

Sublinear Algorithms

LECTURE 22

Last time

- Approximating the distance to sortedness (length of LIS) of 0/1 sequences
- Gap Edit Distance

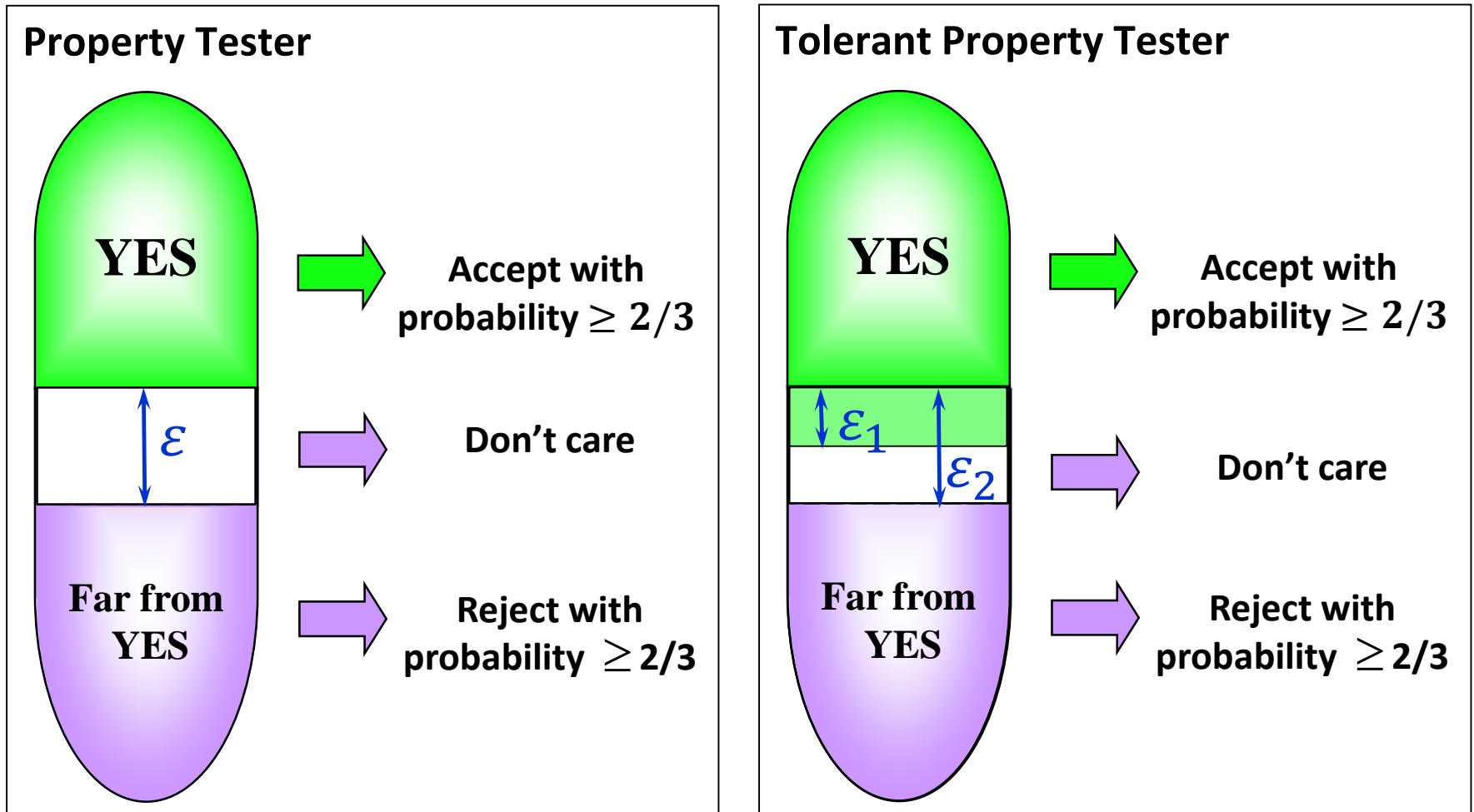
Today

- L_p -testing



Project Reports are due December 3

Testing Models



Two objects are at distance ϵ = they differ in an ϵ fraction of places
Equivalent problem: approximating distance to the property.

Why Hamming Distance?

