

## Tackling the Challenges of Big Data

### Big Data Analytics

### Fast Algorithms I

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology




---



---



---



---



---



---



---



---

## Tackling the Challenges of Big Data

### Big Data Analytics

### Fast Algorithms I:

#### Introduction

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology




---



---



---



---



---



---



---



---

## What to do About REALLY Big Data?



Spiral galaxy NGC 3982. Credit: NASA, ESA, and the Hubble Heritage Team (STScI/AURA)

Read more: <http://www.universetoday.com/30168/galaxies/#ixzz2fegcrxox>

SOURCE: Google




---



---



---



---



---



---



---



---

# No Time

**What can we hope to do without viewing most of the data?**

---



---



---



---



---



---



---

## Small World Phenomenon

- The social network is a graph:
  - “node” is a person
  - “edge” between people that know each other
- “6 degrees of separation”
  - are all pairs of people connected by path of distance at most 6?

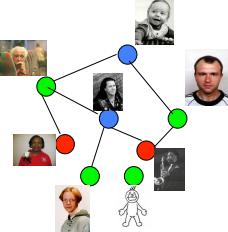


Photo & Image Source: Google

---



---



---



---



---



---



---

## Really Big Data

- Impossible to access all of it?
- Accessible data is too enormous to be viewed by a single individual?
- Once accessed, data can change?

---



---



---



---



---



---



---

## The Gold Standard

- linear time algorithms:
  - for inputs encoded by  $n$  bits/words, allow  $cn$  time steps (constant  $c$ )
- Inadequate...



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

---

## What Can We Hope To Do Without Viewing Most of the Data?

- **Can't answer "for all" or "exactly" type statements:**
  - exactly how many individuals on earth are left-handed?
  - are *all* individuals connected by at most 6 degrees of separation?
- **Compromise?**
  - *approximately* how many individuals on earth are left-handed?
  - is there a *large group* of individuals connected by at most 6 degrees of separation?



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

---

## What Types of Approximation?

- **Property Testing** – distinguish data that has a certain property from data that is *far* from having the property
- **Classical Approximation** – approximate correct output of computational problem



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I:**  
Introduction

**THANK YOU**



**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I**

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I:**  
Property Testing Algorithms

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

## What Types of Approximation?

### Property Testing

Distinguish data that *has* a certain property from data that is *far* from having the property.



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



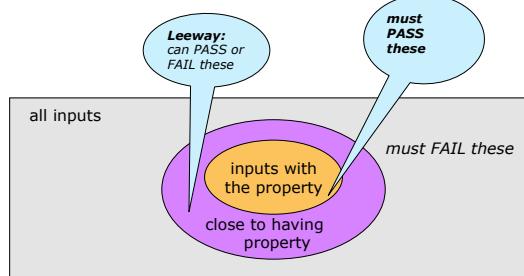
---



---

## Property Testing:

**Quickly** distinguish inputs that *have* the property from those that are *far* from having the property.



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology

---



---



---



---



---



---



---



---



---



---

## Why Property Testing?

- Can often answer such questions much faster
- May be the natural question to ask
  - When some "noise" always present
  - When data constantly changing
  - Gives fast sanity check to rule out very "bad" inputs
  - Model selection problem in machine learning



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Which Properties?

- **Properties of any object, e.g.,**

- *Functions, graphs, strings, matrices, codewords, ...*

- **Examples:**

- *clusterability, small diameter graph, increasing order, close to a codeword, linear or low degree polynomial function*
  - *Lots and lots more...*



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---

## Property Testing

- **Model must specify**

- *representation of object and allowable queries*
  - *notion of close/far, e.g., number of bits/words that need to be changed*



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---

## Examples

- Can test if the social network has 6 degrees of separation in CONSTANT TIME
- Can test if data is clusterable in CONSTANT TIME

[Parnas Ron][Alon Dar Parnas Ron]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---

## Constructing a Property Tester:

- Find characterization of property that is
  - Efficiently testable
  - Robust
    - objects that have the property **satisfy** characterization,
    - and objects **far from having** the property are **unlikely** to PASS

Usually the  
bigger  
challenge



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Example: 6 Degrees of Separation

- **A “bad” testing characterization:**
  - For every node, *all* other nodes within distance 6.
- **Another bad one:**
  - For *most* nodes, *all* other nodes within distance 6.
- **Good characterization:** [Parnas Ron]
  - For *most* nodes, there are *many other* nodes within distance 6.



