



Matemáticas para

Machine Learning

Fundamentos de Ciencia de Datos con Python

Jheser Guzman
Jhohanser Guzman

Matemáticas para Machine Learning

Fundamentos de Ciencia de Datos con Python

Matemáticas para Machine Learning

Fundamentos de Ciencia de Datos con Python

por

Bluminds LLC

Jheser Guzman

Senior Software Development Engineer Amazon-AWS, ex-JPMorgan & Chase, ex Google in USA. Liverpool John Moore's University, University of Chile University of León.

Jhohanser Guzman Dr. (c)

Senior Data Scientist, ex-UNDP, Columbia University, UNINI University, Universidad Europea del Atlántico, Universitat Politècnica de Catalunya, AIU University and Economist.

Matemáticas para Machine Learning

Por *Jheser Guzman & Jhohanser Guzman*

Copyright © 2026 Bluminds LLC. Todos los derechos reservados.

Impreso en Estados Unidos de América

Publicado por Bluminds Books, parte de Bluminds LLC.

Domicilio: 20 S Charles St Ste 403 #878, Baltimore, MD 21201 - USA

Queda prohibida la reproducción, registro o transmisión, total o parcial, de esta publicación por cualquier sistema de recuperación de información, en cualquier formato o medio —electrónico, mecánico, fotoquímico, magnético o electroóptico—, incluyendo fotocopia, grabación o cualquier otro procedimiento, sin la autorización previa y por escrito del editor.

El préstamo, alquiler o cualquier otra modalidad de cesión de uso de este ejemplar requerirá igualmente la autorización expresa del editor o de sus representantes.

Las obras de Bluminds Books pueden adquirirse para fines educativos, empresariales y comerciales. Existen, además, ediciones digitales disponibles en www.bluminds.com. Para información adicional, contacte a nuestro departamento de ventas corporativas: info@bluminds.com

ISBN: 978-1-967406-02-9

Editora y Revisora: Raphaella Dariana Illanes Bequer

Editor de producción: Jheser Guzmán & Jhohanser Guzmán

Primera edición, Enero, 2026

Historial de revisiones de la primera edición: 2024-2026

Diseñador de portada: Jhohanser Guzmán

Diseñador de interiores: Jheser Guzmán & Jhohanser Guzmán

Verificación de originalidad y ausencia de plagio realizada mediante Dupli Checker y Plagiarism Detector.net

Para detalles adicionales sobre el lanzamiento de este libro, consulte info@bluminds.com.

El logotipo de Bluminds constituye una marca registrada de Bluminds LLC. La imagen de portada corresponde a un diseño original generado mediante inteligencia artificial.

Muchas de las denominaciones empleadas por fabricantes y distribuidores para distinguir sus productos están registradas como marcas comerciales.

Ni el editor ni los autores asumen responsabilidad por daños resultantes del mal uso de la información contenida en este libro. Nos reservamos el derecho de rectificar errores u omisiones

“Científico de Datos: El trabajo más sexy del siglo XXI”

— por Harvard Business Review: Data Scientist: The Sexiest Job of the 21st Century.

Agradecimientos

Manifestamos nuestro más profundo reconocimiento a Bluminds LLC, en particular al equipo responsable de la edición y revisión técnica de este volumen. Su dedicación, pericia y compromiso han sido determinantes para asegurar la precisión del contenido presentado. Asimismo, damos nuestro agradecimiento a colegas y colaboradores, cuyas observaciones y sugerencias, enriquecieron de manera significativa la presente obra, al igual que nuestros estudiantes, cuya inquietud intelectual y entusiasmo por el aprendizaje constituyeron un estímulo constante para profundizar y clarificar los conceptos aquí expuestos. A nuestras familias, por su respaldo y aliento, durante el desarrollo de este proyecto; sin estos valores, la culminación de este libro no habrían sido factibles.

Hay que destacar la contribución de la comunidad de software libre, cuyos esfuerzos ha posibilitado la expansión y accesibilidad de herramientas como el Lenguaje Python, y todos los ecosistemas de ciencia de datos. El espíritu de cooperación, que caracteriza a la comunidad open source, ha sido una fuente inagotable de inspiración y un pilar esencial en la gestación de este libro.

Cabe señalar, con especial gratitud, la labor de Raphaella Dariana Illanes Bequer, cuya dedicación y aguda capacidad analítica resultaron determinantes en la revisión de este manuscrito. Su habilidad para identificar matices y formular observaciones enriquecedoras ha elevado de manera notable la calidad textual, garantizando la transmisión precisa y clara de cada idea. Su compromiso con la excelencia académica y su profesionalismo se reflejan en cada página, consolidándose como un soporte para la realización de este proyecto.

Reconocemos a la comunidad global de profesionales en ciencia de datos, quienes, mediante su innovación continua, han propiciado el avance acelerado de este campo, generando un impacto positivo y trascendental en personas y organizaciones a nivel internacional.

Acerca de los Autores



Jheser Guzmán es científico en computación, con una trayectoria de dos décadas en Machine Learning, Procesamiento de Lenguaje Natural (NLP) y en especial, Ingeniería de Software. Hizo sus estudios doctorales en Ciencias de la Computación en la Universidad de Chile y ha completado cuatro maestrías en universidades del Reino Unido, España y Bolivia, además de certificaciones profesionales como Cisco CCNA, AWS Architect/Developer Associate y Oracle Certified Associate (OCA). Ha ejercido funciones en empresas tecnológicas de referencia, tales como *Amazon AWS*, *Google – YouTube*, *Yahoo Research* y *JPMorgan Chase*, donde ha liderado iniciativas en sistemas de backend y arquitecturas de microservicios. Su especialización en Python, para ciencia de datos, ML y NLP, junto con su experiencia industrial, ha favorecido avances en el rendimiento algorítmico y en el desarrollo de softwares. En el ámbito académico, ha impartido docencia en universidades de Bolivia, Chile y Estados Unidos, orientando a estudiantes y profesionales.



Jhohanser Guzman es economista de profesión, científico de datos y ambientalista. Cuenta con quince años de experiencia en los sectores financiero, jurídico y de análisis de datos. Obtuvo la licenciatura en la Universidad Católica y ha desarrollado su perfil profesional con diversas maestrías, entre ellas en Ciencias de Datos, Machine Learning (AIU), Derecho y Estudios Jurídicos Avanzados. Cuenta un certificado profesional en Machine Learning por la Universidad de Columbia en Nueva York. Actualmente cursa un doctorado en Ciencias de la Salud Pública, con énfasis en análisis de datos, en la Universidad Internacional Iberoamericana de Mexico. Ha trabajado en organismos nacionales e internacionales, como el Programa de las Naciones Unidas (PNUD) en Bolivia, y en BeeGrids Social Intelligence en Chile. Donde ha desarrollado análisis geoespacial, evaluaciones estadísticas de redes sociales y estudios sobre comportamiento social ante desastres naturales.

Acerca de Bluminds

Bluminds LLC, es un startup EdTech orientada al desarrollo de soluciones tecnológicas de vanguardia para la democratización del acceso a la educación y el aprendizaje. Especializada en plataformas de aprendizaje digital y en la generación de contenido académico. Bluminds LLC mantiene un compromiso sostenido con la creación de herramientas que posibiliten a instituciones educativas y empresas ofrecer experiencias formativas interactivas, escalables y accesibles.

El equipo de Bluminds LLC integra la experiencia de profesionales en ingeniería de software, ciencia de datos, economía y tecnologías educativas, colaborando con actores de la industria, universidades y entidades gubernamentales. Además, la organización brinda servicios de consultoría y capacitación técnica personalizada para la implementación de soluciones de aprendizaje digital, tanto en modalidad presencial como remota.

Bluminds LLC ha desarrollado, asimismo, libros y recursos educativos destinados a fortalecer comunidades tecnológicas y académicas, abarcando temáticas clave como ciencia de datos, Python, desarrollo de software, entre otros.

Para establecer contacto con Bluminds LLC y sus autores, o para solicitar una propuesta de capacitación técnica presencial o virtual, diríjase a: *info@bluminds.com*

Para información adicional sobre los servicios de capacitación, soluciones educativas, materiales de trabajo y presentaciones de Bluminds LLC, visite: *https://www.bluminds.com*

Individuos o instituciones interesados en adquirir recursos educativos de Bluminds LLC pueden hacerlo a través de: *https://www.amazon.com*

Las solicitudes de adquisición en grandes volúmenes por parte de corporaciones, organismos gubernamentales o instituciones educativas deben gestionarse directamente con Bluminds LLC. Para mayor información, escriba a: *info@bluminds.com*

Índice general

| | |
|---|--------------|
| Prefacio | XIII |
| Antes de Comenzar | XXIII |
| Capítulo 1 Introducción a Machine Learning | 1 |
| 1.1 ¿Qué es Machine Learning? | 1 |
| 1.2 Machine Learning vs Data Mining | 4 |
| 1.3 ¿Qué es la Ciencia de Datos? | 9 |
| 1.4 Matemática y estadística en Machine Learning | 12 |
| 1.5 Python para Machine Learning | 14 |
| 1.6 Herramientas para Machine Learning | 16 |
| 1.7 Proyectos de Machine Learning | 18 |
| 1.8 Resumen | 21 |
| 1.9 Problemas y sus Soluciones | 24 |
| 1.10 Problemas | 27 |
| Capítulo 2 Programación Python para Machine Learning | 33 |
| 2.1 Pasos iniciales de Jupyter Notebook | 33 |
| 2.2 Tipos de datos | 34 |
| 2.3 Operadores lógicos | 41 |
| 2.4 Tablas, listas, diccionarios e indexación | 47 |
| 2.5 Operaciones con DataFrames y bases de datos | 54 |
| 2.6 Tuplas | 67 |
| 2.7 Funciones | 75 |
| 2.8 Argumentos | 86 |
| 2.9 Métodos en objetos | 99 |
| 2.10 Estructuras de control | 110 |
| 2.11 Excepciones | 122 |
| 2.12 Resumen | 132 |
| 2.13 Problemas y sus Soluciones | 137 |
| 2.14 Problemas | 145 |
| Capítulo 3 Python para Análisis de Datos | 149 |
| 3.1 Módulos de Python para análisis de datos | 149 |
| 3.2 NumPy - Python numérico | 151 |
| 3.3 SymPy - Matematica simbolica | 161 |
| 3.4 SciPy - Python científico | 172 |
| 3.5 Pandas - Manipulación y análisis de datos tabulares | 188 |
| 3.6 Matplotlib - Visualización de datos | 201 |

| | | |
|-------------------|---|------------|
| 3.7 | Seaborn - Visualización estadística | 203 |
| 3.8 | Plotly - Visualización interactiva | 204 |
| 3.9 | Bokeh - Visualización interactiva | 205 |
| 3.10 | Dash - Aplicaciones Web interactivas | 212 |
| 3.11 | Folium - Visualización de datos geoespaciales | 218 |
| 3.12 | Resumen | 222 |
| 3.13 | Problemas y sus Soluciones | 225 |
| 3.14 | Problemas | 230 |
| Capítulo 4 | Visualización de Datos | 235 |
| 4.1 | Diagrama de dispersiones | 235 |
| 4.2 | Gráficos de barras | 245 |
| 4.3 | Gráficos de líneas | 265 |
| 4.4 | Histogramas | 278 |
| 4.5 | Diagrama de cajas | 291 |
| 4.6 | Gráficos de sectores | 299 |
| 4.7 | Mapas de calor | 307 |
| 4.8 | Diagrama de curvas de nivel | 317 |
| 4.9 | Diagrama de violín | 327 |
| 4.10 | Diagrama de contornos para mapas | 334 |
| 4.11 | Diagrama de áreas y volúmenes | 354 |
| 4.12 | Diagrama de redes | 365 |
| 4.13 | Diagrama de radar | 383 |
| 4.14 | Diagrama de Marimekko | 388 |
| 4.15 | Gráfico de burbujas | 390 |
| 4.16 | Diagrama de árbol | 392 |
| 4.17 | Gráfica de densidad | 395 |
| 4.18 | Diagrama de Venn | 399 |
| 4.19 | Resumen | 404 |
| 4.20 | Problemas y sus Soluciones | 409 |
| 4.21 | Problemas | 414 |
| Capítulo 5 | Aritmética, Álgebra, Geometría y Trigonometría | 419 |
| 5.1 | Aritmética | 419 |
| 5.2 | Álgebra | 434 |
| 5.3 | Geometría y Trigonometría | 455 |
| 5.4 | Resumen | 498 |
| 5.5 | Problemas y sus Soluciones | 500 |
| 5.6 | Problemas | 505 |
| Capítulo 6 | Álgebra Lineal, Cálculo y Ecuaciones Diferenciales | 509 |
| 6.1 | Álgebra lineal | 509 |

| | | |
|---|--|------------|
| 6.2 | Cálculo | 527 |
| 6.3 | Ecuaciones diferenciales | 553 |
| 6.4 | Resumen | 564 |
| 6.5 | Problemas y sus Soluciones | 566 |
| 6.6 | Problemas | 572 |
| Capítulo 7 Probabilidad y Estadística | | 577 |
| 7.1 | Teoría de probabilidades | 577 |
| 7.2 | Distribuciones de probabilidades | 598 |
| 7.3 | Simulaciones de datos | 607 |
| 7.4 | Estadística descriptiva | 612 |
| 7.5 | Medidas de posición central y no central | 618 |
| 7.6 | Medidas de dispersión | 620 |
| 7.7 | Medidas de forma | 622 |
| 7.8 | Estadística inferencial | 624 |
| 7.9 | Análisis bidimensional | 631 |
| 7.10 | Resumen | 639 |
| 7.11 | Problemas y sus Soluciones | 641 |
| 7.12 | Problemas | 646 |
| Capítulo 8 Módulos de Python para Machine Learning | | 649 |
| 8.1 | Librerías para Machine Learning | 649 |
| 8.2 | Scikit para Machine Learning | 650 |
| 8.3 | TensorFlow | 662 |
| 8.4 | Librería PyTorch | 677 |
| 8.5 | Librería Keras | 686 |
| 8.6 | Librería XGBoost | 695 |
| 8.7 | Librería LightGBM | 699 |
| 8.8 | Librería Statsmodels | 706 |
| 8.9 | Resumen | 713 |
| 8.10 | Problemas y sus Soluciones | 716 |
| 8.11 | Problemas | 722 |
| Capítulo 9 Algoritmos Supervisados | | 727 |
| 9.1 | Algoritmos supervisados | 728 |
| 9.2 | Regresión lineal simple | 729 |
| 9.3 | Regresión Ridge | 734 |
| 9.4 | Regresión Lasso | 739 |
| 9.5 | Regresión Elastic Net | 743 |
| 9.6 | Regresión Kernel | 749 |
| 9.7 | Regresión por procesos Gaussianos | 753 |
| 9.8 | Regresión Probabilística | 756 |

| | | |
|---|---|------------|
| 9.9 | Regresión Bayesiana | 761 |
| 9.10 | Regresión Sparse (con énfasis en LARS) | 764 |
| 9.11 | Clasificación lineal | 767 |
| 9.12 | Clasificación logística | 771 |
| 9.13 | Clasificación Probit | 775 |
| 9.14 | Clasificación Tobit | 778 |
| 9.15 | Clasificación por la aproximación de Laplace | 782 |
| 9.16 | Clasificación por K-Vecinos más cercanos (k-NN) | 786 |
| 9.17 | Clasificación Bayesiana (Naive Bayes) | 788 |
| 9.18 | Clasificación por Support Vector Machines (SVM) | 792 |
| 9.19 | Árboles de decisión | 795 |
| 9.20 | Bosques aleatorios - Random Forest | 799 |
| 9.21 | Remuestreo por Bootstraps | 802 |
| 9.22 | Remuestreo por RANSAC | 805 |
| 9.23 | Ensamblaje: Bagging | 808 |
| 9.24 | Ensamblaje: Boosting | 811 |
| 9.25 | Ensamblaje: Stacking | 814 |
| 9.26 | Ensamblaje: Voting | 817 |
| 9.27 | Reglas de asociación - Apriori | 821 |
| 9.28 | Análisis del Discriminante Lineal (LDA) | 826 |
| 9.29 | Análisis del Discriminante Cuadrático (QDA) | 830 |
| 9.30 | Redes neuronales supervisadas: Perceptrón | 833 |
| 9.31 | Redes neuronales supervisadas: Perceptrón multicapa (MLP) | 836 |
| 9.32 | Resumen | 844 |
| 9.33 | Problemas y sus Soluciones | 851 |
| 9.34 | Problemas | 856 |
| Capítulo 10 Algoritmos No Supervisados | | 861 |
| 10.1 | Introducción | 861 |
| 10.2 | Clustering: K-means | 862 |
| 10.3 | Clustering: K-medoids | 866 |
| 10.4 | Missing Data: Algoritmo de maximización de la esperanza | 869 |
| 10.5 | Dimensionality Reduction: Principal Component Analysis (PCA) | 871 |
| 10.6 | Dimensionality Reduction: Componente Principal Probabilístico | 875 |
| 10.7 | Dimensionality Reduction: Kernel based PCA | 878 |
| 10.8 | Matrix Factorization: SVD | 882 |
| 10.9 | Matrix Factorization: Matriz Factorizada Probabilística | 885 |
| 10.10 | Matrix Factorization: Matriz Factorizada No Negativa | 888 |
| 10.11 | Modelos secuenciales: Cadenas de Markov | 891 |
| 10.12 | Modelos secuenciales: Hidden Markov | 894 |
| 10.13 | Modelos secuenciales: Rankings | 896 |
| 10.14 | Recomendación: Filtro Colaborativo | 899 |

| | | |
|---|---|-------------|
| 10.15 | Recomendación: Tracking por el filtro de Kalman | 903 |
| 10.16 | Búsqueda de tópicos: Asignación Latente de Dirichlet | 907 |
| 10.17 | Redes Neuronales: Autoencoders NS | 910 |
| 10.18 | Redes Neuronales: Restricted Boltzmann Machines (RBM) | 913 |
| 10.19 | Resumen | 919 |
| 10.20 | Problemas y sus Soluciones | 923 |
| 10.21 | Problemas | 928 |
| Bibliografía | | 933 |
| Apéndice A Introducción a Python | | 941 |
| A.1 | ¿Qué es Python? | 941 |
| A.2 | Variables y tipos de datos | 949 |
| A.3 | Operaciones Básicas | 959 |
| A.4 | Estructuras de control de flujo | 977 |
| A.5 | Funciones y programación modular | 988 |
| A.6 | Estructuras de datos en Python | 1000 |
| Apéndice B Respuestas a Preguntas de los Capítulos | | 1011 |
| B.1 | Capítulo 1 - Respuestas | 1011 |
| B.2 | Capítulo 2 - Respuestas | 1011 |
| B.3 | Capítulo 3 - Respuestas | 1020 |
| B.4 | Capítulo 4 - Respuestas | 1021 |
| B.5 | Capítulo 5 - Respuestas | 1021 |
| B.6 | Capítulo 6 - Respuestas | 1022 |
| B.7 | Capítulo 7 - Respuestas | 1023 |
| B.8 | Capítulo 8 - Respuestas | 1023 |
| B.9 | Capítulo 9 - Respuestas | 1036 |
| B.10 | Capítulo 10 - Respuestas | 1036 |
| Índice alfabético | | 1039 |

