

Machine learning in science and society

From automated science to beneficial artificial intelligence

Christos Dimitrakakis

October 25, 2022

Contents

2	Introduction	9
2.1	Introduction to machine learning	9
2.1.1	Data analysis, learning and planning	9
2.1.2	Experiment design	14
2.1.3	Bayesian inference.	15
2.1.4	Book overview	17
3	Privacy	21
3.1	Things we do with data	22
3.2	Statistical disclosure	23
3.3	Algorithmic privacy.	25
3.4	k -anonymity	26
3.5	Differential privacy	27
3.6	The local and central privacy models	34
3.6.1	Properties of differential privacy.	35
3.7	The Laplace mechanism	36
3.8	Interactive data access.	39
3.8.1	Utility of queries	40
3.9	Advanced topics.	41
3.9.1	Relaxations of differential privacy.	41
3.10	Discussion	43
3.11	Exercises	45
9	Fairness	135
9.1	Fairness in machine learning	135
9.2	Concepts of fairness	135
5	Decision and Learning Theory	49
5.1	Hierarchies of decision making problems	49
5.1.1	Simple decision problems	49
5.1.2	Decision rules	52
5.1.3	Statistical testing [*]	53
5.2	Formalising classification problems	61
5.3	Beliefs and probabilities	63
5.3.1	Probability and Bayesian inference	68
5.4	Classification with stochastic gradient descent [*]	72
5.4.1	Neural network models	74
5.5	Nearest neighbours	77

5.6	Naive Bayes classifiers*	82
4	Reproducibility	27
4.1	Reproducibility	27
4.1.1	The human as an algorithm	30
4.1.2	Algorithmic sensitivity	32
4.1.3	Beyond the data you have: simulation and replication	34
6	Causality	97
6.1	Introduction	97
6.1.1	Decision diagrams	100
6.1.2	Common structural assumptions	101
6.2	Interventions	105
6.3	Policy evaluation and optimisation	107
6.4	Individual effects and counterfactuals	109
6.4.1	Disturbances and structural equation models	110
6.4.2	Example: Learning instrumental variables	111
6.4.3	Discussion	112
6.4.4	Exercises	112
7	Experiment design	115
7.1	Introduction	115
7.2	Bandit problems	116
7.2.1	An example: Bernoulli bandits	117
7.2.2	The stochastic n -armed bandit problem	118
7.2.3	Estimation and Robbins-Monro approximation	119
7.2.4	Decision-theoretic bandit process	121
7.2.5	Heuristic algorithms for the n -armed bandit problem	124
7.3	Contextual Bandits	125
7.4	Case study: experiment design for clinical trials	126
7.4.1	Practical approaches to experiment design	127
8.1	Graphical models	124
8.1.1	Inference and prediction in graphical models	128
8.1.2	Testing conditional independence	130
8.1.3	Hierarchical Bayesian models	131

Chapter 3

Privacy

One interpretation of privacy is simply the ability to maintain a personal secret. Is it possible to maintain perfect secrecy? Can the thoughts in our head be perfectly safe, or can they be revealed indirectly through our actions?

While we certainly *can* securely store and transmit data using cryptography, the problems we will discuss in this chapter are *not* solvable solely through cryptography. We are interested in scenarios where we must make a public decision or release some summary statistics that depend on private data, in such a way as to minimise harm to the individuals contributing their data.

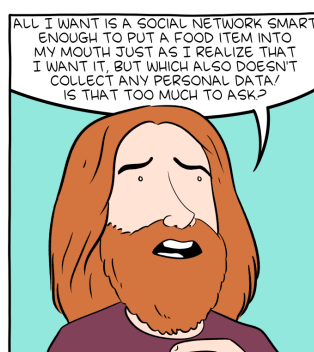
If we never have to reveal any of our computations, cryptography is what is needed. For example, through homomorphic computation, an untrusted party can even perform some computations on encrypted data, returning an encrypted result to us, while learning nothing about the original secret. As long as the data, and the results of any computation on it, are kept under lock and key, our personal information cannot be revealed.

However, sometimes we must make public decisions or release public information based on this data. Then it can be revealed indirectly. For example, you can trust your doctor to maintain confidentiality, but when you go to the pharmacy to get the prescribed medicine, somebody can infer the medical condition you suffer from.

It is even possible to learn personal information from aggregate statistics. Let us say your doctor publishes a list of cases of different diseases every week, together with some other information such as the approximate patient age. Even though your own data is mixed with that of all other patients, it is possible to infer your diagnosis, especially with some additional side-information: If somebody knows you were the only person in that age group visiting the doctor that week, they will learn your diagnosis.

From that point of view, it is not the data itself, but the complete process of data collection, treatment and public release that can be characterised as private or non-private. For that reason, we will emphasise an *algorithmic* view of privacy: participants entrusts their data to an algorithm, which produces a useful output in return. The algorithmic process is typically not fully automated, as it also depends on some human input: One example is a medical study examining different treatments for a disease. While humans select and administer the treatments, they will typically rely on a randomised strategy for assigning treatments to individuals, and use a statistical method to report their results. Given this mixture of ad-hoc decisions and formal algorithmic methods, is it possible to guarantee privacy in any sense? What kind of guarantees can we make?

Generally speaking, an algorithm has good privacy properties, if the amount of information that can be revealed through the algorithm's output about any individual contributing data



Just because they're the problem,
doesn't mean we aren't.

is bounded. In particular, we are interested in how much an adversary can learn about any individual's input to the algorithm from the algorithm's output.

This does not preclude learning general facts about individuals from the output. For example, a study about the use of steroids in sports may show that 90% of sprinters with times under 10 seconds are using steroids, while only 50% of slower sprinters do so. Any sprinter with a time under 10 seconds is thus suspected of using steroids by association. However, it does not matter if their data has been used in the study. The publication of the result does not impact their privacy, but the amount of harm it does to them does not depend in their participation: the same statistical result would have been obtained with or without them.

In this chapter, we will look at two formal concepts of privacy protection: k -anonymity and *differential privacy*. The first is a simple method for anonymising databases. However, it provides only limited resistance to identification. The latter is a more general concept, which provides full information-theoretic protection to individuals. A major problem with any privacy definition and method, however is correct interpretation of the privacy concept used, and correct implementation of the algorithm used.

3.1 Things we do with data

How do we use data in the first place? Sometimes we simply publish the data, perhaps after some initial processing. Publication of datasets is useful for researchers that want to do perform further analysis on the data. Usually though, we collect data in order to calculate specific statistics. For example the census collects data about the number of people in different households, wages, etc, and then publishes tables detailing the average age and number of people per household in different areas of the country. This demographic information is useful for policy makers, urban planning to organise voting centers and the distribution of police and fire stations, as well as other public services.

Fundamentally, privacy in statistics is an *issue of trust*. The analyst, be it a human, or an automated service, will use your data to make decisions. You must also decide who to trust and how much. Do we trust the data analyst? How much privacy are we willing to sacrifice to the analyst? How much to the public at large? What you want out of the service. Is the service important enough to sacrifice significant amounts of privacy? What is an acceptable trade-off between utility and privacy? These are difficult questions and are hard to quantify, hence we assume that we have already decided how to answer them, and we simply want to find an appropriate methodology for achieving a good result.

Consider a researcher wishing to collect data for a statistical analysis. If the analysis is eventually published,¹ this creates two possible scenarios for that may lead to privacy violations.

- Publication of “anonymised” data. Sometimes we may collect data in order to publish the dataset itself for other researchers to use. This is common practice in machine learning, with image classification datasets being a good example. However, there is always some privacy risk even if identifying information such as names are removed.
- Public data analysis. In this setting, we only publish summary statistics or models about the data. A prime example is a national census analysis, which may provide detailed demographic information for all towns and regions in a nation. Although only aggregate data is published, it is theoretically possible to infer personal information. For that reason, the US Census Bureau conducted its analysis in 2020 using differential private algorithms.

Cryptography is not enough.

Cryptography provides secure communication and computation, authentication and verification. These are used to establish secure channels with somebody that we trust. However, the privacy violations we are concerned with relate to *publicly released* outputs of algorithms. It is not important whether or not all the data and computation are encrypted: as long as the algorithm generates a public output, it is theoretically possible for somebody to learn something about the algorithm’s input.

3.2 Statistical disclosure

“Data is everywhere”, said the statistician, “data, data!”

In the past, statistical analysis was performed with laboriously collected and annotated datasets. Even as recently as in the early 21st century, databases for machine learning were limited to a few thousand entries at most. At the time of writing, not only has the size of datasets become extremely large, but the sources of data are much more diverse. Data are collected and commercialised whenever we visit a website and even as we walk around with our phone. To a limited extent, there is a tradeoff between what we can get out of a service and what we pay into it. Many free services such as navigation software rely on collecting user data to perform better: If you can tell that there is a traffic jam in Central Avenue, you can after all try and take another route. As long as informed consent exists, use of private data is generally regarded as unproblematic.²

However, even apparently benign data collected with appropriate consent can lead to serious and unexpected privacy violations. There are three famous examples of this: Firstly, the identification of people in supposedly anonymous health data in the 1990s in the state of Massachusetts, which we will go over in detail in this chapter. Secondly, the identifications of users through anonymised movie ratings in the Netflix dataset. Finally, the ability to discover if any given individual’s data is contained in a pooled genomic study.