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Many More Properties Studied!

- Graphs, functions, point sets, strings, ...
- Amazing characterizations of problems testable in graph and function testing models in constant time!



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I:**  
Property Testing Algorithms

**THANK YOU**



---

---

---

---

---

---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I**

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I:**  
Sublinear Time Approximation Algorithms

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

## “Classical” Approximation

- **Output number close to value of the optimal solution (not enough time to construct a solution)**
  - **Some examples:**
    - Minimum spanning tree,
    - vertex cover,
    - max cut,
    - positive linear program,

MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## What is Close?

- **c-multiplicative approximation**
    - output a number that is within a multiplicative factor of  $c$  of the best solution
      - e.g.,  $output < 2 \times OPTIMAL$
  - **c-additive approximation**
    - output a number that is additively within  $c$  of best solution
      - e.g.,  $output < OPTIMAL + c$

MIF PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Example: Vertex Cover

- Given graph  $G(V,E)$ , a vertex cover (VC) C is a subset of V such that it “touches” every edge.
  - What is minimum size of a vertex cover?
    - NP-complete
    - Polynomial time multiplicative 2-approximation based on relationship of VC and maximal matching

MIT PROFESSIONAL EDUCATION

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## "Classical" Approximation Examples

- Can get **CONSTANT TIME** approximation for vertex cover on sparse graphs!
  - Output  $y$  which is at most  $2 \cdot OPT + \epsilon n$
- **How?**
- **Oracle reduction framework [Parnas Ron]**
  - Construct "oracle" that tells you if node  $u$  in 2-approximate vertex cover
  - Use oracle + standard sampling to estimate size of cover

**But how do you implement the oracle?**



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Implementing The Oracle – Two Approaches:

- Sequentially simulate computations of a fast distributed algorithm [Parnas Ron]
- Figure out what a greedy maximal matching algorithm would do on  $u$  [Nguyen Onak]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Vertex Cover & Maximal Matching

- **Maximal Matching:**
  - $M \subseteq E$  is a matching if no node is involved in more than one edge.
  - $M$  is a maximal matching if adding any edge violates the matching property
- **Well known fact: Nodes in any maximal matching  $M$  give a pretty good Vertex Cover!**
  - That is,  $|M| \leq VC \leq 2 |M|$



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Greedy Algorithm for Maximal Matching

### Algorithm:

- Consider every edge  $(u,v)$  in arbitrary order:
  - If neither of  $u$  or  $v$  matched
    - Add  $(u,v)$  to  $M$
- Output  $M$

### Why is $M$ maximal?

- If  $(u,v)$  not in  $M$  then either  $u$  or  $v$  already matched



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

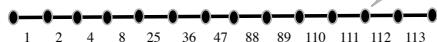
## Implementing the Oracle via Greedy

### To decide if edge $e$ in matching:

- Must know if adjacent edges that come before  $e$  in the ordering are in the matching
- Do not need to know anything about edges coming after

### Arbitrary edge order can have long dependency chains!

Odd or even steps from beginning?



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology

---



---



---



---



---



---



---



---



---



---

## Breaking Long Dependency Chains

### Assign random ordering to edges

- Greedy works under any ordering
- Important fact: random order has short dependency chains

[Nguyen Onak]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology

---



---



---



---



---



---



---



---



---



---

## Better Complexity for VC

- Additional ideas yield query complexity nearly linear in average degree for general graphs [Yoshida Yamamoto Ito] [Onak Ron Rosen Rubinfeld]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Further Work

