

TRABAJO FIN DE MÁSTER

Máster en Ciencia de Datos e Ingeniería de Datos en la Nube

WiDS Dathaton 2024 - Challenge 2: Modelos de regresión para estimación del periodo de diagnóstico metastático

Autor: Luna Jiménez Fernández

Tutor: Juan Carlos Alfaro Jiménez

Junio, 2025

*Dedicado a la gente que, pese a todo,
sigue persiguiendo sus sueños.
Nunca os rindáis.*

Declaración de autoría

Yo, **Luna Jiménez Fernández**, con DNI **47092045M**, declaro que soy la única autora del Trabajo Fin de Master titulado ***“WiDS Dathon 2024 - Challenge 2: Modelos de regresión para estimación del periodo de diagnóstico metastático”***, que el citado trabajo no infringe las leyes en vigor sobre propiedad intelectual, y que todo el material no original contenido en dicho trabajo está apropiadamente atribuido a sus legítimos autores.

Albacete, a ... de **Junio de 2025**

Fdo.: **Luna Jiménez Fernández**

Resumen

TODO RESUMEN AQUI

Abstract

TODO ABSTRACT HERE

Agradecimientos

En primer lugar, quiero agradecer a todos mis compañeros y amigos del grupo de **Sistemas Informáticos y Minería de Datos (SIMD)** - y, especialmente, a mi amigo y director **Juan Carlos Alfaro Jiménez** - por su apoyo, recursos y consejos durante la realización de este trabajo. Aunque ya no sea formalmente parte de este grupo, siempre me sentiré vinculada a él.

Además, quiero agradecer a mis amigos y familia del **Curso de Comic Online de la Escola Joso - Arai, Aina, Arkaitz, Clara, Irene, Martín, Pau, Rafi...** -, con los que compartí un proyecto de gran importancia personal, mi primer comic publicado, y en los que he encontrado un grupo al que pertenecer. Muchas gracias por todo.

Finalmente, quiero agradecer a **mi familia y seres queridos** - tanto los que me acompañan presencialmente como los que se encuentran a distancia. Vuestro apoyo y cariño continuo me ha ayudado a seguir adelante y acabar este trabajo a pesar de todas las dificultades.

Índice general

1	Introducción	1
1.1	Objetivos	1
1.2	Estructura de la memoria	2
2	Revisión de técnicas	3
2.1	Ciencia de datos y el ciclo de vida de los datos	3
2.1.1	Ciencia de datos	3
2.1.2	Cross-Industry Standard Process for Data Mining - CRISP-DM.	5
2.2	Aprendizaje automático y ajuste de modelos	6
2.2.1	Aprendizaje automático	6
2.2.2	Modelos de regresión	7
2.2.3	Selección de modelos: ajuste de hiperparámetros y validación cruzada.	9
3	Estudio del problema	11
3.1	Definición del problema.	11
3.1.1	Atributos del problema	11
3.2	Análisis exploratorio de datos	11
3.2.1	Variable objetivo - distribución y comportamiento.	11
3.2.2	Valores perdidos	11
3.2.3	Atributos categóricos	11
3.2.4	Atributos numéricos	11
3.2.5	Variables geográficas, sociales y económicas	11
4	Preprocesamiento del conjunto de datos.	13
4.1	Selección de atributos	13

4.2	Procesamiento de los datos	13
5	Modelado y experimentación	15
5.1	Selección de modelos	15
5.2	Experimentación	15
5.2.1	<i>Ajuste de hiperparámetros y selección de subconjuntos de atributos</i>	<i>15</i>
5.2.2	<i>Validación y selección de modelo final</i>	<i>15</i>
5.3	Análisis de resultados	15
5.3.1	<i>Rendimiento de los subconjuntos de hiperparámetros</i>	<i>15</i>
5.3.2	<i>Rendimiento de los modelos entrenados</i>	<i>15</i>
5.3.3	<i>Rendimiento del modelo final</i>	<i>15</i>
6	Aplicación web	17
6.1	Aplicación para usuario - predicción individual	17
6.2	Aplicación <i>batch</i> - predicción en grupo	17
7	Conclusiones	19
7.1	Trabajo futuro	19
	Referencia bibliográfica	22
A	Anexo 1	23

Índice de figuras

2.1	Ciclo de vida de los datos [4]	4
2.2	Ciclo de vida de CRISP-DM [8]	5

Índice de tablas

1. Introducción

El acceso equitativo a una **atención sanitaria de calidad** es un problema de gran interés a nivel global, existiendo desigualdades sustanciales en la **calidad y acceso** a dicho servicio entre distintas poblaciones. Estos problemas, además, se pueden llegar a exacerbar por distintos factores: geográficos, socioeconómicos y climáticos.

