

## Comentario Detallado del Código:

### 1. Importaciones:

- `org.apache.spark.{SparkConf, SparkContext}`: Importa las clases necesarias para configurar y crear un contexto de Spark, que es el punto de entrada para cualquier funcionalidad de Spark.
- `org.apache.spark.graphx._`: Importa todas las clases y objetos relacionados con GraphX, la librería de Spark para computación de grafos. Esto incluye las definiciones de `Graph`, `VertexId`, `RDD[Edge[T]]`, etc.
- `org.apache.spark.rdd.RDD`: Importa la clase `RDD` (Resilient Distributed Dataset), la estructura de datos fundamental en Spark que representa una colección distribuida de elementos.
- `org.apache.log4j.{Level, Logger}`: Importa clases para configurar el nivel de registro (logging) de la aplicación.

### 2. Objeto `GraphXConnectedComponents`:

- Define un objeto Scala llamado `GraphXConnectedComponents`, que actuará como contenedor para nuestro programa principal.

### 3. Método `main`:

- `def main(args: Array[String]): Unit = { ... }`: Define el punto de entrada de la aplicación. El parámetro `args` es un array de strings que podría contener argumentos de línea de comandos (aunque en este ejemplo no se utilizan). `Unit` indica que el método no devuelve ningún valor explícito.

### 4. Configuración del Logging:

- `Logger.getLogger("org").setLevel(Level.ERROR)`: Establece el nivel de log para los mensajes generados por las librerías dentro del paquete `org` a `ERROR`. Esto significa que solo los mensajes de error o superiores (como `FATAL`) se mostrarán en la consola, reduciendo la verbosidad de Spark.
- `Logger.getLogger("akka").setLevel(Level.ERROR)`: Similar al anterior, pero para los mensajes generados por Akka, el framework de concurrencia subyacente en Spark.

### 5. Creación de la Configuración de Spark:

- `val conf = new SparkConf()`: Crea una instancia de `SparkConf`, que se utiliza para configurar la aplicación Spark.
- `.setAppName("GraphX Connected Components Example")`: Establece un nombre para la aplicación Spark, que será visible en la interfaz de usuario de Spark.
- `.setMaster("local[*]")`: Establece la URL del maestro de Spark. `"local[*]"` indica que Spark se ejecutará en modo local utilizando todos los núcleos disponibles de la máquina. Esto es común para pruebas y desarrollo.

### 6. Creación del Contexto de Spark:

- `val sc = new SparkContext(conf)`: Crea una instancia de `SparkContext`, utilizando la configuración previamente definida. El `SparkContext` es la puerta de entrada principal para interactuar con la funcionalidad de Spark.

### 7. Bloque `try-catch-finally`:

- Este bloque se utiliza para manejar posibles errores durante la ejecución del código y asegurar que el `SparkContext` se detenga correctamente al final.

### 8. Creación del Grafo de Ejemplo:

- **Vertices:**
  - `val vertices: RDD[(VertexId, String)] = sc.parallelize(Array( ... ))`: Crea un RDD de vértices. Cada vértice está representado por una

tupla (VertexId, String), donde VertexId es un Long que identifica de forma única el vértice, y String es un atributo del vértice (en este caso, una letra).

- Los vértices creados son: (1L, "A"), (2L, "B"), (3L, "C"), (4L, "D"), (5L, "E"), (6L, "F"), (7L, "G"), (8L, "H").

- **Aristas (Edges):**

- `val edges: RDD[Edge[Int]] = sc.parallelize(Array( ... ))`: Crea un RDD de aristas. Cada arista está representada por un objeto `Edge[Int]`, que contiene el ID del vértice de origen, el ID del vértice de destino y un atributo opcional (en este caso, un `Int` con valor 1, que podría representar un peso o simplemente la existencia de la conexión).
- Las aristas creadas son: (1L -> 2L), (2L -> 3L), (3L -> 4L), (5L -> 6L), (6L -> 7L), (7L -> 8L). El grafo es no dirigido implícitamente para el algoritmo de componentes conectados.

- **Creación del Grafo:**

- `val graph: Graph[String, Int] = Graph(vertices, edges)`: Crea un objeto `Graph` utilizando los RDD de vértices y aristas. El tipo del atributo del vértice es `String` y el tipo del atributo de la arista es `Int`.

## 9. Ejecución del Algoritmo Connected Components:

- `val cc = graph.connectedComponents()`: Llama al método `connectedComponents()` en el objeto `graph`. Este algoritmo calcula el ID del componente conectado al que pertenece cada vértice. Dos vértices están en el mismo componente conectado si existe un camino entre ellos. El resultado `cc` es un nuevo grafo donde el atributo de cada vértice es el ID del componente conectado al que pertenece (el ID del vértice más pequeño en ese componente).

