

Safer Data Mining: Algorithmic Techniques in Differential Privacy

Credit: Mary McGlohon

#### Speakers: Moritz Hardt and Sasho Nikolov IBM Almaden and Rutgers SDM 2014

#### Disclaimer

#### This talk is *not* a survey.



#### Disclaimer

This talk is not a history lesson.

Millions of years ago, giant marshmallows ruled the earth.

Having no skeletal system, they left no fossil-record.

#### Focus

# A few widely applicable algorithmic techniques

Algorithms you can run on an actual computer

#### Outline

#### **Part I: Differential Privacy:** Attacks, basic definitions & algorithms

**Part II: Unsupervised Learning** Technique: Singular Value Decomposition

> **Part III: Supervised Learning** Technique: Convex optimization

> > **Part IV: Streaming**

Technique: Streaming data structures

#### Releasing Data: What could go wrong?

#### Wealth of examples

One of my favorites: Genome-Wide Association Studies (GWAS)

> Trust me: You will like this even if you don't like biology



#### GWAS

Typical Setup:

- NIH takes DNA of 1000 test candidates with common disease
- 2. NIH releases minor allele frequencies (MAF) of test population at 100,000 positions (SNPs)
- **Goal:** Find association between SNPs and disease



#### Attack on GWAS data [Homer et al.]

## Can **infer membership in test group** of an individual with known DNA from published data!

MAF 0.02 0.03 0.05 ···· 0.02	SNP	1	2	3	 100000	Test
	MAF	0.02	0.03	0.05	 0.02	ρορυιατιοι

SNP	1	2	3	 100000
MA	NO	NO	YES	 YES

SNP	1	2	3	 100000
MAF	0.01	0.04	0.04	 0.01

Reference population (HapMap data, public)

Moritz's

DNA

#### Attack on GWAS data [Homer et al.]

## Can **infer membership in test group** of an individual with known DNA from published data!



#### Interesting characteristics

- Only innocuous looking data was released
   Data was HIPAA compliant
- Data curator is trusted (NIH)
- Attack uses background knowledge (HapMap data set) available in public domain
- Attack uses unanticipated algorithm
- Curator pulled data sets (now hard to get)

#### How not to solve the problem

- Ad-hoc anonymization (e.g., replace name by random string)
  - Gone wrong: Netflix user data, AOL search logs, MA medical records, linkage attacks
- Allow only innocuous looking statistics (e.g. low-sensitivity counts)

– Gone wrong: GWAS

• Use the wrong crypto tool (e.g., encrypt hard disk, securely evaluate query)

#### Meaningful privacy guarantee

- handles attacks with background information
- powerful composition properties

Intuition: Presence or absence of any individual in the data *cannot* be inferred from output

Two data sets *D,D*' are called *neighboring* if they differ in at most one data record. Example: *D* = {GWAS test population}, *D*' = *D* – {Moritz's DNA}

Informal Definition (Differential Privacy): A randomized algorithm *M*(*D*) is differentially private if for all neighboring data sets *D*,*D*' and all events *S*:

 $\Pr{M(D) \in S} \approx \Pr{M(D') \in S}$ 

Two data sets *D*,*D*' are called *neighboring* if they differ in at most one data record. Example: *D* = {GWAS test population}, *D*' = *D* – {Moritz's DNA}

Definition (Differential Privacy): A randomized algorithm M(D) is  $\varepsilon$ -differentially private if for all neighboring data sets D, D' and all events S:

 $\Pr{M(D) \in S} \le e^{\epsilon} \Pr{M(D') \in S}$ 

Think:  $\varepsilon = 0.01$  and  $e^{\varepsilon} \approx 1 + \varepsilon$ 

Two data sets *D,D*' are called *neighboring* if they differ in at most one data record. Example: *D* = {GWAS test population}, *D*' = *D* – {Moritz's DNA}

Definition (Differential Privacy): A randomized algorithm M(D) is  $(\varepsilon, \delta)$ -differentially private if for all neighboring data sets D, D' and all events S:

