

Testing Monotone Continuous Distributions on High-dimensional Real Cubes

Artur Czumaj

Department of Computer Science &
DIMAP (Centre for Discrete Maths and its Applications)
University of Warwick

Joint work with Michal Adamaszek & Christian Sohler

Testing properties of distributions

- General question:
 - Test if a given probability distribution has a given property

Distribution is available by accessing only samples drawn from the distribution

Examples:

- is given distribution uniform?
- are two distributions identical?
- are two distributions independent?

Testing properties of distributions

Lots of research in statistics

Some recent research in algorithms

- Typical result:
 - Given a probability distribution on n points, we can test if it's uniform after seeing $\sim\sqrt{n}$ random samples

[Batu et al '01]

Testing = distinguish between uniform distribution and distributions which are ϵ -far from uniform

ϵ -far from uniform:

error probab. $\leq 1/3$

$$\sum_{x \in \mathcal{X}} \left| \Pr[x] - \frac{1}{n} \right| \geq \epsilon$$

Testing properties of distributions

- Typical result:
 - Given a probability distribution on n points, we can test if it's uniform after seeing $\sim\sqrt{n}$ random samples
- [Batu et al '01]
- Similar bounds for testing
 - if a distribution is monotone
 - if two distributions are independent
 - ...

Testing properties of distributions

- Typical result:
 - Given a probability distribution on n points, we can test if it's uniform after seeing $\sim\sqrt{n}$ random samples

[Batu et al '01]

Many properties of distributions can be tested in time sublinear in the domain/support size (typically, with $n^{O(1)}$ samples)

Testing properties of distributions

- Typical result:
 - Given a probability distribution on n points, we can test if it's uniform after seeing $\sim\sqrt{n}$ random samples
- [Batu et al '01]
- What if distribution has infinite support?
 - Continuous probability distributions?

Testing properties of continuous distributions

- Typical result:
 - Given a probability distribution on n points, we can test if it's uniform after seeing $\Theta(\sqrt{n})$ random samples
 - $\Theta(\sqrt{n})$ random samples are necessary
- Given a continuous probability distribution on $[0,1]$, can we test if it's uniform?
- Impossible
 - Follows from lower bound for discrete case with $n \rightarrow \infty$

Testing properties of continuous distributions

- What can be tested?
- First question:
test if the distribution is indeed continuous

Testing properties of continuous distributions

- Dual question:
Test if a probability distribution is **discrete**
- Prob. distribution **D** on Ω is **discrete** on **N** points if there is a set $X \subseteq \Omega$, $|X| \leq N$, st. **$\Pr_D[X]=1$**
- **D** is **ϵ -far from discrete** on **N** points if

$$\forall X \subseteq \Omega, |X| \leq N$$

$$\Pr_D[X] \leq 1-\epsilon$$

Testing if distribution is discrete on N points

- We repeatedly draw random points from D
- All what can we see:
 - Count frequency of each point
 - Count number of points drawn

For some D (eg, uniform or close):

- we need $\sim (\sqrt{N})$ to see first multiple occurrence

Gives a hope that can be solved in sublinear-time

Shows that we cannot be better than $\sim (\sqrt{N})$

Testing if distribution is discrete on N points

Raskhodnikova et al '07 (Valiant'08):

Distinct Elements Problem:

- D discrete with each element with prob. $\geq 1/N$
- Estimate the support size

$\Omega(N^{1-o(1)})$ queries are needed to distinguish instances with $\leq N/100$ and $\geq N/11$ support size

Key property:

- two distributions that have identical first $\log^{\Theta(1)}N$ moments
- their expected frequencies up to $\log^{\Theta(1)}N$ are identical

Testing if distribution is discrete on N points

Raskhodnikova et al '07 (Valiant'08):

Distinct Elements Problem:

- D discrete with each element with prob. $\geq 1/N$
- Estimate the support size

$\Omega(N^{1-o(1)})$ queries are needed to distinguish instances with $\leq N/100$ and $\geq N/11$ support size

Corollary: Testing if a distribution is discrete on N points requires $\Omega(N^{1-o(1)})$ samples

Testing if distribution is discrete on N points

- We repeatedly draw random points from D
- All what can we see:
 - Count frequency of each point
 - Count number of points drawn
- Can we get $O(N)$ time?