- Nice probabilistic interpretation
 - probability that two functions differ on a random point in the domain
- Natural measure for
 - algebraic properties (linearity, low degree)
 - properties of graphs and other combinatorial objects
- Motivated by applications to probabilistically checkable proofs (PCPs)
- It is equivalent to other natural distances for
 - properties of Boolean functions

Which stocks grew steadily?



Microsoft



IBM



Data from

<http://finance.google.com>

L_p -Testing

for properties of real-valued data

[Berman Raskhodnikova Yaroslavtsev]

Use L_p -metrics to Measure Distances

- Functions $f, g: D \rightarrow [0,1]$ over (finite) domain D

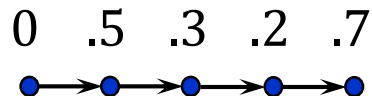
Normalize the values, so they are between 0 and 1

- For $p \geq 1$
$$L_p(f, g) = \|f - g\|_p = \left(\sum_{x \in D} |f(x) - g(x)|^p \right)^{1/p}$$

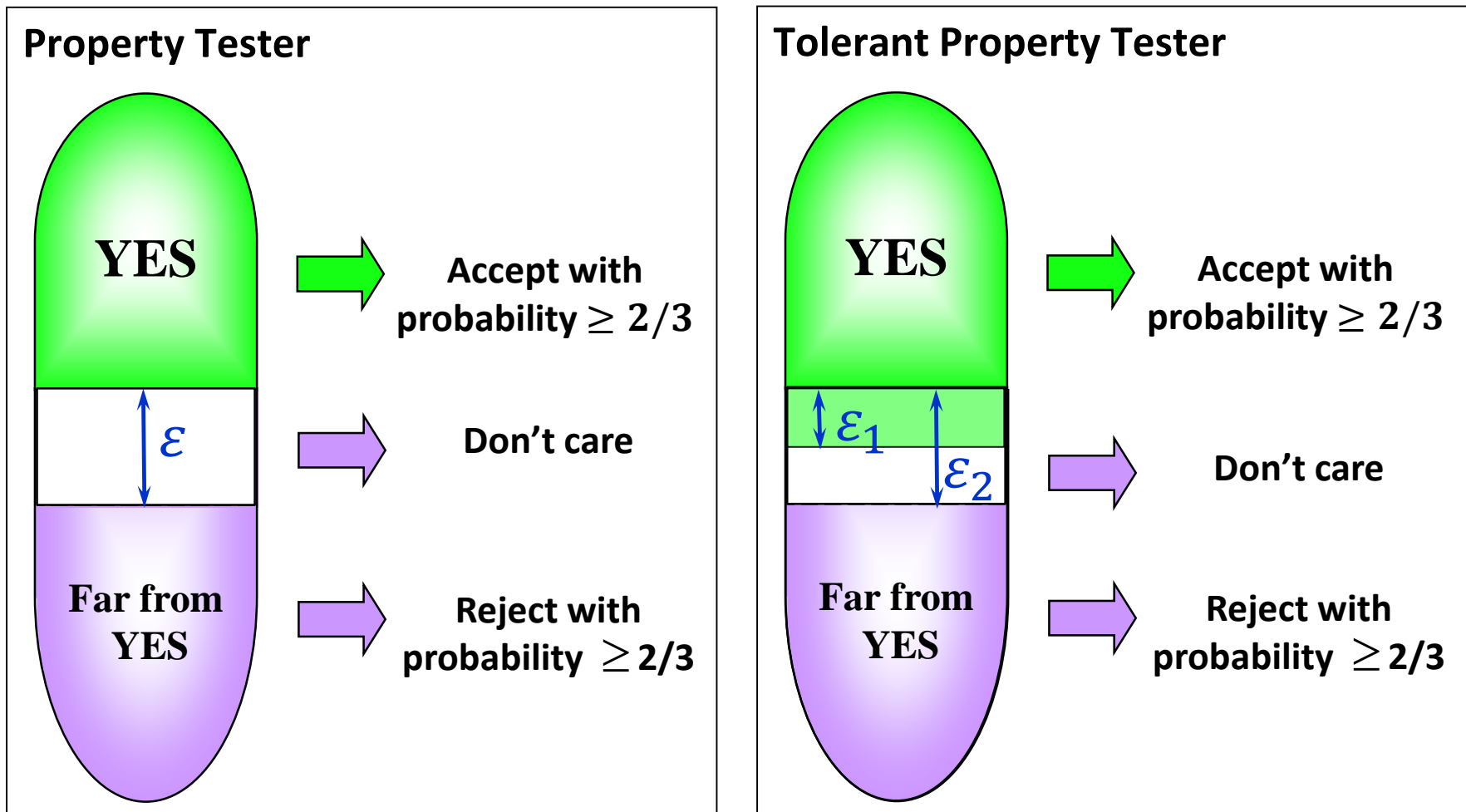
$$L_0(f, g) = \|f - g\|_0 = |\{x \in D: f(x) \neq g(x)\}|$$