Con el fin de estudiar la influencia de dichos factores en la atención sanitaria, la iniciativa *Women in Data Science* propuso en el año 2024 una **competición** [1] con el objetivo de **estimar el tiempo necesario para realizar un diagnóstico de metástasis para cáncer de mama** a partir de un conjunto de datos médico ampliado con información geográfica, socioeconómica y climática - y, a su vez, estudiar como dichos factores pueden influir al tiempo necesario para realizar un diagnóstico.

Por tanto, la meta de este trabajo es la creación de **modelos de regresión** capaces de estimar dicho tiempo de diagnóstico con el menor error posible - utilizando, para ello, el proceso completo de **ciencia de datos**.

1.1. Objetivos

El principal objetivo de este trabajo es el **desarrollo de un modelo de regresión** capaz de resolver el problema propuesto por la competición: la predicción del tiempo necesario para realizar un diagnóstico de metástasis de cáncer de mama, evaluando su rendimiento y dejando disponible el modelo para ser accesible por los hipotéticos usuarios finales.

Para alcanzar dicho objetivo, es necesario llevar a cabo los siguientes pasos, siguiendo el **ciclo de vida de la ciencia de datos**:

1. Análisis exploratorio de los datos disponibles en la competición, para comprender su comportamiento y características.
2. Pre-procesamiento de los datos para la propuesta de subconjuntos de atributos reducidos y preparación posterior para el uso con modelos.
3. Estudio, selección y caracterización de los modelos y sus hiperparámetros a estudiar durante el proceso.

-
4. Experimentación y estudio de los resultados para seleccionar un modelo definitivo a ser utilizado.
 5. Creación de una aplicación web para desplegar el modelo final entrenado, con el fin de ser utilizado por expertos en el campo de la medicina sin experiencia previa en ciencia de datos.

A su vez, este trabajo aborda el segundo objetivo planteado por la propia competición: el **estudio de la influencia de los factores geográficos, socioeconómicos y climáticos** en la calidad de la atención sanitaria.

1.2. Estructura de la memoria

La memoria está dividida en un total de **7** capítulos, como se describen a continuación:

- **Capítulo 1:** En este capítulo se introduce el problema a resolver, los objetivos que se busca cumplir con el trabajo y la estructura general de la memoria.
- **Capítulo 2:** En este capítulo se realiza una breve revisión de las principales técnicas a utilizar durante la memoria: tanto el proceso de ciencia de datos y sus etapas como los modelos a utilizar durante la experimentación - desde los modelos simples como las regresiones lineales y los árboles de decisiones hasta los *ensembles* de modelos simples.
- **Capítulo 3:** En este capítulo se realiza un estudio más exhaustivo del problema: tanto su definición como un análisis exploratorio de los datos disponibles, estudiando el comportamiento de la variable objetivo y la relevancia y correlación de los atributos respecto al tiempo de diagnóstico.
- **Capítulo 4:** En este capítulo se introduce el pre-procesamiento a realizar sobre el conjunto de datos, obteniendo varios subconjuntos de atributos reducidos a ser estudiado posteriormente y preparando *pipelines* automáticos para realizar todas las transformaciones necesarias para el uso de los datos por parte de los modelos.
- **Capítulo 5:** En este capítulo se detalla la experimentación a realizar. Se proponen varios modelos sobre los que se realizará un proceso de ajuste de hiperparámetros y selección de modelos, con el fin de obtener un modelo definitivo a ser utilizado para resolver el problema. Además, se presentan y estudian los resultados de dicha experimentación.
- **Capítulo 6:** En este capítulo se presenta una aplicación web a través de la cual se hace disponible a los usuarios expertos el modelo obtenido en el capítulo anterior - detallando la interfaz gráfica y las distintas funcionalidades ofrecidas.
- **Capítulo 7:** Finalmente, en este capítulo se muestran las conclusiones alcanzadas tras el desarrollo del trabajo, proponiendo posibles líneas de trabajo futuro para ampliarlo.

