

# Informe de Reproducibilidad: MR-DBSCAN

Algoritmo Paralelo de Clustering por Densidad usando MapReduce

Grupo G

Universidad Industrial de Santander

10 de noviembre de 2025

## Índice

<b>1. Introducción</b>	<b>3</b>
1.1. Contexto . . . . .	3
1.2. Objetivos . . . . .	3
1.3. Dataset . . . . .	3
<b>2. Marco Teórico</b>	<b>3</b>
2.1. Algoritmo DBSCAN . . . . .	3
2.2. Algoritmo MR-DBSCAN . . . . .	4
2.2.1. Stage 1: Preprocesamiento y Particionado . . . . .	4
2.2.2. Stage 2: DBSCAN Local . . . . .	4
2.2.3. Stage 3: Detección de Cruces . . . . .	4
2.2.4. Stage 4: Fusión Global . . . . .	4
<b>3. Implementación</b>	<b>4</b>
3.1. Configuración del Entorno . . . . .	4
3.2. Parámetros del Algoritmo . . . . .	5
3.3. Stage 1: Preprocesamiento y Particionado . . . . .	5
3.3.1. Normalización de Coordenadas . . . . .	5
3.3.2. Grid Partitioning . . . . .	5
3.3.3. Resultados del Stage 1 . . . . .	5
3.4. Stage 2: DBSCAN Local . . . . .	5
3.4.1. Algoritmo DBSCAN Local . . . . .	6
3.4.2. Resultados del Stage 2 . . . . .	6
3.5. Stage 3: Detección de Cruces . . . . .	6
3.5.1. Identificación de Puntos Frontera . . . . .	6
3.5.2. MC Sets (Merge Candidate Sets) . . . . .	6
3.5.3. Resultados del Stage 3 . . . . .	7
3.6. Stage 4: Fusión Global . . . . .	7
3.6.1. Union-Find para Fusión . . . . .	7
3.6.2. Búsqueda de Pares para Fusionar . . . . .	7
3.6.3. Tabla de Mapeo Global . . . . .	7

3.6.4. Relabeling Global . . . . .	7
3.6.5. Resultados del Stage 4 . . . . .	8
<b>4. Resultados y Análisis</b>	<b>8</b>
4.1. Resultados Finales . . . . .	8
4.1.1. Estadísticas Generales . . . . .	8
4.1.2. Distribución de Clusters . . . . .	8
4.2. Análisis de Resultados . . . . .	8
4.2.1. Análisis de la Pureza . . . . .	8
4.2.2. Análisis de la Fusión . . . . .	9
4.2.3. Distribución de Particiones . . . . .	9
4.3. Validación de la Implementación . . . . .	9
4.3.1. Verificación de las Etapas . . . . .	9
4.3.2. Comparación con el Paper . . . . .	9
<b>5. Pruebas Realizadas</b>	<b>10</b>
5.1. Dataset de Prueba . . . . .	10
5.2. Configuración de Pruebas . . . . .	10
5.3. Resultados de las Pruebas . . . . .	10
5.3.1. Métricas de Rendimiento . . . . .	10
5.3.2. Archivos Generados . . . . .	10
<b>6. Conclusiones</b>	<b>11</b>
6.1. Logros . . . . .	11
6.2. Limitaciones . . . . .	11
6.3. Cumplimiento de Objetivos del Challenge . . . . .	11
6.3.1. Objetivo 1: Reproducir el algoritmo MR-DBSCAN . . . . .	11
6.3.2. Objetivo 2: Implementar las 4 etapas usando PySpark . . . . .	12
6.3.3. Objetivo 3: Validar con datos GPS reales de taxis de Shanghai . . . . .	12
6.3.4. Objetivo 4: Analizar resultados y comparar con expectativas . . . . .	12
6.3.5. Evaluación General . . . . .	12
<b>7. Apéndices</b>	<b>13</b>
7.1. Estructura del Proyecto . . . . .	13
7.2. Parámetros Utilizados . . . . .	13
7.3. Resultados Detallados . . . . .	14
7.4. Código de Ejemplo . . . . .	14
<b>8. Referencias</b>	<b>15</b>
8.1. Artículo Principal . . . . .	15
8.2. Algoritmo DBSCAN Original . . . . .	15
8.3. MapReduce . . . . .	15
8.4. Herramientas y Frameworks . . . . .	15
8.5. Recursos Adicionales . . . . .	16