 $\Pr\{M(D) \in S\} \le e^{\epsilon} \Pr\{M(D') \in S\} + \delta$ 

Think:  $\varepsilon = 0.01$  and  $e^{\varepsilon} \approx 1 + \varepsilon$  $\delta << 1/|D|$  Definition (Differential Privacy):

A randomized algorithm M(D) is  $\varepsilon$ -differentially private if for all neighboring data sets D,D' and all events S:

### $\Pr{M(D) \in S} \le e^{\epsilon} \Pr{M(D') \in S}$



**Outputs** 

Definition (Differential Privacy): A randomized algorithm M(D) is  $(\varepsilon, \delta)$ -differentially private if for all neighboring data sets D, D' and all events S:

## $\Pr\{M(D) \in S\} \le e^{\epsilon} \Pr\{M(D') \in S\} + \delta$



#### Sensitivity

Assume f maps databases to real vectors **Definition:** 

- L1-sensitivity Δ<sub>1</sub>(f) = max || f(D) f(D') ||<sub>1</sub>
   maximum over all neighboring D,D'
- L2-sensitivity  $\Delta_2(f) = \max || f(D) f(D') ||_2$

Intuition: Low sensitivity implies "compatible with differential privacy"

#### Exercise

#### f maps D to $(q_1(D), q_2(D), ..., q_k(D))$

where each  $q_i(D)$  is a count: "How many people satisfy predicate  $P_i$ ?"

#### L1-sensitivity? k

L2-sensitivity?  $k^{1/2}$ 





Scale noise to L1-sensitivity

Fact 1: Satisfies (*\varepsilon*, *O*)-differential privacy

Fact 2: Expected error  $\Delta_1(f)/\epsilon$  in each coord.

#### Laplacian Mechanism [DMNS'06]





Scale noise to L2-sensitivity

Fact: Satisfies ( $\varepsilon$ ,  $\delta$ )-differential privacy

Fact: Expected error  $\Delta_2(f)\log(1/\delta)/\varepsilon$ 

#### Composition



Differential privacy composes: 1. in parallel 2. sequentially 3. adaptively

## *T*-fold **composition** of $\varepsilon_0$ -differential privacy satisfies:

Answer 1 [DMNS'06]:

 $\varepsilon_0 T$ -differential privacy

Answer 2 [DRV'10]:

 $(\varepsilon, \delta)$ -differential privacy  $\epsilon \leq 4\epsilon_0 \sqrt{T \log(1/\delta)}$ 

Note: for small enough  $\varepsilon$ 

## Part II: Private SVD and Applications

#### We meet the protagonist



A is a real n x d matrix thought of as n data points in d dimensions

#### ...and its singular value decomposition



U,V have orthonormal columns  $\Sigma$  is diagonal  $r \ge r$ , r = rank(A)



#### The left factor

Columns  $u_1, ..., u_r$  are the singular vectors of A  $\Leftrightarrow$  eigenvectors of the  $n \ge n$ Gram matrix  $AA^T$ 

By Courant-Fisher:

$$\|u_1^T A\| = \max \{ \|x^T A\| : \|x\| = 1 \}$$
  
$$\|u_2^T A\| = \max \{ \|x^T A\| : \|x\| = 1, \langle x, u_1 \rangle = 0 \}$$



Columns  $v_1, ..., v_r$  of V are the eigenvectors of the  $d \ge d$ covariance matrix  $A^T A$ 



Diagonal entries  $\sigma_1, ..., \sigma_r$ square roots of the nonzero eigenvalues of  $A^T A$ 

Notation:  $\sigma_i(A) = \sigma_i$ 

Note: Singular values are unique, singular vectors are not (e.g. identity matrix)

#### Why compute the SVD?



# Principal Component Analysis right singular vectors = principal components

#### **Low Rank Approximation**

truncated SVD gives best low rank approximation in Frobenius and spectral norm

**Many more:** Spectral clustering, collaborative filtering, topic modeling

#### **Differential Privacy on Matrices**

How should we define "neighboring"?

Differ in one row

Differ in one row by norm at most 1

Differ in at most one entry by at most 1

#### **Differential Privacy on Matrices**

How should we define "neighboring"?