Anonymisation

¹If somebody knows that the analysis is being conducted, however, they could still learn something private from the fact that the analysis has *not* been published.

²This is actually underscored by the GDPR legislation, which focuses on consent and data use methods.

Data is collected for many reasons. Any typical service you might want to use will require a minimal amount of data. For example, a dating service will at a minimum require your age and location, as shown in the example below.

Birthday	Name	Height	Weight	Age	Postcode	Profession
06/07	Li Pu	190	80	60-70	1001	Politician
06/14	Sara Lee	185	110	70+	1001	Rentier
01/01	A. B. Student	170	70	40-60	6732	Time Traveller

EXAMPLE 1 (Typical relational database in a dating website.). If somebody is a user in a dating website, you expect them to give some minimal personal information to be stored in the site's database. This might include their birthday, location and profession, for example.

When you submit your data to a service, you expect it to be used responsibly. For the dating service, you expect it to use the information to find good matches, based on your preferences and location. Whenever somebody uses the service, they obtain some information about you, at least indirectly. For example, if they make a query for similar singles in their neighbourhood, and they see your profile, then they have learned that you live nearby.

If we wish to publish a database, frequently we need to protect the identities of the people involved. A simple idea is to erase directly identifying information. However, this does not really work most of the time, especially since attackers can have side-information that can reveal the identities of individuals in the original data, when combined with the published dataset.

The simple act of hiding or using random identifiers is called anonymisation. However this is generally insufficient as other identifying information may be used to re-identify individuals in the data. In particular, even if somebody is unable to infer the information of any individual from the published, anonymised, dataset, they may be able to do so via some side-information. All that is needed is another dataset with some columns in common with the dataset they want to target.

Record linkage

As an example, in the 1990s a the governor of Massachussets, decided to publish anonymised information about the health records of individual state employs. They were careful to hide all obviously identifying information such as their name, and thus claimed that there are no potential privacy violations. However, they left in some information that they thought would be useful for researchers: the postcode, the birthdate, and the sex of each individual.

This allowed a PhD student, McSweeney, to perform the following linkage attack. She order a floppy disk containing database of voting records in the state. This contained names and addresses, as well as the postcode, birth date and sex of every voter. In that way, she was able to cross-reference this data, and so obtain the identities of individuals in that database. The first record she obtained was that of the governor himself, Bill Weld. Later, she estimated that approximately 87% of Americans are uniquely identifiable through those three attributes.

Clearly, anonymisation is not enough. So, is there a way to formally guarantee privacy? In the next section we will go over the solution of McSweeney, which provides us with an algorithmic definition of anonymity. While this provides some degree of protection against certain linkage attacks, it is sadly insufficient to protect privacy in general. Later, we will introduce the notion of differential privacy, which protects individuals against statistical disclosure in an information-theoretic manner.

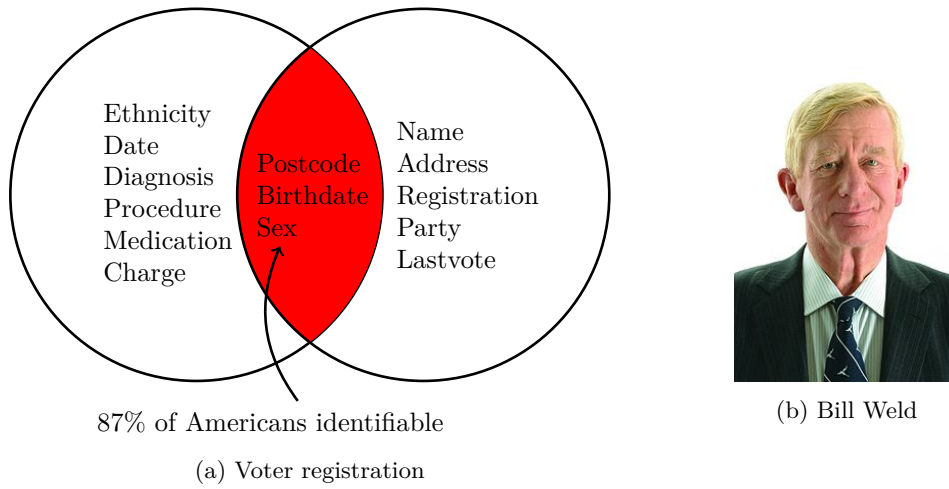


Figure 3.1: Linkage attack on an anonymised patient information database. A misguided attempt to release medical information as public data (left circle in figure (a)) by then-governor Bill Weld resulted in a simple linkage attack. This was done by accessing a party registration database (right circle in figure (b)) which contained three common fields with the medical database.

3.3 Algorithmic privacy.

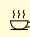
The above discussion should help emphasise that *a dataset* cannot be characterised as private or non-private. We should actually start from the observation that *any data contributed by an individual* should be considered private, because it could be combined with other datasets to obtain sensitive personal information.

For that reason, we wish to ensure that, whenever individuals contribute data for processing, the result of the processing cannot be used to reveal whatever information they had contributed. Thus, the idea of privacy only applies on *algorithms* that are used on personal data. Generally speaking, we wish for algorithms to both be useful and have good privacy properties. Unfortunately, as we shall see in the sequel, there is a distinct trade-off between privacy and utility. However, let us specify more precisely what we mean by an algorithm.

What is an algorithm?

Any functional process is an algorithm. We typically identify inputs with observations or features \mathcal{X} , and outputs with actions \mathcal{A} by the algorithm: i.e. we view the algorithm as taking actions that have an observable effect on its external environment. We generally consider *stochastic* algorithms.

Definition 3.3.1 (Stochastic algorithm π). A stochastic algorithm π with input domain \mathcal{X} and output domain \mathcal{A} is a mapping $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$ from observations \mathcal{X} to distributions over outputs $\Delta(\mathcal{A})$. When \mathcal{A} is finite, we will write $\pi(a \mid x)$ to denote the conditional probability that the algorithm outputs a given input x .

 **The general case.**

It is best to think of $\pi(\cdot|x)$ as a probability measure over \mathcal{A} . Then we write

$$\pi(A | x) \triangleq \mathbb{P}_\pi(a \in A | x), \quad A \subset \mathcal{A},$$

for the conditional probability that algorithm's output is in some set $A \subset \mathcal{A}$ given input x . This allows us to treat the case when \mathcal{A} is continuous or discrete with the same notation.

⚙️ Algorithmic randomness.

We can construct a random algorithm through access to a random coin ω taking values in Ω . We can then define the output through a deterministic function $a_\pi(\omega, x)$. Since ω is random, the output of the function is also random. If P is the probability distribution of ω , then:

$$\pi(A | x) = P(\{\omega \in \Omega \mid a_\pi(\omega, x) \in A\}),$$

i.e. the probability that the algorithm's output is in A is equal to the measure of the values of ω for which the $a_\pi(\omega, x) \in A$.

What is private?

3.4 k -anonymity

k -anonymity



(a) Samarati



(b) Sweeney

The concept of k -anonymity was introduced by Samarati and Sweeney¹⁷ and provides good guarantees against inferring personal information from a single database. This requires the analyst to first determine the variables of interest, and then determine which variables could be potentially used to identify somebody in the database. *It's the analyst's job to define quasi-identifiers.*

Definition 3.4.1 (k -anonymity). A database provides k -anonymity if for every person in the database is indistinguishable from $k - 1$ persons with respect to *quasi-identifiers*.

This hope is that, if the database satisfies k -anonymity it can be safely released, without revealing any private information directly. Let us now walk through an extended example.

***k*-anonymity example**

In particular, let us say that the analyst simply wants to calculate some statistics about how different professions correlate with age, weight, height and where people live. Some areas of the country might produce more politicians, for example. And taller people may be more successful in politics. The initial data collected might look like the table below. It was obvious to the analyst, that even if he did remove all the names, somebody knowing where A. B. Student lived and saw the table would have little trouble recognising them.

Birthday	Name	Height	Weight	Age	Postcode	Profession
06/07	Li Pu	190	80	65	1001	Politician
06/14	Sara Lee	185	110	67	1001	Rentier
06/12	Nikos Karavas	180	82	72+	1243	Politician
01/01	A. B. Student	170	70	52	6732	Time Traveller
05/08	Li Yang	175	72	35	6910	Politician

Table 3.1: 1-anonymity.

After thinking about it for a bit, the analyst decides to remove the birthday and name, and broadly categorise people according to their height in increments of 10cm, the weight in increments of 20cm, and keep just the first digit of the postcode. Now that looked much more reasonable. Still, somebody that knows that Li Yang and A. B. Student are in the table, as well as their postcodes, and they also know that Li Yang is a politician, can infer that A. B. Student is a Time Traveller.

