

# **Actividad de Ciencia de Datos**

## **Creación de Entorno Local**

**María del Carmen Rodríguez Pérez**

### **Ciencia de Datos – Introducción a la Ciencia de Datos**

**30 de noviembre de 2024**

## Índice

1. Introducción
2. Investigación sobre entornos de ciencia de datos
3. Instalación y configuración
4. Creación de un entorno virtual
5. Uso inicial de Jupyter
6. Exploración de datos
7. Documentación del proceso
8. Conclusión
9. Bibliografía

## 1. Introducción

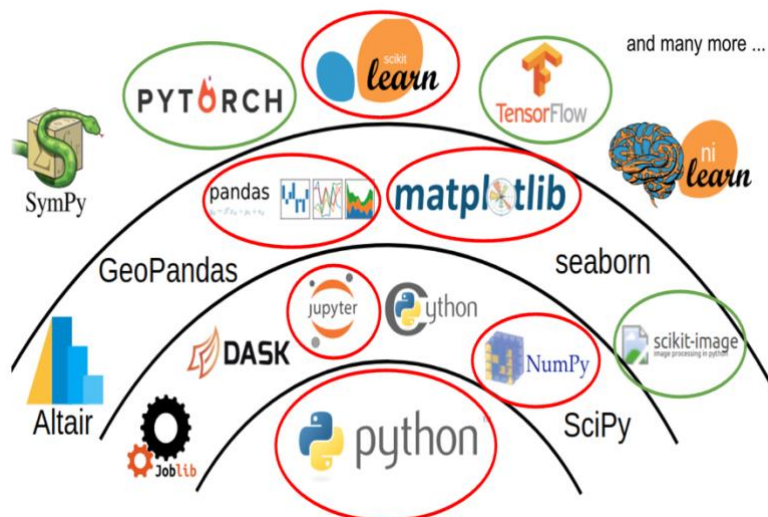
Como estudiante de ciencia de datos, entiendo la importancia de contar con un entorno de trabajo local bien configurado para el desarrollo de proyectos. La capacidad de trabajar de manera independiente en un entorno eficiente y estructurado es fundamental para llevar a cabo análisis y modelados de datos de forma profesional.

Esta actividad me permitirá aprender a configurar un entorno local utilizando herramientas clave como Anaconda o Miniconda, crear entornos virtuales personalizados e instalar librerías esenciales para la ciencia de datos. Además, exploraré datos reales a través de Jupyter Notebook, lo que me ayudará a reforzar mi comprensión práctica de las herramientas y técnicas aprendidas.



## 2. Investigación sobre entornos de ciencia de datos

- **Librerías necesarias:**



- Numpy, Pandas, Matplotlib, Jupyter, Scikit-learn, Seaborn, TensorFlow, PyTorch, entre otras.
- Explica brevemente el uso de cada una:

### 1. **Numpy**

- Uso: Proporciona soporte para cálculos matemáticos complejos y manejo eficiente de arrays y matrices multidimensionales. Es esencial para operaciones algebraicas y numéricas.

### 2. **Pandas**

- Uso: Especializada en la manipulación y análisis de datos estructurados como tablas o series temporales. Permite cargar, limpiar, transformar y analizar grandes conjuntos de datos.

### 3. **Matplotlib**

- Uso: Herramienta fundamental para crear gráficos estáticos, personalizados y de alta calidad. Es ideal para visualizar relaciones y distribuciones en los datos.

### 4. **Jupyter**

- Uso: Proporciona un entorno interactivo basado en notebooks para escribir y ejecutar código, visualizar resultados y documentar el proceso de análisis en tiempo real.

### 5. **Scikit-learn**

- Uso: Biblioteca para el aprendizaje automático, con herramientas para clasificación, regresión, clustering y reducción de dimensionalidad. Ideal para construir modelos predictivos.

### 6. **Seaborn**

- Uso: Extensión de Matplotlib para crear gráficos estadísticos avanzados con líneas de tendencia, mapas de calor y distribuciones con una sintaxis sencilla.

### 7. **TensorFlow**

- Uso: Framework avanzado para construir y entrenar modelos de aprendizaje profundo, redes neuronales y aplicaciones de inteligencia artificial.

## 8. PyTorch

- **Uso:** Alternativa a TensorFlow, popular por su facilidad de uso en investigación y desarrollo de modelos de aprendizaje profundo con soporte dinámico de gráficos computacionales.

## 9. Otras sugerencias:

- **Statsmodels:** Análisis estadístico avanzado y modelado econométrico.
- **Plotly:** Visualización interactiva y atractiva de datos.
- **XGBoost:** Optimización de algoritmos de árboles de decisión para aprendizaje automático.

Estas librerías cubren todas las etapas del flujo de trabajo en ciencia de datos: manipulación de datos, visualización, análisis estadístico y modelado avanzado.