Differ in at most one entry by at most 1

- Focus in this talk, others are reasonable
- Reasonable for sparse data with bounded coeffs
- Many entries? Choose smaller ε.
- Algorithms we'll see also give other guarantees using different noise scaling.

#### Objectives

- Approximating the top singular vector
  - Objective:  $\max_{\substack{x: \|x\|=1}} \|Ax\| \qquad \text{OPT} = \sigma_1(A)$
- Approximating k singular vectors
  - Objective:

$$\min_{B: rk(B)=k} ||A - B||_2 \qquad \text{OPT} = \sigma_{k+1}(A)$$

Goal: Minimize additive error subject to differential privacy

## Baseline Algorithm

1965

#### RANDOMIZED RESPONSE: A SURVEY TECHNIQUE FOR ELIMINATING EVASIVE ANSWER BIAS

STANLEY L. WARNER Claremont Graduate School

For various reasons individuals in a sample survey may prefer not to confide to the interviewer the correct answers to certain questions. In such cases the individuals may elect not to reply at all or to reply with incorrect answers. The resulting evasive answer bias is ordinarily difficult to assess. In this paper it is argued that such bias is potentially removable through allowing the interviewee to maintain privacy through the device of randomizing his response. A randomized response method for estimating a population proportion is presented as an example. Unbiased maximum likelihood estimates are obtained and their mean square errors are compared with the mean square errors of conventional estimates under various assumptions about the underlying population.

#### 1. INTRODUCTION

**F**OR reasons of modesty, fear of being thought bigoted, or merely a reluctance to confide secrets to strangers, many individuals attempt to evade certain questions put to them by interviewers. In survey vernacular, these people become the "non-cooperative" group [5, pp. 235–72], either refusing outright to be surveyed, or consenting to be surveyed but purposely providing wrong answers to the questions. In the one case there is the problem of refusal bias [1, pp. 355–61], [2, pp. 33–6], [5, pp. 261–9]; in the other case there is the Input Perturbation [Warner65, BlumDworkMcSherryNissim05] aka Randomized Response (RR)

Randomized Response(Input A):1.  $N = Lap(1/\varepsilon)^{n \times d}$ 2. B = A + N3. Output SVD(B)

Fact: Satisfies *ɛ*-differential privacy.

How much error does this introduce?

#### Some matrix theory

**Theorem** (Weyl): Let *A*,*N* be *n* x *d* matrices. Then,

#### $|\sigma_i(A+N) - \sigma_i(A)| \le \|N\|_2$

**Theorem**: Let  $N = Lap(1/\varepsilon)^{n \times d}$ . Then,

$$\mathbb{E}\|N\|_2 \le O\left(\frac{\sqrt{n} + \sqrt{d}}{\epsilon}\right)$$
## When does RR Work?

n >> d, d relatively small

Run on d x d covariance matrix

#### **Example:**

Differentially private recommender system evaluated on Netflix data set [McSherryMironov09] Here, n = 480189 and d = 17770

## When RR doesn't work

#### **Problem 1:**

n and d both large

#### Problem 2:

matrix sparse, e.g, ⊿ ones per row

RR generates huge dense noise matrix

 $\sigma_1(A) \leq \sqrt{n\Delta}$ 

Can we improve on RR? No! But we'll do it anyway.

## Why not? Dictator example



n

n

00000000000...00000000000

v<sub>1</sub> = top row (up to scaling)

Fact: Differential privacy requires error  $\sqrt{n}$ 

Error o(n<sup>1/2</sup>) definition of "blatant non-privacy"!

How are we going to get around this?

## Beat **Randomized Response** (1965) with **Power Method** (1929)

#### ZUSAMMENFASSENDE BERICHTE

#### Praktische Verfahren der Gleichungsauflösung.

Von R. v. MISES und H. POLLACZEK-GEIRINGER in Berlin.

