

# Modelo de Ciencia de Datos para predecir ventas en una Empresa

Luisa López Vázquez¹
lisaamazon7@gmail.com
Universidad Americana de Europa

Rodrigo Cadena Martínez <u>profesor@unade.net</u> Universidad Americana de Europa

### **RESUMEN**

En este trabajo de investigación se realizará un modelo de predicción mediante algoritmos de ciencia de datos, que le permita a una empresa del sector de bienes raíces tener una predicción de sus ventas, tomando como base el historial de años anteriores; lo cual permitirá a los directivos y asesores inmobiliarios conocer los meses en los que se venden más inmuebles. Se trabajará con la metodología CRISP-DM la cual permitirá hacer ajustes fácilmente en los cambios que surgen en el desarrollo del proyecto. Posteriormente se identificarán las fuentes de datos de la inmobiliaria, se realizará un análisis de los mismos para clasificarlos en datos estructurados y no estructurados. Se realizará el diseño de la solución mediante la herramienta Jupyter Notebook utilizando algoritmos de Machine Learning, el tipo de aprendizaje utilizado será supervisado. El modelo de Ciencia de Datos permitirá a la inmobiliaria conocer las operaciones que más ganancias generan, identificar las características de los inmuebles más vendidos, así como tener el pronóstico de sus ventas futuras.

Palabras clave: big data; ciencia de datos; análisis; datos.

Correspondencia <u>lisaamazon7@gmail.com</u>

<sup>&</sup>lt;sup>1</sup> Autor Principal

Data Science Model To Predict Sales In A Real Estate Agency.

**ABSTRACT** 

In this research work, a prediction model will be made through data science algorithms, which allows a

company in the real estate sector to have a prediction of its sales, based on the history of previous years;

which will allow managers and real estate consultants to know the months in which more properties are

sold. The CRISP-DM methodology will be used, which will allow easy adjustments to the changes that

arise in the development of the project. Subsequently, the real estate data sources will be identified, an

analysis will be carried out to classify them into structured and unstructured data. The design of the solution

will be carried out using the Jupyter Notebook tool using Machine Learning algorithms, the type of learning

used will be supervised. The Data Science model will allow the real estate company to know the operations

that generate the most profits, identify the characteristics of the best-selling properties, as well as have the

forecast of their future sales.

Keywords: big data; data science; analysis; data.

Artículo recibido 25 julio 2023

Aceptado para publicación: 25 agosto 2023

pág. 9375

### INTRODUCCIÓN

Actualmente el 86 % de empresas en México aún tiene que progresar en su tecnología y procesos de Big Data que les permitan generar valor y mejorar su toma de decisiones esto se debe principalmente a la recopilación de datos más rápido de lo que puede analizar y usar, el 73% se queja de que tiene un exceso de datos que no puede cumplir con los requisitos de seguridad y cumplimiento, y el 68% dice que sus equipos ya están abrumados por los datos que tienen.

El 14 % de empresas que utilizan modelos de ciencia de daros pertenece a los grandes monopolios en el país, sin embargo, esta tecnología no está alcanzando a las Pymes en su desarrollo. Actualmente las plataformas digitales están generando una gran cantidad de información que bien interpretada podría impulsar el desarrollo de las Pymes en México.

El 72% de las empresas encuestadas, tiene la intención de implementar el Machine Learning para automatizar la forma en que detectan datos de anomalías, el 54% planea profundizar en la plataforma de desempeño para reestructurar cómo procesa y usa los datos en los próximos 1 a 3 años.

Becerra Pozas, J. L. (2023, 3 de septiembre). Una paradoja: la mayoría de las empresas en México reconocen que necesitan datos, pero tienen dificultades con la proliferación de éstos. CIO MX. https://cio.com.mx/una-paradoja-la-mayoria-de-las-empresas-en-mexico-reconocen-que-necesitan-datos-pero-tienen-dificultades-con-la-proliferacion-de-estos/

### Hipótesis

El diseño de un modelo de datos predictivo de los consumidores con base a sus patrones de consumo e historial de compras, aumentarán las ventas de las empresas que lo implementen.

