

L_p -Testing

Sofya Raskhodnikova

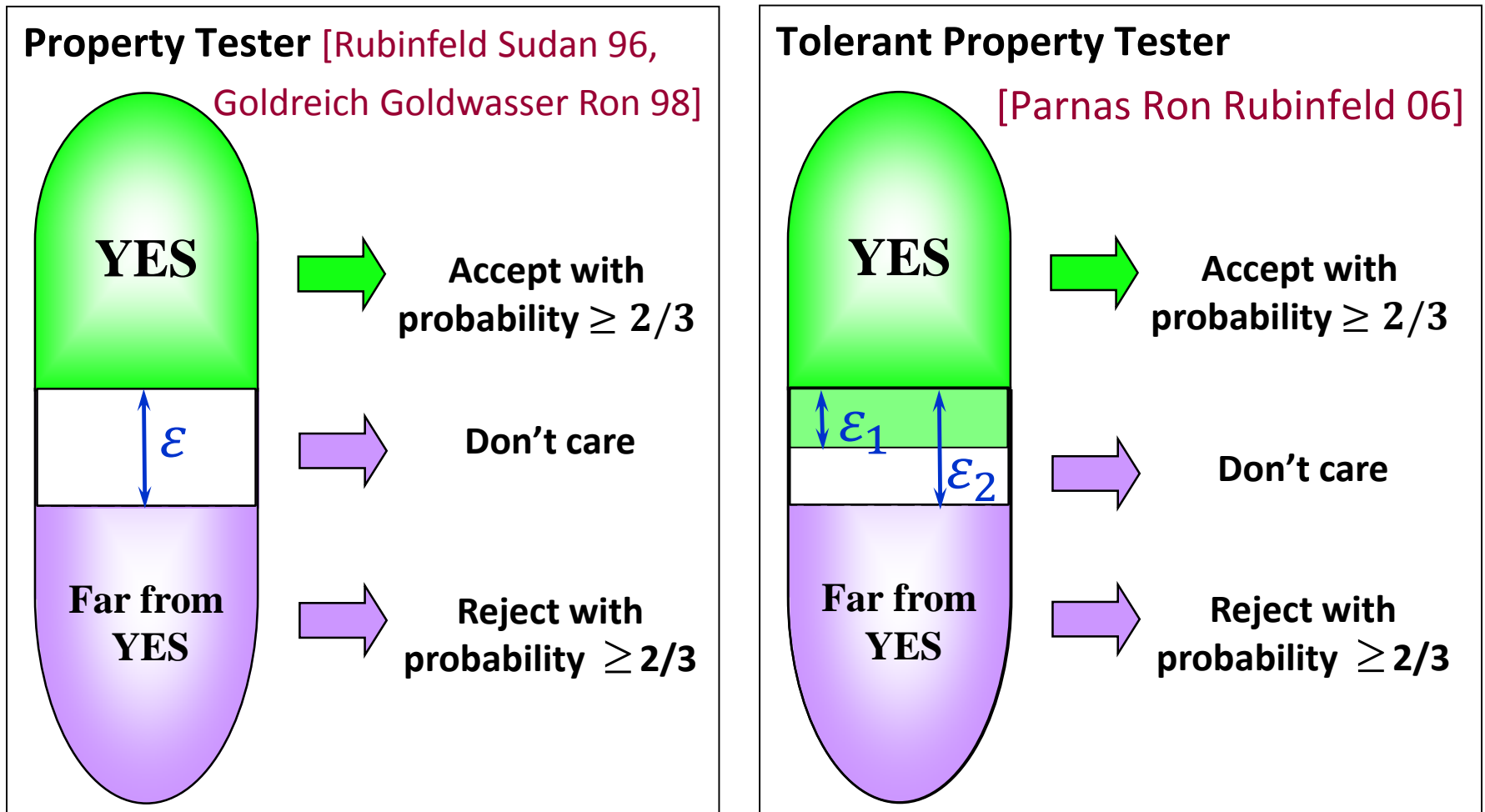
*Penn State University,
visiting Boston University and Harvard*

Joint work with

Piotr Berman (Penn State),

Grigory Yaroslavtsev (Penn State → Brown)