- **Requisitos para un entorno:**
  - Python instalado (versión 3.8 o superior).
  - Gestores de paquetes como pip o conda.
  - Herramientas de virtualización (como Anaconda o Miniconda).
- **Alternativas de entornos:**
  - Anaconda vs Miniconda: Anaconda incluye más herramientas preinstaladas; Miniconda es más ligera y flexible.
  - Comparación con entornos como Docker o IDEs como PyCharm y JupyterLab.

## 3. Instalación y configuración

Pasos para instalar Anaconda o Miniconda:

1. Hay que descargar el instalador desde su página oficial.
2. Instalación seleccionando Python 3.8+.
3. Configuración inicial del entorno base.

## 4. Creación de un entorno virtual

- Comando para crear el entorno:

```
conda create -n entorno_cd python=3.8 numpy pandas  
matplotlib jupyter
```

- Breve descripción de las librerías básicas:
  - **Numpy**: Manejo eficiente de matrices y álgebra lineal.
  - **Pandas**: Análisis y manipulación de datos estructurados.
  - **Matplotlib**: Creación de gráficos estáticos y dinámicos.
  - **Jupyter**: Entorno interactivo para análisis y visualización de datos.
- Librerías adicionales sugeridas:
  - **Scikit-learn**: Modelado y aprendizaje automático.
  - **Seaborn**: Visualizaciones estadísticas avanzadas.
  - **Statsmodels**: Análisis estadístico.
  - **Plotly**: Visualizaciones interactivas.



## 5. Uso inicial de Jupyter

- Hay que abrir Jupyter Notebook y crear un archivo para:
- Importar librerías y mostrar sus versiones:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import jupyter
print(f"Numpy: {np.__version__}")
print(f"Pandas: {pd.__version__}")
print(f"Matplotlib: {plt.__version__}")
print(f"Jupyter: {jupyter.__version__}")
```

- Mostrar la versión de Python:

```
import sys
print(f"Python: {sys.version}")
```

## 6. Exploración de datos

Análisis en el notebook:

- **Representación de características:**

```
df.describe()
```

- **Diagrama de dispersión:**

```
import seaborn as sns
sns.scatterplot(x='columna1', y='columna2', data=df)
```

- **Matriz de correlación:**

```
corr = df.corr()  
sns.heatmap(corr, annot=True)
```

- **Cálculo de estadísticos:**

```
print(df.mean(), df.std(), df.isnull().sum())
```

## 7. Documentación del proceso

- Crear un archivo README.md:
- Pasos para la instalación del entorno.
- Comandos para crear y activar el entorno.
- Descripción de librerías instaladas y su propósito.

## 8. Conclusión

Configurar un entorno local para proyectos de ciencia de datos es un paso fundamental que garantiza eficiencia, reproducibilidad y personalización en el desarrollo de análisis y modelos. Esta configuración permite trabajar en un entorno controlado y adaptado a las necesidades, lo que resulta crucial en un campo tan dinámico como la ciencia de datos.

En resumen, configurar un entorno local no solo es un requisito técnico, sino una habilidad estratégica que facilita el trabajo organizado, eficiente y reproducible en ciencia de datos. Esto no solo reduce la frustración por problemas técnicos, sino que también permite concentrarse en lo más importante: resolver problemas con datos.



## 9. Bibliografía

- **Python:** <https://www.python.org/doc/>  
Tutoriales y referencia oficial para aprender Python y su biblioteca estándar.
- **Anaconda:** <https://docs.anaconda.com/>  
Guía oficial para instalar y gestionar entornos virtuales y paquetes.
- **Jupyter:** <https://jupyter.org/documentation>  
Documentación sobre instalación y uso de Jupyter Notebook y JupyterLab.
- **Real Python:** <https://realpython.com/>  
Tutoriales prácticos sobre Python, ciencia de datos y automatización.
- **Scikit-learn:** <https://scikit-learn.org/stable/documentation.html>  
Documentación de algoritmos de aprendizaje automático y ejemplos prácticos.
- **Pandas:** <https://pandas.pydata.org/docs/>  
Guías para análisis y manipulación de datos estructurados.
- **Kaggle Learn:** <https://www.kaggle.com/learn>  
Cursos interactivos gratuitos sobre ciencia de datos y Python.
- **Matplotlib:** <https://matplotlib.org/stable/users/index.html>  
Referencia para crear visualizaciones estáticas, animadas e interactivas.
- **Miniconda:** <https://docs.conda.io/>  
Guía de instalación para Miniconda, una alternativa ligera a Anaconda.
- **DataCamp:** <https://www.datacamp.com/community/tutorials/tutorial-jupyter-notebook>  
Tutorial detallado para empezar con Jupyter Notebook.