

MACHINE LEARNING BÁSICO

# Introducción al Machine Learning

Conceptos fundamentales para estudiantes de pregrado

---

# Contenido

---

01

## Introducción al Machine Learning

Conceptos básicos, tipos de aprendizaje y flujo de trabajo

02

## Preprocesamiento de Datos

Transformando datos brutos en información útil

03

## Aprendizaje Supervisado: Regresión

Regresión lineal simple, múltiple y evaluación

04

## Aprendizaje No Supervisado

Clustering y algoritmo K-Means

05

## Ética en Machine Learning

Sesgos, implicaciones éticas y responsabilidad



01

# Introducción al Machine Learning

---

Conceptos básicos y fundamentos del aprendizaje automático



# ¿Qué es el Machine Learning?

## Definición

El **Machine Learning** o Aprendizaje Automático es una rama de la Inteligencia Artificial que permite a las computadoras **aprender de datos** sin ser programadas explícitamente para cada tarea.

En lugar de seguir instrucciones fijas, los algoritmos identifican patrones en los datos y toman decisiones basadas en ese aprendizaje.

## Programación Tradicional vs ML

Programación Tradicional:

**Datos** + Reglas → Resultados

Machine Learning:

**Datos** + Resultados → Reglas/Modelo

## Ejemplos Cotidianos

✉ Filtros de spam

🗣 Asistentes virtuales

🎬 Recomendaciones Netflix

🛒 Recomendaciones Amazon

## Tipos de Machine Learning

### 🏷 Supervisado

Aprende de datos **etiquetados** con respuestas correctas.

Ej.: Clasificar emails como spam/no spam

### 🔍 No Supervisado

Encuentra patrones en datos **sin etiquetas**.

Ej.: Agrupar clientes por comportamiento

### 🧠 Por Refuerzo

Aprende mediante **prueba y error** recibiendo recompensas.

Ej.: Entrenar agentes para jugar ajedrez



# Tipos de Aprendizaje Automático



## Supervisado

El algoritmo aprende de datos **etiquetados** que contienen las respuestas correctas. El objetivo es aprender una función que mapee entradas a salidas conocidas.

### Clasificación

Predice **categorías o clases** discretas.

- Email: spam / no spam
- Imagen: gato / perro / pájaro
- Tumor: benigno / maligno

### Regresión

Predice valores **continuos numéricos**.

- Precio de una casa
- Temperatura del día
- Ingresos de una empresa



## No Supervisado

El modelo trabaja con datos **sin etiquetas**, buscando identificar patrones, estructuras o agrupaciones ocultas en los datos.

### Clustering

Agrupar datos **similares** en clústeres.

- Segmentación de clientes
- Agrupar artículos similares
- Identificar comunidades

### Reducción de Dimensionalidad

Reduce **complejidad** manteniendo información.

- Visualización de datos
- Compresión de imágenes
- Selección de características



## Por Refuerzo

Un **agente** aprende mediante **prueba y error** en un entorno, recibiendo recompensas o penalizaciones por sus acciones.

### Ciclo de Aprendizaje

1. Agente observa el estado
2. Toma una acción
3. Recibe recompensa/penalización
4. Ajusta su estrategia
5. Repite el proceso

### Aplicaciones

- Juegos (AlphaGo, Dota 2)
- Robótica y control
- Trading algorítmico

# Flujo de Trabajo Básico de un Proyecto de ML

1

## Identificar el Problema

Definir claramente qué se quiere resolver y los objetivos del proyecto.

2

## Recolección de Datos

Obtener datos relevantes de bases de datos, APIs, sensores, encuestas, etc.

3

## Preprocesamiento

Limpiar datos, manejar valores faltantes, normalizar, codificar variables.

4

## EDA

Análisis Exploratorio: visualizar distribuciones, detectar patrones y anomalías.

5

## Selección del Modelo

Elegir algoritmo adecuado según el tipo de problema y datos disponibles.

6

## Entrenamiento

Ajustar parámetros del modelo usando datos de entrenamiento.

7

## Evaluación

Validar el rendimiento con datos de prueba usando métricas apropiadas.