Property Testing Models



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

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

Use L_p -metrics to Measure Distances

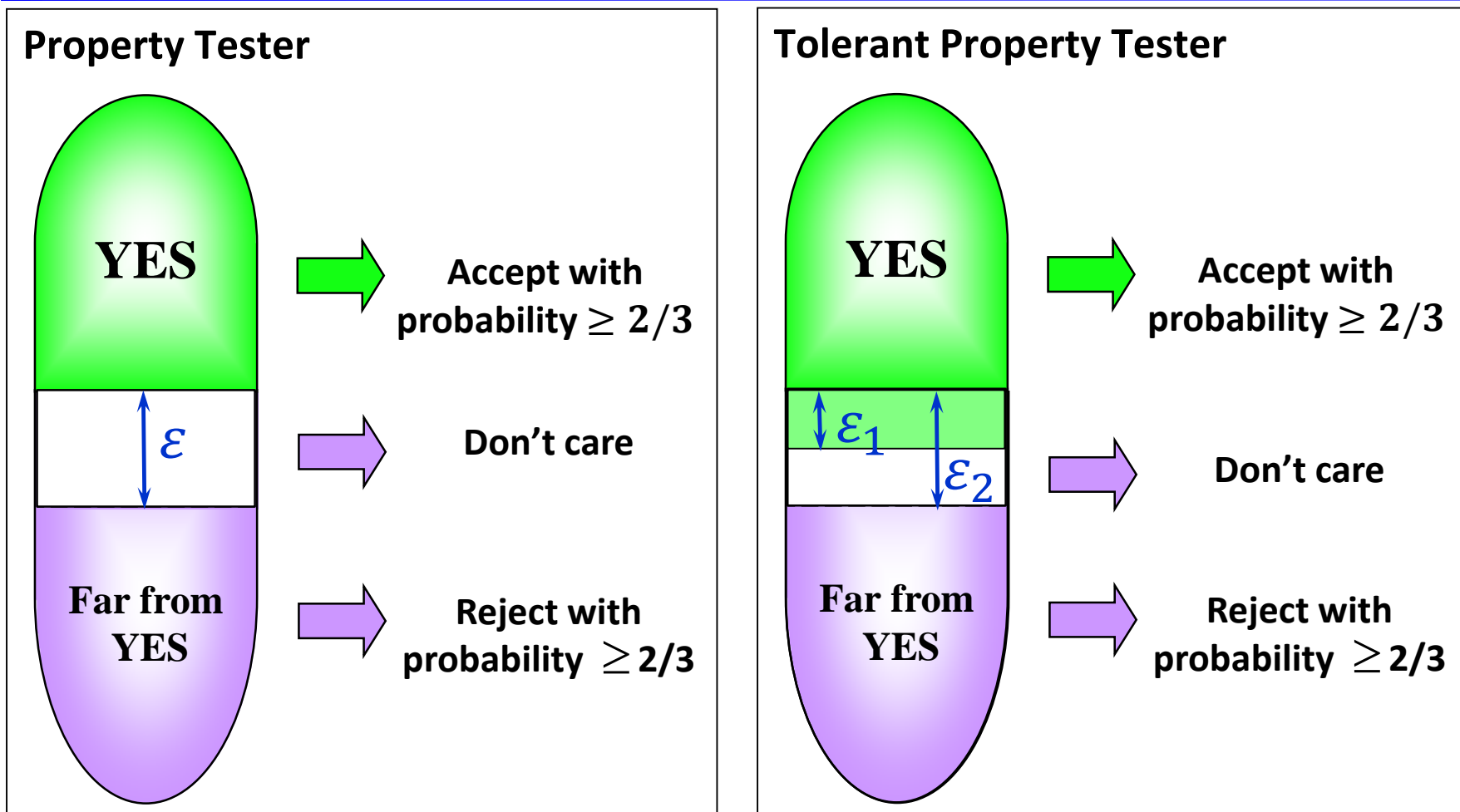
- Functions $f, g: D \rightarrow [0,1]$ over (finite) domain D
- For $p \geq 1$

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

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

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

L_p -Testing and Tolerant L_p -Testing



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

New 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, $d_0(f, g) = d_p(f, g)^p$.

