

guayerd

# Fundamentos IA

## Análisis con Python

Clase 6

En colaboración con  
IBM SkillsBuild





# ¡Bienvenidos!

¿Nos presentamos?

guayerd

- ¿Qué recuerdan de la clase anterior?
- ¿Qué esperan aprender?
- ¿Tienen alguna pregunta?

En colaboración con  
**IBM SkillsBuild**

# Contenidos

Por temas

- 05** • Copilot Chat y prompts  
• Demo asincrónica

- 07** • Estadística aplicada

- 06** • Limpieza y transformación

- 08** • Visualización

# Objetivos de la clase



- Pandas
- Lectura de archivos
- Estructuras principales
- Inspección y limpieza

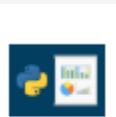
# Análisis con Python

## Limpieza y transformación

guayerd

En colaboración con  
**IBM SkillsBuild**

# Plataforma Skill Build: Python



eLearning  
Data Visualization with Python  
3 horas  1.849  150



eLearning  
Utilizar la IA generativa para el desarrollo de software  
1 hora  34.080  2.316

# Limpieza y transformación de datos

## Etapa 3 del ciclo de vida del dato

Proceso técnico para **preparar datos antes del análisis**.

- Elimina errores, inconsistencias y valores irrelevantes
- Estandariza formatos y estructuras
- Mejora la calidad y utilidad de los datos

Estandarizar = criterios

# Formatos comunes

- **CSV:** datos tabulares separados por comas
- **JSON:** objetos anidados y flexibles
- **Excel:** hojas de cálculo
- **Bases de datos:** estructuras relacionales



**Pandas integra el trabajo con todos estos formatos**

# Pandas

(Panel Data)

Librería basada en NumPy para **análisis y manipulación de datos estructurados**.

- Lectura de archivos CSV, JSON, Excel
- Filtrado, agrupación y ordenamiento
- Transformaciones con DataFrames y Series

## Comandos

- **Instalación:** `pip install pandas`
- **Uso:** `import pandas as pd`

# Estructuras principales

## Series (.s)

- Arreglo unidimensional
- Compuesta por índices y valores
- Similar a una columna

## DataFrame (.df)

- Tabla bidimensional
- Colección de Series alineadas por índice
- Base para análisis y transformación



# Lectura de archivos

## CSV

- `df_csv = pd.read_csv("nombre_archivo.csv")`
- Útil: `encoding='utf-8', sep=''`

## JSON

- `df_json = pd.read_json("nombre_archivo.json")`
- Para listas: `orient='records'`

## Excel

- `df_xls = pd.read_excel("nombre_archivo.xlsx")`
- Especificar hoja: `sheet_name="nombre_hoja"`

### Ejemplos comunes de encoding

Encoding	Características	Uso típico
UTF-8	Universal, soporta todos los idiomas	Web, Python, archivos modernos
latin1	Limitado a caracteres occidentales	Archivos antiguos, Excel en español
ISO-8859-1	Similar a latin1	Sistemas europeos
utf-8-sig	UTF-8 con marca BOM	Archivos exportados desde Excel

Siempre devuelve un DataFrame listo para trabajar

# Inspección inicial

Aspecto	Comando	Descripción
Estructura	<code>df.shape</code>	Dimensiones del conjunto
Tipos	<code>df.dtypes</code>	Variables numéricas y categóricas
Compleitud	<code>df.isnull().sum()</code>	Valores faltantes
Muestra	<code>df.head()</code>	Primeros registros
Resumen	<code>df.info()</code>	Información general

# Desafíos comunes en datos

- **Valores faltantes:** celdas vacías o NaN
- **Duplicados:** registros repetidos exactos
- **Inconsistencias:** formatos diferentes para mismo dato
- **Valores atípicos:** datos extremos que sesgan análisis
- **Tipos incorrectos:** números como texto, fechas mal formateadas