In einer Vorlesung über »Praktische Analysis«, die der erstgenannte Verfasser im Sommer-Semester 1927 hielt, kamen eine Reihe von Rechnungsverfahren zur Sprache, die teils neu, teils wenig bekannt sind. Es erschien zweckmäßig, einiges davon mit kurzer Begründung und Hinweisen auf die wichtigsten Anwendungen zusammenzustellen. Dieser Arbeit hat sich der zweitgenannte Verfasser — zunächst hinsichtlich der Methoden der Gleichungsauflösung — unterzogen; von ihm rühren auch wesentliche Ergänzungen und nähere Ausführungen sowie die Zahlenbeispiele her.

Im ersten Abschnitt des vorliegenden Berichtes wird ein Verfahren zur Lösung einer beliebigen Gleichung mit einer Unbekannten entwickelt, das sich von den bisher bekannten dadurch unterscheidet, daß es stets zum Ziele führt (immer konvergiert) und in der Regel weniger Rechenarbeit verlangt. Die Aufsuchung der Wurzeln von Gleichungssystemen wird in den beiden folgenden Abschnitten nur für den Sonderfall der linearen Gleichungen behandelt. Für das »Iterationsverfahren in Gesamtschritten«, welches die unmittelbare Uebertragung des Verfahrens für Eine Unbekannte darstellt, werden in Abschnitt 2 die zum Teil bekannten, meist nicht gehörig beachteten Konvergenzbedingungen zusammengestellt sowie Fehlerabschätzungen gegeben.

Besondere Bedeutung kommt dem in Abschnitt 3 behandelten, als »Iteration in Einzelschritten« bezeichneten, von Ph. Seidel stammenden Ansatz zu. Dieser ist — nach

## From here on

#### Assume A is symmetric and n x n

Wlog, consider:



## Recap: Power Method

**Input:** *n* x *n* matrix A, parameter T

Pick random unit vector x<sub>0</sub>

**For** *t* = 1 *to T*:

$$y_t = Ax_{t-1}$$
$$x_t = y_t / |y_t|_2$$

**Output**  $X_T$ 

# **Fact:** *x*<sub>*T*</sub> **converges to top SV** for now all we care about!

## Senstivity of Matrix-Vector Product

- Suppose A, A' differ in one entry by 1
- Assume x is a unit vector

Fact:  $||(A-A')x||_2 \le ||x||_{\infty}$ where  $||x||_{\infty} = \max_i |x_i|$ 

## **Noisy** Power Method

**Input:** *n* x *n* matrix A, parameters T,  $\varepsilon, \delta > 0$ Pick random unit vector x<sub>0</sub> For t = 1 to T:  $g_t = N(0, T \log(1/\delta)/\varepsilon^2 |x_t|^2$  $y_{t} = Ax_{t-1} + g_{t}$  $1 \infty$ )<sup>n</sup> largest squared  $x_{t} = y_{t} / |y_{t}|_{2}$ entry of  $X_{t-1}$ **Output**  $X_T$ in [1/n,1]

**Fact:** Satisfies  $(\varepsilon, \delta)$ -differential privacy

## Bounding the largest entry of x<sub>t</sub>

Let  $v_1, \dots, v_n$  be singular vectors of A Put  $B = \max_{i=1}^{n} \|v_i\|_{\infty}$ **Easy bound :**  $||x_t||_{\infty} \leq \sqrt{nB}$ Can we do better? Lemma [H-Roth13]:

$$\Pr\left\{\|x_t\|_{\infty} \geq 4B\sqrt{\log n}\right\} \leq 1/n^3$$

## A pleasant surprise

$$\mu(A) = n \cdot \max_{i=1}^{r} \|v_i\|_{\infty}^2$$

known as *coherence* of A and widely studied

Theoretically and empirically often small

polylog(n) in random models [CandesTao09] much less than *n* for real world data

Previous lemma:

$$: ||x_t||_{\infty} \lesssim \sqrt{\frac{\mu(A)}{n}}$$
  
and hence  $||g_t|| \lesssim \sqrt{\mu(A)}$ 

## Performance of Power Method

#### **Theorem** [*H*-Roth'13]:

The Noisy Power Method satisfies  $(\varepsilon, \delta)$ -differential privacy and with  $T = O(\log(n))$  steps returns a unit vector x such that

$$\|Ax\| \ge \sigma_1(A) - O\left(\epsilon^{-1}\sqrt{\mu(A)\log(1/\delta)}\log n\right)$$

provided A satisfies "singular value smaration."

