

Private Release and Learning of Thresholds

Mark Bun (3rd year Ph.D., supported by NDSEG Fellowship)

Joint work with Kobbi Nissim, Uri Stemmer, and Salil Vadhan



Privacy Tools
for Sharing Research Data

A National Science Foundation
Secure and Trustworthy Cyberspace Project



MOTIVATING QUESTION

How many **data samples** do we need to achieve both **differential privacy** and **statistical accuracy**?

i.e. How big a study do we need to conduct to answer our questions and preserve privacy?

PRIVATE QUERY RELEASE

Counting queries: What fraction of rows in a database satisfy property q ?

e.g. $q(x) = \text{Age}(x) \geq 42$?

DarkSide?	Age	Home	Weight	
0	896	Dagobah	17	$q(x_1)=1$
0	19	Alderaan	49	$q(x_2)=0$
0	19	Tatooine	77	$q(x_3)=0$
1	42	Tatooine	136	$q(x_4)=1$

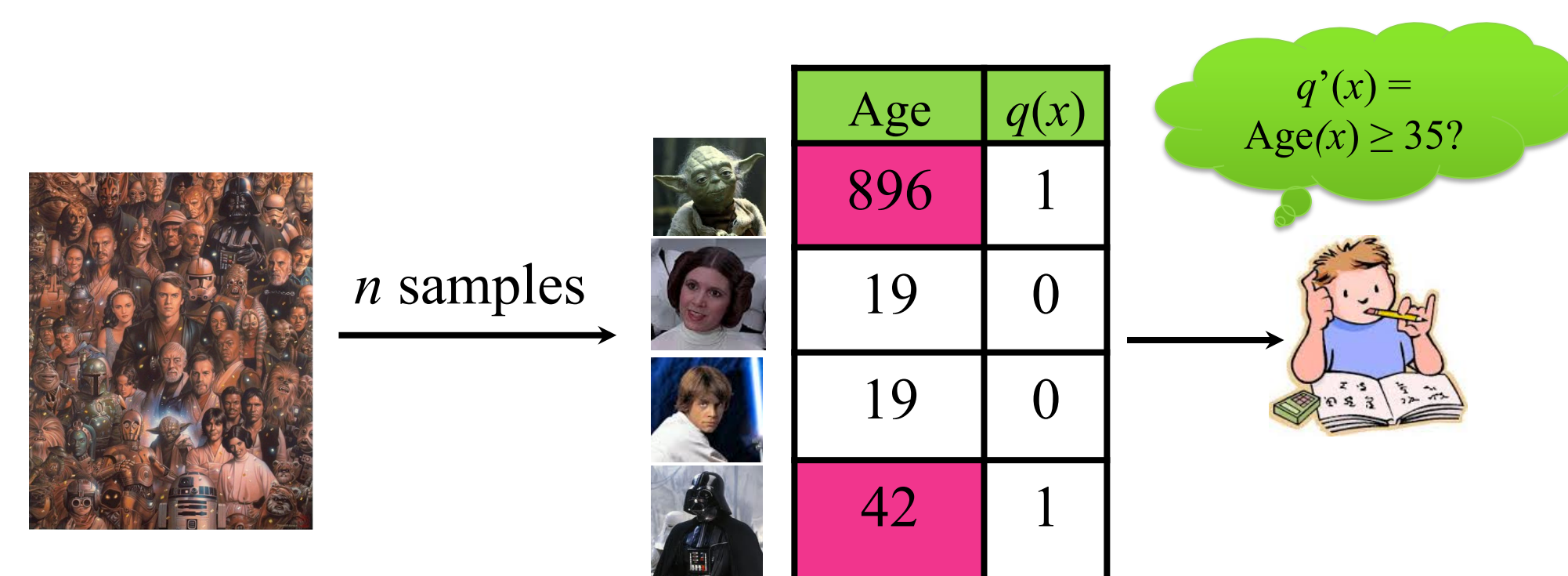
$q(D)=1/2$

Goal: Privately answer *all* $q \in Q$ to within 0.05 error

PRIVATE (PROPER) LEARNING

Examples drawn from a distribution and labeled by an unknown predicate $q \in Q$

e.g. $q(x) = \text{Age}(x) \geq 42$?



Goal: Output $q' \in Q$ that classifies new examples with 95% accuracy

THRESHOLD FUNCTIONS

$$\text{THRESH}^R(x) = \begin{cases} 1 & \text{if } x \geq y \\ 0 & \text{otherwise} \end{cases}$$



For any (even infinite) domain, $\text{VC}(\text{THRESH}) = 1$
 $\Rightarrow \text{THRESH}$ can be learned (non-privately) with $O(1)$ samples

QUESTION FOR THIS WORK

Can we *privately* release/learn thresholds over infinite domains?

