

Guía de Laboratorio: Introducción a Pandas

Curso: Inteligencia Artificial

Docente: Carlos R. P. Tovar

1. Introducción

Pandas es una biblioteca de Python especializada en el manejo y análisis de datos. Proporciona estructuras de datos flexibles y eficientes que facilitan trabajar con datos tabulares y series temporales. Esta guía te introducirá a los conceptos fundamentales de pandas.

2. Objetivos de Aprendizaje

Al finalizar esta guía, serás capaz de:

- Crear y manipular Series y DataFrames
- Cargar datos desde diferentes formatos (CSV, Excel, JSON)
- Realizar operaciones básicas de limpieza y transformación de datos
- Filtrar, seleccionar y agrupar datos
- Realizar análisis exploratorio básico con pandas

3. Requisitos Previos

- Conocimientos básicos de Python
- Instalación de Python y pandas (`pip install pandas`)
- Editor de código o entorno de desarrollo (Jupyter Notebook recomendado)

4. Ejercicios Prácticos

4.1. Parte 1: Fundamentos de Pandas

1. Creación de Series y DataFrames

```
1 import pandas as pd
2 import numpy as np
3
4 # Crear una Serie
```

```

5 serie_ejemplo = pd.Series([10, 20, 30, 40], index=['a', 'b', 'c', 'd'])
6 print("Serie ejemplo:")
7 print(serie_ejemplo)
8
9 # Crear un DataFrame desde un diccionario
10 datos = {
11     'Nombre': ['Ana', 'Juan', 'María', 'Carlos'],
12     'Edad': [25, 32, 28, 35],
13     'Ciudad': ['Lima', 'Bogotá', 'Lima', 'Quito']
14 }
15 df_personas = pd.DataFrame(datos)
16 print("\nDataFrame personas:")
17 print(df_personas)
18

```

2. Atributos básicos de un DataFrame

```

1 # Explorar el DataFrame
2 print("Dimensiones:", df_personas.shape)
3 print("\nInformación del DataFrame:")
4 print(df_personas.info())
5 print("\nPrimeras 3 filas:")
6 print(df_personas.head(3))
7 print("\nEstadísticas descriptivas:")
8 print(df_personas.describe())
9

```

4.2. Parte 2: Carga y Exploración de Datos

1. Carga de datos desde CSV

```

1 # Cargar datos desde un archivo CSV
2 # Nota: Asegurate de tener el archivo 'ejemplo_datos.csv' en tu
   directorio
3 df = pd.read_csv('ejemplo_datos.csv')
4
5 # Explorar los datos cargados
6 print("Primeras filas:")
7 print(df.head())
8 print("\nColumnas disponibles:", df.columns.tolist())
9 print("\nTipos de datos:")
10 print(df.dtypes)
11

```

2. Manejo de valores nulos

```

1 # Identificar valores nulos
2 print("Valores nulos por columna:")
3 print(df.isnull().sum())
4
5 # Estrategias para manejar valores nulos
6 # Opción 1: Eliminar filas con valores nulos
7 df_sin_nulos = df.dropna()
8
9 # Opción 2: Rellenar valores nulos
10 df_rellenado = df.fillna({

```

```

11     'columna_numerica': df['columna_numerica'].mean(),
12     'columna_categorica': 'Desconocido'
13 })
14
15 print(f"DataFrame original: {df.shape}")
16 print(f"DataFrame sin nulos: {df_sin_nulos.shape}")
17

```

4.3. Parte 3: Selección y Filtrado de Datos

1. Selección de columnas

```

1 # Seleccionar una columna (devuelve una Serie)
2 nombres = df['Nombre']
3 print("Columna Nombre:")
4 print(nombres.head())
5
6 # Seleccionar mltiples columnas (devuelve un DataFrame)
7 subconjunto = df[['Nombre', 'Edad', 'Ciudad']]
8 print("\nSubconjunto de columnas:")
9 print(subconjunto.head())
10

```

2. Filtrado de filas

```

1 # Filtrar por condici n
2 mayores_30 = df[df['Edad'] > 30]
3 print("Personas mayores de 30 a os:")
4 print(mayores_30)
5
6 # Mltiples condiciones
7 limenos_mayores_25 = df[(df['Ciudad'] == 'Lima') & (df['Edad'] >
8     25)]
9 print("\nLime os mayores de 25 a os:")
10 print(limenos_mayores_25)
11
12 # Filtrar con isin()
13 ciudades_seleccionadas = df[df['Ciudad'].isin(['Lima', 'Bogot '])]
14 print("\nPersonas de Lima o Bogot :")
15 print(ciudades_seleccionadas)