2. Revisión de técnicas

En este capítulo se describen los procesos y algoritmos utilizados a lo largo del trabajo descrito en esta memoria. Concretamente, se comienza explicando el concepto de la **ciencia de datos** y su ciclo de vida, haciendo énfasis en **CRISP-DM** como metodología utilizada a lo largo del proyecto para resolver el problema propuesto. Tras esto, se estudian conceptos de aprendizaje automático como los **modelos de regresión** - haciendo especial énfasis en los modelos de *ensemble* basados en técnicas de **Gradient Boosting** - o la búsqueda de hiperparámetros.

2.1. Ciencia de datos y el ciclo de vida de los datos

2.1.1. Ciencia de datos

La **ciencia de datos** es el estudio de la extracción de conocimiento útil a partir de datos, y de la generalización de dicho proceso a cualquier problema [2]. Dicho proceso incluye la recolección y almacenamiento, mantenimiento, procesamiento, análisis y visualización de enormes cantidades de datos heterogéneos - asociados a un gran abanico de aplicaciones y dominios en muchas ocasiones multidisciplinarios [3].

Desde su origen, la ciencia de datos ha evolucionado como un campo interdisciplinar que integra conocimientos y técnicas de otras disciplinas afines como el análisis de datos, la estadística o la minería de datos [4]. Ahora bien, la principal diferencia con estos campos se encuentra en el fin: el aprendizaje a partir de los datos [2] y la capacidad de adquirir nuevo conocimiento capaz de ser utilizado para la toma de decisiones y la predicción [3].

Por definición, la ciencia de datos depende de los datos sobre los que se está trabajando. Por esto, el proceso de trabajo de la ciencia de datos depende generalmente del **ciclo de vida de los datos**: las distintas etapas por las que pasa un conjunto de datos desde su recolección e investigación hasta su uso final [5]. Como se observa en la **Figura 2.1**, este ciclo está tradicionalmente dividido en **cinco** apartados [4]:

1. **Adquisición:** En la actualidad, los datos se generan en cantidades masivas - del orden de **exabytes por hora** [6]. Por tanto, el primer paso del ciclo consiste en la adquisición y almacenamiento eficiente de los datos necesarios para el proceso.

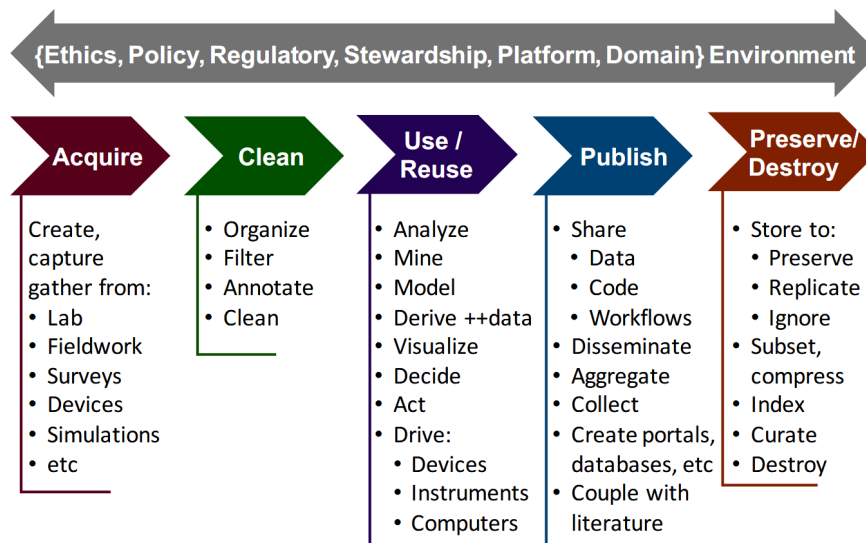


Figura 2.1: Ciclo de vida de los datos [4]