Testing if distribution is discrete on N points

- Testing if a distribution is discrete on N points:

• Draw a sample $S = (s_1, \dots, s_t)$ with $t = 2N/\epsilon$
• If S has more than N distinct elements
then REJECT
else ACCEPT

- If D is discrete on N points then we will accept D
- We only have to prove that
 - if D is ϵ -far from discrete on N points, then we will reject with probability $>2/3$

Testing if distribution is discrete on N points

- Testing if a distribution is discrete on N points:

- Draw a sample $S = (s_1, \dots, s_t)$ with $t = 2N/\epsilon$
- If S has more than N distinct elements then **REJECT** else **ACCEPT**

D is ϵ -far from discrete on N points, then reject with prob $> 2/3$

D is ϵ -far from discrete on N points \Rightarrow

D is ϵ -far from discrete on N points iff $\forall X \subseteq \Omega$, if $|X| \leq N$ then $Pr_D[\Omega \setminus X] \geq \epsilon$

- Assuming that we haven't chosen n points yet, we choose a new point with probability at least ϵ

After $(1 + o(1))N/\epsilon$ samples, we choose $N + 1$ points with prob. ≥ 0.99

Testing if distribution is discrete on N points

- Testing if a distribution is discrete on N points:

- Draw a sample $S = (s_1, \dots, s_t)$ with $t = 2N/\epsilon$
- If S has more than N distinct elements
then REJECT
else ACCEPT

Can we do better (if we only count distinct elements)?

D: has 1 point with prob. $1-4\epsilon$ and $2N$ points with prob. $2\epsilon/N$

D is ϵ -far from discrete on N points

We need $\Omega(N/\epsilon)$ samples to see at least N points

Open problem

What is the complexity of testing if
distribution is discrete on N points?

Upper bound: $O(N/\epsilon)$

Lower bound: $\Omega(N^{1-o(1)})$

Open problem: close the gap

Testing continuous probability distributions

- What can we test efficiently?
 - Complexity for discrete distributions should be “independent” on the support size
- Uniform distribution ... under some conditions
- Rubinfeld & Servedio'05:
 - testing monotone distributions for uniformity

Testing uniform distributions (discrete)

Rubinfeld & Servedio'05:

- Testing monotone distributions for uniformity

D: distribution on **n-dimensional** cube; $D:\{0,1\}^n \rightarrow \mathbf{R}$

$x, y \in \{0,1\}^n$, $x \preceq y$ iff $\forall i: x_i \leq y_i$

D is monotone if $x \preceq y \rightarrow \Pr[x] \leq \Pr[y]$

Goal: test if a monotone distribution is uniform

Rubinfeld & Servedio'05:

Testing if a monotone distribution on n-dimensional binary cube is uniform:

- Can be done with $O(n \log(1/\epsilon)/\epsilon^2)$ samples
- Requires $\Omega(n/\log^2 n)$ samples

Testing continuous distributions

Rubinfeld & Servedio'05:

- Testing monotone distributions for uniformity

D: distribution on n-dimensional cube; $D:\{0,1\}^n \rightarrow \mathbf{R}$

$x, y \in \{0,1\}^n$, $x \preceq y$ iff $\forall i: x_i \leq y_i$

D is monotone if $x \preceq y \rightarrow \Pr[x] \leq \Pr[y]$

Goal: test if a monotone distribution is uniform

D: distribution on n-dimensional cube;

density function $f:[0,1]^n \rightarrow \mathbf{R}$

$x, y \in [0,1]^n$, $x \preceq y$ iff $\forall i: x_i \leq y_i$

D is monotone if $x \preceq y \rightarrow f(x) \leq f(y)$

Testing continuous distributions

Lower bounds holds for n-dimensional real cubes
Upper bound: ???

Rubinfeld & Servedio'05:

Testing if a monotone distribution on n-dimensional
binary cube is uniform:

- Can be done with $O(n \log(1/\epsilon)/\epsilon^2)$ samples
- Requires $\Omega(n/\log^2 n)$ samples

Testing monotone distributions for uniformity

D is ϵ -far from uniform if $\frac{1}{2} \int_{x \in \mathcal{X}} |f(x) - 1| dx \geq \epsilon$

L_1 distance between f and uniform distribution

To test uniformity, we need to characterize monotone distributions that are ϵ -far from uniform

On the high level:

- we follow approach of Rubinfeld & Servedio'05;
- details are different

Testing monotone distributions for uniformity

D is ϵ -far from uniform if $\frac{1}{2} \int_{x \in \mathcal{X}} |f(x) - 1| dx \geq \epsilon$

