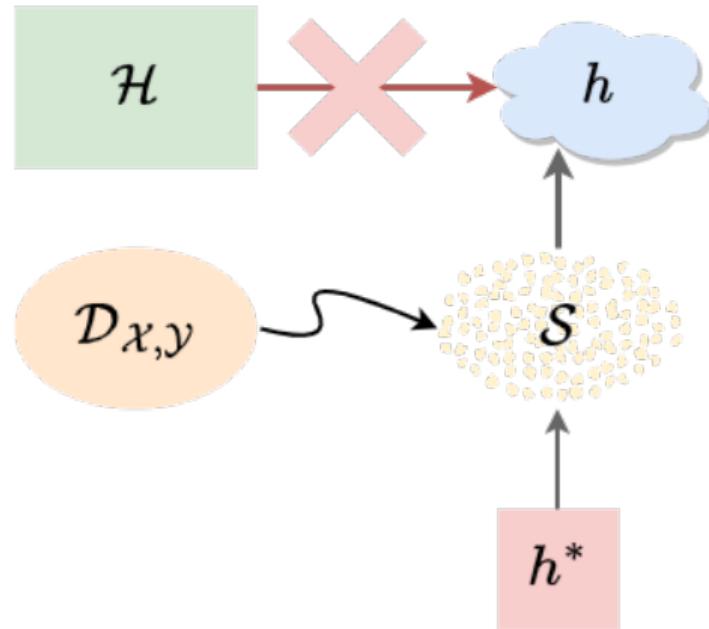
## Open Problem

PAC learning with zero inductive bias









# PAC-Fourier Auditor

## Estimating all properties at once

## Appendix

Ommited technical details



## Gram-Schmidt orthogonalization

Assumption 2: Generalizing to Agnostic Fourier expansion to the distribution

For an unknown distribution  $\mathcal{D}$ , the set of parity functions are not necessarily orthogonal. Fix the following subsets of  $\{1, \dots, n\}$  in the following order:

$$\{\emptyset\}, \{1\}, \{2\}, \{1, 2\}, \{3\}, \{1, 3\}, \{2, 3\}, \{1, 2, 3\}, \dots, \{1, 2, \dots, n\}$$

- Apply Gram-Schmidt process on the monomials of parity functions (to make them orthogonal) with the above ordering.
- The first element of the basis is trivially:  $\psi_\emptyset = 1$
- The  $j$ th basis function corresponding to  $S_j$  is obtained from the following operation.

$$\tilde{\psi}_{S_j} = \chi_{S_j} - \sum_{l=1}^{j-1} \langle \psi_{S_l}, \chi_{S_j} \rangle_{\mathcal{D}} \psi_{S_l}$$

$$\psi_{S_j} = \begin{cases} \frac{\tilde{\psi}_{S_j}}{\|\tilde{\psi}_{S_j}\|_{2,\mathcal{D}}} & \text{if } \|\tilde{\psi}_{S_j}\|_{2,\mathcal{D}} > 0, \\ 0 & \text{Otherwise} \end{cases}$$

- The new Fourier expansion is given by the formula :

$$\frac{h(x)}{\mathcal{D}_{\mathcal{X}}(x)} = \sum_{S \subseteq [n]} \hat{f}_{\mathcal{D}}(S) \chi_S(x)$$

- The Fourier coefficients are given by:

$$\forall S \subseteq [n] : \hat{f}(S) = \langle f, \chi_S \rangle$$

- All the results obtained from uniform distribution remain the same considering this approach. However, in practice we don't have access to  $\mathcal{D}$ , considering an empirical distribution could lead to unsatisfying results.

Let  $G$  be a finite abelian group. We consider the general case:  $G = \prod_{i=1}^k \mathbb{F}_{p_i}$

### Group character

A map  $\chi : G \rightarrow \mathbb{C} - \{0\}$  is called a character of  $G$  if it is a group homomorphism, that is:

$$\chi(0) = 1$$

$$\forall a, b \in G : \chi(a + b) = \chi(a) + \chi(b)$$

The constant map 1 is always a character for any abelian group, and is called the principal character of  $G$ .

Let  $G$  be a finite abelian group. We consider the general case:  $G = \prod_{i=1}^k \mathbb{F}_{p_i}$

For every observation  $a = (a_1, \dots, a_k) \in G$ , define  $\chi_a \in \mathcal{L}_2(G)$ :

$$\chi_a : x \rightarrow \prod_{j=1}^k e^{\frac{i2\pi}{p_j} a_j x_j}$$

## Theorem

If  $G$  is a finite Abelian group, then the characters of  $G$  form an orthonormal basis for  $\mathcal{L}_2(G)$ . Furthermore, we have  $G \cong \tilde{G}$ .

$\tilde{G}$  is called the Pontryagin dual of  $G$  and it is the group of characters of  $G$  together with the usual point-wise product of complex-valued functions.

## Extension to finite groups

Fourier expansion in finite domain as cartesian product of different cardinals

Let  $G$  be a finite abelian group. We consider the general case:  $G = \prod_{i=1}^k \mathbb{F}_{p_i}$

For every observation  $a = (a_1, \dots, a_k) \in G$ , define  $\chi_a \in \mathcal{L}_2(G)$ :

$$\chi_a : x \rightarrow \prod_{j=1}^k e^{\frac{i2\pi}{p_j} a_j x_j}$$

### Theorem

The Fourier transform of a function  $f : G \rightarrow \mathbb{C}$  is the unique function  $\hat{f} : \hat{G} \rightarrow \mathbb{C}$  defined as  $\hat{f}(a) = \langle f, \chi_a \rangle = \mathbb{E}[f(x)\bar{\chi}_a(x)]$ .

It follows from the fact that the characters form an orthonormal basis for  $\mathcal{L}_2(G)$  that  $f = \sum_{a \in G} \hat{f}(a) \chi_a$

# Relaxation: Assumptions on the distribution

Partial answer

- Distribution dependent rates [Cohen, K., "Local Glivenko-Cantelli", COLT 2023].
- Bayes consistency, No rate [Steve Hanneke et al. Universal Bayes consistency in metric spaces", AOS 2021.]

Estimating all  
properties at once!