

Análisis Exploratorio del Dataset *ActiveMatter* (The Well)

Equipo 2 – Project 02: Revealing Hidden Order

22 de octubre de 2025

1. Introducción

El dataset *ActiveMatter*, parte de la colección **The Well** ([arXiv:2412.00568](https://arxiv.org/abs/2412.00568)), contiene simulaciones de materia activa que presentan transiciones de fase desde estados isotrópicos hacia estados con alineamiento colectivo (nemático o polar). Este conjunto de datos es ideal para estudiar estructuras latentes de baja dimensión mediante autoencoders o VAEs, en el contexto del proyecto *Revealing Hidden Order*.

A continuación se detallan los pasos para instalar las dependencias, descargar el conjunto de datos, cargarlo en Python y realizar un primer análisis exploratorio (EDA).

2. Instalación de librerías necesarias

Antes de comenzar, se deben instalar las librerías requeridas para la descarga y el análisis del conjunto de datos:

```
1 # --- instalar el loader oficial ---
2 pip install the_well
3
4 # --- librerias para analisis y visualizacion ---
5 pip install numpy matplotlib h5py torch tqdm scikit-learn
```

Estas librerías permiten descargar los datos, manipularlos en memoria y generar visualizaciones básicas.

3. Carga del dataset *ActiveMatter*

El paquete `the_well` proporciona una interfaz directa para acceder a los datos desde Hugging Face o descargarlos localmente.

Opción A: descarga local

```
1 the-well-download --base-path ./data \
2           --dataset ActiveMatter \
3           --split train
```

Esto descargará los archivos HDF5 en la ruta `./data/ActiveMatter/train/`.

Opción B: acceso por streaming (sin descarga completa)

```
1 from the_well.data import WellDataset
2
3 trainset = WellDataset(
4     well_base_path="hf://datasets/polymathic-ai/",
5     well_dataset_name="ActiveMatter",
6     well_split_name="train"
7 )
```

Esta opción permite leer los datos directamente desde la nube sin almacenarlos completamente en disco.

4. Exploración inicial de la estructura

```
1 import numpy as np
2 from torch.utils.data import DataLoader
3
4 # Crear DataLoader
5 loader = DataLoader(trainset, batch_size=1, shuffle=True)
6
7 # Tomar una muestra
8 batch = next(iter(loader))
9
10 print(type(batch))
11 print(batch.shape)
```

La estructura típica es de la forma $(1, 81, 256, 256, n_{\text{fields}})$, donde los 81 pasos corresponden al tiempo y la malla es de 256×256 .

5. Identificación de los campos físicos

Los campos registrados en el dataset son los siguientes:

Campo	Tipo	Descripción
c	Escalar	Concentración
U	Vectorial	Velocidad (componentes u_x, u_y)
Q	Tensorial	Tensor de orientación
S	Tensorial	Tasa de deformación

Cuadro 1: Campos físicos disponibles en el conjunto *ActiveMatter*.

Ejemplo de separación de los campos en Python:

```
1 c_field    = batch[0, :, :, :, 0]    # concentración
2 u_x_field = batch[0, :, :, :, 1]    # velocidad x
3 u_y_field = batch[0, :, :, :, 2]    # velocidad y
```

6. Visualización de un instante

Para graficar una instantánea del campo de concentración:

```
1 import matplotlib.pyplot as plt
2
3 t = 0 # primer paso temporal
4 plt.figure(figsize=(6,5))
5 plt.imshow(c_field[t], cmap='viridis')
6 plt.title(f"Concentración - paso {t}")
7 plt.colorbar()
8 plt.show()
```

Visualización del campo de velocidad (submuestreado para claridad):

```
1 plt.figure(figsize=(6,6))
2 plt.quiver(u_x_field[t][::8, ::8], u_y_field[t][::8, ::8])
3 plt.title("Campo de velocidad (submuestreado)")
4 plt.show()
```

Estas visualizaciones permiten identificar la evolución de estructuras coherentes y patrones de alineamiento en el sistema de materia activa.