Key Technical Lemma:

Let $g: [0,1]^n \rightarrow \mathbf{R}$ be a monotone function with $\int_{\mathcal{X}} g(x) dx = 0$ then

$$\int_{\mathcal{X}} \|x\|_1 \cdot g(x) dx \geq \frac{1}{4} \int_{\mathcal{X}} |g(x)| dx$$

Key Lemma follows from Key Technical Lemma with $g(x) = f(x) - 1$

Key Lemma:

If D is a monotone distribution on $[0,1]^n$ with density function f and which is ϵ -far from uniform then

$$E_f[\|x\|_1] = \int_{\mathcal{X}} \|x\|_1 \cdot f(x) dx \geq \frac{n}{2} + \frac{\epsilon}{2}$$

Testing monotone distributions for uniformity

Key Lemma:

If D is a monotone distribution on $[0,1]^n$ with density function f and which is ϵ -far from uniform then

$$E_f[\|x\|_1] = \int_x \|x\|_1 \cdot f(x) dx \geq \frac{n}{2} + \frac{\epsilon}{2}$$

Uniform distribution:

If D is uniform on $[0,1]^n$ with density function f then

$$E_f[\|x\|_1] = \int_x \|x\|_1 \cdot f(x) dx = \frac{n}{2}$$

Testing monotone distributions for uniformity

Key Lemma:

If D is a monotone distribution on $[0,1]^n$ with density function f and which is ϵ -far from uniform then

$$E_f[\|x\|_1] = \int_x \|x\|_1 \cdot f(x) dx \geq \frac{n}{2} + \frac{\epsilon}{2}$$

$$s = cn/\epsilon^2$$

Repeat 20 times

Draw a sample $S=(x_1, \dots, x_s)$ from $[0,1]^n$

If $\sum_i \|x_i\|_1 \geq s (n/2 + \epsilon/4)$ then REJECT and exit

ACCEPT

Testing monotone distributions for uniformity

Theorem:

The algorithm below tests if D is uniform.

Its complexity is $O(n/\epsilon^2)$.

Slightly better bound than the one by RS'05

$$s = cn/\epsilon^2$$

Repeat 20 times

Draw a sample $S=(x_1, \dots, x_s)$ from $[0,1]^n$

If $\sum_i \|x_i\|_1 \geq s(n/2 + \epsilon/4)$ then REJECT and exit

ACCEPT

Testing monotone distributions for uniformity

$$s = cn/\epsilon^2$$

Repeat 20 times

Draw a sample $S=(x_1, \dots, x_s)$ from $[0,1]^n$

If $\sum_i \|x_i\|_1 \geq s(n/2 + \epsilon/4)$ then REJECT and exit

ACCEPT

Lemma 1: If D is uniform then

$$\Pr[\sum_i \|x_i\|_1 \geq s(n/2 + \epsilon/4)] \leq 0.01$$

Easy application of Chernoff bound

Lemma 2: If D is ϵ -far from uniform then

$$\Pr[\sum_i \|x_i\|_1 < s(n/2 + \epsilon/4)] \leq 12/13$$

By Key Lemma + Feige lemma

Testing monotone distributions for uniformity

Key Technical Lemma:

Let $g:[0,1]^n \rightarrow \mathbf{R}$ be a monotone function with $\int_{\mathbf{x}} g(\mathbf{x}) d\mathbf{x} = 0$ then

$$\int_{\mathbf{x}} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} \geq \frac{1}{4} \int_{\mathbf{x}} |g(\mathbf{x})| d\mathbf{x}$$

Why such a bound:

Tight for $g(\mathbf{x}) = \text{sgn}(x_1 - 1/2)$

$$\int_{x_1 > 1/2} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} = \frac{1}{2} \int_{x_1 > 1/2} (x_1 + \dots + x_n) d\mathbf{x} = \frac{1}{2} \left(\frac{3}{4} + \frac{1}{2} + \dots + \frac{1}{2} \right) dx = \frac{n}{4} + \frac{1}{8} .$$