- More complicated arguments for sparse maximum matching, set cover, positive Linear Programming... [Parnas Ron + Kuhn Moscibroda Wattenhofer] [Nguyen Onak]
- Even better results for special classes of graphs [Hassidim Kelner Nguyen Onak] [Newman Sohler]
  - e.g., planar, hyperfinite



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Tackling the Challenges of Big Data Big Data Analytics Fast Algorithms I: Sublinear Time Approximation Algorithms

**THANK YOU**




---



---



---



---



---



---



---



---



---

**Tackling the Challenges of Big Data****Big Data Analytics****Fast Algorithms I****Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

**Tackling the Challenges of Big Data****Big Data Analytics****Fast Algorithms I:****Local Computation Algorithms****Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

# Big inputs

# Big outputs



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

## When We Don't Need to See All the Output...

do we need to see all the input?



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Locally (list-)Decodable Codes

[Sudan-Trevisan-Vadhan, Katz-Trevisan,...]

Input

Encoding of big message

Output of Decoding Algorithm

Original big message

What is the ith bit?

Can design codes so that few queries needed to compute answer!!!



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Local Decompression Algorithms

[Muthukrishnan Strauss Zheng][Chandar Shah Wornell][Sadakane Grossi][Gonzalez Navarro][Ferragina Venturini][Kreft Navarro][Billie Landau Raman Sadakane Satti Weimann][Dutta Levi Ron Rubinfeld]

Input

Compression of big data

Output of Decompression Algorithm

Original big data

What is the ith bit?

design compression scheme so that compress well and few queries needed to compute answer



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



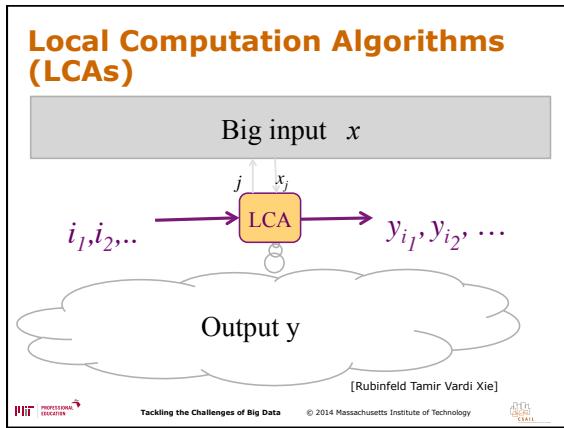
---



---



---




---

---

---

---

---

---

**What More Can be Done in This Model?**

- **Optimization problems** [Rubinfeld Tamir Vardi Xie]  
[Alon Rubinfeld Vardi Xie]
  - Maximal independent set, certain constraint satisfaction problems,...
- **Graph sparsification** [Campagna Guo Rubinfeld]  
[Levi Ron Rubinfeld]

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

**Even more:**

- **Local algorithms for**
  - Ranking webpages
  - Graph partitioning

[Andersen Borgs Chayes Hopcroft Mirrokni Teng]  
[Andersen Chung Lang] [Spielman Teng] [Borgs Brautbar Chayes Lucier]
- **Property preserving data reconstruction**  
[Ailon Chazelle Comandur Liu] [Comandur Saks] [Jha Raskhodnikova][Campagna Guo R.] [Awasthi Jha Molinaro Raskhodnikova] ...