**Contrast with** 

for Randomized Response even if  $\mu(A)=1$ 

**Theorem:** Nearly matching lower bound for every setting of  $\mu(A)$ .

#### When utility and privacy benefit from the same principle



## Generalization: Subspace Iteration

Input: n x n matrix A symmetric, target rank k
X<sub>0</sub> random orthonormal matrix
For t = 1 to T:
– Pick Gaussian perturbation G<sub>t</sub>

$$-Y_t = AX_{t-1} + G_t$$

 $-X_t = Orthonormalize(Y_t)$ 

**Output**  $X_{T}$  (approx top k singular vectors)

# **Principle Angle Between Subspaces** Let U, X subspaces of dimension k k=1 cos $\Theta(U,X) = |U^TX|$ **In general** $\cos \Theta(U,X) = \sigma_{\min}(U^T X)$ $\sin \Theta(U,X) = \sigma_{\max}(V^T X)$ where V orthog. complement of U $\tan \Theta(U,X) = \sin \Theta(U,X) / \cos \Theta(U,X)$

## Main Convergence Lemma

Let U be spanned by top k singular vectors of A.



## **Application 1: Spectral Clustering**

#### Planted multisection model:

c clusters, intra-cluster edge probability pinter-cluter probability q, q < p



#### How to recover a cluster?

Simple approach: Cheeger Cut

- 1. Compute second eigenvector  $v_2$  of the graph G
- 2. Sort coordinates in ascending order
- 3. Pick vertices corresponding to first n/c coordinates

## **Application 1: Spectral Clustering**

Graph with n = 20,000 vertices, p = 0.2



#### **Open Problem:**

Explore more sophisticated spectral clustering techniques.

## Application 2: Matrix approximation

Upper left 40 x 100 corner of a 1000 x 1000 matrix.

## **Differentially Private Approximation**

One step Subspace (Non-)Iteration, noisy projection, rounding



In this example: Matrix has rank 2 and coherence 2 making nearly exact recovery feasible even under differential privacy.

## Enter Graph Cuts

 Given graph G=(V,E), cut query is a subset S of V, answer is E<sub>G</sub>(S,S<sup>c</sup>)



**Goal:** Come up with weighted graph G' that satisfies differential privacy and approximates all cuts in G.

## Synthetic Data for Cuts: What's known

Goal: Preserve all cuts of size s

**Randomized Response** 

Johnson-Lindenstrauss [BlockiBlumDattaSheffet12]

Error  $O(sn^{1/2})$ 

Error  $O(s^{1.5})$ 

Better for s << n

MW+EM (*inefficient*) gives error:



## JL Approximation of the Laplacian

Graph Laplacian  $L_G$  Cut query:  $\mathbb{1}_{S}^{T}L_G\mathbb{1}_{S}$ 

Fact: 
$$L_G = E_G^T E_G$$
  $E_G$  is the weighted  $\binom{n}{2} \times n$   
edge-vertex incidence matrix

Suppose we pick random Gaussian  $r \times {n \choose 2}$ matrix M and put:  $\widetilde{L}_G = \frac{1}{r} E_G^T M^T M E_G$ 

By JL Theorem,  $r = \alpha^{-2} \log(1/\beta)$ preserves single cut up to factor (1± $\alpha$ ) with pr 1- $\beta$ 

## Adding in Privacy [BBDM12]

$$H = \frac{w}{n}K_n + \left(1 - \frac{w}{n}\right)G$$

**Theorem** [BBDM]: For  $w = O\left(\epsilon^{-1}\sqrt{r}\log(1/\delta)\log(r/\delta)\right)$ the JL approximation of  $L_H$  satisfies  $(\varepsilon, \delta)$ -differential privacy.