Height	Weight	Age	Postcode	Profession
180-190	80+	60+	1*	Politician
180-190	80+	60+	1*	Rentier
180-190	80+	60+	1*	Politician
170-180	60-80	20-60	6*	Time Traveller
170-180	60-80	20-60	6*	Politician

Table 3.2: 2-anonymity: the database can be partitioned in sets of at least 2 records

However, with enough information, somebody may still be able to infer something about the individuals. In the example above, it remains true that if somebody knows that both A. B. Student and Li Yang are in the database, as well as their postcodes, as well as that Li Yang is a politician, they can infer A. B. Student's profession. Fortunately, there is a way to protect individual information from adversaries with arbitrary side-information. This is given by differential privacy.

3.5 Differential privacy

This section introduces one of the main tools for giving formal guarantees about the privacy of any algorithm ran on a dataset, *differential privacy*. This will provide individual-level privacy, in the sense that an algorithm that is differentially private guarantees that no adversary can significantly increase their knowledge about any particular individual by observing the algorithm's output. This is independent of the adversary's existing knowledge, or computational power.

While *k*-anonymity can protect against specific re-identification attacks when used with care, it says little about what to do when the adversary has a lot of knowledge. For example, if the

adversary knows the data of everybody that has participated in the database, it is trivial for them to infer what our own data is. For some particularly sensitive datasets, we may want for the adversary to be unable to tell whether or not your data was part of the base. Differential privacy offers protection against adversaries with unlimited side-information or computational power. Informally, an algorithmic computation is differentially-private if an adversary cannot distinguish two “neighbouring” database based on the result of the computation. Informally, two databases are neighbours when they are identical apart from the data of one person. A differentially private algorithm, because of its randomness, makes it impossible for somebody to tell from the algorithm’s output whether any specific individual’s data was in the database.

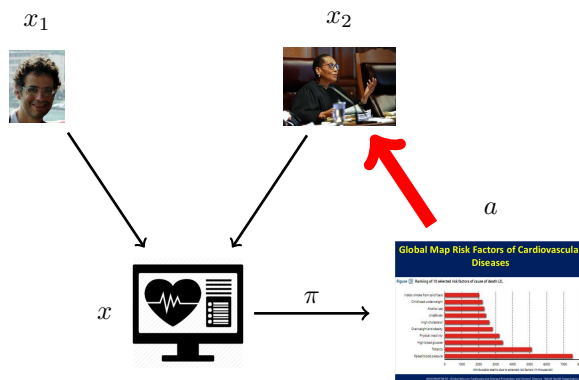


Figure 3.3: If two people contribute their data $x = (x_1, x_2)$ to a medical database, and an algorithm π computes some public output a from x , then it should be hard infer anything about the data from the public output.

Consider the example given in Figure 3.3, where two people contribute their data to a medical database. The i -th individual contributes data x_i , and the complete dataset is $x = (x_1, x_2)$. The algorithm π defines a distribution over the set of possible outputs $a \in \mathcal{A}$, with $\pi(a|x)$ being the probability that the algorithm outputs a if the data is x . If the algorithm was deterministic, then it might be possible for an adversary to invert the computation and obtain x from a . But even if that is not possible, maybe they can learn *something* about the data from the output. In the section below, we will formalise this notion.

Privacy desiderata

Consider a scenario where n persons give their data x_1, \dots, x_n to an analyst. This analyst then performs some calculation on the data and publishes the result through a randomised algorithm π , where for any output a , and any dataset x , $\pi(a|x)$ is the probability that the algorithm generates a . The following properties are desirable from a general standpoint.

Anonymity. Individual participation in the study remains a secret. From the release of the calculations results, nobody can significantly increase their probability of identifying an individual in the database.

Secrecy. The data of individuals is not revealed. The release does not significantly increase the probability of inferring individual’s information x_i .

Side-information. Even if an adversary has arbitrary side-information, he cannot use that to amplify the amount of knowledge he would have obtained from the release.

Utility. The released result has, with high probability, only a small error relative to a calculation that does not attempt to safeguard privacy.

⚙️ The prevalence of drugs in sport

Let's say you need to perform a statistical analysis of the drug-use habits of athletes. Obviously, even if you promise the athlete not to reveal their information, you still might not convince them. Yet, you'd like them to be truthful. The trick is to allow them to randomly change their answers, so that you can't be *sure* if they take drugs, no matter what they answer.

🎲 Algorithm for randomising responses about drug use

1. Flip a coin.
2. If it comes heads, respond truthfully.
3. Otherwise, flip another coin and respond **yes** if it comes heads and **no** otherwise.

📊 Calculating the true rate of responses.

Assume that the observed rate of positive responses in a sample is p , that everybody follows the protocol, and the coin is fair. Then, what is the true rate q of drug use in the population?

The problem with this approach, of course, is that we are effectively throwing away half of our data sources. In particular, if we repeated the experiment with a coin that came heads at a rate ϵ , then our error bounds would scale as $O(1/\sqrt{\epsilon n})$ for n data points.

This algorithm is very specific: it assumes binary responses, and it uses a fair coin, which introduces a lot of noise. Since the coin flips make the responses noisy, we may want to have some way of controlling it. In general, we want to consider an algorithm that takes data x_1, \dots, x_n from n users transforms it randomly to a_1, \dots, a_n using the following mapping.

The binary randomised response mechanism

Definition 3.5.1 (Randomised response). The i -th user, whose data is $x_i \in \{0, 1\}$, responds with $a_i \in \{0, 1\}$ with probability

$$\pi(a_i = j \mid x_i = k) = p, \quad \pi(a_i = k \mid x_i = k) = 1 - p,$$

where $j \neq k$.

Given the complete data x , the algorithm's output is $a = (a_1, \dots, a_n)$. Since the algorithm independently calculates a new value for each data entry, the output probability is

$$\pi(a \mid x) = \prod_i \pi(a_i \mid x_i)$$

This mechanism satisfies the formal notion of ϵ -differential privacy, which will be given in Definition 3.5.2. In a more general setting, we may have multiple possible responses. While the algorithm can be trivially generalised to n -ary outputs, the special case of when the outputs are integers is deferred until later.

⚙️ The original randomised response mechanism¹⁹

As first proposed by Warner, the mechanism distributes spinners to people. The spinner has a probability p of landing on A and $1 - p$ of landing on B . This can be easily arranged by having the corresponding areas to have the appropriate proportions. The interesting thing about this mechanism is that the responders must merely say whether or not the spinner landed on the group they identify with. They do *not* have to reveal a group. This makes the mechanism feel like you are revealing less, even if it is just a special case of a general randomised response mechanism.

What can we learn from the output?

In the Bayesian setting, we can think of the adversary as having some prior belief $\beta(x_i)$, expressed as a probability distribution over the secret value of each individual.

After the adversary observes the output, they can form a posterior belief $\beta(x_i | a_i)$, representing the information they collected. The following example quantifies the amount of information gained by the adversary.

EXAMPLE 2. For simplicity, consider only one individual, and let the adversary have a prior $\beta(x = 0) = 1 - \beta(x = 1)$ over the values of the true response of an individual. We use the randomised response mechanism with parameter p and the adversary observes the randomised data $a = 1$ for that individual, then what is $\mathbb{P}_\beta^\pi(x = 1 | a = 1)$, assuming the adversary knows the mechanism?

Proof. Bayes's theorem states that³

$$\mathbb{P}_\beta^\pi(x | a) = \pi(a | x)\beta(x) / \mathbb{P}_\beta^\pi(a),$$

where

$$\mathbb{P}_\beta^\pi(a) = \sum_{x'} \pi(a | x')\beta(x').$$

Let $q = \beta(x = 1)$. Then we have:

$$\mathbb{P}_\beta^\pi(x = 1 | a = 1) = pq / [pq + (1 - p)(1 - q)].$$

It is particularly interesting to consider the case where $q = 1/2$. Then

$$\mathbb{P}_\beta^\pi(x = 1 | a = 1) = p,$$

so we only have limited evidence for whether $x = 1$. Now consider the case where we have some arbitrary prior q , and $p = 1/2$. This means that the output of the algorithm is completely random. Consequently:

$$\beta(x = 1 | a = 1) = q = \beta(x = 1).$$

So, in this scenario we learn nothing from the algorithm's output. \square

³If there are too many symbols for you, you can write Bayes theorem simply with $P(x|a) = P(a|x)P(x)/P(a)$.

Differential privacy.

Now let us take a look at a way to characterise the the inherent privacy properties of algorithms. This is called differential privacy, and it can be seen as a bound on the information an adversary with arbitrary power or side-information could extract from the result of a computation π on the data. For reasons that will be made clear later, this computation has to be stochastic.

