# $L_p$-Testing

## Grigory Yaroslavtsev



Joint work with Piotr Berman and Sofya Raskhodnikova

# Testing Big Data

- **Q**: How to understand properties of large data looking only at a small sample?

- **Q**: How to ignore noise and outliers?

- **Q**: How to minimize assumptions about the sample generation process?

- **Q**: How to optimize running time?
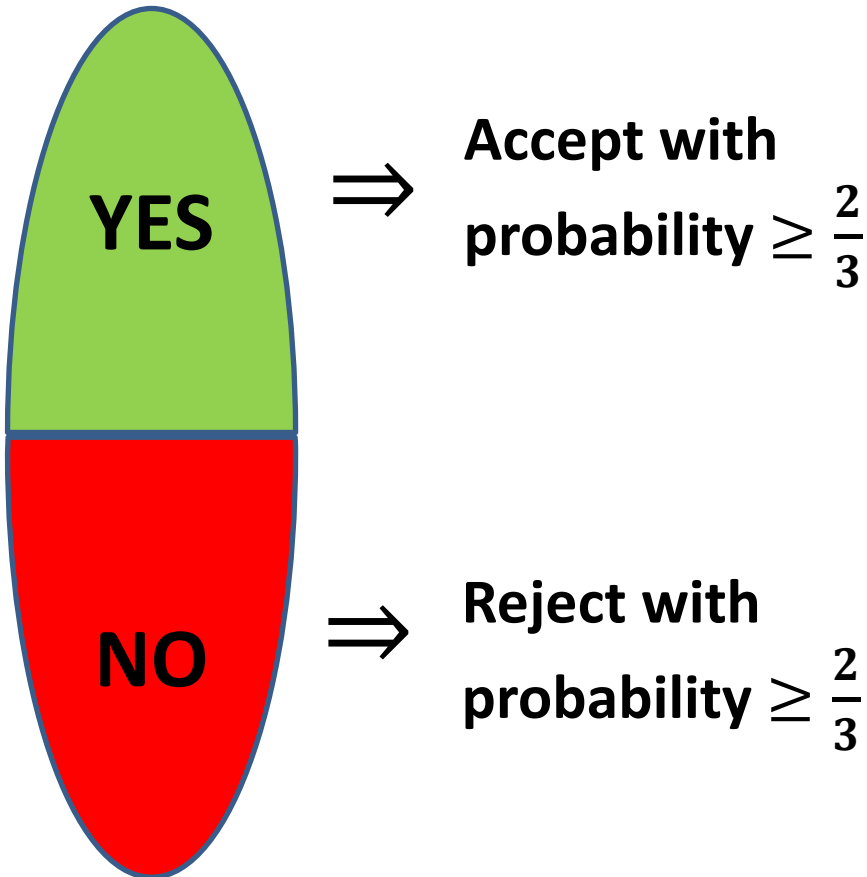
# Which stocks were growing steadily?



Data from http://finance.google.com

# Property Testing

[Goldreich, Goldwasser, Ron; Rubinfeld, Sudan]

**Randomized Algorithm**

**YES** $\Longrightarrow$ **Accept with probability** $\geq \frac{2}{3}$

**NO** $\Longrightarrow$ **Reject with probability** $\geq \frac{2}{3}$

**Property Tester**

**YES** $\Longrightarrow$ **Accept with probability** $\geq \frac{2}{3}$

$\epsilon$**-close** $\Longrightarrow$ **Don't care**

**NO** $\Longrightarrow$ **Reject with probability** $\geq \frac{2}{3}$

$\epsilon$**-close :** $\leq \epsilon$ fraction has to be changed to become **YES**