If not, does the sample complexity depend on R ?

Answer: Thresholds require $n > \log^* R$

PRIOR WORK & OUR RESULTS

No privacy

Q arbitrary $Q = \text{POINT}_R$ $Q = \text{THRESH}_R$

$n = \Theta(\text{VC}(Q))$ [e.g. Vap98]	$n = \Theta(1)$ [Vap98]	$n = \Theta(1)$ [Vap98]
--	----------------------------	----------------------------

(ϵ, δ) -diff. priv.

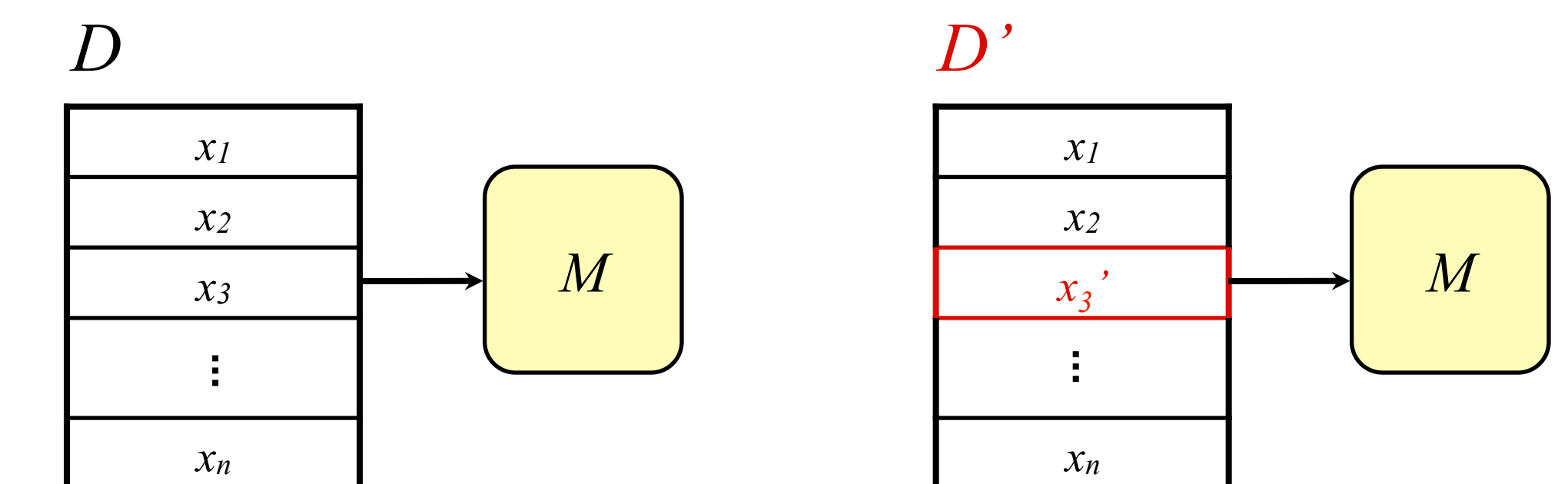
Query release

$\Theta(\log Q \cdot (\log R)^{1/2})$ [HR10, BUV14]	$\Theta(1)$ [BNS13]	$n < 8 \log^* R$ [BNS13] OUR WORK: $\log^* R < n < 2 \log^* R$
---	------------------------	--

Proper learning

$O(\log Q)$ [KLNRS08] ...but no general lower bounds	$\Theta(1)$ [BNS13]	$n < 8 \log^* R$ [BNS13] OUR WORK: $\log^* R < n < 2 \log^* R$
---	------------------------	--

DIFFERENTIAL PRIVACY



D and D' are neighbors if they differ only on one user's data

An algorithm M is (ϵ, δ) -differentially private if for all neighbors D, D' and every $S \subseteq \text{Range}(M)$,

$$\Pr[M(D) \in S] \leq e^\epsilon \Pr[M(D') \in S] + \delta$$

Think of $\epsilon = \Theta(1)$ and $\delta = o(1/n)$

CONCLUSIONS

- Releasing/learning thresholds requires sample complexity growing with R
- Separates private release/learning from non-private cases, even for $\text{VC}(Q) = 1$
- Open questions: Can we characterize the difference between private/non-private sample complexity? Extend results to improper learning?