- $$d_p(f, g) = \frac{\|f - g\|_p}{\|1\|_p}$$

Example:



L_p -Testing and Tolerant L_p -Testing



Functions $f, g: D \rightarrow [0,1]$ are at distance ϵ if $d_p = \frac{\|f-g\|_p}{\|1\|_p} = \epsilon$.

L_p -Testing Model for Real-Valued Data

- Generalizes standard L_0 -testing
- For $p > 0$ still have a nice probabilistic interpretation:
distance $d_p(f, g) = (\mathbf{E}[|f - g|^p])^{1/p}$
- Compatible with existing PAC-style learning models
(preprocessing for model selection)
- For Boolean functions, $d_0(f, g) = d_p(f, g)^p$.

Plan

1. Relationships between L_p -testing models
2. L_p -testing monotonicity

Relationships between L_p -Testing Models

Relationships Between L_p -Testing Models

$C_p(P, \epsilon)$ = complexity of L_p -testing property P
with distance parameter ϵ

- e.g., query or time complexity
- for general or restricted (e.g., nonadaptive) tests

For all properties P

- L_1 -testing is no harder than Hamming testing

$$C_1(P, \epsilon) \leq C_0(P, \epsilon)$$

- L_p -testing for $p > 1$ is close in complexity to L_1 -testing

$$C_1(P, \epsilon) \leq C_p(P, \epsilon) \leq C_1(P, \epsilon^p)$$

Relationships Between L_p -Testing Models

$C_p(P, \varepsilon)$ = complexity of L_p -testing property P
with distance parameter ε

- e.g., query or time complexity
- for general or restricted (e.g., nonadaptive) tests

For properties of Boolean functions $f: D \rightarrow \{0,1\}$

- L_1 -testing is equivalent to Hamming testing

$$C_1(P, \varepsilon) = C_0(P, \varepsilon)$$

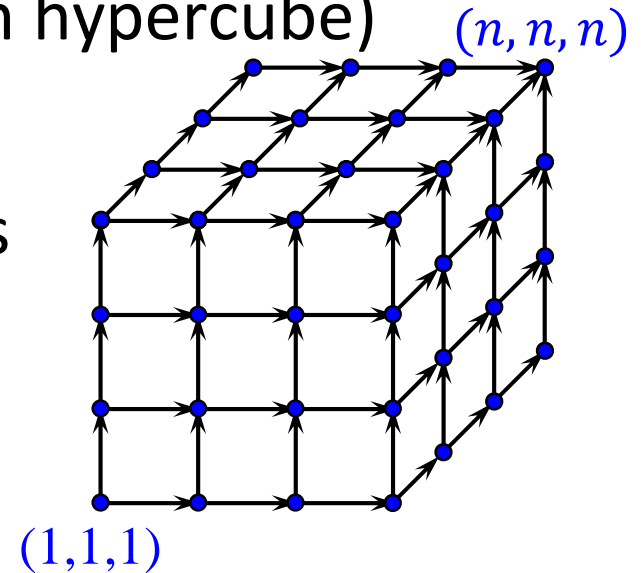
- L_p -testing for $p > 1$ is equivalent to L_1 -testing
with appropriate distance parameter

$$C_p(P, \varepsilon) = C_1(P, \varepsilon^p)$$

Property: Monotonicity of Functions

Monotonicity

- Domain $D=[n]^d$ (vertices of d -dim hypercube)
- A function $f: D \rightarrow \mathbb{R}$ is **monotone** if increasing a coordinate of x does not decrease $f(x)$.
- Special case $d = 1$
 $f: [n] \rightarrow \mathbb{R}$ is monotone $\Leftrightarrow f(1), \dots, f(n)$ is sorted.