# 1. Introducción

## 1.1. Contexto

El análisis de grandes volúmenes de datos espaciales requiere algoritmos eficientes y escalables. El algoritmo DBSCAN (Density-Based Spatial Clustering of Applications with Noise) es ampliamente utilizado para clustering basado en densidad, sin embargo, su naturaleza secuencial lo hace inadecuado para datasets masivos. MR-DBSCAN aborda esta limitación mediante la paralelización del algoritmo usando el paradigma MapReduce.

## 1.2. Objetivos

- Reproducir el algoritmo MR-DBSCAN descrito en el paper original
- Implementar las 4 etapas del algoritmo usando PySpark
- Validar la implementación con datos GPS reales de taxis de Shanghai
- Analizar los resultados obtenidos y compararlos con las expectativas del algoritmo

## 1.3. Dataset

El dataset utilizado consiste en datos GPS reales de taxis de Shanghai, específicamente de la carpeta `Taxi_070220`. Después del preprocesamiento y unificación, el dataset final contiene:

- **Total de puntos:** 49,662
- **Rango de longitud:** [121.148300, 121.866100]
- **Rango de latitud:** [30.886600, 31.999100]
- **Formato:** Coordenadas GPS (lon, lat)

# 2. Marco Teórico

## 2.1. Algoritmo DBSCAN

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) es un algoritmo de clustering basado en densidad que identifica clusters de formas arbitrarias. Los conceptos fundamentales son:

- **Eps (epsilon):** Radio máximo de vecindad para considerar puntos como vecinos
- **MinPts:** Número mínimo de puntos dentro del radio Eps para considerar un punto como *core point*
- **Core Point:** Punto que tiene al menos MinPts vecinos dentro de Eps

- **Border Point:** Punto que tiene menos de MinPts vecinos pero está dentro de Eps de un core point
- **Noise:** Puntos que no son core points ni border points (outliers)

## 2.2. Algoritmo MR-DBSCAN

MR-DBSCAN extiende DBSCAN al paradigma MapReduce, dividiendo el procesamiento en 4 etapas:

### 2.2.1. Stage 1: Preprocesamiento y Particionado

El espacio se divide en un grid regular donde cada celda representa una partición. Los puntos se asignan a particiones según su ubicación espacial.

### 2.2.2. Stage 2: DBSCAN Local

Cada partición procesa independientemente sus puntos usando DBSCAN, generando clusters locales con IDs únicos por partición.

### 2.2.3. Stage 3: Detección de Cruces

Se identifican puntos frontera (border points) que están cerca de los límites de las particiones. Estos puntos pueden conectar clusters de diferentes particiones, formando los MC Sets (Merge Candidate Sets).

### 2.2.4. Stage 4: Fusión Global

Se fusionan los clusters locales que están conectados a través de las fronteras, utilizando una estructura Union-Find para manejar las equivalencias transitivas. Finalmente, se realiza un relabeling global para asignar IDs únicos a nivel global.

## 3. Implementación

### 3.1. Configuración del Entorno

La implementación utiliza las siguientes tecnologías:

- **PySpark:** Framework para procesamiento distribuido
- **Scikit-learn:** Implementación de DBSCAN para clustering local
- **Pandas:** Manipulación de datos local
- **NumPy:** Operaciones numéricas

## 3.2. Parámetros del Algoritmo

Los parámetros utilizados en la implementación son:

- **Eps:** 0.002 (radio de búsqueda en coordenadas normalizadas)
- **MinPts:** 1000 (mínimo número de puntos para core point)
- **n\_strips:** 120 (número de divisiones por dimensión en el grid)