# ¿Qué harías en cada situación?



- Dataset con 10% de valores faltantes distribuidos aleatoriamente
- Registros de clientes con emails duplicados pero datos diferentes
- Precios con valores negativos en sistema de inventario
- Fechas en formatos: "2024-01-15", "15/01/2024", "Jan 15, 2024"

# Tratamiento de valores faltantes

Estrategia	Comando(s)	Descripción
Detección	<code>df.isnull()</code> , <code>df.isna()</code>	Identificar valores faltantes
Eliminación	<code>df.dropna()</code>	Eliminar registros con pocos casos
Valor constante	<code>df.fillna(0)</code>	Rellenar con cero o texto
Valor estadístico	<code>df.fillna(df['columna'].median())</code>	Rellenar con promedio o mediana

# Eliminación de duplicados

Estrategia	Comando(s)	Descripción
Detección	<code>df.duplicated()</code>	Marca filas duplicadas
Eliminación completa	<code>df.drop_duplicates()</code>	Remueve todos los duplicados
Por columnas específicas	<code>df.drop_duplicates(subset=['col1', 'col2'])</code>	Detecta duplicados según columnas elegidas
Conservar primera/última	Parámetro <code>keep</code>	Mantiene el primer o último registro

# Inconsistencias de formato

Tipo de dato	Comando(s) / Acción	Descripción
Texto	<code>.str.lower(), .str.strip()</code>	Normalizar mayúsculas/minúsculas y espacios
Fechas	<code>pd.to_datetime()</code>	Convertir a formato fecha uniforme
Números	<code>pd.to_numeric()</code>	Ajustar separadores decimales
Categorías	—	Estandarizar variaciones del mismo valor

# Manejo de valores atípicos

**Valores extremos que se alejan significativamente** del resto de los datos.

- **Detección visual:** boxplots muestran valores extremos
- **Filtrado por rango:** remover valores fuera de límites lógicos
- **Criterio de dominio:** usar conocimiento del área
- **Verificación manual:** confirmar si son errores o valores reales

# Tipos de datos incorrectos

Estrategia	Comando(s)	Descripción
Verificar tipos	<code>df.dtypes</code>	Identificar tipos actuales
Conversión manual	<code>df.astype()</code>	Cambiar tipo de columna
Fechas	<code>pd.to_datetime()</code>	Convertir a formato fecha
Numéricos	<code>pd.to_numeric()</code>	Forzar conversión numérica

# Transformaciones básicas

Operación	Comando(s)	Descripción
Filtrado	<code>df[condición]</code>	Seleccionar subconjuntos específicos
Agrupación	<code>df.groupby()</code>	Calcular estadísticas por categorías
Ordenamiento	<code>df.sort_values()</code>	Organizar registros por columnas
Selección	<code>df[['col1', 'col2']]</code>	Elegir columnas específicas

# Café de barrio



1. Calcular correlación entre temperatura y ventas
2. Identificar el mes con mejor retorno publicitario
3. Analizar relación personal vs satisfacción cliente
4. Proponer estrategia basada en datos

Mes	Ventas (\$)	Temp (°C)	Publicidad (\$)	Personal	Satisfacción
Ene	15,000	18	800	4	4.2
Feb	22,000	25	1,200	5	4.5
Mar	18,000	22	900	4	4.1
Abr	28,000	28	1,500	6	4.8
May	25,000	30	1,300	5	4.6

# Proyecto

## Tienda Aurelion

- **Documentación:** notebook Markdown
- **Desarrollo técnico:** programa Python
- **Visualización de datos:** dashboard en Power BI
- **Presentación oral:** problema, solución y hallazgos