Prefacio

El avance acelerado de la inteligencia artificial y el aprendizaje automático ha generado una demanda creciente por una base matemática aplicada que sustente el desarrollo de modelos y algoritmos en estos campos. Este libro se propone ofrecer una exposición sistemática y profunda de las Matemáticas para Machine Learning, proporcionando al lector los fundamentos teóricos y metodológicos necesarios para abordar los desafíos actuales en la investigación y la práctica profesional. En este prefacio, se introduce el enfoque interdisciplinario que caracteriza la obra, articulando la relación entre aritmética, entre álgebra, entre cálculo y otras disciplinas matemáticas para el progreso del aprendizaje automático y la inteligencia artificial.¹

Históricamente, las matemáticas han sido las herramientas más importantes para comprender y modelar todo lo que nos rodea, con base sólida y acompañada de lógica, formulando una base robusta para las demás ciencias. Desde ecuaciones que describen el movimiento de los planetas en el universo hasta las fórmulas que describen las leyes de la física, la naturaleza y la mecánica cuántica, las matemáticas han sido el lenguaje del universo, y ahora las aplicaremos para crear y modelar el comportamiento en sistemas computacionales.²

Machine Learning, o Aprendizaje Automático en español, se ha convertido en un pilar importante en la revolución tecnológica actual. Así están los algoritmos de recomendación en plataformas de comunicaciones, hasta los sistemas de reconocimiento que están en nuestros teléfonos móviles. Así, Machine Learning está presente en diferentes dispositivos, como en la mayor parte de los dispositivos electrónicos que tenemos en el hogar y en la oficina. Detrás de la magia de estos algoritmos, existen fundamentos matemáticos que dan vida a estas aplicaciones.³

En este libro se examina el papel central que desempeñan las matemáticas como fundamento teórico y metodológico para el desarrollo y la comprensión de los modelos de Machine Learning, considerados como una de las áreas más útiles dentro de la inteligencia artificial.

El cálculo, en su forma más avanzada, se convierte en una herramienta muy útil para comprender cómo los modelos matemáticos pueden evolucionar a medida

¹Edwards, C.H. Jr & Penney David. (1994). “Cálculo con Geometría Analítica”. México DF, México. Prentice Hall Hispanoamericana S.A.

²Chávez Reyes, Carmen & León Quintanar. (2007). “La Biblia de las Matemáticas”. México D.F., México. Editorial Letrarte S.A.

³Baldor. A. (2003). “Aritmética: Teoría Practica”. México DF, México. Publicaciones Cultural Decima Octava Impresión.

que se enfrentan a nuevas características o a nuevos tipos de datos.⁴ Las derivadas parciales, que antes eran solo unos conceptos abstractos en ciencias, ahora se transforman en las piedras angulares de los gradientes y optimizadores que impulsan el aprendizaje de las máquinas, es decir, que el cálculo se encarga de optimizar cada uno de los algoritmos que necesitan ser optimizados para después aplicarlos a casos concretos. Además del cálculo, también veremos álgebra lineal, con conceptos como las matrices y los vectores, se convierten en herramientas útiles, para representar y manipular datos en el contexto de ML.

En cada capítulo de este libro, destacaremos la aplicabilidad de estos conceptos matemáticos avanzados en el desarrollo y mejora de los algoritmos de Machine Learning. Aquí, no solo nos quedamos con la teoría, sino que también hemos de mostrar ejemplos prácticos que permitirán a los lectores experimentar directamente con la implementación de estos modelos algorítmicos y matemáticos en un entorno de programación en el lenguaje Python.⁵

Hablando de Python, este libro continuará abogando por su utilidad y versatilidad en el ámbito de las matemáticas, estadística y Machine Learning. Aquí, Python, con su sintaxis y sus amplias bibliotecas (librerías) de herramientas específicas para Machine Learning, se erige como el compañero ideal para aquellos que desean traducir conceptos matemáticos abstractos en código ejecutable. A lo largo de los capítulos, proporcionaremos Jupyter Notebooks que permitirán a los lectores realizar experimentos, verificar resultados y explorar nuevas ideas de manera interactiva.

Este libro tiene como objetivo el ser más que un manual académico; busca ser una guía para la aplicación de las matemáticas en Machine Learning usando el lenguaje Python, trascendiendo la mera teoría. Su contenido explora los modelos matemáticos, la formulación de predicciones y la implementación de estas en proyectos del campo.

El Lenguaje Python

El lenguaje de programación Python, se ha vuelto en uno de los lenguajes más populares en las últimas décadas, especialmente por su facilidad de uso, codificación y su amplia gama de aplicaciones en ciencias. Python es el lenguaje líder en ciencia de datos y es la herramienta principal que utiliza un científico de datos. Su simplicidad en la manipulación de datos, curva de aprendizaje accesible, versatilidad de aplicación y gran diversidad de módulos, lo han convertido en el lenguaje para la

⁴Haaser, N. La Salle, J. & Sullivan, J. (2000). “Análisis Matemático: Curso de Introducción Volumen 1”. México DF., México. Octavo Edición.

⁵Deitel, P & Deitel, H. (2015). “Java: How to Program”. New Jersey, United States of America. Pearson Education Inc.

Ingeniería de Datos, Ingeniería de Software, Ciencias de la Computación y Ciencia de Datos.

Python tiene una gran variedad de módulos de uso especializado, que se utilizan desde la construcción de bases de datos, procesamiento de datos, creación de modelos para su uso en Machine Learning y visualización de todo tipo de datos, como: NumPy, Pandas, Seaborn, TensorFlow y Scikit-learn. Las bibliotecas y frameworks para Python también proporcionan funcionalidades para el procesamiento de datos no estructurados, ya sean texto, audio, imágenes o videos.

La facilidad para aprender Python, lo ha convertido en una opción popular tanto entre programadores como entre no programadores, para desarrollar aplicaciones web, aplicaciones predictivas y servicios diversos, que a menudo se utilizan para obtener resultados de ciencia de datos a otros usuarios o sistemas.

A quien va dirigido el libro

Este libro está diseñado para estudiantes, profesionales y entusiastas de la ciencia de datos utilizando Python, pero para su uso exclusivo en Machine Learning. No es necesario tener un amplio conocimiento previo en programación o estadística, aunque contar con nociones básicas es beneficioso. El objetivo es guiar al lector desde los fundamentos hasta conceptos más avanzados, proporcionando las herramientas necesarias para desarrollarse como un científico de datos competente.

Aquellos que aspiren a una carrera en Ciencia de Datos, Ingeniería de Datos, Análisis de Datos o áreas afines o similares, encontrarán en este libro una base sólida. También es útil para profesionales de otras disciplinas que buscan aplicar técnicas de ciencia de datos en sus respectivos campos, como negocios, finanzas, salud, marketing, entre otros.

Este libro es apropiado para cursos de pregrado y posgrado en ciencias de la computación y disciplinas relacionadas que desean complementar sus conocimientos con Machine Learning, siguiendo las recomendaciones curriculares de la ACM/IEEE para programas de computación.⁶ Además, es una excelente opción para cursos en instituciones de educación superior, incluyendo universidades y colegios de dos años, que buscan preparar a los estudiantes para programas avanzados en ciencia de datos.

Asimismo, es relevante para estudiantes de secundaria interesados en programación y Ciencia de Datos. Cada vez más escuelas secundarias están incorporando cursos de ciencia de datos en sus planes de estudio, reconociendo la importancia

⁶Association for Computing Machinery (ACM). Curricula Recommendations. Recuperado de <https://www.acm.org/education/curricula-recommendations>.

de estas habilidades en la actualidad.⁷ Según expertos, enseñar Ciencia de Datos en niveles educativos tempranos prepara a los estudiantes para las carreras del futuro y refleja hacia dónde se dirige la tecnología y el mercado laboral. No olvidemos, además, que programar se ha vuelto en una herramienta indispensables para los profesionales del futuro.

Organización del libro

El libro, es una recopilación detallada y práctica en Python para el análisis de la matemática aplicada a los algoritmos de ML, donde:

Capítulo 1: Introducción a Machine Learning

El primer capítulo del documento proporciona una introducción comprensiva al campo del Machine Learning (ML), destacando su importancia creciente en la era digital y su aplicación en una variedad de disciplinas y sectores. El ML es un subcampo de la Inteligencia Artificial, que se centra en el desarrollo de algoritmos que hace a las máquinas, aprender y hacer proyecciones y decisiones basadas en datos. Se muestra la capacidad de ML para automatizar tareas analíticas, lo que antes requería intervención del hombre, o demasiados conocimientos especializados. Además, se señala cómo ML está revolucionando sectores como la educación, las ciencias y la manufacturación de productos. El capítulo describe tres tipos principales de aprendizaje en ML: supervisado, no supervisado y por refuerzo. El aprendizaje supervisado se trata de entrenar modelos en datos etiquetados; el no supervisado busca patrones en datos sin etiquetar, y el aprendizaje por refuerzo, consiste en tomar decisiones paso a paso con el objetivo, de maximizar alguna noción de recompensa. Las aplicaciones prácticas del ML son diversas y se extienden a través de numerosos dominios. El texto aborda ejemplos como sistemas de recomendación en plataformas de películas por internet, diagnósticos médicos, vehículos automatizados y que se conduzcan solos, y algunos análisis de sentimientos en diferentes tipos de redes sociales. Estos ejemplos ilustran la versatilidad y la relevancia del ML en resolver problemas complejos y en ofrecer soluciones innovadoras en el mundo moderno. El capítulo también plantea reflexiones sobre los desafíos asociados con el ML, que incluye la calidad de los datos y la cantidad de los mismos, la interpretabilidad de los modelos entrenados, y los retos computacionales tales como el big data. Además, se discuten las implicaciones éticas y sociales del ML, como el sesgo en los datos y en los algoritmos, la privacidad de los datos, y la toma de decisiones automatizada.

Capítulo 2: Programación en Python para Machine Learning

⁷“Newark School of Data Science & Information Technology”. Recuperado de <https://www.nps.k12.nj.us/data-science/>

El segundo capítulo se adentra en el uso de Python como herramienta para la búsqueda de soluciones relacionadas con ML. Python se presenta como un lenguaje de programación dominante en el campo del Machine Learning, gracias a su flexibilidad y la vasta disponibilidad de bibliotecas especializadas. El capítulo argumenta que Python facilita la experimentación y la implementación rápida de modelos de ML, haciendo del lenguaje una opción popular entre los profesionales y los académicos. Se hace hincapié en bibliotecas como NumPy, para cálculos numéricos eficientes; Pandas, para la manipulación y análisis de datos; y Scikit-Learn, una biblioteca integral para el ML. Cada una de estas herramientas se describe en términos de sus funcionalidades y cómo habilitan diferentes aspectos del flujo de trabajo en ML, desde la preparación de datos hasta la construcción y evaluación de modelos. El texto resalta la compatibilidad de Python con otras tecnologías y plataformas, lo que permite una integración fluida en sistemas más grandes y en flujos de trabajo de ciencia de datos. Se mencionan ejemplos de cómo Python se utiliza en conjunción con las bases de datos, también con las aplicaciones diseñadas para la web y plataformas de uso en la nube. Sin embargo, en la parte final del libro, hay un apartado de programación Python, mas extensa y complementaria para aprender los conceptos mas profundos de manera directa y concreta.

Capítulo 3: Python para Análisis de Datos

En este capítulo, el foco se pone en los módulos específicos de Python utilizados en el análisis de datos. Se detallan las características y usos de módulos como NumPy y Pandas, mostrando cómo facilitan operaciones complejas con datos. NumPy se presenta como la columna principal para operaciones matemáticas de alto rendimiento, mientras que Pandas ofrece estructuras de datos y herramientas de análisis flexibles y potentes. El capítulo también aborda herramientas de visualización como Matplotlib y Seaborn, destacando su importancia para explorar y entender los datos. Se discute cómo estas herramientas ayudan a crear representaciones gráficas intuitivas de los datos, la exploración, el diagnóstico de problemas encontrados y el señalamiento de resultados. El texto muestra ejemplos de cómo estas herramientas se aplican en escenarios de análisis de datos, desde el preprocesamiento y preparación de datos hasta su análisis y visualización. Se enfatiza cómo el dominio de estas herramientas es vital para cualquier profesional del ML y la ciencia de datos.

Capítulo 4: Visualización de Datos

Este capítulo se centra en las técnicas de visualización de datos en ML. La visualización se describe como un aspecto crítico del análisis de datos, que permiten descubrir patrones y anomalías. El capítulo argumenta que una visualización efectiva puede ofrecer insights que de otro modo podrían pasar desapercibidos, y puede ser una herramienta útil para extraer conocimiento de los datos. Se explora una

variedad de técnicas de visualización, desde gráficos básicos como histogramas y gráficos de dispersión hasta visualizaciones más avanzadas tales como los útiles mapas de calor y gráficos relacionados con las redes. Se discuten las herramientas específicas de Python que facilitan estas visualizaciones, proporcionando ejemplos y aplicaciones prácticas. El capítulo también considera cómo la visualización se integra en el flujo de trabajo, ayudando en la comprensión de los datos y en la interpretación de los resultados de los modelos. Se examina cómo las visualizaciones pueden ayudar en la evaluación de modelos y en la presentación de resultados a una audiencia no técnica.

Capítulo 5: Aritmética, Álgebra, Geometría y Trigonometría

El quinto capítulo sirve como una introducción a los conceptos matemáticos claves necesarios en ML, ofreciendo una base sólida para entender los algoritmos y técnicas avanzadas. Se comienza con los fundamentos de la aritmética y el álgebra, explicando cómo estos conceptos se aplican en ML. Se destacan las operaciones aritméticas básicas, las propiedades algebraicas, y su importancia en la formulación y comprensión de los algoritmos de ML. El capítulo aborda la relevancia de la geometría y la trigonometría en ML, particularmente en la comprensión de las relaciones espaciales y las transformaciones en los datos. Se explica cómo estos conceptos ayudan en la visualización y el análisis de datos multidimensionales. Se proporcionan ejemplos de cómo los conceptos aritméticos, algebraicos, geométricos y trigonométricos se aplican en problemas de ML. Esto incluye el uso de estas matemáticas en la construcción de modelos, la búsqueda de la optimización de algoritmos y la búsqueda de resultados. Así, proporciona una visión general detallada de los primeros cinco capítulos del documento, abarcando desde la introducción y las bases del Machine Learning hasta las herramientas y matemáticas necesarias para su implementación y comprensión.

Capítulo 6: Álgebra Lineal, Cálculo y Ecuaciones Diferenciales

El álgebra lineal es utilizada para procesar datos en Machine Learning. Se utilizan vectores y matrices para representar datos y realizar transformaciones. Operaciones como la multiplicación de matrices son utilizadas comúnmente para preprocesamiento y normalización de datos. Los sistemas de ecuaciones lineales, son centrales en muchos algoritmos de ML, especialmente en técnicas como la regresión lineal y la clasificación. El cálculo proporciona herramientas para comprender y manejar cambios y tasas de cambio en los datos, muy usados en la optimización de modelos. Conceptos como gradientes y derivadas son relevantes en el contexto en algoritmos de optimización como el Gradiente Descendiente. Las ecuaciones diferenciales permiten modelar la evolución de sistemas complejos y fenómenos que cambian con el tiempo, como la propagación de información en redes neuronales o la simulación

de epidemias.

Capítulo 7: Probabilidad y Estadística

En este capítulo se introducen los conceptos de la teoría de probabilidades, que son usados especialmente en estadística y Machine Learning. Se usa Python, con bibliotecas como NumPy y SciPy para aplicar estos conceptos en ejemplos prácticos. Se exploran distribuciones de probabilidades, tanto discretas como continuas, enfocándose en binomiales y normales. La estadística descriptiva ayuda a describir datos en diferentes tablas, y la estadística inferencial permite inferenciar sobre un conjunto de datos poblacionales, donde estos son a partir de muestras.

Capítulo 8: Módulos de Python para Machine Learning

Se exploran diversas librerías clave en Python para Machine Learning, como Scikit-Learn, TensorFlow, PyTorch, Keras, XGBoost, LightGBM y Statsmodels. Estas herramientas facilitan desde la construcción y entrenamiento de modelos hasta el análisis estadístico avanzado.

Capítulo 9: Algoritmos Supervisados en Machine Learning

Se abordan diversos algoritmos supervisados, incluyendo múltiples formas de regresión (lineal, ridge, LASSO, entre otros) y clasificación (lineal, logística, SVM, etc.). Se exploran decisiones mediante los algoritmos de árboles de decisiones, mediante los bosques aleatorios y técnicas de ensamblaje, como bagging y boosting.