#### Justificación.

# ¿Por qué es relevante que se resuelva este problema?

La relevancia de la investigación consiste en ayudar a una empresa en México a restablecer su economía y con ello reducir el desempleo debido a que actualmente en la raíz de la pandemia se encuentran pasando por una crisis económica, por lo cual pretendo generar modelos de datos que ayuden a conocer a su consumidor final aumentando sus ventas.

### ¿Qué viabilidad tiene, en cuanto a implementación, realización, costo-beneficio, esfuerzo?

Para garantizar la viabilidad de la implantación de mi propuesta, se realizará una evaluación de los datos disponibles para análisis; el conocimiento que puede obtenerse de ese análisis; y los recursos disponibles para definir, diseñar, crear y desplegar un modelo de Ciencia de Datos.

Será sumamente importante que la empresa cuente con fuentes de datos necesarias para poder hacer un modelo, así como con la infraestructura tecnológica para la implantación.

La implementación se realizará utilizando herramientas Open Source, esto permitirá a la empresa reducir los costos de licenciamiento. Se impartirá una capacitación al personal de TI para que pueda administrar e interpretar los resultados del modelo de Ciencia de Datos.

El beneficio que obtendrá la empresa será contar con un modelo predictivo de ventas en información de valor que le ayude a tomar sus decisiones, además de capacitar a su personal para estar actualizado en las herramientas de Ciencia de Datos. El esfuerzo que tendría la empresa será la de proporcionar los insumos necesarios para que el proyecto de lleve acabó.

# ¿A quién beneficiará?

Beneficiará a una empresa del sector privado, que desea incrementar sus ventas mejorando sus estrategias y la toma de daciones incrementando sus clientes y mejorando sus servicios que les ofrece a los mismos.

### ¿Qué beneficio aportará en específico?

Una transformación digital para obtener valor de sus datos aportando soluciones innovadoras que impacten positivamente en su negocio.

Aumento de las ventas de productos en una empresa.

Conocer las preferencias de los consumidores, para ofrecerles productos a sus necesidades.

Las utilidades de la empresa van a incrementar con el análisis de datos.

# ¿Por qué es mejor que las soluciones que actualmente existen?

Actualmente las soluciones de Ciencia de Datos se enfocan en las grandes empresas, el modelo que pretendo generar está enfocada en las pymes para ayudarles en su toma de decisiones, cada día el consumidor genera

gran cantidad de información, el modelo de entrenamiento le dará un valor agregado a dicha información permitiendo a los directivos generar estrategias que les ayuden a mejorar sus servicios al cliente.

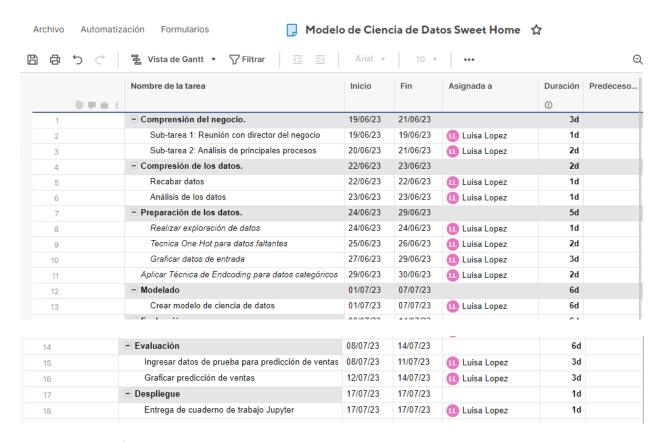
### Objetivo General.

 Diseño un modelo de datos que permita a los directivos tomar decisiones en las estrategias a seguir para la venta de sus productos.