Our Contributions

1. Relationships between L_p -testing models
2. Algorithms
 - L_p -testers for $p \geq 1$
 - monotonicity, Lipschitz, convexity
 - Tolerant L_p -tester for $p \geq 1$
 - monotonicity in 1D (aka sortedness)
 - ❖ Our L_p -testers beat lower bounds for L_0 -testers
 - ❖ Simple algorithms backed up by involved analysis
 - ❖ Uniformly sampled (or easy to sample) data suffices
3. Nearly tight lower bounds

Implications for L_0 -Testing

Some techniques/observations/results carry over to L_0 -testing

- Improvement on Levin's work investment strategy
Gives improvements in run time of testers for
 - Connectivity of bounded-degree graphs [Goldreich Ron 02]
 - Properties of images [R 03]
 - Multiple-input problems [Goldreich 13]
- First example of monotonicity testing problem where adaptivity helps
- Improvements to L_0 -testers for Boolean functions

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

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

Relationships: *Tolerant* L_p -Testing Models

$C_p(\mathbf{P}, \varepsilon_1, \varepsilon_2)$ = complexity of tolerant L_p -testing property \mathbf{P} with distance parameters $\varepsilon_1, \varepsilon_2$

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

For all properties \mathbf{P}

- No obvious relationship between tolerant L_1 -testing and tolerant Hamming testing
- L_p -testing for $p > 1$ is close in complexity to L_1 -testing

$$C_1(\mathbf{P}, \varepsilon_1^p, \varepsilon_2) \leq C_p(\mathbf{P}, \varepsilon_1, \varepsilon_2) \leq C_1(\mathbf{P}, \varepsilon_1, \varepsilon_2^p)$$

Relationships: *Tolerant* L_p -Testing Models

$C_p(P, \varepsilon_1, \varepsilon_2)$ = complexity of *tolerant* L_p -testing property P with distance parameters $\varepsilon_1, \varepsilon_2$

- 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

$$C_1(P, \varepsilon_1, \varepsilon_2) = C_0(P, \varepsilon_1, \varepsilon_2)$$

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

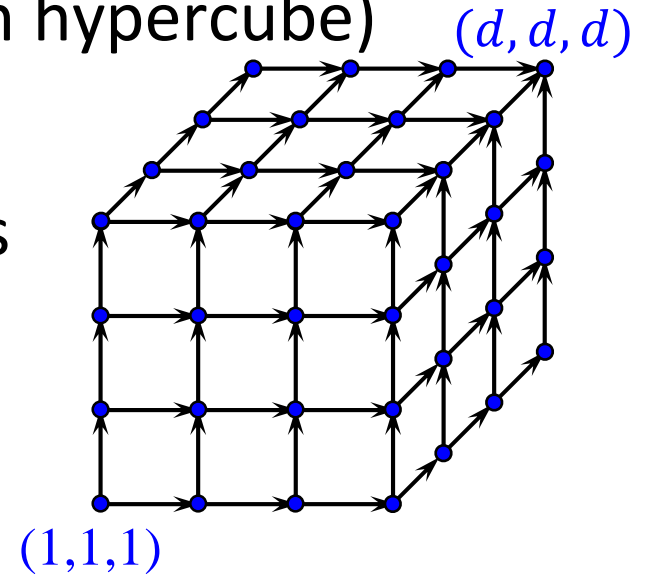
$$C_p(P, \varepsilon_1, \varepsilon_2) = C_1(P, \varepsilon_1^p, \varepsilon_2^p)$$

Our Results

Property: Monotonicity

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.