Capítulo 10: Algoritmos no Supervisados en Machine Learning

Se estudian técnicas de clustering como K-means y K-medoids, y métodos de reducción de dimensionalidad, incluyendo PCA y SVD. Se centra en descubrir todos aquellos patrones y estructuras, en base a datos no estructurados, abordando técnicas de análisis de tópicos, filtros colaborativos y autoencoders.

Cada capítulo aporta conocimientos base y prácticos que son para la comprensión del campo del Machine Learning, abarcando desde los aspectos teóricos y matemáticos hasta las aplicaciones prácticas con Python y su ecosistema de librerías. Juntos, proporcionan una comprensión integral de las herramientas, técnicas y teorías que sustentan la práctica moderna del Machine Learning.

Software utilizado en el libro

Todo el software que necesitará para este libro está disponible para Windows, macOS y Linux y es gratuito para descargar desde Internet. Escribimos los ejemplos del libro utilizando la distribución gratuita de Python Anaconda. Incluye la mayoría de las bibliotecas de Python, visualización y ciencia de datos que necesitará, así como Python, el intérprete IPython, Jupyter Notebooks y Spyder, considerado uno de los mejores entornos de desarrollo integrado (IDE) para ciencia de datos en Python—en el libro usamos solo IPython y Jupyter Notebooks para el desarrollo de programas. La sección “**Antes de Comenzar**” discute la instalación de Anaconda y otros elementos que necesitará para trabajar con nuestros ejemplos.

Documentación de Python y Recursos de Ciencia de Datos

Encontrará la siguiente documentación especialmente útil mientras avanza en el libro:

- **La biblioteca estándar de Python:**
<https://docs.python.org/3/library/index.html>
- **La referencia del lenguaje Python:**
<https://docs.python.org/3/reference/index.html>
- **Lista de documentación de Python:**
<https://docs.python.org/3/>

Además, los siguientes recursos y documentación son especialmente relevantes para la ciencia de datos:

- **Documentación de NumPy** (computación numérica con Python):
<https://numpy.org/doc/>
- **Documentación de Pandas** (análisis y manipulación de datos):
<https://pandas.pydata.org/docs/>
- **Documentación de Matplotlib** (visualización de datos en 2D):
<https://matplotlib.org/stable/contents.html>
- **Documentación de SciPy** (algoritmos y herramientas matemáticas):
<https://docs.scipy.org/doc/scipy/reference/>
- **Documentación de Scikit-learn** (aprendizaje automático en Python):
https://scikit-learn.org/stable/user_guide.html
- **Documentación de Seaborn** (visualización estadística de datos):
<https://seaborn.pydata.org/>
- **Documentación de Plotly** (visualización interactiva de datos):
<https://plotly.com/python/>
- **Documentación de TensorFlow** (ML y redes neuronales):
<https://www.tensorflow.org/overview>
- **Documentación de PyTorch** (aprendizaje profundo y redes neuronales):

<https://pytorch.org/docs/stable/index.html>

- **Documentación de Keras** (API para redes neuronales):

<https://keras.io/>

- **Documentación de StatsModels** (modelado estadístico y econométrico):

<https://www.statsmodels.org/stable/index.html>

Estos recursos le proporcionarán información detallada y ejemplos sobre las bibliotecas y herramientas más utilizadas en ciencia de datos con Python.

Cómo contactarnos

Por favor, sus comentarios y preguntas sobre este libro al editor:

Bluminds LLC 20 S Charles St Ste 403 #878
Baltimore, MD 21201 - USA

Se tiene una página web del libro, donde listamos erratas y cualquier información adicional. Puede acceder a esta página en *<https://books.bluminds.com>*

Para realizar preguntas técnicas sobre este libro, envíe un correo electrónico a *info@bluminds.com*

Más información sobre nuestros libros y cursos visite nuestro sitio web en *<https://www.bluminds.com>*

- Síguenos en Twitter (x.com): *<https://x.com/blumindsllc>*
- Vea nuestros videos en YouTube: *<https://www.youtube.com/@bluminds>*
- Código Fuente: *<https://github.com/Bluminds/Book-MMLv1-CodeExamples>*

Antes de Comenzar

Con la aparición de mayores y mejores conceptos y aplicaciones de Inteligencia Artificial, los profesionales se están quedando fuera del mercado laboral si es que no conocen sobre de la programación y los lenguajes de programación,⁸ y Python es el lenguaje más sencillo y potente para que un profesional amplie y mejore sus habilidades.

Aquí asumimos el gran desafío de seguir desarrollando habilidades matemáticas integradas con programación con Python!

No es necesario saber programar para usar este libro, pero si es necesario conocer el entorno de Jupyter Notebooks en Anaconda, ya que ahí se han de desarrollar los scripts para programar con Python y obtener los resultados que se buscan. Aquí, se muestran los pasos principales que se deben saber antes de empezar con la lectura del libro, donde el lector debe procurar llegar al final de este apartado y poder realizar los ejercicios que aquí se muestra antes de empezar con el capítulo 1. En la parte final de este libro se podrán ver algunos conceptos de Programación Python, que el lector debe ver para aprender más sobre los conceptos de Python, pero conceptos genéricos; existe un apartado llamado “*Apéndice A: Introducción a Python*”.

Así, empecemos con la instalación de Python.

Instalación de Anaconda y Python

La instalación de Anaconda y Python constituye la etapa inicial indispensable para integrarse de manera efectiva en los ámbitos académico, profesional y de la programación científica. En el contexto de este libro, dicha instalación habilita el acceso al entorno computacional idóneo para el aprendizaje de matemáticas aplicadas a la inteligencia artificial.

Anaconda, una distribución de Python, simplifica significativamente la configuración del entorno de desarrollo al proporcionar una colección completa de herramientas y bibliotecas esenciales que son fáciles de utilizar. Aquí se explorará detalladamente el proceso de instalación de Anaconda y sus beneficios, asegurando una experiencia de desarrollo coherente y eficiente. Anaconda incluye Jupyter, una aplicación web que permite crear documentos que combinan código, texto y visualizaciones. Esto es útil para el análisis de datos y la presentación de resultados. La aplicación de Anaconda está disponible para diversos sistemas operativos

⁸Deitel, P & Deitel, H. (2015). “Java: How to Program”. New Jersey, United States of America. Pearson Education Inc.

como Windows, MacOS y Linux, lo que garantiza una experiencia de desarrollo consistente sin importar la plataforma utilizada.

¿Qué es Anaconda?

Anaconda distribuye el lenguaje Python, que incluye el lenguaje de programación en sí, y también un conjunto de herramientas y bibliotecas específicamente diseñadas para el análisis de datos y la programación científica. Desarrollada por Anaconda Inc., esta distribución más moderna, tiene como objetivo, la gestión de paquetes y entornos virtuales, haciendo que el proceso de desarrollo sea más accesible para usuarios de todos los niveles.

Pasos para la instalación:

La instalación de Anaconda y Python es un proceso sencillo que se puede realizar en pocos pasos, pero dependiendo de la computadora, esta puede tardar menos o más. Aquí, supondremos que su computador no tiene instalado Anaconda, por lo que le proporcionaremos una guía, paso a paso junto con referencias, para una instalación exitosa.

Descargar Anaconda:

Visita el sitio web oficial de Anaconda para descargar la versión más reciente:

(<https://www.anaconda.com/products/distribution>)

Selecciona (*DOWNLOAD*) la versión adecuada para tu sistema operativo (Windows, macOS o Linux) y no olvides que tu computador debe cumplir ciertos requisitos mínimos. De todas maneras, existen varias versiones que, si tu computador no cumple algunos, otra versión puede cumplirlas. Aquí los requisitos mínimos son:

DEL COMPUTADOR

- Se recomienda un procesador moderno de 1 GHz para un rendimiento optimo.
- Se recomienda tener 2 GB de RAM, pero con mayor memoria, se mejora la velocidad de las ejecuciones.
- Se recomienda tener un mínimo de 3 GB de espacio libre en el disco duro, más espacio para procesos de ejecución donde se preste memoria el disco.

DEL WINDOWS

- Anaconda es compatible con Windows 7 o superiores.
- Se recomienda un espacio libre de 10 GB en el disco para las instalaciones y las creaciones de entornos virtuales.

DEL MACOS

- Anaconda es compatible con macOS 10.9 o superiores.
- Se recomienda espacio libre de al menos 10 GB en el disco para instalaciones y creaciones de entornos virtuales.

DEL LINUX

- Anaconda es compatible con varias distribuciones de Linux, Ubuntu, Fedora y otros.
- Se recomienda espacio libre de al menos 10 GB en el disco para instalaciones, pero puede variar según la distribución.

Instalación de Anaconda:

- Ejecuta el instalador de Anaconda, descargado en tu PC.
- Seguir todas las instrucciones del instalador, además de la aceptando de los términos y condiciones de uso.
- Puedes optar por agregar Anaconda al PATH durante la instalación, lo que facilitará la ejecución de comandos desde la línea de comandos.

Una vez completada la instalación, se debe verificar que Anaconda y Python se hayan instalado correctamente. Es posible realizarlo abriendo una terminal, o mediante líneas de comandos, y después ejecutar el siguiente comando:

```
conda --version
```

Este comando muestra la versión de Anaconda, que es el administrador de paquetes incluido en Anaconda.

Iniciar la aplicación:

Hacer doble clic en la aplicación.



Figura 1: Icono de Anaconda

Y esperar hasta que se inicie el programa.

Manejo de Anaconda y Python

Una vez iniciado se tendrá un entorno de Anaconda, que es similar al que se muestra en la siguiente figura.

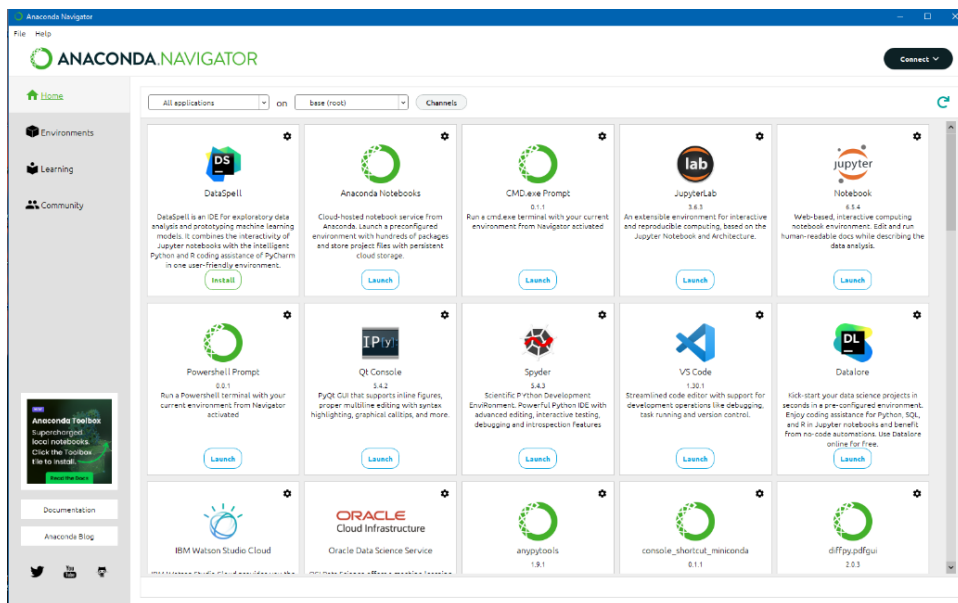


Figura 2: Ventana de Anaconda Navigator

Una vez que se tiene esta imagen, se ejecuta la aplicación que dice Jupyter, que en la figura está en la parte superior derecha, mediante “Launch”.

Aquí, es el administrador de paquetes de la aplicación Anaconda, donde simplifica la instalación, la actualización y la gestión de librerías y sus dependencias, haciendo que sea mucho más fácil tener el lenguaje en óptimas condiciones para su uso. Además, permite la creación de entornos virtuales para proyectos específicos, asegurando la consistencia en el desarrollo. Anaconda incluye una variedad de bibliotecas esenciales para el análisis de datos, para la programación, y para el aprendizaje de otras ciencias, o como en este caso, para aprender más sobre Python.⁹

Configuración del Entorno de Trabajo

Para asegurar la compatibilidad y el funcionamiento óptimo de todas las bibliotecas y códigos utilizados en este libro, es recomendable crear un entorno virtual específico denominado bluminds con Python 3.12. Los entornos virtuales permiten mantener las dependencias del proyecto aisladas y evitar conflictos entre diferentes versiones de bibliotecas.

⁹Venero B. (1999). “Análisis Matemático 2”. Lima, Perú. Ediciones Gemar.

La creación de este entorno se realiza desde la línea de comandos o terminal. Para acceder a la terminal en Windows, se puede usar el Anaconda Prompt, mientras que en macOS y Linux se utiliza la terminal del sistema. Una vez abierta la terminal, se ejecutan los siguientes comandos:

Crear el entorno bluminds:

Para crear el entorno virtual con Python 3.12, ejecute el siguiente comando:

```
conda create --name bluminds python=3.12
```

Este comando creará un nuevo entorno virtual llamado `bluminds` con Python 3.12 instalado. Durante el proceso, el sistema solicitará confirmación para proceder con la instalación.

Activar el entorno:

Una vez creado el entorno, debe activarse antes de instalar las bibliotecas o ejecutar código Python. Para activar el entorno, use:

```
conda activate bluminds
```

Cuando el entorno esté activo, notará que el nombre (`bluminds`) aparece al inicio del prompt de la terminal, indicando que está trabajando dentro del entorno virtual.

Verificar la instalación:

Para verificar que Python 3.12 se instaló correctamente en el entorno, ejecute:

```
python --version
```

Este comando debe mostrar una versión de Python 3.12.x.

Instalar Jupyter en el entorno:

Para poder usar Jupyter Notebooks dentro del entorno `bluminds`, es necesario instalar Jupyter:

```
conda install jupyter
```

Registrar el entorno como kernel de Jupyter:

Para que el entorno bluminds aparezca como opción al crear nuevos notebooks, debe registrarse como kernel:

```
python -m ipykernel install --user --name bluminds
--display-name "Python 3.12 (bluminds)"
```

Una vez completados estos pasos, el entorno bluminds estará disponible en Jupyter Notebook, y podrá seleccionarse al crear un nuevo notebook, asegurando que todos los ejercicios y códigos del libro funcionen correctamente con las versiones específicas de las bibliotecas.

Desactivar el entorno:

Cuando termine de trabajar, puede desactivar el entorno ejecutando:

```
conda deactivate
```

Creación de un Jupyter Notebook

Una vez que se haya presionado “Launch”, se ha de desplegar una imagen, que es similar a la que se muestra en la siguiente figura.



Figura 3: Inicio de Jupyter Notebook

Ahora, en la parte derecha de la pantalla de la aplicación, buscamos la palabra “NEW” y escogemos la opción que dice “Python 3 (ipykernel)”, presionándolo para que despliegue un primer entorno de Jupyter Notebooks. Se espera unos segundos y se tiene un entorno inicial de Jupyter Notebooks.

Jupyter Notebooks

El primer entorno de Jupyter Notebooks es similar al que se muestra en la siguiente figura.

Ahora está listo para realizar sus primeros programas, no se preocupe por los demás botones y opciones, eso lo ira aprendiendo en el transcurso de la lectura de este libro, o en libros más especializados, donde desea conocer más sobre el entorno de Jupyter Notebooks, para lo cual le recomendamos el libro de nuestra colección de Programación Python.

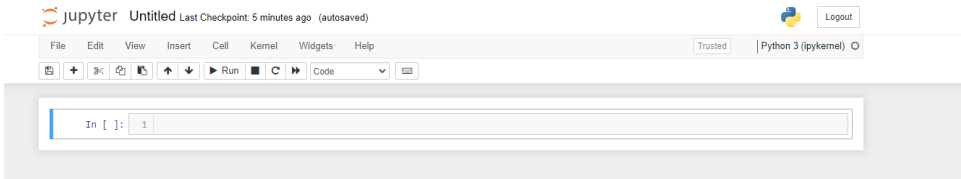


Figura 4: Ventana Práctica de Jupyter Notebook

Cinco ejercicios básicos con Jupyter Notebooks

Ahora hemos de realizar las siguientes operaciones matemáticas simples, con el objetivo de ver y comprobar que Python está ejecutando el código correctamente. No olvide colocar los datos correctamente, respetando los espacios, los signos y símbolos, así:

Ejercicio 1:

Haremos un primer ejercicio de sumas y restas de números complejos, colocando los scripts en la celda de color gris claro, y una vez que se haya terminado, se presiona **RUN**.

In [1]:

```
1 # Ejercicio 1
2 # Suma y Resta de Números Complejos
3 num1 = 3 + 4j
4 num2 = 1 - 2j
5
6 # Suma
7 suma = num1 + num2
8
9 # Resta
10 resta = num1 - num2
11
12 suma, resta
```

Out [1]:

```
((4+2j), (2+6j))
```

Ejercicio 2:

Haremos un segundo ejercicio de resolución de ecuaciones lineales, colocando los scripts en la celda gris claro y una vez terminado, se presiona **RUN**.

In [2]:

```
1 # Ejercicio 2
2 # Resolución de Ecuaciones Lineales
3 import numpy as np
4 import warnings
5
6 warnings.filterwarnings("ignore")
7
8 coeficientes = np.array([[3, 2], [4, -1]])
9 constantes = np.array([9, 5])
10
11 solucion = np.linalg.solve(coeficientes, constantes)
12 solucion
```

Out [2]:

```
array([1.72727273, 1.90909091])
```

Ejercicio 3:

Haremos un tercer ejercicio de polinomios cuadráticos, colocando los scripts en la celda y una vez terminado, se vuelve a presionar **RUN**.