2. **Limpieza:** Tras la adquisición, el segundo paso del ciclo consiste en la transformación de los datos originales en datos utilizables posteriormente - a través de procesos de limpieza, imputación, formateo...
3. **Uso y re-uso:** El tercer paso del ciclo consiste en el uso de los datos procesados con el fin de adquirir conocimiento y tomar decisiones a partir de éstos. Éste apartado se puede dividir, a su vez, en tres subapartados [6]:
 - (a) **Análisis exploratorio:** El estudio del comportamiento de los datos con el fin de plantear hipótesis para guiar el resto del ciclo de datos [7].
 - (b) **Modelado:** El uso de técnicas computacionales y estadísticas para extraer conocimiento y predicciones a partir del conjunto de datos.
 - (c) **Visualización, interpretación y actuación:** La representación gráfica de los resultados del uso de los datos, con el fin de facilitar la toma de decisiones posterior a las personas.
4. **Publicación:** El cuarto paso del ciclo consiste en la disseminación de los resultados del proceso - con el fin de que el conocimiento creado pueda ser conocido y reutilizado por el mayor número de personas posible.
5. **Preservación o destrucción:** El quinto y último paso del ciclo consiste en la preservación o destrucción de los datos utilizados - cumpliendo con otros factores como pueden ser las consideraciones éticas o regulatorias.

Con el fin de regularizar, estandarizar y hacer reproducible el proceso completo de la ciencia de datos - desde la adquisición de los conjuntos de datos hasta la distribución de los resultados -, se han propuesto varias ampliaciones y adaptaciones del ciclo de datos estudiado, conocidas como **ciclos de vida de la ciencia de datos** [5].

Aunque actualmente no existe un ciclo estandarizado, uno de los procesos más utilizados para ciencia de datos es el **Cross-Industry Standard Process for Data Mining (CRISP-DM)**, propuesto originalmente para el campo de la minería de datos pero adaptado a las necesidades de la ciencia de datos [8] - siendo el proceso utilizado a lo largo del trabajo descrito en esta memoria.

2.1.2. Cross-Industry Standard Process for Data Mining - CRISP-DM

Cross-Industry Standard Process for Data Mining (abreviado como *CRISP-DM*) es una metodología desarrollada con el fin de ofrecer un proceso de trabajo completo de principio a fin para la minería de datos; independientemente del campo, las herramientas o la aplicación final de los datos [8]. Si bien fue propuesto originalmente en el año 2000, en la actualidad sigue siendo uno de los procesos más utilizados tanto en minería de datos como en ciencia de datos [9].

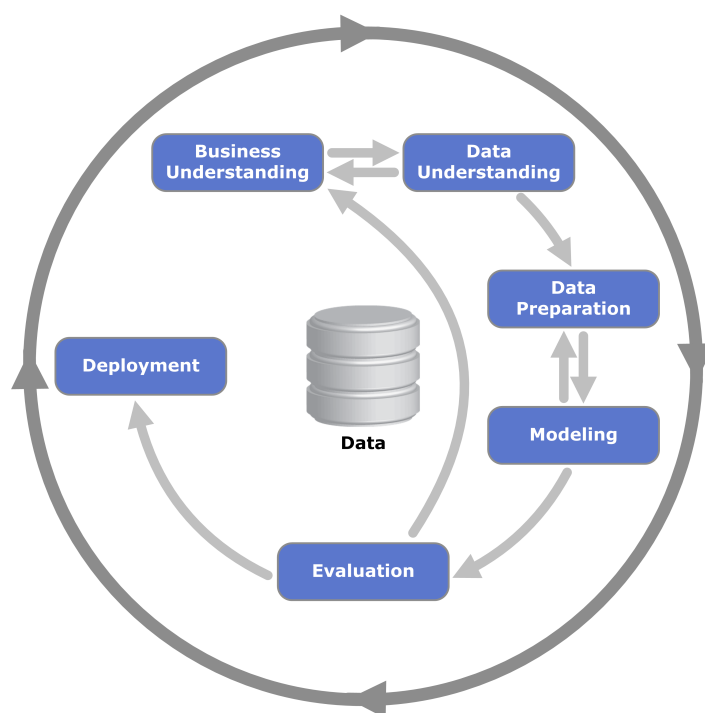


Figura 2.2: Ciclo de vida de CRISP-DM [8]

