

Apéndice A

Código de Python

En este apéndice se detalla la implementación en Python del algoritmo de refinamiento de conjuntos de datos propuesto en este trabajo. El objetivo es implementar un algoritmo que refine conjuntos de datos eliminando ruido y puntos atípicos, basándose en principios de Análisis Topológico de Datos (TDA) para preservar la estructura intrínseca de los datos. Vamos a hacer un resumen del código que hemos usado en Python para filtrar nuestros conjuntos de datos.

A.1. Librerías Importantes

Definimos las librerías más relevantes para ejecutar el algoritmo y visualizar sus resultados

```
1 import numpy as np # Computacion numerica
2 import matplotlib.pyplot as plt # Crear graficos y visualizaciones
3 from ripser import ripser # Calcular la homologia persistente de una nube de
  puntos
4 from persim import plot_diagrams # Creacion diagramas de persistencia
5 from persim import bottleneck # Calculo de distancias bottleneck
```

A.2. Generación del Conjunto de Datos de Prueba

A continuación, crearemos una nube de puntos siguiendo la forma de un toroide, al que posteriormente añadiremos los outliers. La construcción se basa en las ecuaciones paramétricas canónicas del toro, que describen la superficie como el resultado de hacer girar un círculo de radio menor r alrededor de un eje que se encuentra a una distancia R (radio mayor) de su centro.

```
1 def toro(n_puntos=500, ruido=0.2, n_outliers=10):
2
3     # Definimos radios
4     R, r = 2, 1
5
6     # Creamos array de distribucion uniforme en [0,2*pi) para el angulo en el
  plano xy y el angulo de la seccion circular
7     theta = 2 * np.pi * np.random.rand(n_puntos)
8     phi = 2 * np.pi * np.random.rand(n_puntos)
9
10    # Calculamos las coordenadas a traves de las ecuaciones parametricas y las
  unimos
```

```

11 x = (R + r * np.cos(theta)) * np.cos(phi)
12 y = (R + r * np.cos(theta)) * np.sin(phi)
13 z = r * np.sin(theta)
14 toro = np.vstack([x, y, z]).T
15
16 # Generamos una matriz de ruido del mismo tamaño que la matriz del toro
17 ruido_gaussiano = np.random.normal(scale=ruido, size=toro.shape)
18
19 # Sumamos el ruido a las coordenadas de cada punto, desplazándolo
20 ligeramente de su posición original en la superficie ideal.
21 toro_ruidoso = toro + ruido_gaussiano
22
23 # Creamos puntos con coordenadas distribuidas uniformemente dentro de un
24 cubo definido por los límites [-6,6)
25 outliers = np.random.uniform(low=-6, high=6, size=(n_outliers, 3))
26 # Obtenemos el toro con el ruido en distribución normal más los outliers
27 return np.vstack([toro_ruidoso, outliers])

```

A.3. Filtrado Topológico

Esta función implementa el algoritmo de filtrado topológico. Su propósito es refinar un conjunto de datos eliminando iterativamente puntos que tienen un impacto insignificante en su estructura homológica fundamental.

```

1 def filtrar_por_impacto_topologico(puntos, epsilon, iteraciones):
2
3     # Creamos una copia de nuestro conjunto a filtrar
4     S = puntos.copy()
5
6     for iter in range(iteraciones):
7
8         # Calculamos el diagrama de persistencia de dimensión k=0,1,2
9         D = ripser(S)['dgms'][k]
10
11        # Creamos la lista donde iremos introduciendo la posición de los puntos
12        # a eliminar
13        eliminar = []
14
15        for i in range(len(S)):
16
17            # Creamos un array donde eliminamos el punto en la posición i de S y
18            # calculamos su diagrama de persistencia de dimensión k
19            S_i = np.delete(S, i, axis=0)
20            D_i = ripser(S_i)['dgms'][k]
21
22            # Calculamos la distancia bottleneck entre el diagrama de
23            # persistencia nuestro conjunto original y el del conjunto sin el
24            # punto en la posición i
25            d = bottleneck(D, D_i)

```

```

23         # Si dicha distancia es menor que nuestro umbral, agregamos la
           posicion i a eliminar
24         if d < epsilon:
25             eliminar.append(i)
26         # Si eliminar sigue vacia rompemos el bucle
27         if not eliminar:
28             break
29
30         # Eliminamos de S los puntos en las posiciones guardadas en eliminar
31         S = np.delete(S, eliminar, axis=0)
32     return S