In [3]:

```
1 # Ejercicio 3
2 # Factorización de Polinomios Cuadráticos
3 from sympy import symbols, factor
4
5 x = symbols("x")
6 polinomio = x**2 - 5 * x + 6
7
8 factorizado = factor(polinomio)
9 factorizado
```

Out [3]:

```
(x - 3)*(x - 2)
```

Ejercicio 4:

Ahora un cuarto ejercicio de operaciones con matrices, colocando los scripts en la celda gris claro y una vez terminado, se vuelve a presionar **RUN**.

In [4]:

```
1 # Ejercicio 4
2 # Operaciones con Matrices
3 matriz_a = np.array([[1, 2], [3, 4]])
4 matriz_b = np.array([[5, 6], [7, 8]])
5
6 # Multiplicación de Matrices
7 multiplicación = np.dot(matriz_a, matriz_b)
8
9 # Transposición de Matrices
10 transpuesta_a = matriz_a.T
11
12 multiplicación, transpuesta_a
```

Out [4]:

```
(array([[19, 22],
        [43, 50]]),
 array([[1, 3],
        [2, 4]]))
```

Ejercicio 5:

Haremos un quinto ejercicio de despejar variables en ecuaciones, colocando los scripts en la celda gris claro y una vez terminado, se vuelve a presionar **RUN**.

In [5]:

```
1 # Ejercicio 5
2 # Despejar Variables en Ecuaciones
3 from sympy import symbols, solve
4
5 x, y = symbols("x y")
6 ecuacion = 2 * x + 3 * y - 12
7
8 # Despejar y
9 solucion_y = solve(ecuacion, "y")
10 solucion_y
```

Out [5]:

```
[4 - 2*x/3]
```

Una vez que realizados los anteriores pasos correctamente, se encuentra listo para empezar a usar Python con este libro, caso opuesto vuelva a realizar las operaciones hasta lograr los mismos resultados.

Como actualizar Anaconda, Python y sus librerías

Actualizar Anaconda, Python y sus librerías desde un Jupyter Notebook, se realiza utilizando las celdas y poniendo scripts en ella, por lo que, para actualizar Anaconda se coloca en una de las celdas, lo siguiente:

```
!conda update --all
```

Para actualizar Python se coloca:

```
!conda install python=3.12 --update-deps
```

Y para actualizar una biblioteca especifica se coloca

```
!pip install -- upgrade <MODULO>
```

Librerías que deben Instalarse en Python

Las librerías utilizadas en el presente libro, deben estar instaladas en Jupyter Notebook, a continuacion se listan los comandos para su instalación desde una celda en Jupyter Notebook.

```
!pip install numpy
!pip install sympy
!pip install scipy
!pip install pandas
!pip install matplotlib
!pip install numdifftools

!pip install openpyxl
!pip install sqlalchemy

!pip install statsmodels
!pip install scikit-learn
!pip install tensorflow
!pip install keras
!pip install torch
!pip install mlxtend
!pip install xgboost
!pip install lightgbm
!pip install opencv-python
```

```
!pip install torchvision

!pip install seaborn
!pip install bokeh
!pip install plotly
!pip install dash
!pip install dash-core-components
!pip install dash-html-components

!pip install cartopy
!pip install folium
!pip install networks
!pip install basemap
!pip install geopandas==0.14.4

!pip install hmmlearn
!pip install filterpy
!pip install squarify
```

... aqui se procede a la instalacion de cada una de las librerias con pip



Buena Práctica

Las versiones de las librerías de Python cambian en el tiempo, y algunas funciones pueden quedar obsoletas (deprecadas) debido a mejoras en rendimiento, seguridad, o inclusión de nuevas funcionalidades. Es una buena practica mantener las dependencias (librerías o módulos) actualizados en el sistema.

Otras buenas prácticas respecto al manejo de dependencias y actualizaciones se tienen:

- **Instalar versiones específicas:** Cuando sea necesario usar una versión particular de una librería (como `geopandas==0.14.4`), especifique la versión exacta para garantizar compatibilidad.
- **Verificar versiones instaladas:** Use `pip list` o `conda list` para ver todas las librerías instaladas y sus versiones actuales en su entorno.
- **Consultar versiones disponibles:** Visite <https://pypi.org> para verificar las versiones más recientes de cualquier librería de Python y revisar el historial de cambios.
- **Manejo de deprecaciones:** Si encuentra advertencias sobre funciones deprecadas, consulte la documentación oficial de la librería para encontrar las alternativas recomendadas.
- **Actualización gradual:** Actualice las librerías una por una y pruebe su código

después de cada actualización para identificar posibles incompatibilidades.

- **Entornos virtuales:** Use entornos virtuales separados, para diferentes proyectos y evitar conflictos entre versiones entre librerías.

Resumen y Preguntas

Cada uno de los capítulos de este libro, contienen un apartado llamado “*Resumen*”, para que el lector pueda verificar de manera resumida los avances del capítulo, y además contiene otros apartados de ejercicios que están resueltos (ver el final del libro), y otro de ejercicios no resueltos, para que el lector pueda practicar.

Este libro contiene “Matemática Aplicada para Machine Learning”, por lo que es útil tener conocimientos de programación en Python para un mejor aprendizaje.

Capítulo 1 Introducción a Machine Learning

Objetivos

- ❑ Conceptos básicos de Machine Learning e Inteligencia Artificial (IA).
- ❑ Tipos de aprendizaje: supervisado, no supervisado y por refuerzo.
- ❑ Aplicaciones contemporáneas de ML e IA generativa.
- ❑ Distinción conceptual entre Machine Learning y Data Mining.
- ❑ Rol de la ciencia de datos, las matemáticas y la estadística en ML.
- ❑ Bibliotecas esenciales para el desarrollo en Python.
- ❑ Fases metodológicas de un proyecto basado en ML.

1.1 ¿Qué es Machine Learning?

Machine Learning, también denominado *aprendizaje automático o aprendizaje de máquinas* en español, constituye una vertiente moderna y transformadora de la Inteligencia Artificial. En este marco, la Inteligencia Artificial se erige como el concepto central, dentro del cual Machine Learning ocupa un lugar preponderante, sustentado en sólidos fundamentos teóricos, particularmente en el ámbito matemático.

El campo de Machine Learning se orienta al desarrollo de algoritmos y modelos que capacitan a los sistemas computacionales para aprender patrones a partir de datos y ejecutar tareas sin requerir una programación explícita.¹ A diferencia de los enfoques tradicionales, en los que las máquinas siguen instrucciones específicas, Machine Learning habilita a los sistemas a mejorar su rendimiento mediante la experiencia y la adaptación continua.



Información Relevante

¿Por qué estudiar Machine Learning? Esta tecnología impulsa innovaciones cotidianas como los asistentes virtuales (Alexa, Siri), sistemas de recomendación (Netflix, Spotify), diagnósticos médicos automatizados, vehículos autónomos y análisis de mercados financieros. Comprender ML te posiciona en la vanguardia tecnológica actual.

En este contexto, Machine Learning se clasifica convencionalmente en tres

¹ Deitel, P & Deitel, H. (2015). “Java: How to Program”. New Jersey, United States of America. Pearson Education Inc.

áreas principales de aprendizaje:

1. Aprendizaje supervisado:

En el aprendizaje supervisado, el modelo se fundamenta en algoritmos específicos y se entrena utilizando conjuntos de datos etiquetados, es decir, datos en los que tanto las variables dependientes como las independientes están claramente identificadas. Cada entrada de datos se asocia con una etiqueta que representa la respuesta esperada. Por ejemplo, en un modelo de clasificación de imágenes, cada imagen podría estar etiquetada como “perro” o “gato”, y las variables independientes corresponderían a características propias de la imagen. El modelo matemático utiliza estas etiquetas para aprender a predecir la respuesta correcta en nuevas instancias no observadas durante el entrenamiento.

2. Aprendizaje no supervisado:

En contraste, el aprendizaje no supervisado aborda conjuntos de datos no etiquetados, en los cuales las categorías o respuestas no están previamente definidas. En este escenario, el modelo debe identificar patrones y estructuras subyacentes de manera autónoma, sin la guía de etiquetas predefinidas. Un ejemplo paradigmático es el análisis de agrupamiento, donde el modelo organiza los datos en grupos según su similitud, sin conocer de antemano las categorías a las que pertenecen.

3. Aprendizaje por refuerzo:

El aprendizaje por refuerzo se caracteriza por la toma de decisiones del modelo en función de su interacción con el entorno, recibiendo retroalimentación en forma de recompensas o penalizaciones. Este enfoque es análogo al proceso de aprendizaje humano, donde las acciones se refuerzan mediante experiencias positivas o negativas. Por ejemplo, en un entorno de videojuegos, un agente de aprendizaje por refuerzo puede recibir recompensas al tomar decisiones que favorecen el progreso, y penalizaciones por acciones desfavorables.

Cabe preguntarse si existen otros tipos de aprendizaje además de los tres mencionados. La respuesta es afirmativa; sin embargo, la literatura científica aún no ha establecido una clasificación específica y universalmente aceptada para otros enfoques emergentes. Por el momento, el análisis se centrará en estos tres paradigmas principales.

Para fomentar la comprensión y el aprendizaje de los algoritmos, resulta pertinente ilustrar las aplicaciones prácticas de estos tipos de aprendizaje. Algunos ejemplos representativos incluyen:

- **Reconocimiento de voz:** El desarrollo de sistemas como asistentes virtuales (por ejemplo, Siri, Alexa) se apoya en Machine Learning para optimizar la precisión en el reconocimiento de voz, a medida que el sistema se expone a volúmenes crecientes de datos. Esta área se estudia bajo el campo de “*Procesamiento del Lenguaje Natural (NLP)*”.
- **Recomendaciones en línea:** Plataformas como Netflix y Amazon emplean algoritmos de Machine Learning para analizar el comportamiento de los usuarios y ofrecer recomendaciones personalizadas de películas, series, productos o servicios. Esta aplicación se enmarca en la disciplina de “*Sistemas de Recomendación*”.
- **Diagnósticos médicos:** En el ámbito de la salud, los modelos de Machine Learning permiten analizar datos médicos para asistir en el diagnóstico de enfermedades y la predicción de tratamientos. Esta línea de investigación se denomina “*Procesamiento de Señales Médicas*”.
- **Conducción autónoma:** En la industria automotriz, Machine Learning se emplea en sistemas de conducción autónoma para interpretar y responder a situaciones en tiempo real, como la detección de obstáculos y la toma de decisiones en rutas. Esta área se conoce como “*Dispositivos Autónomos*”.
- **Reconocimiento de imágenes y rostros:** Esta área, denominada “*Visión por Computadoras*”, que se caracteriza por utilizar algoritmos de Machine Learning para dotar a las máquinas de capacidades de interpretación y comprensión visual. El campo abarca desde tareas elementales como la detección de bordes hasta procesos complejos, como los reconocimientos faciales y tridimensionales o en secuencias de video.
- **Redes neuronales artificiales:** El Deep Learning constituye una rama especializada de Machine Learning, centrada en el entrenamiento de modelos basados en redes neuronales profundas, tanto supervisadas como no supervisadas y de refuerzo. Este enfoque ha desarrollado múltiples campos de la investigación, ya que las redes profundas permiten el aprendizaje de representaciones jerárquicas de características, que hace más fácil la identificación de patrones complejos en datos visuales.
- **Robótica:** En el ámbito de la automatización y la mecánica, Machine Learning se aplica para potenciar las capacidades de los robots, permitiéndoles percibir su entorno, tomar decisiones y ejecutar tareas de manera autónoma. Esto abarca la planificación de movimientos, la navegación, la interacción con objetos y personas, así como la adaptación a entornos dinámicos. Los algoritmos de Machine Learning posibilitan que los robots aprendan de la experiencia y mejoren sus capacidades progresivamente.
- **Biomimética:** En este contexto, Machine Learning se inspira en principios observados en la naturaleza para el desarrollo de algoritmos y modelos matemáticos de aprendizaje. Esto incluye el diseño de soluciones basadas en comportamientos de organismos vivos. La rama correspondiente, “*Algorit-*

mos Genéticos”, busca imitar el proceso de evolución biológica para hallar soluciones óptimas, previamente optimizadas por la naturaleza.

- **Preferencias individuales:** Machine Learning se emplea en el desarrollo de interfaces de usuario inteligentes y adaptativas, orientadas a optimizar la interacción entre humanos y computadoras mediante la comprensión de patrones de comportamiento en la interacción humano-máquina.
- **Generative IA:** La “*Inteligencia Artificial Generativa*” se orienta a la creación de datos novedosos y originales, en contraposición a tareas tradicionales como el reconocimiento de patrones o la toma de decisiones. La IA generativa puede emplearse para la generación de contenido textual original (poesía, historias, noticias), mediante modelos de lenguaje como GPT (Generative Pre-trained Transformer), entrenados para producir texto coherente y relevante. En los ámbitos del arte y el diseño, la IA generativa posibilita la creación de obras artísticas, ilustraciones y productos de diseño, aprendiendo estilos y generando creaciones únicas. Asimismo, puede componer música original, crear imágenes realistas y animaciones mediante modelos como GANs (Redes Generativas Adversarias), y desarrollar chatbots avanzados capaces de mantener conversaciones contextuales. En aplicaciones de simulación, la IA generativa permite la creación de escenarios realistas para entrenamiento en contextos militares o médicos, y en la investigación científica, facilita el diseño de nuevos materiales y moléculas con propiedades específicas. Al igual que en Machine Learning, la IA generativa contribuye a la mejora de interfaces de usuario, promoviendo diseños inteligentes y adaptativos para optimizar la interacción humano-sistema.

En síntesis, Machine Learning constituye un campo amplio y dinámico que sigue transformando múltiples industrias. Desde la toma de decisiones inteligente hasta la automatización de tareas complejas, las aplicaciones de Machine Learning son vastas y se encuentran en continua expansión.

1.2 Machine Learning vs Data Mining

Machine Learning (ML) y Data Mining (DM), términos que suelen emplearse de manera intercambiable, aunque presentan diferencias conceptuales sustanciales, poseen objetivos y enfoques propios. Estos contribuyen de manera diferenciada al procesamiento y análisis de datos. Machine Learning se caracteriza por su capacidad para desarrollar algoritmos y modelos que permiten a los sistemas computacionales aprender patrones a partir de datos, sin requerir una programación explícita de reglas y decisiones. En este paradigma, el modelo se adapta y mejora conforme se expone a nuevos datos, lo que lo vincula estrechamente con la formación de modelos predictivos basados en algoritmos.

Por otro lado, Data Mining constituye una rama de la Ciencia de Datos orientada

a la exploración exhaustiva de conjuntos de datos, desde pequeños volúmenes hasta grandes bases de datos, con el propósito de descubrir patrones, tendencias y relaciones que no son evidentes a simple vista.² A diferencia de ML, DM se enfoca principalmente en la extracción de información valiosa y conocimiento implícito, utilizando para ello metodologías y herramientas como Python para la identificación de patrones y tendencias.

En cuanto a la clasificación de algoritmos, Machine Learning abarca diversas categorías según la naturaleza del problema abordado, tales como el aprendizaje supervisado, no supervisado y por refuerzo.³ Estos enfoques son relevantes para la construcción de modelos predictivos que ejecutan tareas específicas como clasificación, regresión o agrupamiento, basándose en patrones aprendidos a partir de los datos de entrenamiento.⁴

En contraste, Data Mining recurre a técnicas estadísticas y matemáticas para analizar grandes volúmenes de datos. Entre estas técnicas se incluyen la minería de asociaciones para identificar patrones frecuentes, la clasificación para organizar datos en categorías específicas y la regresión para modelar relaciones entre variables.⁵ El objetivo de DM es comprender en profundidad los datos y fundamentar la toma de decisiones informadas a partir de la información descubierta.

Objetivos principales:

- **Machine Learning:** Se orienta a la construcción de modelos predictivos capaces de realizar tareas específicas de manera autónoma, permitiendo que el sistema aprenda y ejecute dichas tareas sin intervención programática directa.
- **Data Mining:** Se enfoca en la exploración y el descubrimiento de patrones e información relevante en grandes conjuntos de datos, utilizando técnicas que facilitan la identificación de tendencias y relaciones útiles para la ciencia y la teoría.

Enfoque:

- **Machine Learning:** Se centra en el desarrollo de algoritmos y modelos que

²Provost, Foster & Fawcett, T. (2013). "Data Science for Business". United States of America. O'Reilly Media, Inc. Primera Edición.

³Tan, Pang-Ning. Steinbach, Michael. Karpadne, Anuj & Kumar, Vipin. (2005). "Introduction to Data Mining". New York, United States of America. Pearson Education.

⁴Deitel, P & Deitel, H. (2015). "Java: How to Program". New Jersey, United States of America. Pearson Education Inc.

⁵Nisbet, R. Elder, J. & Miner, G. (2009). "Handbook of Statistical Analysis & Data Mining Applications". Canada. Elsevier Inc.

mejoran progresivamente a partir de la experiencia y la incorporación de nuevos datos.

- **Data Mining:** Hace uso de técnicas estadísticas y matemáticas para explorar, identificar y revelar patrones subyacentes en los datos.

Aplicaciones:

- **Machine Learning:** Se aplica en la creación de sistemas de recomendación, diagnósticos médicos, reconocimiento de voz e imágenes, así como en pronósticos financieros y bursátiles, entre otros ámbitos.
- **Data Mining:** Se utiliza para la identificación de tendencias de mercado, segmentación de clientes, detección de fraudes y otras aplicaciones orientadas al análisis exploratorio de grandes volúmenes de datos.