Como se observa en la **Figura 2.2**, el ciclo de CRISP-DM está dividido en **seis fases** [8], similares al ciclo de datos estudiado:

1. **Conocimiento del campo (*Business Understanding*):** El primer paso del ciclo consiste en entender el problema y los objetivos a resolver - estudiando la situación actual y estableciendo los pasos para alcanzar las metas propuestas.

-
2. **Conocimiento de los datos (*Data Understanding*):** El segundo paso del ciclo consiste en adquirir y estudiar los datos - tanto de forma superficial como en un análisis exploratorio más profundo -, además de verificar que los datos disponibles son útiles para los objetivos propuestos.
 3. **Preparación de los datos (*Data Preparation*):** Tras la adquisición de conocimiento, el tercer paso consiste en preparar los datos obtenidos para su uso posterior - seleccionando las instancias relevantes, limpiando los datos para eliminar valores perdidos, enriqueciendo los datos con información externa...
 4. **Modelado (*Modeling*):** Con los datos preparados, la cuarta fase del ciclo consiste en el uso y calibración de modelos de aprendizaje automático a aplicar sobre los datos - definiendo los estudios y experimentos a realizar sobre los modelos, y evaluando el rendimiento final de éstos.
 5. **Evaluación (*Evaluation*):** Antes de desplegar el modelo final, la quinta fase del ciclo consiste en evaluar si los resultados obtenidos satisfacen los objetivos propuestos y si el proceso de ciencia de datos se ha aplicado de forma adecuada.
 6. **Despliegue (*Deployment*):** La última fase del ciclo es el despliegue y diseminación de los resultados obtenidos - haciendo disponible el modelo y los resultados a los usuarios finales.

Es importante destacar que, como indican las flechas de la **Figura 2.2**, la metodología propuesta no es lineal, sino que el flujo entre los distintos pasos se puede ver alterado:

- Las fases tienen dependencias entre sí - los descubrimientos en algunas fases pueden producir que sea necesario volver a fases anteriores para perfeccionar el proceso.
- El proceso es **cíclico** - los conocimientos adquiridos durante las distintas fases se aplican para refinar futuros procesos, ya sean sobre el mismo conjunto de datos o datos nuevos.

2.2. Aprendizaje automático y ajuste de modelos

2.2.1. Aprendizaje automático

El **aprendizaje automático** (también conocido en inglés como *Machine Learning*) es una rama de la inteligencia artificial que consiste en la creación de programas capaces de **aprender** - es decir, de mejorar su rendimiento en una tarea - a través de la experiencia y de la información que se les aporta [10].

También se puede definir el término como el conjunto de métodos capaces de detectar patrones en los datos de forma autónoma, y de utilizar dichos patrones para predecir datos futuros [11] - siendo este último de mayor interés a los principios de la ciencia de datos.

Generalmente, los algoritmos de aprendizaje automático se dividen en dos grandes familias, en función del tipo de datos e información que se aporta a los algoritmos [12] [11]:

- **Aprendizaje supervisado:** El objetivo del algoritmo es aprender una función capaz de, dados unos datos de entrada X , predecir una salida Y . Esta función se aprende a partir de un conjunto de datos $D = (x_i, y_i)_{i=1}^N$ donde a cada instancia x_i del conjunto de datos D se le asocia un valor esperado y_i .

Este tipo de aprendizaje se puede dividir a su vez en dos categorías dependiendo del tipo de salida Y que se espera [12]:

- **Clasificación:** El algoritmo busca obtener para cada entrada x_i un valor concreto **dentro de un conjunto finito de posibles valores**.
 - **Regresión:** El algoritmo busca obtener, para cada entrada x_i , un **valor numérico continuo**.
- **Aprendizaje no supervisado:** El objetivo del algoritmo es aprender patrones subyacentes de los datos de entrada X ofrecidos, sin buscar predecir una salida. Esta función se aprende a partir de un conjunto de datos $D = x_{i=1}^N$ donde no se ofrece ningún tipo de etiqueta a cada instancia x_i .

