







Anonimización de datos de salud













Características de los datos en salud



+ Compartir datos de salud puede ser muy beneficioso para evitar sesgos de recolección, sesgo del investigador y generar descubrimientos por fuera de los ensayos clínicos

- La filtración de datos personales en salud puede ser muy perjudicial tanto a nivel individual como a nivel general











Ley 25.326 – Protección de datos personales

- Tipos de datos
- Datos personales: Información de cualquier tipo referida a personas físicas o de existencia ideal determinadas o determinables.
- Datos sensibles: Datos personales que revelan origen racial y étnico, opiniones políticas, convicciones religiosas, filosóficas o morales, afiliación sindical e información referente a la salud o a la vida sexual.
- ARTICULO 7° (Categoría de datos).
- 1. Ninguna persona puede ser obligada a proporcionar datos sensibles.
- 2. Los datos sensibles sólo pueden ser recolectados y objeto de tratamiento cuando medien razones de interés general autorizadas por ley. También podrán ser tratados con finalidades estadísticas o científicas cuando no puedan ser identificados sus titulares.
- 3. Queda prohibida la formación de archivos, bancos o registros que almacenen información que directa o indirectamente revele datos sensibles. Sin perjuicio de ello, la Iglesia Católica, las asociaciones religiosas y las organizaciones políticas y sindicales podrán llevar un registro de sus miembros.
- ARTICULO 8° (Datos relativos a la salud).
- Los establecimientos sanitarios públicos o privados y los profesionales vinculados a las ciencias de la salud pueden recolectar y tratar los datos personales relativos a la salud física ó mental de los pacientes que acudan a los mismos o que estén o hubieren estado bajo tratamiento de aquéllos, respetando los principios del secreto profesional.
- ARTICULO 9° (Seguridad de los datos).
- 1. El responsable o usuario del archivo de datos debe adoptar las medidas técnicas y organizativas que resulten necesarias para garantizar la seguridad y confidencialidad de los datos personales, de modo de evitar su adulteración, pérdida, consulta o tratamiento no autorizado, y que permitan detectar desviaciones, intencionales o no, de información, ya sea que los riesgos provengan de la acción humana o del medio técnico utilizado.
- 2. Queda prohibido registrar datos personales en archivos, registros o bancos que no reúnan condiciones técnicas de integridad y seguridad







Ley 26.529 - Derechos del Paciente en su Relación con los Profesionales e Instituciones de la Salud.



Art 2.

c) **Intimidad.** Toda actividad médico - asistencial tendiente a obtener, clasificar, utilizar, administrar, custodiar y transmitir información y documentación clínica del paciente debe observar el estricto respeto por la dignidad humana y la autonomía de la voluntad, así como el debido **resguardo de la intimidad del mismo y la confidencialidad de sus datos sensibles**, sin perjuicio de las previsiones contenidas en la Ley N° 25.326









Anonimización



• La anonimización es un proceso que depende principalmente de dos factores:

- o La estructura de los datos en cuestión:
 - Cuán identificables son nuestros registros
- o El contexto en el que se comparten esos datos:
 - Qué otras bases de datos disponibles hay
 - Qué capacidad técnica hay
 - Qué interés tiene la información para los demás
 - Qué implicancias tiene perder la anonimización de esa información en particular















Tipos de variables



- <u>Cuasi identificadores:</u> Localidad, Edad, Sexo, Nivel educativo, Situación laboral, etc
 - Combinados estos datos pueden ser utilizados para la identificación de personas. Mientras más única sea la combinación de variables, más riesgo hay de reidentificación.
- <u>Atributos sensibles</u>: Diagnóstico médico, Procedimiento, Salario, Orientación política, etc
 - o Generalmente es la información que presenta interés para otras partes. Muchas veces a mayor diversidad, el valor decrece.











Medición de riesgo



- K-anonimity: Por cada combinación de variables cuasi identificadoras, cuántas tienen menos de k registros
- <u>L-diversity</u>: Por cada combinación de variables cuasi identificadoras, cuántas tienen menos de L cantidad de categorías de variables sensibles.
- <u>T-closeness:</u> La distribución de la variable sensible en cada combinación de variables cuasi identificadoras debe ser t parecida a la distribución en el total general

Software:

- o ARX
- o BIBLIOTECA sdcMicro en R
 - Ambos dan cuenta de cuántos registros tienen riesgo de poder ser identificados una vez suprimidos los identificadores directos









Estrategias de anonimización





- Agregación (ej. convertir edad en una variable de intervalos cada 5 años)
- Muestreo (ej. aclarar que de una base determinada sólo se muestran el 50% de los casos)
- Aleatorización de datos (ej. mezclar los datos de edad y sexo)
- Introducción de "ruido" (ej. sumarle o restarle al año de nacimiento un factor aleatorio conocido o desconocido)









Base de egresos inicial



- Sin datos de fecha de nacimiento, localidad y DNI, la base sigue teniendo un riesgo de más del 90% de casos identificables
- Definir cuál es la información de mayor valor para la publicación (sexo, edad, hospital, diagnóstico por ejemplo) y realizar la evaluación sólo con esas variables
- Introducir alguna estrategia de anonimización para esos datos en particular











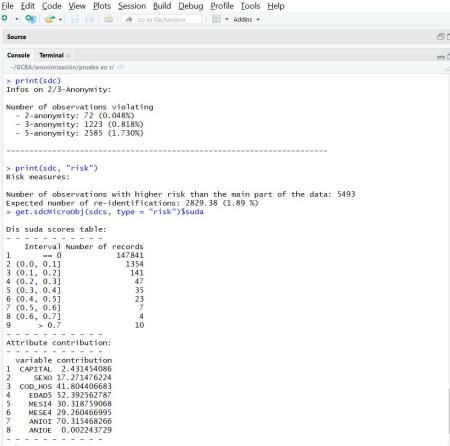
Base de egresos modificada

- <u>Variables cuasi identificadoras:</u> Sexo, Edad en intervalos de 5 años, Código de hospital de internación, Trimestre y Año de ingreso, Trimestre y Año de egreso
- <u>Variables sensibles</u>: Código de diagnóstico, Código de procedimiento, Días de internación
- <u>Variables suprimidas:</u> Fecha de nacimiento, ID, Educación,
 Ocupación, Localidad, Especialidad, todas las relativas a gestación y parto, Segundo diagnóstico o procedimiento









Es posible reducir el riesgo sustancialmente pero es necesario evaluar el valor de la base con la estructura resultante









Procesos de anonimización de Historia Clínica



- Algoritmos que buscan y remueven o alteran nombres, fechas, localidades, DNIs, números de teléfono, nacionalidad, IDs, correo electrónico, direcciones
- Algoritmos que buscan términos médicos y el resto lo dejan como asterisco (text scrubbers)
- Algoritmos de clasificación de frases con información personal (PHI) o no
- En general todo está desarrollado en inglés
- ---- Otro tema es el de la encriptación









HIPAA (Ley de Responsabilidad y Transferibilidad de Seguros Médicos de EEUU)

- Exige excluir de la información de los pacientes:
 - Nombres
 - Toda subdivisión geográfica menor a Estado o los primeros tres dígitos del código postal si la unidad es mayor a 20,000 habitantes
 - o Todas las fechas salvo año
 - Todas las edades arriba de 89 años (se pueden agrupar en 90 y más)
 - Números de teléfono, Correos electrónicos, Número de identificación, Número de histórica clínica, etc
 - Identificadores biométricos y fotos de rostro
 - Cualquier otro identificador que pueda ser utilizado para reconocimiento





