

Exponential Mechanism, Private Data Release Katrina Ligett

differential privacy

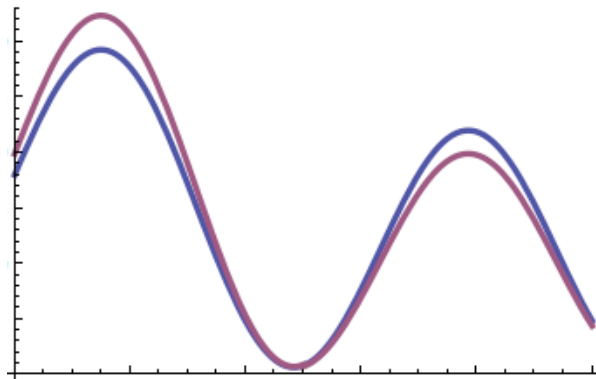
[DinurNissim03, DworkNissimMcSherrySmith06, Dwork06]

ϵ -Differential Privacy for algorithm M :

for any two neighboring data sets x_1, x_2 , differing by the addition or removal of a single row

any $S \subseteq \text{range}(M)$,

$$\Pr[M(x_1) \in S] \leq e^\epsilon \Pr[M(x_2) \in S]$$



yesterday

DP definition, properties

Randomized Response

Laplace Mechanism

reportNoisyMax

Ok, but I wanted to use my data for a scenario where direct noise addition doesn't make sense

selecting from among discrete set of alternatives

small perturbation in outcome space could be disastrous for outcome quality

exponential mechanism

[McSherryTalwar07]

Output an element $t \in \text{range}(M)$ with probability $\sim \exp(\varepsilon u(x, t) / (2 \Delta u))$

where u is a “scoring function”

privacy pretty straightforward...

utility... depends

Thm. The exponential mechanism preserves $(\epsilon, 0)$ -differential privacy.

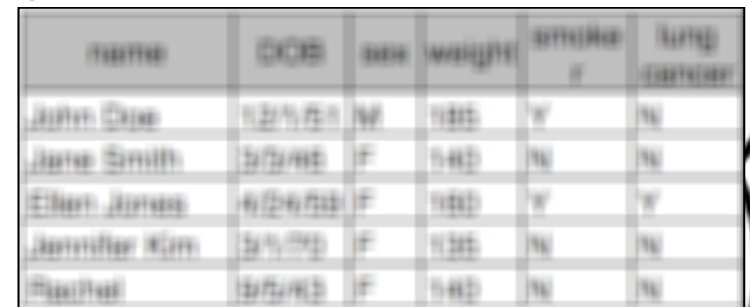
What about general accuracy theorem for Exponential Mechanism?

handling an exponential number of queries

what fraction of males over age 50? what fraction smoke and have lung cancer? what fraction of males over 150 lbs?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

...



name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel	9/5/43	F	140	N	N

Data can be “big” in two dimensions: more rows makes privacy easier (lower sensitivity); more columns makes it harder (more queries to preserve)

idea: leverage structure/overlap between queries

What if we asked the *same* count query k times?

Naively, Laplace mechanism would add fresh zero-mean noise to each instance of the query, so noise must grow with k , to avoid mean of the noisy answers converging to true mean.

Should have re-used answers.

Error scales with k . Can we get $\log k$ instead?

synthetic data (“offline case”)

succinctly represents answers to many queries

publish once and for all; anyone can run any statistic on it, as many times as you like...

...(but only certain computations will have guaranteed accuracy)

[BlumLigettRoth08]

Idea: use the Exponential Mechanism to sample a small database that answers all queries in set Q of interest well.

Will need to know there exists such a small database (sample complexity bounds from learning theory)

Will need to show we get a good output with high probability

BLR mechanism

utility function of candidate output
database y given true database x

$$u(x, y) = - \max_{f \in Q} |f(x) - f(y)|$$

Sampling bounds

Lemma. For any $x \in \mathbb{N}^{|X|}$ and any collection of linear queries Q , there exists a database y of size $\log |Q| / \alpha^2$ s.t.

$$\max_{f \in Q} |f(x) - f(y)| \leq \alpha$$

[BLR08]