Implementación práctica en Python:

Para ilustrar la diferencia entre Machine Learning y Data Mining, se presenta una implementación práctica en Python. Se emplea la biblioteca Scikit-Learn para construir un modelo predictivo representativo de Machine Learning, y Pandas para realizar un análisis exploratorio típico de Data Mining. Aunque ambos enfoques comparten técnicas de análisis de datos, sus objetivos difieren: Machine Learning se orienta a la creación de modelos capaces de realizar predicciones o tomar decisiones automáticas, mientras que Data Mining se enfoca en descubrir patrones ocultos o relaciones significativas en grandes volúmenes de datos. Se presentan dos ejemplos:⁶

El primer ejemplo, implementado en Python, muestra la exactitud de un modelo de árboles aleatorios (random forest) utilizando datos reales. El análisis detallado de los algoritmos se abordará en secciones posteriores.

In [1]:

```
1 # Ejercicio 1
2 import pandas as pd
3 import warnings
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import accuracy_score
7 from sklearn.preprocessing import LabelEncoder
8
9 warnings.filterwarnings("ignore")
10
11 # Crear un conjunto de datos de ejemplo
12 data = {
```

⁶Venero B. (1999). “Análisis Matemático 2”. Lima, Perú. Ediciones Gemar.

```

13     "Feature1": [1, 2, 3, 4, 5],
14     "Feature2": [10, 20, 30, 40, 50],
15     "Label": ["A", "B", "A", "B", "A"],
16 }
17 df = pd.DataFrame(data)
18
19 # Convertir etiquetas categóricas a numéricas
20 le = LabelEncoder()
21 df["Label"] = le.fit_transform(df["Label"])
22
23 # Dividir el conjunto de datos en características (X) y etiquetas (y)
24 X = df[["Feature1", "Feature2"]]
25 y = df["Label"]
26
27 # Dividir el conjunto de datos en entrenamiento y prueba
28 X_train, X_test, y_train, y_test = train_test_split(
29     X, y, test_size=0.2, random_state=42
30 )
31
32 # Construir un modelo de Machine Learning (Random Forest en este caso)
33 model = RandomForestClassifier()
34 model.fit(X_train, y_train)
35
36 # Hacer predicciones en el conjunto de prueba
37 predictions = model.predict(X_test)
38
39 # Evaluar la precisión del modelo
40 accuracy = accuracy_score(y_test, predictions)
41 print(f"Precisión del modelo: {accuracy}")

```

Out [1]:

```
Precision del modelo: 0.0
```

El modelo obtenido es resultado de un proceso de Machine Learning; sin embargo, desde la perspectiva de Data Mining, la precisión del pronóstico puede ser insuficiente, lo que indica la ausencia de patrones relevantes en los datos. Así, mientras Machine Learning se orienta a la obtención de modelos, Data Mining se encarga de la interpretación y el análisis de los mismos.

El segundo ejemplo ilustra el uso de la estadística descriptiva para analizar la distribución de los datos y la matriz de correlaciones entre variables. Los algoritmos subyacentes serán abordados con mayor profundidad en secciones posteriores.

In [2]:

```

1 # Ejercicio 2
2 import pandas as pd

```



```
3
4 # Crear un conjunto de datos de ejemplo
5 data = {
6     "Feature1": [1, 2, 3, 4, 5],
7     "Feature2": [10, 20, 30, 40, 50],
8     "Label": ["A", "B", "A", "B", "A"],
9 }
10
11 df = pd.DataFrame(data)
12
13 # Explorar y descubrir información en el conjunto de datos
14 # Por ejemplo, calcular estadísticas descriptivas (solo columnas numé-
15     ricas)
16 descriptive_stats = df.describe()
17
18 # Obtener la distribución de las etiquetas
19 label_distribution = df["Label"].value_counts()
20
21 # Encontrar correlaciones entre características numéricas solamente
22 numeric_columns = df.select_dtypes(include=[float, int]).columns
23 correlation_matrix = df[numeric_columns].corr()
24
25 # Mostrar resultados
26 print("Estadísticas descriptivas:")
27 print(descriptive_stats)
28
29 print("\nDistribución de etiquetas:")
30 print(label_distribution)
31
32 print("\nMatriz de correlación (solo columnas numéricas):")
33 print(correlation_matrix)
```

Out [2]:

Estadísticas descriptivas:

| | Feature1 | Feature2 |
|-------|----------|-----------|
| count | 5.000000 | 5.000000 |
| mean | 3.000000 | 30.000000 |
| std | 1.581139 | 15.811388 |
| min | 1.000000 | 10.000000 |
| 25% | 2.000000 | 20.000000 |
| 50% | 3.000000 | 30.000000 |
| 75% | 4.000000 | 40.000000 |
| max | 5.000000 | 50.000000 |

Distribucion de etiquetas:

| Label | |
|-------|---|
| A | 3 |
| B | 2 |

```
Name: count, dtype: int64
```

```
Matriz de correlacion (solo columnas numericas):
```

| | Feature1 | Feature2 |
|----------|----------|----------|
| Feature1 | 1.0 | 1.0 |
| Feature2 | 1.0 | 1.0 |

Habiendo introducido el tema de la Ciencia de Datos, es pertinente profundizar en esta disciplina, de la cual Machine Learning constituye una subárea relevante.

1.3 ¿Qué es la Ciencia de Datos?

La Ciencia de Datos se ha consolidado como un campo interdisciplinario en la era digital, desempeñando un papel central en la transformación de las ciencias y la industria. En un contexto donde se generan volúmenes masivos de datos en todas las áreas profesionales, surge la necesidad de procesar, analizar y transformar estos datos en información relevante que contribuya a la mejora de la calidad profesional y de los productos y servicios ofrecidos por las organizaciones. La Ciencia de Datos se sitúa en la intersección de la estadística, las matemáticas, la informática y el conocimiento del dominio, con el propósito de analizar datos complejos y extraer conocimientos significativos.⁷

En el contexto del Machine Learning, la Ciencia de Datos ha desempeñado y continúa desempeñando un papel significativo al proporcionar los fundamentos para la toma de decisiones informadas y la construcción de modelos predictivos precisos. Estos modelos, presentes en productos cotidianos como el reconocimiento de imágenes y rostros en dispositivos móviles, no serían posibles sin la integración de Machine Learning y las herramientas desarrolladas en el ámbito de la Ciencia de Datos.⁸

1.3.1 Dominio de conocimiento relacionados a la Ciencia de Datos

1. La estadística:

La estadística constituye un pilar esencial en la Ciencia de Datos, al proveer los métodos y herramientas necesarios para describir, analizar e interpretar datos, tanto en muestras pequeñas como en grandes volúmenes almacenados en tablas o bases

⁷McKinney, W. (2013). “Python for Data Analysis”. United States of America. O’Reilly Media Inc. Primera Edición.

⁸Provost, Foster & Fawcett, T. (2013). “Data Science for Business”. United States of America. O’Reilly Media, Inc. Primera Edición.

de datos. Se emplean técnicas de estadística descriptiva, estadística inferencial y pruebas de hipótesis, para sustentar científicamente los análisis realizados. De este modo, la estadística permite a los científicos de datos extraer información relevante y confiable de conjuntos de datos y modelos, dotándolos de un repertorio robusto de herramientas para el aprendizaje y la toma de decisiones.⁹

2. La matemática:

Las matemáticas son muy utilizadas en la Ciencia de Datos, especialmente en la descripción de fenómenos complejos y en la creación de modelos de Machine Learning.¹⁰ Áreas como el álgebra lineal, el cálculo, la aritmética y la teoría de la probabilidad resultan esenciales para comprender y desarrollar algoritmos de Machine Learning por su uso en la manipulación y transformación de datos. La modelización matemática permite representar fenómenos del mundo real de manera precisa y eficiente, facilitando la construcción de modelos robustos y generalizables.¹¹

3. La ciencia de la computación:

La ciencia de la computación es determinante en el procesamiento eficiente de grandes volúmenes de datos, mediante el uso de computadoras para la automatización de cálculos y la programación de soluciones. Lenguajes como Python y R son ampliamente empleados en Ciencia de Datos por su versatilidad y la disponibilidad de bibliotecas especializadas. Asimismo, la gestión de bases de datos y la optimización de procesos computacionales son aspectos clave para garantizar la eficiencia en la manipulación de datos a gran escala.¹²

4. Dominio del negocio o dominio de las ciencias conexas:

La comprensión del contexto o dominio de aplicación es esencial para el éxito de los proyectos de Ciencia de Datos. Los científicos de datos colaboran estrechamente con expertos del área específica para identificar problemas relevantes y diseñar soluciones que generen impacto real. Este conocimiento contextual asegura que los resultados obtenidos sean pertinentes y valiosos en el entorno profesional o

⁹Grus, J. (2015). “Data Science From Scratch”. United States of America. O’Reilly Media Inc.

¹⁰Chávez Reyes, Carmen & León Quintanar. (2007). “La Biblia de las Matemáticas”. México D.F., México. Editorial Letrarte S.A.

¹¹Edwards, C.H. Jr & Penney David. (1994). “Cálculo con Geometría Analítica”. México DF, México. Prentice Hall Hispanoamericana S.A.

¹²Deitel, P & Deitel, H. (2015). “Java: How to Program”. New Jersey, United States of America. Pearson Education Inc.

científico correspondiente.¹³

1.3.2 Aplicaciones de la Ciencia de Datos en ML

La Ciencia de Datos provee el marco conceptual y metodológico indispensable para la implementación exitosa de proyectos de Machine Learning. Su influencia se extiende a lo largo de las distintas etapas del proceso de Machine Learning, como las definidas en la metodología CRISP-DM (Cross-Industry Standard Process for Data Mining), que incluyen:

1. Entendimiento del problema:

En la fase inicial de un proyecto de Machine Learning, los científicos de datos emplean técnicas estadísticas y de análisis exploratorio de datos para comprender la naturaleza y las características del problema en estudio, así como para identificar variables clave, analizar relaciones entre ellas y anticipar posibles desafíos. Esta etapa permite la búsqueda de soluciones viables y adecuadas a los recursos disponibles.

2. Preparación de datos:

La preparación de datos constituye una etapa crítica en el proceso de Machine Learning. La Ciencia de Datos se encarga de la limpieza, gestión de valores atípicos y transformación de los datos para adecuarlos a los requerimientos de los algoritmos de ML. Asimismo, se aplican técnicas de ingeniería de características para mejorar la calidad de los datos y facilitar la construcción de modelos robustos.¹⁴

3. Selección de modelos:

En la etapa de selección de modelos, la Ciencia de Datos orienta la elección de los algoritmos y enfoques más adecuados para el problema específico. Esta decisión se fundamenta en una evaluación rigurosa de los datos y en la comprensión de las fortalezas y limitaciones de los distintos algoritmos disponibles.

4. Evaluación y validación:

La estadística desempeña un papel central en la evaluación y validación de los modelos de Machine Learning. Se emplean métricas específicas para medir el rendimiento

¹³McKinney, W. (2013). “Python for Data Analysis”. United States of America. O’Reilly Media Inc. Primera Edición.

¹⁴Provost, Foster & Fawcett, T. (2013). “Data Science for Business”. United States of America. O’Reilly Media, Inc. Primera Edición.

del modelo y garantizar su capacidad de generalización a nuevos conjuntos de datos. Técnicas como la validación cruzada y otros métodos estadísticos se utilizan para prevenir el sobreajuste y asegurar la robustez de los resultados obtenidos.¹⁵

5. Interpretación de resultados:

La Ciencia de Datos facilita la interpretación de los resultados obtenidos mediante la aplicación de técnicas de visualización de datos y análisis estadístico. Estas herramientas permiten comunicar de manera efectiva las conclusiones a los interesados, asegurando una comprensión clara y fundamentada para la toma de decisiones informadas.¹⁶

La Ciencia de Datos se erige como una disciplina importante en la actualidad, donde la masiva generación de datos requiere un enfoque interdisciplinario para extraer conocimientos significativos. En el contexto de Machine Learning, la Ciencia de Datos tiene un papel relevante en las etapas del proceso, desde el entendimiento del problema hasta la interpretación de los resultados.

1.4 Matemática y estadística en Machine Learning

La matemática y la estadística constituyen los cimientos sobre los cuales se edifica la estructura del Machine Learning (ML). Estas disciplinas resultan indispensables para comprender, diseñar y optimizar los algoritmos que impulsan la Inteligencia Artificial, permitiendo a las máquinas aprender patrones a partir de datos y mejorar su desempeño de manera autónoma.

1.4.1 Rol de la matemática en Machine Learning

Desde los conceptos de álgebra lineal y cálculo hasta la aplicación de la estadística, estas herramientas matemáticas son necesarias para los profesionales de Machine Learning. Entre las principales herramientas matemáticas se destacan:

1. Álgebra lineal:

El álgebra lineal es una rama de las matemáticas que resulta especialmente relevante en Machine Learning, en particular para la representación de datos y la formulación de modelos. Conceptos como matrices y operaciones lineales permiten describir y

¹⁵McKinney, W. (2013). “Python for Data Analysis”. United States of America. O’Reilly Media Inc. Primera Edición.

¹⁶Bowers, N. Jr. Gerber, H. Hickman, J. Jones, D. & Nesbitt, C. (1997). “Actuarial Mathematics”. United States of America. The Society of Actuaries.

datos de manera eficiente. En la implementación de algoritmos, el álgebra lineal se utiliza para definir relaciones entre variables, representar transformaciones en el espacio de características y, posteriormente, contribuir a la construcción de redes neuronales artificiales.

2. Cálculo:

El cálculo, a través de sus ramas de derivadas e integrales, son muy utilizados para la optimización de modelos en Machine Learning. Algoritmos como el Gradiente Descendiente, empleados para minimizar funciones de pérdida y ajustar parámetros del modelo, se fundamentan en conceptos de cálculo. La capacidad de optimizar modelos y encontrar mínimos o máximos locales resulta clave para mejorar el rendimiento y la eficacia de los algoritmos de ML.¹⁷

3. Probabilidades:

Las probabilidades permiten comprender la incertidumbre y la toma de decisiones en ML. Muchos modelos, especialmente en aprendizaje supervisado, se basan en probabilidades para realizar predicciones y clasificaciones. La probabilidad permite cuantificar el grado de certeza o incertidumbre asociado a las predicciones de un modelo, lo que resulta esencial en la evaluación y validación de su desempeño. No obstante, es igualmente relevante considerar la estadística, ya que al construir modelos es necesario contar con criterios científicos que permitan determinar si un algoritmo está óptimamente ajustado y cumple con los requisitos para su implementación.¹⁸

1.4.2 Rol de la estadística en Machine Learning

La estadística también ocupa un lugar destacado en la validación, inferencia y toma de decisiones en el contexto del Machine Learning. Entre los aspectos más relevantes se encuentran:

1. Inferencia estadística:

La inferencia estadística permite realizar afirmaciones sobre una población a partir de un conjunto de datos muestral. En Machine Learning, esto se traduce en la capacidad de generalizar conclusiones sobre un conjunto más amplio de datos a partir de un subconjunto de entrenamiento. Herramientas como las pruebas de

¹⁷Edwards, C.H. Jr & Penney David. (1994). “Cálculo con Geometría Analítica”. México DF, México. Prentice Hall Hispanoamericana S.A.

¹⁸Nisbet, R. Elder, J. & Miner, G. (2009). “Handbook of Statistical Analysis & Data Mining Applications”. Canada. Elsevier Inc.

hipótesis e intervalos de confianza son empleadas para inferir propiedades de la población y validar la robustez de los modelos.

2. Validación de modelos:

La validación estadística de modelos permite evaluar el desempeño de los modelos con datos no observados previamente. Técnicas como la validación cruzada proporcionan estimaciones robustas del rendimiento del modelo al evaluarlo en diferentes subconjuntos de datos, lo que contribuye a prevenir el sobreajuste y asegurar que el modelo generalice adecuadamente a nuevos datos.

3. Métricas de evaluación:

Las métricas empleadas para evaluar modelos, tales como la precisión, la sensibilidad y la especificidad, tienen fundamentos estadísticos. Estas métricas permiten cuantificar la calidad del modelo y ofrecen información relevante sobre su capacidad para realizar predicciones precisas y fiables.

Así, la matemática y la estadística son pilares en la construcción y evaluación de modelos de Machine Learning. Desde la aplicación de álgebra lineal y cálculo para entender la estructura de los datos y la optimización de modelos, hasta la aplicación de las probabilidades y la estadística para cuantificar la incertidumbre, y evaluar la precisión del modelo, estas disciplinas sirven en la creciente disciplina de la Inteligencia Artificial.

1.5 Python para Machine Learning

En el ámbito del Machine Learning, la elección del lenguaje de programación es determinante para el desarrollo eficiente y efectivo de algoritmos. En este libro, se empleará exclusivamente Python para modelar todo tipo de algoritmos y procesos relacionados con las matemáticas.¹⁹ Python se ha consolidado como el lenguaje de programación predominante en Machine Learning, gracias a su facilidad de uso, sintaxis simple, extensas bibliotecas especializadas y una comunidad activa y comprometida. Este apartado explora cómo Python se ha convertido en la herramienta preferida en el campo, destacando las bibliotecas clave y la flexibilidad que ofrece a desarrolladores y científicos de datos. Antes de profundizar, se deben señalar algunas características que hacen de Python una opción destacada en Machine Learning:

1. Sintaxis simple:

¹⁹Chávez Reyes, Carmen & León Quintanar. (2007). “La Biblia de las Matemáticas”. México D.F., México. Editorial Letrarte S.A.

La sintaxis simple de Python es una de las principales razones de su adopción en Machine Learning, especialmente entre profesionales sin formación previa en programación. Esta característica facilita la comprensión de los scripts y promueve la colaboración efectiva en equipos multidisciplinarios. La filosofía de diseño del lenguaje, orientada a la legibilidad y simplicidad, resulta especialmente adecuada para abordar la complejidad inherente de los algoritmos de Machine Learning.²⁰ No obstante, en entornos empresariales puede ser conveniente considerar otros lenguajes más eficientes para implementaciones en dispositivos electrónicos.

