Differential privacy and generalization

Germán Bassi

School of Electrical Engineering and Computer Science KTH Royal Institute of Technology

November 18, 2019

Diving into DP

Generalization and holdout

Appendix



Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

Diving into DP

Generalization and holdout

Appendix 0000



Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix



What is privacy?

- It is the 1940s in the US; a large medical study reveals that smoking regularly has a strong connection to higher risks of developing lung cancer.
- Bob, who regularly smokes, participates in the study.
- He later sees that, as a result of this medical study, his health-insurance premiums rise.
- Clearly, the study revealed some unknown information about Bob. But, was this a breach of privacy?

Privacy in the context of DP

Differential privacy addresses the paradox of learning nothing about an individual while learning useful information about a population.¹

In the context of differential privacy (DP), we consider that there is no leakage of private information if the "impact" on the smoker is the same independent of whether or not he was in the study.



 $^{^1 \}text{C}.$ Dwork and A. Roth. "The algorithmic foundations of differential privacy." Foundations and Trends in Theoretical Computer Science 9.3–4 (2014): 211–407.

Privacy in the context of DP

Differential privacy addresses the paradox of learning nothing about an individual while learning useful information about a population.¹

In the context of differential privacy (DP), we consider that there is no leakage of private information if the "impact" on the smoker is the same independent of whether or not he was in the study.

Bob's participation should not affect the result of the study

 $^{^{1}\}text{C}.$ Dwork and A. Roth. "The algorithmic foundations of differential privacy." Foundations and Trends in Theoretical Computer Science 9.3–4 (2014): 211–407.

The basic model in DP

Differential privacy is a definition of privacy tailored to the problem of privacy-preserving data analysis:

- We assume that the information is stored in a database *D*.
- The database has *n* rows, each from a different individual.
- The database is maintained by a trusted and trustworthy curator.
- A data analyst, who does not have direct access to *D*, asks queries in order to perform some statistical analysis of the database as a whole (aggregated data).

If the data analyst does not have access to *D*, how can there be a breach of Bob's privacy?

G. Bassi (KTH EECS)

The basic model in DP

Differential privacy is a definition of privacy tailored to the problem of privacy-preserving data analysis:

- We assume that the information is stored in a database *D*.
- The database has *n* rows, each from a different individual.
- The database is maintained by a trusted and trustworthy curator.
- A data analyst, who does not have direct access to *D*, asks queries in order to perform some statistical analysis of the database as a whole (aggregated data).

If the data analyst does not have access to D, how can there be a breach of Bob's privacy?

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

A differencing attack

Person	# of cig.	Person	# of cig.
Alice	3	Alice	3
Bob	5	Charles	1
Charles	1	Diana	0
:	÷	:	÷
Zoe	4	Zoe	4
Database 1		Database 2	

- Assume you have two identical databases except for Bob's data, which it only appears in the first one.
- The innocent-looking queries "total number of participants" and "average cigarette consumption", when performed on both databases, allow for perfectly characterizing Bob's daily cigarette consumption.

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

A differencing attack

Person	# of cig.	Person	# of cig.
Alice	3	Alice	3
Bob	5	Charles	1
Charles	1	Diana	0
:	:	:	÷
Zoe	4	Zoe	4
Database 1		Database 2	

- Adding zero-mean independent random noise to the queries' answers helps us protect the participants but degrades the quality of the statistical analysis. There is a trade-off.
- If the analyst performs multiple identical queries, the effect of the noise can be averaged out.

Lessons learned

- The most detrimental thing (from the point of view of DP) is to have two almost identical databases, i.e., two databases that differ in just one entry.
 - We need to focus on databases with "distance 1".
- Some noise needs to be added to the queries' answers in order to assure some level of privacy.
 - There is a natural and unavoidable trade-off between statistical accuracy and privacy.
- Allowing multiple queries to the same database inherently decreases the level of privacy any method can achieve.
 - We cannot release unlimited number of results. Clever strategies need to be devise to release the most amount of useful information.

G. Bassi (KTH EECS)

Lessons learned

- The most detrimental thing (from the point of view of DP) is to have two almost identical databases, i.e., two databases that differ in just one entry.
 - We need to focus on databases with "distance 1".
- Some noise needs to be added to the queries' answers in order to assure some level of privacy.
 - There is a natural and unavoidable trade-off between statistical accuracy and privacy.
- Allowing multiple queries to the same database inherently decreases the level of privacy any method can achieve.
 - We cannot release unlimited number of results. Clever strategies need to be devise to release the most amount of useful information.

G. Bassi (KTH EECS)

Lessons learned

- The most detrimental thing (from the point of view of DP) is to have two almost identical databases, i.e., two databases that differ in just one entry.
 - We need to focus on databases with "distance 1".
- Some noise needs to be added to the queries' answers in order to assure some level of privacy.
 - There is a natural and unavoidable trade-off between statistical accuracy and privacy.
- Allowing multiple queries to the same database inherently decreases the level of privacy any method can achieve.
 - We cannot release unlimited number of results. Clever strategies need to be devise to release the most amount of useful information.

Generalization and holdout

Appendix



Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

First definitions

- A database x is a collection of records from a universe \mathcal{X} .
 - In our example about the daily cigarette consumption, $\mathcal{X} = \{0, 1, 2, \ldots\}$
- The database is represented by a vector of occurrences of each type in *X*, i.e., *x* ∈ N^{|X|}.
 - In our example, x = [10, 2, 4, 13, 44, 35] means that the database contains 10 participants that declared 0 cigarettes per day, 2 participants that smoke only 1 cigarette per day, ...
 - It is similar to the type in information theory but without normalization.
- A (privacy) mechanism is an algorithm that takes as inputs a database, a query, and other parameters, and returns the corresponding answer to the query in the database (hopefully with some guarantees on the level of privacy).

G. Bassi (KTH EECS) Differential privacy and generalization

First definitions (cont.)

 With our previous definition of database, a useful measure is the ℓ₁ norm. Given two databases x and y:

$$\underbrace{\|x\|_1 = \sum_{i=1}^{|\mathcal{X}|} |x_i|}_{\text{size of } x} \qquad \underbrace{\|x - y\|_1 = \sum_{i=1}^{|\mathcal{X}|} |x_i - y_i|}_{\text{number of records which are different in } x \text{ and } y}$$

 Going back to our example of the differencing attack, x represents the database which has Bob's data (he smokes 5 cigarettes a day) and y is the other almost identical database:

$$\begin{array}{ll} \mathsf{x} = [10, \, 2, \, 4, \, 13, \, 44, \, 35] \\ \mathsf{y} = [10, \, 2, \, 4, \, 13, \, 44, \, 34] \end{array} \longrightarrow \quad \|\mathsf{x} - \mathsf{y}\|_1 = 1. \end{array}$$

G. Bassi (KTH EECS)

Diving into DP

Appendix 0000

First definitions (cont.)

Randomized Algorithm

A randomized algorithm \mathcal{M} with domain A and discrete range B is associated with a mapping $M : A \to \Delta(B)$, where $\Delta(B)$ represents the probability simplex on B.

On input $a \in A$, the algorithm \mathcal{M} outputs $\mathcal{M}(a) = b$ with probability $(\mathcal{M}(a))_b$ for each $b \in B$.

Example: In our example about the daily cigarette intake, let $\mathcal{M}(x)$ calculate the (noisy) empirical mean of the $||x||_1 = n$ integer values from the database.

G. Bassi (KTH EECS)

First definitions (cont.)

<u>Remark</u>: for simplicity ignore border effects. If $\mathcal{X} = \{0, 1, 2, 3\}$ the mean should not be larger than 3. This example ignores this fact.

We take $B = \left\{\frac{k}{n}\right\}_{k=0}^{\infty}$, i.e., a quantization of the space.

The true mean is $\mu(x) = \frac{1}{n} \sum_{i=1}^{|\mathcal{X}|} (i-1) x_i$.

However, ${\mathcal M}$ outputs a random value close to the true mean:

$$\mathsf{Pr}\{\mathcal{M}(x) = b\} = \begin{cases} \frac{1}{2L+1} & \text{if } b \in \left[\mu(x) - \frac{L}{n}, \mu(x) + \frac{L}{n}\right] \\ 0 & \text{else} \end{cases}$$



G. Bassi (KTH EECS)

Definition and properties of DP

Differential Privacy

A randomized algorithm \mathcal{M} with domain $\mathbb{N}^{|\mathcal{X}|}$ is (ε, δ) -differentially private if for all $\mathcal{S} \subseteq \text{Range}(\mathcal{M})$ and for all $x, y \in \mathbb{N}^{|\mathcal{X}|}$ such that $||x - y||_1 \leq 1$:

$$\mathsf{Pr}\{\mathcal{M}(x)\in\mathcal{S}\}\leq \exp(arepsilon)\,\mathsf{Pr}\{\mathcal{M}(y)\in\mathcal{S}\}+\delta,$$

where the probability space is over the coin flips of the mechanism \mathcal{M} . If $\delta = 0$, we say that \mathcal{M} is ε -differentially private (sometimes referred to as pure DP).