 ϵ -Differential Privacy

Definition 3.5.2. A stochastic algorithm $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$, where \mathcal{X} is endowed with a neighbourhood relation N , is said to be ϵ -differentially private if

$$\left| \ln \frac{\pi(A | x)}{\pi(A | x')} \right| \leq \epsilon, \quad \forall xNx', \quad A \subset \mathcal{A}. \quad (3.5.1)$$

This form is related to standard notions of statistical distance such as the KL divergence $\sum_a \ln \frac{\pi(a|x)}{\pi(a|x')} \pi(a|x)$. However, the above inequality can be equivalently rewritten as

$$\pi(A | x) \leq e^\epsilon \pi(A | x').$$

Frequently (and particularly when \mathcal{A} is finite) we can also use $\pi(a | x)$ without any technical difficulties.

Typically, algorithms are applied to datasets $x = (x_1, \dots, x_n)$ composed of the data of n individuals. Thus, all privacy guarantees relate to the data contributed by these individuals.

Neighbourhoods

Differential privacy guarantees that it is hard to distinguish neighbouring datasets. Hence, the definition of neighbourhood we use reflects what we want to protect. It makes sense to define neighbourhoods in terms of changes in one individual's data, because then an adversary cannot learn about the values of any particular individual.

In this book we will use two definitions of neighbourhoods. The first is constructed so that the adversary cannot distinguish whether or not any particular individual's information is in the dataset. If not, then they cannot infer the individual's presence from the output of the algorithm.

Definition 3.5.3 (Addition/deletion neighbourhood). If two datasets x, x' are neighbours, then we write xNx' , and it holds that

$$x = (x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n), \quad x' = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

for some i , i.e. if one dataset is missing an element.

This is an important concept if the participation in the dataset is itself sensitive. For example, if somebody is enlisted in study for experimental treatment of a rare disease, then the mere fact that they are part of the study is strong evidence that they have the disease.

The second definition is slightly weaker. It determines whether or not we can distinguish between different *values* of individual data. If not, then they cannot infer whether the data submitted by the individual has a particular value.

Definition 3.5.4 (Edit neighbourhood). We say that two datasets x, x' are neighbours, and write xN_ex' if

$$x = (x_1, \dots, x_i, \dots, x_n), \quad x' = (x_1, \dots, x'_i, \dots, x_n), \quad x_i \neq x'_i.$$

i.e. if one dataset has an altered element.

If x, x' are 1-neighbours under the second definition, then they are 2-neighbours under the first definition. To see this, create a new dataset \hat{x} from x , without the data of person i . Then $xN\hat{x}$ and $\hat{x}Nx'$, since we can change the data of one person by removing the original data x_i and then re-inserting the altered data x'_i . This is illustrated in the example below.

⚙️ Neighbourhood example

In this example, we have three datasets, x, x' and $hat{x}$. In the second dataset, \hat{x} , the highlighted row in x is missing. In the third dataset, x' , another row with the same name is added, but the height and weight are changed.

Birthday	Name	Height	Weight
06/07	Li Pu	190	80
06/14	Sara Lee	185	110
06/12	John Smith	170	82
01/01	A. B. Student	170	70
05/08	Li Yang	175	72

Table 3.3: Data x

In fact, \hat{x} is obtained through a deletion from x . Of course, the operation can be reversed: x can be obtained from an addition to \hat{x} .

Birthday	Name	Height	Weight
06/07	Li Pu	190	80
06/14	Sara Lee	185	110
01/01	A. B. Student	170	70
05/08	Li Yang	175	72

Table 3.4: \hat{x} , 1-Neighbour of x .

We can now instead add another row to \hat{x} , which will be similar to the row we had removed. This will have the effect of altering the original data in x .

Birthday	Name	Height	Weight
06/07	Li Pu	190	80
06/14	Sara Lee	185	110
06/12	John Smith	180	80
01/01	A. B. Student	170	70
05/08	Li Yang	175	72

Table 3.5: x' , 2-Neighbour of x .

Since x, x' only differ in the contents of a single row, they are edit-neighbours.

As we hinted earlier, the randomised response mechanism satisfies ϵ -DP. The ϵ parameter is dependent on p , with higher values giving a smaller ϵ , and thus better privacy protection.

Remark 3.5.1. The randomised response mechanism with $p \leq 1/2$ is $(\ln \frac{1-p}{p})$ -DP with respect to the edit neighbourhood N_e .

Proof. Consider $x = (x_1, \dots, x_j, \dots, x_n)$, $x' = (x_1, \dots, x'_j, \dots, x_n)$. Then

$$\begin{aligned} \pi(a \mid x) &= \prod_i \pi(a_i \mid x_i) \\ &= \pi(a_j \mid x_j) \prod_{i \neq j} \pi(a_i \mid x_i) \\ &\leq \frac{1-p}{p} \pi(a_j \mid x'_j) \prod_{i \neq j} \pi(a_i \mid x_i) \\ &= \frac{1-p}{p} \pi(a \mid x') \end{aligned}$$

$\pi(a_j = k \mid x_j = k) = 1 - p$ so the ratio is $\max\{(1-p)/p, p/(1-p)\} \leq (1-p)/p$ for $p \leq 1/2$. \square

Moving to a new neighbourhood

Is the randomised-response mechanism for it to be differentially private with respect to the insertion-neighbourhood definition? If not, is it possible to modify it in order to satisfy that privacy definition? *Hint: The mechanism must be able to hide the participation of a single individual in the database.*

Finally, it may be convenient to look at neighbourhoods in terms of a distance between datasets. Let $\mathbb{N}^{\mathcal{X}}$ be the set of all possible dataset histograms, i.e. counts of different possible rows in each dataset. Then the distance between two datasets is simply the total difference in counts:

Definition 3.5.5 (Hamming distance between datasets).

$$\|\mathbf{x} - \mathbf{x}'\|_1 = \sum_j \left| \sum_{i=1}^n \mathbb{I}\{x_i = j\} - \sum_{i=1}^{n'} \mathbb{I}\{x'_i = j\} \right| = \sum_j |N_j(\mathbf{x}) - N_j(\mathbf{x}')|, \quad (3.5.2)$$

where $N_j(\mathbf{x})$ is the number of elements in \mathbf{x} equal to j .

Let us see how the Hamming distance relates to the two neighbourhoods we defined. In particular, $\|\mathbf{x} - \mathbf{x}'\|_1 = 1$ if and only if $\mathbf{x} N \mathbf{x}'$.

Taking the first definition, if \mathbf{x}' has one row less than \mathbf{x} , but is otherwise identical, then there is exactly one value j for which $N_j(\mathbf{x}) = N_j(\mathbf{x}') + 1$, with $N_j(\mathbf{x}) = N_j(\mathbf{x}')$ otherwise. Consequently, $\|\mathbf{x} - \mathbf{x}'\|_1 = 1$.

On the other hand, for the second neighbourhood definition, things are slightly different. There, we assume that $x_i \neq x'_i$ for one i . Without loss of generality, let us say that $x_i = j$ and $x'_i = k$, with $j \neq k$. Then $N_j(\mathbf{x}) = N_j(\mathbf{x}') + 1$ and $N_k(\mathbf{x}') = N_k(\mathbf{x}) + 1$. Consequently $\|\mathbf{x} - \mathbf{x}'\|_1 = 2$.

The reverse direction follows by contradiction: if the Hamming distance is one, then the two databases must differ in at least one element. If they differ in more, then their distance has to be larger than one.

3.6 The local and central privacy models

So far, we have only seen the concept of differential privacy applied as a method to generate a privatised version of a dataset. In general, however, we want to perform some specific calculations on the data. Is it possible to compute things privately in the first place? Let us assume that you have defined a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, which you wish to compute on arbitrary data $x \in \mathcal{X}$. If you wish to preserve privacy, however, you need the computation to satisfy some formal guarantee like differential privacy. The simplest way to do this is by simply generating a dataset a with a DP algorithm and then processing the data with the function we have already specified. This corresponds to the local privacy model.

The local privacy model.

Given a function $f : \mathcal{X} \rightarrow \mathcal{Y}$, and a set of individual data $\mathbf{x} = (x_i)_{i=1}^n$, use a differentially private algorithm $\pi(\mathbf{a}|\mathbf{x})$ to generate $\mathbf{a} \in \mathcal{X}$. Then output $f(\mathbf{a})$.

The advantage of the local model is that the calculation f does not need to be modified. The individuals do not have to trust anybody, as they can modify their data locally. However, they must take care that the DP calculation is performed correctly.

Typically, in the local privacy model, the i -th individual's data x_i is used to generate a private response a_i . This means that individual data is only available to a local private mechanism. This model allows us to publish a complete dataset of private responses. In the central privacy model, the data x is collected and the result a is published by a *trusted curator*.

The central privacy model.

A trusted curator obtains the data \mathbf{x} of all individuals and selects a DP policy $\pi(\mathbf{a} | \mathbf{x})$ to generate the output a , so that $\mathbf{a} \approx f(\mathbf{x})$.

The main advantage of the centralised model is that approximating f with a differentially-private version can be much more accurate than simply using the original function with noisy data. However, obtaining such an approximation is not always easy.

The dependency diagrams in Figures 3.4 shows how the output depends on the individual data, the non-private and private algorithm.

