



## **Análisis del Contexto y Normatividad**

**Materia:** Inteligencia Artificial Avanzada  
Para la Ciencia de Datos I

**Actividad:** Portafolio Análisis

Axel Amos Hernández Cardenas - A00829837

## 1. Normativa Asociada al Tipo de Datos Utilizados

Para el entregable 2 del módulo de machine learning, se utilizó el Pima Indians Diabetes Dataset, disponible bajo la licencia de Creative Commons Zero (CC0 1.0), que permite el uso libre de los datos, aspecto que se abordará en la segunda sección de este documento. (UCI, 2016) Sin embargo, el uso de datos de salud conlleva la necesidad de cumplir con normativas que buscan proteger la privacidad y seguridad de la información, aún y cuando esta información se encuentra bajo una licencia abierta. Por lo anterior, es importante comprender las normativas o leyes vigentes que protegen precisamente este tipo de información.

En Estados Unidos, la Health Insurance Portability and Accountability Act (HIPAA) regula el manejo de la información médica de los pacientes, establece controles estrictos sobre la recopilación, uso y el intercambio de los datos. De acuerdo con la normativa, los datos de salud solo pueden ser utilizados con el consentimiento del paciente o institución, y deben ser tratados con confidencialidad. (U.S HHS, 2022) Además, HIPAA ve por la implementación de medidas de seguridad para evitar el mal uso de los datos. (INAI, S.F)

En México, la Ley Federal de Protección de Datos Personales en Posesión de los Particulares (LFPDPPP) también regula el tratamiento de información sensible, incluida la relacionada con la salud. Esta ley establece que los datos de salud son considerados información sensible y, por lo tanto, requieren un tratamiento especial. Al igual que la normativa HIPAA, se exige un consentimiento explícito del titular de los datos para su uso y se deben adoptar medidas adicionales para garantizar su confidencialidad y evitar su uso inadecuado. (INAI, 2016)

Aunque la licencia Creative Commons Zero otorga flexibilidad para el uso libre del dataset, es de gran importancia seguir cumplir con las normativas aplicables al dataset, ya que de esta manera se garantiza que el manejo de datos sea ético y responsable asegurando que, aunque los datos estén anonimizados, la privacidad de los individuos representados en ellos siga siendo respetada,

## 2. Uso de los datos

Para el uso de los datos del dataset, se utilizó la plataforma de Kaggle. Como ya se mencionó, el dataset se encuentra bajo licencia de Creative Commons Zero (CC0 1.0), lo cual otorga total libertad para usar, modificar y distribuir los datos sin restricciones. (Creative Commons, 2024) Bajo esta licencia, se utilizó el dataset de la siguiente manera:

- 1. Para copiar, modificar e interpretar el dataset con propósitos educativos:** Se extrajeron y manipularon los datos con el fin de utilizar técnicas de machine learning para el entrenamiento de modelos.
- 2. Para construir un modelo de red neuronal eficiente sobre la información contenida en el dataset:** Los datos fueron procesados y utilizados para entrenar un modelo de red neuronal que busca predecir la aparición de diabetes basándose en la información proporcionada.

3. **Para distribuir contenido con fines educativos y científicos relacionados al punto anterior:** Los resultados del análisis y entrenamiento se distribuyeron en un documento educativo para su posterior revisión y evaluación.

De igual manera, los datos se usaron bajo las normativas descritas en la sección 1 asegurando el cumplimiento de:

- **Confidencialidad:** Tanto HIPAA como la LFPDPPP exigen el respeto a la confidencialidad de los datos, garantizando que los individuos no puedan ser identificados. Desde el principio, los datos utilizados estaban anonimizados, lo que significa que la privacidad de los individuos está completamente salvaguardada. (U.S HHS, 2022) (INAI, 2016)
- **Consentimiento:** La licencia CC0 1.0 asegura que los datos fueron proporcionados con consentimiento y sin necesidad de atribución, permitiendo su uso sin restricciones adicionales. (Creative Commons, 2024) Lo anterior está alineado con los términos de ambas normativas, al exigir que los datos solo se utilicen con autorización. (U.S HHS, 2022) (INAI, 2016)

### 3. La Red Neuronal y los Datos

La red neuronal construida sigue las consideraciones en la licencia y normativas anteriormente mencionadas de la siguiente manera:

- **Confidencialidad:** Naturalmente, la red neuronal se entrenó con datos anonimizados, lo que asegura que en ningún momento se trate información identificable de un individuo, lo cual está alineado con las normativas de HIPAA y la LFPDPPP.
- **Seguridad:** El modelo fue desarrollado y probado en un entorno local seguro, impidiendo el acceso no autorizado a los datos durante su desarrollo. Aunque estos siguen siendo libres de uso (por CC0 1.0), cualquier modificación sin conocimiento previo podría comprometer la integridad de los mismos. Por lo anterior, se cumplió con los requisitos mínimos de seguridad para ambas normativas.