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

---

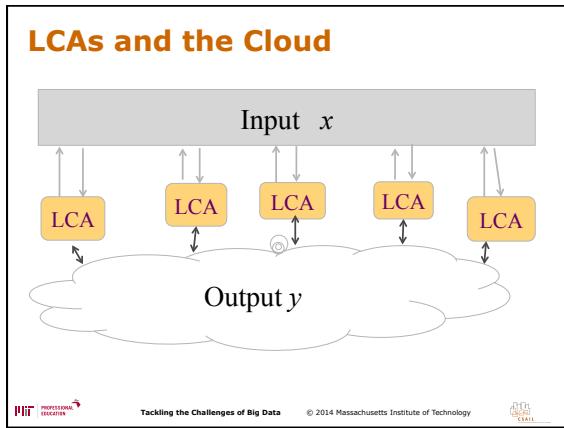
---

---

---

---

---




---



---



---



---



---



---




---



---



---



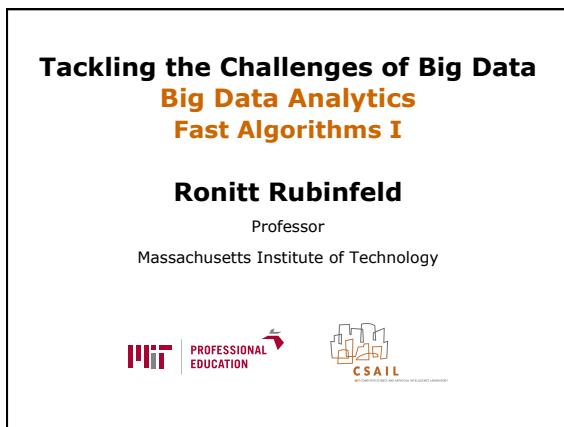
---



---



---




---



---



---



---



---



---

**Tackling the Challenges of Big Data**

**Big Data Analytics**

**Fast Algorithms I**

**Big Distributions**

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology





---



---



---



---



---



---



---

# No Samples

What if data only accessible via random samples?









---



---



---



---



---



---



---

## Distributions

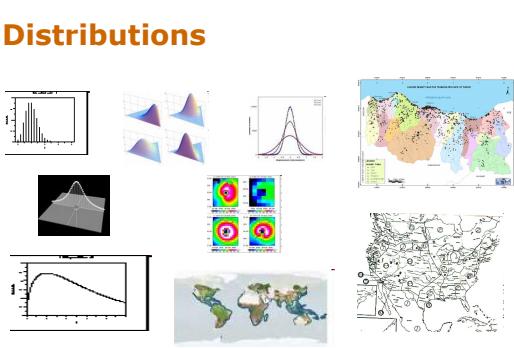


Photo & Image Source: Google









---



---



---



---



---



---



---

A collage of lottery tickets and scratch-off cards. The top left image shows a large pile of various lottery tickets from different states. The top right image shows a collection of scratch-off cards with visible numbers and patterns. The bottom image shows four specific lottery tickets standing upright: 'Georgia Power', 'Florida State Lottery', 'Michigan Lottery', and 'Oregon Lottery'.

---

---

---

---

---

---

---

---

---

---

---

# Is The Lottery Unfair?

- From [Hitlotto.com](#): Lottery experts agree, past number histories can be the key to predicting future winners.

---

---

---

---

---

---

---

# True Story!



- **Polish lottery Multilotek**
  - Choose "uniformly" at random distinct 20 numbers out of 1 to 80.
  - Initial machine biased
    - e.g., probability of 50-59 too small
- **Past results:** [http://serwis.lotto.pl:8080/archiwum/wyniki\\_wszystkie.php?id\\_gra=2](http://serwis.lotto.pl:8080/archiwum/wyniki_wszystkie.php?id_gra=2)

---

---

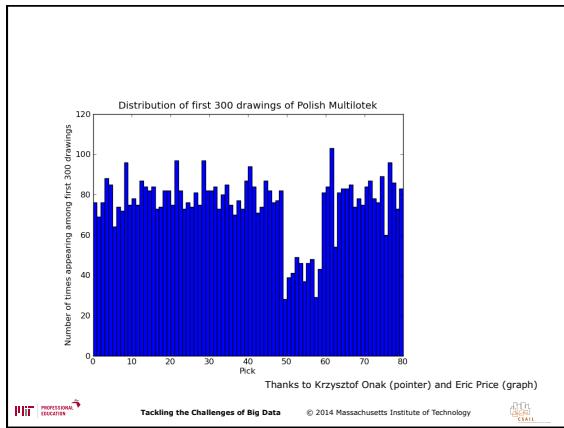
---

---

---

---

---



## New Jersey Pick 3, 4 Lottery

#### ▪ New Jersey Pick k ( =3,4) Lottery.

- Pick k digits in order.
  - $10^k$  possible values.
  - Assume lottery draws iid

**Data:**

  - Pick 3 - 8522 results from 5/22/75 to 10/15/00.
    - $\chi^2$ -test gives 42% confidence
  - Pick 4 - 6544 results from 9/1/77 to 10/15/00.
    - fewer results than possible values
    - $\chi^2$ -test gives no confidence



## Distributions on BIG domains

- Given samples of a distribution, need to know, e.g.,

- entropy
  - number of distinct elements
  - “shape” (monotone, bimodal,...)
  - closeness to uniform, Gaussian, Zipfian...

