



La privacidad diferencial cómo solución a la amenaza de la privacidad en los datos biológicos

Albert Márquez Navarro-Soto
Grado en ingeniería informática
Área de seguridad

Consultor: Carlos Bonet Rapell
Profesor: Jorge Miguel Moneo

28 de diciembre de 2021



Esta obra está sujeta a una licencia de
Reconocimiento-NoComercial-SinObraDerivada
[3.0 España de Creative Commons](https://creativecommons.org/licenses/by-nc-nd/3.0/es/)

FICHA DEL TRABAJO FINAL

Título del trabajo:	<i>La privacidad diferencial cómo solución a la amenaza de la privacidad en los datos biológicos</i>
Nombre del autor:	<i>Albert Márquez Navarro-Soto</i>
Nombre del consultor/a:	<i>Carlos Bonet Rapell</i>
Nombre del PRA:	<i>Jorge Miguel Moneo</i>
Fecha de entrega (mm/aaaa):	<i>12/2021</i>
Titulación:	<i>Ingeniería informática</i>
Área del Trabajo Final:	<i>Seguridad</i>
Idioma del trabajo:	<i>Español</i>
Palabras clave:	<i>Ruido, Gauss, Laplace, aleatorio, CSV, seguridad, criptografía, privacidad, diferencial, criptografía, salud, paciente, usuario</i>
Resumen del Trabajo (máximo 250 palabras):	
<p><i>Debido al creciente uso de dispositivos inteligentes los cuales miden las constantes vitales del usuario, este trabajo de fin de grado propone un análisis sobre la situación de la privacidad diferencial, la cual tiene relación directa con el funcionamiento de dichos dispositivos.</i></p> <p><i>Inicialmente se hace una presentación de la privacidad diferencial, sus tipos así cómo compañías conocidas que hacen uso de la misma. Una vez realizado el marco teórico y analizada la información obtenida, se ha procedido a crear un apartado práctico, el cual se ha basado en añadir ruido a un conjunto de datos.</i></p> <p><i>Se han aplicado 2 tipos de adiciones de ruido, la aleatoria, mediante la función Jitter, y la Gaussiana, a través de la distribución de probabilidad de Gauss. El resultado obtenido ha sido favorable en la aplicación de la distribución de probabilidad Gaussiana, pero de falta de utilidad para la aleatoria dado que no se podían recuperar los datos.</i></p> <p><i>A continuación se ha desarrollado un modelo básico de sistema de la información, el cual se basa en un modelo ya existente y certificado por el Gobierno Canario, pero modificándolo y adaptándolo en base a las necesidades de este documento y enfoque de este trabajo de fin de grado.</i></p> <p><i>Finalmente, después de obtener todos los resultados, tanto teóricos cómo prácticos y habiendo creado un sistema de información básico, se han desarrollado las conclusiones asociadas a estos 3 apartados y se ha hecho una reflexión sobre la situación</i></p>	

Abstract (in English, 250 words or less):

Due to the increase of use of smart devices which measure the biological data of the user, this final degree project proposes an analysis of the situation of differential privacy, which is directly related to the operation and behaviour of said devices.

Initially, a presentation about differential privacy is made, its types and also well-known companies that makes use of it.

Once the theoretical part had been completed and the information obtained analyzed, a practical section was created, which was based on adding noise to a data set.

Two types of noise additions have been applied, the random one (by using the Jitter function) and the Gaussian (by using the Gaussian probability distribution). The result obtained has been positive in the application of the Gaussian probability distribution, but of lack of utility for the random one, because the data wasn't possible to recover.

A basic IS model has been developed below, which is based on an existing model certified by the Canary Islands Government, but modified and adapted based on the needs of this document and the focus of this final degree project.

Finally, after obtaining all the results, both theoretical and practical and having created a basic IS, the conclusions associated with these 3 sections have been developed and a reflection has been made on the situation.

Cita

”Si crees que la tecnología puede solucionar tus problemas de seguridad, entonces no entiendes los problemas y tampoco la tecnología”

Bruce Schneier

Agradecimientos

Gracias a todas las personas que han creído en mi en este largo camino

Índice

1. Introducción	9
1.1. Contexto y justificación del Trabajo	9
1.2. Objetivos del Trabajo	9
1.3. Enfoque y método seguido	10
1.4. Planificación del Trabajo	11
1.4.1. Planificación de entregas	11
1.4.2. Recursos	12
1.4.3. Requisitos	12
1.5. Breve resumen de productos obtenidos	12
1.6. Breve descripción de los otros capítulos de la memoria	12
2. Parte teórica	13
2.1. Definición de privacidad diferencial	13
2.2. Relación con la criptografía	15
2.2.1. Computación multiparte segura	16
2.2.2. Privacidad diferencial	17
2.3. Tipos de privacidad diferencial	18
2.4. Compañías conocidas	20
2.4.1. Apple	20
2.4.2. Microsoft	21
2.4.3. Google	22
2.4.4. Breve comparativa entre compañías y ejemplo	22
2.5. Adición de ruido	24
2.5.1. Distribución Gauss	24
2.5.2. Distribución Laplace	25
2.5.3. Comparativa entre distribuciones	26
3. Parte práctica	28
3.1. Creación de un archivo CSV en Numbers	28
3.2. Análisis de la información en RStudio	32
3.2.1. Adición de ruido aleatorio	37
3.2.2. Adición de ruido Gaussiano	40
4. Parte aplicación SI	51
4.1. Fase 1 - Protección de los datos	54
4.2. Fase 2 - Normativa y registro de actividades de tratamiento	55
4.3. Fase 3 - Evaluación de riesgos y medidas de seguridad	56
4.4. Fase 4 - Medidas técnicas y de seguridad	59
4.5. Fase 5 - Denuncias	59

5. Conclusiones	60
5.1. Apartado teórico	60
5.2. Apartado práctico	60
5.3. Apartado SI	61
5.4. Conclusión global	61
5.5. Lineas de trabajo futuras	62
6. Glosario	63
7. Bibliografía	64
8. Anexos	67

Índice de figuras

1.	Privacidad vs utilidad	13
2.	Computación multiparte	16
3.	MPC - Machine Learning	17
4.	Privacidad diferencial	18
5.	Privacidad diferencial (Laplace)	19
6.	Apple - Marketing iPhone	20
7.	Privacidad global vs local	23
8.	Distribución de Gauss	24
9.	Distribución de probabilidad de Gauss	25
10.	Distribución de Laplace	26
11.	Gauss vs Laplace (distribuciones - I)	26
12.	Gauss vs Laplace (distribuciones - II)	27
13.	Numbers (MacOS)	29
14.	Defunciones por motivos respiratorios	31
15.	CSV final	32
16.	RStudio - Comando lectura datos	33
17.	RStudio - Resultado carga CSV	34
18.	Rstudio - Información de las tablas	35
19.	RStudio - Detección gráfica de pacientes inmunes	36
20.	RStudio - Función Jitter	37
21.	RStudio - Aplicación Jitter	38
22.	RStudio - Representación gráfica Jitter	38
23.	RStudio - Filtro Inmunidad (Jitter usado)	39
24.	RStudio - Creación secuencia	40
25.	RStudio - Secuencia numérica (Gráfico)	41
26.	RStudio - Comando dnorm	42
27.	RStudio - Distribución Gaussiana	42
28.	RStudio - Asignación a variable de pacientes inmunes	44
29.	RStudio - Pacientes inmunes (nueva variable)	44
30.	RStudio - Asignación nueva secuencia numérica	44
31.	RStudio - Parámetro inmunidad modificado	45
32.	RStudio - Parámetro inmunidad -1 y +1	45
33.	RStudio - Intercambio valores inmunidad (1)	45
34.	RStudio - Intercambio valores inmunidad (2)	46
35.	RStudio - Creación Columna Número inmunes	47
36.	RStudio - Instalación paquete encryptr	47
37.	RStudio - Encriptación archivo CSV	48
38.	RStudio - Desencriptación archivo CSV	48
39.	RStudio - Proceso encriptación y desencriptación	49
40.	Comparativa archivos	50
41.	Planteamiento	52
42.	Mapa de riesgo	57

43. Funcionamiento EIPD 58

1. Introducción

1.1. Contexto y justificación del Trabajo

Actualmente la sociedad gestiona tanto su vida personal como profesional mediante infinidad de dispositivos, ya sean ordenadores, portátiles, tabletas, teléfonos móviles, relojes inteligentes, etc. En los últimos años se ha hecho especial hincapié en los datos biológicos.

A lo largo del grado, se han ofrecido diferentes asignaturas dónde se utilizan conjuntos de datos para obtener resúmenes de información. Es en la asignatura de minería de datos dónde realmente se hace más hincapié en la relevancia de la información y de su tratamiento, ya que se pueden descubrir patrones y metodologías de descubrimiento de datos.

Un ejemplo es el que se encontró Latanya Sweeney (investigador de la universidad de Harvard) hace 5 años. Mediante la información recogida de los diarios, pudo identificar hasta un 43 % de los pacientes de la base de datos del estado de Washington D.C, Estados Unidos de América. El problema no residía en algún hospital en concreto, sino en cómo se gestionaban los datos.

1.2. Objetivos del Trabajo

Los objetivos en este proyecto han sido:

1. Establecer cómo mantener los datos visibles y utilizables únicamente por el personal sanitario y el paciente
2. Definir un procedimiento para que, en casos como el anteriormente citados, no se repitan
3. Definir un método de comunicación entre el personal sanitario y el paciente sin que otras partes puedan usar esos datos tanto de forma directa como indirecta
4. Hacer una comparativa de la vulnerabilidad de los datos usando privacidad diferencial y sin ella
5. Hacer un estudio acorde referente a cómo la privacidad diferencial se puede aplicar en otros campos además del sanitario y cómo puede afectar

El objetivo general consiste en diseñar un modelo en que la información se quede en el círculo del paciente y entidad sanitaria, sin que los intermediarios o agentes externos puedan sacar provecho de dicha información así como identificar los beneficios de la privacidad diferencial no solamente en el ámbito sanitario.

1.3. Enfoque y método seguido

La estrategia y metodología seguidas para el desarrollo de este documento han sido las siguientes:

1. **Definición de un marco teórico de trabajo:** mediante la búsqueda de documentación científica acreditada, entrevistas y noticias, se analizado y probado el funcionamiento de la privacidad diferencial.

Los pilares sobre los que se cimienta este marco teórico han sido:

- Definición de privacidad diferencial y sus tipos
 - Diferenciación y utilidad de la privacidad diferencial local y global, así cómo su uso
 - Explicación de varias técnicas criptográficas
 - Breve explicación de la adición de ruido a datos mediante las funciones de Gauss y Laplace
 - Análisis básico de varias compañías que implementan la privacidad diferencial hoy en día
2. **Implementación de caso práctico:** al finalizar la definición del marco teórico, se ha realizado una aplicación práctica, la cual ha consistido en implementar en un archivo .CSV la adición de ruido aleatorio y posteriormente ruido mediante la función de Gauss. Se ha realizado una comparación tanto de efectividad cómo de utilidad entre ambos modelos.
 3. **Creación de un modelo SI sanitario:** en este apartado se ha realizado un modelo de sistema de la información, basado en el documento del gobierno de Canarias y Confederación Canaria de empresarios "Pasos prácticos para la implementación de un sistema de gestión en privacidad de la información", con un enfoque al sector sanitario, el cual es en el que se basa este documento. Así mismo, este proceso se ha hecho basándose en la ISO/IEC 27701.

Esta parte se ha usado de nexo entre el apartado teórico y práctico anteriores para dar coherencia y justificación de utilidad a este trabajo.

4. **Conclusiones:** al obtener tanto unos resultados teóricos en el primer apartado como prácticos en el segundo y tercer apartados, se ha procedido a exponer las conclusiones obtenidas. En este apartado se ha pretendido poner de manifiesto cómo gracias a la privacidad diferencial, los datos son mucho más seguros.

1.4. Planificación del Trabajo

La planificación de entregas, recursos y requisitos han sido:

1.4.1. Planificación de entregas

1. **PAC1 (15/09/2021 - 28/09/2021):** comunicación con el profesor dónde se ha tratado el tema del trabajo así como del enfoque del mismo. Esta comunicación se ha realizado mediante correo electrónico y la compartición de un documento de Google Drive con los puntos principales e ideas.
2. **PAC2 (29/09/2021 - 26/10/2021):** trabajo intensivo en el apartado teórico. El objetivo ha sido poder tener toda la base teórica finalizada al acabar esta PAC. Así mismo, en base a los resultados obtenidos, se ha dejado preparado el apartado práctico.
3. **PAC3 (27/10/2021 - 23/11/2021):** focalización en el apartado práctico y de creación del modelo SI. Debido a que ya se tenían los resultados del apartado teórico, se ha procedido a realizar el estudio práctico y la creación del modelo SI. La premisa de esta parte ha sido que bajo ningún concepto se finalizaría sin tener el apartado teórico correctamente enfocado.
4. **PAC4 (24/11/2021 - 28/12/2021):** preparación de la memoria y producto resultante. En esta PAC se ha realizado el trabajo de maquetación de todo el trabajo realizado hasta la fecha. Dicho proceso se ha realizado mediante LaTeX, en la página web de Overleaf.
5. **Entrega de la presentación virtual (29/12/2021 - 04/01/2022):** realización y entrega de la presentación virtual. Mediante un vídeo de una duración máxima de 15 minutos, se han explicado los puntos clave del trabajo. En dicha presentación se ha hecho una explicación del estudio realizado y los resultados tanto teóricos como prácticos obtenidos.
6. **Defensa (10/01/2022 - 14/01/2022):** periodo enfocado a dar respuesta a las preguntas del tribunal en menos de 24h.