<u>Remark</u>: The smaller ε and δ are, the more private the algorithm is.

G. Bassi (KTH EECS)

Definition and properties of DP (cont.)

DP:
$$\forall x, y \in \mathbb{N}^{|\mathcal{X}|}$$
 s.t. $||x - y||_1 \le 1$,
 $\Pr{\mathcal{M}(x) \in \mathcal{S}} \le \exp(\varepsilon) \Pr{\mathcal{M}(y) \in \mathcal{S}} + \delta$

If $\delta = 0$ we see that

$$\varepsilon \geq \ln \frac{\Pr{\{\mathcal{M}(x) \in \mathcal{S}\}}}{\Pr{\{\mathcal{M}(y) \in \mathcal{S}\}}},$$

i.e., the log-likelihood ratio should be bounded.

For every run of the mechanism $\mathcal{M}(x)$, the output observed is (almost) equally likely to be observed on every neighboring database.

Moreover,
$$\Pr{\{\mathcal{M}(x) \in \mathcal{S}\}} = 0 \iff \Pr{\{\mathcal{M}(y) \in \mathcal{S}\}} = 0.$$

Definition and properties of DP (cont.)

DP:
$$\forall x, y \in \mathbb{N}^{|\mathcal{X}|}$$
 s.t. $||x - y||_1 \le 1$,
 $\Pr{\mathcal{M}(x) \in \mathcal{S}} \le \exp(\varepsilon) \Pr{\mathcal{M}(y) \in \mathcal{S}} + \delta$

If $\delta \neq 0$ we see that, $\forall S \subseteq \mathsf{Range}(\mathcal{M}) \text{ s.t. } \mathsf{Pr}\{\mathcal{M}(x) \in S\} \geq \delta$,

$$arepsilon \geq \ln rac{\Pr\{\mathcal{M}(x) \in \mathcal{S}\} - \delta}{\Pr\{\mathcal{M}(y) \in \mathcal{S}\}}.$$

For very small δ , it is extremely unlikely that the observed value $\mathcal{M}(x)$ will be much more or much less likely to be generated when the database is x than when the database is y.

However, it is possible that $0 < \Pr{\mathcal{M}(x) \in S} < \delta$ and $\Pr{\mathcal{M}(y) \in S} = 0$.

G. Bassi (KTH EECS)

Definition and properties of DP (cont.)

Post-processing

Let $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \to R$ be a randomized algorithm that is (ε, δ) -DP. Let $f : R \to R'$ be an arbitrary randomized mapping. Then $f \circ \mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \to R'$ is (ε, δ) -DP.

Proof: Let x and y s.t. $||x - y||_1 \le 1$.

Assume f is deterministic, fix $S \subseteq R'$ and let $T = \{r \in R : f(r) \in S\}$. Then, we have that:

$$Pr\{f(\mathcal{M}(x)) \in S\} = Pr\{\mathcal{M}(x) \in T\}$$

$$\leq exp(\varepsilon)Pr\{\mathcal{M}(y) \in T\} + \delta$$

$$= exp(\varepsilon)Pr\{f(\mathcal{M}(y)) \in S\} + \delta.$$

Stochastic f follows as a convex combination of deterministic functions.

G. Bassi (KTH EECS)

Composition of DP mechanisms

What happens when we combine the outputs of two DP mechanisms (on the same database)?

As we saw in the first part, the strength of the privacy guarantee degrades by performing several queries.

We show first the following simple result:

Composition of two pure DP mechanisms

Let $\mathcal{M}_1 : \mathbb{N}^{|\mathcal{X}|} \to R_1$ be an ε_1 -DP algorithm and let $\mathcal{M}_2 : \mathbb{N}^{|\mathcal{X}|} \to R_2$ be an ε_2 -DP algorithm. Then their combination, defined to be $\mathcal{M}_{1,2} : \mathbb{N}^{|\mathcal{X}|} \to R_1 \times R_2$ by the mapping: $\mathcal{M}_{1,2}(x) = (\mathcal{M}_1(x), \mathcal{M}_2(x))$ is $(\varepsilon_1 + \varepsilon_2)$ -DP.

Composition of DP mechanisms (cont.)

Proof: Let $x, y \in \mathbb{N}^{|\mathcal{X}|}$ be such that $||x - y||_1 \leq 1$. Fix any $(r_1, r_2) \in R_1 \times R_2$. Then,

$$\frac{\Pr\{\mathcal{M}_{1,2}(x) = (r_1, r_2)\}}{\Pr\{\mathcal{M}_{1,2}(y) = (r_1, r_2)\}} = \frac{\Pr\{\mathcal{M}_1(x) = r_1\}}{\Pr\{\mathcal{M}_1(y) = r_1\}} \frac{\Pr\{\mathcal{M}_2(x) = r_2\}}{\Pr\{\mathcal{M}_2(y) = r_2\}} \\ \leq \exp(\varepsilon_1) \exp(\varepsilon_2) \\ = \exp(\varepsilon_1 + \varepsilon_2).$$

By symmetry,

$$\begin{aligned} \frac{\Pr\{\mathcal{M}_{1,2}(x) = (r_1, r_2)\}}{\Pr\{\mathcal{M}_{1,2}(y) = (r_1, r_2)\}} &\geq \exp(-\varepsilon_1)\exp(-\varepsilon_2)\\ &= \exp(-(\varepsilon_1 + \varepsilon_2)), \end{aligned}$$

which proves the statement.

G. Bassi (KTH EECS)

Composition of DP mechanisms (cont.) What if the algorithms are not *pure* DP?

Composition of several DP mechanisms

Let $\mathcal{M}_i : \mathbb{N}^{|\mathcal{X}|} \to R_i$ be an $(\varepsilon_i, \delta_i)$ -DP algorithm for $i \in [k]$. If $\mathcal{M}_{[k]} : \mathbb{N}^{|\mathcal{X}|} \to \prod_{i=1}^k R_i$ is defined to be $\mathcal{M}_{[k]}(x) = (\mathcal{M}_1(x), \dots, \mathcal{M}_k(x))$, then $\mathcal{M}_{[k]}$ is $(\sum_{i=1}^k \varepsilon_i, \sum_{i=1}^k \delta_i)$ -DP.

What if the algorithms are chosen adaptively?

Adaptive composition of several DP mechanisms Let $\mathcal{M}_1 : \mathbb{N}^{|\mathcal{X}|} \to R_1$ be an $(\varepsilon_1, \delta_1)$ -DP algorithm and, for $2 \le i \le k$, let $\mathcal{M}_i : \mathbb{N}^{|\mathcal{X}|} \times \prod_{j=1}^{i-1} R_j \to R_i$ be an $(\varepsilon_i, \delta_i)$ -DP algorithm for every input $(r_1, \ldots, r_{i-1}) \in \prod_{j=1}^{i-1} R_j$. Then their adaptive composition, defined to be $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \to R_k$, is $(\sum_{i=1}^k \varepsilon_i, \sum_{i=1}^k \delta_i)$ -DP.

G. Bassi (KTH EECS)

Composition of DP mechanisms (cont.) What if the algorithms are not *pure* DP?

Composition of several DP mechanisms

Let $\mathcal{M}_i : \mathbb{N}^{|\mathcal{X}|} \to R_i$ be an $(\varepsilon_i, \delta_i)$ -DP algorithm for $i \in [k]$. If $\mathcal{M}_{[k]} : \mathbb{N}^{|\mathcal{X}|} \to \prod_{i=1}^k R_i$ is defined to be $\mathcal{M}_{[k]}(x) = (\mathcal{M}_1(x), \dots, \mathcal{M}_k(x))$, then $\mathcal{M}_{[k]}$ is $(\sum_{i=1}^k \varepsilon_i, \sum_{i=1}^k \delta_i)$ -DP.

What if the algorithms are chosen adaptively?

Adaptive composition of several DP mechanisms

Let $\mathcal{M}_1 : \mathbb{N}^{|\mathcal{X}|} \to R_1$ be an $(\varepsilon_1, \delta_1)$ -DP algorithm and, for $2 \leq i \leq k$, let $\mathcal{M}_i : \mathbb{N}^{|\mathcal{X}|} \times \prod_{j=1}^{i-1} R_j \to R_i$ be an $(\varepsilon_i, \delta_i)$ -DP algorithm for every input $(r_1, \ldots, r_{i-1}) \in \prod_{j=1}^{i-1} R_j$.

Then their adaptive composition, defined to be $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \to R_k$, is $\left(\sum_{i=1}^k \varepsilon_i, \sum_{i=1}^k \delta_i\right)$ -DP.