- No assumptions on shape of distribution

- i.e., smoothness, monotonicity, Normal distribution,...



## Distributions on BIG domains, ctd.

- Considered in statistics, information theory, machine learning, databases, algorithms, physics, biology,...



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Key Question

- How many samples do you need in terms of domain size?

- Do you need to estimate the probabilities of each domain item?
- Can sample complexity be sublinear in size of the domain?

Rules out standard statistical techniques,  
learning distribution



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Our Aim:

Algorithms with **sublinear** sample complexity



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I**  
Big Distributions

**THANK YOU**



---

---

---

---

---

---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I**

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I:**  
Similarity of Big Distributions?

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

## Our Aim:

Algorithms for understanding distributions with  
**sublinear** sample complexity



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Similarities of Distributions

- Are two distributions  $p$  and  $q$  close or far?

- $p$  is given via samples
- $q$  is either
  - known to the tester (e.g. uniform)
  - given via samples



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Is $p$ uniform?



**Theorem:** ([Goldreich Ron] [Batu Fortnow Rubinfeld Smith White] [Paninski]) Sample complexity of distinguishing

$$p = U \text{ from } ||p - U||_1 > \epsilon \text{ is } \Theta(n^{1/2})$$



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Upper Bound for L<sub>2</sub> Distance

- L<sub>2</sub> distance:

$$\begin{aligned} \|\mathbf{p}-\mathbf{U}\|_2^2 &= \sum_{i \in [1..n]} (p_i - 1/n)^2 \\ &= \sum p_i^2 - 2\sum p_i/n + \sum 1/n^2 \\ &= \sum p_i^2 - 1/n \end{aligned}$$

- Estimate collision probability to estimate L<sub>2</sub> distance from uniform

[Goldreich Ron]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Testing Uniformity

- **Upper bound: Estimate collision probability and use known relation between L<sub>1</sub> and L<sub>2</sub> norms**

- Issues:
  - Collision probability of uniform is 1/n
  - Use O(sqrt(n)) samples via recycling
  - Comment: [Paninski] uses different estimator

- **Easy lower bound:  $\Omega(n^{1/2})$  samples needed**

- Can improve to  $\Omega(n^{1/2}/\epsilon^2)$  [P]

[Goldreich Ron, Batu Fortnow Rubinfeld Smith White]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

## Back to The Lottery...

plenty of samples!



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---

**Is  $p$  Uniform?**

**Theorem:** Sample complexity of distinguishing  $p=U$  from  $|p-U|_1 > \varepsilon$  is  $\Theta(n^{1/2})$

[Goldreich Ron][Batu Fortnow Rubinfeld Smith White] [Paninski]

Nearly same complexity to test if  $p$  is any known distribution  
[Batu Fischer Fortnow Kumar Rubinfeld White]

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology MIT PROFESSIONAL EDUCATION CAAI

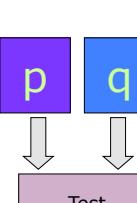
# Testing Closeness of Two Distributions

Transactions of  
20-30 yr olds

Transactions of  
30-40 yr olds

trend change?

# Testing Closeness



**Theorem:** Sample complexity of distinguishing  $p=q$  from  $\|p-q\|_1 > \epsilon$  is  $\Theta(n^{2/3})$

[BFRSW] [P. Valiant] [Chan Diakonikolas Valiant Valiant]

Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology 

## Why So Different?

- Collision statistics are all that matter
- Collisions on “heavy” elements can hide collision statistics of rest of the domain



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Additively Estimate Distance?