1.4.2. Recursos

Los recursos de software utilizados han sido: Apple Pages, Overleaf (página web de edición LaTeX) y RStudio.

1.4.3. Requisitos

A nivel de software no se ha previsto ningún requisito en especial. Las consideraciones han sido principalmente:

1. Disponer de las versiones estables más actuales de los programas indicados
2. A nivel de documentación, el requisito indispensable ha sido el de tener documentación oficial acreditada

1.5. Breve resumen de productos obtenidos

En este trabajo se han obtenido principalmente 3 productos:

1. **Apartado teórico:** se ha obtenido un producto netamente académico para posteriormente proceder al apartado práctico
2. **Apartado práctico:** se ha obtenido un archivo CSV encriptado, el cual a su vez tiene implementada privacidad diferencial en el campo de inmunidad
3. **Apartado sistema de la información:** se ha obtenido un modelo adaptado al sector sanitario junto a la privacidad diferencial y teniendo en cuenta la ISO/IEC 27701

1.6. Breve descripción de los otros capítulos de la memoria

No se han previsto capítulos adicionales a los indicados con anterioridad y definidos en la plantilla proporcionada.

2. Parte teórica

En este primer apartado se tratará la introducción a la privacidad diferencial mediante su definición, indicación de sus diferentes tipos, relación con la criptografía así como las compañías más conocidas en la aplicación de este modelo matemático.

Finalmente se mostrará un ejercicio teórico, el cual permitirá acceder a la siguiente sección con la información necesaria.

2.1. Definición de privacidad diferencial

A continuación se muestra la definición formal de privacidad diferencial:

Definición de privacidad diferencial (Universidad de Harvard)

La privacidad diferencial es la definición matemática rigurosa de privacidad. En la configuración más simple, considerar un algoritmo que analiza un conjunto de datos y calcula estadísticas sobre él (como la media, la varianza, la mediana, etc. de los datos).

Se dice que un algoritmo de este tipo es diferencialmente privado si, al observar el resultado, no se puede saber si los datos de algún individuo se incluyeron en el conjunto de datos original o no.

Cuándo se hace referencia a la privacidad diferencial, el enfoque está en cuantificar dicho espacio privado de forma plenamente objetiva. En este documento se considera el uso de la privacidad diferencial cómo la adición de ruido a un contenido digital (más específicamente datos de salud) de un dispositivo informático.

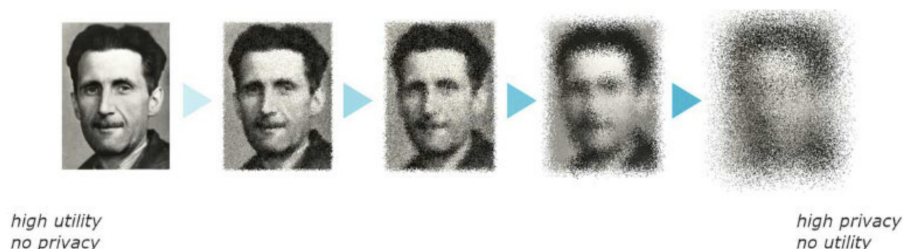


Figura 1: Privacidad vs utilidad

Un ejemplo interesante es la figura 1, dónde se pone en una balanza la utilidad de una imagen y la privacidad de la persona que sale en ella. Aunque este documento se centrará en datos numéricos (enfocados principalmente al ámbito sanitario) y no contenido audiovisual.

Esta imagen ejemplifica cuán importante es utilizar debidamente la matemática en este campo, ya que dependiendo de ello, si no se aplica correctamente se tendrán unos datos con muy poca privacidad, o bien insertables.

Dado que la privacidad diferencial se fundamenta en algoritmos matemáticos, ello conlleva varios aspectos a tratar:

1. **Tiempo de procesamiento:** el hecho de que se utilice uno o varios algoritmos matemáticos en beneficio de la privacidad de un contenido determinado, conlleva a un aumento del tiempo de procesamiento. Esto es especialmente importante en el transporte de paquetes en redes.
2. **Nivel de adición de ruido:** aún y teniendo unos tiempos de procesamiento correctos, hay que tener en cuenta el nivel de ruido que se añade a los datos, para que estos no sean completamente insertables o bien queden expuestos completamente.
3. **Almacenamiento:** el hecho de añadir ruido a un contenido, implicará también más espacio requerido. Dependiendo de dónde se almacenen dichos datos y el tipo de contenido de los mismos, también conllevará a un incremento significativo del espacio solicitado.
4. **Legislación:** el hecho de que las compañías tecnológicas implementen privacidad diferencial (y lo que ello conlleva), resulta en que éstas sean obligadas por ley a proteger los datos a un mayor nivel, lo que conlleva a una mayor responsabilidad con los clientes.
5. **Coste:** en pequeñas cantidades de datos no se apreciará, pero cuándo hablamos de miles o millones de usuarios, el coste de almacenamiento y energía aumenta exponencialmente.

El objetivo de este documento no es el desarrollo matemático de los algoritmos de la privacidad diferencial ni la explicación profunda de los mismos, no obstante, a continuación se mostrará la definición formal que tiene la privacidad diferencial:

Definición matemática general de privacidad diferencial (Silveri Fu)

$$Pr[K(D_1) \in S] \leq e^\epsilon Pr[K(D_2) \in S] \quad (1)$$

En la expresión matemática (1) se puede observar que una función aleatoria K proporciona un valor ϵ de privacidad diferencial si para dos conjuntos de datos (D_1 y D_2) difieren como máximo en un elemento para cualquier conjunto de salidas S pertenecientes al rango de K .

Así pues, K retorna el mismo resultado para D_1 y D_2 con una probabilidad muy semejante (debido a la explicación anterior). A continuación se define la sensibilidad de la función f . Tal y como se ha indicado antes, la variación en ambos conjuntos de datos, será máximo de 1 elemento, por lo tanto, se calculará el valor máximo de la diferencia de los 2 conjuntos

Definición matemática de la sensibilidad de la función f (Silveri Fu)

$$\Delta_f = \max \| f(D_1) - f(D_2) \| \quad (2)$$

En la expresión matemática (2) se puede observar cómo se define la variación entre los dos conjuntos de datos mediante la función de resta.

2.2. Relación con la criptografía

Dadas las definiciones matemáticas anteriores, éstas pueden hacer surgir la duda sobre qué diferencia/s hay entre la criptografía misma y la privacidad diferencial, dado que el enfoque de ambas es el de proteger los datos. La respuesta a ello es directa, la criptografía es la base sobre la cuál se fundamenta la privacidad diferencial, ya que es en sí misma una técnica criptográfica.

Aunque técnicas criptográficas como la computación multiparte segura o la privacidad diferencial existen desde hace bastantes años (principalmente debido a que surgen de conceptos matemáticos aplicados a la tecnología), la realidad es que debido al coste computacional que representan, no ha sido hasta ahora que se han podido implementar con garantías.

El motivo principal de su implementación es debido a que la potencia y capacidad de procesamiento ha aumentado drásticamente gracias a tecnologías como el multinúcleo y el multihilo así como la reducción del tamaño

de los procesadores, lo que permite poder implementar más núcleos e hilos en los procesadores.

De las 2 técnicas criptográficas, este documento se centrará en el de la privacidad diferencial. No obstante, a continuación se ofrece una explicación sobre ellas y unas imágenes que permiten entender de una forma más clara el concepto de ambas técnicas criptográficas.

2.2.1. Computación multiparte segura

MPC () es una técnica criptográfica, la cual permite a un grupo de personas o dispositivos electrónicos acordar una función para computar de forma conjunta. Para ello, acuerdan utilizar un mismo protocolo MPC.

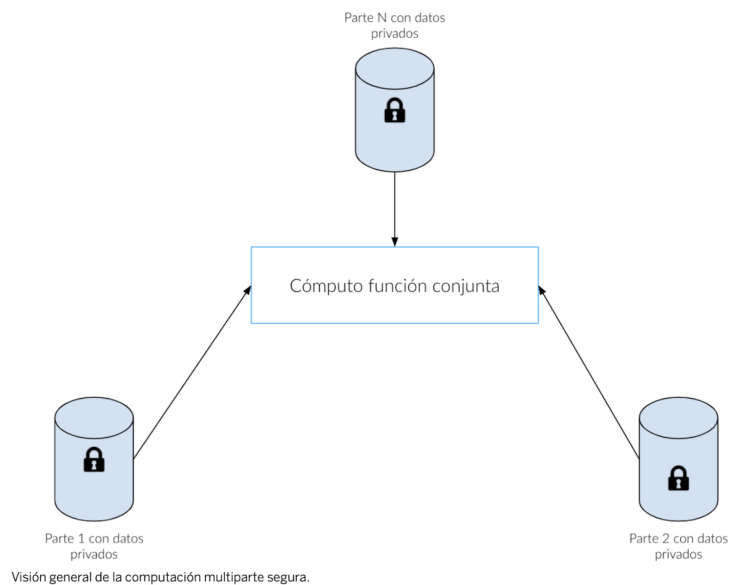


Figura 2: Computación multiparte

En la figura 2 se pueden observar como hay 3 dispositivos conectados entre sí. Cada dispositivo contiene una información privada, pero a nivel computacional se unen para ello. La ventaja de este tipo de modelo criptográfico es que el límite no lo pone el dispositivo, sino el colectivo con el que trabaja.

En la figura 3 se muestra una situación mejorada respecto la anterior figura 2, ya que se proporciona un modelo de cliente-servidor donde

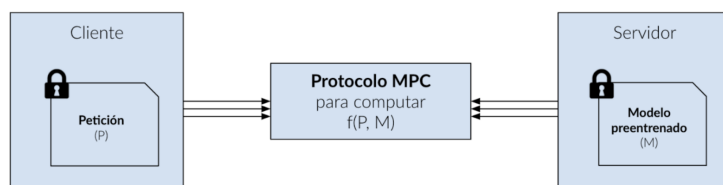


Figura 3: MPC - Machine Learning

el servidor tiene habilitado el machine learning, lo que permite mejorar el rendimiento de esta técnica.

Así pues, si se conectan varios dispositivos con el machine learning habilitado, ya sea en fase de entrenamiento (aprendizaje) o inferencia (funcionamiento real), el rendimiento aumentará respecto el caso anterior, pues a una misma petición, si se ha entrenado un servidor, la respuesta será más rápida.

2.2.2. Privacidad diferencial

La privacidad diferencial provee un paradigma diferente al que ofrece MPC, ya que aquí el objetivo es la seguridad de la información en si misma.

Mientras que en el caso anterior, la finalidad es que cada parte mantenga su información privada y que se pueda computar por otros dispositivos pero sin que estos sepan de que va, aquí el objetivo es que la información tratada, en caso de análisis por terceros, no puedan proporcionar datos privados.

Ejemplo de privacidad diferencial en un juego (Albert Márquez)

Ana y Pedro están jugando a tirar la moneda y adivinar si sale cara o cruz. Ellos, no obstante, no quieren que los demás sepan que resultado les ha salido, por lo tanto, pintan la cara de azul y la cruz de rojo. Al llegar Albert, les pregunta que hacen y ellos le dicen que ha ganado Ana por 6 a 4 en las 10 tiradas que han hecho y le dice que ella jugaba por azul y Pedro por rojo.

Aunque Albert sabe de que va el juego, no sabe que opción es realmente cara y cual cruz, solamente sabe que Ana ha ido a por azul y Pedro por rojo.

Se ha añadido ruido a los datos.

En la figura que se muestra a continuación se ve un ejemplo más elaborado que el anterior. En este caso, se presentan 3 personas (todas ellas con datos sensibles) y que comparten sus datos con un hospital y/o una empresa. El objetivo es que dicha información no salga de ahí y en caso de que así fuese, no se pueda interpretar.

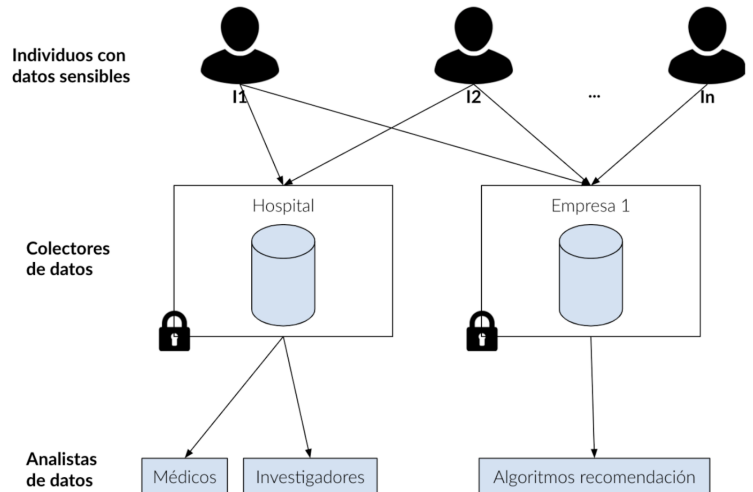


Figura 4: Privacidad diferencial

Para proceder a aplicar esta técnica criptográfica, hacen falta mecanismos. En este documento se tratará de forma teórica el de Gauss y Laplace, y en el apartado práctico únicamente Gauss, debido a las limitaciones de extensión. El objetivo de utilizar dichos mecanismos es el de aplicar ruido a los datos con los que se trabaja. De la misma forma que en el ejemplo anterior, Ana y Pedro pintaban las monedas, aquí se ha de hacer lo equivalente (aplicar ruido) a los datos.

2.3. Tipos de privacidad diferencial

En la privacidad diferencial, hay dos grandes tipos:

- Privacidad diferencial local:** en este tipo, se aplica ruido sobre los datos originales en el mismo punto en que se recopilan. Esto implica por tanto que, si se desea tener un conjunto de datos limpio no es posible, ya que se obtiene y añade el ruido directamente.