# Objetivos específicos.

- Investigación de las metodologías para el análisis de la información.
- Investigación de los diferentes algoritmos para predecir ventas con Machine Learning.
- Investigar cómo se implementan los modelos de Ciencia de Datos.
- Investigar las herramientas para procesar las fuentes de datos y crear laboratorio para hacer pruebas y análisis de resultados.
- Identificar una empresa para desarrollar el proyecto de Ciencia de Datos.
- Identificar las fuentes de datos.
- Realizar un preprocesamiento de datos de una empresa.
- Desarrollo de un modelo de datos que genere un valor agregado en la toma de decisiones de los directivos.
- Implementar modelo de Ciencia de Datos.
- Obtener resultados y conclusiones de implantación de modelo de Ciencia de Datos.

#### CRONOGRAMA DE ACTIVIDADES



#### METODOLOGÍA

El desarrollo del proyecto se basa en la metodología CRISP-DM con las siguientes etapas:

### Conocimiento del negocio.

#### Planteamiento del problema.

La empresa del sector **inmobiliario Asesoría Sweet Home**, aspira a crecer en el mercado de los bienes raíces.

Los directivos han decidido implementar un modelo de ciencia de datos para predecir las comisiones de los próximos años en las ventas y rentas de los inmuebles. Específicamente, quieren comprender los meses en los cuales se vende más en el mercado de los bienes raíces, ya que pueden variar de acuerdo al periodo del año.

La empresa desea obtener una predicción de ventas y las comisiones que obtendrá por las mismas para los próximos 3 años (2024,2025 y 2026).

- Se deberá identificar los meses con mayor número de comisiones.
- El promedio de comisiones mensuales.
- Identificar la operación que genera el monto mayor de comisiones.
- Identificar los meses en los que se obtiene el monto mayor de comisiones.

# Adquisición y comprensión de los datos.

Los datos que se utilizarán para el modelo de entrenamiento se proporcionan por la inmobiliaria Asesoría **Sweet Home** en la cual describe los inmuebles más vendidos de los últimos tres años(2021, 2022 y 2023). Se solicitarán en un formato Layout ventas.xlsx con las siguientes columnas:

Columna	Descripción
Consecutivo	Número de operación realizada
Fecha	Fecha en la que se realizó la operación
Año	Año en que se realizó la operación
Operación	Columna que indica si la operación fue una Venta / Renta / Tramite Infonavit
Tipo de Inmueble	Categoriza los inmuebles en Casa / Departamento / Oficina
Monto	Monto de la comisión ganada

# Preparación de los datos.

Se importan los datos de las ventas.

Consecutivo	Fecha	Año	Operacion	TipoInmueble	Monto
1	2021-05-01	2021	Venta	Casa	45000.0
2	2021-06-01	2021	Venta	Casa	70500.0
3	2021-07-01	2021	Venta	Casa	25500.0
4	2021-08-01	2021	Venta	Casa	62500.0
5	2021-09-01	2021	Venta	Casa	144800.0

Se idéntica el tamaño de la muestra: 53 registros \* 6 columnas.

Se identifica el tipo de dato de las columnas.

Column	Non-Null Count	Dtype
Consecutivo	48 non-null	int64
Fecha	48 non-null	datetime64[ns]
Año	48 non-null	int64
Operacion	48 non-null	object
TipoInmueble	48 non-null	object
Monto	48 non-null	float64

Se obtiene un resumen estadístico de los montos de las comisiones.

count	53.000000
mean	29680.528302
std	29454.000313
min	1080.000000
25%	10000.000000
50%	18500.000000
75%	42000.000000
max	144800.000000

Se observa que la media de las comisiones es \$29,680.52.