Additive error on one cut of size s: O(ws)

Geometric Intuition Sample  $g \sim N(0, 1)^{\binom{n}{2}}$ 

**Fact:**  $E_H^T g$  Gaussian variable with covariance matrix  $L_H$ 

Sparse cuts  $\Leftrightarrow$  directions of small variance

Mixing in (w/n)K<sub>n</sub> gives variance w<sup>2</sup> in every unit direction (thus hiding single edge change in *G*)





# Part II

## But wait, there's more...

## Part III: Supervised Learning (Empirical Risk Minimization)

## Supervised Learning: Classification



## Supervised Learning: Regression



## **Generalized Linear Model**

- <u>Data</u>  $X_n = \{(x_{1,}y_1), (x_{2,}y_2), ..., (x_{n,}y_n)\}$ -  $(x_{i,}y_i)$  sampled IID from distribution D - data point  $x_i$  is d-dimensional, label  $y_i$  is real
- <u>Goal</u>: predict y for new x, w
- <u>Hypothesis class:</u> H
- Loss function: l(<w, x>; y)
  - loss of hypothesis w on (x,y)
  - assumption: y can be predict measurement of x
  - I will be *convex* in *w*



## **Risk Minimization**

• The Risk Minimization Problem:

 $-w^* = \arg \min_{w} E_{(x, y) \sim D} [I(\langle w, x \rangle; y)]$ 

- How to do that based on  $X_n$ ?
- Minimize empirical risk:

$$R_n(w) = \frac{1}{n} \sum_i \ell(\langle w, x_i \rangle; y_i)$$

 Uniform convergence: if everything is "nice", empirical risk minimizer w → w\* as n → ∞

## Regularization

- Often there are many solutions to the ERM problem
  - many hypotheses fit the data
- Occam's Razor: pick the "simplest" hypothesis
- Regularized ERM:

$$J_n(w) = R_n(w) + \lambda r(w)$$

• Regularizer *r*(*w*): "complexity" of *w* 

## Strong Convexity

- Need regularizer to be strongly convex:
  - unique optimal w
  - <u>robustness</u> to data perturbation
  - helps <u>privacy</u>: output does not depend on any data point too much

$$r(z) \ge r(w) + \left\langle \nabla r(w), z - w \right\rangle + \left\| z - w \right\|_{2}^{2}$$

## **Strong Convexity**



## Examples

- SVM:
  - $-I(\langle w, x \rangle; y) = \max\{0, 1-y \langle w, x \rangle\};$  $-r(w) = 0.5 |w|^2$
- Logistic Regression

   l(<w, x>; y) = log(1 exp(-y<w,x>);
   r(w) = 0.5 |w|<sup>2</sup>
- Ridge Regression

$$-|(\langle w, x \rangle; y) = (y - \langle w, x \rangle)^{2};$$
  
- r(w) = 0.5 | w |<sup>2</sup>

### **Private Algorithm: Output Perturbation**

- 1. Compute minimizer w of  $J_n(w)$
- 2. Sample noise  $b \sim N(0, c(\varepsilon, \delta)/\lambda^2 n^2)^d$
- 3. Output  $w^{\sim} = w + b$
- (ε,δ)-DP if:
  - I is 1-Lipschitz:  $|I(\langle z, x \rangle; y) - I(\langle w, x \rangle; y)| \le |\langle z - w, x \rangle|$  $- |x| \le 1$
- [Chaudhuri, Monteleoni, Sarwate '11]
#### Sensitivity Analysis

• *Intuition*: strong convexity -> low sensitivity

$$X = \{(x_{1}, y_{1}), (x_{2}, y_{2}), ..., (x_{n}, y_{n})\}$$

$$X' = \{(x_{1}, y_{1}), (x_{2}, y_{2}), ..., (x', y')\}$$

$$J'_{n}(w): \text{risk for } X; \qquad J'_{n}(w): \text{risk for } X'$$

$$w: \text{minimizer of } J_{n}(w); \qquad w': \text{minimizer of } J'_{n}(w);$$

Sensitivity: 
$$|w - w'| \le 2/(\lambda n)$$

Privacy follows as usual.