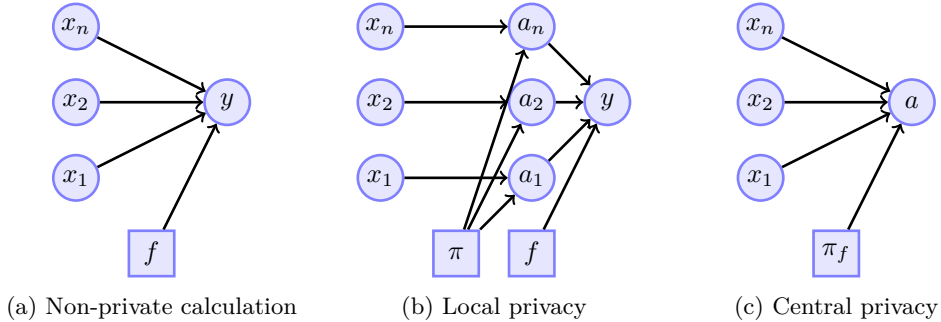


Figure 3.4: Non-privacy, local privacy, and centralised privacy models. In the first case, the value of y depends directly on the function f to be calculated, as well as the data x . In the local privacy model, noise is added to the individual data with a differentially private algorithm π , and the function is calculated on the output. In the central model, a stochastic version π_f of the original function f is calculated on the data.

Trust models.

Google uses a local privacy model to collect data from Android phone users. How does this compare to Google gathering data in the clear and then performing private computations? Who do the users have to trust? What are the privacy risks in either case?

3.6.1 Properties of differential privacy.

Remark 3.6.1. Any differentially private algorithm must be stochastic.

To prove that this is necessary, consider the example of counting how many people take drugs in a competition. If the adversary only doesn't know whether you in particular take drugs, but knows whether everybody else takes drugs, it's trivial to discover your own drug habits by looking at the total. This is because in this case, $f(x) = \sum_i x_i$ and the adversary knows x_i for all $i \neq j$. Then, by observing $f(x)$, he can recover $x_j = f(x) - \sum_{i \neq j} x_i$. Consequently, it is not possible to protect against adversaries with arbitrary side information without stochasticity.

Randomness is also necessary in cryptography. In that setting, Alice wants to communicate a message x to Bob. In a secret key cryptographic scheme, information theoretic security is achieved by making sure that the encrypted message a comes from a uniform distribution $\pi(a|x)$. This is achieved by uniformly selecting a key and selecting an appropriate hash function

Depending on the DP scheme, each query answered may leak privacy. In particular, if we always respond with an ϵ -DP mechanism, after T queries our privacy guarantee is $T\epsilon$. There exist mechanisms that do not respond to each query independently, which can reduce the total privacy loss, but those are outside the scope of this chapter.

Definition 3.6.1 (T -fold composition). In this privacy model, we compose a mechanism out of T differential private mechanisms π_1, \dots, π_T . This composition is fixed ahead of time.

Theorem 3.6.1. For any $\epsilon > 0$, the class of ϵ -differentially private mechanism satisfy $T\epsilon$ -differential privacy under T -fold composition. More generally, if each mechanism is ϵ_i -DP, the composed mechanism is $\sum_{i=1}^T \epsilon_i$ -DP.

Theorem 3.6.2 (Post-processing). Let mechanism $\pi(a | x)$ be ϵ -DP. Applying any transformation $f : A \rightarrow Y$ to the output of the mechanism to obtain $y = f(a)$, results in another ϵ -DP mechanism.

The composition theorem is a very useful tool, as it allows us to create new mechanisms that are composed of simpler parts. As a first example, we can look at how we can use it to create a randomised response algorithm when the respondents provide multiple attributes.

Randomised response for multiple attributes.

Up to now we have been discussing the case where each individual only has one attribute. However, in general each individual t contributes multiple data $x_{t,i}$, which can be considered as a row \mathbf{x}_t in a database. Then the mechanism can release each $a_{t,i}$ independently.

For n users and k attributes, if the release of each attribute i is ϵ -DP then the data release is $k\epsilon$ -DP. Thus to get ϵ -DP overall, we need ϵ/k -DP per attribute.

The result follows immediately from the composition theorem. We can see each attribute release as the result of an individual query. More generally, if each attribute i is released with an ϵ_i -DP mechanism, the overall mechanism is $\sum_i \epsilon_i$ -DP.

✍ Adversary knowledge

Assume that the adversary knows that the data is either \mathbf{x} or \mathbf{x}' . For concreteness, assume the data is either

$$\mathbf{x} = (x_1, \dots, x_j = 0, \dots, x_n)$$

where x_i indicates whether or not the i -th person takes drugs, or

$$\mathbf{x}' = (x_1, \dots, x_j' = 1, \dots, x_n).$$

In other words, the adversary knows the data of all people apart from one, the j -th person. We can assume that the adversary has some prior belief

$$\beta(\mathbf{x}) = 1 - \beta(\mathbf{x}')$$

for the two cases. Assume the adversary knows the output a of a mechanism π . What can we say about the posterior distribution of the adversary $\beta(\mathbf{x} \mid a, \pi)$ after having seen the output, if π is ϵ -DP? How does it depend on ϵ ?

Solution. We can write the adversary posterior as follows.

$$\mathbb{P}_\beta^\pi(\mathbf{x} \mid a) = \frac{\pi(a \mid \mathbf{x})\beta(\mathbf{x})}{\pi(a \mid \mathbf{x})\beta(\mathbf{x}) + \pi(a \mid \mathbf{x}')\beta(\mathbf{x}')} \quad (3.6.1)$$

$$\geq \frac{\pi(a \mid \mathbf{x})\beta(\mathbf{x})}{\pi(a \mid \mathbf{x})\beta(\mathbf{x}) + \pi(a \mid \mathbf{x})e^\epsilon\beta(\mathbf{x}')} \quad (\text{from DP definition})$$

$$= \frac{\beta(\mathbf{x})}{\beta(\mathbf{x}) + e^\epsilon\beta(\mathbf{x}')} \quad (3.6.2)$$

Note that $\pi(a \mid \mathbf{x}) \leq e^\epsilon \pi(a \mid \mathbf{x}')$ and conversely $\pi(a \mid \mathbf{x}') \leq e^\epsilon \pi(a \mid \mathbf{x})$. We can also then bound the quantity from above:

$$\mathbb{P}_\beta^\pi(\mathbf{x} \mid a) \leq \frac{\beta(\mathbf{x})}{\beta(\mathbf{x}) + e^{-\epsilon}\beta(\mathbf{x}')}.$$

Consequently, $\lim_{\epsilon \rightarrow 0} \mathbb{P}_\beta^\pi(\mathbf{x} \mid a) = \beta(\mathbf{x})$, hence the information gained by the adversary is bounded by the prior and the information loss ϵ . \square

3.7 The Laplace mechanism

Many times we already have a function $f : \mathcal{X} \rightarrow \mathbb{R}$ we want to calculate on data, and we would like to make the function preserve privacy. This clearly falls within the *central privacy* model: The analyst has access to the data and function, and wishes for the algorithm generating the final output to be private. One solution, that can satisfy differential privacy, is to add noise to the function's output. We start by first calculating the value of the function for the data we have, $f(x)$, and then we add some random noise ω , hence our calculation is random:

$$a = f(x) + \omega.$$

The amount and type of noise added, together with the smoothness of the function f , determine the amount of privacy we have. One of the simplest noise-adding mechanisms is to add zero-mean Laplace noise. Then we write $\omega \sim \text{Laplace}(\lambda)$.⁴ The mechanism is defined below.

⁴When unspecified, the mean parameter is assumed to be zero.

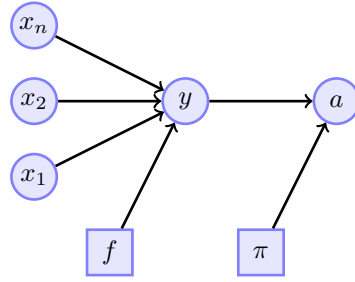


Figure 3.5: Laplace mechanism

Definition 3.7.1 (The Laplace mechanism). For any function $f : \mathcal{X} \rightarrow \mathbb{R}$, the output a of the mechanism is sampled from a Laplace distribution with mean $f(x)$ and scaling parameter λ , i.e.

$$\pi(a | x) = \text{Laplace}(f(x), \lambda). \quad (3.7.1)$$

The probability density function of the Laplace distribution with mean μ and scaling λ is given by:

$$p(\omega | \mu, \lambda) = \frac{1}{2\lambda} \exp\left(-\frac{|\omega - \mu|}{\lambda}\right).$$

and has mean μ and variance $2\lambda^2$.

⚙️ The Laplace mechanism for averages

Here we have n individuals for which we wish to calculate the average salary.

- The i -th person receives salary x_i
- We wish to calculate the average salary in a private manner.

We can do this in two ways. By using a DP mechanism on each individual salary and then calculating the average, or by first calculating the average and then applying a DP mechanism to the result. In particular, we can add Laplace noise in both cases.