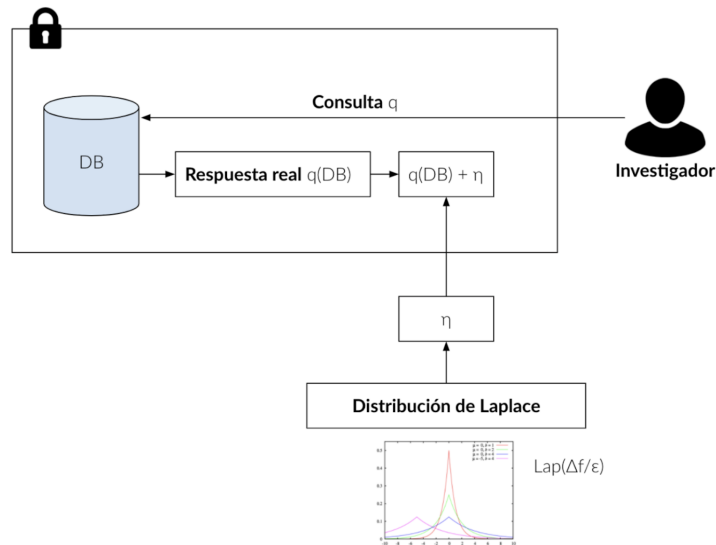


Figura 5: Privacidad diferencial (Laplace)

- **Privacidad diferencial global:** para este tipo, se dispone de un repositorio de datos limpios sin ruido. Al querer usar dichos datos, posteriormente se le añade ruido, pero se dispone de la fuente original.

Mostradas estas dos definiciones, puede dar pie a confusión y pensar que la global es mejor que la local, ya que permite trabajar mejor con la información. No obstante, hay que valorar dicha clasificación no cómo cuál es mejor, sino cuál es más adecuada según la situación.

El modelo global, permite disponer de una versión limpia de los datos y posteriormente una ya con ruido, esto permite que la precisión sobre la cual se realizan los análisis de los datos sea superior. Esto es especialmente interesante si los datos con los que se trata, requieren gran exactitud, cómo por ejemplo cálculos matemáticos con gran cantidad de decimales.

En cambio, el modelo local, al no disponer de los datos limpios, puede aportar una ventaja respecto el modelo global y es que los datos están más asegurados. Esto es debido a que al estar con ruido ya incorporado de serie, esto permitirá que ante un atacante, éste no pueda obtener fácilmente la información.

Por lo tanto, si lo que se requiere es trabajar con mayor precisión, el modelo global tiene claras ventajas, pero si lo que se desea es asegurar al

máximo la información, la local aporta una clara ventaja.

2.4. Compañías conocidas

Dentro del campo de la privacidad diferencial, hay 3 grandes compañías conocidas por la gran mayoría de la población, las cuales han mostrado un claro apoyo a ésta. Estas compañías son Apple, Google i Microsoft.

2.4.1. Apple

Apple ha hecho una fuerte apuesta en la privacidad diferencial del tipo local. Esto lo indican constantemente bajo el lema de “Lo que sucede en tu iPhone, se queda en tu iPhone”. Esto se aplica a todos sus dispositivos, ya que lo implementan a nivel de ecosistema.



Figura 6: Apple - Marketing iPhone

El punto dónde más se ha focalizado la compañía ha sido en los datos biológicos recogidos en su reloj inteligente Apple Watch. Este dispositivo puede recoger valores de oxígeno en sangre, pulso, arritmias, frecuencia cardíaca durante el ejercicio, forma de pisada, etc.

En estos casos, la compañía indica explícitamente que estos datos, residirán en tu cuenta y que Apple no hará uso de ellos. Dichos datos se pueden consultar en otros dispositivos cómo por ejemplo el iPhone, ya que mediante el servicio iCloud se pueden consultar.

Así pues, en el Apple Watch se generan los datos con ruido y los únicos dispositivos que son capaces de interpretar dichos datos son los dispositivos asociados a la cuenta de usuario iCloud.

No obstante, si se desea colaborar con la compañía, ésta ofrece que se puedan compartir datos con ellos (sin incluir identificadores o IPs) y se compromete a borrar la recolección de datos pasados 3 meses.

El valor usado en las sugerencias de búsquedas es de $\epsilon=4$, mientras que en el ámbito de salud $\epsilon=2$. En la detección de la intención de reproducción automática en Safari se usa una $\epsilon=8$.

Para realizar todo ello, Apple utiliza 2 técnicas principalmente:

- **Count Mean Sketch:** codificado con SHA256 y después de codificarla, cada valor se invierte bajo una probabilidad de:

Probabilidad usada (Apple)	
$\frac{1}{(1 + e_{\epsilon/2})}$	(3)

- **Hadamard Count Mean Sketch:** muy semejante al anterior tipo (definición matemática 3), pero usando la

2.4.2. Microsoft

Microsoft por su parte, tiene en funcionamiento un proyecto de código abierto llamado SmartNoise el cual contiene los componentes necesarios para crear sistemas con privacidad diferencial de tipo global. Este proyecto se compone de 2 grandes bloques:

1. **SmartNoise Core:** contiene los mecanismos necesarios para la implementación de un sistema de privacidad diferencial. Tiene 4 componentes:
 - Análisis: descripción de gráfico de cálculos arbitrarios
 - Enlaces: bibliotecas auxiliares para compilar análisis
 - Tiempo de ejecución: tiempo medio donde se realizará la ejecución del análisis
 - Validador: biblioteca de Rust que permite que los análisis sean privados

2. **SmartNoise SDK:** proporciona acceso a datos (mediante procesamiento SQL e implementación con Python, C++, R y más lenguajes), servicio (que proporciona una REST y evaluador (de tipo estocástico, el cual verifica las infracciones de privacidad así como determinar la precisión y el sesgo).

En este caso, Microsoft enfoca esta tecnología para realizar análisis así como machine learning.

2.4.3. Google

Google tiene muchos proyectos en funcionamiento, por lo que es complicado, dadas las limitaciones de extensión que tiene este documento abordarlas todas, pero hay un proyecto en colaboración con OpenMined, el cual tiene a la comunidad expectante ya que tiene como objetivo el de replicar la infraestructura de privacidad diferencial que tiene Google y ponerla a la disposición de los desarrolladores de Python.

2.4.4. Breve comparativa entre compañías y ejemplo

En base a los datos aportados, se puede apreciar que, las 3 compañías utilizan la privacidad diferencial, pero con objetivos diferenciados y con tipos diferentes. Mientras que Apple se focaliza en la local y alrededor del usuario, Microsoft lo hace con un objetivo claramente empresarial y Google a nivel de desarrolladores.

A continuación se mostrará un gráfico dónde se puede ver un ejemplo de los dos tipos de privacidad diferencial tratados en este documento y aplicados al ámbito sanitario. Es un ejemplo muy básico y sencillo, pero muy claro y directo, el cual muestra el enfoque principal de este documento.

Dada las explicaciones anteriores, en la figura (7) , la cual se mostrará a continuación, se ejemplifica un posible funcionamiento de la privacidad diferencial en el ámbito sanitario.

En este caso, se pone de manifiesto cómo en el caso de la privacidad global, se tiene de intermediario un trabajador sanitario acreditado en el sistema, el cual intercambia datos, mientras que en el lado de la privacidad diferencial local se coloca un personal sanitario no identificada, el cual perturba la información (aunque ésta ya esté con ruido añadido).

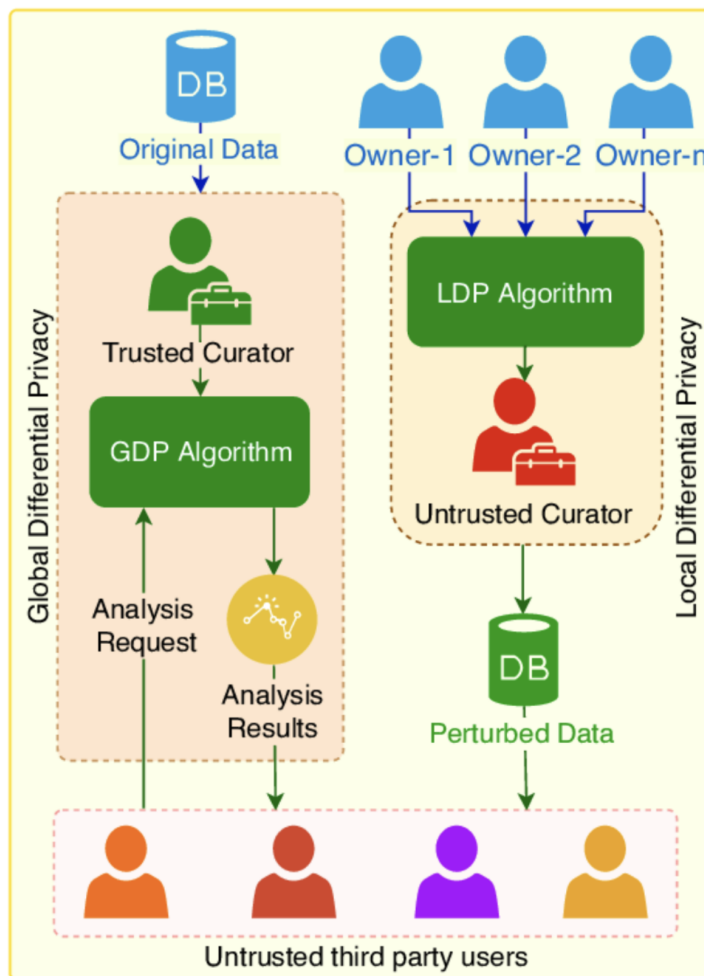


Figura 7: Privacidad global vs local

2.5. Adición de ruido

Tal y como se ha indicado en el apartado anterior, el fundamento de la privacidad diferencial es el de añadir ruido a los datos, para que así no sean reconocibles de forma directa. Aunque hay muchos mecanismos de aplicación de privacidad diferencial mediante algoritmos de todo tipo, hay 2 implementaciones que sobresalen, las cuales son el mecanismo de Gauss y el mecanismo de Laplace.

2.5.1. Distribución Gauss

La distribución de Gauss (también llamada distribución normal) se utiliza especialmente en el ámbito de la estadística. Dicha distribución tiene como principal ventaja que se puede analizar fácilmente, lo que permite poder obtener resultados de manera mucho más rápida. Sus principales características residen en una media 0 y una desviación típica o estándar de 1 (la estándar).

A continuación se muestra un ejemplo de una distribución de Gauss:

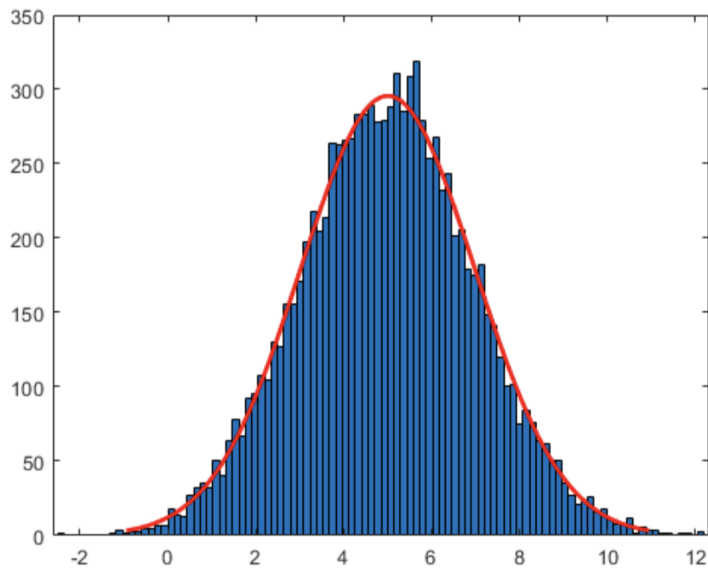


Figura 8: Distribución de Gauss

En la figura 8, se puede observar cómo hay una serie de resultados de un histograma y si se siguen los puntos de resultado de los mismos, se obtiene una distribución de Gauss. En este apartado, no se tratará la media ni la

desviación típica, sino el concepto que sirve para aplicarse en la privacidad diferencial.

A continuación, lo que se va a mostrar es la distribución de probabilidad. Esta permitirá ver cómo se distribuye la probabilidad para que un suceso determinado ocurra.

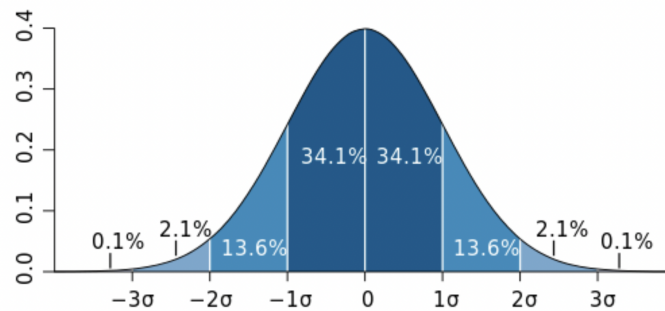


Figura 9: Distribución de probabilidad de Gauss

Tal y como se puede observar, a medida que el valor del eje de ordenadas aumenta, la probabilidad de que un evento ocurra sube. Así pues, los eventos que se encuentran en los extremos, son altamente improbables que sucedan, mientras que los centrales son los que tienen mayor probabilidad.

2.5.2. Distribución Laplace

La distribución de Laplace es semejante a la de Gauss pero con un comportamiento mucho más abrupto. A continuación se puede ver un ejemplo de la misma:

Tal y cómo se puede ver, esta función crece mucho más rápido que la de Gauss. En este caso, la probabilidad de que un suceso ocurra estará menos repartido a lo largo del eje de abscisas, por lo tanto, éste tendrá un margen menor que en el caso de la distribución de probabilidad de Gauss.