Similarly,

$$\int_{x_1 < 1/2} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} = \frac{1}{2} \left(\frac{1}{4} + \frac{1}{2} + \dots + \frac{1}{2} \right) = \frac{n}{4} - \frac{1}{8} ,$$

and hence,

$$\int_{\mathbf{x}} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} = \int_{x_1 > 1/2} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} - \int_{x_1 < 1/2} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} = \frac{1}{4} = \frac{1}{4} \cdot \int_{\mathbf{x}} |g(\mathbf{x})| d\mathbf{x} .$$

Testing monotone distributions for uniformity

Key Technical Lemma:

Let $g:[0,1]^n \rightarrow \mathbf{R}$ be a monotone function with $\int_{\mathbf{x}} g(\mathbf{x}) d\mathbf{x} = 0$ then

$$\int_{\mathbf{x}} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} \geq \frac{1}{4} \int_{\mathbf{x}} |g(\mathbf{x})| d\mathbf{x}$$

Testing monotone distributions for uniformity

Let $P = \{\mathbf{x} : g(\mathbf{x}) \geq 0\}$ and $N = \{\mathbf{x} : g(\mathbf{x}) < 0\}$. Consider:

$$\int_{\mathbf{x} \in N} \int_{\mathbf{y} \in P} |g(\mathbf{x}) - g(\mathbf{y})| \, d\mathbf{y} \, d\mathbf{x} .$$

For $g(\mathbf{x}) < 0 \cdot g(\mathbf{y})$, we have $|g(\mathbf{x}) - g(\mathbf{y})| = |g(\mathbf{x})| + |g(\mathbf{y})|$.

Moreover $\int_{\mathbf{x} \in N} |g(\mathbf{x})| \, d\mathbf{x} = \int_{\mathbf{y} \in P} |g(\mathbf{y})| \, d\mathbf{y} = \frac{1}{2} \int_{\mathbf{x}} |g(\mathbf{x})| \, d\mathbf{x}$.

Hence:

$$\begin{aligned} &= \int_{\mathbf{x} \in N} \int_{\mathbf{y} \in P} (|g(\mathbf{x})| + |g(\mathbf{y})|) \, d\mathbf{y} \, d\mathbf{x} = \int_{\mathbf{y} \in P} \int_{\mathbf{x} \in N} |g(\mathbf{x})| \, d\mathbf{x} \, d\mathbf{y} + \int_{\mathbf{x} \in N} \int_{\mathbf{y} \in P} |g(\mathbf{y})| \, d\mathbf{y} \, d\mathbf{x} \\ &= \frac{1}{2} \int_{\mathbf{y} \in P} \int_{\mathbf{x}} |g(\mathbf{x})| \, d\mathbf{x} \, d\mathbf{y} + \frac{1}{2} \int_{\mathbf{x} \in N} \int_{\mathbf{y}} |g(\mathbf{y})| \, d\mathbf{y} \, d\mathbf{x} = \frac{1}{2} \int_{\mathbf{y}} \int_{\mathbf{x}} |g(\mathbf{x})| \, d\mathbf{x} \, d\mathbf{y} = \frac{1}{2} \int_{\mathbf{x}} |g(\mathbf{x})| \, d\mathbf{x} . \end{aligned}$$

Since every pair (\mathbf{x}, \mathbf{y}) can satisfy at most one of the conditions $(\mathbf{x}, \mathbf{y}) \in P \times N$ and $(\mathbf{x}, \mathbf{y}) \in N \times P$, we obtain:

$$\int_{\mathbf{x} \in N} \int_{\mathbf{y} \in P} |g(\mathbf{x}) - g(\mathbf{y})| \, d\mathbf{y} \, d\mathbf{x} \cdot \frac{1}{2} \int \int_{\mathbf{x}, \mathbf{y}} |g(\mathbf{x}) - g(\mathbf{y})| \, d\mathbf{y} \, d\mathbf{x} .$$

Hence:

$$\frac{1}{2} \int_{\mathbf{x}} |g(\mathbf{x})| \, d\mathbf{x} = \int_{\mathbf{x} \in N} \int_{\mathbf{y} \in P} |g(\mathbf{x}) - g(\mathbf{y})| \, d\mathbf{x} \, d\mathbf{y} \cdot \frac{1}{2} \int \int_{\mathbf{x}, \mathbf{y}} |g(\mathbf{x}) - g(\mathbf{y})| \, d\mathbf{x} \, d\mathbf{y} .$$