Advanced composition

A more sophisticated argument yields significant improvement on the scaling of ε when we allow for a degradation in terms of δ .

Advanced adaptive composition

For all $\varepsilon, \delta, \delta' \ge 0$, the adaptive composition of k arbitrary (ε, δ) -DP algorithms is $(\varepsilon', k\delta + \delta')$ -DP, where

$$arepsilon' = \sqrt{2k \ln(1/\delta')} \, arepsilon + k arepsilon (extbf{e}^arepsilon - 1).$$

Assuming $\varepsilon \ll 1$, $e^{\varepsilon} \approx 1 + \varepsilon$, thus

$$\varepsilon' \approx \sqrt{2k \ln(1/\delta')} \varepsilon + k \varepsilon^2 \approx \sqrt{2k \ln(1/\delta')} \varepsilon.$$

By allowing a small δ' , this result states that $\varepsilon' = \mathcal{O}(\sqrt{k}\varepsilon)$ instead of $\varepsilon' = \mathcal{O}(k\varepsilon)$ as in the simple argument.

Advanced composition (cont.)

Parameter optimization for composition of DP mechanisms² Given target privacy parameters $\varepsilon', \delta' > 0$, to ensure ($\varepsilon', k\delta + \delta'$) cumulative privacy loss over k mechanisms, it suffices that each

mechanism is (ε, δ) -DP, where

$$\varepsilon \leq \frac{\varepsilon'}{\sqrt{2k\ln(1/\delta')}}.$$

Example: Bob is a member of $k = 10,000 \varepsilon$ -DP (uncoordinated) databases. What should be the value of ε so that Bob's cumulative privacy loss is bounded by $\varepsilon' = 1$ with probability at least $1 - e^{-32}$? Taking $\delta' = e^{-32} \approx 10^{-14}$ it suffices to have $\varepsilon \leq 1/800$. Compare this to $\varepsilon \leq 1/10,000$ from the simple composition bound, which assumes $\delta' = 0$.

 $^{^{2}}$ As we saw, this is an approximation and should be handled with care. For simplicity, we will take the inequality as true for the rest of the presentation.

Advanced composition (cont.)

The results on the composition of k (ε , δ)-DP mechanisms has two interpretations:

- One user that participates in k (ε , δ)-DP databases.
- One database that answers k (ε , δ)-DP queries.

In the following part, we will see that ε is related to the amount of noise that we need to add before answering the queries, and thus it is related to the expected distortion of the answers.

It is therefore convenient (if possible) to use the advanced composition result such that the final privacy level is $\varepsilon' = \mathcal{O}(\sqrt{k}\varepsilon)$ instead of $\varepsilon' = \mathcal{O}(k\varepsilon)$.

0000000000000000

Take-away messages

- A randomized algorithm \mathcal{M} is (ε, δ) -DP if for all databases x and $y \in \mathbb{N}^{|\mathcal{X}|}$ such that $||x - y||_1 < 1$: $\Pr{\mathcal{M}(x) \in S} \leq \exp(\varepsilon) \Pr{\mathcal{M}(y) \in S} + \delta.$
- Post-processing: differential privacy is preserved under any arbitrary randomized mapping.
- The composition of k (ε, δ)-DP mechanisms is:
 - a $(k\varepsilon, k\delta)$ -DP mechanism; or,
 - a $(\varepsilon', k\delta + \delta')$ -DP mechanism, where $\varepsilon' = \mathcal{O}(\sqrt{k}\varepsilon)$ and $\delta' > 0$.

Diving into DP

Diving into DP 000000000000000 Generalization and holdout

Appendix



Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

Diving into DP

 Generalization and holdout

Appendix

Preliminaries

The most fundamental type of queries are numeric queries, i.e., functions $f : \mathbb{N}^{|\mathcal{X}|} \to \mathbb{R}^k$.

 ℓ_1 -sensitivity

The ℓ_1 -sensitivity of a function $f : \mathbb{N}^{|\mathcal{X}|} \to \mathbb{R}^k$ is:

$$\Delta f = \max_{x,y \in \mathbb{N}^{|\mathcal{X}|} \text{ s.t. } \|x-y\|_1 = 1} \|f(x) - f(y)\|_1$$

The ℓ_1 -sensitivity of a function f captures the magnitude by which a single individual's data can change the function f the most.

This is related to the uncertainty in the response that we must introduce in order to hide the participation of a single individual.

Diving into DP

 Generalization and holdout

Appendix 0000

Preliminaries (cont.)

Examples:

- Counting query: "How many users smoke 5 cigarettes a day?" Removing one user can only affect the count by one, i.e., Δf = 1.
- *m* (general) counting queries: 1) "How many users smoke 5 cigarettes a day? and 2) "How many users are older than 30?" In the worst case, a user affects all counts, i.e., Δ*f* = *m*.
- *m* disjoint counting queries: "How many users smoke *i* cigarettes a day?" (repeat for *i* = 0 to 10³)
 A user can only affect the count of only one question by one, i.e., Δ*f* = 1.

Diving into DP 0000000000000000 DP mechanisms

Generalization and holdout

Appendix

Laplace mechanism

Laplace distribution

The Laplace distribution (centered at 0) with scale b is the distribution with probability density function:

$$Lap(x|b) = \frac{1}{2b} \exp\left(-\frac{|x|}{b}\right).$$

The variance of this distribution is $2b^2$.

For simplicity, we abuse notation and write Lap(b) when clear.

<u>Remark</u>: The Laplace distribution is a symmetric version of the exponential distribution.

Diving into DP

DP mechanisms

Generalization and holdout

Appendix 0000

Laplace mechanism (cont.)

Laplace mechanism

Given any function $f : \mathbb{N}^{|\mathcal{X}|} \to \mathbb{R}^k$, the Laplace mechanism is defined as:

$$\mathcal{M}_L(x, f, \varepsilon) = f(x) + (Y_1, \ldots, Y_k),$$

where Y_i are i.i.d. random variables drawn from Lap(σ), and $\sigma = \Delta f / \varepsilon$ is sometimes called the noise rate.

<u>Remark</u>: The larger the ℓ_1 -sensitivity of f ($\Delta f \rightarrow \infty$) or the more stringent the privacy requirement ($\varepsilon \rightarrow 0$), the larger the variance of the added noise.

G. Bassi (KTH EECS)

 Generalization and holdout

Appendix

Laplace mechanism (cont.)

Theorem

The Laplace mechanism preserves (ε , 0)-differential privacy.

Proof: Let $x, y \in \mathbb{N}^{|\mathcal{X}|}$ be such that $||x - y||_1 \leq 1$, and let $f(\cdot)$ be some function $f : \mathbb{N}^{|\mathcal{X}|} \to \mathbb{R}^k$.

Let p_x and p_y denote the probability density functions of $\mathcal{M}_L(x, f, \varepsilon)$ and $\mathcal{M}_L(y, f, \varepsilon)$, respectively. We compare the two at some arbitrary point $z \in \mathbb{R}^k$

$$\frac{p_{x}(z)}{p_{y}(z)} = \prod_{i=1}^{k} \frac{\exp\left(-\frac{\varepsilon|f(x)_{i}-z_{i}|}{\Delta f}\right)}{\exp\left(-\frac{\varepsilon|f(y)_{i}-z_{i}|}{\Delta f}\right)}$$
$$= \prod_{i=1}^{k} \exp\left(\frac{\varepsilon(|f(y)_{i}-z_{i}|-|f(x)_{i}-z_{i}|)}{\Delta f}\right)$$

G. Bassi (KTH EECS)
Diving into DP

DP mechanisms

Generalization and holdout

Appendix

Laplace mechanism (cont.)

$$\frac{p_{x}(z)}{p_{y}(z)} = \dots \leq \prod_{i=1}^{k} \exp\left(\frac{\varepsilon(|f(x)_{i} - f(y)_{i}|)}{\Delta f}\right)$$
$$= \exp\left(\frac{\varepsilon ||f(x) - f(y)||_{1}}{\Delta f}\right)$$
$$\leq \exp(\varepsilon).$$

Following similar steps, we can prove that $\frac{p_x(z)}{p_y(z)} \ge \exp(-\varepsilon)$.

G. Bassi (KTH EECS)

Accuracy of queries under the Laplace mechanism

Simple counting query: As we saw, for counting queries we have that $\Delta f = 1$, and thus ε -DP can be achieved for counting queries by the addition of noise drawn from Lap $(1/\varepsilon)$. The expected distortion, or error, is

$$\mathbb{E}[|f(x) - \mathcal{M}_L(x, f, \varepsilon)|] = \mathbb{E}[|Y|] = \frac{1}{\varepsilon},$$

which is independent of the size of the database.

<u>*m*</u> (general) counting queries: In this case, $\Delta f = m$, and thus we need to add noise drawn from Lap (m/ε) to the answer of each query in order to achieve ε -DP.