The local model. In this case, $\pi(a | x)$ is obtained by independent Laplace noise ω_i for each individual:

- Obtain $y_i = x_i + \omega_i$, where $\omega_i \sim \text{Laplace}(\lambda)$.
- Return $a = n^{-1} \sum_{i=1}^n y_i$.

The centralised model. In this case, $\pi(a | x)$ is obtained by averaging first and adding noise later.

- Calculate $y = n^{-1} \sum_{i=1}^n x_i$.
- Return $a = y + \omega$, where $\omega \sim \text{Laplace}(\lambda')$.

We must tune λ, λ' appropriately in to obtain the needed ϵ -DP guarantee.

In the centralised privacy model, the non-private calculation directly outputs $y = f(x)$. We can approximate this with a DP calculation with distribution $\pi(a | x)$. The Laplace mechanism does so by first calculating $y = f(x)$ and then generating $\pi(a | y)$ so that $\mathbb{E}_\pi[a | x] = f(x)$.

Let us now talk about how the Laplace mechanism that can be used both in the centralised and local model when $\mathcal{X} \subset \mathbb{R}^n$ and $\mathcal{Y} \subset \mathbb{R}$.

DP properties of the Laplace mechanism

To use the Laplace mechanism in the centralised setting, we must relate it to a specific function f that we wish to compute. In particular, the mechanism works by adding noise to the output of the function f in a carefully calibrated manner so that we achieve exactly ϵ -differential privacy.

Definition 3.7.2 (Sensitivity). The sensitivity of a function f is

$$\mathbb{L}(f) \triangleq \sup_{xNx'} |f(x) - f(x')|$$

If we define a metric d , so that $d(x, x') = 1$ for xNx' , then:

$$|f(x) - f(x')| \leq \mathbb{L}(f) d(x, x'),$$

i.e. f is $\mathbb{L}(f)$ -Lipschitz with respect to d .

EXAMPLE 3. If $f : \mathcal{X} \rightarrow [0, B]$, e.g. $\mathcal{X} = \mathbb{R}$ and $f(x) = \min\{B, \max\{0, x\}\}$, then $\mathbb{L}(f) = B$.

EXAMPLE 4. If $f : [0, B]^n \rightarrow [0, B]$ is $f = \frac{1}{n} \sum_{t=1}^n x_t$, then $\mathbb{L}(f) = B/n$.

Proof. Consider two neighbouring datasets x, x' differing in example j . Then

$$f(x) - f(x') = \frac{1}{n} [f(x_j) - f(x'_j)] \leq \frac{1}{n} [B - 0]$$

□

Theorem 3.7.1. The Laplace mechanism on a function f with sensitivity $\mathbb{L}(f)$, ran with $\text{Laplace}(\lambda)$ is $\mathbb{L}(f)/\lambda$ -DP. Consequently, if the Laplace mechanism is ran with $\lambda = \mathbb{L}(f)/\epsilon$, then it is ϵ -DP.

Proof.

$$\frac{\pi(a | x)}{\pi(a | x')} = \frac{e^{|a-f(x')|/\lambda}}{e^{|a-f(x)|/\lambda}} \leq \frac{e^{|a-f(x)|/\lambda + \mathbb{L}(f)/\lambda}}{e^{|a-f(x)|/\lambda}} = e^{\mathbb{L}(f)/\lambda}$$

The first step follows from the definition of the Laplace mechanism. The inequality follows from the fact that for xNx' , we have $|f(x) - f(x')| \leq \mathbb{L}(f)$. In particular,

(a) If $f(x') - a \geq 0$ then

$$|f(x') - a| = f(x') - a \leq \mathbb{L}(f) + f(x) - a \leq L + |a - f(x)|,$$

as $f(x') \leq f(x) + \mathbb{L}(f)$ from the Lipschitz property.

(b) If $f(x') - a < 0$ then

$$|f(x') - a| = a - f(x') \leq \mathbb{L}(f) - f(x) + a \leq L + |a - f(x)|,$$

as $-f(x') \leq \mathbb{L}(f) - f(x)$.

Replacing into the exponential gives us the required inequality. The final result is obtained with elementary algebra. □

⚙️ DP properties of the average in the local and central model.

What is the effect of applying the Laplace mechanism in the local versus centralised model? Let us continue the average example. Here let us assume $x_i \in [0, B]$ for all i .

The Laplace mechanism in the local privacy model

The sensitivity of the individual data is B , so to obtain ϵ -DP we need to use $\lambda = B/\epsilon$. The variance of each component is $2(B/\epsilon)^2$, so the total variance is $2B^2/\epsilon^2 n$.

The Laplace mechanism in the centralised privacy model

The sensitivity of f is B/n , so we only need to use $\lambda = \frac{B}{n\epsilon}$. The variance of a is $2(B/\epsilon n)^2$.

Thus the two models have a significant difference in the variance of the estimates obtained, for the same amount of privacy. While the central mechanism has variance $O(n^{-2})$, the local one is $O(n^{-1})$ and so our estimates will need much more data to be accurate under this mechanism. In particular, we need square the amount of data in the local model as we need in the central model. Nevertheless, the local model may be the only possible route if we have no specific use for the data.

3.8 Interactive data access.

In a lot of applications, we cannot pre-define the computation that we want to perform. After we have collected the data, we need to perform an adaptive data analysis. This requires performing a sequence of computations on the data, where the result of one computation determines what computation we will do next.

We can think of this as performing queries to the database. You may be familiar with database access languages such as SQL, where you ask a question such as “what is the sum of the attribute `age` in table `students`?” and obtain the exact sum back. It is possible to answer such queries through a differentially private mechanism. However, the more queries you answer, the more you reveal about the original data. In addition, since an adversary may be cleverly adjusting the queries to more efficiently discover a secret, we need the concept of adaptive composition.

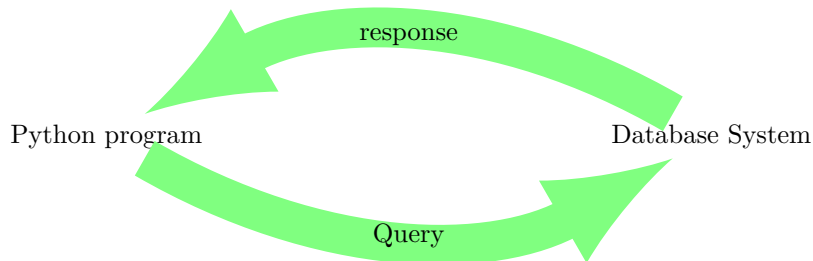


Figure 3.6: Database access model

3.8.1 Utility of queries

Rather than saying that we wish to calculate a private version of some specific function f , sometimes it is more useful to consider the problem from the perspective of the utility of different answers to queries. More precisely, imagine the interaction between a database system and a user:

Interactive queries

- System has data x .
- At time t , user asks query $q_t = q$.
- System responds with $a_t = a$.
- There is a common utility function $U : \mathcal{X}, \mathcal{A}, \mathcal{Q} \rightarrow \mathbb{R}$.

We wish to give the most useful answer, by return a that maximises U , i.e. $a = \arg \max_{a'} U(x, a', q)$, but are constrained by the fact that we also want to preserve privacy.

The utility $U(x, a, q)$ describes how appropriate each answer a given by the system for a query q is given the data x . It can be seen as how useful the response is⁵ It allows us to quantify exactly how much we would gain by replying correctly. The exponential mechanism, described below is a simple differentially private mechanism for responding to queries while trying to maximise utility for *any possible* utility function.

The Exponential Mechanism.

Here we assume that we can answer queries q , whereby each possible answer a to the query has a different utility to the DM: $U(q, a, x)$. The idea is that the best answer to the query should have the highest utility. The further away from the correct answer the response is, the lower the utility should be. This allows us to quantify how much we value correct answers to a query.

As an example, if the optimal response to a query is given by a real-valued function $f(q, x)$, then one possible utility function is $U(q, a, x) = -[f(q, x) - a]^2$. This results in the correct response having a utility of zero, and responses whose value is further away will have negative values.

The exponential mechanism allows us to answer queries in a “soft” manner, by using this utility function. It assigns lower probability to answers with lower utility, hence the better responses have a higher likelihood of being selected.

Definition 3.8.1 (The Exponential mechanism). For any utility function $U : \mathcal{Q} \times \mathcal{A} \times \mathcal{X} \rightarrow \mathbb{R}$, define the policy, which implicitly depends on U and ϵ , as:

$$\pi(a \mid x, q) \triangleq \frac{e^{\epsilon U(q, a, x) / 2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x) / 2\mathbb{L}(U(q))}}, \quad (3.8.1)$$

where $\mathbb{L}(U(q)) \triangleq \max_a \sup_{x, N, x'} |U(q, a, x) - U(q, a, x')|$ denotes the sensitivity of a query.

Clearly, when $\epsilon \rightarrow 0$, this mechanism is uniformly random. When $\epsilon \rightarrow \infty$ the action maximising $U(q, a, x)$ is always chosen (see Exercise 4). Although the exponential mechanism can be used to describe some known DP mechanisms (see Exercise 5), its best use is in settings where there is a natural utility function.