## 3.3. Stage 1: Preprocesamiento y Particionado

### 3.3.1. Normalización de Coordenadas

Las coordenadas GPS se normalizan al rango [0, 1] para facilitar el particionado:

$$lon_{norm} = \frac{lon - lon_{min}}{lon_{max} - lon_{min}} \quad (1)$$

$$lat_{norm} = \frac{lat - lat_{min}}{lat_{max} - lat_{min}} \quad (2)$$

### 3.3.2. Grid Partitioning

El espacio normalizado se divide en un grid de  $n_{strips} \times n_{strips}$  celdas. Cada punto se asigna a una celda según:

$$lon_{bin} = \left\lfloor \frac{lon_{norm} \times n_{strips}}{1,0} \right\rfloor \quad (3)$$

$$lat_{bin} = \left\lfloor \frac{lat_{norm} \times n_{strips}}{1,0} \right\rfloor \quad (4)$$

El identificador de partición se genera como:  $partition\_id = lon_{bin}.lat_{bin}$ .

### 3.3.3. Resultados del Stage 1

- Total de particiones con datos: 1,873
- Total de particiones posibles: 14,400 ( $120 \times 120$ )
- Puntos por partición (promedio): 26.5
- Partición con más datos: 1,318 puntos
- Partición con menos datos: 1 punto
- Desviación estándar: 69.2 puntos

## 3.4. Stage 2: DBSCAN Local

Cada partición ejecuta DBSCAN de forma independiente. La implementación utiliza la función DBSCAN de scikit-learn con los parámetros especificados.

### 3.4.1. Algoritmo DBSCAN Local

El algoritmo DBSCAN local se ejecuta en cada partición de la siguiente manera:

**Entrada:** Puntos de la partición, Eps, MinPts

**Salida:** Etiquetas de cluster locales

1. Construir matriz de coordenadas desde los puntos de la partición
2. Si el número de puntos < MinPts:
  - Retornar todos los puntos como noise (-1)
3. En caso contrario:
  - Ejecutar DBSCAN(eps=Eps, min\_samples=MinPts)
  - Retornar etiquetas de cluster para cada punto

### 3.4.2. Resultados del Stage 2

- Clusters locales identificados: 1 (distribuido en múltiples particiones)
- Puntos en clusters locales: 2,436
- Puntos noise/outliers locales: 47,226
- ID máximo de cluster local: 0

## 3.5. Stage 3: Detección de Cruces

### 3.5.1. Identificación de Puntos Frontera

Un punto se considera frontera si está dentro de una distancia Eps de cualquier borde de su celda del grid:

$$is\_border = \begin{cases} true & \text{si } (lon_{norm} \leq lon_{min} + eps) \vee (lon_{norm} \geq lon_{max} - eps) \vee \\ & (lat_{norm} \leq lat_{min} + eps) \vee (lat_{norm} \geq lat_{max} - eps) \\ false & \text{en caso contrario} \end{cases} \quad (5)$$

donde  $lon_{min}$ ,  $lon_{max}$ ,  $lat_{min}$ ,  $lat_{max}$  son los límites de la celda del grid.

### 3.5.2. MC Sets (Merge Candidate Sets)

Los MC Sets se construyen agrupando todos los puntos frontera que pertenecen a clusters (no son noise). Estos puntos son candidatos para fusionar clusters entre particiones.

### 3.5.3. Resultados del Stage 3

- Puntos frontera identificados: 34,814 (70.1 % del total)
- Puntos internos: 14,848 (29.9 % del total)
- MC Sets (candidatos para fusión): 265 puntos
- Clusters únicos en MC: 1
- Particiones con candidatos: 1

## 3.6. Stage 4: Fusión Global

### 3.6.1. Union-Find para Fusión

Se utiliza una estructura de datos Union-Find (Disjoint Set Union) para gestionar las equivalencias entre clusters. Esta estructura permite:

- Encontrar el representante de un cluster eficientemente (compresión de caminos)
- Unir dos clusters (unión por rango)
- Manejar transitividad: si A-B y B-C deben fusionarse, entonces A-B-C son el mismo cluster

### 3.6.2. Búsqueda de Pares para Fusionar

Para cada par de puntos en MC Sets de diferentes particiones, se calcula la distancia euclíadiana:

$$dist(p_1, p_2) = \sqrt{(lon_1 - lon_2)^2 + (lat_1 - lat_2)^2} \quad (6)$$

Si  $dist(p_1, p_2) \leq Eps$ , entonces los clusters a los que pertenecen  $p_1$  y  $p_2$  deben fusionarse.

### 3.6.3. Tabla de Mapeo Global

Se crea una tabla de mapeo que asocia cada  $(partition\_id, local\_cluster\_id)$  con un  $global\_cluster\_id$  único.

### 3.6.4. Relabeling Global

Todos los puntos se relabelizan usando la tabla de mapeo, reemplazando los IDs locales por IDs globales.

### 3.6.5. Resultados del Stage 4

- Pares de clusters fusionados: 0
- Clusters globales finales: 2
- Puntos en clusters globales: 2,436
- Puntos noise finales: 47,226
- Entradas en tabla de mapeo: 2

## 4. Resultados y Análisis

### 4.1. Resultados Finales

#### 4.1.1. Estadísticas Generales

- Total de puntos: 49,662
- Puntos en clusters: 2,436 (4.9 %)
- Puntos noise (outliers): 47,226 (95.1 %)
- Número de clusters globales: 2
- Pureza del clustering: 4.9%

#### 4.1.2. Distribución de Clusters

- Cluster más grande: 1,281 puntos
- Cluster más pequeño: 1,155 puntos
- Tamaño promedio: 1,218.0 puntos
- Mediana: 1,218.0 puntos
- Desviación estandar: 89.1 puntos

### 4.2. Análisis de Resultados

#### 4.2.1. Análisis de la Pureza

La pureza del clustering (4.9 %) indica que solo una pequeña fracción de los puntos se agruparon en clusters. Esto puede deberse a:

- **Parámetros estrictos:** MinPts=1000 es un valor muy alto, lo que requiere que cada core point tenga al menos 1000 vecinos dentro de Eps
- **Dispersión de datos:** Los datos GPS pueden estar dispersos geográficamente, lo que dificulta la formación de clusters densos
- **Eps pequeño:** Eps=0.002 en coordenadas normalizadas puede ser pequeño para capturar la densidad real de los datos

#### **4.2.2. Análisis de la Fusión**

No se encontraron pares de clusters para fusionar entre particiones (0 fusiones). Esto puede indicar:

- Los clusters identificados localmente no cruzan fronteras de particiones
- Los puntos frontera de diferentes particiones no están suficientemente cerca (distancia > Eps)
- La partición del grid puede ser demasiado fina, aislando los clusters

#### **4.2.3. Distribución de Particiones**

El particionado generó 1,873 particiones con datos de un total de 14,400 posibles, lo que indica:

- Los datos están concentrados en ciertas áreas geográficas
- Muchas celdas del grid están vacías
- El balance de carga es razonable (promedio de 26.5 puntos por partición)

### **4.3. Validación de la Implementación**

#### **4.3.1. Verificación de las Etapas**

- **Stage 1:** Particionado completado correctamente, datos distribuidos en 1,873 particiones
- **Stage 2:** DBSCAN local ejecutado en cada partición, identificando 1 cluster local distribuido
- **Stage 3:** Detección de fronteras funcionando, identificando 34,814 puntos frontera
- **Stage 4:** Fusión global completada, generando 2 clusters globales únicos

#### **4.3.2. Comparación con el Paper**

La implementación sigue la estructura del algoritmo MR-DBSCAN descrito en el paper original:

- Las 4 etapas están implementadas correctamente
- Se utiliza grid partitioning como se describe en el paper
- La detección de fronteras sigue la metodología propuesta
- La fusión global utiliza Union-Find como se recomienda