8

## Ajuste

Optimizar hiperparámetros para mejorar el rendimiento del modelo.

9

## Deployment

Implementar el modelo en producción para hacer predicciones en tiempo real.

10

## Monitoreo

Vigilar el rendimiento continuo y reentrenar si es necesario.



**Proceso Iterativo:** El flujo de trabajo de ML no es lineal. Según los resultados de evaluación, es común volver a etapas anteriores (preprocesamiento, selección de modelo, etc.) para mejorar el rendimiento.

# 02

## Preprocesamiento de Datos

---

Transformando datos brutos en información útil para el modelo



# Preprocesamiento de Datos

## ¿Por Qué Es Importante?

Los datos del mundo real son **sucios, incompletos y inconsistentes**. Sin preprocesamiento, incluso los algoritmos más avanzados producirán resultados engañosos o inexactos.

El preprocesamiento **mejora la calidad de los datos**, lo que se traduce en modelos más precisos y confiables.

### Problemas Comunes

- Valores faltantes
- Datos duplicados
- Formatos inconsistentes
- Outliers (valores atípicos)
- Escalas diferentes

### Beneficios

- Mayor precisión del modelo
- Mejor rendimiento
- Reducción de sesgos
- Menos sobreajuste
- Entrenamiento más rápido

## Limpieza de Datos

- ✓ **Manejar valores faltantes:** Imputar con media/mediana/moda o eliminar registros
- ✓ **Eliminar duplicados:** Identificar y remover registros repetidos
- ✓ **Corregir errores:** Tipos, formatos incorrectos, valores fuera de rango

## Transformación de Datos

### Normalización y Escalado

Ajustar variables a una escala común (0-1 o media=0, std=1)

**Importancia:** Algoritmos sensibles a la escala como KNN, SVM, redes neuronales

### Codificación de Variables Categóricas

**One-Hot Encoding:** Crear columnas binarias para cada categoría

**Label Encoding:** Asignar números enteros a categorías

### Detección de Outliers

Identificar valores atípicos que pueden distorsionar el modelo usando métodos estadísticos

## Ingeniería de Características

Crear nuevas características a partir de las existentes para mejorar el rendimiento del modelo.

- Extraer información de fechas (día, mes, año)
- Crear variables de relación entre características
- Agregar datos de fuentes externas



03

# Aprendizaje Supervisado

## Regresión

---

Predicción de valores continuos usando relaciones entre variables



# Regresión Lineal Simple

## Concepto Fundamental

La regresión lineal simple es el modelo más básico de ML supervisado. Predice una variable dependiente (y) a partir de una variable independiente (x) asumiendo una **relación lineal**.

### Fórmula Matemática

$$y = \beta_0 + \beta_1 x$$

$\beta_0$  (Intercepto):

Valor de y cuando x = 0

$\beta_1$  (Pendiente):

Cambio en y por cada unidad de x

### ¿Cómo Funciona?

El algoritmo encuentra **la línea recta que mejor se ajusta** a los puntos de datos, minimizando la suma de los errores cuadráticos entre los valores reales y predichos (Método de Mínimos Cuadrados).

### Supuestos del Modelo

- ✓ Relación lineal entre variables
- ✓ Independencia de observaciones
- ✓ Normalidad de residuales
- ✓ Homoscedasticidad (varianza constante)

## Ejemplo Práctico

### Predicción de Precios de Vivienda

Queremos predecir el precio de una casa basándonos en su tamaño en metros cuadrados.

Variable independiente (x): Tamaño en m<sup>2</sup>

Variable dependiente (y): Precio en USD

### Datos de Entrenamiento

Casa de 80 m<sup>2</sup> → \$120,000  
Casa de 120 m<sup>2</sup> → \$180,000  
Casa de 150 m<sup>2</sup> → \$225,000  
Casa de 200 m<sup>2</sup> → \$300,000

### Resultado del Modelo

$$\text{Precio} = 1500 \times (\text{Tamaño})$$

Predicción: Casa de 100 m<sup>2</sup> ≈ \$150,000

### Otras Aplicaciones

**Ventas vs. Publicidad:** Predecir ventas según inversión en marketing

**Temperatura vs. Altitud:** Relación entre altura y temperatura



# Regresión Lineal Múltiple

## Ampliando el Concepto

La regresión lineal múltiple extiende el concepto simple usando **múltiples variables independientes** ( $x_1, x_2, \dots, x_n$ ) para predecir la variable dependiente ( $y$ ).