The expected distortion for each query is now

$$\mathbb{E}[|f(x)_i - \mathcal{M}_L(x, f, \frac{\varepsilon}{m})_i|] = \mathbb{E}[|Y_i|] = \frac{m}{\varepsilon},$$

G. Bassi (KTH EECS)

Accuracy of queries under the Laplace mechanism

Simple counting query: As we saw, for counting queries we have that $\Delta f = 1$, and thus ε -DP can be achieved for counting queries by the addition of noise drawn from Lap $(1/\varepsilon)$. The expected distortion, or error, is

$$\mathbb{E}[|f(x) - \mathcal{M}_L(x, f, \varepsilon)|] = \mathbb{E}[|Y|] = \frac{1}{\varepsilon},$$

which is independent of the size of the database.

<u>*m* (general) counting queries</u>: In this case, $\Delta f = m$, and thus we need to add noise drawn from Lap (m/ε) to the answer of each query in order to achieve ε -DP.

The expected distortion for each query is now

$$\mathbb{E}[|f(x)_i - \mathcal{M}_L(x, f, \frac{\varepsilon}{m})_i|] = \mathbb{E}[|Y_i|] = \frac{m}{\varepsilon},$$

G. Bassi (KTH EECS)

<u>*m*</u> (general) counting queries (advanced): Remember that by accepting a small $\delta > 0$, we can ensure (ε, δ) cumulative privacy loss over <u>*m*</u> mechanisms if each mechanism is ε_0 -DP, where

$$\varepsilon_0 \approx rac{arepsilon}{\sqrt{2m\ln(1/\delta)}}$$

Therefore, the expected distortion for each query is now

$$\mathbb{E}[|f(x)_i - \mathcal{M}_L(x, f, \varepsilon_0)_i|] = \mathbb{E}[|Y_i|] pprox rac{\sqrt{2m\ln(1/\delta)}}{arepsilon},$$

which has a lower scaling factor with respect to m.

G. Bassi (KTH EECS)

If our database has n entries, how many counting queries m can we ask without losing too much accuracy?

Usually, counting queries are used to estimate the percentage of certain statement/property on a population:

 $\mathsf{Percentage} = \frac{\mathsf{Count}}{n}.$

Therefore, given that the expected distortion is $\mathcal{O}(\sqrt{m})$, we may ask $o(n^2)$ queries with non-trivial accuracy.

Asking o(n) queries reduces the error introduced by the DP mechanism below the level of sampling error, i.e., $o(1/\sqrt{n})$.

If the Laplacian noises are not i.i.d. this can be further improved. To achieve this, we need to coordinate the answers.

G. Bassi (KTH EECS)

If our database has n entries, how many counting queries m can we ask without losing too much accuracy?

Usually, counting queries are used to estimate the percentage of certain statement/property on a population:

 $\mathsf{Percentage} = \frac{\mathsf{Count}}{n}.$

Therefore, given that the expected distortion is $\mathcal{O}(\sqrt{m})$, we may ask $o(n^2)$ queries with non-trivial accuracy.

Asking o(n) queries reduces the error introduced by the DP mechanism below the level of sampling error, i.e., $o(1/\sqrt{n})$.

If the Laplacian noises are not i.i.d. this can be further improved. To achieve this, we need to coordinate the answers.

G. Bassi (KTH EECS)

If our database has n entries, how many counting queries m can we ask without losing too much accuracy?

Usually, counting queries are used to estimate the percentage of certain statement/property on a population:

$$\mathsf{Percentage} = \frac{\mathsf{Count}}{n}.$$

Therefore, given that the expected distortion is $\mathcal{O}(\sqrt{m})$, we may ask $o(n^2)$ queries with non-trivial accuracy.

Asking o(n) queries reduces the error introduced by the DP mechanism below the level of sampling error, i.e., $o(1/\sqrt{n})$.

If the Laplacian noises are not i.i.d. this can be further improved. To achieve this, we need to coordinate the answers.

G. Bassi (KTH EECS)

If our database has n entries, how many counting queries m can we ask without losing too much accuracy?

Usually, counting queries are used to estimate the percentage of certain statement/property on a population:

$$\mathsf{Percentage} = \frac{\mathsf{Count}}{n}.$$

Therefore, given that the expected distortion is $\mathcal{O}(\sqrt{m})$, we may ask $o(n^2)$ queries with non-trivial accuracy.

Asking o(n) queries reduces the error introduced by the DP mechanism below the level of sampling error, i.e., $o(1/\sqrt{n})$.

If the Laplacian noises are not i.i.d. this can be further improved. To achieve this, we need to coordinate the answers.

G. Bassi (KTH EECS)

If our database has n entries, how many counting queries m can we ask without losing too much accuracy?

Usually, counting queries are used to estimate the percentage of certain statement/property on a population:

$$\mathsf{Percentage} = \frac{\mathsf{Count}}{n}.$$

Therefore, given that the expected distortion is $\mathcal{O}(\sqrt{m})$, we may ask $o(n^2)$ queries with non-trivial accuracy.

Asking o(n) queries reduces the error introduced by the DP mechanism below the level of sampling error, i.e., $o(1/\sqrt{n})$.

If the Laplacian noises are not i.i.d. this can be further improved. To achieve this, we need to coordinate the answers.

Diving into DP

 Generalization and holdout

Appendix

Report Noisy Max

Imagine we want to determine which of m counting queries has the highest value when the queries are not disjoint.

Releasing all the counts and letting the data analyst find the maximum is suboptimal. Remember that the vector of counts has high ℓ_1 -sensitivity, specifically, $\Delta f = m$.

A better approach is to (internally) add independently generated Laplace noise $Lap(1/\varepsilon)$ to each count and return the index of the largest noisy count.

The Report Noisy Max algorithm is $(\varepsilon, 0)$ -differentially private.

G. Bassi (KTH EECS)

Diving into DP

 Generalization and holdout

Appendix

Report Noisy Max

Imagine we want to determine which of m counting queries has the highest value when the queries are not disjoint.

Releasing all the counts and letting the data analyst find the maximum is suboptimal. Remember that the vector of counts has high ℓ_1 -sensitivity, specifically, $\Delta f = m$.

A better approach is to (internally) add independently generated Laplace noise $Lap(1/\varepsilon)$ to each count and return the index of the largest noisy count.

The Report Noisy Max algorithm is $(\varepsilon, 0)$ -differentially private.

G. Bassi (KTH EECS)

Diving into DP

 Generalization and holdout

Appendix

Report Noisy Max

Imagine we want to determine which of m counting queries has the highest value when the queries are not disjoint.

Releasing all the counts and letting the data analyst find the maximum is suboptimal. Remember that the vector of counts has high ℓ_1 -sensitivity, specifically, $\Delta f = m$.

A better approach is to (internally) add independently generated Laplace noise $Lap(1/\varepsilon)$ to each count and return the index of the largest noisy count.

The Report Noisy Max algorithm is $(\varepsilon, 0)$ -differentially private.

G. Bassi (KTH EECS)

Diving into DP

 Generalization and holdout

Appendix

Report Noisy Max

Imagine we want to determine which of m counting queries has the highest value when the queries are not disjoint.

Releasing all the counts and letting the data analyst find the maximum is suboptimal. Remember that the vector of counts has high ℓ_1 -sensitivity, specifically, $\Delta f = m$.

A better approach is to (internally) add independently generated Laplace noise $Lap(1/\varepsilon)$ to each count and return the index of the largest noisy count.

The Report Noisy Max algorithm is $(\varepsilon, 0)$ -differentially private.

Sparse vector mechanism

- The Laplace mechanism can be used to answer adaptively chosen low-sensitivity queries.
 - The privacy loss increases proportionally with the number of queries answered (or its square root)
 - What do we do if we want to answer a very large number of queries with reasonable accuracy?
- In some situations, we may only care on identifying the queries that lie above a certain threshold.
 - If we only report when the answer to a query exceeds the threshold, we can show that privacy degrades only with the number of answered queries, rather than with the total.
 - This can be a huge savings if we know that the set of queries that lie above the threshold is much smaller than the total number of queries that is, if the answer vector is sparse.

G. Bassi (KTH EECS)

Sparse vector mechanism

- The Laplace mechanism can be used to answer adaptively chosen low-sensitivity queries.
 - The privacy loss increases proportionally with the number of queries answered (or its square root)
 - What do we do if we want to answer a very large number of queries with reasonable accuracy?
- In some situations, we may only care on identifying the queries that lie above a certain threshold.
 - If we only report when the answer to a query exceeds the threshold, we can show that privacy degrades only with the number of answered queries, rather than with the total.
 - This can be a huge savings if we know that the set of queries that lie above the threshold is much smaller than the total number of queries that is, if the answer vector is sparse.