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Michael Ray	3/2/81	M	200	Y	N
Fran Michaels					
Rachel Kim					
Michelle Lo					
Nira Waters	9/5/43	F	140	N	N
Jennifer Kim	3/1/70	F	135	N	N
Lisa Smith	9/5/43	F	140	N	N
Tony Miller	12/1/51	M	210	Y	N
Eve Casey	3/3/46	F	140	N	N
Paul Larson	4/24/59	F	160	Y	Y
Noelle Mason	3/1/70	F	130	N	N
Rachel Waters	9/5/43	F	140	Y	N
Shirley Wu	3/1/70	F	150	N	N
Rachel Waters	9/5/43	F	140	N	Y
Lawrence Vay	12/1/51	M	185	Y	N
Laura Rich	3/3/46	F	140	N	N

Size $O(\log |Q| / \epsilon)$

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y
YYY	1/11/74	F	130	N	Y
ZZZ	10/5/44	F	150	N	N

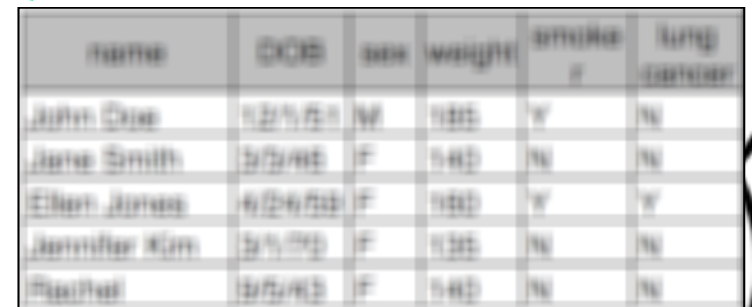
name	DOB	se x	weight	smoker	lung cancer
WWW	4/13/48	F	135	N	N
XXX	4/22/61	M	165	Y	Y

handling an exponential number of queries

what fraction of males over age 50? what fraction smoke and have lung cancer? what fraction of males over 150 lbs?

name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel Waters	9/5/43	F	140	N	N

...



name	DOB	sex	weight	smoker	lung cancer
John Doe	12/1/51	M	185	Y	N
Jane Smith	3/3/46	F	140	N	N
Ellen Jones	4/24/59	F	160	Y	Y
Jennifer Kim	3/1/70	F	135	N	N
Rachel	9/5/43	F	140	N	N

Algorithm 4 The Small Database Mechanism

SmallDB($x, Q, \varepsilon, \alpha$)

Let $\mathcal{R} \leftarrow \{y \in \mathbb{N}^{|\mathcal{X}|} : \|y\|_1 = \frac{\log |Q|}{\alpha^2}\}$

Let $u : \mathbb{N}^{|\mathcal{X}|} \times \mathcal{R} \rightarrow \mathbb{R}$ be defined to be:

$$u(x, y) = - \max_{f \in Q} |f(x) - f(y)|$$

Sample And Output $y \in \mathcal{R}$ with the exponential mechanism

$\mathcal{M}_E(x, u, \mathcal{R})$

BLR privacy

Thm. The BLR mechanism is $(\epsilon, 0)$ -differentially private.

BLR accuracy

Proposition 4.4. Let \mathcal{Q} be any class of linear queries. Let y be the database output by $\text{SmallDB}(x, \mathcal{Q}, \varepsilon, \alpha)$. Then with probability $1 - \beta$:

$$\max_{f \in \mathcal{Q}} |f(x) - f(y)| \leq \alpha + \frac{2 \left(\frac{\log |\mathcal{X}| \log |\mathcal{Q}|}{\alpha^2} + \log \left(\frac{1}{\beta} \right) \right)}{\varepsilon \|x\|_1}.$$

BLR accuracy

Theorem 4.5. By the appropriate choice of α , letting y be the database output by $\text{SmallDB}(x, \mathcal{Q}, \varepsilon, \frac{\alpha}{2})$, we can ensure that with probability $1 - \beta$:

$$\max_{f \in \mathcal{Q}} |f(x) - f(y)| \leq \left(\frac{16 \log |\mathcal{X}| \log |\mathcal{Q}| + 4 \log \left(\frac{1}{\beta} \right)}{\varepsilon \|x\|_1} \right)^{1/3}. \quad (4.2)$$

Equivalently, for any database x with

$$\|x\|_1 \geq \frac{16 \log |\mathcal{X}| \log |\mathcal{Q}| + 4 \log \left(\frac{1}{\beta} \right)}{\varepsilon \alpha^3} \quad (4.3)$$

with probability $1 - \beta$: $\max_{f \in \mathcal{Q}} |f(x) - f(y)| \leq \alpha$.

notes on the BLR mechanism

Can replace $\log |Q|$ with VCdimension
computational efficiency...