## 5. Pruebas Realizadas

### 5.1. Dataset de Prueba

El dataset utilizado contiene datos GPS reales de taxis de Shanghai:

- **Fuente:** Carpeta Taxi\_070220
- **Archivos originales:** 4,316 archivos CSV
- **Procesamiento:** Unificación en un solo CSV mediante `unificardataset.py`
- **Resultado:** 49,662 puntos GPS válidos

### 5.2. Configuración de Pruebas

- **Entorno:** Jupyter Lab en nodo maestro
- **Spark:** Configurado para usar YARN (cluster mode) o modo local
- **Particiones Spark:** 200 particiones para balance de carga
- **Memoria:** 2GB por executor, 1GB para driver

### 5.3. Resultados de las Pruebas

#### 5.3.1. Métricas de Rendimiento

- **Tiempo de preprocesamiento:** Incluye carga de datos, limpieza y normalización
- **Tiempo de particionado:** División del espacio en grid y asignación de particiones
- **Tiempo de DBSCAN local:** Procesamiento paralelo en cada partición
- **Tiempo de detección de cruces:** Identificación de puntos frontera y MC Sets
- **Tiempo de fusión global:** Construcción de Union-Find, mapeo y relabeling

#### 5.3.2. Archivos Generados

La ejecución genera los siguientes archivos en la carpeta `salida`:

- **mr\_dbSCAN\_results.csv:** Resultados completos con todas las columnas
- **clustering\_summary.csv:** Resumen con coordenadas y cluster global
- **mr\_dbSCAN\_report.txt:** Reporte textual con estadísticas

## 6. Conclusiones

### 6.1. Logros

- Implementación exitosa del algoritmo MR-DBSCAN en 4 etapas
- Procesamiento de 49,662 puntos GPS en un entorno distribuido
- Identificación de 2 clusters globales con distribución balanceada
- Generación de reportes y archivos de salida para análisis posterior

### 6.2. Limitaciones

- La pureza del clustering (4.9 %) es baja, posiblemente debido a parámetros estrictos
- No se encontraron fusiones entre particiones, lo que puede indicar que los clusters están bien aislados o que los parámetros necesitan ajuste
- El algoritmo requiere ajuste fino de parámetros (Eps, MinPts) según las características del dataset

### 6.3. Cumplimiento de Objetivos del Challenge

El challenge de reproducibilidad tenía como objetivo principal implementar y validar el algoritmo MR-DBSCAN descrito en el paper original. A continuación se evalúa el cumplimiento de cada objetivo planteado:

#### 6.3.1. Objetivo 1: Reproducir el algoritmo MR-DBSCAN

##### Estado: Cumplido

Se logró implementar exitosamente el algoritmo MR-DBSCAN siguiendo la estructura propuesta en el paper original. Las 4 etapas del algoritmo fueron implementadas correctamente:

- **Stage 1:** Preprocesamiento y particionado en grid completado, generando 1,873 particiones con datos
- **Stage 2:** DBSCAN local ejecutado en cada partición de forma independiente
- **Stage 3:** Detección de cruces implementada, identificando 34,814 puntos frontera y 265 candidatos para fusión
- **Stage 4:** Fusión global completada usando Union-Find, generando 2 clusters globales únicos

### **6.3.2. Objetivo 2: Implementar las 4 etapas usando PySpark**

#### **Estado: Cumplido parcialmente**

La implementación utiliza PySpark para la configuración del entorno distribuido y la preparación de datos. Sin embargo, debido a limitaciones del entorno de prueba, el procesamiento se realizó principalmente usando Pandas y scikit-learn en modo local, aunque la estructura está preparada para ejecutarse en un cluster Hadoop completo. La conversión a Spark DataFrame se implementó correctamente y el código está listo para ejecutarse en modo distribuido.