De cara a cumplir el objetivo propuesto por el trabajo descrito en esta memoria - la creación de un modelo capaz de **predecir el tiempo de diagnóstico** -, se van a trabajar con modelos supervisados de **regresión**. Por esto, resulta de interés describir los principales modelos a utilizar y el ajuste que se va a realizar sobre ellos.

2.2.2. Modelos de regresión

Un **modelo** es el resultado del proceso de aprendizaje automático: una función capaz de predecir una salida para una entrada dada, y cuyos parámetros han sido ajustados a través de un entrenamiento sobre un conjunto de datos para **minimizar un error** [13].

En el caso de la **regresión**, el objetivo del modelo es aprender una función capaz de predecir un valor numérico continuo para cada instancia de datos de entrada [11]. Dicha función se ajusta buscando encontrar el conjunto de parámetros que **minimiza** la diferencia entre los valores predichos por la función y los valores reales asociados a los datos de entrada [13].

En el caso de los modelos de regresión, el objetivo por tanto es encontrar los parámetros que **reducen al mínimo la diferencia** entre los valores predichos por la función y los valores reales esperados.

Se han propuesto y estudiado un gran número de modelos de regresión, con parametrizaciones y funcionamientos diversos, en la bibliografía [14]. Pese a esta variedad, es posible dividir todas estas familias de modelos en dos grandes grupos: modelos **tradicionales** y modelos de **conjuntos o ensembles** [12]

Modelos tradicionales - regresión lineal, árboles de decisiones y máquinas de vectores de soporte

Si bien no hay una definición consensuada sobre su definición, se puede entender como **modelo tradicional** a un modelo que entrena una única función con el fin de realizar predicciones sobre la salida esperada para cada entrada de datos [12].

Existe una gran cantidad de familias de modelos con una larga trayectoria en la bibliografía existente [13]. Ahora bien, el estudio realizado en la memoria se centra en las siguientes tres familias de modelos utilizadas en el trabajo:

Regresión lineal: Los modelos más simples, trabajando con la suposición de que **existe una correlación lineal** entre los atributos de entrada y la salida del modelo [11]. Por lo general, la salida y para una entrada x se predice utilizando la siguiente fórmula:

$$y(x) = \sum_{j=1}^D w_j x_j$$

Donde x_j representa cada atributo de la entrada y w_j el peso asignado a cada atributo, siendo el objetivo de estos modelos ajustar los pesos asignados a cada atributo para minimizar el error cuadrado [12]. Ahora bien, cuando se trabaja con conjuntos de datos de gran dimensionalidad, el gran número de atributos puede afectar de forma negativa al rendimiento del modelo, causando un sobreajuste al conjunto de entrenamiento [15].

Para evitar este problema, se proponen técnicas de **regularización** - penalizaciones aplicadas a la fórmula del error con el objetivo de conseguir modelos menos complejos y más generalizables [16]. Los tres modelos de regularización más utilizados son los siguientes:

- **Ridge (L2) [17]:** Como factor de penalización, se utiliza $\sum_{j=1}^D (w_j)^2$ - la suma de los pesos cuadrados del modelo, buscando reducir de forma generalizada la influencia de los atributos para evitar sobreajustes y correlaciones.
- **Lasso (L1) [18]:** Como factor de penalización, se utiliza $\sum_{j=1}^D |w_j|$ - la suma del valor absoluto de los pesos del modelo, buscando eliminar los atributos irrelevantes reduciendo su peso a 0.
- **Elastic-Net [16]:** Como factor de penalización, se utiliza $\lambda \left(\sum_{j=1}^D (w_j)^2 \right) + (1-\lambda) \left(\sum_{j=1}^D |w_j| \right)$ - utilizando a la vez las regularizaciones L1 y L2 de forma ponderada, buscando aunar los beneficios de ambas aproximaciones.