Monotonicity Testers: Running Time

| f | L_0 | L_p |
|--------------------------------|---|--|
| $[n]$ $\rightarrow [0,1]$ | $\Theta\left(\frac{\log n}{\varepsilon}\right)$ [Ergün Kannan Kumar Rubinfeld Viswanathan 00, Fischer 04] | $\Theta\left(\frac{1}{\varepsilon^p}\right)$ |
| $[n]^d$ $\rightarrow [0,1]$ | $\Theta\left(\frac{d \cdot \log n}{\varepsilon}\right)$ [Chakrabarty Seshadhri 13] | $O\left(\frac{d}{\varepsilon^p} \log \frac{d}{\varepsilon^p}\right)$ $\Omega\left(\frac{1}{\varepsilon^p} \log \frac{1}{\varepsilon^p}\right)$ for $d = 2$ nonadaptive 1-sided error |

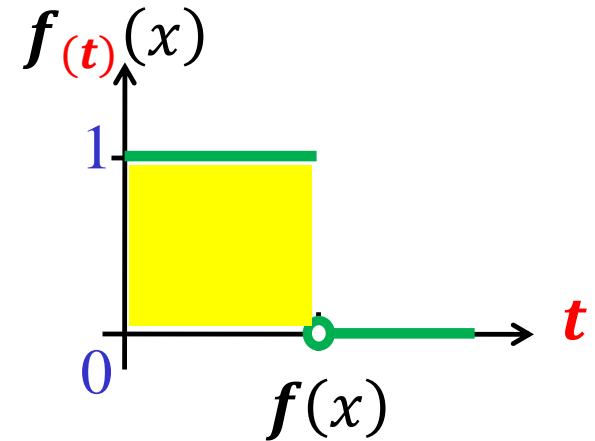
* Hiding some $\log 1/\varepsilon$ dependence

L_1 -Testing of Monotonicity

Monotonicity: Reduction to Boolean Functions

Boolean threshold function $f_{(t)}: D \rightarrow \{0,1\}$

$$f_{(t)}(x) = \begin{cases} 1 & \text{if } f(x) \geq t \\ 0 & \text{otherwise} \end{cases}$$

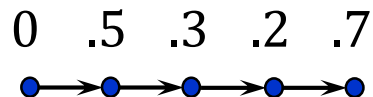


- Decomposition: $f(x) = \int_0^1 f_{(t)}(x) dt$
- M = class of monotone functions

Characterization Theorem

$$L_1(f, M) = \int_0^1 L_1(f_{(t)}, M) dt$$

Example:



Characterization Theorem: One Direction

$$L_1(\mathbf{f}, M) \leq \int_0^1 L_1(\mathbf{f}(t), M) dt$$

- $\forall t \in [0,1]$, let g_t = closest monotone (Boolean) function to $\mathbf{f}(t)$.
- Let $\mathbf{g} = \int_0^1 g_t dt$. Then \mathbf{g} is monotone, since g_t are monotone.

$$L_1(\mathbf{f}, M) \leq \|\mathbf{f} - \mathbf{g}\|_1$$

\mathbf{g} is monotone

$$= \left\| \int_0^1 \mathbf{f}(t) dt - \int_0^1 g_t dt \right\|_1$$

Decomposition & definition of \mathbf{g}

$$= \left\| \int_0^1 (\mathbf{f}(t) - g_t) dt \right\|_1$$

$$\leq \int_0^1 \|\mathbf{f}(t) - g_t\|_1 dt$$

Triangle inequality

$$= \int_0^1 L_1(\mathbf{f}(t), M) dt$$

Definition of g_t

Characterization Theorem: the Other Direction

$$L_1(\mathbf{f}, M) \geq \int_0^1 L_1(\mathbf{f}(t), M) dt$$

- Let \mathbf{h} be closest monotone function to \mathbf{f} .
- Then $\mathbf{h}(t)$ is monotone for all $t \in [0,1]$.

$$L_1(\mathbf{f}, M) = \|\mathbf{f} - \mathbf{h}\|_1$$

Because \mathbf{h} is monotone

$$= \left\| \int_0^1 (\mathbf{f}(t) - \mathbf{h}(t)) dt \right\|_1$$

Decomposition

$$\mathbf{f}(x) \geq \mathbf{h}(x)$$

\Leftrightarrow

$$\mathbf{f}(t) \geq \mathbf{h}(t)$$

$$\forall t \in [0,1]$$

$$= \sum_{x:\mathbf{f}(x) \geq \mathbf{h}(x)} \int_0^1 (\mathbf{f}(t) - \mathbf{h}(t)) dt + \sum_{x:\mathbf{f}(x) < \mathbf{h}(x)} \int_0^1 (\mathbf{h}(t) - \mathbf{f}(t)) dt$$