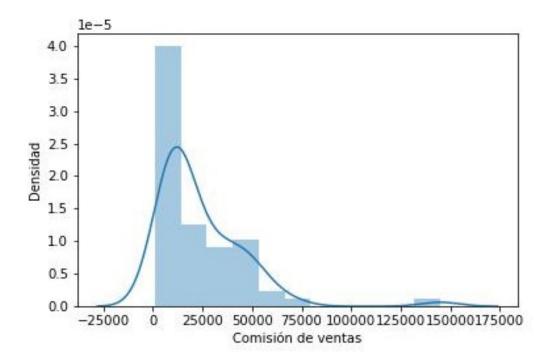
Identificar que valores faltantes. Se puede observar que no hay valores perdidos.

Consecutivo	0
Fecha	0
Año	0
Operacion	0
TipoInmueble	0
Monto	0

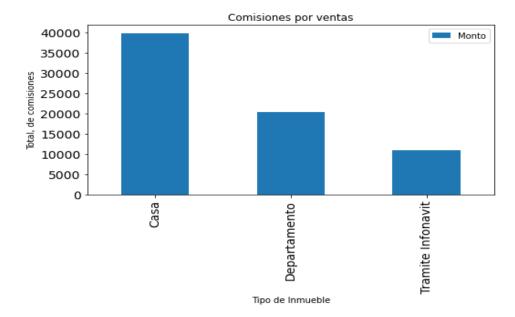
Visualización inicial de los datos.

Se realiza un análisis exploratorio de los datos:

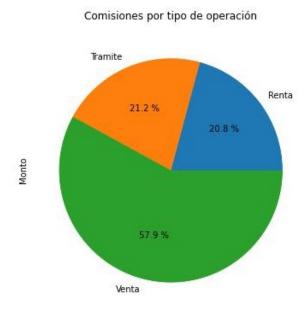
**Densidad de las comisiones obtenidas.** Se identifica que el rango del monto que se tienen más comisiones pro las ventas esta entre 0 y \$50,000.



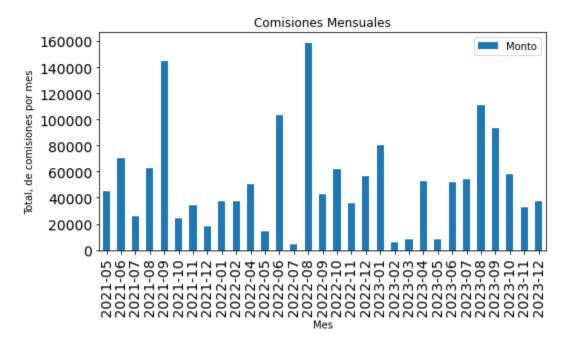
**Comisiones obtenidas.** Se identifica la venta de casas genera el mayor monto de comisiones, le sigue la venta de departamentos y tramites de créditos Infonavit.



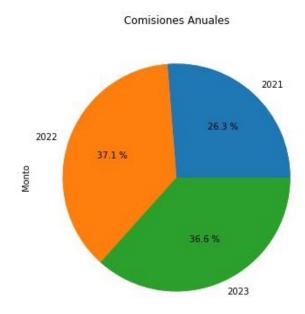
**Comisiones por tipo de operación.** Se identifica a la operación de ventas como la que genera el mayor monto de comisiones, le sigue tramites de créditos Infonavit y rentas.



Total, de comisiones mensuales. Se observa que los meses en que se vende más son julio y agosto.



**Total, de comisiones Anuales.** Se observa que el año con más ventas es 2022, aquí consideramos que el año 2021 se estuvo e pandemia por lo cual las ventas bajaron.



# Procesamiento de los datos.

Los modelos **ARIMA/SARIMA** son algunas de las mejores técnicas de modelado que se utilizan para el análisis de series temporales. Estos modelos requieren el paso de parámetros que deben conocerse para crear un modelo preciso. Hay diferentes métodos que se pueden utilizar para encontrar los parámetros más óptimos.

### Justificación l Uso del Modelo ARIMA.

Como las ventas pasadas afectan los valores actuales o futuros se puede predecir tendencias futuras basadas en fluctuaciones recientes. En ese caso, el pronóstico de series de tiempo es la solución para tal problema de regresión. Varios otros modelos de pronóstico de series temporales se basan en la incorporación de cambios sucesivos o desarrollos más recientes en los datos para predecir tendencias futuras.