Testing monotone distributions for uniformity

By considering all the possible relative placements of \mathbf{x} and \mathbf{y} within $[0, 1]^n$ and splitting the domain accordingly, one can prove that

$$\int \int_{\mathbf{x}, \mathbf{y}} |g(\mathbf{x}) - g(\mathbf{y})| dy dx \cdot \int \int_{\mathbf{x} \prec \mathbf{y}} \left(\sum_{(\mathbf{u}, \mathbf{v}) \in D(\mathbf{x}, \mathbf{y})} |g(\mathbf{u}) - g(\mathbf{v})| \right) dy dx ,$$

where $D(\{0, 1\}^n)$ is the set of all **main diagonals** of **discrete** cube $\{0, 1\}^n$:

$$D(\{0, 1\}^n) = \{(\mathbf{x}, \mathbf{y}) \in \{0, 1\}^n \times \{0, 1\}^n : x_i = 1 - y_i \text{ for every } i\}$$

Testing monotone distributions for uniformity

Key Technical Lemma:

Let $g: [0,1]^n \rightarrow \mathbf{R}$ be a monotone function with $\int_{\mathbf{x}} g(\mathbf{x}) d\mathbf{x} = 0$ then

$$\int_{\mathbf{x}} \|\mathbf{x}\|_1 \cdot g(\mathbf{x}) d\mathbf{x} \geq \frac{1}{4} \int_{\mathbf{x}} |g(\mathbf{x})| d\mathbf{x}$$

Key inequalities in the proof:

$$\begin{aligned} & \frac{1}{4} \int_{\mathbf{x}} |g(\mathbf{x})| d\mathbf{x} \cdot \frac{1}{4} \int \int_{\mathbf{x}, \mathbf{y}} |g(\mathbf{x}) - g(\mathbf{y})| d\mathbf{x} d\mathbf{y} \\ & \cdot \frac{1}{4} \int \int_{\mathbf{x} < \mathbf{y}} \left(\sum_{(\mathbf{u}, \mathbf{v}) \in D(\mathbf{x}, \mathbf{y})} |g(\mathbf{u}) - g(\mathbf{v})| \right) d\mathbf{x} d\mathbf{y} \\ & \cdot \frac{1}{2} \sum_{i=1}^n \int \int_{\mathbf{x} < \mathbf{y}} \left(\sum_{(\mathbf{u}, \mathbf{v}) \in E_i(\mathbf{x}, \mathbf{y})} |g(\mathbf{u}) - g(\mathbf{v})| \right) d\mathbf{x} d\mathbf{y} \\ & \cdot \frac{1}{2} \sum_{i=1}^n \int_{\mathbf{x}} (2x_i - 1) g(\mathbf{x}) d\mathbf{x} \\ & \cdot \int_{\mathbf{x}} \|\mathbf{x}\|_1 g(\mathbf{x}) d\mathbf{x} \end{aligned}$$

Testing monotone continuous distributions

Rubinfeld & Servedio'05:

Testing if a monotone distribution on n -dimensional

binary cube is uniform:

- Can be done with $O(n \log(1/\epsilon)/\epsilon^2)$ samples
- Requires $\Omega(n/\log^2 n)$ samples

Here:

Testing if a monotone distribution on n -dimensional

continuous cube is uniform :

- Can be done with $O(n/\epsilon^2)$ samples
- (Requires $\Omega(n/\log^2 n)$ samples)

Testing monotone continuous distributions

Further extension/application:

Testing if a monotone distribution on n-dimensional discrete cube $\{0,1,2,\dots,k\}^n$ is uniform:

- Can be done with $O(n/\epsilon^2)$ samples

Here:

Testing if a monotone distribution on n-dimensional **continuous** cube is uniform :

- Can be done with $O(n/\epsilon^2)$ samples
- (Requires $\Omega(n/\log^2 n)$ samples)

Conclusions

- Testing distributions on infinite/uncountable support is different from testing discrete distributions
 - Continuous distributions are harder
- Challenge: understand when it's possible to test
 - Usually some additional conditions are to be imposed
- Tight(er) bounds?

Conclusions

- Continuous distributions are harder
- Is the L_1 -norm the right one?
 - It doesn't work well for continuous distributions
- Earth mover norm?
 - How much mass has to be moved and how far to obtain a given distribution
 - Ba, Nguyen, Nguyen, Rubinfeld 2009:
 - Testing uniformity on $[0,1]$ can be done in time $f(1/\epsilon)$
 - Framework (holds for a variety of properties):
reduction to the problem on the support of size $f(1/\epsilon)$