#### Sensitivity Analysis Sketch

• Optimality of w, w' and strong convexity:

$$J'_{n}(w) \ge J'_{n}(w') + \frac{\lambda}{2} \|w - w'\|^{2}$$
$$J_{n}(w') \ge J_{n}(w) + \frac{\lambda}{2} \|w - w'\|^{2}$$

• Only one data point changed:

$$J_n(w) - J'_n(w) = \frac{1}{n} (\ell(\langle w, x \rangle; y) - \ell(\langle w, x' \rangle; y'))$$

• Combine + Lipschitz<sup>"</sup>-ness:

$$\lambda \|w - w'\|^2 \leq \frac{1}{n} |\langle w - w', x \rangle| + \frac{1}{n} |\langle w - w', x' \rangle| \leq \frac{2}{n} \|w - w'\|$$

#### **Generalization Error**

• Expected risk relates to the optimal risk as:  $E[\ell(\langle w^{\sim}, x \rangle; y) - \ell(\langle w^{*}, x \rangle; y)] \leq \frac{c'(\varepsilon, \delta)}{\sqrt{n}}$ 

when  $\lambda = n^{-1/2}$ .

- Proof idea [Jain, Thakurta '13]
  - Lipschitzness:

$$\left|\ell(\langle w^{\sim}, x \rangle; y) - \ell(\langle w, x \rangle; y)\right| \leq \left|\langle b, x \rangle\right|$$

 $- \langle b, x \rangle$  is Gaussian with variance  $c(\varepsilon, \delta)/(\lambda^2 n^2)$ 

- bounds  $J_n(w^{\sim}) - J_n(w)$ . Suffices by convergence results.

#### Variants and Extensions

- Objective Perturbation [Chaudhuri, Monteleoni, Sarwate '11], [Kifer, Smith, Thakurta '12], [Jain, Thakurta '13]
  - minimize  $J_n(w) + \langle b, w \rangle$  for random Gaussian b
  - improved guarantees for "nice" data (adapts to convexity of the instance)
- Analysis can be extended to:
  - ERM with more general loss function
  - Structural constraints (sparse regression)
  - But usually a dependence on *d* creeps in

### **Other Approaches**

- Exploiting robustness
  - Private algorithms from learning algorithms robust to perturbations of the input
  - [Smith, Thakurta '13] Private Lasso with optimal sampling complexity
- Online learning: data arrives online, minimize regret
  - [Jain, Kothari, Thakurta '12] Sensitivity analysis
  - [Smith, Thakurta '13] Follow the approximate leader

# Part IV: Streaming Models

#### The Streaming Model

1, 4, 5, 19, 145, 14 , 5, 5, 16, 4

+,-,+, -, +, +,-,+, -, +

- Underlying *frequency vector* A = A [1], ..., A[n]
  - start with A[i] = 0 for all *i*.

#### • We observe an <u>online sequence of updates:</u>

- Increments only (cash register):
  - Update is  $i_t \rightarrow A[i_t] := A[i_t] + 1$
- Fully dynamic (turnstile):
  - Update is  $(i_t, \pm 1) \rightarrow A[i_t] := A[i_t] \pm 1$

#### • <u>Requirements</u>: compute statistics on A

- Online, O(1) passes over the updates
- Sublinear space, polylog(n,m)

### **Typical Problems**

- Frequency moments: F<sub>k</sub> = |A[1]|<sup>k</sup> + ... + |A[n]|<sup>k</sup>
   related: L<sub>p</sub> norms
- Distinct elements:  $F_0 = \#\{i: A[i] \neq 0\}$
- k-Heavy Hitters: output all *i* such that  $A[i] \ge F_1/k$
- Median: smallest *i* such that  $A[1] + ... + A[i] \ge F_1/2$ - Generalize to Quantiles
- Different models:
  - Graph problems: a stream of edges, increments or dynamic
    - matchings, connectivity, triangle count
  - Geometric problems: a stream of points
    - various clustering problems

#### When do we need this?

- The universe size *n* is *huge*.
- Fast arriving stream of updates:
  - IP traffic monitoring
  - Web searches, tweets
- Large unstructured data, external storage:
  - multiple passes make sense
- Streaming algorithms can provide a *first rough approximation* 
  - decide whether and when to analyze more
  - fine tune a more expensive solution
- Or they can be the *only feasible solution*

#### A taste: the AMS sketch for F<sub>2</sub> [Alon Matias Szegedy 96]



 $h:[n] \rightarrow \{\pm 1\}$  is 4-wise independent:

 $E[X^2] = F_2$   $E[X^4]^{1/2} \le O(F_2)$ 

#### The Median of Averages Trick



Average: reduces variance by  $\alpha^2$ .