Por lo general, cuando se implementan modelos ARIMA/SARIMA, es necesario conocer los parámetros p, d y q para el modelado ARIMA y p, d, q, P, D, Q y m para el modelado SARIMA.

p es el orden autorregresivo de la tendencia

d es el orden de diferenciación de tendencia

q es el orden promedio móvil de tendencia

P es el orden autorregresivo estacional

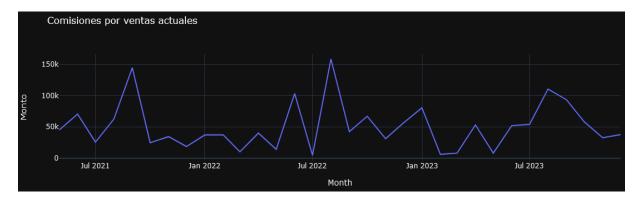
D es el orden de diferenciación estacional

Q es el orden promedio móvil estacional

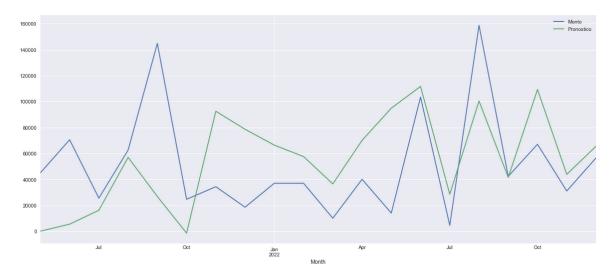
m es el número de pasos en un período estacional

pmdarima.auto\_arima() es una herramienta útil que realizará una búsqueda en cuadrícula para encontrar los mejores parámetros posibles. El único que no encontrará es el parámetro m, esto es algo que tendremos que resolver e ingresar nosotros mismos. Podemos resolver esto haciendo una descomposición\_estacional para encontrar el número de pasos en un período estacional. (También es posible conocer este parámetro simplemente sabiendo si sus datos son recolectados mensual, trimestral o anualmente)

**Total, de comisiones por ventas actuales.** Podemos observar u aumento de ventas entre los meses de julio y octubre del 2022.



Comparando las comisiones actuales contra las predicciones realizadas por el modelo ARIMA. Se observa que las predicciones van coincidiendo con las ventas reales que se tienen para los tres primeros años (2021, 2022 y 2023).



#### Evaluación.

Se ejecuta la predicción del modelo ARIMA y se obtiene las predicciones de comisiones para los próximos tres años 2024, 2025 y 2026, se observa que de continuar las estrategias actuales se tendrá un alta de las ventas, esto considerando que al ya no estar en pandemia las ventas aumentaron considerablemente durante 2022 y 2023.

### Implementación.

Se entregará a la inmobiliaria el repositorio que contiene los layouts de información y el cuaderno de Jupyter Noobook que contiene el modelo de ciencia de datos. Adicional se brindará una capacitación sobre como ejecutar el modelo de ciencia de datos.

También se entregará un Data storytelling, la cual permitirá a los directivos interpretar de una manera sencilla las ventas actuales y el comportamiento que tendrán en los próximos 3 años.

Se exportan las predicciones en un archivo prediccion\_comisiones.csv el cual también se entrega en el repositorio.



# RESULTADOS Y DISCUSIÓN.

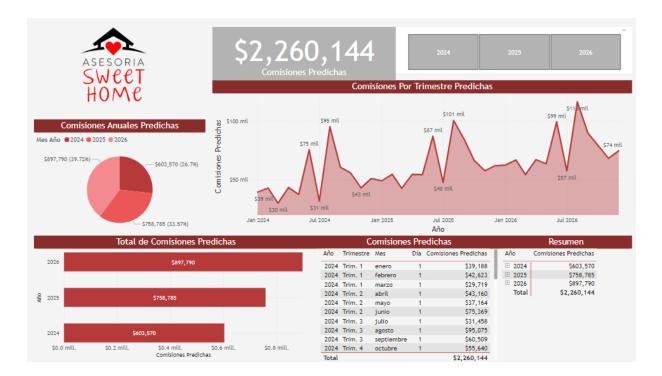
Histórico de comisiones 2021,2022 y 2023.