```

Para el caso del último tipo de filtrado presentado en nuestro trabajo, dado que hay que comparar respecto tanto de H_1 como de H_2 , tenemos un algoritmo ligeramente distinto a los anteriores.

```

1 def filtrar_por_impacto_topologico(puntos, epsilon, epsilon2, iteraciones):
2
3     # Creamos una copia de nuestro conjunto a filtrar
4     S = puntos.copy()
5     for _ in range(iteraciones):
6
7         # Calculamos el diagrama de persistencia de dimension 1 y 2
8         resultados = ripser(S, maxdim=2)
9         D1 = resultados['dgms'][1]
10        D2 = resultados['dgms'][2]
11
12        # Creamos la lista donde iremos introduciendo la posicion de los puntos
           a eliminar
13        eliminar = []
14
15        for i in range(len(S)):
16
17            # Creamos un array donde eliminamos el punto en la posicion i de S y
           calculamos sus diagramas de persistencia de dimension 1 y 2
18            S_i = np.delete(S, i, axis=0)
19            resultados_i = ripser(S_i, maxdim=2)
20            D1_i = resultados_i['dgms'][1]
21            D2_i = resultados_i['dgms'][2]
22
23            # Calculamos las distancias bottleneck entre los diagramas de
           persistencia nuestro conjunto original y el del conjunto sin el
           punto en la posicion i
24            d1 = bottleneck(D1, D1_i)
25            d2 = bottleneck(D2, D2_i)
26
27            # Si dichas distancias son menores que nuestros umbrales, agregamos
           la posicion i a eliminar
28            if d1 < epsilon and d2 < epsilon2:
29                eliminar.append(i)
30
31            # Si eliminar sigue vacia rompemos el bucle
32            if not eliminar:
33                break

```

```

34
35     # Eliminamos de S los puntos en las posiciones guardadas en eliminar
36     S = np.delete(S, eliminar, axis=0)
37
38     return S

```

A.4. Análisis y Comparación de Resultados

La función *comparar persistencia* es una herramienta de validación y análisis diseñada para evaluar la efectividad del algoritmo de filtrado. Su objetivo principal es comparar la estructura topológica de dos nubes de puntos (típicamente, la original y la filtrada) y proporcionar una medida tanto cuantitativa como cualitativa de su similitud.

```

1 def comparar_persistencia(puntos1, puntos2, dim):
2
3     # Creamos los diagramas
4     dgms1 = ripser(puntos1, maxdim=2)['dgms']
5     dgms2 = ripser(puntos2, maxdim=2)['dgms']
6     dist = bottleneck(dgms1[dim], dgms2[dim])
7
8     # Representacion visual
9     plt.figure(figsize=(10,4))
10    plt.subplot(1,2,1)
11    plt.title('Diagrama persistencia original')
12    plt.scatter(dgms1[dim][:,0], dgms1[dim][:,1])
13    plt.plot([0, max(dgms1[dim][:,1])], [0, max(dgms1[dim][:,1])], 'k--')
14    plt.subplot(1,2,2)
15    plt.title('Diagrama persistencia filtrado')
16    plt.scatter(dgms2[dim][:,0], dgms2[dim][:,1])
17    plt.plot([0, max(dgms2[dim][:,1])], [0, max(dgms2[dim][:,1])], 'k--')
18    plt.show()
19
20    print(f"Distancia bottleneck dimension {dim}: {dist:.4f}")

```

A.5. Visualización de Resultados

La función *main* sirve como el punto de entrada principal para ejecutar el flujo de trabajo completo del experimento. Su rol es orquestar la secuencia de operaciones: desde la generación de los datos hasta la aplicación del algoritmo de filtrado y la posterior visualización y análisis de los resultados.

```

1 def main():
2     puntos = toro()
3     print(f"Puntos originales: {len(puntos)}")
4
5     # Filtramos el toro
6     puntos_filtrados = filtrar_por_impacto_topologico(puntos,0.02,100)
7     print(f"Puntos tras filtrado: {len(puntos_filtrados)}")
8
9     # Muestra nubes de puntos 3D

```

```

10     fig = plt.figure(figsize=(12, 6))
11
12     ax1 = fig.add_subplot(121, projection='3d')
13     ax1.scatter(puntos[:,0], puntos[:,1], puntos[:,2], s=10)
14     ax1.set_title("Conjunto original")
15
16     ax2 = fig.add_subplot(122, projection='3d')
17     ax2.scatter(puntos_filtrados[:,0], puntos_filtrados[:,1], puntos_filtrados
18    [:,2], s=10, c="#ff7f0e")
19     ax2.set_title("Conjunto filtrado")
20
21     plt.show()
22
23     # Calculamos diagramas de persistencia
24     resultados = ripser(puntos, maxdim=2)
25     diagramas = resultados['dgms']
26
27     # Muestra diagramas de persistencia
28     plot_diagrams(diagramas, show=True)
29
30     # Muestra distancias bottleneck de cada dimension
31     comparar_persistencia(puntos, puntos_filtrados, dim=0)
32     comparar_persistencia(puntos, puntos_filtrados, dim=1)
33     comparar_persistencia(puntos, puntos_filtrados, dim=2)
34
35 if __name__ == "__main__":
36     main()

```