Median: reduces probability of large error to  $\delta$ .

### **Defining Privacy for Streams**

- We will use *differential privacy*.
- The database is represented by a stream
  - online stream of transactions
  - offline large unstructured database
- Need to define *neighboring inputs:* 
  - Entry (event) level privacy: differ in a single update

1, 4, 5, 19, 145, 14 , 5, 5, 16, 4

1, 1, 5, 19, 145, 14 , 5, 5, 16, 4

- <u>User level privacy</u>: replace some updates to *i* with updates to *j* 

1, 4, 5, 19, 145, 14 , 5, 5, 16, 4

1, 4, 3, 19, 145, 14, 3, 5, 16, 4

 We also allow the modified updates to be placed somewhere else

## Streaming & DP?

- Large unstructured database of transactions
- Estimate how many distinct users initiated transactions?

- i.e.  $F_0$  estimation

- Can we satisfy <u>both</u> the streaming and privacy constraints?
  - $-F_0$  has sensitivity 1 (under user privacy)
  - Computing  $F_0$  exactly takes  $\Omega(n)$  space
  - Classic sketches from streaming may have large sensitivity

### Flajolet Martin Sketch for F<sub>0</sub>

- Store a *bit map B* of *L* = O(log *n*) bits.
   One computer word
- Randomly hash update to L bits
- Bitmap: information about least significant 1 in hashed values



Estimate: k = index of lowest 0; Output f(S) = 2<sup>k</sup>
 k = 3; Output 8

### **Oblivious Sketch**

- <u>Accuracy:</u>
  - $-F_0/2 \le f(S) \le 2F_0$  with constant probability
- <u>Obliviousness</u>: distribution of f(S) is *entirely* determined by  $F_0$ 
  - similar to <u>functional privacy</u> [Feigenbaum Ishai Malkin Nissim Strauss Wright 01]
- Why it helps:
  - Pick noise  $\eta$  from discretized Lap(1/ $\epsilon$ )
  - Create new stream S' to feed to f:
    - If  $\eta < 0$ , ignore first  $\eta$  distinct elements
    - If  $\eta > 0$ , insert elements n+1, ...,  $n+\eta$
- Distribution of f(S') is a function of  $max{F_0 + \eta, 0}$ :  $\epsilon$ -DP (user)
- Error:  $F_0/2 O(1/\epsilon) \le f(S) \le 2F_0 + O(1/\epsilon)$
- Space:  $O(1/\epsilon + \log n)$ 
  - can make log *n* w.h.p. by first inserting  $O(1/\epsilon)$  elements

### **Continual Observation**

- In an online stream, often need to *track* the value of a statistic.
  - number of reported instances of a viral infection
  - sales over time
  - number of likes on Facebook
- <u>Privacy under continual observation</u> [Dwork Naor Pitassi Rothblum 10]:
  - At each time step the algorithm outputs the value of the statistic
  - The entire sequence of outputs is ε-DP (usually event level)
- Results:
  - A single counter (number of 1's in a bit stream) [DNPR10]
  - Time-decayed counters [Bolot Fawaz Muthukrishnan Nikolov Taft 13]
  - Online learning [DNPR10] [Jain Kothari Thakurta 12] [Smith Thakurka 13]
  - Generic transformation for monotone algorithms [DNPR10]

#### Binary Tree Technique [DPNR10], [Chan Shi Song 10]



#### **Binary Tree Technique**



#### **Continuous Counter**

- Achieves polylog(*m*) error per time step
- Simple variations:
  - the value of *m* is unknown
  - other statistics decomposable over time intervals
- Improve error for time-decayed statistics:
   vary the noise on different levels of the tree
- Applications to online learning
  - continuous counters  $\rightarrow$  track gradient of risk function

# Thank

# you!