One of the most studied properties in property testing

[Ergün Kannan Kumar Rubinfeld Viswanathan, Goldreich Goldwasser Lehman Ron, Dodis Goldreich Lehman R Ron Samorodnitsky, Batu Rubinfeld White, Fischer Lehman Newman R Rubinfeld Samorodnitsky, Fischer, Halevy Kushilevitz, Bhattacharyya Grigorescu Jung R Woodruff, ..., Chakrabarty Seshadhri, Blais R Yaroslavtsev, Chakrabarty Dixit Jha Seshadhri]

Monotonicity Testers: Running Time

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

Monotonicity Testers: Running Time

f	L_0	L_p
$[n]$ $\rightarrow \{0,1\}$	$\Theta\left(\frac{1}{\epsilon}\right)$	$\Theta\left(\frac{1}{\epsilon^p}\right)$
$[n]^d$ $\rightarrow \{0,1\}$	$\Theta\left(\frac{d}{\epsilon} \cdot \log^3 \frac{d}{\epsilon}\right)$ [Dodis Goldreich Lehman R Samorodnitsky 99]	$O\left(\frac{d}{\epsilon^p} \log \frac{d}{\epsilon^p}\right)$ $\Omega\left(\frac{1}{\epsilon^p} \log \frac{1}{\epsilon^p}\right)$ for $d = 2$ nonadaptive 1-sided error
		$\Theta\left(\frac{1}{\epsilon^p}\right)$ for constant d adaptive 1-sided error

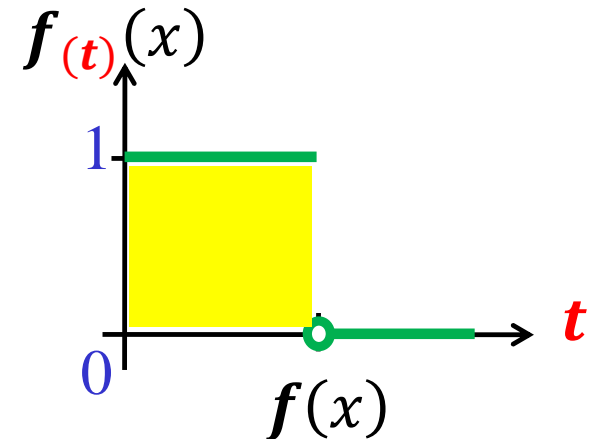
L_1 -Testing of Monotonicity

Monotonicity: Reduction to Boolean Functions

Given $f: D \rightarrow [0,1]$, a 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

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

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.

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

Because \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

Monotonicity: Using Characterization Theorem

Characterization Theorem

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

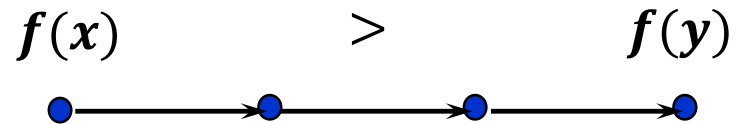
We use Characterization Theorem
to get monotonicity L_1 -testers
and tolerant testers from
standard property testers for Boolean functions.

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

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

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

Distance Approximation and Tolerant Testing

Approximating L_1 -distance to monotonicity $\pm\delta$ w. $p. \geq 2/3$

f	L_0	L_1
$[n]$ $\rightarrow [0,1]$	$\text{polylog } n \cdot \left(\frac{1}{\delta}\right)^{O(1/\delta)}$ [Saks Seshadhri 10]	$\Theta\left(\frac{1}{\delta^2}\right)$

- Time complexity of tolerant L_1 -testing for monotonicity is

$$O\left(\frac{\epsilon_2}{(\epsilon_2 - \epsilon_1)^2}\right).$$

L_1 -Testers for Other Properties

Via combinatorial characterization of L_1 -distance to the property

- Lipschitz property $f: [n]^d \rightarrow [0,1]$:

$$\Theta\left(\frac{d}{\epsilon}\right) \text{ (tight)}$$

Via (implicit) proper learning: approximate in L_1 up to error ϵ , test approximation on a random $O(1/\epsilon)$ -sample

- Convexity $f: [n]^d \rightarrow [0,1]$:

$$O\left(\epsilon^{-\frac{d}{2}} + \frac{1}{\epsilon}\right) \text{ (tight for } d \leq 2)$$

- Submodularity $f: \{0,1\}^d \rightarrow [0,1]$

$$2^{\tilde{O}\left(\frac{1}{\epsilon}\right)} + \text{poly}\left(\frac{1}{\epsilon}\right) \log d \text{ [Feldman Vondrak 13]}$$

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?