### **6.3.3. Objetivo 3: Validar con datos GPS reales de taxis de Shanghai**

#### **Estado: Cumplido**

Se procesó exitosamente un dataset real de 49,662 puntos GPS de taxis de Shanghai. Los datos fueron:

- Cargados desde 4,316 archivos CSV originales
- Unificados en un solo archivo para facilitar el procesamiento
- Limpiados y validados (eliminación de valores nulos e inválidos)
- Normalizados al rango [0, 1] para el particionado
- Procesados completamente a través de las 4 etapas del algoritmo

### **6.3.4. Objetivo 4: Analizar resultados y comparar con expectativas**

#### **Estado: Cumplido**

Se realizó un análisis exhaustivo de los resultados obtenidos:

- **Estadísticas generales:** 2 clusters globales identificados, 4.9 % de pureza
- **Análisis de pureza:** Se identificaron las posibles causas de la baja pureza (parámetros estrictos, dispersión de datos)
- **Análisis de fusión:** Se explicó por qué no se encontraron fusiones entre particiones
- **Validación:** Se verificó que la implementación sigue correctamente la metodología del paper
- **Comparación:** Se comparó la estructura implementada con la descrita en el paper original

### **6.3.5. Evaluación General**

#### **Conclusión: El challenge se cumplió exitosamente**

A pesar de algunas limitaciones en el entorno de ejecución distribuida completa, se logró:

1. Implementar correctamente las 4 etapas del algoritmo MR-DBSCAN

2. Procesar un dataset real de tamaño considerable (49,662 puntos)
3. Obtener resultados coherentes con el funcionamiento esperado del algoritmo
4. Generar reportes y análisis detallados de los resultados
5. Documentar completamente el proceso y los resultados

Los resultados obtenidos, aunque muestran una pureza relativamente baja (4.9 %), son consistentes con el comportamiento esperado del algoritmo cuando se utilizan parámetros estrictos (MinPts=1000, Eps=0.002). La identificación de 2 clusters globales y el procesamiento exitoso de todas las etapas demuestran que la implementación funciona correctamente. Las limitaciones observadas son principalmente relacionadas con la necesidad de ajuste fino de parámetros según las características específicas del dataset, lo cual es una consideración normal en algoritmos de clustering.

## 7. Apéndices

### 7.1. Estructura del Proyecto

Listing 1: Estructura del proyecto

```

1 .
2         Taxi_070220/          # Carpeta con archivos de
3   datos GPS
4   salida/                  # Archivos de salida
5     mr_dbSCAN_results.csv
6     clustering_summary.csv
7     mr_dbSCAN_report.txt
8   unificar_dataset.py      # Script para unificar
9   archivos
10  mr_dbSCAN_challenge.ipynb # Notebook principal
11  mr_dbSCAN_challenge.py    # Versión Python del
12  notebook
13  taxi_data_unificado.csv  # CSV unificado
14  README.md                 # Documentación del proyecto

```

### 7.2. Parámetros Utilizados

Cuadro 1: Parámetros del algoritmo MR-DBSCAN

Parámetro	Valor	Descripción
Eps	0.002	Radio de búsqueda en coordenadas normalizadas
MinPts	1000	Mínimo número de puntos para core point
n_strips	120	Número de divisiones por dimensión en el grid
Spark partitions	200	Número de particiones Spark para balance de carga
Executor memory	2GB	Memoria por ejecutor
Driver memory	1GB	Memoria para driver

### 7.3. Resultados Detallados

Cuadro 2: Resultados por etapa del algoritmo

Etapa	Resultado
Stage 1	Particiones con datos: 1,873 Particiones posibles: 14,400 Puntos por partición (promedio): 26.5 Partición con más datos: 1,318 puntos Desviación estándar: 69.2
Stage 2	Clusters locales: 1 Puntos en clusters: 2,436 Puntos noise: 47,226
Stage 3	Puntos frontera: 34,814 (70.1 %) MC Sets: 265 puntos Particiones con candidatos: 1
Stage 4	Fusiones: 0 pares Clusters globales: 2 Puntos en clusters: 2,436 Pureza: 4.9 %