⁵This is essentially the utility to the user that asks the query, but it could be the utility to the person that answers. In either case, the motivation does not matter the action should maximise it, but is constrained by privacy.

Theorem 3.8.1. *The exponential mechanism is ϵ -differentially private.*

Proof. We only need to look at the ratio between two different distributions of answers to queries. Then we have:

$$\begin{aligned}
\frac{\pi(a \mid x, q)}{\pi(a \mid x', q)} &= \frac{e^{\epsilon U(q, a, x)/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/2\mathbb{L}(U(q))}}{e^{\epsilon U(q, a, x')/2\mathbb{L}(U(q))}} \\
&= e^{\epsilon U(q, a, x)/2\mathbb{L}(U(q)) - \epsilon U(q, a, x')/2\mathbb{L}(U(q))} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \\
&= e^{\epsilon [U(q, a, x) - U(q, a, x')]/2\mathbb{L}(U(q))} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \\
&\leq e^{\epsilon/2} \times \frac{\sum_{a'} e^{\epsilon U(q, a', x')/2\mathbb{L}(U(q))}}{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \\
&\leq e^{\epsilon/2} \times \frac{\sum_{a'} e^{\epsilon [(U(q, a', x) + \mathbb{L}(U(q)))/2\mathbb{L}(U(q))]} }{\sum_{a'} e^{\epsilon U(q, a', x)/2\mathbb{L}(U(q))}} \leq e^{\epsilon}
\end{aligned}$$

□

Interactive differential privacy.

So far we only defined differential privacy in terms of a fixed mechanism. However, in general the mechanism must respond to sequential queries q_1, \dots, q_t , with answers a_1, \dots, a_t . The answer can depend arbitrarily on the sequence of responses and queries. So, privacy must hold no matter what the previous sequence of questions and answers. This gives rise to the following definition.

Definition 3.8.2 (Interactive DP). An algorithm π satisfies interactive $\{\epsilon_1, \dots, \epsilon_t\}$ differential privacy with respect to N if

$$\ln \frac{\pi(a_t \mid x, q_1, \dots, q_t, a_1, \dots, a_{t-1})}{\pi(a_t \mid x', q_1, \dots, q_t, a_1, \dots, a_{t-1})} \leq \epsilon_t, \quad \forall x N x'. \quad (3.8.2)$$

Since the queries are not defined in advance, this means that q_t may also depend on all the previous answers a_1, \dots, a_{t-1} .

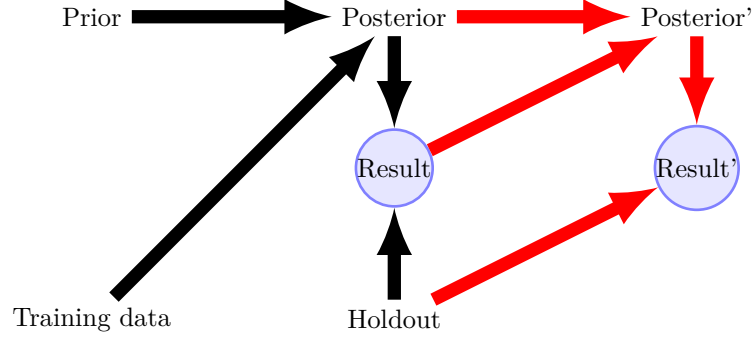
Answering queries independently in the interactive setting is sufficient for achieving differential privacy ϵ_t at each step t . The total privacy loss is at most $\sum_{i=1}^t \epsilon_i$. However, it is possible to bound the privacy loss more tightly with advanced composition theorems.

3.9 Advanced topics.

In this section, we give a brief overview of some more advanced topics on privacy. The interested reader is urged to look into the given references.

3.9.1 Relaxations of differential privacy.

In practice, satisfying differential privacy with a small ϵ results in mechanisms with low utility. It is possible to relax the definition of differential privacy to include an additive term, as shown below, which greatly improves the privacy-utility trade-off we can achieve.



Definition 3.9.1. Approximate differential privacy. A stochastic algorithm $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{A})$, where \mathcal{X} is endowed with a neighbourhood relation N , is said to be (ϵ, δ) differentially private if

$$\pi(A \mid x) \leq \pi(A \mid x')\epsilon + \delta, \quad \forall xNx', \quad A \subset \mathcal{A}. \quad (3.9.1)$$

One simple interpretation of this definition is that the mechanism is $(\epsilon, 0)$ -DP with probability $1 - \delta$. Conversely, this implies that with a very small probability δ , the complete dataset might be revealed. For that reason, the parameter δ must be as small as possible.

Approximate differential privacy is satisfied by the *Gaussian* mechanism. This is similar to the Laplace mechanism, but scales noise to the ℓ_2 sensitivity of the function.

⚙️ The Gaussian mechanism

Given a function $f : \mathcal{X} \rightarrow \mathbb{R}^k$, the Gaussian mechanism with parameter σ outputs $a \sim \mathcal{N}(f(x), \sigma)$.

Let ℓ_2 sensitivity of f be

$$\mathbb{L}(f)_2 = \max \{ \|f(x) - f(y)\|_2 \mid xNy \}.$$

For $\sigma \geq 2 \ln(1.25/\delta) \frac{\Delta_2(f)}{\epsilon}$, then the Gaussian mechanism is (ϵ, δ) -differentially private.

⚖️ Renyi differential privacy.

The careful reader will have noticed that differential privacy is essentially a bound on the distributional distance of mechanism outputs for similar inputs. In particular, the distance used for ϵ -DP is $D_\infty(P \parallel Q) \triangleq \sup_A |\ln P(A)/Q(A)|$. We can use any other divergence between distributions, such as the KL divergence, as long as this becomes unbounded for distributions with unequal support.

Reproducibility

Training and testing, overfitting on the test set.

The unfortunate practice of adaptive analysis

In the ideal data analysis, we start from some prior hypothesis, then obtain some data, which we split into training and holdout. We then examine the training data and obtain a posterior

that corresponds to our conclusions. We can then measure the quality of these conclusions in the independent holdout set.

However, this is not what happens in general. Analysts typically use the same holdout repeatedly, in order to improve the performance of their algorithms. This can be seen as indirectly using the holdout data to obtain a new posterior, and so it is possible that you can overfit on the holdout data, even if you never directly see it. It turns out we can solve this problem if we use differential privacy, so that the analyst only sees a differentially private version of queries.

The reusable holdout¹⁰⁶

One idea to solve this problem is to only allow the analyst to see a private version of the result. In particular, the analyst will only see whether or not the holdout result is τ -close to the training result.

Algorithm parameters

- Performance measure f .
- Threshold τ . How close do we want f to be on the training versus holdout set?
- Noise σ . How much noise should we add?
- Budget B . How much are we allowed to learn about the holdout set?

Algorithm idea

Run algorithm λ on data D_T and get e.g. classifier parameters θ .
Run a DP version of the function $f(\theta, D_H) = \mathbb{I}\{U(\theta, D_T) \geq \tau U(\theta, D_H)\}$.

So instead of reporting the holdout performance at all, you just see if you are much worse than the training performance, i.e. if you're overfitting. The fact that the mechanism is DP also makes it difficult to learn the holdout set. See the thresholdout link for more details.

3.10 Discussion

The definition of differential privacy

- First rigorous mathematical definition of privacy.
- Relaxations and generalisations possible.
- Connection to learning theory and reproducibility.

Current uses

- Apple. DP is used internally in the company to “protect user privacy”. It is not clear exactly what they are doing but their efforts seem to be going in the right direction.
- Google. The company has a DP API available based on randomised response, RAPPOR.

¹⁰⁶Also see <https://ai.googleblog.com/2015/08/the-reusable-holdout-preserving.html>

- Uber. Elastic sensitivity for SQL queries, which is available as open source. This is a good thing, because it is easy to get things wrong with privacy.
- US 2020 Census. It uses differential privacy to protect the confidentiality of responders' information while maintaining data that are suitable for their intended uses.

Open problems

- Complexity of differential privacy.
- Verification of implementations and queries.

Available privacy toolboxes

Differential privacy

- Open DP <<https://opendp.org/>>
- <https://github.com/bmcmenamin/thresholdOut-explorations> Threshold out
- <https://github.com/steven7woo/Accuracy-First-Differential-Privacy> Accuracy-constrained DP
- <https://github.com/menisadi/pydp> Various DP algorithms
- <https://github.com/haiphanNJIT/PrivateDeepLearning> Deep learning and DP

k -anonymity

- <https://github.com/qiyuangong/Mondrian> Mondrian k -anonymity

Learning outcomes

Understanding

- Linkage attacks and k -anonymity.
- Inferring data from summary statistics.
- The local versus centralised differential privacy model.
- False discovery rates.

Skills

- Make a dataset satisfy k -anonymity with respect to identifying attributes.
- Apply the randomised response and Laplace mechanism to data.
- Apply the exponential mechanism to simple decision problems.

