## Statistical-Computational Trade-offs for Hypothesis Selection

Anders Aamand, Alexandr Andoni, Justin Y. Chen, Piotr Indyk, Shyam Narayanan, Sandeep Silwal, Haike Xu



 $p_4$ 

*p*<sub>1</sub>

## Searching over Distributions

**Optimizing Samples or Time** 

Preprocess: k discrete distributions over [n]

 $p_1, p_2, ..., p_k$ 



Samples:  $O(\log k/\epsilon^2)$ Existing solutions use  $\Omega(k)$  time (e.g., via a tournament)

Time:  $O(k^{\rho})$  query time

With  $O(n/\epsilon^2)$  samples and  $O(k^{1+\rho})$ 

Query: samples from unknown distribution  $p_i$ 





Trivial solutions if we allow superpolynomial preprocessing/space

Goal: output  $p_i$  that is  $\epsilon$ -close to  $p_i$  (*TV* distance)

Nearest neighbor search over distributions

Sublinear samples (in *n*) and query time (in k)?

Our Results: Novel Statistical-Computational Trade-offs

1.000

0.975

0.950

0.925

0.900 -

0.875-

UPPER BOUND [AACINS'23]: Doubly sublinear data structure using hierarchy of LSH-based data structures on "heavy"/"light" elements.

LOWER BOUND [AACINSX'24]: Any data structure in the list of points model with samples must have runtime



- Insufficient samples to *learn* unknown  $p_i$
- Insufficient runtime to scan over all distributions
- But both facts above are *barely* true!

Can we get truly sublinear samples and query time?  $n^{0.9}$  samples and  $k^{0.9}$  time?

[AK '20]: Thomas D Ahle and Jakob BT Knudsen. Subsets and supermajorities: Optimal hashing-based set similarity search. FOCS '20. [AACINS '23]: Anders Aamand, Alexandr Andoni, Justin Y. Chen, Piotr Indyk,