2. Amplias bibliotecas:

Python sobresale por su ecosistema de bibliotecas especializadas, principalmente en el desarrollo de aplicaciones de Machine Learning. Entre las más relevantes se encuentran:

- **NumPy:** Ofrece soporte para arreglos multidimensionales y funciones matemáticas, permitiendo la manipulación eficiente de datos y la realización de operaciones matriciales. En Machine Learning, donde el manejo de datos en forma de matrices es habitual, NumPy resulta esencial. Proporciona funciones para operaciones matriciales y algebraicas, simplificando la implementación de algoritmos.²¹
- **Pandas:** Proporciona estructuras de datos avanzadas para el análisis y manipulación de datos. Facilita la limpieza, transformación y preparación de los datos para su uso en algoritmos de Machine Learning. Permite la carga eficiente de conjuntos de datos, su manipulación y preprocesamiento, asegurando que estén listos para su utilización en modelos de ML.
- **Scikit-Learn:** Esta biblioteca comprende una extensa variedad de algoritmos y funciones orientadas a la validación de modelos. Scikit-Learn se erige como una herramienta integral para la construcción, evaluación y selección de modelos en Machine Learning. Ofrece un repertorio amplio de algoritmos, facilitando así la implementación rigurosa de soluciones de Machine Learning en Python.

3. Comunidad activa:

La comunidad de Python se distingue por su dinamismo, permitiendo a desarrolladores y científicos de datos, acceder rápidamente a recursos y soluciones a problemas recurrentes. Esta comunidad impulsa la innovación continua y la evolución de las

²⁰McKinney, W. (2013). “Python for Data Analysis”. United States of America. O’Reilly Media Inc. Primera Edición.

²¹Bressert, E. (2013). “SciPy and NumPy”. United States of America. O’Reilly Media Inc. Primera Edición.

herramientas y prácticas en el ámbito del Machine Learning.

Flexibilidad y versatilidad de Python:

La flexibilidad de Python constituye uno de sus principales activos en el contexto del Machine Learning. Esta característica permite a los profesionales implementar algoritmos de manera eficiente y realizar experimentos de forma ágil. La versatilidad de Python se manifiesta en varias características:

- **Diversidad de algoritmos:** Python es compatible con una amplia gama de algoritmos de Machine Learning, desde los más básicos hasta los más avanzados. Esta diversidad permite seleccionar el enfoque metodológico más adecuado para cada problema específico.
- **Experimentación rápida:** La facilidad de implementación que ofrece Python favorece la experimentación expedita con diferentes algoritmos y metodologías. Los científicos de datos pueden ajustar y validar modelos de manera eficiente, acelerando el ciclo de desarrollo y optimizando la búsqueda de soluciones óptimas.
- **Integración con tecnologías externas:** Python se integra eficazmente con diversas tecnologías y herramientas, facilitando la construcción de soluciones integrales (end-to-end) en Machine Learning. Desde la conexión con bases de datos hasta la implementación de modelos en sistemas productivos, Python proporciona una integración fluida y robusta.

El uso de Python en Machine Learning no es arbitrario; responde a su sintaxis simple, la disponibilidad de bibliotecas especializadas y su flexibilidad inherente. Herramientas como NumPy, Pandas y Scikit-Learn simplifican la construcción y evaluación de modelos, permitiendo abordar desde tareas sencillas hasta desafíos complejos en el ámbito del aprendizaje automático.²²

1.6 Herramientas para Machine Learning

El desarrollo de modelos de Machine Learning exige el uso de herramientas especializadas que simplifican tareas complejas y fomentan la colaboración interdisciplinaria. Si bien Python es el lenguaje predominante, el ecosistema de Machine Learning se enriquece con una variedad de herramientas y entornos diseñados para abordar los desafíos del aprendizaje automático. Se destacan algunas de las herramientas más relevantes, como TensorFlow, PyTorch y Jupyter Notebooks, subrayando su contribución al desarrollo eficiente de modelos y a la creación de documentos que integran código, visualización y análisis reproducible.

²²Deitel, P & Deitel, H. (2015). “Java: How to Program”. New Jersey, United States of America. Pearson Education Inc.

- **TensorFlow:** Esta biblioteca de código abierto, desarrollada por Google, se ha consolidado como un pilar en el ámbito del aprendizaje profundo. Su arquitectura flexible, permite la creación de modelos complejos, desde redes neuronales convolucionales hasta redes recurrentes. TensorFlow destaca por su capacidad para gestionar datos de alto nivel de abstracción y por su escalabilidad, lo que facilita la implementación de modelos avanzados.
- **PyTorch:** Es un módulo donde su estructura dinámica permite modificar y construir modelos durante la ejecución, lo que resulta especialmente útil en investigación y desarrollo de nuevas arquitecturas.
- **Jupyter Notebooks:** Es una herramienta hecha para científicos de datos y profesionales de Machine Learning. Jupyter Notebooks permite combinar código ejecutable, visualizaciones y texto explicativo en un único documento, lo que facilita la exploración de datos, la experimentación con modelos y la presentación de resultados de manera clara y reproducible. Su entorno interactivo promueve la documentación transparente del proceso analítico y la colaboración efectiva entre equipos multidisciplinares.²³

1.6.1 Facilitando la experimentación y colaboración

Estas herramientas no solo son potentes de forma individual, sino que también se integran de manera sinérgica para potenciar la experimentación y la colaboración en proyectos de Machine Learning. Python, TensorFlow, PyTorch y Jupyter Notebooks conforman un ecosistema cohesivo que permite a los equipos trabajar de manera eficiente, compartir resultados y reproducir experimentos. La experimentación rápida y estas herramientas proporcionan la flexibilidad y agilidad necesarias para probar distintos enfoques y algoritmos. La documentación interactiva mediante Jupyter Notebooks facilita la comprensión, la comunicación y la replicabilidad de los resultados obtenidos.

1.6.2 Integración con TensorFlow y PyTorch

La integración de TensorFlow y PyTorch con otras herramientas y bibliotecas del ecosistema Python permite optimizar el flujo de trabajo en Machine Learning y Deep Learning. Por ejemplo, TensorFlow y Keras (una interfaz de alto nivel para TensorFlow) permiten construir y entrenar modelos de Deep Learning de manera eficiente y modular. Estas herramientas proporcionan a los profesionales de datos capacidades avanzadas para construir, entrenar y evaluar modelos, facilitando la experimentación rápida y la colaboración efectiva. El uso de Python como lenguaje principal añade flexibilidad y compatibilidad, permitiendo aprovechar la amplia variedad de librerías y frameworks disponibles en el ecosistema científico y de ingeniería de datos.

²³Grus, J. (2015). “Data Science From Scratch”. United States of America. O’Reilly Media Inc.

1.7 Proyectos de Machine Learning

El desarrollo de un proyecto de Machine Learning requiere la adopción de un enfoque sistemático que abarque desde la comprensión profunda del problema hasta la implementación y evaluación rigurosa del modelo. Se describen las etapas en el ciclo de vida de un proyecto de Machine Learning, subrayando la importancia de definir con precisión el problema, preparar los datos de manera adecuada, seleccionar el algoritmo idóneo y realizar una evaluación exhaustiva del rendimiento. Asimismo, se destaca la relevancia de la iteración y la optimización como componentes clave para alcanzar resultados robustos y reproducibles.

1. Definición del problema:

La primera etapa consiste en comprender y definir con claridad el problema a abordar. Esto implica establecer los objetivos del proyecto, identificar las variables relevantes y analizar el contexto en el que se aplicará la solución. Una definición precisa del problema orienta adecuadamente en el resto del proceso y maximiza el impacto de la solución propuesta.

2. Recopilación y preparación de datos:

La calidad de los datos es determinante para el éxito de cualquier proyecto de Machine Learning. Esta etapa abarca la recopilación de datos relevantes, su limpieza, transformación y preparación para el entrenamiento del modelo. Un preprocesamiento adecuado garantiza que los datos sean representativos y aptos para el análisis posterior.

3. Selección del algoritmo:

La selección del algoritmo más adecuado depende del tipo de tarea de Machine Learning (supervisada, no supervisada o por refuerzo), la naturaleza de los datos y los objetivos específicos del proyecto. Existen múltiples opciones, como algoritmos de regresión, clasificación, agrupamiento o redes neuronales, y la elección debe basarse en un análisis cuidadoso de las características del problema y los datos disponibles.

4. División del conjunto de datos:

Esta etapa implica dividir el conjunto de datos en subconjuntos de entrenamiento y prueba (o validación). El conjunto de entrenamiento se utiliza para ajustar el modelo, mientras que el conjunto de prueba permite evaluar su rendimiento en datos no vistos previamente. Esta práctica es esencial para evitar el sobreajuste, un

fenómeno en el que el modelo se ajusta excesivamente a los datos de entrenamiento y pierde capacidad de generalización.

5. Entrenamiento del modelo:

En esta fase, el modelo se entrena utilizando el conjunto de datos de entrenamiento. El objetivo es ajustar los parámetros del modelo para minimizar la diferencia entre las predicciones y los valores reales, optimizando así su capacidad predictiva.

6. Evaluación del rendimiento:

La evaluación del rendimiento del modelo permite determinar su eficacia a un problema dado. Se emplean métricas específicas según la naturaleza de la tarea: precisión, recall y F1-score para clasificación; error cuadrático medio y coeficiente de determinación para regresión, entre otras. La validación cruzada y otras técnicas permiten obtener una evaluación más robusta y confiable del desempeño del modelo.

7. Iteración y optimización:

La iteración y la optimización son procesos continuos en el ciclo de vida de un proyecto de Machine Learning. Tras evaluar el rendimiento inicial, se pueden realizar ajustes en la selección de características, la arquitectura del modelo o los hiperparámetros para mejorar su desempeño. Este proceso iterativo se repite hasta alcanzar los resultados deseados, promoviendo la mejora continua del modelo.

1.7.1 Antes de empezar con Machine Learning

- **Importancia de la división del conjunto de datos:** Antes de iniciar cualquier proyecto de Machine Learning, se debe realizar la división correcta del conjunto de datos en subconjuntos de entrenamiento y prueba. Esta práctica es esencial para evaluar la capacidad de generalización del modelo y evitar el sobreajuste, fenómeno en el que el modelo memoriza los datos de entrenamiento pero falla al enfrentarse a nuevas instancias. Python, mediante bibliotecas como Scikit-Learn, ofrece herramientas eficientes para realizar esta división de manera sistemática.
- **Iteración y optimización en la práctica:** La iteración y optimización constituyen procesos continuos en el desarrollo de proyectos de Machine Learning. Tras la evaluación inicial del modelo, los profesionales pueden realizar ajustes y mejoras para incrementar la precisión y robustez de los resultados. Herramientas como Grid Search y Random Search facilitan la búsqueda sistemática de hiperparámetros óptimos.
- **Integración de herramientas:** La integración de TensorFlow, PyTorch y

Jupyter Notebooks en el flujo de trabajo de desarrollo de modelos en Python aporta flexibilidad y eficiencia adicionales. Estas herramientas no solo facilitan la implementación de algoritmos avanzados, sino que también promueven la documentación interactiva y la colaboración entre profesionales. El uso de Jupyter Notebooks para documentar y compartir código y resultados permite a los equipos comprender y reproducir el flujo de trabajo de extremo a extremo. Además, TensorFlow y PyTorch se integran fácilmente con otras herramientas del ecosistema, maximizando las capacidades de estas bibliotecas en proyectos de Machine Learning.

En síntesis, la adopción de un enfoque sistemático en proyectos de Machine Learning—desde la definición precisa del problema hasta la iteración y optimización—es clave para el éxito. La adecuada división del conjunto de datos, el entrenamiento riguroso del modelo y la evaluación exhaustiva del rendimiento constituyen pasos esenciales. Python, junto con bibliotecas especializadas como Scikit-Learn, TensorFlow y PyTorch, proporciona las herramientas necesarias para abordar estos desafíos de manera eficiente y reproducible.

1.8 Resumen

En resumen, se tienen los siguientes puntos relevantes del capítulo:

1.1. ¿Qué es Machine Learning?

Definición: Machine Learning, es una rama de la IA, centrada en el desarrollo de algoritmos y modelos que permiten a las máquinas y a las computadoras, aprender de los datos y realizar tareas, y así lograr objetivos sin programación explícita.

Áreas de aprendizaje:

- Aprendizaje Supervisado: Entrenamiento con datos etiquetados.
- Aprendizaje no Supervisado: Descubrimiento de patrones en datos sin etiquetar.
- Aprendizaje por refuerzo: Son todos aquellos modelos que aprenden a través de recompensas y castigos en base a sus acciones.

Aplicaciones:

Reconocimiento de voz, recomendaciones en línea, diagnósticos médicos, conducción autónoma, reconocimiento de imágenes y otros.

Tecnologías relacionadas:

Redes neuronales artificiales, Robótica, Biomimética, Generative IA.

1.2. Machine Learning vs Data Mining

Diferencias entre Machine Learning y Data Mining:

- Machine Learning: Se enfoca en construir modelos predictivos basados en patrones de datos.
- Data Mining: Se enfoca en explorar grandes conjuntos de datos para encontrar patrones y tendencias.²⁴

Objetivos y enfoques:

²⁴Tan, Pang-Ning. Steinbach, Michael. Karpatne, Anuj & Kumar, Vipin. (2005). “Introduction to Data Mining”. New York, United States of America. Pearson Education.

- Machine Learning: Mejora de la precisión de las predicciones a través de la experiencia y datos adicionales.
- Data Mining: Descubrimiento de información valiosa y conocimiento en los datos.

1.3. ¿Qué es la Ciencia de Datos?

Definición: Campo interdisciplinario que utiliza estadísticas, matemáticas, informática y conocimiento del negocio para analizar datos complejos y extraer conocimientos valiosos.

Componentes:

- Estadística: Descripción, análisis e interpretación de datos.
- Matemáticas: Modelización funcional de fenómenos complejos.
- Ciencias de la Computación: Para el manejo y procesamiento de grandes volúmenes de datos (big data) mediante el uso de la programación.
- Conocimiento del Negocio: Comprensión del contexto empresarial o científico.

1.4. El uso de la matemática y la estadística en Machine Learning

Importancia: En el diseño y optimización de algoritmos.

Herramientas:

- Álgebra Lineal: Representación y manipulación de datos.
- Cálculo: Optimización de modelos.
- Probabilidades y estadística: Comprensión de la incertidumbre y evaluación del rendimiento de los modelos.

1.5. El uso de Python para Machine Learning

Ventajas de Python:

- Sintaxis Simple: Facilita la comprensión y colaboración.
- Bibliotecas especializadas: NumPy, Pandas, Scikit-Learn.
- Comunidad Activa: Recursos y soporte continuo.

1.6. Herramientas para Machine Learning

Herramientas Clave:

- TensorFlow y PyTorch: Bibliotecas para Deep Learning.
- Jupyter Notebooks: Entorno interactivo para documentar y compartir código y resultados.

1.7. El proyecto de Machine Learning**Proceso sistemático:**

- Definición del Problema: Lograr entender y poder definir el problema.
- Recopilación y preparación de datos: Incluye la limpieza y preparación de datos que han de ser utilizados.
- Selección del algoritmo: Elección del algoritmo más adecuado.
- División del conjunto de datos: Entrenamiento y prueba para evitar el sobreajuste.
- Entrenamiento del modelo: Ajuste de parámetros.
- Evaluación del rendimiento: Uso de métricas específicas para evaluar la eficacia del modelo.
- Iteración y optimización: Mejora continua del modelo.

El capítulo aborda aspectos relevantes de Machine Learning, desde su definición y áreas de aprendizaje hasta las herramientas y procesos necesarios para desarrollar proyectos exitosos en este campo.²⁵

²⁵Tan, Pang-Ning. Steinbach, Michael. Karpatne, Anuj & Kumar, Vipin. (2005). “Introduction to Data Mining”. New York, United States of America. Pearson Education.

1.9 Problemas y sus Soluciones

Se tienen los siguientes problemas y soluciones:

Pregunta 1: ¿Qué es Machine Learning?

- A) Un tipo de software antivirus.
- B) Una rama de la Inteligencia Artificial.
- C) Un lenguaje de programación.
- D) Un sistema operativo.
- E) Un hardware específico para computadoras.
- F) Una herramienta de análisis estadístico.
- G) Un modelo económico.

Pregunta 2: ¿Cuál es un ejemplo de aprendizaje supervisado?

- A) Agrupamiento de datos.
- B) Juegos de estrategia.
- C) Identificación de imágenes de caras y rostros.
- D) Descubrimiento de patrones en datos.
- E) Navegación autónoma.
- F) Diseños inspirados en la naturaleza.
- G) Reconocimiento de patrones de comportamiento.

Pregunta 3: ¿Qué es el aprendizaje no supervisado?

- A) Algoritmos entrenados con datos etiquetados.
- B) Aprendizaje de modelos a través de recompensas y castigos.
- C) Algoritmos que encuentran patrones en datos no etiquetados.
- D) Sistemas que requieren programación explícita.
- E) Modelos basados en la interacción humana.
- F) Algoritmos que imitan el comportamiento humano.
- G) Sistemas que solo se utilizan en robótica.

Pregunta 4: ¿Qué describe mejor el aprendizaje por refuerzo?

- A) Modelos que aprenden sin datos.
- B) Algoritmos basados en etiquetas predefinidas.
- C) Sistemas que mejoran su rendimiento a través de recompensas y castigos.
- D) Técnicas centradas exclusivamente en la estadística.
- E) Modelos que no requieren interacción con el entorno.
- F) Sistemas basados en la previsión del mercado.

G) Algoritmos utilizados únicamente en finanzas.

Pregunta 5: ¿Cuál es una aplicación de Machine Learning en la industria automotriz?

- A) Desarrollo de juegos.
- B) Diagnósticos médicos.
- C) Análisis de mercado.
- D) Conducción autónoma.
- E) Traducción de idiomas.
- F) Educación en línea.
- G) Diseño gráfico.