### Fórmula Matemática

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

Cada coeficiente  $\beta_i$  representa el cambio en  $y$  por cada unidad de cambio en  $x_i$ , manteniendo las demás variables constantes (ceteris paribus).

### Ventajas

- + Captura relaciones más **complejas y realistas**
- + Mejora la **precisión predictiva** al usar más información
- + Permite identificar qué variables tienen **mayor impacto**

## Ejemplo Práctico Detallado

### Predicción Avanzada de Precios de Vivienda

Ahora predecimos el precio usando múltiples características:

$x_1$ : Tamaño (m<sup>2</sup>)

$x_2$ : Ubicación (1-10)

$x_3$ : N° Habitaciones

$x_4$ : Antigüedad (años)

### Modelo Resultante (Ejemplo)

Precio = 10,000 + 1,200×(Tamaño)  
+ 15,000×(Ubicación)  
+ 8,000×(Habitaciones)  
- 500×(Antigüedad)

### Interpretación de Coeficientes

- Cada m<sup>2</sup> adicional aumenta precio en **\$1,200**
- Cada punto de ubicación suma **\$15,000**
- Cada habitación extra agrega **\$8,000**
- Cada año de antigüedad resta **\$500**



# Evaluación de Modelos de Regresión

## ¿Por Qué Evaluar?

Necesitamos métricas objetivas para medir qué tan bien nuestro modelo predice valores nuevos. La evaluación nos permite **comparar diferentes modelos** y seleccionar el mejor.

## Métricas Principales

### MAE - Error Absoluto Medio

Promedio de diferencias absolutas entre valores reales y predichos

**MAE = 5000**

Interpretación: En promedio, nos equivocamos en \$5,000

### RMSE - Raíz del Error Cuadrático Medio

Penaliza errores grandes (eleva al cuadrado antes de promediar)

**RMSE = 7500**

Interpretación: Error típico de \$7,500, da más peso a predicciones muy erróneas

## R<sup>2</sup> - Coeficiente de Determinación

Indica el **porcentaje de varianza** en la variable dependiente que es explicada por el modelo. Varía entre 0 y 1.




**0.85**

R<sup>2</sup> = 85%

El modelo explica el 85% de la variabilidad

El 15% restante es error/noise

### Guía de Interpretación

-  R<sup>2</sup> < 0.5: Modelo débil
-  0.5 ≤ R<sup>2</sup> < 0.8: Modelo aceptable
-  R<sup>2</sup> ≥ 0.8: Modelo fuerte

## Mejores Prácticas

- ✓ **Usar múltiples métricas:** Ninguna métrica es suficiente por sí sola
- ✓ **Validación cruzada:** Evaluar en datos no usados en entrenamiento
- ✓ **Contexto:** Interpretar métricas según el problema específico



04

# Aprendizaje No Supervisado

## Clustering

---

Descubriendo grupos naturales en datos sin etiquetas



# Clustering y Algoritmo K-Means

## ¿Qué es el Clustering?

El clustering es una técnica de aprendizaje no supervisado que **agrupa datos similares** en conjuntos llamados **clústeres**, sin necesidad de etiquetas previas.

### Objetivo Principal

Maximizar la **similitud dentro del clúster** (intra-cluster) y minimizar la similitud entre **diferentes clústeres** (inter-cluster).

## Tipos de Clustering

### Particionado

Divide datos en K grupos no superpuestos (K-Means)

### Jerárquico

Construye árbol de clústeres anidados

## Algoritmo K-Means

El algoritmo más popular. El usuario especifica **K** (número de clústeres).