```

4.4. Parte 4: Transformación de Datos

1. Operaciones sobre columnas

```

1 # Crear una nueva columna
2 df['Edad_en_decadas'] = df['Edad'] / 10
3 print("DataFrame con nueva columna:")
4 print(df.head())
5
6 # Aplicar una funci n a una columna
7 df['Nombre_mayusculas'] = df['Nombre'].str.upper()
8 print("\nNombres en may sculas:")
9 print(df[['Nombre', 'Nombre_mayusculas']].head())
10

```

2. Agrupación y agregación

```
1 # Agrupar por ciudad y calcular estadísticas
2 agrupado_ciudad = df.groupby('Ciudad')
3 print("Conteo por ciudad:")
4 print(agrupado_ciudad.size())
5
6 print("\nEdad promedio por ciudad:")
7 print(agrupado_ciudad['Edad'].mean())
8
9 # Múltiples agregaciones
10 resumen = df.groupby('Ciudad').agg({
11     'Edad': ['mean', 'min', 'max', 'count'],
12     'Nombre': 'count'
13 })
14 print("\nResumen estadístico por ciudad:")
15 print(resumen)
16
```

4.5. Parte 5: Ordenamiento y Manipulación de Índices

1. Ordenamiento de datos

```
1 # Ordenar por una columna
2 df_ordenado_edad = df.sort_values('Edad')
3 print("Ordenado por edad (ascendente):")
4 print(df_ordenado_edad[['Nombre', 'Edad']].head())
5
6 # Ordenar por múltiples columnas
7 df_ordenado_ciudad_edad = df.sort_values(['Ciudad', 'Edad'],
8     ascending=[True, False])
9 print("\nOrdenado por ciudad (A-Z) y edad (descendente):")
10 print(df_ordenado_ciudad_edad[['Nombre', 'Ciudad', 'Edad']].head())
```

2. Manejo de índices

```
1 # Establecer una columna como índice
2 df_indexado = df.set_index('Nombre')
3 print("DataFrame con Nombre como índice:")
4 print(df_indexado.head())
5
6 # Resetear el índice
7 df_reseteado = df_indexado.reset_index()
8 print("\nDataFrame con índice reseteado:")
9 print(df_reseteado.head())
10
```

5. Ejercicios de Práctica

5.1. Ejercicio 1: Análisis de Datos de Ventas

Carga el siguiente dataset de ventas y realiza las operaciones solicitadas:

```

1 datos_ventas = {
2     'Fecha': ['2023-01-01', '2023-01-01', '2023-01-02', '2023-01-02', '2023-01-03'],
3     'Producto': ['A', 'B', 'A', 'C', 'B'],
4     'Cantidad': [10, 5, 8, 12, 6],
5     'Precio_Unitario': [25.0, 40.0, 25.0, 30.0, 40.0]
6 }
7 df_ventas = pd.DataFrame(datos_ventas)

```

1. Calcula el total de ventas por producto (Cantidad * Precio Unitario)
2. Encuentra el producto con mayores ventas totales
3. Calcula el promedio de ventas por día

5.2. Ejercicio 2: Limpieza de Datos

Dado el siguiente dataset con problemas, realiza las operaciones de limpieza:

```

1 datos_problema = {
2     'Nombre': ['Ana', 'Juan', None, 'Mar a', 'Carlos'],
3     'Edad': [25, 32, 28, None, 35],
4     'Puntuacion': [85, 92, 78, 88, None]
5 }
6 df_problema = pd.DataFrame(datos_problema)

```

1. Identifica los valores nulos en cada columna
2. Rellena los valores nulos en Edad con la mediana y en Puntuacion con la media
3. Elimina las filas donde el Nombre es nulo

6. Recursos Adicionales

- Documentación oficial de Pandas: <https://pandas.pydata.org/docs/>
- Pandas Cheat Sheet: https://pandas.pydata.org/Pandas_Cheat_Sheet.pdf
- Tutoriales interactivos: <https://www.w3schools.com/python/pandas/default.asp>
- Datasets para practicar: <https://www.kaggle.com/datasets>

7. Conclusión

Pandas es una herramienta esencial para cualquier persona que trabaje con datos en Python. Esta guía ha cubierto los conceptos fundamentales, pero hay muchas más funcionalidades por explorar. La práctica constante es clave para dominar esta biblioteca.