Pregunta 6: ¿En qué se diferencia Machine Learning de Data Mining?

- A) Machine Learning se enfoca en programar explícitamente las computadoras.
- B) DM utiliza algoritmos para aprender de los datos.
- C) ML se centra en construir modelos predictivos.
- D) DM se usa exclusivamente en medicina.
- E) ML no utiliza estadísticas.
- F) DM no permite el análisis en big data.
- G) ML es una subcategoría de Data Mining.

Pregunta 7: ¿Qué lenguaje de programación es ampliamente utilizado en Machine Learning?

- A) C++.
- B) Java.
- C) Python.
- D) Ruby.
- E) PHP.
- F) Swift.
- G) JavaScript.

Pregunta 8: ¿Cuál es un uso común de las redes neuronales artificiales en Machine Learning?

- A) Procesamiento de textos.
- B) Automatización de tareas domésticas.
- C) Desarrollo de sistemas operativos.
- D) Creación de música.
- E) Análisis de datos financieros.

- F) Aprendizaje profundo.
- G) Seguridad informática.

Pregunta 9: ¿Qué representa el término 'Biomimética' en el contexto de Machine Learning?

- A) Aprendizaje de idiomas.
- B) Modelos inspirados en la naturaleza.
- C) Robótica industrial.
- D) Simulaciones climáticas.
- E) Desarrollo de software educativo.
- F) Diseño de videojuegos.
- G) Marketing digital.

Pregunta 10: ¿Qué herramienta de Machine Learning se utiliza para la construcción y entrenamiento de modelos complejos?

- A) Excel.
- B) TensorFlow.
- C) Adobe Photoshop.
- D) AutoCAD.
- E) Microsoft Word.
- F) WordPress.
- G) Oracle Database.

1.10 Problemas

Se tienen los siguientes problemas:

Pregunta 1: ¿Qué es una característica clave de Python, en el contexto de Machine Learning?

- A) Sus capacidades de diseño gráfico.
- B) Programación de bajo nivel.
- C) Su sintaxis simple y legible.
- D) Su principal uso en juegos.
- E) Su incompatibilidad con la mayoría de las bibliotecas de datos.
- F) Su lento tiempo de ejecución.
- G) Su uso limitado en la comunidad científica.

Pregunta 2: ¿Qué librería de Python proporciona estructuras flexibles para análisis de datos?

- A) TensorFlow.
- B) PyTorch.
- C) Matplotlib.
- D) Pandas.
- E) Django.
- F) Flask.
- G) SQLAlchemy.

Pregunta 3: ¿Cuál es un ejemplo de aplicación del Machine Learning en el campo de la salud?

- A) Desarrollo de videojuegos.
- B) Producción de música.
- C) Diagnósticos médicos.
- D) Diseño de interiores.
- E) Creación de contenido en redes sociales.
- F) Construcción de edificios.
- G) Agricultura.

Pregunta 4: ¿Qué describe mejor el término "*Procesamiento del Lenguaje Natural (NLP)*"

- A) Algoritmos para videojuegos.
- B) Modelos para el reconocimiento de voz.

- C) Técnicas para la conducción autónoma.
- D) Sistemas para el diseño gráfico.
- E) Herramientas para la edición de video.
- F) Métodos para la agricultura.
- G) Protocolos para redes sociales.

Pregunta 5: ¿Qué herramienta se utiliza comúnmente para la experimentación y documentación interactiva en Machine Learning?

- A) Microsoft Excel.
- B) Adobe Illustrator.
- C) Jupyter Notebooks.
- D) Unity.
- E) AutoCAD.
- F) WordPress.
- G) Tableau.

Pregunta 6: ¿Cuál es un ejemplo de Aprendizaje Profundo (Deep Learning)?

- A) Programación básica.
- B) Redes Neuronales Artificiales.
- C) Cálculo simple.
- D) Diseño web.
- E) Gestión de bases de datos.
- F) Marketing digital.
- G) Análisis literario.

Pregunta 7: ¿Qué factor es importante para evitar en el entrenamiento de modelos de Machine Learning?

- A) Aprendizaje rápido.
- B) Sobreajuste.
- C) Uso de Python.
- D) Colaboración.
- E) Documentación.
- F) Integración con otras tecnologías.
- G) Experimentación.

Pregunta 8: ¿Qué biblioteca de Python se especializa en operaciones matriciales y algebraicas?

- A) Flask.

- B) Django.
- C) PyTorch.
- D) Matplotlib.
- E) Pandas.
- F) TensorFlow.
- G) NumPy.

Pregunta 9: ¿Cuál es un enfoque común en Data Mining?

- A) Construcción de redes sociales.
- B) Desarrollo de sistemas operativos.
- C) Minería de asociaciones para identificar patrones.
- D) Programación de videojuegos.
- E) Creación de contenido digital.
- F) Diseño arquitectónico.
- G) Educación en línea.

Pregunta 10: ¿Qué representa el término “Aprendizaje Supervisado”?

- A) Modelos que aprenden sin datos.
- B) Algoritmos que no requieren entrenamiento.
- C) Sistemas que funcionan exclusivamente con datos etiquetados.
- D) Que se aplican solo en robótica.
- E) Modelos que no utilizan estadísticas.
- F) Algoritmos diseñados para juegos.
- G) Sistemas que solo se aplican en finanzas.

Pregunta 11: ¿En qué se diferencia principalmente el Aprendizaje no Supervisado del Aprendizaje Supervisado?

- A) Uso de datos etiquetados.
- B) Aplicación en la salud.
- C) Uso de Python.
- D) Enfoque en robótica.
- E) Integración con redes neuronales.
- F) Uso de tecnologías externas.
- G) Enfoque en juegos.

Pregunta 12: ¿Qué tarea es un ejemplo de Aprendizaje por Refuerzo en Machine Learning?

- A) Desarrollo de bases de datos.

- B) Creación de contenido web.
- C) Toma de decisiones en un entorno de juego.
- D) Diseño de software antivirus.
- E) Construcción de edificios.
- F) Producción de películas.
- G) Marketing en redes sociales.

Pregunta 13: ¿Qué lenguaje de programación es frecuentemente utilizado en Data Mining?

- A) C++.
- B) JavaScript.
- C) Ruby.
- D) Python.
- E) PHP.
- F) Swift.
- G) Kotlin.

Pregunta 14: ¿Qué área de Machine Learning se centra en análisis de datos visuales?

- A) Procesamiento del Lenguaje Natural.
- B) Biomimética.
- C) Visión por Computadoras.
- D) Análisis literario.
- E) Diseño de videojuegos.
- F) Marketing digital.
- G) Desarrollo web.

Pregunta 15: ¿Qué técnica de Machine Learning es utilizada para el reconocimiento de voz?

- A) Aprendizaje supervisado.
- B) Aprendizaje no supervisado.
- C) Aprendizaje por refuerzo.
- D) Biomimética.
- E) Análisis de redes sociales.
- F) Desarrollo de videojuegos.
- G) Diseño gráfico.

Pregunta 16: ¿Qué herramienta de Python es útil para construir y evaluar

modelos?

- A) Django.
- B) Flask.
- C) Scikit-Learn.
- D) TensorFlow.
- E) PyTorch.
- F) NumPy.
- G) Pandas.

Pregunta 17: ¿Qué método se utiliza para evaluar el rendimiento de un modelo de Machine Learning en diferentes subconjuntos de datos?

- A) Sobreajuste.
- B) Validación cruzada.
- C) Aprendizaje supervisado.
- D) Aprendizaje no supervisado.
- E) Aprendizaje por refuerzo.
- F) Minería de asociaciones.
- G) Procesamiento del Lenguaje Natural.

Pregunta 18: ¿Cuál es una aplicación común del Machine Learning en el comercio electrónico?

- A) Creación de contenido en redes sociales.
- B) Diseño de software antivirus.
- C) Sistemas de recomendación.
- D) Desarrollo de sistemas operativos.
- E) Producción de películas.
- F) Educación en línea.
- G) Construcción de edificios.

Pregunta 19: ¿Qué describe mejor la Ciencia de Datos en relación con Machine Learning?

- A) Una herramienta para el diseño gráfico.
- B) Una disciplina centrada en la creación de juegos.
- C) Un campo que proporciona la base para tomar decisiones informadas.
- D) Un método para el desarrollo de software antivirus.
- E) Una tecnología centrada en la producción de música.
- F) Una técnica utilizada solo en marketing.
- G) Un protocolo para redes sociales.

Pregunta 20: ¿Qué se usa en la Ciencia de Datos para analizar e interpretar grandes volúmenes de datos?

- A) Juegos.
- B) Diseño gráfico.
- C) Estadística.
- D) Producción de música.
- E) Desarrollo de software antivirus.
- F) Marketing digital.
- G) Educación en línea.

Índice alfabético

- Acceder a elementos, 47, 50
- Adaptación continua, 1
- Ajuste de modelos, 434, 805
- Álgebra lineal, 12–14, 150, 162, 165, 172, 174, 175, 182, 443, 474, 509, 511, 512, 515, 527, 563
- Álgebra simbólica, 161
- Algoritmos
 - de clasificación, 649, 771, 819
 - de maximización de la esperanza, 862
 - de regresión, 155
 - no supervisados, 821, 869
 - supervisados, 728
- Amplia variedad de modelos, 680
- Amplias bibliotecas, 15
- Anidación de condiciones, 46
- Análisis
 - bidimensional, 631, 633
 - de componentes principales, 461, 518, 827, 878
 - de datos científicos, 184
 - de frecuencia, 482
 - del componente principal por kernel, 862, 878, 881
 - exploratorio, 6, 11, 194, 195, 245, 395
- Aprendizaje, 96, 151, 553, 591, 662, 669, 677, 686, 687, 692, 695
 - no supervisado, 2, 913, 914, 918
 - por refuerzo, 2
 - supervisado, 2, 5, 13, 786, 792
- Árboles de decisión, 659, 695, 697, 699, 727, 795, 799, 808, 809, 811, 812, 816
- Área bajo la curva, 426
- Argumentos por palabra clave, 88, 89, 93, 94
- Argumentos posicionales, 86, 87, 93
- Aritmética, 10
- Asignación latente de dirichlet, 862, 907, 909
- Asimetría, 577, 615, 618, 622–624
- Autoencoders ns, 862
- Automatización de tareas, 4, 112, 169
- Backends, 686, 687, 690, 692
- Bagging, 808–811
- Barras, 105, 203, 212, 235, 245, 246, 248, 249, 251, 252, 254–256, 258–260, 262, 263, 279–281, 287, 388
- Biblioteca sqlalchemy, 64
- Bokeh, 150, 205, 207, 208, 210, 212
- Booleanos, 36–38
- Boosting, 649, 650, 695, 697, 699, 706, 727, 799, 809, 811–815
- Bosques aleatorios, 659, 727, 799, 800, 802, 808, 811
- Broadcasting, 154
- Bucles, 33, 110, 112, 115, 119, 120, 122
- Burbujas, 206, 235, 390, 391
- Búsqueda de tópicos, 907
- Cadenas, 35, 37, 48, 99, 101, 189
- Cadenas de markov, 861, 862, 891–894
- Cajas, 235, 291–298, 327
- Cambio, 310, 534, 536, 546, 730, 768, 803, 863, 870
- Características, 3, 11, 13, 19, 48, 54, 67, 84, 85, 269, 327, 345, 350, 358, 383, 423, 443, 449, 457, 465, 526, 532, 597, 609, 612, 624, 673, 677, 686, 695, 697, 699, 702, 706, 710, 727, 739, 743, 745, 749, 753, 757, 761, 764–768, 786, 789–792, 795, 796, 800, 802, 815, 821, 826, 827, 830, 831, 833, 861, 878, 880, 886, 889, 896, 905, 910, 911, 914, 918
- Carga de datos, 124, 689
- Ciencia de datos, 1, 4, 9–12, 212, 307, 365, 390, 695, 699
- Claridad, 88, 94, 395
- Clasificación
 - lineal, 767–769, 771, 778
 - logística, 771–774
 - por k-vecinos más cercanos, 786, 788
 - por la aproximación de laplace, 782, 785, 786
 - por support vector machines, 792
 - probit, 775, 778, 779
 - tobit, 778, 781
- Claves en diccionarios, 71, 74
- Clustering, 468, 474, 649, 650, 652
- Coefficiente de correlación, 634, 899
- Coefficiente de determinación, 19, 736, 739
- Coefficientes, 598, 710, 727, 730, 731, 733–736, 738–740, 743–745, 753, 762, 764–766, 772
- Compatibilidad, 17, 690

- Compatibilidad con versiones anteriores, 94
- Componente principal probabilístico, 862
- Componentes, 215, 365, 452, 560, 666, 761, 831, 866, 869, 872, 873, 876, 878–880, 882, 885, 888, 889, 907
- Comunidad activa, 14, 15, 662, 681
- Condicionales, 33, 110, 111, 120, 122, 586, 594, 595, 757, 769
- Condiciones, 41, 42, 46, 110, 120, 520, 526, 558, 589, 608, 750, 754, 784, 795
- Configuración de modelos, 96
- Configuración flexible, 95, 97
- Conjunto de datos, 2, 12, 13, 18–20, 54, 84, 85, 104, 115, 119, 126, 127, 160, 249, 251, 255, 256, 262, 265, 268–270, 275, 278, 281, 283, 284, 286–288, 290, 292, 297–299, 307, 328, 421, 598, 609, 612, 613, 619, 621, 654, 703, 727, 735, 757, 765, 792, 796, 800, 802, 803, 805, 809, 812, 864, 869, 873, 885, 897
- Conjuntos, 5, 10, 15, 54, 84, 112, 119, 120, 124, 126, 151, 153, 202, 204, 235, 240, 245, 265, 292, 308, 334, 383, 389, 390, 392, 395, 399, 400, 419, 420, 443, 481, 532, 598, 605, 612, 624, 631, 673, 680, 686, 697, 699, 705, 756, 765, 768, 786, 792, 799, 803, 805, 806, 808, 809, 815, 818, 821, 822, 824, 830, 831, 871, 875, 880–882, 889, 910, 913
- Contorno para mapas, 334
- Creación de aplicaciones web interactivas, 204
- Curtosis, 615, 622–624
- Curvas de nivel, 235, 312, 317, 320–323
- Cálculo, 10, 12–14, 83, 419, 422, 424–426, 525, 527, 529, 532, 535, 536, 539, 553, 563, 806, 864, 883
- Cálculo diferencial, 163
- Cálculos probabilísticos, 577, 598
- Dashboards, 150, 204, 205
- Data mining, 1, 4–7, 11, 789, 797, 872
- Dataframe, 54, 103, 150, 189, 190, 196–199
- Datos categóricos, 49, 702, 704, 706
- Datos faltantes, 192, 197
- Declaración de funciones, 72
- Definición del problema, 18
- Densidad, 203, 235, 290, 296, 307, 308, 319, 327, 328, 334, 365, 395, 397, 398, 590, 754, 757, 869
- Derivadas, 13, 150, 527, 529, 532, 536, 538, 544, 761, 762
- Desarrollo de software, 80, 168, 171
- Desarrollo web, 205, 210, 220
- Descomposición de valores singulares, 862, 884
- Desigualdad, 42
- Desigualdad de markov, 592
- Despliegue, 217, 677
- Desviación estándar, 85, 588, 597, 598, 614, 618, 620–622, 873
- Devolución de múltiples valores, 74
- Diagrama
 - de marimekko, 388, 389
 - de venn, 400, 401
 - de violín, 328
 - de áreas, 354–356, 358, 359, 361
- Diagrama de venn, 402
- Diccionarios, 33, 47–49, 51, 54, 70, 72
- Dispersión, 202, 205, 235–241, 291, 612–614, 616, 620–622, 764, 769
- Distribuciones continuas, 577
- Distribuciones de probabilidades, 280–282
- División, 492, 797, 800
- División del conjunto de datos, 18, 19, 85
- Documentación, 75, 202
- Ecuaciones diferenciales, 161, 175, 184, 186, 509, 553, 555, 556, 558, 559, 561, 563, 892
- Ecuaciones lineales, 510, 514
- Eficiencia, 10, 20, 120, 151–154, 160, 181, 650, 695, 699, 703, 768, 772, 790, 822, 827
- Elif, 110
- Enfoque, 2, 3, 5, 12, 18, 20, 122, 212, 218, 607, 734, 750, 753, 754, 761, 764, 766, 769, 785, 830, 833, 861, 885, 890, 894, 896, 899, 900, 905, 906, 911, 912
- Enfoque estadístico, 650, 712
- Entendimiento del problema, 11
- Entero, 36, 39, 447
- Entrenamiento, 84, 85, 96, 97, 125, 126, 426, 440, 532, 649, 650, 654, 662, 666, 673, 677, 680, 686, 687, 690, 705, 728, 757, 769, 786, 787, 791, 792, 796, 803, 805, 809, 834, 841, 880, 897, 914, 915
- Entrenamiento del modelo, 18, 19
- Escalabilidad, 649, 662, 899
- Espacios de características, 474, 475