# Limpieza de datos

Trabajo en equipo

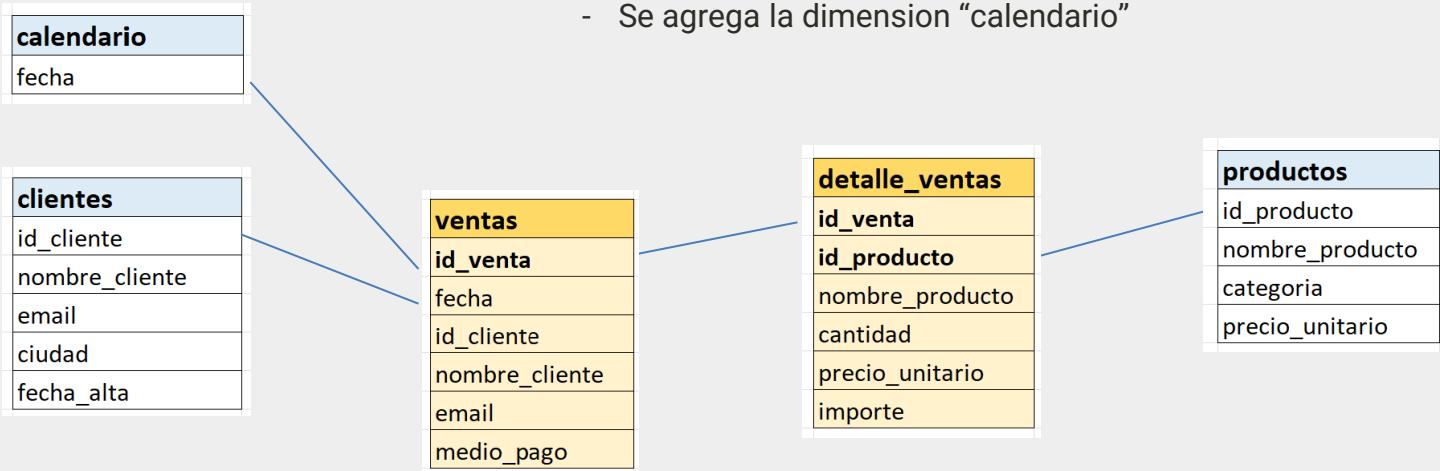


1. Usar **Copilot** para analizar problemas con dataset
2. Limpiar la base de datos
3. Documentar con **Copilot** cada paso aplicado

# Modelo de datos

## Copo de Nieve

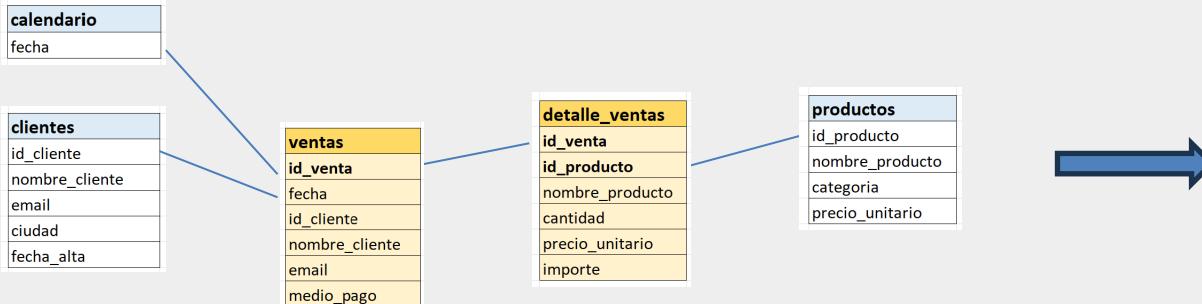
- A partir de los archivos disponibles
- No esta normalizado
- Se agrega la dimension “calendario”



# Modelo de datos

## Preparacion para ML

Es necesario pasar del modelo relacional inicial a un dataframe adecuado para trabajar con un modelo de ML



Este procedimiento requiere verificaciones de integridad y Calidad de los datos

DataFrame
<i>id_hecho</i>
<i>id_venta</i>
fecha
año
mes
día
<i>id_cliente</i>
nombre_cliente
ciudad
email
fecha_alta
<i>id_producto</i>
nombre_producto
categoria
cantidad
precio_unitario
importe
medio_pago



# Retro

¿Cómo nos vamos?

- ¿Qué fue lo más útil de la clase?
- ¿Qué parte te costó más?
- ¿Qué te gustaría repasar o reforzar?