Máquinas de vectores de soporte (SVM): Las máquinas de vectores de soporte se pueden entender como una evolución de los modelos de regresión lineal donde, en vez de buscar la línea que mejor se ajusta al conjunto de datos, se busca el **hiperplano** capaz de ajustarse al conjunto de datos con el **mayor margen** [12].

En regresión, esto se traduce en la búsqueda de la función representando al mejor hiperplano capaz de ajustarse a todas las instancias del conjunto de datos a la vez que es capaz de mantener una distancia inferior a un margen ϵ con todos los puntos [19].

La principal utilidad de estos modelos radica en las dos siguientes características [12]:

- **Funciones kernel:** Un problema de los modelos lineales es que los conjuntos de datos no siempre son linealmente separables. Para solventar este problema, las máquinas de vectores de soporte son capaces de utilizar **funciones kernel** para transformar los datos a una mayor dimensionalidad - donde si es posible ajustar un hiperplano con mayor margen.

- **Vectores de soporte:** Para definir el modelo no es necesario almacenar información sobre el conjunto de datos completo, sino que es suficiente con almacenar información sobre los **puntos que definen la frontera entre el hiperplano y el margen** - conocidos como los vectores de soporte.

Árboles de decisión: Los árboles de decisión son modelos de reglas representando su función a través de grafos dirigidos [13] donde la predicción se obtiene realizando una serie de comprobaciones secuenciales empezando desde la raíz, ramificando hasta llegar a una hoja final [12]. Estos árboles se dividen en los siguientes componentes:

- **Nodos:** Nodos internos del árbol donde se realiza una comprobación sobre el valor de un atributo. Dependiendo del resultado de la comprobación, el nodo se **ramifica** a otros nodos u hojas.
- **Hojas:** El valor final predicho para una entrada, alcanzado tras una serie de comprobaciones en nodos.

El objetivo del modelo es, por tanto, aprender el conjunto de reglas que minimiza el error del modelo para el conjunto de datos dado. Ahora bien, estos modelos tienden a **sobreajustar** - intentando representar patrones presentes en el conjunto de datos no representativos de la distribución real - creando árboles de gran profundidad [12].

Modelos de ensemble - bagging y boosting

Como contraste a los modelos tradicionales, un **modelo de conjunto o ensembles** es un modelo que, durante su entrenamiento, aprende un **conjunto de funciones o modelos** - por lo general, una agrupación de modelos tradicionales -, agrupando las predicciones de todos éstos para obtener una predicción general de la salida esperada para cada entrada de datos [12].

Bagging

- **Random Forest:**
- **Extremely Random Trees:**

Boosting

- **Adaptive Boosting:**
- **Extreme Gradient Boosting:**
- **Categorical Boosting:**
- **Light Gradient-Boosting Model:**
- **Histogram-based Gradient Boosting:**