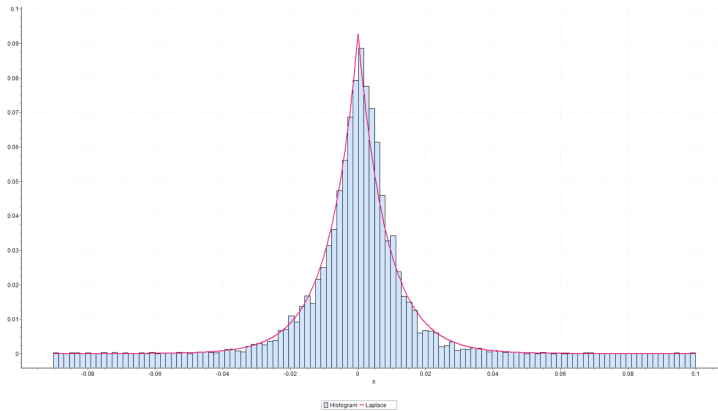


Figura 10: Distribución de Laplace

2.5.3. Comparativa entre distribuciones

Mediante las figuras 8 y 10, se puede observar que ambas tienen una forma semejante, pero mientras que la distribución de Gauss tiene forma cóncava y un progreso suave, la distribución de Laplace es mucho más estrecha y puntiaguda, lo que repercutirá en el resultado de la aplicación de ruido en los datos.

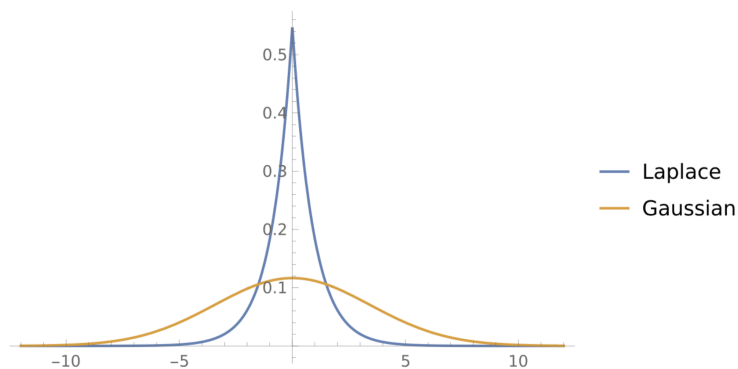


Figura 11: Gauss vs Laplace (distribuciones - I)

En las 2 figuras anteriormente mostradas (11 y 12 respectivamente) se pueden ver 2 casos opuestos, pero que permiten entender mejor la situación.

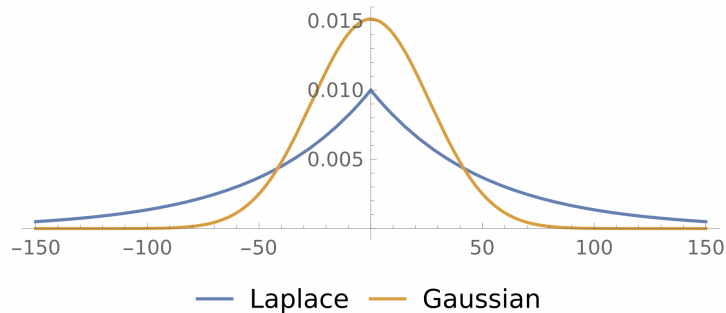


Figura 12: Gauss vs Laplace (distribuciones - II)

La diferencia de uso entre ambas distribuciones se basará en según cuantas estadísticas individuales pueden influir en el resultado. Si una persona o evento puede influir en una estadística únicamente, la distribución de Laplace será mejor, pues tendrá un valor más elevado (caso de la figura 11).

Si por lo contrario, una persona o evento puede influir en varias estadísticas, entonces Gauss será mejor opción, ya que tal y cómo se puede observar en la gráfica, tendrá un valor más elevado (figura 12).

Ejemplo de uso (Damien Desfontaines)

Se dispone de una base de datos con un conjunto de datos definidos por «ID, especialista» donde el ID corresponde a un paciente (identificador único) y especialista al tipo de médico al que visita el paciente.

Si un paciente visita al cardiólogo 10 veces, solo se indicará y contemplará cómo 1 especialista y 1 registro. En caso de que dicho paciente haya visitado mas especialistas, entonces la lista lo reflejará.

Si se desea obtener una estadística del primer caso, la distribución de Laplace permitirá una mejor privacidad diferencial, dado que se trata de una estadística en base a una única información (figura 11). Si por el contrario, se desea valorar todos los especialistas visitados, Gauss proporcionará un mejor resultado (figura 12).

3. Parte práctica

En esta parte, se ha realizado un ejemplo práctico de cómo aplicar adición de ruido. Para ello se ha utilizado el software Numbers para generar un archivo CSV y así dar veracidad académica a los datos generados junto al RStudio para poder trabajar con los datos así como para añadir el ruido a los mismos. Todas las imágenes mostradas en esta sección, son obra del autor de este documento.

Los motivos principales de la utilización de estos programas han sido principalmente 2:

1. Numbers: es la aplicación predeterminada de hojas de cálculo en el sistema operativo MacOS y permite poder generar los archivos CSV de forma rápida y efectiva.
2. RStudio: es una aplicación ampliamente utilizada en el mundo profesional y también académico, sobretodo en el análisis de datos.

Dado que se ha realizado un ejemplo práctico de tipo académico, el archivo usado se ha justificado de forma directa en este documento. El objetivo ha sido generar un archivo CSV (el cual posteriormente se ha usado en RStudio) mediante Numbers. Para ello, se ha decidido usar los siguientes campos: ID, Nombre, Apellido, Nacimiento, Grupo sanguíneo, Estatura, Peso y Enfermedad. En total hay 1000 registros.

Se ha decidido usar 1000 registros para que a nivel computacional sea fácil trabajar con el conjunto de datos así como para verificar el contenido del mismo. Los campos utilizados han sido los anteriormente indicados dado que son datos habituales y de fácil análisis.

3.1. Creación de un archivo CSV en Numbers

Inicialmente se parte de una tabla completamente vacía. En ella primeramente se ha añadido en la columna ID unos registros de 6 dígitos completamente aleatorios (entre 000001 y 999999). Para realizarlo, se ha usado la función:

Función para creación de columna ID (Numbers)

```
aleatori.entre(valor_minim,valor_màxim)
```

El problema de esta función, es que cada vez que se edita cualquier celda en una hoja de cálculo, varia su contenido. Para ello, se ha copiado

toda la columna ID y a continuación se ha seleccionado ”pegar resultados de fórmula” para dejar así bloqueado el contenido.

Las dos columnas siguientes que se han rellenado han sido las de Nombre y Apellido (para que el conjunto de datos sea más fácil de crear y analizar posteriormente, se ha omitido el segundo apellido). Para añadir dicha información, se ha accedido al portal INE y se han buscado 2 estadísticas las cuales corresponden a los 200 nombres (100 de hombre y 100 de mujer) y 100 apellidos más utilizados en España.

Dado que este ejemplo es de tipo académico y que el nombre y apellido no son relevantes en el estudio que se ha realizado, no afecta que se hayan repetido nombres con apellidos (p.e. que hayan aparecido 3 veces Alberto García). Para ello, en Numbers se han generado 2 tablas adicionales. La primera contiene los 200 nombres de hombres y mujeres de forma seguida, así cómo los apellidos también.

La tercera tabla se ha utilizado de forma auxiliar, ya que en ella se ha utilizado la función aleatori (indicada anteriormente) y se ha generado un número aleatorio entre 0 y 1 en la columna AUX2 (en AUX1 se han almacenado los 200 nombres repetidos 5 veces). A continuación se ha realizado una ordenación de forma ascendente de la columna AUX2. El resultado ha sido el siguiente:

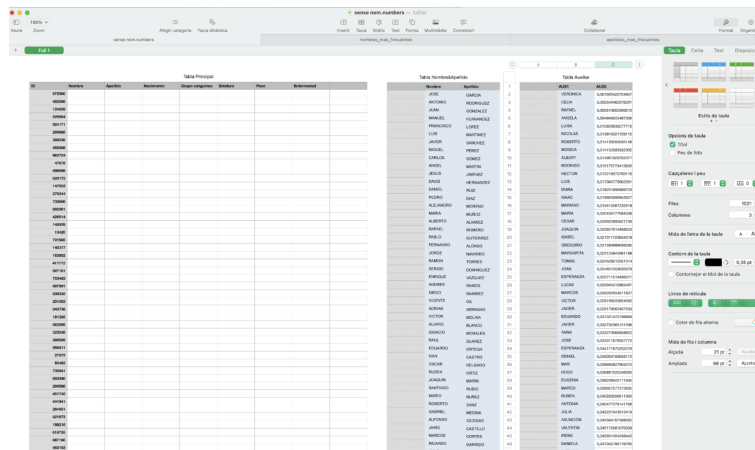


Figura 13: Numbers (MacOS)

Para el campo Apellido se ha actuado de la misma manera. Dado que la columna AUX2 no se había bloqueado mediante el pegado de fórmula explicado anteriormente, si se realizaba un copiado del apellido, los valores de la columna AUX2 y AUX4 variaban. Nuevamente se han ordenado de forma ascendente en base a AUX4. De esta manera se han obtenido los 1000 nombres y 1000 apellidos (los nombres repetidos 5 veces y los apellidos 10) repartidos de forma aleatoria.

El siguiente campo ha sido Nacimiento. Aquí se ha realizado una creación aleatoria de valores comprendidos entre 1930 y 2015. Para el caso de la estatura ha sido entre 150 y 200 (expresado en centímetros).

El peso se ha colocado en base a la estatura, multiplicado por 0,44 y con una suma aleatoria de entre -0,8 y +1,2 del valor anterior para que se tengan unos valores acordes, con ligeras variaciones y no todo el mundo tenga la misma proporción de peso respecto estatura.

Para el grupo sanguíneo, se ha consultado la web de Cruz Roja, para conocer la distribución de grupos sanguíneos en España. Dado que lo indican en porcentajes con un valor decimal y el conjunto de datos creado son 1000 registros, se ha añadido de forma directa a la tabla para generar una ordenación aleatoria a la hora de asignar a las personas. Por ejemplo, el grupo AB- representa el 0.5 % de la población, al tener 1000 registros, se ha asignado a 5 personas.

A continuación se han añadido las enfermedades. Para este caso, se ha utilizado cómo base inicial una noticia del hospital Vall d' Hebron dónde indican que el asma podría no ser un factor de riesgo para el desarrollo del. En dicha noticia se indica que el estudio no es del todo concluyente debido a falta del número de pacientes analizados.

Inicialmente se ha consultado la web del INE para tener una fuente de información respecto a la distribución de muertes según enfermedades del sistema respiratorio. El total, para el mes de enero de 2019, ha sido de 6993 personas, la distribución de las cuales ha sido:

Estadística de defunciones según la causa de muerte
Defunciones por año, enfermedades del sistema respiratorio y mes de defunción.
 Unidades: Número de defunciones

Tabla	Gráfico
	Enero
2019	
062-067 X. Enfermedades del sistema respiratorio	6.693
062 Influenza (gripe) (incluye gripe aviar y gripe A)	417
063 Neumonía	1.337
064 Enfermedades crónicas de las vías respiratorias inferiores (excepto asma)	1.662
065 Asma	129
066 Insuficiencia respiratoria	242
067 Otras enfermedades del sistema respiratorio	2.906

Figura 14: Defunciones por motivos respiratorios

Si se hacen los valores porcentuales, los resultados son los siguientes:

- Influenza de la gripe: 6,23 %
- Neumonía: 19,98 %
- Enfermedades crónicas (asma no incluida): 24,83 %
- Asma: 1,93 %
- Insuficiencia respiratoria: 3,62 %
- Otras: 43,42 %

La metodología de aplicación de estos datos en el archivo CSV ha sido la misma que con los grupos sanguíneos. Dado que en los asmáticos hay un 1.93 %, se han puesto en la tabla 20 pacientes (se ha redondeado al alza para que posteriormente no hubiesen registros sin enfermedades).

A continuación se muestra el archivo CSV resultante:

Figura 15: CSV final

Dado que el asma es el que tiene menos porcentaje, es el que se usará en este archivo CSV como piedra angular. Por ello, se establece una situación hipotética (con objetivo meramente académico):

Situación hipotética

“El 100 % de las personas que tienen asma y son del grupo A-, son inmunes al SARS-CoV2”.

Así pues, se ha creado una novena columna llamada inmunidad, la cual contiene 0 si es negativo en inmunidad y 1 si es positivo en ella. Con todo ello, se ha dispuesto de un archivo CSV en el que posteriormente se han realizado los pertinentes estudios de privacidad diferencial con RStudio.

3.2. Análisis de la información en RStudio

En este programa, se realizarán varias verificaciones para contrastar la veracidad de los datos creados anteriormente. Así mismo, el objetivo del uso de este programa es el de la adición de ruido.

¿Por qué añadir ruido?

La adición de ruido proporciona numerosas ventajas en el tratamiento de los datos así como de su seguridad. En este caso, dado que se ha tratado con datos de pacientes inventados, se ha podido preservar la información aunque hubiese acontecido algún problema de seguridad y este archivo hubiese salido de este contexto académico. La función principal de la adición de ruido en este documento, ha sido la de preservar el contenido del campo inmunidad a toda costa.