**Total, de comisiones Anuales.** Se observa que el año con más ventas es 2022, aquí hay que tomar en cuenta que el año en curso es 2023 por lo cual al no contar con la información del segundo semestre se obtuvo con el promedio de los años 2021 y 2022.

La operación que más comisiones genera es la venta de inmuebles, el tipo de inmueble más vendido es departamento.



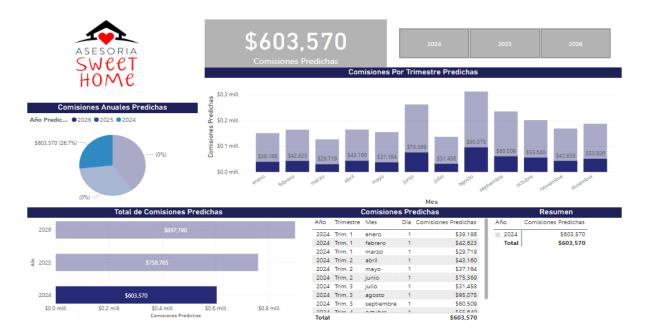
**Predicción de comisiones 2024, 2025 y 2026**. Se observa que los meses de junio y septiembre son los picos más altos en los que se encontrará el mayor número de ventas.



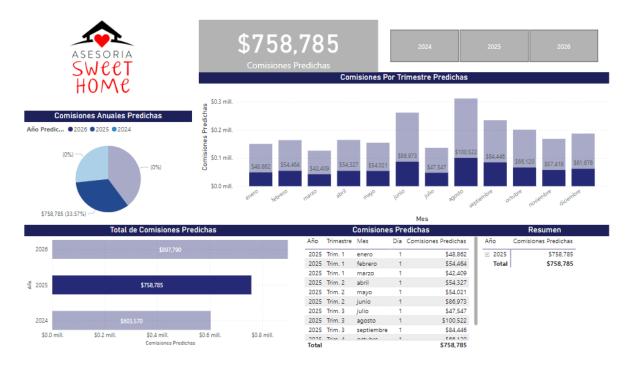
# **CONCLUSIONES**

De acuerdo a los resultados del modeo de predicción ARIMA por series de tiempo.

Para el año 2024 de obtendra un total de \$603,570.00 pesos mexicanos en ganancias.



Para el año 2025 de obtendrá un total de \$ 758,785.00 pesos mexicanos en ganancias.



Para el año 2026 de obtendrá un total de \$ 897,790.00 pesos mexicanos en ganancias. En este año se proyectan las ventas más altas.



En conclusión sobre el modelo generado se observa que es posible aplicar ciencia de datos Para las PYMES en México ya que a pesar de contar con poca información se logró realizar un modelo bastante asertado en la predicción de ventas de los proxismos 3 años para la empresa Asesoria Sweet Home.

Se observa qu existe una estacionalidad en los meses de junio y agosto como los marcados con las ventas más altas, asi mismo se extendie al periode de septeimbre, octubre, nomvienbre y diciembre.

Como siguientes pasos se recomienda a la inmobiliaria Asesoria Sweet Home, aumentar el número de exclusivas para ventas ya que es el tipo de peración que genera más ganancias.

Se recomienda tambien la contratación de un CRM de gestión de propiedades que pueda estar alimentando el modelo de predicción de ventas.

Se recomienda la capacitación del personal en el areá estadística para que sean capaces de leer e interpretar información grafica que les permita una adecuada toma de desiciones basadas en datos.