2.2.3. Selección de modelos: ajuste de hiperparámetros y validación cruzada

[11] PARA DEFINICION DE SELECCION

[20] PARA AJUSTE DE HIPERPARAMETROS

3. Estudio del problema

3.1. Definición del problema

3.1.1. Atributos del problema

3.2. Análisis exploratorio de datos

3.2.1. Variable objetivo - distribución y comportamiento

3.2.2. Valores perdidos

3.2.3. Atributos categóricos

3.2.4. Atributos numéricos

3.2.5. Variables geográficas, sociales y económicas

4. Preprocesamiento del conjunto de datos

4.1. Selección de atributos

4.2. Procesamiento de los datos

5. Modelado y experimentación

5.1. Selección de modelos

5.2. Experimentación

5.2.1. Ajuste de hiperparámetros y selección de subconjuntos de atributos

5.2.2. Validación y selección de modelo final

5.3. Análisis de resultados

5.3.1. Rendimiento de los subconjuntos de hiperparámetros

5.3.2. Rendimiento de los modelos entrenados

5.3.3. Rendimiento del modelo final

6. Aplicación web

6.1. Aplicación para usuario - predicción individual

6.2. Aplicación *batch* - predicción en grupo

7. Conclusiones

7.1. Trabajo futuro

Bibliografía

- [1] Women in Data Science, *WiDS Datathon 2024 Challenge 2*, 2024.
- [2] D. D. and, “50 Years of Data Science”, *Journal of Computational and Graphical Statistics*, vol. 26, n° 4, págs. 745-766, 2017. DOI: 10 . 1080 / 10618600 . 2017 . 1384734. eprint: <https://doi.org/10.1080/10618600.2017.1384734>.
- [3] V. Dhar, “Data science and prediction”, *Commun. ACM*, vol. 56, n° 12, págs. 64-73, dic. de 2013, ISSN: 0001-0782. DOI: 10 . 1145/2500499.
- [4] F. Berman et al., “Realizing the potential of data science”, *Commun. ACM*, vol. 61, n° 4, págs. 67-72, mar. de 2018, ISSN: 0001-0782. DOI: 10 . 1145/3188721.
- [5] V. Stodden, “The data science life cycle: a disciplined approach to advancing data science as a science”, *Commun. ACM*, vol. 63, n° 7, págs. 58-66, jun. de 2020, ISSN: 0001-0782. DOI: 10 . 1145/3360646.
- [6] J. M. Wing, “The Data Life Cycle”, *Harvard Data Science Review*, vol. 1, n° 1, jul. de 2019, <https://hdsr.mitpress.mit.edu/pub/577rq08d>.
- [7] M. Komorowski, D. Marshall, J. Saliccioli e Y. Crutain, “Exploratory Data Analysis”, en sep. de 2016, págs. 185-203, ISBN: 978-3-319-43740-8. DOI: 10 . 1007 / 978 - 3 - 319 - 43742 - 2_15.
- [8] C. Shearer, “The CRISP-DM model: the new blueprint for data mining”, *Journal of data warehousing*, vol. 5, n° 4, págs. 13-22, 2000.
- [9] J. Saltz, *CRISP-DM is Still the Most Popular Framework for Executing Data Science Projects - Data Science PM — datascience-pm.com*, <https://www.datascience-pm.com/crisp-dm-still-most-popular/>, [Accessed 28-05-2025], 2024.
- [10] T. M. Mitchell, *Machine learning*. McGraw-hill New York, 1997, vol. 1.
- [11] K. P. Murphy, *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012, ISBN: 0262018020.
- [12] S. Russell y P. Norvig, *Artificial Intelligence: A Modern Approach*, 3rd. USA: Prentice Hall Press, 2009, ISBN: 0136042597.
- [13] A. Burkov, *The Hundred-Page Machine Learning Book*. 2019.

-
- [14] Y. Tai, *A Survey Of Regression Algorithms And Connections With Deep Learning*, 2021. arXiv: 2104.12647 [cs.LG].
- [15] A. Y. Ng, "Feature selection, L1 vs. L2 regularization, and rotational invariance", en *Proceedings of the Twenty-First International Conference on Machine Learning*, ép. ICML '04, Banff, Alberta, Canada: Association for Computing Machinery, 2004, pág. 78, ISBN: 1581138385. DOI: 10.1145/1015330.1015435.
- [16] H. Zou y T. Hastie, "Regularization and Variable Selection Via the Elastic Net", *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 67, n° 2, págs. 301-320, mar. de 2005, ISSN: 1369-7412. DOI: 10.1111/j.1467-9868.2005.00503.x. eprint: https://academic.oup.com/jrsssb/article-pdf/67/2/301/49795094/jrsssb_67_2_301.pdf.
- [17] A. E. Hoerl y R. W. Kennard, "Ridge Regression: Biased Estimation for Nonorthogonal Problems", *Technometrics*, vol. 42, n° 1, págs. 80-86, 2000, ISSN: 00401706.
- [18] R. Tibshirani, "Regression Shrinkage and Selection via the Lasso", *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58, n° 1, págs. 267-288, 1996, ISSN: 00359246.
- [19] A. Smola y B. Schölkopf, "A tutorial on support vector regression", *Statistics and Computing*, vol. 14, págs. 199-222, ago. de 2004. DOI: 10.1023/B%3ASTC0.0000035301.49549.88.
- [20] I. Goodfellow, Y. Bengio y A. Courville, *Deep Learning*. MIT Press, 2016, Book in preparation for MIT Press.

A. Anexo 1