- Espacios vectoriales, 449, 454, 459, 460
- Especificar excepciones, 131
- Estadística descriptiva, 7, 577, 612, 613, 616, 618
- Estadística inferencial, 10, 624, 629
- Estadísticas, 5, 6, 11, 12, 14, 150, 172, 174, 278, 525, 584, 586, 598, 607, 620, 622, 624, 629, 636, 750, 775, 778, 781, 797
- Estructuras de control de flujo, 119
- Estructuras de datos, 15, 150, 188, 189, 204, 432, 786
- Estructuras de datos en python, 54
- Estética, 205, 207, 210, 212, 301
- Evaluación, 328, 588
- Evaluación de modelos, 14, 16, 422, 424, 650, 687, 690, 805
- Evaluación del rendimiento, 19
- Evaluación y validación, 11, 13
- Excepciones, 122–124, 126, 128, 129
- Experimentación rápida, 17
- Exploración de datos, 34, 205, 636
- Explorando, 808, 811, 814, 817, 821
- FileNotFoundError, 128
- Filtrado de datos, 190
- Filtro colaborativo, 861, 862, 899, 903
- Flexibilidad, 14, 16, 17, 20, 34, 86, 93, 96, 120, 190, 220, 649, 657, 677, 686, 690, 743, 745, 757, 793, 815, 830, 833, 879, 899
- Folium, 149, 218, 220
- For, 735, 744, 755, 762, 866, 872, 914
- Frecuencia y amplitud, 481
- Funcionalidades, 107, 150, 171, 172, 202, 662, 695
- Funciones, 149, 150, 161, 174, 175, 188, 198, 201, 203–205, 326, 355, 357, 749, 782, 827, 837
- Funciones anónimas, 45
- Funciones de pérdida, 13, 422, 426, 427, 440, 536, 591, 695, 897
- Funciones trigonométricas inversas, 496
- Geometría, 161, 164, 165, 419, 455, 456, 458, 460, 461, 471, 474, 475, 478
- Gradientes, 359, 509, 532
- Gráficos de dispersión, 204, 235, 633
- Gráficos en 3d, 204
- Herramientas, 9, 10, 12, 16, 17, 19, 20, 128, 149–151, 161, 172, 174, 181, 184, 188, 201, 202, 205, 207, 210, 235, 245, 265, 279, 307, 308, 334, 365, 383, 388, 392, 393, 395, 399, 445, 455, 478, 512, 553, 563, 598, 607, 616, 620, 624, 629, 636, 649, 650, 659, 677, 680, 699, 706, 710, 735, 805, 811, 814, 821, 826, 833, 869, 875, 896, 898, 910, 913
- Hidden markov, 862
- Histogramas, 201–203, 235, 278, 279, 616, 622, 624
- Homoscedasticidad, 731
- Igualdad, 42, 622, 818, 819, 830
- Independencia, 521, 731, 790, 791
- Indexación, 47, 50–52, 54, 69
- Inferencia estadística, 13, 586
- Integración, 531, 532, 540, 749, 750, 790
- Integración numérica, 175, 180
- Inteligencia artificial, 1, 12, 14, 670
- Interactividad, 212
- Interfaz, 17, 203, 204, 218, 649, 650, 680, 686, 687, 690, 692, 695, 699
- Interpolación, 172, 174, 178, 398, 469
- Interpretabilidad, 734, 767, 771, 818
- Interpretación de resultados, 12, 588
- Interpretación detallada, 650, 710
- Investigación, 3, 657, 735
- Investigación científica, 168, 188
- Iteración y optimización, 18–20
- Jupyter notebooks, 16, 17, 20, 202, 208, 210
- Keras, 17, 649, 686, 687, 690, 692, 695
- La estadística, 9–14, 612, 624, 629, 871
- Lambda, 45, 597, 735, 736, 740, 744, 761, 762, 880
- Legibilidad, 15, 75, 81, 83, 88, 94, 107, 153, 392
- Ley de los grandes números, 588
- Lightgbm, 649, 650, 699, 702, 705, 706
- Linealidad, 521, 731, 827
- Lista anidada, 40
- Lista de cadenas, 40
- Listas, 33, 47–49, 51, 52, 54, 67–70, 99, 152
- Logueo de errores, 129
- Líneas, 161, 202, 203, 205, 220, 235, 265–267, 270, 273–275, 277, 284, 316, 334, 336, 339, 346, 348, 351, 357, 361, 386, 471

- Machine Learning, 1, 3–7, 9–20, 33, 34, 41, 45, 47, 48, 54, 67, 72–75, 84, 96, 102, 106, 110, 112, 115, 119, 120, 124–126, 128, 131, 149–151, 153–155, 157, 160, 188, 193–195, 202–204, 212, 218, 220, 221, 235, 307, 354, 419, 420, 422, 431, 432, 434, 436, 442, 449, 455, 460, 463, 468, 471, 474, 478, 479, 481, 485, 493, 495, 497, 509, 511, 514, 525, 527, 531, 532, 536, 541, 546, 553, 555, 558, 563, 577, 586, 595, 605, 607, 608, 629, 649, 650, 654, 657, 659, 662, 669, 677, 686, 692, 695, 699, 706, 727, 728, 734, 743, 745, 749, 753, 756, 761, 764, 767, 771, 778, 785, 792, 795, 799, 802, 805, 808, 809, 811, 812, 814, 817, 821, 830, 833, 836, 861, 865, 871, 882, 899, 907
- Manejo de excepciones, 33, 122, 128, 131
- Manejo eficiente, 124
- Manipulaciones simbólicas, 446
- Manipulación de datos, 10, 41, 54, 149, 151, 160, 161, 181, 422, 423, 455, 478, 563, 584, 711
- Manipulación eficiente de datos, 201, 686
- Manipulación simbólica, 161, 165, 445
- Manipulación y análisis de datos, 150
- Mantenimiento, 80, 82, 107
- Mapas de calor, 205, 235, 307
- Matplotlib, 104, 105, 150, 201–204, 235, 236, 245, 265, 279, 291, 299, 307, 334, 354, 365, 383, 388, 390, 392, 393, 395, 399, 463, 497, 577, 629, 631, 635, 650
- Matrices, 12, 15, 48, 49, 151, 161, 440, 443, 452, 455, 474, 509, 511–513, 524, 525, 828, 830, 831, 841, 861, 882, 883, 886, 888, 889, 891
- Matrix factorization, 885
- Matriz factorizada no negativa, 862, 888, 891
- Matriz factorizada probabilística, 862
- Mayor o igual que, 43
- Mayor que, 43, 431, 592
- Media, 5, 9, 11, 159, 423, 703, 739, 753, 754, 757, 828, 831, 863, 869, 873, 876, 886, 900
- Mediana, 111, 291, 328, 612, 618–620
- Medidas de dispersión, 577
- Medidas de posición central, 618, 620
- Mejora de la legibilidad, 82
- Memoryerror, 127
- Menor o igual que, 44
- Menor que, 43, 627, 889
- Mensajes explicativos, 130, 131
- Modelado de patrones, 482
- Modelado simple, 515
- Modelos estadísticos, 649, 650, 706, 710, 711
- Modelos secuenciales, 891, 894
- Modularidad, 105, 107, 109
- Modularización del código, 81
- Multipliación, 511, 513
- Mutables, 48
- Métodos, 174, 474, 727, 729, 730, 735, 736, 743, 750, 753, 757, 761, 767, 778, 793, 805, 808, 864, 897
- Métodos en cadenas, 101
- Métodos en listas, 99, 100
- Métodos personalizados, 105, 107
- Métricas de evaluación, 14, 827
- Módulos, 80, 109, 432, 914
- Módulos de python, 149, 657
- Normalidad, 731, 734
- Normalización de datos, 84, 119, 422, 423
- Numpy, 15, 48, 124, 149–155, 157, 160, 161, 172, 174, 181, 182, 188, 354, 563, 577, 584, 598, 620, 629, 631, 634, 650, 884
- Objetivo principal, 205, 612, 775
- Operaciones
 - aritméticas, 37, 419, 421, 422, 426, 429, 434
 - básicas, 157, 160, 201, 420
 - de agrupación, 191, 196
 - elementales, 181
 - en listas, 41
 - estadísticas, 152, 153
 - matriciales, 15, 150, 153, 154, 160, 837
- Operaciones básicas, 677
- Operaciones matriciales, 124
- Operadores lógicos, 41, 42, 45–47
- Optimización, 120, 150, 153, 172–175, 181, 184, 591, 697, 703, 727, 756, 769, 772, 793, 827, 838, 861, 890, 897
- Optimización de modelos, 13, 14, 426, 440
- Paletas de colores, 203, 390
- Pandas, 6, 16, 48, 54, 59, 64, 67, 102, 103, 149, 150, 188–190, 193, 195, 196,

- 200–202, 204, 631, 670, 686, 706, 711
- Parámetro de penalización, 736, 740
- Parámetros, 13, 19, 75, 77, 78, 80, 93, 94, 96, 98, 106, 240, 436, 440, 515, 521, 532, 536, 543, 684, 692, 730, 745, 753, 756, 757, 761, 769, 775, 783, 793, 795, 803, 805, 806, 812, 869, 870, 877, 879, 886, 887, 895, 910, 914
- Personalización, 781
- Ploteo de vectores, 463
- Plotly, 150, 204, 205, 212, 215, 218, 389
- Potencia, 692
- Precisión, 3, 14, 19, 122, 167, 194, 318, 424, 733, 797, 799, 802, 803, 809, 811–815, 821, 827, 830, 888
- Predict, 789
- Preparación de datos, 11, 150, 193, 195
- Probabilidades, 13, 14, 283, 285, 286, 396, 401, 577–579, 584, 586, 594, 595, 598, 605, 607, 611, 757, 758, 769, 771, 772, 774, 775, 789, 818, 830, 831, 892, 894, 895
- Procesamiento de datos, 94, 96, 97, 419, 431, 434
- Procesamiento de imágenes, 429, 471, 888, 910
- Procesamiento de señales, 150, 174, 175, 177, 481, 485, 509, 531, 540, 882
- Proyectos, 11, 17, 19, 20, 33, 47, 54, 81, 107, 109, 119, 120, 124, 126, 128, 131, 150, 151, 167, 202, 205, 212, 218, 220, 307, 563, 649, 650, 657, 662, 669, 686, 695, 699, 888
- Pruebas de hipótesis, 10, 14, 624, 629, 706, 710, 731, 775
- Pruebas estadísticas, 636, 650
- Punto flotante, 35, 36
- Python, 1, 5, 6, 10, 14–17, 19, 20, 33–35, 37–39, 41, 42, 45, 47–50, 52, 54, 67, 75, 86, 93, 99, 102, 105, 106, 110, 120, 122, 149–152, 161, 165, 171, 172, 181, 188, 201, 205, 212, 215, 217, 218, 235, 236, 245, 279, 291, 299, 307, 334, 354, 365, 383, 388, 390, 392, 393, 395, 399, 420, 433–435, 445, 455, 479, 497, 509, 527, 563, 577, 584, 598, 607, 611–613, 619, 620, 622, 624, 629, 631, 636, 649, 650, 670, 680, 686, 706, 712, 729, 743, 745, 750, 753, 764, 773, 821, 826, 830, 833, 836, 865, 869, 875, 881, 884, 893
- Pytorch, 17, 20, 649, 677, 680, 685, 686
- Rango, 46, 113, 114, 119, 123, 279, 451, 613, 624, 734, 771, 882, 889
- Rankings, 862, 896–899
- Recomendación, 6, 821, 822, 885, 886, 888, 896, 898–900, 903, 906, 907
- Recopilación y preparación de datos, 18
- Redes neuronales no supervisadas, 910, 913
- Reducción de dimensionalidad, 649, 650, 653, 740, 830, 872, 910, 913
- Regla de adición, 581
- Regla del producto, 581, 583
- Regresión
 - bayesiana, 761, 762, 764
 - elastic net, 743, 745
 - kernel, 749, 750, 752, 753
 - lasso, 739, 740, 743, 745
 - lineal simple, 729–731, 734
 - por procesos gaussianos, 753, 756
 - probabilística, 756, 757, 760, 761, 772
 - ridge, 734–736, 738, 739, 743, 744
 - sparse, 764, 767
- Remuestreo por bootstraps, 802
- Remuestreo por ransac, 805, 808
- Representación de coordenadas, 73
- Representación de datos, 12, 220, 458, 513
- Representación de datos inmutables, 71
- Representación gráfica, 165, 616, 620, 622
- Resolución de ecuaciones diferenciales, 164, 165
- Resta, 107
- Reutilización de código, 75, 82, 107, 215
- Scikit-learn, 6, 15, 16, 19, 20, 103, 104, 154, 155, 157, 649, 650, 654, 657, 659, 662, 699, 739, 778, 781, 814, 821, 830, 865, 875, 881
- Scipy, 149, 150, 172, 174, 175, 181, 184, 188, 392, 584, 598, 607, 624
- Seaborn, 150, 203, 204, 328, 577, 635
- Sectores, 235, 299–305
- Selección
 - de modelos, 11, 15, 659
 - de variables, 744
 - del algoritmo, 18
- Series, 3, 190, 744, 761, 762, 782, 811
- Series temporales, 49, 235, 485, 490, 495, 536, 538, 650, 706

- Señales, 172, 184, 429, 430, 481, 487, 495, 910
- Simplicidad, 15, 205, 212, 218, 657, 692, 771, 789, 790, 822
- Simplificación de expresiones, 162, 165, 167, 494
- Simulaciones de datos, 607, 608, 611
- Simulación, 4, 489, 553, 563, 736, 841, 894
- Simulación de lanzamiento, 607
- Simulación numérica, 184
- Sintaxis, 122, 686
- Sintaxis simple, 14–16, 215
- Sistemas de ecuaciones, 161, 167, 169, 174, 419, 436, 438, 451, 509, 512, 513, 515, 523
- Sistemas dinámicos, 509, 903
- Stacking, 814, 815, 817
- Statsmodels, 649, 650, 706, 710–712
- Suma, 84, 106, 113, 114, 116, 117, 447, 449, 473, 740, 744, 745, 749, 772, 812, 816, 833, 834, 837, 838, 863, 864, 867, 869, 892
- Svd, 509, 510, 524, 841, 861, 862, 882–884
- Sympy, 149, 161, 165, 167, 171, 445, 446, 455, 563
- Tablas, 9, 48, 50, 818
- Tendencia central, 612, 613, 618, 619
- Tensorflow, 16, 17, 19, 20, 154, 155, 157, 649, 662, 665, 666, 669, 670, 673, 677, 680, 686, 687, 689, 690, 692, 695
- Teorema
 - de bayes, 577, 761, 788, 789
 - de euclides, 432
 - de fermat, 447
 - de gauss-markov, 520, 521
 - de green, 551
 - de la diagonalización, 519
 - de mercer, 526
 - de wilson, 433
 - del coseno, 493
 - del núcleo, 449
 - del seno, 492
 - del valor extremo, 543
 - del valor medio, 546
- Teoría de probabilidades, 577, 581, 584, 598, 629
- Tipos de argumentos, 86, 91–93, 96
- Tipos de datos, 34, 37, 41, 47, 861, 914
- Tracking por filtro de kalman, 862
- Transformación lineal, 419, 449, 454, 471, 474, 511–513, 519
- Transformadas, 175, 179, 183, 184
- Transformadas trigonométricas, 485
- Trigonometría, 165, 419, 478, 479, 481, 485, 497
- Tuplas, 67, 70–74
- Typeerror, 127
- Uso de funciones, 75, 81
- Uso de métodos, 103, 104
- Validación de modelos, 14, 15
- Valueerror, 126
- Varianza, 521, 613, 614, 618, 620–622, 731, 735, 736, 745, 750, 754, 757, 808, 811, 827, 828, 866, 869, 872, 873, 876, 886, 903
- Variedad, 896
- Variedad de capas, 695
- Vectores, 419, 455, 456, 458, 460, 465, 792, 793, 878, 883, 886, 915
- Velocidad, 546, 699, 705
- Ventajas, 153, 205, 220, 692, 710, 735, 771, 775, 837, 866, 872, 875, 879, 882, 885, 889
- Verificación de resultados, 168, 171
- Versatilidad, 10, 16, 86, 151, 204, 650, 662, 699, 893, 903
- Visualizaciones estadísticas, 710
- Visualizaciones estadísticas, 203
- Visualización, 17, 33, 104, 392, 395, 399, 636, 872
 - de datos, 12, 235, 463, 635
 - de datos geoespaciales, 150
 - de distribuciones, 203
- Voting, 727, 817, 818, 821
- Xgboost, 695–697, 699, 702

Matemáticas para

Machine Learning

Fundamentos de Ciencia de Datos con Python

Este libro establece un puente entre los conceptos matemáticos teóricos y la programación en Python, proporcionando un tratamiento sistemático de los fundamentos matemáticos para las aplicaciones de machine learning e inteligencia artificial. Se presenta una progresión estructurada desde los principios fundamentales hasta las metodologías avanzadas, incorporando tanto exposición teórica como implementación práctica.

Temas centrales:

Fundamentos computacionales:

- ✓ Programación en Python con ecosistemas NumPy, Sympy, Pandas y Matplotlib,
- ✓ **Álgebra Lineal:** Espacios vectoriales y matriciales, descomposición de valores propios y transformaciones lineales.
- ✓ **Cálculo multivariable:** Derivadas parciales, gradientes y teoría de optimización .
- ✓ **Probabilidad y estadística:** Modelado probabilístico, inferencia estadística y pruebas de hipótesis.
- ✓ **Aprendizaje supervisado:** Algoritmos de clasificación y regresión con análisis matemático.
- ✓ **Aprendizaje no supervisado:** Algoritmos de clustering y técnicas de reducción de dimensionalidad.
- ✓ **Redes neuronales:** Fundamentos matemáticos de arquitecturas de deep learning.
- ✓ **Aplicaciones especializadas:** Sistemas de recomendación y procesamiento de lenguaje natural.

Características pedagógicas:

- ✓ **Ejemplos comprensivos:** Más de 500 ejemplos de código con documentación.
- ✓ **Metodología progresiva:** Ejercicios graduales con soluciones analíticas completas
- ✓ **Aplicaciones del mundo real:** Estudios de caso y conjuntos de datos derivados de la industria
- ✓ **Análisis Visual:** Visualización matemática utilizando bibliotecas de gráficos profesionales

Experiencia de los autores

Los autores cuentan con amplia experiencia en academia y en industria, habiendo trabajado en empresas tecnológicas líderes incluyendo Amazon Web Services, Google y JPMorgan Chase en Estados Unidos y en universidades alrededor del mundo.

Jheser Guzman
Jhohanser Guzman