#### LISTA DE REFERENCIAS.

- Becerra Pozas, J. L. (2023, 3 de septiembre). Una paradoja: la mayoría de las empresas en México reconocen que necesitan datos, pero tienen dificultades con la proliferación de éstos. CIO MX. <a href="https://cio.com.mx/una-paradoja-la-mayoria-de-las-empresas-en-mexico-reconocen-que-necesitan-datos-pero-tienen-dificultades-con-la-proliferacion-de-estos/">https://cio.com.mx/una-paradoja-la-mayoria-de-las-empresas-en-mexico-reconocen-que-necesitan-datos-pero-tienen-dificultades-con-la-proliferacion-de-estos/</a>
- Dor. (2021, 7 octubre). Forecasting time series with Auto-Arima Data Science portfolio. <a href="https://www.alldatascience.com/time-series/forecasting-time-series/">https://www.alldatascience.com/time-series/forecasting-time-series/</a>
- Kyle, S. (2022, 7 enero). Quick intro: Auto\_Arima from PMDarima Package Steven Kyle Medium.
  Medium. <a href="https://stevenkyle2013.medium.com/quick-intro-auto-arima-from-pmdarima-package-e7aab5e8dfb8">https://stevenkyle2013.medium.com/quick-intro-auto-arima-from-pmdarima-package-e7aab5e8dfb8</a>

Kaggle: your machine learning and data science community. (s. f.). <a href="https://www.kaggle.com/">https://www.kaggle.com/</a>

- Gonzalez, L. (2022, 23 septiembre). K Vecinos más Cercanos Teoría. Aprende IA. <a href="https://aprendeia.com/algoritmo-k-vecinos-mas-cercanos-teoria-machine-learning/">https://aprendeia.com/algoritmo-k-vecinos-mas-cercanos-teoria-machine-learning/</a>
- Sanz, F. (2022, 17 noviembre). K-Means Clustering: algoritmo, aplicaciones y desventajas. The Machine Learners. <a href="https://www.themachinelearners.com/k-means/">https://www.themachinelearners.com/k-means/</a>

- Arce, J. I. B. (2022, 14 junio). Light GBM vs XGBoost . ¿Cuál es mejor el algoritmo? Juan Barrios. https://www.juanbarrios.com/light-gbm-vs-xgboost-cual-es-mejor-algoritmo/
- J. (2021, 9 diciembre). Minería de Datos #1 JAYWRKR. Medium. https://blog.jaywrkr.com/data-mining-1-735be1532599
- Alfredo Bueno Solano. (2018). Diseño conceptual del modelo de big data para el IMT. 27/11/2021, de secretaria de Comunicaciones y Trasportes Sitio web: https://imt.mx/archivos/Publicaciones/PublicacionTecnica/pt537.pdf
- Data.Barcelona. (2021). Cómo desarrollar modelos de negocio basados en datos. 28/11/2021, de Barcelona Sitio web: https://data.barcelona/2021/05/28/como-desarrollar-modelos-de-negocio-basados-endatos/
- Fernandez, R. (2020, 9 abril). Fases de un proyecto de Data Science ▷ Cursos de programación de 0 a experto © garantizados. ▷ Cursos de Programación de 0 a Experto © Garantizados. https://unipython.com/fases-de-un-proyecto-de-data-science/
- ARIMA S-ARIMA. (2021, 21 marzo). Ricardo R. Palma. Recuperado 8 de agosto de 2023, de https://themys.sid.uncu.edu.ar/rpalma/MBA/Evaluaciones%202022/Arima/Arima2.html
- Bernardes, L. (2021). Guía del Data Storytelling: Cómo cautivar a tu audiencia con historias basadas en datos valiosos. Rock Content ES. https://rockcontent.com/es/blog/data-storytelling/
- PowerBI: Qué es y cómo transforma el análisis de datos. (s. f.). https://www.capacitarte.org/blog/nota/Power-BI-analisis-de-datos
- INTELIGENCIA DE NEGOCIOS y ANALITICA DE DATOS. (s. f.). Librerías Gandhi. https://www.gandhi.com.mx/inteligencia-de-negocios-y-analitica-de-atos?gclid=CjwKCAjw29ymBhAKEiwAHJbJ8o\_U-fsYcuwjpdN-Qo\_i301qIRdAiH7TW58I8nq5oKYXhlVUOy81EhoCis4QAvD\_BwE
- Ciencia de datos con aplicaciones en R y Python : Hernández López, Eymard, Cruz Diosdado, Leonardo David, Wences Nájera, Giovanni Arquímedes: Amazon.com.mx: Libros. (s. f.). https://www.amazon.com.mx/Ciencia-aplicaciones-python-Eymard-