$$= \int_0^1 \left(\sum_{x:\mathbf{f}(x) \geq \mathbf{h}(x)} (\mathbf{f}(t) - \mathbf{h}(t)) + \sum_{x:\mathbf{f}(x) < \mathbf{h}(x)} (\mathbf{h}(t) - \mathbf{f}(t)) \right) dt$$

$$= \int_0^1 \|\mathbf{f}(t) - \mathbf{h}(t)\|_1 dt$$

Triangle inequality

$$\geq \int_0^1 L_1(\mathbf{f}(t), M) dt$$

$\mathbf{h}(t)$ is monotone

Monotonicity: Using Characterization Theorem

Characterization Theorem

$$d_1(\mathbf{f}, M) = \int_0^1 d_1(\mathbf{f}_{(t)}, M) dt$$

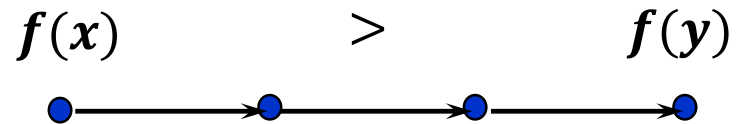
We can use Characterization Theorem
to get monotonicity L_1 -testers
and tolerant testers from
standard property testers for Boolean functions.

L_1 -Testers from Testers for Boolean Ranges

A nonadaptive, 1-sided error L_0 -test for monotonicity of $f: D \rightarrow \{0,1\}$ is also an L_1 -test for monotonicity of $f: D \rightarrow [0,1]$.

Proof:

- A **violation** (x, y) :



- A nonadaptive, 1-sided error test queries a random set $Q \subseteq D$ and rejects iff Q contains a violation.
- If $f: D \rightarrow [0,1]$ is monotone, Q will not contain a violation.
- If $d_1(f, M) \geq \varepsilon$ then $\exists t^*: d_0(f_{(t^*)}, M) \geq \varepsilon$
- W.p. $\geq 2/3$, set Q contains a violation (x, y) for $f_{(t^*)}$

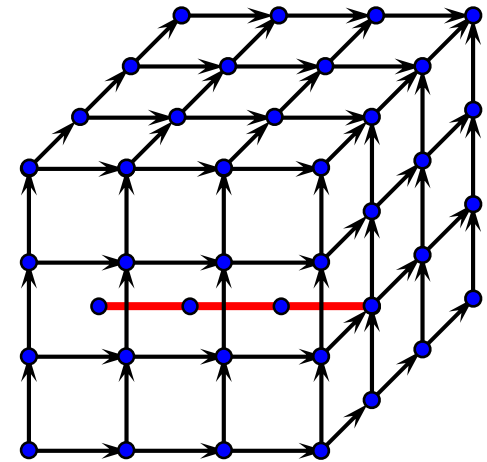
$$f_{(t^*)}(x) = 1, f_{(t^*)}(y) = 0$$

\Downarrow

$$f(x) > f(y)$$

L_0 -Testing Monotonicity of $f: [n]^d \rightarrow \{0, 1\}$

- Idea:
1. Pick axis-parallel lines ℓ .
 2. Sample points from each ℓ ,
and check for violations of $f|_{\ell}$.



[DGLRRS 99]

- **Testing sortedness:** If $f: [n] \rightarrow \{0,1\}$ is ε -far from sorted then $O\left(\frac{1}{\varepsilon}\right)$ samples are sufficient to find a violation w/ const. prob.
- **Dimension reduction:** For $f: [n]^d \rightarrow \{0,1\}$

$$\mathbb{E}[d_0(f|_{\ell}, M)] \geq \frac{d_0(f, M)}{2d}.$$

How many lines should we sample?

How many points form each line?

General Work Investment Problem [Goldreich 13]

- Algorithm needs to find “evidence” (e.g., **a violation**).
- It can select an element from distr. Π (e.g., **a uniform line**).
- Elements e have different quality $q(e) \in [0,1]$
(e.g., $d_0(f_{|\ell}, M)$).
- Algorithm must invest more work into e with lower $q(e)$ to extract evidence from e (e.g., **need $\Theta\left(\frac{1}{q(e)}\right)$ samples**).
- $\mathbb{E}_{e \leftarrow \Pi}[q(e)] \geq \mu$.