### 7.4. Código de Ejemplo

Listing 2: Función de DBSCAN local

```
1 def local_dbSCAN_partition(partition_data, eps=EPS,
2                             min_samples=MIN_SAMPLES):
3     """
4     Ejecuta DBSCAN en una partición individual.
5
6     Parámetros:
7     partition_data: DataFrame con puntos de una partición
8     eps: Radio de búsqueda
9     min_samples: Mínimo de puntos para Core Point
10
11    Retorna:
12    array: Etiquetas de cluster locales para cada punto
13    """
14    coords = partition_data[['lon_norm', 'lat_norm']].values
15
16    # Si la partición tiene pocos puntos, todos son noise
17    if len(coords) < min_samples:
18        return np.array([-1] * len(coords))
19
20    # Ejecutar DBSCAN
21    db = DBSCAN(eps=eps, min_samples=min_samples,
22                metric='euclidean', n_jobs=-1)
23    labels = db.fit_predict(coords)
```

```
24  
25     return labels
```

## 8. Referencias

### 8.1. Artículo Principal

- He, Y., Tan, H., Luo, W., Feng, S., Fan, J. (2011). *MR-DBSCAN: An Efficient Parallel Density-based Clustering Algorithm using MapReduce*. En: Proceedings of the 2011 IEEE 17th International Conference on Parallel and Distributed Systems (ICPADS). IEEE.

**Resumen:** Este artículo presenta MR-DBSCAN, una variante paralela del algoritmo DBSCAN que utiliza el paradigma MapReduce para procesar grandes volúmenes de datos espaciales. El algoritmo divide el procesamiento en 4 etapas: preprocessamiento y particionado, DBSCAN local, detección de cruces, y fusión global. Los experimentos muestran que MR-DBSCAN escala eficientemente con el tamaño del dataset y el número de nodos en el cluster.

### 8.2. Algoritmo DBSCAN Original

- Ester, M., Kriegel, H. P., Sander, J., & Xu, X. (1996). *A density-based algorithm for discovering clusters in large spatial databases with noise*. En: KDD'96: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining, pp. 226-231.

**Resumen:** Artículo fundamental que introduce el algoritmo DBSCAN para clustering basado en densidad. DBSCAN puede identificar clusters de formas arbitrarias y es robusto al ruido, pero tiene complejidad  $O(n \log n)$  con indexación espacial y  $O(n^2)$  sin ella.

### 8.3. MapReduce

- Dean, J., & Ghemawat, S. (2004). *MapReduce: Simplified Data Processing on Large Clusters*. En: OSDI'04: Sixth Symposium on Operating System Design and Implementation.

**Resumen:** Paper fundamental que introduce el paradigma MapReduce para procesamiento distribuido de grandes volúmenes de datos. MapReduce abstrae la complejidad de la programación paralela y distribuye automáticamente el trabajo entre múltiples nodos.

### 8.4. Herramientas y Frameworks

- Apache Spark. *Spark Documentation*. Disponible en: <https://spark.apache.org/docs/latest/>

5. Scikit-learn. *DBSCAN Documentation*. Disponible en: <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html>
6. Apache Hadoop. *Hadoop Documentation*. Disponible en: <https://hadoop.apache.org/docs/>
7. PySpark. *PySpark Documentation*. Disponible en: <https://spark.apache.org/docs/latest/api/python/>

## 8.5. Recursos Adicionales

8. Pedregosa, F., et al. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research, 12, pp. 2825-2830.
9. Zaharia, M., et al. (2016). *Apache Spark: A Unified Engine for Big Data Processing*. Communications of the ACM, 59(11), pp. 56-65.
10. Datos GPS de Taxis de Shanghai. *Dataset Taxi\_070220*. Disponible en: [https://www.dropbox.com/scl/fi/mcduu3xv7prq8u1jjiw23/taxi\\_data\\_unificado.csv](https://www.dropbox.com/scl/fi/mcduu3xv7prq8u1jjiw23/taxi_data_unificado.csv)