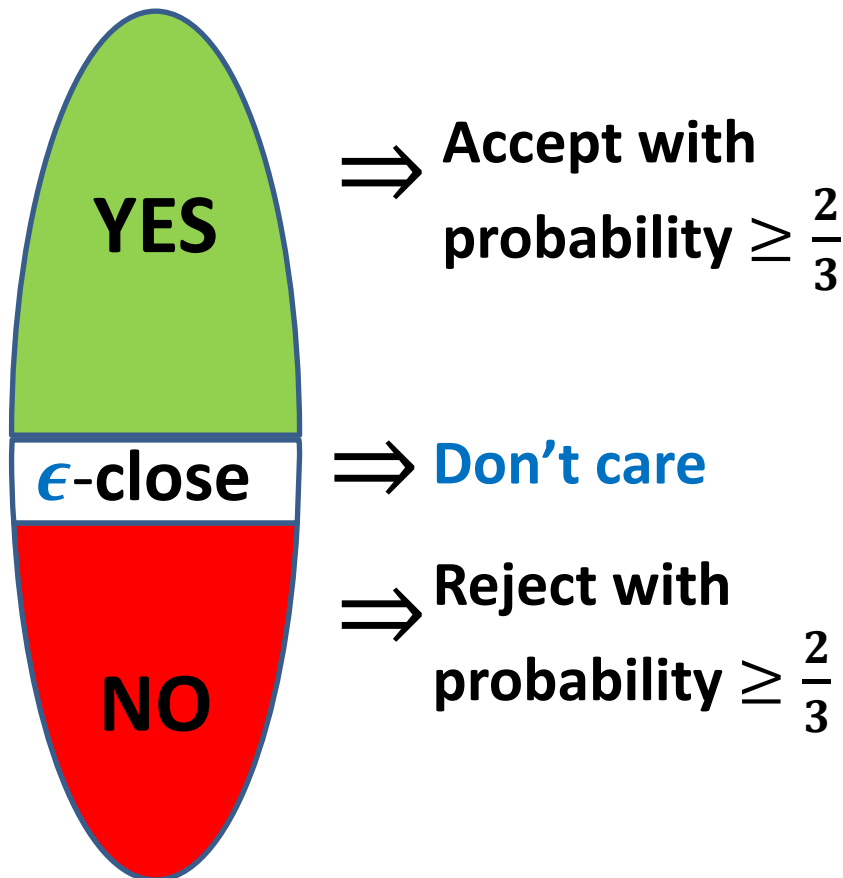# Tolerant Property Testing

[Parnas, Ron, Rubinfeld]

**Property Tester**

**YES** $\Rightarrow$ **Accept with** probability $\geq \frac{2}{3}$

$\epsilon$-**close** $\Rightarrow$ **Don't care**

**NO** $\Rightarrow$ **Reject with** probability $\geq \frac{2}{3}$

**Tolerant Property Tester**

**YES** $\Rightarrow$ **Accept with** probability $\geq \frac{2}{3}$

$\epsilon_1$-**close**

$(\epsilon_1, \epsilon_2)$-**close** $\Rightarrow$ **Don't care**

**NO** $\Rightarrow$ **Reject with** probability $\geq \frac{2}{3}$

$\epsilon$-**close** : $\leq \epsilon$ fraction has to be changed to become **YES**

# Which stocks were growing steadily?



Data from http://finance.google.com

# Tolerant "$L_1$ Property Testing"

- $f : \{1, \dots, n\} \to [0,1]$
- $P$ = class of monotone functions
- $dist_1(f, P) = \dfrac{\min\limits_{g \in P} |f - g|_1}{n}$
- $\epsilon$-close: $dist_1(f, P) \leq \epsilon$

**Tolerant "$L_1$ Property Tester"**

YES $\Rightarrow$ **Accept with probability** $\geq \dfrac{2}{3}$

$\epsilon_1$-close

$(\epsilon_1, \epsilon_2)$-close $\Rightarrow$ **Don't care**

NO $\Rightarrow$ **Reject with probability** $\geq \dfrac{2}{3}$

# New $L_p$-Testing Model for Real-Valued Data

- **Generalizes** standard Hamming testing

- For $p > 0$ still have a **probabilistic interpretation**:
$$d_p(f, g) = (\mathbf{E}[|f - g|^p])^{1/p}$$

- 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 (sublinear algorithm for isotonic regression)

   ❖ Our $L_p$-testers **beat lower bounds** for Hamming testers
   ❖ **Simple algorithms** backed up by involved analysis
   ❖ Uniformly sampled (or **easy to sample**) data suffices

3. Nearly tight lower bounds

# Implications for Hamming Testing

Some techniques/results carry over to Hamming testing

- Improvement on **Levin's work investment strategy**
  - Connectivity of bounded-degree graphs [Goldreich, Ron '02]
  - Properties of images [Raskhodnikova '03]
  - Multiple-input problems [Goldreich '13]

- First example of **monotonicity testing** problem where **adaptivity helps**
- Improvements to Hamming testers for Boolean functions

# Definitions

- $f: D \to [0,1]$ (D = finite domain/poset)

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

- $\|f\|_0 =$ Hamming weight (# of non-zero values)

- Property $P$ = class of functions (monotone, convex, linear, Lipschitz, …)

- $dist_p(f, P) = \dfrac{\min\limits_{g \in P} \|f - g\|_p}{\|1\|_p}$

# Relationships: $L_p$-Testing

$Q_p(P, \epsilon)$ = query complexity of $L_p$-testing property $P$ at distance $\epsilon$

- $Q_1(P, \epsilon) \leq Q_0(P, \epsilon)$
- $Q_1(P, \epsilon) \leq Q_2(P, \epsilon)$ (Cauchy-Shwarz)
- $Q_1(P, \epsilon) \geq Q_2(P, \sqrt{\epsilon})$

Boolean functions $f: D \to \{0,1\}$

$Q_0(P, \epsilon) = Q_1(P, \epsilon) = Q_2(P, \sqrt{\epsilon})$

# Relationships: Tolerant $L_p$-Testing

$Q_p(P, \epsilon_1, \epsilon_2)$ = query complexity of tolerant $L_p$-testing property $P$ with distance parameters $\epsilon_1, \epsilon_2$