- There are *m* total number of sensitivity-1 queries, which may be chosen adaptively.
- There is a single threshold *T* fixed in advance but it is possible for each query to have its own threshold.
- We add noise to query values and compare the results to T.
- We expect a small number *c* of noisy values to exceed the threshold, and we only release the noisy values above the threshold.
- The algorithm stops after releasing c values.

- There are *m* total number of sensitivity-1 queries, which may be chosen adaptively.
- There is a single threshold *T* fixed in advance but it is possible for each query to have its own threshold.
- We add noise to query values and compare the results to T.
- We expect a small number *c* of noisy values to exceed the threshold, and we only release the noisy values above the threshold.
- The algorithm stops after releasing c values.

- There are *m* total number of sensitivity-1 queries, which may be chosen adaptively.
- There is a single threshold *T* fixed in advance but it is possible for each query to have its own threshold.
- We add noise to query values and compare the results to T.
- We expect a small number *c* of noisy values to exceed the threshold, and we only release the noisy values above the threshold.
- The algorithm stops after releasing c values.

- There are *m* total number of sensitivity-1 queries, which may be chosen adaptively.
- There is a single threshold *T* fixed in advance but it is possible for each query to have its own threshold.
- We add noise to query values and compare the results to T.
- We expect a small number *c* of noisy values to exceed the threshold, and we only release the noisy values above the threshold.
- The algorithm stops after releasing c values.

- There are *m* total number of sensitivity-1 queries, which may be chosen adaptively.
- There is a single threshold *T* fixed in advance but it is possible for each query to have its own threshold.
- We add noise to query values and compare the results to T.
- We expect a small number *c* of noisy values to exceed the threshold, and we only release the noisy values above the threshold.
- The algorithm stops after releasing *c* values.

Detour: AboveThreshold algorithm

The following is a special case with c = 1 released values.

procedure AboveThreshold $(D, \{f_i\}, T, \varepsilon)$ 1: $\hat{T} \leftarrow T + \text{Lap}(\frac{2}{2})$ > Add noise to the threshold 2: for all f; do 3 $\nu_i \leftarrow Lap(\frac{4}{c})$ \triangleright Noise for the query value 4: if $f_i(D) + \nu_i > \hat{T}$ then 5: output $a_i = \top$ 6: halt ▷ We stop after a positive match 7: else 8. output $a_i = \bot$ 9: end if 10: end for 11: 12: end procedure

Detour: AboveThreshold algorithm (cont.)

The AboveThreshold algorithm is $(\varepsilon, 0)$ -differentially private.

Sketch of Proof: It is a sequential proof, i.e., we assume that we are on round *i* and the previous noises ν_1, \ldots, ν_{i-1} are fixed.

We then calculate the probability of halting at round *i* (by integrating over the p.d.f. of \hat{T} and ν_i) for two almost-equal databases.

The proof looks like the composition of two DP mechanisms: one hiding the effect of previous rounds (a 1-sensitivity query) and one hiding the effect of the present round (a 2-sensitivity query).

The parameters of the Laplacian noises of \hat{T} and ν_i are selected such that the composition produces a ε -DP mechanism.

Sparse vector mechanism (again)

The Sparse algorithm is constructed as follows:

- Each received query is forwarded to the AboveThreshold algorithm.
- Each time a query produces a result exceeding the noisy threshold (and AboveThreshold halts), Sparse restarts AboveThreshold and continues feeding the remaining queries.
- The algorithm halts after AboveThreshold has been restarted *c* times.
- Since each instance of AboveThreshold is *ε*-DP, the composition results apply:
 - we may obtain a $(c\varepsilon, 0)$ -DP mechanism; or,
 - we may obtain a ($\varepsilon',\delta')\text{-}\mathsf{DP}$ mechanism with the advanced composition result.

G. Bassi (KTH EECS)

Sparse vector mechanism (again)

The Sparse algorithm is constructed as follows:

- Each received query is forwarded to the AboveThreshold algorithm.
- Each time a query produces a result exceeding the noisy threshold (and AboveThreshold halts), Sparse restarts AboveThreshold and continues feeding the remaining queries.
- The algorithm halts after AboveThreshold has been restarted *c* times.
- Since each instance of AboveThreshold is *ε*-DP, the composition results apply:
 - we may obtain a $(c\varepsilon, 0)$ -DP mechanism; or,
 - we may obtain a ($\varepsilon',\delta')\text{-}\mathsf{DP}$ mechanism with the advanced composition result.

G. Bassi (KTH EECS)

Sparse vector mechanism (again)

The Sparse algorithm is constructed as follows:

- Each received query is forwarded to the AboveThreshold algorithm.
- Each time a query produces a result exceeding the noisy threshold (and AboveThreshold halts), Sparse restarts AboveThreshold and continues feeding the remaining queries.
- The algorithm halts after AboveThreshold has been restarted *c* times.
- Since each instance of AboveThreshold is *ε*-DP, the composition results apply:
 - we may obtain a $(c\varepsilon, 0)$ -DP mechanism; or,
 - we may obtain a ($\varepsilon',\delta')\text{-}\mathsf{DP}$ mechanism with the advanced composition result.

G. Bassi (KTH EECS)

Sparse vector mechanism (again)

The Sparse algorithm is constructed as follows:

- Each received query is forwarded to the AboveThreshold algorithm.
- Each time a query produces a result exceeding the noisy threshold (and AboveThreshold halts), Sparse restarts AboveThreshold and continues feeding the remaining queries.
- The algorithm halts after AboveThreshold has been restarted *c* times.
- Since each instance of AboveThreshold is *ε*-DP, the composition results apply:
 - we may obtain a $(c\varepsilon, 0)$ -DP mechanism; or,
 - we may obtain a ($\varepsilon',\delta')\text{-}\mathsf{DP}$ mechanism with the advanced composition result.

DP mechanisms Sparse vector mechanism (cont.) 1: procedure Sparse(D, { f_i }, T, c, ε , δ) if $\delta = 0$ then $\sigma \leftarrow 2\frac{c}{c}$ else $\sigma \leftarrow 2\frac{\sqrt{2c\ln(1/\delta)}}{c}$ 2: Select noise level 3. $\hat{T}_0 \leftarrow T + \mathsf{Lap}(\sigma)$ 4: *count* \leftarrow 0 Initialize the counter 5: for all f; do 6: $\nu_i \leftarrow Lap(2\sigma)$ if $f_i(D) + \nu_i \geq \hat{T}_{count}$ then 7: 8: output $a_i = \top$ 9. $count \leftarrow count + 1$ Increase the counter $T_{count} \leftarrow T + Lap(\sigma)$ Restart AboveThreshold 10: 11: else 12: output $a_i = \bot$ end if 13: 14: if count > c then 15: halt \triangleright Halt if we have produced *c* positive results 16: end if 17: end for 18: end procedure

DP mechanisms 00000000000000000000

Take-away messages

- Given a function f with ℓ_1 -sensitivity equal to Δf , we can construct a $(\varepsilon, 0)$ -DP mechanism to answer it by adding $Lap(\Delta f/\varepsilon)$ noise to the result.
- The composition of k queries, each one being $(\varepsilon, 0)$ -DP, is:
 - (kε, 0)-DP; or

$$DP$$
 with $c' \sim \sqrt{2k \ln t}$

- (ε', δ') -DP with $\varepsilon' \approx \sqrt{2k \ln(1/\delta')} \varepsilon$ for $\delta' > 0$.
- We can answer more queries if we reveal less information about them, e.g., by only saying that the result is above a certain threshold.

Diving into DP

Generalization and holdout

Appendix 0000



Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

The problem of generalization in data analysis

- It is commonly assumed that each analysis procedure operates on a freshly sampled dataset, or at least, it is validated on a freshly sampled holdout (or testing) set.
- However, in practice the holdout dataset is rarely used only once.
- If the set of all tested hypotheses is known and independent of the holdout set, then it is easy to account for such multiple testing.
- However, such approach does not apply if the estimates or hypotheses tested on the holdout are chosen adaptively: that is, if the choice of hypotheses depends on previous analyses performed on the dataset.



Model description

- Adaptive data analysis is seen as a sequence of procedures: $\mathcal{A}_1 \rightarrow \mathcal{A}_2 \rightarrow \ldots \rightarrow \mathcal{A}_m$.
- These steps analyze a fixed dataset S = (x₁,..., x_n) drawn from a distribution D over Xⁿ.
- Each step is an algorithm $\mathcal{A}_i : \mathcal{X}^n \times \prod_{j=1}^{i-1} \mathcal{Y}_j \to \mathcal{Y}_i$ that takes the fixed dataset S and the previous results as inputs.
- It is assumed that each individual algorithm generalizes well when executed on a fresh dataset sampled from \mathcal{D} for every input vector (y_1, \ldots, y_{i-1}) .

Model description (cont.)