Encontrar una tabla con unas características mínimamente semejantes y de acceso público en Internet respecto la que se ha tratado en el anterior apartado no es posible ya que se hubieran visto comprometidos los datos de los pacientes y se hubiera puesto en peligro el secreto médico cómo tal. Es por ello que se ha decidido crear una por cuenta del propio autor de este documento.

Para empezar a trabajar con este archivo, primeramente se ha realizado la importación archivo CSV en RStudio. Para ello, se ha ejecutado el siguiente código:

```
library(readr)
Taula <- read_delim("Taula.csv", delim = ";",
  escape_double = FALSE, trim_ws = TRUE)
View(Taula)
```

Figura 16: RStudio - Comando lectura datos

Así pues, en un tiempo no superior a 10 minutos, un atacante ha podido detectar el patrón y obtener la información deseada.

Debido a que adicionalmente se tienen los nombres, primer apellido e identificadores, para el atacante no es difícil poder obtener información confidencial. Ha sido una obtención rápida y efectiva de datos, dado que éstos estaban directamente expuestos. Si el parámetro de inmunidad no hubiese sido de tipo booleano, quizás la obtención de datos hubiese sido más difícil.

A nivel gráfico, a continuación se muestra la relación de pacientes positivos en inmunidad.

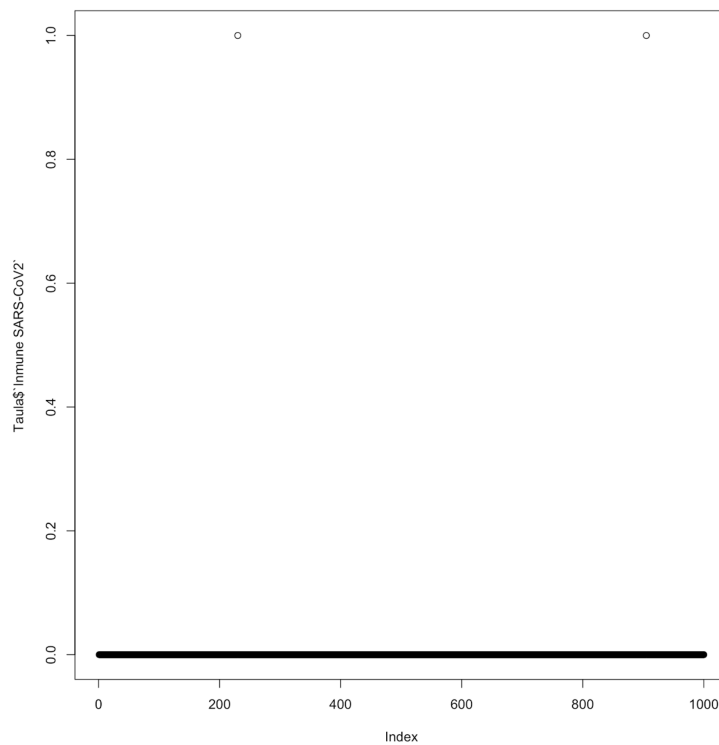


Figura 19: RStudio - Detección gráfica de pacientes inmunes

Tal y cómo se observa, tanto a nivel analítico cómo a nivel gráfico, detectar los casos excepcionales, ha sido muy fácil y rápido.

Dada la naturaleza de la variable, pues un paciente puede ser inmune o vulnerable (no existen medios estados), esto obliga a establecer la variable como un tipo booleano.

Se podría usar una segunda variable la cuál indicase la severidad de la infección (en caso de no ser inmune), pero el hecho es que la inmunidad estaría presente o no, dependiendo del sistema inmunitario del paciente.

El gran problema de ello, es que es muy fácil identificar la información. No obstante, para ello, en este documento se han tratado 2 métodos para poder solucionar este problema. Dado que el valor de inmunidad se puede obtener gracias a los campos de grupo sanguíneo y el de enfermedad, éste primero se puede alterar mediante la adición de ruido.

Que hayan pacientes con el grupo A- y asmáticos es proporcional a la densidad de población, pero la relación de ello con la inmunidad, al añadirle el ruido, no se verá tan clara como antes. Este documento se basa en un ejemplo académico, pero su aplicabilidad es plenamente factible en otros casos.

3.2.1. Adición de ruido aleatorio

La primera solución a estudiar ha sido la de la función Jitter. Esta función añade una parte de ruido a los datos. La sintaxis de la misma es:

```
jitter(x, factor = 1, amount = NULL)
```

Figura 20: RStudio - Función Jitter

Por ello, lo que se ha realizado es aplicar dicha función a la columna de inmunidad. Los parámetros de la función Jitter indican lo siguiente:

- `x` → El vector de datos numéricos con el que se trabaja
- `factor` → Un valor numérico con el que se trabaja en el siguiente parámetro
- `amount` → De forma predeterminada es $\text{factor} * d / 5$ donde d es la menor distancia entre los valores del vector `x`.

Así pues, se ha aplicado la función Jitter a la columna de inmunidad en RStudio:

```
Taula$'Inmune SARS-CoV2' <-  
jitter(Taula$'Inmune SARS-CoV2', factor = 1)
```

Figura 21: RStudio - Aplicación Jitter

Y a continuación se hace su representación gráfica, se obtiene el siguiente resultado:

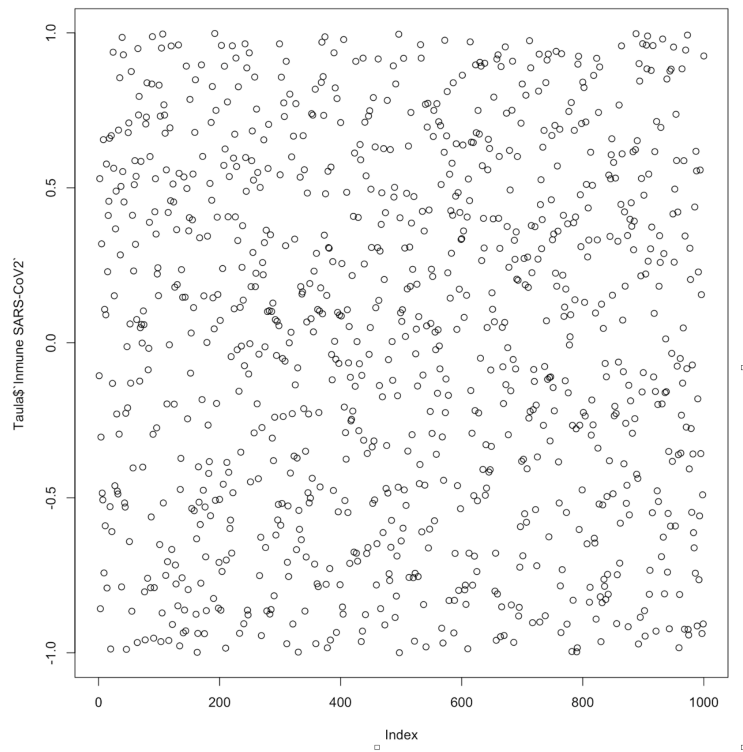


Figura 22: RStudio - Representación gráfica Jitter

Tal y cómo se puede observar, los 1000 pacientes tienen un estado de inmunidad comprendido entre -1.0 y 1.0. La ventaja de este método es que todos los valores asociados a la inmunidad son únicos. En este caso, si se realiza un filtrado (de forma gráfica para que quede todo mejor indicado), se puede observar, cómo a pesar de tener inmunidad, los dos pacientes tienen un valor ahora muy dispar:

	ID	Nombre	Apellido	Nacimiento	Grupo sanguíneo	Estatura	Peso	Enfermedad	Inmune SARS-CoV2
1	425942	ROBERTO	PRIETO	1967	A-	174	78	Asma	0.3298672
2	437707	JAIME	GOMEZ	1948	A-	156	70	Asma	0.9603746

Figura 23: RStudio - Filtro Inmunidad (Jitter usado)

La principal desventaja de este método es que si se quiere recuperar la información, dada la aleatoriedad aplicada a los datos, no es posible y por lo tanto, el campo de inmunidad pasa a ser inútil. Se ha destruido información en el camino, aunque se pueda deducir gracias al grupo sanguíneo y la enfermedad, pero esto solo lo sabe la persona que genera el CSV.

Es por lo tanto, una adición de ruido altamente ineficiente para este caso.

Este tipo de adición de ruido podría ser realmente útil si de lo que se trata es de detectar el comportamiento de una función por ejemplo. Si sabemos que una función oscila entre -10 y 10 y otra en cambio entre 200 y 300, se puede utilizar esta función para detectar el comportamiento de cada una de ellas añadiendo puntos aleatorios comprendidos entre esos márgenes de valores.

Así mismo, también se puede usar esta función a modo de padding para realizar el relleno de campos vacíos si estos requieren información obligatoriamente (aunque esta resulte irrelevante).

3.2.2. Adición de ruido Gaussiano

Para poder solucionar el caso anterior, se ha trabajado con el ruido Gaussiano. Este tipo de adición, a diferencia de la anterior, proporciona una situación mucho más controlada.

Para ello, lo que se ha realizado es aplicar sobre el campo de inmunidad un cambio de valores, los cuales han sido una secuencia definida entre un valor $-n$ y $+n$ (esto es el eje de abscisas sobre el cuál se construye la distribución Gaussiana). A continuación se ha definido el eje de ordenadas mediante el comando `dnorm`, en el cual se indica la variable x (anteriormente explicada), la media y la desviación.

Métodos para trabajar sobre este sistema hay varios, por ejemplo definir un vector en la columna de inmunidad con los parámetros X e Y o bien dejar en la columna de inmunidad solamente la variable X con los valores generados y tener que trabajar con la Y para analizarlos posteriormente.

En este ejemplo, en la columna de inmunidad, se han dejado los valores generados de X , para poder realizar la representación de la distribución Gaussiana de forma más práctica simplemente modificando los parámetros del comando `dnorm`, el cual permite realizar la representación.

Para este caso y poder mostrar de la forma más visual posible este concepto, se ha decidido utilizar una secuencia comprendida entre los valores -500 y $+500$ (ambos incluidos).

Primeramente se ha definido dicha secuencia mediante el comando:

```
x <- seq(-500, 500, by = 1)
```

Figura 24: RStudio - Creación secuencia

Con este comando se indica la secuencia y de cuánto serán los saltos.

A continuación se muestra la gráfica asociada a dicha secuencia:

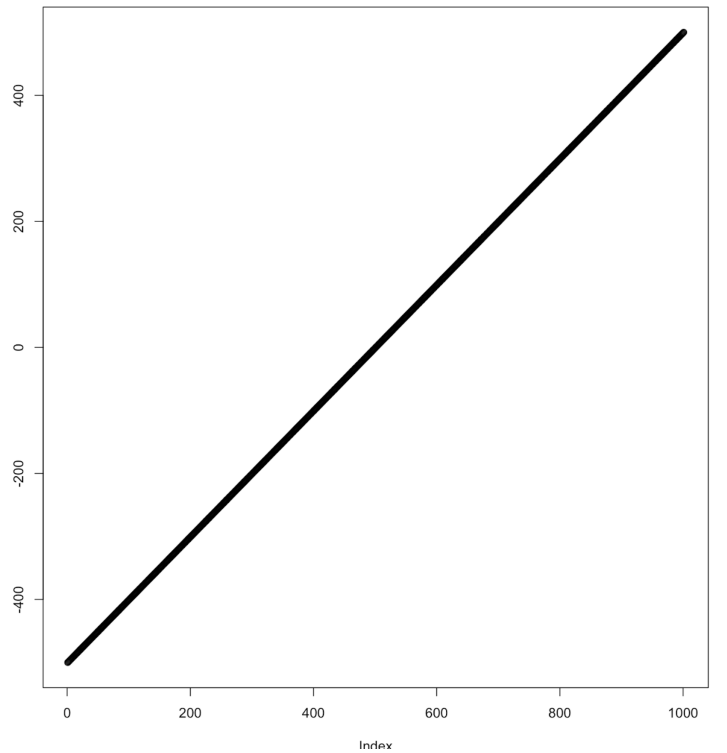


Figura 25: RStudio - Secuencia numérica (Gráfico)

Tal y cómo se muestra en dicha gráfica, los valores están asignados de forma correcta, pues siguen una gráfica con una progresión claramente lineal y definida. Así mismo, también se puede ver que no hay valores repetidos (aspecto importante el cual se tratará más adelante).

El siguiente paso ha sido el de dar la altura de la distribución de probabilidad para cada punto definido en el eje de abscisas en base a una mediana y una desviación estándar. Esto se obtiene mediante el comando `dnorm`. El objetivo es obtener una situación en la que el resultado sea lo más uniforme posible y que no se pueda detectar patrones e información relevante de forma directa.

El comando aplicado ha sido el siguiente:

```
y <- dnorm(x, mean = 0, sd = 100)
```

Figura 26: RStudio - Comando dnorm

Al realizar la gráfica en base a los valores X e Y, se obtiene el siguiente resultado:

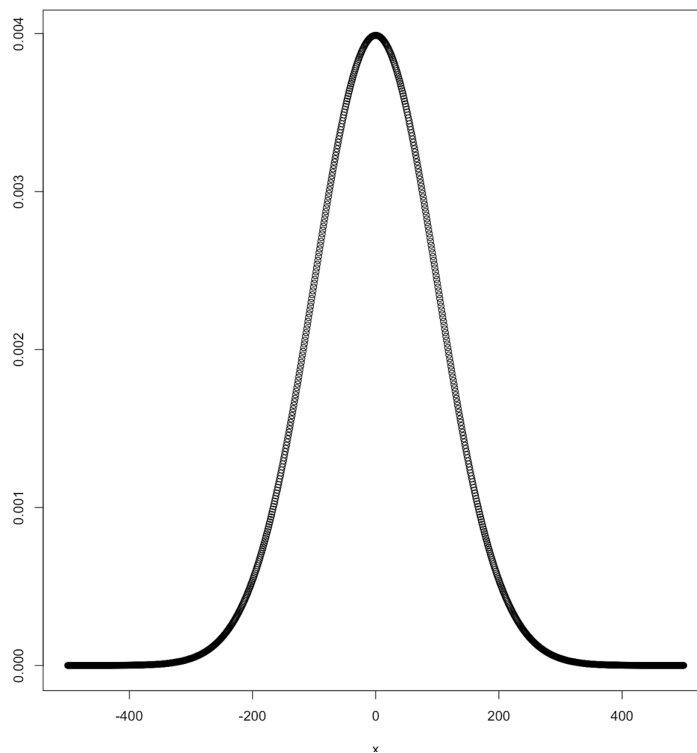


Figura 27: RStudio - Distribución Gaussiana

Tal y cómo se muestra en esta imagen, se puede observar una distribución Gaussiana. El fundamento de ello es que los valores de inmunidad se mantengan ocultos y no salgan cómo unos y ceros.