- No general relationship between tolerant $L_1$-testing and tolerant Hamming testing

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

$$Q_1(P, \varepsilon_1^p, \varepsilon_2) \leq Q_p(P, \varepsilon_1, \varepsilon_2) \leq Q_1(P, \varepsilon_1, \varepsilon_2^p)$$

For Boolean functions $f: D \to \{0,1\}$

$$Q_0(P, \varepsilon_1, \varepsilon_2) = Q_1(P, \varepsilon_1, \varepsilon_2) = Q_p(P, \epsilon_1^{1/p}, \varepsilon_2^{1/p})$$

# Our Results: Testing Monotonicity

- Hypergrid ($D = [n]^d$)

| | $L_0$ | $L_1$ |
|---|---|---|
| Upper bound | $O\left(\dfrac{d \log n}{\epsilon}\right)$ <br><br> [Dodis et al. '99,..., Chakrabarti, Seshadhri '13] | $O\left(\dfrac{d}{\epsilon} \log \dfrac{d}{\epsilon}\right)$ |
| Lower bound | $\Omega\left(\dfrac{d \log n}{\epsilon}\right)$ <br><br> [Dodis et al.'99..., Chakrabarti, Seshadhri '13] | $\Omega\left(\dfrac{1}{\epsilon} \log \dfrac{1}{\epsilon}\right)$ <br><br> Non-adaptive 1-sided error |

- $2^{O(d)}/\epsilon$ **adaptive** tester for Boolean functions

# Monotonicity: Key Lemma

- M = class of monotone functions
- Boolean slicing operator $\boldsymbol{f_y}: D \rightarrow \{0,1\}$

$$\boldsymbol{f_y}(x) = 1, \text{ if } \boldsymbol{f}(x) \geq \boldsymbol{y},$$

$$\boldsymbol{f_y}(x) = 0, \text{ otherwise.}$$

- **Theorem:**

$$dist_1(\boldsymbol{f}, M) = \int_0^1 dist_0(\boldsymbol{f_y}, M)d\boldsymbol{y}$$

# Proof sketch: slice and conquer

1) Closest monotone function with **minimal $L_1$-norm** is **unique** (can be denoted as an operator $M_f^1$).

2) $\left\| f - g \right\|_1 = \int_0^1 \left\| f_y - g_y \right\| dy$

3) $M_f^1$ and $f_y$ commute: $\left( M_f^1 \right)_y = M^1_{(f_y)}$

$$dist_1(f, M) = \overset{\textbf{1)}}{\frac{\left\| f - M_f^1 \right\|_1}{|D|}} = \overset{\textbf{2)}}{\frac{\int_0^1 \left\| f_y - (M_f^1)_y \right\|_1 dy}{|D|}} = \overset{\textbf{3)}}{}$$

$$= \frac{\int_0^1 \left\| f_y - M^1_{(f_y)} \right\|_1 dy}{|D|} = \int_0^1 dist_0(f_y, M) dy$$

# $L_1$-Testers from Boolean Testers

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

Proof:

$$f(x) \qquad > \qquad f(y)$$

- 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 $f: D \to [0,1]$ 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 $f_{(t^*)}$

$$f_{(t^*)}(x) = 1, f_{(t^*)}(y) = 0$$
$$\Downarrow$$
$$f(x) > f(y)$$

# Distance Approximation and Tolerant Testing

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

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

- Time complexity of tolerant $L_1$-testing for monotonicity is

$$O\left(\frac{\varepsilon_2}{(\varepsilon_2 - \varepsilon_1)^2}\right)$$

  – Better dependence than what follows from distance appoximation for $\epsilon_2 \ll 1$

  – Improves $\tilde{O}\left(\dfrac{1}{\delta^2}\right)$ adaptive distance approximation of [Fattal,Ron'10] for Boolean functions

# $L_1$-Testers for Other Properties

Via combinatorial characterization of $L_1$-distance to the property

- Lipschitz property $\boldsymbol{f} \colon [\boldsymbol{n}]^{\boldsymbol{d}} \to [0,1]$:

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

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

- Convexity $\boldsymbol{f} \colon [\boldsymbol{n}]^{\boldsymbol{d}} \to [0,1]$:

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

- Submodularity $\boldsymbol{f} \colon \{0,1\}^{\boldsymbol{d}} \to [0,1]$

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

# Open Problems

- All our algorithms for for $p > 1$ were obtained directly from $L_1$-testers.

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

- Our complexity for $L_p$-testing convexity grows exponentially with d

Is there an $L_p$-testing algorithm for convexity with subexponential dependence on the dimension?

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

Is there a better adaptive tester?

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

Tolerant testers for higher dimensions?

Other properties?