- For a function $\phi : \mathcal{X} \to \mathbb{R}$:
 - If we fix a dataset S, we denote the average value of φ on S as:
 *ε*_S[φ] = ¹/_n Σⁿ_{i=1} φ(x_i).
 - If we take a distribution \mathcal{P} over \mathcal{X} , we denote the expected value of ϕ as: $\mathcal{P}[\phi] = \mathbb{E}_{x \sim \mathcal{P}}[\phi(x)]$.
- In the (common) situation that the samples are taken i.i.d. according to a distribution *P* over *X*, we say that the random dataset *S* is drawn from a distribution *Pⁿ* over *Xⁿ*.
- When we say that an "algorithm generalizes well when executed on a fresh dataset," we mean that $\mathcal{E}_{S}[\phi]$ and $\mathcal{P}[\phi]$ should be close, for ϕ and S chosen independently.
- The random dataset S and the randomly selected function φ might not be independent in an adaptive setting, and thus *ε_s*[φ] and *P*[φ] might not be close.

G. Bassi (KTH EECS)

Unrealistic example

- Assume that the dataset S contains n standard Gaussian RVs, i.e, $X \sim \mathcal{N}(0, 1)$.
- Further assume that the dataset is partitioned in n/ln(n) subsets, each one containing ln(n) samples. For simplicity, all these values are assumed to be integers.
- The query φ_j is the sample mean of the jth subset, whose answer is a zero-mean Gaussian RV with variance 1/ln(n).
- If we were to randomly select $\sqrt{n/\ln(n)}$ subsets and average the queries' responses (final query $\overline{\phi}$), we would have an unbiased estimator of the true mean with a smaller variance $(1/\sqrt{n\ln(n)})$.
- But what if we do not choose these subsets randomly?

Unrealistic example (cont.)

- Assume that we randomly group the $n/\ln(n)$ subsets in $\sqrt{n/\ln(n)}$ groups, each containing $\sqrt{n/\ln(n)}$ subsets.
- After receiving the queries' responses, we pick the maximum from each group and we create the final query as the sample mean on those particular groups.

G. Bassi (KTH EECS)

Unrealistic example (cont.)

• In contrast to choosing a random subset, whose sample mean behaves as expected: $\phi_j \sim \mathcal{N}(0, \frac{1}{\ln(n)})$, choosing the maximum among the sample means in the group has a particular p.d.f. that does not concentrate around the true mean of X.



Unrealistic example (cont.)

Therefore, this (unrealistic) method for improving the sample convergence does not result in *E*_S[φ̄] being close to *P*[φ̄], even though *E*_S[φ_j] is close to *P*[φ_j] for every fixed φ_j.


Generalization via differential privacy

Generalization via DP for i.i.d. datasets

Let us define \mathcal{A} , \boldsymbol{S} , and \boldsymbol{Y} , where

- $\mathcal{A}: \mathcal{X}^n \to \mathcal{Y}$ and it is an ε -DP algorithm,
- **S** is a RV drawn from a distribution \mathcal{P}^n over \mathcal{X}^n , and
- $\boldsymbol{Y} = \mathcal{A}(\boldsymbol{S})$ is the corresponding output distribution.

Assume that for each element $y \in \mathcal{Y}$ there is a subset $R(y) \subseteq \mathcal{X}^n$ such that $\max_{y \in \mathcal{Y}} \Pr\{\mathbf{S} \in R(y)\} \leq \beta$. Then, for $\varepsilon \leq \sqrt{\frac{\ln(1/\beta)}{2n}}$ we have $\Pr\{\mathbf{S} \in R(\mathbf{Y})\} \leq 3\sqrt{\beta}$.

Interpretation: In the statement, R(y) denotes the set of datasets for which y is a "bad" output of the algorithm. If setting A to output a fix y is a bad idea with probability at most β , then the output $\mathbf{Y} = \mathcal{A}(\mathbf{S})$ of an ε -DP algorithm has a bounded probability of being bad for the particular input \mathbf{S} .

G. Bassi (KTH EECS)

Generalization via differential privacy (cont.)

A statistical query is defined by a function $\phi : \mathcal{X} \to [0, 1]$ and a tolerance τ . For a distribution \mathcal{P} over \mathcal{X} a valid response to such a query is any value v such that $|v - \mathcal{P}[\phi]| \leq \tau$.

By Hoeffding's inequality, a fixed query function verifies that

$$\Pr\{\left|\mathcal{P}[\phi] - \mathcal{E}_{\boldsymbol{S}}[\phi]\right| > \tau\} \le 2\exp(-2n\tau^2).$$

Corollary for statistical queries

Let \mathcal{A} be an ε -DP algorithm that outputs a function from \mathcal{X} to [0,1]. For a distribution \mathcal{P} over \mathcal{X} , let \boldsymbol{S} be a random variable distributed according to \mathcal{P}^n and let $\phi = \mathcal{A}(\boldsymbol{S})$.

Then for any
$$au > 0$$
, setting $\varepsilon \leq \sqrt{ au^2 - \frac{\ln(2)}{2n}}$ ensures that

$$\mathsf{Pr}\big\{\big|\mathcal{P}[\phi] - \mathcal{E}_{\boldsymbol{S}}[\phi]\big| > \tau\big\} \le 3\sqrt{2}\exp\big(-n\tau^2\big).$$

Appendix 0000

Generalization via differential privacy (cont.)

Proof: Apply general theorem with

•
$$R(\phi) = \{ S \in \mathcal{X}^n : |\mathcal{P}[\phi] - \mathcal{E}_S[\phi]| \ge \tau \}$$
 and
• $\beta = 2 \exp(-2n\tau^2).$

This results in

•
$$\varepsilon \leq \sqrt{\frac{\ln(1/\beta)}{2n}} = \sqrt{\tau^2 - \frac{\ln(2)}{2n}}$$
 and
• $\Pr{\{\mathbf{S} \in R(\mathbf{Y})\}} = \Pr{\{|\mathcal{P}[\phi] - \mathcal{E}_{\mathbf{S}}[\phi]| > \tau\}}$
 $\leq 3\sqrt{\beta}$
 $= 3\sqrt{2}\exp(-n\tau^2).$

<u>Remark</u>: Both papers show the same result but with different wording.

Unrealistic example (revisited)

We revisit the Gaussian example where now the $n/\ln(n)$ queries ϕ_j are answered by a DP mechanism following the previous guidelines. Once the appropriate subsets are selected, the final query $\bar{\phi}$ is answered by a non-DP mechanism. This is possible since there will be no further adaptivity.



Max-divergence

For two random variables X and Y over the same domain \mathcal{X} , the max-divergence of X from Y is defined³ as

$$D_{\infty}(\boldsymbol{X} \| \boldsymbol{Y}) = \log \max_{x \in \mathcal{X}} \frac{\Pr{\{\boldsymbol{X} = x\}}}{\Pr{\{\boldsymbol{Y} = x\}}}$$

while the $\delta\text{-approximate}$ max-divergence is defined as

$$D^{\delta}_{\infty}(\boldsymbol{X} \| \boldsymbol{Y}) = \log \max_{\mathcal{O} \subseteq \mathcal{X}, \Pr\{\boldsymbol{X} \in \mathcal{O}\} \geq \delta} \frac{\Pr\{\boldsymbol{X} \in \mathcal{O}\} - \delta}{\Pr\{\boldsymbol{Y} \in \mathcal{O}\}}$$

A randomized algorithm \mathcal{A} with domain \mathcal{X}^n is (ε, δ) -DP if for all pairs of datasets that differ in a single element $S, S' \in \mathcal{X}^n$:

$$D^\delta_\inftyig(\mathcal{A}(\mathcal{S})\|\mathcal{A}(\mathcal{S}')ig) \leq (\log e)\,arepsilon.$$

³In the following, log and In denote the logarithm on base 2 and *e*, respectively.

G. Bassi (KTH EECS)

Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

Max-information

Max-information

Let X and Y be jointly distributed random variables. Then, the max-information between X and Y is defined as

$$I_{\infty}(\boldsymbol{X};\boldsymbol{Y}) = D_{\infty}((\boldsymbol{X},\boldsymbol{Y}) \| \boldsymbol{X} imes \boldsymbol{Y}),$$

while the β -approximate max-information is defined as:

$$I^{eta}_{\infty}(oldsymbol{X};oldsymbol{Y})=D^{eta}_{\infty}ig((oldsymbol{X},oldsymbol{Y})\|oldsymbol{X} imesoldsymbol{Y}ig).$$

Some (simple) properties:

- If $I_{\infty}(\boldsymbol{X}; \boldsymbol{Y}) = k$, then $\Pr{\{\boldsymbol{X} = x \mid \boldsymbol{Y} = y\}} \le 2^k \Pr{\{\boldsymbol{X} = x\}}$.
- The max-information is an upper bound of the mutual information: *I*(*X*; *Y*) ≤ *I*_∞(*X*; *Y*).