Es ahora cuándo surge la cuestión de que se puede realizar con los datos creados. Es en este punto dónde se muestra la propuesta académica sobre ello.

Tal y cómo se puede observar en la gráfica generada, a medida que el eje X se acerca a 0, el valor asociado al eje Y aumenta. El objetivo de la distribución de probabilidad es obtener un valor fehaciente de que un cierto evento ocurra.

En el caso a tratar, el evento real es que un paciente presente inmunidad (1) o no (0), la variación que se ha realizado es que en vez de disponer dos únicos valores, lo que se ha hecho es crear una secuencia numérica, lo que permite camuflar dicha información. No obstante, nos encontramos con el mismo problema que con el ruido aleatorio y es que se pierde el sentido de los datos. Por ello, la proposición es la siguiente:

Propuesta de solución para la recuperación de datos

Dado que se tienen X casos de pacientes inmunes y una secuencia numérica comprendida entre $-n/2$ y $+n/2$ (dónde n es el número total de pacientes, que es 1000), se realizará una distribución de valores encubierta respecto la inmunidad.

Así pues, dado que se tienen 2 pacientes que cumplen las condiciones indicadas, se utilizarán los valores opuestos -1 y +1 del eje X. Si fuese un caso impar (p.e 3 pacientes) se usaría -1, 0 y +1. Con 4 pacientes: -2, -1, +1, +2. Con 5 pacientes: -2, -1, 0, +1, +2.

Así pues, el procedimiento a seguir para llevar a cabo esta idea ha sido el siguiente:

1. Previamente a la asignación de valores mediante a la secuencia numérica, se identifican los casos que cumplen la condición de que inmunidad = 1. Se guardan las posiciones de dichos registros. Esto se ha realizado mediante el comando:

```
selected <- Taula[Taula$
'Immune SARS-CoV2'==1,]
```

Figura 28: RStudio - Asignación a variable de pacientes inmunes

Así pues, si se verifica posteriormente su contenido, se puede observar cómo almacena los datos de dichos registros.

```
> selected <- Taula[Taula$`Immune SARS-CoV2`==1,]
> selected
# A tibble: 2 x 9
  ID Nombre Apellido Nacimiento `Grupo sanguine... Estatura Peso Enfermedad `Immune SARS-Co...
  <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <chr> <dbl>
1 425942 ROBERTO PRIETO 1967 A- 174 78 Asma 1
2 437707 JAIME GOMEZ 1948 A- 156 70 Asma 1
> |
```

Figura 29: RStudio - Pacientes inmunes (nueva variable)

2. Se ha aplicado la asignación de la secuencia numérica. Cómo detalle, RStudio cuando se generan los valores secuenciales, crea 1001 registros dado que contempla el 0, por lo tanto hay que borrar un registro, así pues quedará del -499 al 500 o bien eliminar un registro de los 1000, para que queden 999 y así se puedan distribuir en -499, 0 y 499. Se usará el primer caso

```
x <- seq(-500, 500, by = 1)
x <- tail (x, -1)
Taula$`Immune SARS-CoV2`<- x
```

Figura 30: RStudio - Asignación nueva secuencia numérica

3. A continuación, se han buscado los pacientes con el el ID indicado en el punto 1) los cuales tenían valor 1 en inmunidad:

```

> Taula[Taula$ID ==425942,]
# A tibble: 1 × 9
  ID Nombre Apellido Nacimiento `Grupo sanguineo` Estatura Peso Enfermedad `Inmune SARS-CoV2`
  <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <chr> <dbl>
1 425942 ROBERTO PRIETO 1967 A- 174 78 Asma -270
> Taula[Taula$ID == 437707,]
# A tibble: 1 × 9
  ID Nombre Apellido Nacimiento `Grupo sanguineo` Estatura Peso Enfermedad `Inmune SARS-CoV2`
  <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <chr> <dbl>
1 437707 JAIME GOMEZ 1948 A- 156 70 Asma 405
>

```

Figura 31: RStudio - Parámetro inmunidad modificado

Ahora el valor de inmunidad se ha visto modificado por -270 y 405 respectivamente.

- El siguiente paso ha sido buscar los pacientes cuyo valor de inmunidad es -1 y +1

```

> Taula[Taula$`Inmune SARS-CoV2`==-1,]
# A tibble: 1 × 9
  ID Nombre Apellido Nacimiento `Grupo sanguineo` Estatura Peso Enfermedad `Inmune SARS-CoV2`
  <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <chr> <dbl>
1 584205 ESTEBAN VELASCO 1995 0+ 160 71 Insuficiencia -1
> Taula[Taula$`Inmune SARS-CoV2`==1,]
# A tibble: 1 × 9
  ID Nombre Apellido Nacimiento `Grupo sanguineo` Estatura Peso Enfermedad `Inmune SARS-CoV2`
  <dbl> <chr> <chr> <dbl> <chr> <dbl> <dbl> <chr> <dbl>
1 99161 GREGORIO ROJAS 1965 A+ 182 81 Crónica 1
>

```

Figura 32: RStudio - Parámetro inmunidad -1 y +1

- El paso final, ha sido intercambiar los valores de inmunidad entre estos 4 pacientes. Dado que el ID es único en todos los pacientes, el procedimiento aplicado ha sido buscar dicho identificador y la fila asociada al valor de inmunidad, modificarlo:

A los pacientes que no tienen inmunidad y que tenían en la sucesión asignada el -1 y +1 se les ha cambiado el valor por el de los pacientes que realmente sí tienen inmunidad:

```

Taula[Taula$ID==584205,9] = -270
Taula[Taula$ID==437707,9] = 405

```

Figura 33: RStudio - Intercambio valores inmunidad (1)

Finalmente a los pacientes que sí que tienen inmunidad, se les ha aplicado el +1 explicado con anterioridad. Si hubiesen sido 3 pacientes, el tercer paciente sería el que tuviese valor 0.

```
Taula[Taula$ID==425942,9] = 1
Taula[Taula$ID==437707,9] = -1
```

Figura 34: RStudio - Intercambio valores inmunidad (2)

Así pues, se ha dejado ya definida la inmunidad correctamente. No obstante, ahora quedaría dejar un indicador conforme cuántos pacientes inmunes hay. Si un médico A le indica a un médico B que hay pacientes inmunes, pero no cuántos, de nada servirá realizar este procedimiento. Así pues, se ha añadido una columna adicional a la tabla de datos, donde se pondrá un indicador de cuántos pacientes inmunes hay.

Dado que hay información explícita cómo el nombre, apellido, grupo sanguíneo, enfermedad, año de nacimiento, etc. y que se repiten varias veces (pues hay un número finito bastante limitado (dadas las fuentes consultadas y número de registros) una opción hubiese sido aplicar cifrados comunes a éstos datos cómo por ejemplo el cifrado de César, pero dada la repetición de los datos, sería altamente peligroso, ya que el atacante podría detectar los patrones mediante minería de datos aplicando reglas de asociación.

Es por ello, se ha decidido aplicar un cifrado a nivel de archivo. Dado que el campo clave que se desea preservar a toda costa, es el de inmunidad, se ha aplicado una segunda capa, la cual será a nivel global de archivo.

Se ha aplicado por lo tanto privacidad diferencial sobre el campo inmunidad (mediante adición de ruido Gaussiano) y cifrado a nivel de archivo para preservar la demás información de los pacientes.

Ahora mismo, si un médico A, envía a un médico B, éste no sabrá cuántos pacientes habrá con inmunidad, ya que se ha modificado el campo de inmunidad mediante la adición de ruido gaussiano. Por ello, se generará una columna adicional la cual se llamará "Patients.group" (en este caso su valor será 2).

El procedimiento completo es el siguiente:

1. Se ha añadido la columna de Patients_group con el valor de 2 (pacientes con inmunidad pero que no se quiere indicar explícitamente) y se ha exportado la tabla a un nuevo archivo CSV. Esto se ha realizado mediante los comandos:

```
Taula$Patients_Group = 2  
write.csv(Taula, "patients.csv")
```

Figura 35: RStudio - Creación Columna Número inmunes

2. Ahora que ya se tiene el fichero completado (con el campo de inmunidad con adición de ruido y con los cambios de parámetros indicados anteriormente), la columna de Patients_Group creada y todo analizado, se ha procedido a cifrar el archivo. Lo primero, si no se tiene ya instalado, es poner a disposición de RStudio la librería encryptr y activarla.

```
install.packages("encryptr")  
library(encryptr)
```

Figura 36: RStudio - Instalación paquete encryptr

3. Dado que se ha cifrado un archivo, esto ha implicado la generación de una clave privada y una pública, en el siguiente apartado se explicará más en detalle el funcionamiento de todo el sistema. Esto se logra mediante el comando genkeys().

Este proceso requiere crear una contraseña. Un buen procedimiento es ponerle un identificador interno que solamente conozcan el emisor y el receptor. Para este ejemplo se ha usado como contraseña TFG21AMNS1990 (código interno del proyecto el cual contiene la asinatura, año, iniciales del autor y fecha de nacimiento).

4. El paso siguiente es el cifrado del archivo. Para ello se ha usado el siguiente comando:

```
encrypt_file("patients.csv")
```

Figura 37: RStudio - Encriptación archivo CSV

5. Este último proceso ha generado un archivo llamado patients.csv.encryptr.bin. Este archivo es el que se ha de enviar a otra persona con la que se acuerda analizar la tabla. Debido a la generación de claves y cifrado del archivo original, si no se tiene al alcance ello, será muy complicado acceder a la información.
6. La persona que reciba el archivo, deberá de disponer de la clave para poder descifrarlo. Para ello, se debe aplicar el comando:

```
decrypt_file("patients.csv.encryptr.bin",  
file_name = "output_file.csv")
```

Figura 38: RStudio - Desencriptación archivo CSV

Es ahora cuándo se puede visualizar el archivo finalmente. A continuación se muestra el proceso en RStudio:

The image shows a screenshot of the RStudio interface. The top pane displays a data table with the following columns: Nombre, Apellido, Nacimiento, Grupo sanguineo, Estatura, Peso, Enfermedad, Inmune SARS-CoV2, and Patients_Group. The table contains 34 rows of patient data. The bottom pane shows the R console with the following code and output:

```

R 4.1.2 ~-~
> install.packages("encryptr")
trying URL 'https://cran.rstudio.com/bin/macosx/contrib/4.1/encryptr_0.1.3.tgz'
Content type 'application/x-gzip' length 85534 bytes (83 KB)
=====
downloaded 83 KB

The downloaded binary packages are in
  /var/folders/Gs/sws9s3xn6rx3y418t3ht2zym000gn/T//RtmpFWsW9z/downloaded_packages
> library(encryptr)
> Taula$Patients_Group = 2
> write.csv(Taula, "patients.csv")
> genkeys()
Private key written with name 'id_rsa'
Public key written with name 'id_rsa.pub'
> encrypt_file("patients.csv")
Encrypted file written with name 'patients.csv.encryptr.bin'
> decrypt_file("patients.csv.encryptr.bin", file_name = "output_file.csv")
Decrypted file written with name 'output_file.csv'

```

Figura 39: RStudio - Proceso encriptación y desencriptación

Si a continuación se analizan el archivo creado y el encriptado, se puede observar cómo la diferencia de espacio que representan es 1KB, que tratándose de archivos de 65 y 64KB respectivamente, representa un incremento del 1,54%.

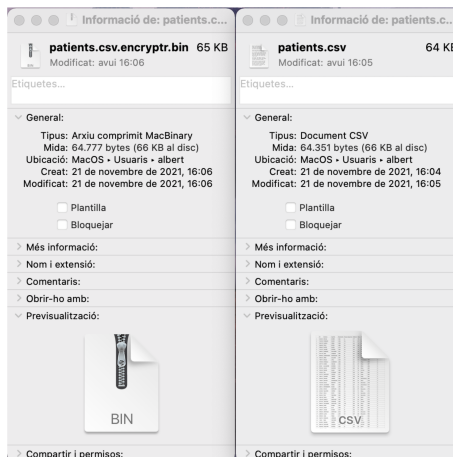


Figura 40: Comparativa archivos

Cuestión: ¿Por qué no criptografía de clave pública para los casos planteados?

Si bien es cierto que el sistema de claves implementado en el archivo de pacientes (figura 39), es una muy buena idea y añade una capa adicional de seguridad, no obstante, este sistema implica un intercambio de llaves. Esto comporta que haya un intercambio de credenciales entre 2 o varios usuarios.

Esto conlleva a que si estas llaves están generadas con un cifrado muy débil o hay un usuario interceptando las comunicaciones (Man-in-the-middle), esto hace que la información se vuelva altamente vulnerable. Para dar solución a estos casos, se debe de aplicar con anterioridad SSL/TLS.