- Use differential privacy to improve reproducibility.

Reflection

- How can potentially identifying attributes be chosen to achieve k -anonymity?
- How should the parameters of the two ideas, ϵ -DP and k -anonymity be chosen?
- Does having more data available make it easier to achieve privacy?

Further reading

- k -anonymity¹⁷
- Randomness, privacy, and the US 2020 census¹³
- The paper introducing differential privacy⁹
- Differential privacy book⁸
- Bayesian inference and privacy⁶
- Local differential privacy and statistics⁷
- Local differential privacy and applications²⁰
- The exponential mechanism¹⁶

3.11 Exercises

EXERCISE 1. Show that the randomised response mechanism, as originally defined, is not differentially private with respect to the addition/deletion neighbourhood.

Proof. Let $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$ be three datasets with $\mathbf{x}_1 N \mathbf{x}_2 N \mathbf{x}_3$ and $\mathbf{x}_1 N_e \mathbf{x}_3$. First of all, an algorithm that is ϵ -DP with respect to the N , is 2ϵ -DP with respect to N_e . This follows from the fact that $\pi(a|\mathbf{x}_1) \leq e^\epsilon \pi(a|\mathbf{x}_2) \leq e^{2\epsilon} \pi(a|\mathbf{x}_3)$. What about the converse? Consider the case where \mathbf{x}_1 has n entries and \mathbf{x}_2 has $n+1$ entries. Let A_n be the set of all responses with n entries. Then clearly $\pi(A_n|\mathbf{x}_1) = 1$ and $\pi(A_n|\mathbf{x}_2) = 0$. Consequently $|\ln \frac{\pi(A_n|\mathbf{x}_1)}{\pi(A_n|\mathbf{x}_2)}|$ is not bounded by any constant, and differential privacy is not satisfied with respect to insertions and deletions. \square

EXERCISE 2. Define a variant of the binary randomised response mechanism that satisfies differential privacy with respect to the insertion and deletion neighbourhood. More precisely, it should hold that for any dataset pair $\mathbf{x} = (x_1, \dots, x_n)$, and $\mathbf{x}' = (x_1, \dots, x_n)$, $\mathbf{x} = (x_1, \dots, x_n, x_{n+1})$, the mechanism satisfies

$$\pi(A|\mathbf{x}) \leq e^\epsilon \pi(A|\mathbf{x}'), \quad \forall A \subset \mathcal{A}.$$

Proof. First of all, for the mechanism to satisfy ϵ -DP, it must be the case that an output size must have a non-zero probability: If with an n -size input we can only generate $\leq n$ -size outputs, then we cannot generate a $n+1$ -size output. Consequently, the randomised response mechanism defined for N -neighbourhoods makes sense only when we know the maximum possible value of

n . This is normally the case when we are surveying a population whose size is assumed to be common knowledge.

In particular, consider a mechanism which samples person i with probability q so that $s_i = 1$ if a person is selected. If a person is selected, then a_i is drawn from a standard randomised response mechanism (which flips x_i with probability p), otherwise $a_i = \perp$. Then, we can factorise the probability of a specific output as follows:

$$\pi(a_i = k | x_i = k) = \pi(a_i = k | s_i = 1, x_i = k)$$

□

EXERCISE 3. The mean estimator $\hat{\theta}$ for the randomised response mechanism with flipping probability p , where the true data generation process is Bernoulli with parameter θ , is defined as

$$\hat{\theta}(a) = \frac{\bar{a} - p}{1 - 2p},$$

where $\bar{a} = \frac{1}{T} \sum_{t=1}^T a_t$ is the mean of the observed answers. Prove that this estimator is unbiased, i.e. that

$$\mathbb{E}[\hat{\theta}] = \theta,$$

where the expectation is taken over the unobserved data randomness and the randomness of the mechanism.

EXERCISE 4. Prove that the exponential mechanism is uniform when $\epsilon \rightarrow 0$ and deterministically returns the maximising action when $\epsilon \rightarrow \infty$.

EXERCISE 5. Prove that the exponential mechanism, when we used to calculate a noisy version of the function $q(x)$ with $q : \mathcal{X} \rightarrow \mathbb{R}$ and utility $U(a, q, x) = -|q(x) - a|^p$, results in the Laplace mechanism for $p = 1$ and the Gaussian mechanism for $p = 2$.

EXERCISE 6. Consider the following relaxation of differential privacy to KL divergences. A mechanism is ϵ -KL-private if

$$\sum_a \ln \frac{\pi(a|x)}{\pi(a|x')} \pi(a|x) \leq \epsilon$$

Prove that two-folded composition of such a mechanism is 2ϵ -KL-private.

Proof. For simplicity, we compose the same mechanism twice over a finite alphabet. Then $\pi(a, a'|x) = \pi(a|x)\pi(a'|x)$ because the mechanism is not interactive.

$$\begin{aligned} \sum_{a, a'} \ln \frac{\pi(a, a'|x)}{\pi(a, a'|x')} \pi(a, a'|x) &= \sum_{a, a'} \ln \frac{\pi(a|x)\pi(a'|x)}{\pi(a|x')\pi(a'|x')} \pi(a|x)\pi(a'|x) \\ &= \sum_{a, a'} \ln \frac{\pi(a|x)}{\pi(a|x')} \pi(a|x)\pi(a'|x) + \ln \frac{\pi(a'|x)}{\pi(a'|x')} \pi(a|x)\pi(a'|x) \\ &= \sum_{a'} \pi(a'|x) \sum_a \ln \frac{\pi(a|x)}{\pi(a|x')} \pi(a|x) + \sum_a \pi(a|x) \sum_{a'} \ln \frac{\pi(a'|x)}{\pi(a'|x')} \pi(a'|x) \\ &\leq \sum_{a'} \pi(a'|x)\epsilon + \sum_a \pi(a|x)\epsilon = 2\epsilon. \end{aligned}$$

□

Bibliography

- [1] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2/3):235–256, 2002.
- [2] Craig M. Bennett, George L. Wolford, and Michael B. Miller. The principled control of false positives in neuroimaging. *Social cognitive and affective neuroscience*, 4 4:417–22, 2009. URL <https://pdfs.semanticscholar.org/19c3/d8b67564d0e287a43b1e7e0f496eb1e8a945.pdf>.
- [3] Craig M Bennett, Abigail A Baird, Michael B Miller, and George L Wolford. Journal of serendipitous and unexpected results. *Journal of Serendipitous and Unexpected Results (jsur. org)-Vol*, 1(1):1–5, 2012. URL <https://teenspecies.github.io/pdfs/NeuralCorrelates.pdf>.
- [4] Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996. URL citeseer.nj.nec.com/breiman96bagging.html.
- [5] Philip Dawid. The decision-theoretic approach to causal inference. *Causality: Statistical perspectives and applications*, pages 25–42, 2012. URL <https://arxiv.org/pdf/1405.2292.pdf>.
- [6] Christos Dimitrakakis, Blaine Nelson, Zuhe Zhang, Aikaterini Mitrokotsa, and Benjamin I. P. Rubinstein. Differential privacy for bayesian inference through posterior sampling. *Journal of Machine Learning Research*, 18(1):343–381, 2017.
- [7] John C Duchi, Michael I Jordan, and Martin J Wainwright. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 429–438. IEEE, 2013.
- [8] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [9] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [10] Cynthia Dwork, Vitaly Feldman, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Aaron Roth. The reusable holdout: Preserving validity in adaptive data analysis. *Science*, 349 (6248):636–638, 2015.
- [11] Dean Eckles and Maurits Kaptein. Thompson sampling with the online bootstrap. *arXiv preprint arXiv:1410.4009*, 2014.

- [12] Evelyn Fix and Joseph L Hodges Jr. Discriminatory analysis-nonparametric discrimination: consistency properties. Technical report, California Univ Berkeley, 1951.
- [13] Simson L. Garfinkel and Philip Leclerc. Randomness concerns when deploying differential privacy, 2020.
- [14] Jason Hartford, Greg Lewis, Kevin Leyton-Brown, and Matt Taddy. Counterfactual prediction with deep instrumental variables networks. Technical Report 1612.09596, arXiv, 2016.
- [15] PN Lee. Relation between exposure to asbestos and smoking jointly and the risk of lung cancer. *Occupational and environmental medicine*, 58(3):145–153, 2001.
- [16] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE, 2007.
- [17] Pierangela Samarati and Latanya Sweeney. Protecting privacy when disclosing information: k-anonymity and its enforcement through generalization and suppression. Technical report, technical report, SRI International, 1998.
- [18] W.R. Thompson. On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of two Samples. *Biometrika*, 25(3-4):285–294, 1933.
- [19] Stanley L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309):63–69, 1965. ISSN 01621459. URL <http://www.jstor.org/stable/2283137>.
- [20] Mengmeng Yang, Lingjuan Lyu, Jun Zhao, Tianqing Zhu, and Kwok-Yan Lam. Local differential privacy and its applications: A comprehensive survey. *arXiv preprint arXiv:2008.03686*, 2020.