• For
$$\beta \geq 0$$
, $I_{\infty}^{\beta}(\boldsymbol{X}; \boldsymbol{Y}) = I_{\infty}^{\beta}(\boldsymbol{Y}; \boldsymbol{X}).$

The max-information of an algorithm

Max-information of an algorithm

We say that a (randomized) algorithm \mathcal{A} has β -approximate max-information of value k if for every distribution \mathcal{D} over \mathcal{X}^n , $I_{\infty}^{\beta}(\boldsymbol{S}; \mathcal{A}(\boldsymbol{S})) \leq k$, where \boldsymbol{S} is a dataset chosen randomly according to \mathcal{D} . We denote this by $I_{\infty}^{\beta}(\mathcal{A}, n) \leq k$.

For the particular case of $\beta = 0$, we can alternatively define the (pure) max-information of algorithm A as

$$I_{\infty}(\mathcal{A}, n) = \max_{S, S' \in \mathcal{X}^n} D_{\infty}(\mathcal{A}(S) \| \mathcal{A}(S')),$$

where S and S' are any two datasets.

G. Bassi (KTH EECS)

The max-information of an algorithm (cont.)

Post-processing

Let $\mathcal{A}: \mathcal{X}^n \to \mathcal{Y}$ and $\mathcal{B}: \mathcal{Y} \to \mathcal{Y}'$ be two randomized algorithms. Then, the algorithm $\mathcal{B} \circ \mathcal{A}$ with domain \mathcal{X}^n and range \mathcal{Y}' satisfies that, for every random variable \boldsymbol{S} over \mathcal{X}^n and every $\beta \geq 0$, $I_{\infty}^{\beta}(\boldsymbol{S}; \mathcal{B} \circ \mathcal{A}(\boldsymbol{S})) \leq I_{\infty}^{\beta}(\boldsymbol{S}; \mathcal{A}(\boldsymbol{S})).$

Composition of algorithms

Let $\mathcal{A} : \mathcal{X}^n \to \mathcal{Y}$ be an algorithm such that $I_{\infty}^{\beta_1}(\mathcal{A}, n) \leq k_1$, and let $\mathcal{B} : \mathcal{X}^n \times \mathcal{Y} \to \mathcal{Z}$ be an algorithm such that, for every $y \in \mathcal{Y}$, $I_{\infty}^{\beta_2}(\boldsymbol{S}; \mathcal{B}(\boldsymbol{S}, y)) \leq k_2$. Let $\mathcal{C} : \mathcal{X}^n \to \mathcal{Z}$ be defined such that $\mathcal{C}(S) = \mathcal{B}(S, \mathcal{A}(S))$. Then, $I_{\infty}^{\beta_1 + \beta_2}(\mathcal{C}, n) \leq k_1 + k_2$. For a longer composition, we have that $I_{\infty}^{\sum_j \beta_j}(\mathcal{C}', n) \leq \sum_i k_j$.

The max-information of an algorithm (cont.)

Usefulness of max-information

Let **S** be a random dataset in \mathcal{X}^n and \mathcal{A} be an algorithm with range \mathcal{Y} such that for some $\beta \geq 0$, $I_{\infty}^{\beta}(\mathcal{A}, n) = k$. Then for any event $\mathcal{O} \subseteq \mathcal{X}^n \times \mathcal{Y}$,

$$\mathsf{Pr}\big\{\big(\boldsymbol{S},\mathcal{A}(\boldsymbol{S})\big)\in\mathcal{O}\big\}\leq 2^k\mathsf{Pr}\{\boldsymbol{S}\times\mathcal{A}(\boldsymbol{S})\in\mathcal{O}\}+\beta.$$

In particular,

$$\Pr\left\{\left(\boldsymbol{S}, \mathcal{A}(\boldsymbol{S})\right) \in \mathcal{O}\right\} \leq 2^{k} \max_{y \in \mathcal{Y}} \Pr\left\{\left(\boldsymbol{S}, y\right) \in \mathcal{O}\right\} + \beta.$$

As before, this allows us to bound the probability of any joint event between the random dataset and random output of the algorithm, i.e., an adaptive behavior, with respect to the independent selection of \boldsymbol{S} and $\mathcal{A}(\boldsymbol{S})$.

The max-information of DP algorithms

The (approximate) max-information of \mathcal{A} , i.e., $I_{\infty}^{\beta}(\boldsymbol{S}; \mathcal{A}(\boldsymbol{S})) = I_{\infty}^{\beta}(\mathcal{A}, n)$, is a measure of how much the random dataset \boldsymbol{S} affects the distribution of the algorithm's output. In other words, how much information from \boldsymbol{S} is leaked to the output of \mathcal{A} .

Pure max-information of pure DP algorithms

Let \mathcal{A} be an ε -DP algorithm. Then $I_{\infty}(\mathcal{A}, n) \leq (\log e) \varepsilon n$.

Proof: Given that any two datasets S and S' might differ in at most n elements, for every output y of A we have that

$$\Pr{\{\boldsymbol{Y} = y \mid \boldsymbol{S} = S\}} \le \exp(\varepsilon n) \Pr{\{\boldsymbol{Y} = y \mid \boldsymbol{S} = S'\}},$$

which implies that $D_{\infty}(\mathcal{A}(S) \| \mathcal{A}(S')) \leq (\log e) \varepsilon n$. Therefore,

$$I_{\infty}(\mathcal{A}, n) = \max_{S, S' \in \mathcal{X}^n} D_{\infty}(\mathcal{A}(S) \| \mathcal{A}(S')) \leq (\log e) \varepsilon n.$$

G. Bassi (KTH EECS)

The max-information of DP algorithms (cont.)

The previous result applied to general distributions of datasets but we can obtain a tighter bound for datasets with i.i.d. samples.

Approx. max-information of DP algorithms for i.i.d. datasets Let us define A, S, and Y, where

- $\mathcal{A}: \mathcal{X}^n \to \mathcal{Y}$ and it is an ε -DP algorithm,
- **S** is a RV drawn from a distribution \mathcal{P}^n over \mathcal{X}^n , and
- $\boldsymbol{Y} = \mathcal{A}(\boldsymbol{S})$ is the corresponding output distribution.

Then for any $\beta > 0$, $I_{\infty}^{\beta}(\mathcal{A}, n) \leq (\log e) (\varepsilon^2 n/2 + \varepsilon \sqrt{n \ln(2/\beta)/2}).$

For small ε and fixed $\beta > 0$, the bound on $I_{\infty}^{\beta}(\mathcal{A}, n)$ is $\mathcal{O}(\varepsilon \sqrt{n})$. For comparison, remember that the composition of $k \varepsilon$ -DP mechanisms results in a (ε', δ') -DP mechanism with $\varepsilon' = \mathcal{O}(\varepsilon \sqrt{k})$.

Generalization via max-information

Let \mathcal{A} be an algorithm that outputs a function $f : \mathcal{X}^n \to \mathbb{R}$ of sensitivity c and define the "bad event" \mathcal{O}_{τ} when the empirical estimate of f is far from its expected value, i.e.,

$$\mathcal{O}_{ au} = ig\{(\mathcal{S}, f) : ig| f(\mathcal{S}) - \mathcal{D}[f] ig| \geq au ig\},$$

where \mathcal{D} is a general distribution over \mathcal{X}^n .

By McDiarmid's inequality, if the samples are i.i.d. $(\mathcal{D} = \mathcal{P}^n)$, then for every possible function f we have that

$$\mathsf{Pr}\{(\boldsymbol{S},f)\in\mathcal{O}_{\tau}\}=\mathsf{Pr}\{\left|f(\boldsymbol{S})-\mathcal{P}^{n}[f]\right|\geq\tau\}\leq 2\exp\left(-\frac{2\tau^{2}}{nc^{2}}\right).$$

G. Bassi (KTH EECS)

Generalization via max-information (cont.)

Generalization via max-information for i.i.d. datasets Let us define A, S, and f, where

- \mathcal{A} is an algorithm that outputs a *c*-sensitive function $f: \mathcal{X}^n \to \mathbb{R}$,
- \boldsymbol{S} is a random dataset drawn according to \mathcal{P}^n , and
- $\boldsymbol{f} = \mathcal{A}(\boldsymbol{S})$ is the corresponding random output of \mathcal{A} .

If for
$$\beta \ge 0$$
 and $\tau > 0$, $I_{\infty}^{\beta}(\mathcal{A}, n) \le (\log e) \frac{\tau^2}{nc^2}$, then
 $\Pr\{|\mathbf{f}(\mathbf{S}) - \mathcal{P}^n[\mathbf{f}]| \ge \tau\} \le 2\exp\left(-\frac{\tau^2}{nc^2}\right) + \beta.$

In particular, if \mathcal{A} is a $\left(\frac{\tau^2}{n^2c^2}\right)$ -DP algorithm, then $\Pr\{\left|\boldsymbol{f}(\boldsymbol{S}) - \mathcal{P}^n[\boldsymbol{f}]\right| \ge \tau\} \le 2\exp\left(-\frac{\tau^2}{nc^2}\right).$

Generalization via max-information (cont.)