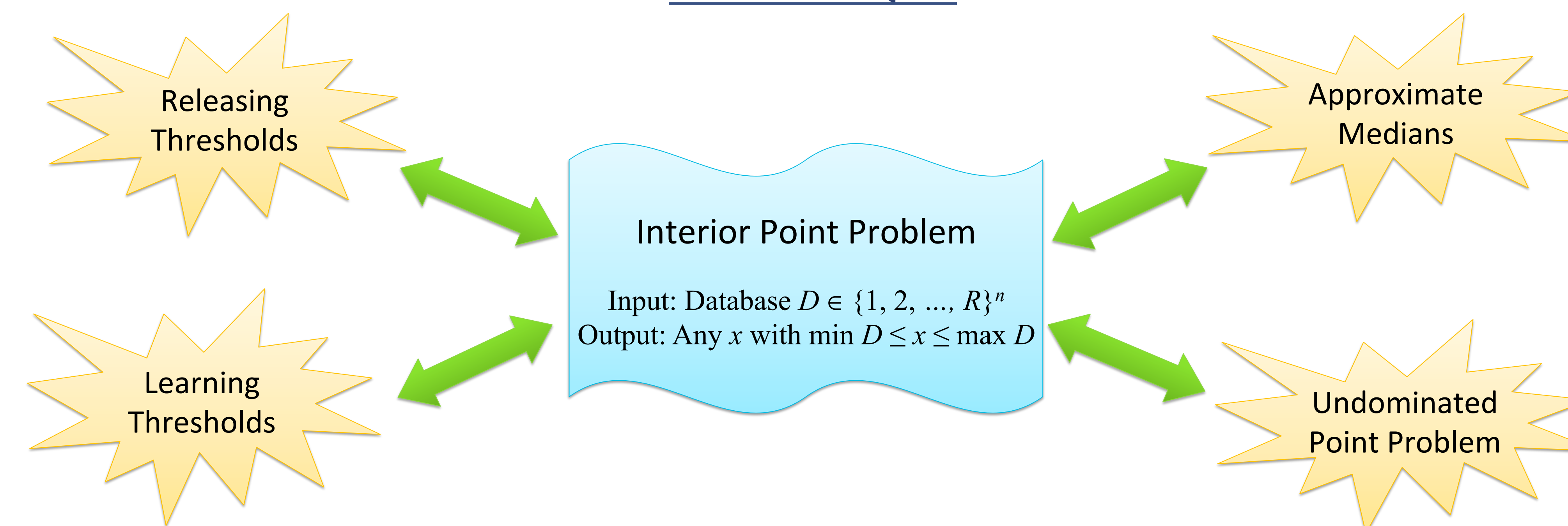
REFERENCES

- [BUV14] Mark Bun, Jonathan Ullman, and Salil Vadhan. Fingerprinting codes and the price of approximate differential privacy. In *STOC*, 2014.
- [BNS13] Amos Beimel, Kobbi Nissim, and Uri Stemmer. Private learning and sanitization: pure vs. approximate differential privacy. In *RANDOM*, 2013.
- [BNSV14] Mark Bun, Kobbi Nissim, Uri Stemmer, and Salil Vadhan. Differentially private release and learning of thresholds. *Manuscript*, 2014.
- [HR10] Moritz Hardt and Guy Rothblum. A multiplicative weights mechanism for privacy-preserving data analysis. In *FOCS*, 2010.
- [KLNRS08] Shiva Kasiviswanathan, Homin Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. What can we learn privately? In *FOCS* 2008.
- [Vap98] Vladimir Vapnik. Statistical Learning Theory. 1998.

CONTACT

mbun@seas.harvard.edu
Harvard School of Engineering and Applied Sciences
Maxwell Dworkin 138
33 Oxford St.
Cambridge, MA 02138

OUR TECHNIQUES



INTERIOR POINT LOWER BOUND

Hard distribution for domain size $R(n)$

$x_1 = 2$
$x_2 = 5$
$x_3 = 8$
\vdots
$x_n = 4$

Hard distribution for domain size $R(n+1) = 2^{R(n)}$

$y_0 = 01011101100110$
$y_1 = 01101001001001$
$y_2 = 01011010100101$
$y_3 = 010111101001001$
\vdots
$y_n = 010101101001100$

Random point

Agrees w/ y_0 in x_1 indices

INTERIOR POINT UPPER BOUND

Database with domain size $R(n)$

$y_1 = 01011101100110$
$y_2 = 01101001001001$
$y_3 = 01011010100101$
$y_4 = 010111101001001$
\vdots
$y_n = 010101101001100$

Random pairings

Database with domain size $R(n/2) = \log R(n)$

$x_1 = 6$
$x_2 = 2$
\vdots
$x_{n/2} = 8$

agreements