### Proceso Iterativo

1. Inicializar K centroides aleatoriamente
2. Asignar cada punto al centroide más cercano
3. Recalcular centroides como promedio de puntos
4. Repetir hasta convergencia

## Aplicaciones Reales

### Marketing

Segmentación de clientes para campañas dirigidas

### Finanzas

Detección de transacciones fraudulentas

### Entretenimiento

Sistemas de recomendación (Netflix, Spotify)

### Logística

Optimización de rutas de entrega

### Manufactura

Mantenimiento predictivo de maquinaria

## Ventajas y Limitaciones

### Ventajas

- Simple y rápido
- Escala bien

### Limitaciones

- Necesita especificar K
- Sensibilidad a inicialización

REGISTER

LOGIN

ABOUT

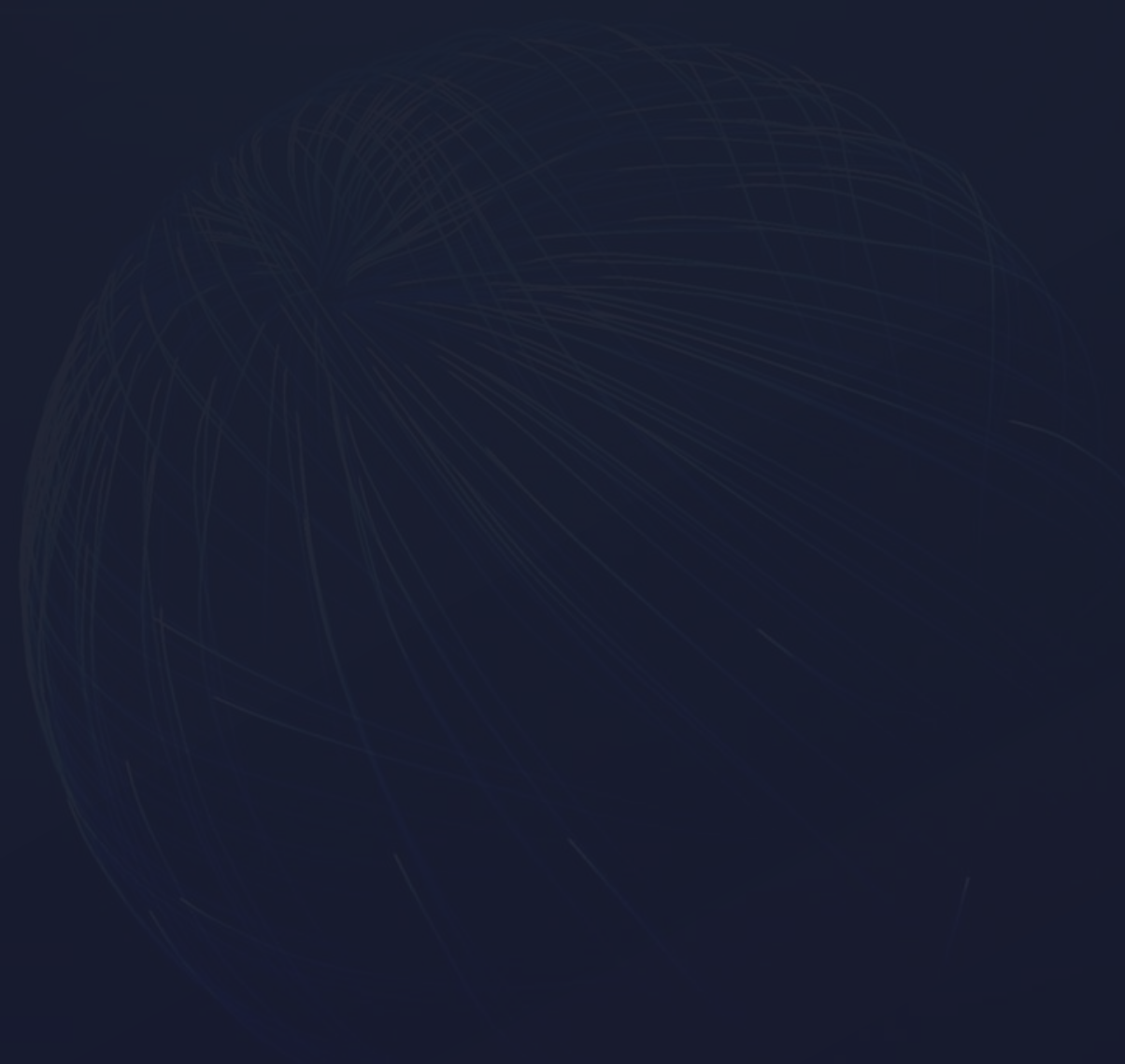
FAQ



# 05

## Ética en Machine Learning

Construyendo sistemas de IA responsables y justos





# Sesgos en Datos y Algoritmos

## ¿Qué es el Sesgo en ML?

El **sesgo** es un **error sistemático** que produce resultados injustos para ciertos grupos, perpetuando o amplificando discriminaciones existentes.

### Origen de los Sesgos

#### 1. Sesgo en los Datos

Datos desbalanceados o históricos

#### 2. Sesgo Algorítmico

Elecciones de diseño que amplifican sesgos

## Tipos de Sesgo en Datos

### Sesgo de Representación

Grupos subrepresentados en datos de entrenamiento

### Sesgo Histórico

Datos reflejan prácticas discriminatorias del pasado

### Sesgo de Medición

Variables usadas como proxy no representan bien el concepto

## Ejemplos Reales Documentados

### COMPAS - Justicia Criminal

Sistema de predicción de reincidencia usado en tribunales de EE.UU.

**Sesgo encontrado:** Falsos positivos 2x más altos para afroamericanos

### Amazon - Reclutación

Herramienta de IA para filtrar currículos (2018)

**Sesgo encontrado:** Penalizaba currículos con palabras relacionadas con mujeres

### Reconocimiento Facial

Sistemas de grandes tecnológicas