If f is a statistical query $\phi : \mathcal{X} \to [0, 1]$, we have that:

•
$$f(S) = \frac{1}{n} \sum_{i=1}^{n} \phi(x_i) = \mathcal{E}_S[\phi]$$

- $c = \frac{1}{n}$, and
- $\mathcal{P}^n[f] = \frac{1}{n} \sum_{i=1}^n \mathcal{P}^n[\phi(x_i)] = \mathcal{P}[\phi].$

Therefore, the previous result may be rewritten as follows:

If
$$\mathcal{A}$$
 is a τ^2 -DP algorithm, then $I_{\infty}(\mathcal{A}, n) \leq (\log e) n\tau^2$ and
 $\Pr\{|\mathcal{E}_{\boldsymbol{S}}[\phi] - \mathcal{P}[\phi]| \geq \tau\} \leq 2\exp(-n\tau^2).$

Compare to our former bound that for \mathcal{A} being an ε -DP algorithm, setting $\varepsilon \leq \sqrt{\tau^2 - \frac{\ln(2)}{2n}}$ ensures that $\Pr\{|\mathcal{E}_{\mathbf{S}}[\phi] - \mathcal{P}[\phi]| \geq \tau\} \leq 3\sqrt{2}\exp(-n\tau^2).$

G. Bassi (KTH EECS)

Generalization via max-information (cont.)

Given that the previous result requires a very restrictive privacy level of $\varepsilon = \tau^2$, we can try to use the other (tighter) bound on generalization via max-information.

Generalization via max-information for i.i.d. datasets (take 2) Let us define \mathcal{A} , \boldsymbol{S} , and \boldsymbol{f} as before. If \mathcal{A} is $(\frac{\tau}{nc})$ -DP algorithm, then $\Pr\{|\boldsymbol{f}(\boldsymbol{S}) - \mathcal{P}^{n}[\boldsymbol{f}]| \geq \tau\} \leq 2\exp\left(-\frac{3}{4}\frac{\tau^{2}}{nc^{2}}\right)$, for large $\frac{\tau^{2}}{nc^{2}}$ (more than 10 seems OK).

Proof: Let $\beta = 2 \exp\left(-\frac{\tau^2}{nc^2}\right)$, then using the tighter bound on the approx. max-information we have $I_{\infty}^{\beta}(\mathcal{A}, n) \leq (\log e) \frac{\tau^2}{nc^2} \left(\frac{1}{2} + \frac{1}{\sqrt{2}}\right)$. Applying this to the bound on the prob. of joint events we obtain $\Pr\left\{\left|\boldsymbol{f}(\boldsymbol{S}) - \mathcal{P}^n[\boldsymbol{f}]\right| \geq \tau\right\} \leq 2 \exp\left(\left(\frac{1}{\sqrt{2}} - \frac{3}{2}\right)\frac{\tau^2}{nc^2}\right) + 2 \exp\left(-\frac{\tau^2}{nc^2}\right)$.

G. Bassi (KTH EECS)

Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

Reusable holdout

In this last part, we present one differentially private algorithm to answer adaptively chosen statistical queries on a holdout set: the Thresholdout algorithm which is based on the "Sparse Vector Algorithm."

The algorithm accepts any statistical query $\phi : \mathcal{X} \to [0, 1]$ and its goal is to provide an estimate of $\mathcal{P}[\phi]$ using a holdout set for validation.

Thresholdout

Thresholdout is given access to:

- the training dataset S_t,
- the holdout dataset S_h, and
- a budget limit B.

Given a function ϕ , Thresholdout checks if

 $\left|\mathcal{E}_{S_t}[\phi] - \mathcal{E}_{S_h}[\phi]\right| \leq T + \eta,$

where T is a fixed threshold and η is Laplace noise:

- If below the threshold, then the algorithm returns $\mathcal{E}_{S_t}[\phi]$.
- If above the threshold, then the algorithm returns $\mathcal{E}_{S_h}[\phi] + \xi$, for another Laplacian noise variable ξ , and the budget B is reduced by one.

Thresholdout (cont.)

Given a budget limit B > 0 and a noise rate $\sigma > 0$, the Thresholdout algorithm is:

• $\left(\frac{2B}{\sigma n}\right)$ -differentially private, or

•
$$\left(\frac{2\sqrt{2B}\ln(2/\delta)}{\sigma n},\delta\right)$$
-differentially private for any $\delta > 0$.

Sketch of proof: Thresholdout is an instance of the Sparse Vector Algorithm (when we check if above or below the threshold) together with the Laplace mechanism (when we release the average on S_h).

Our previous analysis of Sparse was for counting queries, which are 1-sensitive, but now we deal with empirical averages, which are $\frac{1}{n}$ -sensitive queries.



Thresholdout (cont.)

In our instance, the holdout set S_h is the private data set, and each function ϕ corresponds to the following query:

$$f_{\phi}(S_h) = \big| \mathcal{E}_{S_h}[\phi] - \mathcal{E}_{S_t}[\phi] \big|.$$

Thresholdout then is equivalent to the following procedure:

- We run the sparse vector algorithm with c = B, queries f_{ϕ} for each function ϕ , and noise rate 2σ .
- Whenever an above-threshold query is found, we release its value using the Laplace mechanism with noise rate *σ*.

Each mechanism is simultaneously:

•
$$\left(\frac{B}{\sigma n}\right)$$
-differentially private, or

•
$$\left(\frac{\sqrt{2B \ln(2/\delta)}}{\sigma n}, \frac{\delta}{2}\right)$$
-differentially private⁴ for any $\delta > 0$.

⁴This value differs from the one in the article by a factor of 2. The reason is that I am being consistent with our definition of the best ε used in the *advanced* composition of *B* mechanisms.

into DP DP mechanisms

Generalization and holdout



Take-away messages

- Differentially private algorithms can be used to obtain bounds on the probability of overfitting in adaptive data analysis given functions that "generalizes well" in a non-adaptive setting.
- The max-information of an algorithm is an interesting concept that generalizes DP in the context of data analysis.
 - The max-information exhibits similar properties as DP.
 - DP algorithms have bounded max-information.
- Thresholdout is a specific DP mechanism that allows the reuse of a holdout dataset with rigorous generalization guarantees.

Preliminaries

Diving into DP

Generalization and holdout

Appendix •000



Preliminaries

Diving into DP

DP mechanisms

Generalization and holdout

Appendix

Appendix

Notes on the unrealistic example

Let $X_i \sim \mathcal{N}(0, \sigma)$ and let $Y = \max_{i \in [n]} X_i$, then

$$\mathsf{Pr}\{Y \le y\} = \Phi\left(\frac{y}{\sigma}\right)^n,$$

where Φ is the CDF of the standard Gaussian distribution. Consequently, the PDF of Y is

$$f_{\mathbf{Y}}(\mathbf{y}) = \frac{n}{\sigma} \phi\left(\frac{\mathbf{y}}{\sigma}\right) \Phi\left(\frac{\mathbf{y}}{\sigma}\right)^{n-1},$$

where ϕ is the PDF of the standard Gaussian distribution.

Notes on the unrealistic example (cont.)

To create the Laplace mechanism for this example, I did as follows:

- Each query ϕ_j is the sample mean of the *j*th subset, i.e., $\phi_j \sim \mathcal{N}(0, 1/\ln(n)).$
- The subsets are disjoint, and thus the $n/\ln(n)$ queries are disjoint. This helps reduce the amount of noise added.
- With probability $1 2Q(3) \approx 0.9973$, each sample satisfies $|X_i| \leq 3$. Therefore, I assumed that $\Delta \phi_j = \frac{3}{\ln(n)}$ (with high probability) and neglected the effects of $|X_i| > 3$.
- Given that $\Pr\{|\mathcal{E}_{\mathbf{S}}[\phi_j]| \ge t\} = 2Q(t\sqrt{\ln(n)}) = \beta$, according to the generalization result, I took $\varepsilon = \sqrt{\frac{\ln(1/\beta)}{2\ln(n)}}$ and added independent noises $\operatorname{Lap}\left(\frac{\Delta\phi_j}{\varepsilon}\right)$ to every element of the subset before choosing the maximum. Only the index was reported.



References

- C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Foundations and Trends in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth, "Preserving Statistical Validity in Adaptive Data Analysis," arXiv:1411.2664 [cs], Mar. 2016.
- C. Dwork, V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. Roth, "Generalization in Adaptive Data Analysis and Holdout Reuse," arXiv:1506.02629 [cs], Jun. 2015.

G. Bassi (KTH EECS)