New goal: Output  $\|p-q\|_1 \pm \epsilon$

*much harder, but still sublinear!*

need  $\theta(n/\log n)$  samples [G. Valiant P. Valiant]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Collisions Tell All

### ▪ Algorithms:

- Algorithms use collisions to determine “wrong” behavior
- E.g., too many collisions implies far from uniform [GR,BFSRW]
- Use Linear Programming to determine if there is a distribution with the right collision probabilities and the right property [G. Valiant P. Valiant]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Collisions Tell All

### ▪ Lower bounds:

- For symmetric properties, collision statistics are only relevant information [BFRSW] (see also [Orlitsky Santhanam Zhang] [Orlitsky Santhanam Viswanthan Zhang])
  - Need new analysis tools since not independent
    - Central limit theorem for generalized multinomial distributions [G. Valiant P. Valiant]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



# Tackling the Challenges of Big Data

## Big Data Analytics

### Fast Algorithms I

### Big Distributions

THANK YOU



# Tackling the Challenges of Big Data

## Big Data Analytics

### Fast Algorithms I

Ronitt Rubinfeld

Professor

Massachusetts Institute of Technology



**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I**  
**Big Distributions – Final Words**

**Ronitt Rubinfeld**  
 Professor  
 Massachusetts Institute of Technology

---

---

---

---

---

---

**Information Theoretic Quantities**

Entropy  
 Support size

---

---

---

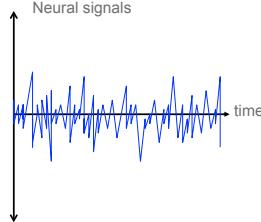
---

---

---

**Information in Neural Spike Trails**

[Strong, Koberle, de Ruyter van Steveninck, Bialek '98]



- Each application of stimuli gives sample of signal (spike train)
- Entropy of (discretized) signal indicates which neurons respond to stimuli

---

---

---

---

---

---

A large pile of trash on the left and a compact roll of trash on the right, separated by a red arrow.

---

---

---

---

---

---

---

---

---

---

# Can We Get Multiplicative Approximations?

- In general, no....
- What if entropy is at least some constant?
  - Can  $\gamma$ -multiplicatively approximate the entropy with  $\tilde{O}(n^{1/\gamma})$  samples (when entropy  $> 2g/\epsilon$ ) [Batu Dasgupta Rubinfeld Kumar]
  - requires  $\Omega(n^{1/\gamma})$  [Valiant]
  - better bounds when support size is small [Brautbar Samorodnitsky]
  - Similar bounds for estimating support size [Raskhodikova Ron Rubinfeld Smith] [Raskhodnikova Ron Shpilka Smith]

---

---

---

---

---

---

# Additive Approximations for Entropy and Support Size

- need  $\Theta(n/\log n)$  samples [Raskhodnikova, Ron, Shpilka, Smith] [Valiant Valiant]

---

---

---

---

---

---

---

## More Properties:

- Independence and Limited Independence [Batu Fischer Fortnow Kumar Rubinfeld White] [Alon Andoni Kaufman Matulef R. Xie] [Haviv Langberg]
- K-histogram distributions [Levi Indyk Rubinfeld]
- K-modal distributions [Daskalakis Diakonikolas Servedio]
- Poisson Binomial Distributions [Daskalakis Diakonikolas Servedio]
- Monotonicity over general posets [Batu Kumar Rubinfeld] [Bhattacharyya Fischer Rubinfeld P. Valiant]
- Properties of multiple distributions [Levi Ron Rubinfeld]



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Conclusion

- **For many problems, we need a lot less time and samples than one might think!**
- **Many cool ideas and techniques have been developed**
- **Lots more to do!**



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---



---



---



---



---



---



---



---



---



---

## Tackling the Challenges of Big Data Big Data Analytics Fast Algorithms I Big Distributions – Final Words

**THANK YOU**




---



---



---



---



---



---



---



---



---



---

**Tackling the Challenges of Big Data**  
**Big Data Analytics**  
**Fast Algorithms I**

**Ronitt Rubinfeld**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

---