**Sesgo encontrado:** Error 35% mayor para mujeres de color vs. hombres blancos

## Consecuencias del Sesgo



**Impacto humano:** Decisiones injustas afectan vidas y oportunidades



**Riesgos legales:** Demandas y regulaciones como el AI Act europeo



**Pérdida de confianza:** Usuarios rechazan sistemas percibidos como injustos



# Implicaciones Éticas y Responsabilidad

## Transparencia

Los usuarios deben entender **cómo** se toman las decisiones automáticas.

**Derecho a la explicación:** Poder explicar por qué el modelo tomó una decisión específica

## Equidad

Los algoritmos deben tratar a todos los grupos de manera **justa**.

**Igualdad de oportunidades:** Mismo rendimiento para todos los grupos demográficos

## Rendición de Cuentas

Debe haber **mecanismos** para cuestionar decisiones incorrectas.

**Canales de apelación:** Procesos para corregir errores del sistema

## Privacidad

Protección de datos **personales** y sensibles.

**Protección de datos:** Cumplir regulaciones como GDPR

## Beneficencia

La IA debe **maximizar beneficios** y minimizar daños.

**Impacto positivo:** Usar IA para resolver problemas sociales reales

## ¿Por Qué la Ética Es Estratégica?

### Construye Confianza

Usuarios confían en sistemas percibidos como justos y transparentes

### Mitiga Riesgos

Evita daños reputacionales, legales y financieros

### Ventaja Competitiva

Diferencia la empresa como líder en IA responsable

## Mejores Prácticas

### Auditar Regularmente

Evaluar modelos en producción para detectar sesgos emergentes

### Datos Diversos

Usar datasets representativos de todas las poblaciones

### Equidad como Objetivo

Incorporar métricas de fairness en el entrenamiento

### Diversidad de Equipos

Incluir perspectivas diversas en el diseño y desarrollo



# Construyendo el Futuro de la **Inteligencia Artificial**

---

El Machine Learning es una herramienta **poderosa** que está transformando nuestro mundo.

Como futuros profesionales, tienen la **responsabilidad** de desarrollar sistemas no solo inteligentes, sino también **éticos, justos y beneficiosos** para toda la sociedad.

La clave del éxito está en combinar **conocimiento técnico** con **pensamiento crítico** y **conciencia ética**.

---

¡Gracias por su atención!

---