What’s a good work investment strategy?

Used in [Levin 85, Goldreich Levin 89], testing connectedness of a graph [Goldreich Ron 97], testing properties of images [R 03], multi-input testing problems [G13]

Work Investment Strategies

- “Reverse” Markov Inequality

For a random variable $X \in [0,1]$ with expectation $\mathbb{E}[X] \geq \mu$,

$$\Pr \left[X \geq \frac{\mu}{2} \right] \geq \frac{\mu}{2}.$$

Proof: $\mu \leq \mathbb{E}[X] \leq \Pr \left[X \geq \frac{\mu}{2} \right] \cdot 1 + \Pr \left[X < \frac{\mu}{2} \right] \cdot \frac{\mu}{2}.$

“Reverse” Markov Strategy:

1. Sample $\Theta \left(\frac{1}{\mu} \right)$ lines.
2. Sample $\Theta \left(\frac{1}{\mu} \right)$ points from each line.

Cost: $\Theta \left(\frac{1}{\mu^2} \right)$ queries.

Work Investment Strategies

Bucketing idea [Levin, Goldreich 13]:

Invest in elements of quality $q(e) \geq \frac{1}{2^i}$ separately.

Bucketing Inequality [Berman R Yaroslavtsev 14]

For a random variable $X \in [0,1]$ with $\mathbb{E}[X] \geq \mu$, let

$$p_i = \Pr \left[X \geq \frac{1}{2^i} \right] \text{ and } k_i = \Theta \left(\frac{1}{2^i \mu} \right).$$

Then $\prod_{i=1}^{\log 4/\mu} (1 - p_i)^{k_i} \leq 1/3$.

Bucketing Strategy: For each bucket $i \in \left[\log \frac{4}{\mu} \right]$

1. Sample $k_i = \Theta \left(\frac{1}{2^i \mu} \right)$ lines.
2. Sample $\Theta(2^i)$ points from each line.

Cost: $\Theta \left(\frac{1}{\mu} \log \frac{1}{\mu} \right)$ queries (for monotonicity, $\mu = \frac{\varepsilon}{2d}$)

Monotonicity Testers: Running Time

| f | L_0 | L_p |
|----------------------------------|--|--|
| $[n]$ $\rightarrow \{0,1\}$ | $\Theta\left(\frac{1}{\epsilon}\right)$ | $\Theta\left(\frac{1}{\epsilon^p}\right)$ |
| $[n]^d$ $\rightarrow \{0,1\}$ | $O\left(\frac{d}{\epsilon} \cdot \log \frac{d}{\epsilon}\right)$ | $O\left(\frac{d}{\epsilon^p} \log \frac{d}{\epsilon^p}\right)$ $\Omega\left(\frac{1}{\epsilon^p} \log \frac{1}{\epsilon^p}\right)$ for $d = 2$ nonadaptive 1-sided error |
| | | $\Theta\left(\frac{1}{\epsilon^p}\right)$ for constant d adaptive 1-sided error |

Monotonicity Testers: Running Time

| f | L_0 | L_p |
|--------------------------------|---|--|
| $[n]$ $\rightarrow [0,1]$ | $\Theta\left(\frac{\log n}{\varepsilon}\right)$ [Ergün Kannan Kumar Rubinfeld Viswanathan 00, Fischer 04] | $\Theta\left(\frac{1}{\varepsilon^p}\right)$ |
| $[n]^d$ $\rightarrow [0,1]$ | $\Theta\left(\frac{d \cdot \log n}{\varepsilon}\right)$ [Chakrabarty Seshadhri 13] | $O\left(\frac{d}{\varepsilon^p} \log \frac{d}{\varepsilon^p}\right)$ $\Omega\left(\frac{1}{\varepsilon^p} \log \frac{1}{\varepsilon^p}\right)$ for $d = 2$ nonadaptive 1-sided error |

* Hiding some $\log 1/\varepsilon$ dependence

Open Problems

- Our L_1 -tester for monotonicity is nonadaptive, but we show that adaptivity helps for Boolean range.

Is there a better adaptive tester?

- All our algorithms for L_p -testing for $p \geq 1$ were obtained directly from L_1 -testers.

Can one design better algorithms by working directly with L_p -distances?

- We designed tolerant tester only for monotonicity ($d=1,2$).

Tolerant testers for higher dimensions?

Other properties?