Hern%C3%A1ndez/dp/B08QRPLP8Y/ref=asc\_df\_B08QRPLP8Y/?tag=gledskshopmx-20&linkCode=df0&hvadid=547161748294&hvpos=&hvnetw=g&hvrand=604067717251175855 7&hvpone=&hvptwo=&hvqmt=&hvdev=c&hvdvcmdl=&hvlocint=&hvlocphy=9142899&hvtarg id=pla-1128590135224&psc=1

Libro Python con aplicaciones a las matematicas, ingenieria y finanzas, David Lopez Baez, ISBN 9786076226735. Comprar en Buscalibre. (s. f.). https://www.buscalibre.com.mx/libro-python-con-aplicaciones-a-las-matematicas-ingenieria-y-finanzas-david-lopez-baez-alfaomega-grupo-editor/9786076226735/p/48965793?bmkt\_source=google&bmkt\_campaign=16070390338&gclid=CjwKCAjw29ymBhAKEiwAHJbJ8gcUQxN-LEdDdn-3xu8P3h6bB-eEg-rCKI86eqNJkVTNAoxBgUlmUxoCLfQQAvD\_BwE

Estadística práctica para ciencia de datos con R y Python: Bruce, Peter, Bruce, Andrew, Gedeck, Peter:

Amazon.com.mx: Libros. (s. f.). https://www.amazon.com.mx/Estad%C3%ADstica-pr%C3%A1ctica-ciencia-datos-

Python/dp/842673443X/ref=asc\_df\_842673443X/?tag=gledskshopmx-20&linkCode=df0&hvadid=547307517812&hvpos=&hvnetw=g&hvrand=295031379937828749

3&hvpone=&hvptwo=&hvqmt=&hvdev=c&hvdvcmdl=&hvlocint=&hvlocphy=9142899&hvtarg

id=pla-1628525925251&psc=1

Ingeniería de software: Sommerville: Amazon.com.mx: libros. (s. f.).

<a href="https://www.amazon.com.mx/Ingenieria-Software-Sommerville/dp/6073206038">https://www.amazon.com.mx/Ingenieria-Software-Sommerville/dp/6073206038</a>

EPJ Data Science. (2023, 11 agosto). SpringerOpen. https://epjdatascience.springeropen.com/articles

Modelo predictivo del proceso de ventas utilizando inteligencia de negocios y data analitics en la empresa
centro textil De la Matta S.A.C. (2023). https://repositorio.uss.edu.pe/. Recuperado 8 de agosto de
2023,

 $\frac{https://repositorio.uss.edu.pe/bitstream/handle/20.500.12802/10636/Carre\%C3\%B1o\%20Guerrer}{o\%2C\%20Santiago\%20An\%C3\%ADbal.pdf?sequence=1\&isAllowed=y}$ 

Implementación de Machine Learning en el área de ventas de la. (2022, 21 febrero).

http://repositorio.ucsg.edu.ec/. Recuperado 8 de agosto de 2023, de <a href="http://repositorio.ucsg.edu.ec/bitstream/3317/18337/1/T-UCSG-PRE-CEAC-CNI-21.pdf">http://repositorio.ucsg.edu.ec/bitstream/3317/18337/1/T-UCSG-PRE-CEAC-CNI-21.pdf</a>