Además de lo anterior, también es importante definir como la improbabilidad de la red neuronal a incurrir en sesgos éticos.

- **Selección de Features:** Aunque no se contaba con features categóricos, el uso único de valores numéricos asegura que el modelo no incorpore sesgos relacionados con la etnicidad de la comunidad Pima.
- **Validación del Modelo:** Sin embargo, aun y solamente utilizando variables numéricas, es necesario comprobar que el modelo no haga overfitting para no favorecer predicciones sesgadas hacia los individuos cuyas características se repiten más y, por ende, sería injusto además de éticamente incorrecto señalar en gran medida a esos individuos de la comunidad.

## 4. Escenarios de Falta Ética

Al desarrollar una herramienta de machine learning basada en datos de salud, es importante considerar los riesgos éticos que podrían surgir si se utiliza de manera inapropiada. Dado que el análisis de datos de salud implica el manejo de información sensible, se debe corroborar que el este genere resultados precisos y que se pueda utilizar bajo un buen estándar ético.

### 4.1. Malicia

Algunos escenarios donde la solución podría utilizarse con fines malintencionados se describen a continuación.

- **Discriminación:** El modelo podría ser utilizado para formar estereotipos discriminatorios hacia la comunidad indígena Pima al sugerir que la diabetes es una enfermedad muy presente entre ellos debido a sus características, lo que podría llevar a negación de servicios de salud debido a la alta presencia de la condición.
- **Manipulación de Datos:** Modificar las predicciones o datos del modelo con el fin de obtener un beneficio propio sería malicioso. Por ejemplo, modificar los datos para exagerar el riesgo de diabetes con el fin de vender tratamientos o seguros sería malicioso, además de fraudulento y dañino para las personas.

### 4.2. Negligencia

Así mismo, se debe ser consciente del uso negligente de la herramienta, ya que, incluso sin malicia, las decisiones basadas en predicciones incorrectas o incompletas pueden ser graves. A continuación se describen algunos ejemplos de un uso negligente del modelo creado.

- **Mala Interpretación de Resultados:** Si los usuarios de la herramienta no toman en cuenta las limitaciones del dataset o del modelo, pueden generar conclusiones incorrectas que afecten el diagnóstico o tratamiento de pacientes. Por ejemplo, utilizar los datos de Pima, una población específica, a otra población, podría llevar a tratamientos incorrectos, ya que las características entre los dos grupos de personas no necesariamente son iguales o similares.
- **Modelo Desactualizado:** Los modelos deben actualizarse para mantenerse eficientes. Un acto negligente implicaría no mantener el modelo actualizado con nuevos datos ya que si la red neuronal se basa solamente en información obsoleta, las predicciones generadas pueden ser inexactas, ocasionando que se generen tratamientos inadecuados que empeoren la situación de un individuo de la comunidad.

## 5. Referencias

Creative Commons. (2024). *CC0 1.0 Universal Código Legal*. CC.

<https://creativecommons.org/publicdomain/zero/1.0/legalcode.en>

Creative Commons. (2024). *CC0 1.0 Universal - Deed*. CC.

<https://creativecommons.org/publicdomain/zero/1.0/>

INAI. (S.F) 4.20. *HIPAA Health Insurance Portability and Accountability Act*. GobMX.

[https://home.inai.org.mx/wp-content/documentos/DocumentosSectorPrivado/4\\_20\\_HIPAA.pdf](https://home.inai.org.mx/wp-content/documentos/DocumentosSectorPrivado/4_20_HIPAA.pdf)

INAI. (2016). *Introducción a la Ley Federal de Protección de Datos Personales de los Particulares*. GobMX.

<https://inicio.inai.org.mx/CalendarioCapacitacion/Manual%20ILFPDPPP.pdf>

UCI Machine Learning. (2016). *Pima Indians Diabetes Database*. Kaggle.

<https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>

U.S Department of Health and Human Services. (2022). *Summary of the HIPAA Security Rule*. HHS.

<https://www.hhs.gov/hipaa/for-professionals/security/laws-regulations/index.html>