Así mismo, con la privacidad diferencial se altera el contenido (adición de ruido) y si un atacante accede a dicha información, no será capaz de leerla a no ser que conozca el modelo matemático del cual se ha tratado dicha información.

4. Parte aplicación SI

Tal y cómo se ha podido observar en el apartado anterior, aplicar ruido gaussiano a un campo específico, proporciona beneficios, los cuales además se pueden ver aumentados si se cifra el archivo que contiene la información.

Debido a la situación acontecida por el SARS-CoV2, esto ha provocado que el sector sanitario haya tenido que realizar una fuerte inversión en el ámbito de las tecnologías, más específicamente en las comunicaciones entre personal sanitario y paciente. El hecho de que haya ocurrido una pandemia y que haya afectado tan rápido, causa que el ámbito SI/TI del sector sanitario haya tenido que evolucionar más rápido de lo previsto.

Para poder implementar un sistema tecnológico sanitario correcto y eficiente, la opción más segura es certificar las tecnologías que lo forman. Es por ello que debido al campo que se trata en este documento, el cual es la seguridad en los datos biológicos, se propone implementar un modelo propio basado en la ISO/IEC 27701.

El principal motivo de ello es dado que dicha certificación se integra con la RGPD (Reglamento General de la Protección de Datos) y dado que este documento ha girado en torno al concepto de la seguridad de los datos biológicos, el crear un modelo adaptado a las necesidades citadas, será una oportunidad para poder llevar a cabo dicha idea.

Dado el alcance que representa una ISO y la RGPD en sí misma, y las limitaciones de extensión de este documento, se realizará un modelo personalizado y adaptado a las necesidades que se han tratado con anterioridad.

Para llevarlo a cabo en su plenitud, se tendría que hablar de las bases legales, reglas y excepciones, tipos de sanciones, etc. En este documento se cubrirán los apartados de consentimiento, seudonimización y violación de datos. Se han seleccionado estos 3 apartados dado que son los que representan más implicación de cara al personal sanitario y al paciente (los dos grandes bloques que interactúan en este documento).

Por este motivo, el modelo que se presenta en este documento se basará en la ISO/IEC 27701, la RGPD y un manual elaborado por la Confederación Canaria de Empresarios el cual implementa una serie de pasos prácticos para implementar un sistema de gestión de privacidad de la información.

Primeramente, se muestra a continuación un esquema básico del planteamiento que se realiza sobre el tratamiento de la información en este documento (imagen creada con Apple Pages por el autor de este documento):

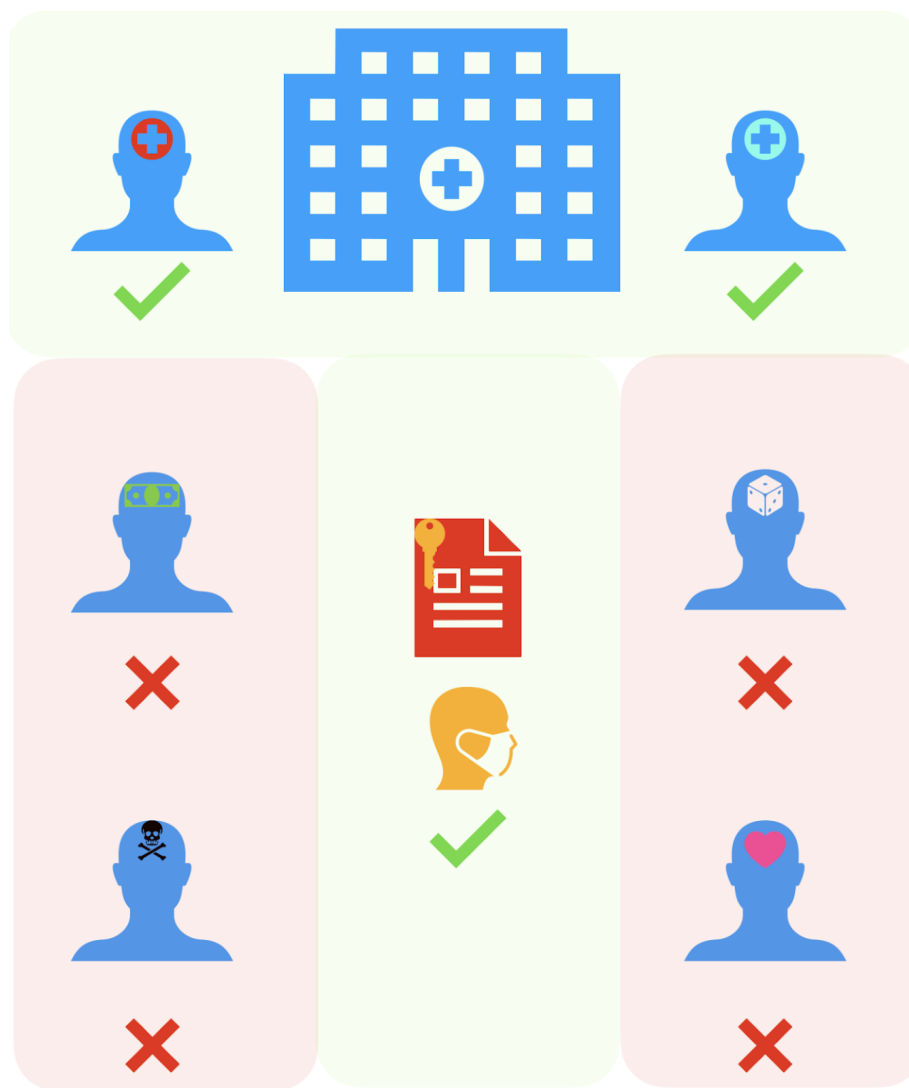


Figura 41: Planteamiento

Este esquema representa el principio fundamental de funcionamiento en un hospital de forma ideal:

1. Únicamente el paciente y el personal sanitario tendrán acceso a la información biológica.
2. Personal no sanitario, gestores del hospital, atacantes, personas al azar o seres allegados no tienen porqué tener acceso a esta información
3. La documentación en ningún caso estará abierta (es decir, el archivo no estará en texto plano o sin cifrar de forma directa).
4. La comunicación se establecerá entre los diferentes integrantes del personal sanitario y el paciente. El resto de secciones no se contempla su acceso a dicha información.

El punto de partida en esta situación se basa en 3 pilares fundamentalmente:

- Mucha información se gestiona en papel. Este problema conlleva que sea altamente accesible en caso de pérdida.
- Hoy en día se tiene acceso a un vasto mundo de información y aunque una persona obtenga una información que no entiende, puede llegar a traducirla con los medios adecuados. El objetivo es que, si esto ocurre, la información no se pueda descifrar. Hay que securizar la información relativa a un paciente para que eso no acontezca.
- Dado que hay muchos puntos de acceso a una información, la mejor forma de que se reduzca su rango de alcance, será que únicamente el personal acreditado pueda acceder a ella.

Por lo tanto, hay que crear 3 pilares: digitalización, securización y acreditación. Para ello, se crearán una serie de fases, las cuales permitirán cumplir dicho objetivo.

4.1. Fase 1 - Protección de los datos

En esta primera fase. El objetivo es saber qué formación, en materia de seguridad tiene el personal y el centro. Para ello, una buena práctica, es realizar un check-list asociado a dichos conocimientos. Los apartados que han de contener dicho check-list son los siguientes:

- Tecnología utilizada en la recolección de datos del paciente y metodología aplicada
- Conocimiento de las ventajas y desventajas de los diferentes tipos de tecnologías requeridas
- Finalidad de la recolección de datos
- Tratamiento de los datos del paciente
- Cesión de datos del paciente entre el diferente personal autorizado
- Actualización, almacenamiento y mantenimiento de los datos del paciente
- Copia de seguridad de los datos
- Protocolo en caso de pérdida de datos

El siguiente paso será realizar 3 nombramientos:

1. **Delegado de protección de datos:** Esta persona deberá de tener la formación y experiencia suficientes en relación a la RGPD y su implementación correcta.
2. **Delegado responsable de seguridad:** Esta persona deberá de tener la formación técnica necesaria para poder llevar a cabo las medidas requeridas por el delegado de protección de datos.
3. **Delegado técnico sanitario:** Dado que los dos anteriores delegados poseen formación en ámbito legal y técnico, falta una figura del ámbito sanitario para que indique la información cómo interpretarla debidamente.

Una vez que dichos delegados y los trabajadores asociados a la organización tengan clara la situación a afrontar, se firmará un documento que fijará la gestión de la privacidad de la información del paciente. Los aspectos a incluir en dicho documento son:

- **Ámbito de aplicación**
- **Acreditación del consentimiento**
- **Datos que se protegerán**
- **Registro de actividad**
- **Derechos del paciente, personal sanitario y organización en el uso de dichos datos**
- **Mecanismos utilizados**

4.2. Fase 2 - Normativa y registro de actividades de tratamiento

Esta fase se centra en el aspecto legislativo y de tratamiento. Es por ello, que se definirán 3 tipos de normativas:

1. **Normativa general:** la que dicta la Unión Europea respecto el tratamiento de datos
2. **Normativa específica:** en base a los Real Decreto aplicables
3. **Normativa interna:** propia del centro

La definición de estos 3 tipos de normativas permitirá conocer el alcance del desarrollo de las actividades y tratamiento con los datos.

El hecho de incumplir estas normativas, puede acarrear sanciones de tipo económico y/o laboral, a estipular según los estamentos asociados y el propio centro si es el caso.

A continuación se gestionará el registro de las actividades de tratamiento. En ella, se especificará el tipo de material (imágenes, vídeos, muestras biológicas, etc.) y se adjuntará con ello:

- **Responsable del tratamiento** (personal sanitario con su respectivo identificador)
- **Finalidad del tratamiento** (por ejemplo analíticas, revisiones, etc.)
- **Categoría de los datos**
- **Plazo de vigencia y existencia de dichos datos**
- **Medidas de seguridad aplicadas a dichos datos**

4.3. Fase 3 - Evaluación de riesgos y medidas de seguridad

En base a la herramienta proporcionada proporcionada por la AEPD (Asociación Española de Protección de Datos) la cual es llamada EIPD (Evaluaciones de impacto de protección de datos) se realizarán varios análisis de riesgo así como evaluaciones de impacto respecto los datos tratados que sean considerados especialmente protegidos.

Un caso es el ejemplo práctico mostrado anteriormente en este mismo documento, en el cuál hay varios datos, pero uno en especial, el de inmunidad, es vital mantenerlo a salvo.

En esta fase se considerarán los casos en los que:

- Hayan filtraciones de información biológica o personal
- Violaciones contractuales por parte del paciente, el personal sanitario o la organización
- Error en la información proporcionada

En estos casos, se definirá la probabilidad de que sucedan estos casos así como el impacto asociado.

Una vez identificados los posibles casos, se creará una matriz de riesgos cómo resultado de la evaluación de riesgos. Esta matriz será una herramienta visual muy práctica para identificar rápidamente la situación. Para hacerlo lo más clara posible, se usará una gama cromática que comprenderá el verde, amarillo, naranja y rojo. Así mismo se muestra también un resumen del funcionamiento de la EIPD.

		PROBABILIDAD				
		Raro	Poco probable	Posible	Muy probable	Casi seguro
Despreciable	Bajo	Bajo	Bajo	Bajo	Medio	Medio
Menores	Bajo	Bajo	Bajo	Medio	Medio	Medio
Moderadas	Medio	Medio	Medio	Medio	Alto	Alto
Mayores	Medio	Medio	Medio	Alto	Alto	Muy alto
Catastróficas	Medio	Alto	Alto	Alto	Muy alto	Muy alto

Figura 42: Mapa de riesgo

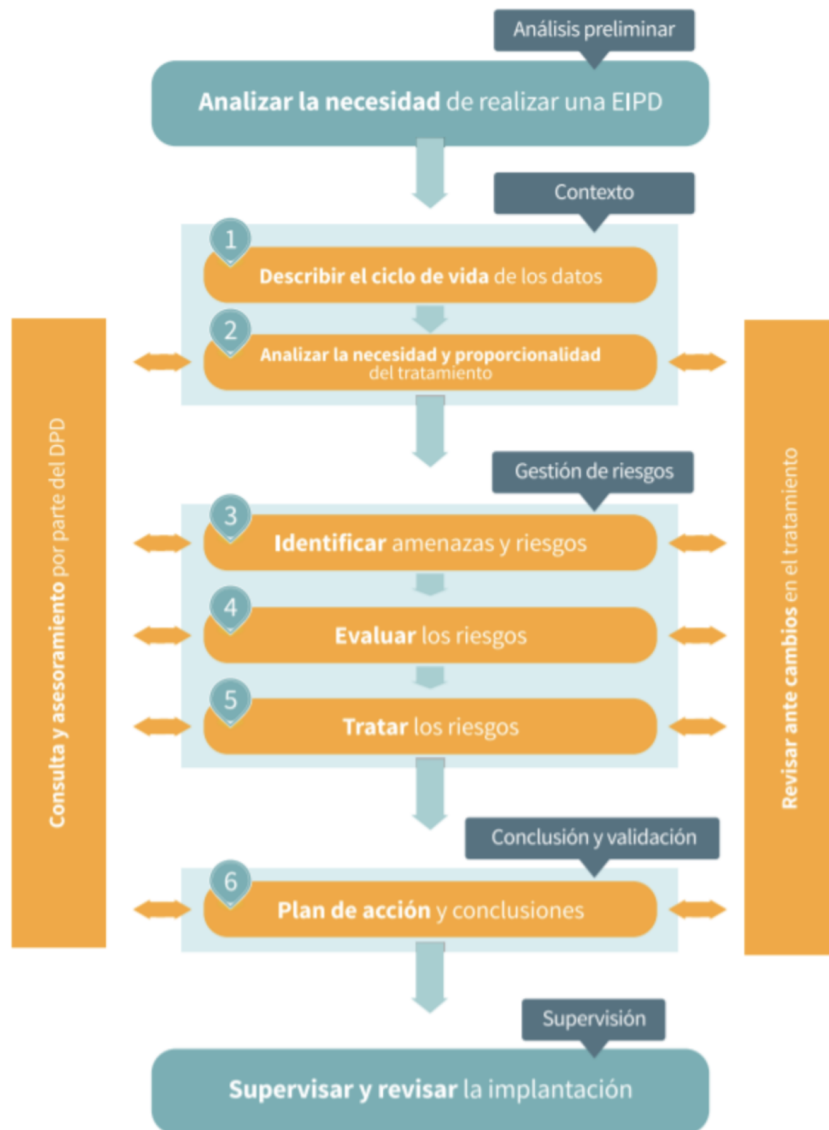


Figura 43: Funcionamiento EIPD

4.4. Fase 4 - Medidas técnicas y de seguridad

Esta fase define el marco legal, ámbito de aplicación, términos, políticas y controles.

4.5. Fase 5 - Denuncias

Esta fase está orientada a dar cobertura a cualquier conducta o situación contraria a lo que se estipula en la normativa establecida anteriormente. Los puntos principales sobre los cuales se fundamenta son:

1. Prevenir conductas, ya sean de pacientes, trabajadores o de la propia organización que no sean acordes a la normativa establecida
2. Identificar debidamente a los responsables de las conductas erróneas, ilícitas o incorrectas.
3. Promover el modelo establecido en la normativa así cómo intentar mejorarlo

De igual manera, esta fase también estará destinada a mejorar el sistema y procedimientos, únicamente cuándo estas mejoras representen objetivamente una clara ventaja para el paciente, trabajador y organización.

5. Conclusiones

En este apartado se tratan las conclusiones obtenidas según el apartado asociado. Así mismo, al final de todo, se indican las líneas de trabajo futuras que potencialmente se desearían haber implementado.

5.1. Apartado teórico

En el aspecto teórico se ha tratado la definición de privacidad diferencial, así como diferentes técnicas criptográficas las cuales permiten asegurar la información en una comunicación. Entre los tipos de privacidad diferencial, se han tratado el tipo local y la global, las cuales han demostrado tener enfoques muy diferentes entre sí.

Dado el trato que hacen a los datos estos dos tipos de privacidad diferencial, se concluye que la privacidad diferencial local es más adecuada para datos que residen en el propio dispositivo (p.e. los datos biológicos de un reloj inteligente), mientras que la global es altamente recomendable para trabajos donde la precisión es requisito indispensable (p.e. cálculo distribuido).

Así mismo, las 3 compañías mencionadas (Apple, Google y Microsoft) aunque difieren en el camino que emprenden para proporcionar al usuario y empresa un producto que implemente privacidad diferencial, la conclusión obtenida de ello es que los tres enfoques son completamente válidos ya que se centran en 3 puntos diferentes pero de interés común, los cuales son usuario, empresa e investigación.

5.2. Apartado práctico

El apartado práctico ha consistido en analizar la adición de ruido en un archivo CSV mediante la función Jitter la cual proporciona ruido aleatorio y la función Gaussiana.

Las conclusiones han sido directas, mientras que con la adición de ruido Gaussiano se puede aplicar privacidad diferencial, en el caso del ruido aleatorio no se pueden obtener buenos resultados ya que el conjunto de datos pierde utilidad, dado que el campo de inmunidad, el cual es donde vira todo el apartado práctico, queda inservible, ya que no se puede recuperar la información.

En el caso de adición de ruido Gaussiano se ha podido comprobar cómo efectivamente se puede implementar privacidad diferencial, no obstante, hace falta crear un mecanismo de recuperación válido y específico para ello. El

resultado ha sido satisfactorio. Como opción o mejora, se plantea el hecho de incluir el valor del eje de ordenadas también en el documento.

5.3. Apartado SI

En este bloque se ha definido un marco de trabajo oficial con el cual se puede poner en marcha un sistema de información. Se han definido 5 fases, las cuales se basan en la protección de datos, normativa, evaluación de riesgos, medidas técnicas y canal de denuncias.

Si bien es cierto que se ha enfocado en el ámbito sanitario, principalmente focalizado en la protección de datos del paciente, este apartado se puede extrapolar e implementar en otros sectores como por ejemplo el comercial (respetando los datos de las partes compradora y vendedora) o de investigación.

5.4. Conclusión global

Este documento por lo tanto ha ofrecido, desde tres puntos de vista diferentes (teórico, práctico y de SI), una propuesta que da solución al problema en la comunicación entre paciente y personal sanitario.

- Mediante el apartado teórico se ha ofrecido un punto de partida dónde seleccionar el tipo de privacidad diferencial más idónea según el caso deseado.
- En el apartado práctico se ha puesto en práctica la viabilidad de las alternativas que se habían valorado desde el principio de este proyecto
- En el SI definido, se ha indicado el procedimiento a aplicar a nivel de servicios y normativa

El seguimiento y la planificación del proyecto han sido satisfactorios y la comunicación entre el autor de este documento y el profesor ha sido correcta. Mediante la evaluación del profesor en cada PAC, se han ido perfilando las carencias del documento y proyecto. El principal punto de mejora tratado ha sido la inclusión de definiciones y glosario para el correcto seguimiento de la explicación del documento.

5.5. Líneas de trabajo futuras

Las principales líneas de trabajo futuras son:

- La creación de una aplicación que simule la aplicación de privacidad diferencial (automatizar el proceso de la parte práctica)
- El diseño y edición de documentación en el apartado SI que sirva como plantilla para la puesta en marcha del mismo
- La consideración de más métodos de aplicación de privacidad diferencial además de las funciones aleatorias, Gauss y Laplace

6. Glosario

- **AEPD (Asociación Española de Protección de Datos):** entidad que vela por el cumplimiento de la ley orgánica de protección de datos (LOPD)
- **Conexión REST:** Arquitectura que permite estandarizar la comunicación entre sistemas de computación web
- **Cruz Roja:** Institución humanitaria que trabaja para prestar servicios de sangre así como a las víctimas de violencia armada
- **CSV:** Documento de texto cuyo contenido está separado por comas y para diferenciar los registros, se utiliza un salto de línea. Es un formato ampliamente usado en estadística
- **EIPD:** evaluación de impacto de protección de datos. Estudio relacionado con la seguridad y privacidad de los datos.
- **INE:** Organismo oficial español el cual tiene como finalidad recopilar estadísticas de tipo demográficas, económicas y sociales
- **MPC (Multi-party computation):** Técnica criptográfica que permite a un colectivo de usuarios o máquinas realizar un cálculo sin que ninguna de las partes implicadas conozcan los datos privados de los participantes
- **Padding:** Método aplicado sobretodo en criptografía para introducir información irrelevante con un objetivo definido
- **RGPD (Reglamento General de Protección de Datos):** conjunto de norma y leyes de ámbito europeo con el objetivo de unificar criterios y la protección de datos de los ciudadanos europeos
- **SARS-CoV-2:** Virus que provoca una enfermedad respiratoria llamada enfermedad por coronavirus de 2019 (COVID-19)
- **SHA256:** Función hash criptográfica la cual permite generar identificadores únicos a partir de la información proporcionada. Utilizada ampliamente por su buen compromiso entre rendimiento y seguridad
- **Transformada Hadamard:** Función matemática basada en el desarrollo en serie de funciones básicas cuyo resultado es -1 o $+1$

7. Bibliografía

Referencias

- [1] Harvard University Privacy Tools Project [En línea] [consulta: 20/12/2021] Disponible en <https://privacytools.seas.harvard.edu/differential-privacy>
- [2] Sartor, Nicolas. 2019. Privacy Budget – Don’t You Spend It All At Once! [imagen web] [consulta: 24/09/2021]. Disponible en: <https://aircloak.com/explaining-differential-privacy/>
- [3] Fu, Silveri. 2017. Differential Privacy [en línea] [consulta 27/09/2021] Disponible en <https://inst.eecs.berkeley.edu/~cs261/fa17/subscribe/diffprivacy.pdf>
- [4] Sainz, Nerea; Cuadrado Saez, Jorge; Ignacio Escribano, Jose. 2021. Privacidad en machine learning: un concepto ahora más formal y tangible [imagen web] [consulta: 16/10/2021]. Disponible en: <https://www.bbvanexttechnologies.com/pills/privacidad-en-machine-learning-un-concepto-ahora-mas-formal-y-tangible>
- [5] Miller, Chance. 2019. Ahead of CES, Apple touts ‘what happens on your iPhone, stays on your iPhone’ with privacy billboard in Las Vegas [imagen web]. [consulta: 23/12/2021]. Disponible en: https://i2.wp.com/9to5mac.com/wp-content/uploads/sites/6/2019/01/DwGoq2uV4AA_Aov.jpg-large.jpeg?w=2500&quality=82&strip=all&ssl=1
- [6] Apple (Privacidad) [en línea]. [Consulta 10/10/2021]. Disponible en: <https://www.apple.com/la/privacy/control/>
- [7] Apple (Privacidad diferencial) [en línea] [Consulta 10/10/2021] Disponible en: https://www.apple.com/privacy/docs/Differential_Privacy_Overview.pdf
- [8] ¿Qué es la privacidad diferencial en el aprendizaje automático (versión preliminar?) [en línea] [Consulta 11/10/2021] Disponible en: <https://docs.microsoft.com/es-es/azure/machine-learning/concept-differential-privacy>
- [9] Cómo Google anonimiza los datos [en línea] [consulta 11/10/2021] Disponible en <https://policies.google.com/technologies/anonymization?hl=es>
- [10] How we’re helping developers with differential privacy [en línea] [consulta 12/10/2021] Disponible en: <https://developers.googleblog.com/2021/01/how-were-helping-developers-with-differential-privacy.html>

- [11] Khalil, Ibrahim; Camtepe Seyit; Bertok, P; Liu, D; Chamikara, M.A.P. 2019. Global vs. Local differential privacy [imagen web] Consulta [22/10/2021]. Disponible en: <https://www.researchgate.net/profile/Map-Chamikara/publication/335095429/figure/fig1/AS:790188276056065@1565406974646/Global-vs-Local-differential-privacy.png>
- [12] Conexión a los Servicios de datos de REST [en línea] [consulta 26/12/2021] Disponible en: https://help.highbond.com/helpdocs/analytics/142/user-guide/es/Content/analytics/defining_importing_data/data_access_window/connecting_to_rest.htm
- [13] Gaussian Distribution and Maximum Likelihood Estimate Method (Step-by-Step) [en línea] [consulta 8/10/2021] Disponible en: <https://medium.com/swlh/gaussian-distribution-and-maximum-likelihood-estimate-method-step-by-step-e4f6014fa83e>
- [14] The Laplace distribution and financial returns [en línea] [consulta 8/10/2021] Disponible en: <https://businessforecastblog.com/the-laplace-distribution-and-financial-returns/>
- [15] The magic of Gaussian noise [en línea] [consulta 13/10/2021] Disponible en: <https://desfontain.es/privacy/gaussian-noise.html>
- [16] INE [en línea][consulta 24/10/2021] Disponible en: <https://www.ine.es>
- [17] Cruz Roja, donar sangre [en línea] [consulta 30/10/2021] Disponible en: <https://www.donarsangre.org>
- [18] El asma podría proteger frente a la COVID-19 en un grupo concreto de pacientes asmáticos [en línea] [consulta 05/10/2021] Disponible en: <https://www.vallhebron.com/es/noticias/el-asma-podria-proteger-frente-la-covid-19-en-un-grupo-concreto-de-pacientes-asmaticos>
- [19] Definición SARS-CoV-2 [en línea] [consulta 15/12/2021] Disponible en: <https://www.cancer.gov/espanol/publicaciones/diccionarios/diccionario-cancer/def/sars-cov-2>
- [20] Documentación Overleaf [en línea] [consulta 1/10/2021] Disponible en: <https://www.overleaf.com/learn>
- [21] Pasos prácticos para la implementación de un sistema de gestión en privacidad de la información[En línea] [Consulta: 11/10/2021] Disponible en: <https://www.ccelpa.org/wp-content/uploads/2019/12/Pasos-pr%C3%A1cticos-para-la-implementaci%C3%B3n-de-un-Sistema-de-Gesti%C3%B3n-en-Privacidad-de-la-Informaci%C3%B3n.pdf>

- [22] Sweeney, Latanya. Matching Known Patients to Health Records in Washington State Data. [en línea] [Consulta 22/09/2021]. Disponible en: <https://privacytools.seas.harvard.edu/files/privacytools/files/1089-1.pdf>

8. Anexos

Este documento no contempla anexos, dado que las tablas de los archivos indicados a lo largo de este trabajo, son de una excesiva extensión y supondría un tamaño no justificado del mismo. No obstante, se indicarán los archivos adjuntos y su finalidad:

- Archivo "Taulas.numbers". Este archivo contiene las tablas creadas en el programa Numbers (archivo en bruto)
- Archivo "Taula.csv". Archivo con el cual se trabaja inicialmente en RStudio
- Archivo "patients.csv". Archivo ya modificado en el propio RStudio el cual ya contempla la privacidad diferencial
- Archivo "id_rsa". Archivo con la clave privada
- Archivo "id_rsa.pub". Archivo con la clave pública (compartida entre usuarios)
- Archivo "patients.csv.encrypt.bin". Archivo CSV encriptado mediante las claves indicadas anteriormente