## 10. Mostrar los Componentes Conectados:

- `println("Connected Components:")`: Imprime un encabezado.
- `cc.vertices.collect().sortBy(_._1).foreach { case (id, component) => ... }`:
  - `cc.vertices`: Accede al RDD de vértices del grafo resultante `cc`.
  - `.collect()`: Recopila todos los vértices del RDD en el nodo del driver (se debe usar con precaución en datasets grandes).
  - `.sortBy(_._1)`: Ordena los vértices por su ID original.
  - `.foreach { case (id, component) => ... }`: Itera sobre cada vértice y su ID de componente.
  - `val vertexName = graph.vertices.filter(_._1 == id).first()._2`: Busca el nombre original del vértice en el grafo original usando su ID.
  - `println(s"Vertex $id ($vertexName) belongs to component $component")`: Imprime el ID del vértice, su nombre original y el ID del componente al que pertenece.

## 11. Contar el Número de Componentes Conectados:

- `val numComponents = cc.vertices.map(_._2).distinct().count()`:
  - `cc.vertices.map(_._2)`: Extrae solo los IDs de los componentes de los vértices.
  - `.distinct()`: Elimina los IDs de componentes duplicados, dejando solo un ID único por cada componente conectado.
  - `.count()`: Cuenta el número de IDs de componentes únicos, que es el número total de componentes conectados.
- `println(s"\nNumber of connected components: $numComponents")`: Imprime el número de componentes conectados.

## 12. Calcular el Tamaño de Cada Componente:

- o `val componentSizes = cc.vertices.map { case (_, component) => (component, 1) }`: Crea un nuevo RDD donde cada vértice se transforma en una tupla `(componentId, 1)`. El 1 representa una cuenta para ese componente.
- o `.reduceByKey(_ + _)`: Agrupa las tuplas por el `componentId` y suma las cuentas (`1 + 1 + ...`) para obtener el tamaño de cada componente.
- o `.collect()`: Recopila los resultados en el nodo del driver.
- o `.sortBy(_._1)`: Ordena los componentes por su ID.
- o `println("\nSize of each component:")`: Imprime un encabezado.
- o `componentSizes.foreach { case (component, size) => println(s"Component $component: $size vertices") }`: Itera sobre los tamaños de los componentes e imprime el ID del componente y su tamaño.

## 13. Encontrar el Componente Más Grande:

- o `val largestComponent = componentSizes.maxBy(_._2)`: Encuentra la tupla en `componentSizes` con el valor más grande en la segunda posición (el tamaño del componente).
- o `println(s"\nLargest component: Component ${largestComponent._1} with ${largestComponent._2} vertices")`: Imprime el ID y el tamaño del componente más grande.

## 14. Bloque **finally**:

- o `sc.stop()`: Detiene el `SparkContext`, liberando los recursos utilizados por la aplicación Spark. Esto es crucial para una correcta gestión de recursos.

## Resultado del Código:

Dado el grafo de ejemplo definido en el código, el algoritmo de componentes conectados identificará dos grupos de vértices que están interconectados por aristas. Los vértices {1, 2, 3, 4} forman un componente conectado, y los vértices {5, 6, 7, 8} forman otro. No hay aristas entre estos dos grupos.

Por lo tanto, el resultado de la ejecución del código será similar al siguiente:

```
Connected Components:
Vertex 1 (A) belongs to component 1
Vertex 2 (B) belongs to component 1
Vertex 3 (C) belongs to component 1
Vertex 4 (D) belongs to component 1
Vertex 5 (E) belongs to component 5
Vertex 6 (F) belongs to component 5
Vertex 7 (G) belongs to component 5
Vertex 8 (H) belongs to component 5

Number of connected components: 2

Size of each component:
Component 1: 4 vertices
Component 5: 4 vertices

Largest component: Component 1 with 4 vertices
```

## Explicación del Resultado:

- **Connected Components:** Muestra para cada vértice su ID original y el ID del componente conectado al que pertenece. Observamos que los vértices 1, 2, 3 y 4 tienen el mismo ID de

componente (1), y los vértices 5, 6, 7 y 8 tienen el mismo ID de componente (5). El ID del componente suele ser el ID del vértice con el valor más pequeño dentro de ese componente.

- **Number of connected components:** Indica que se encontraron dos componentes conectados distintos en el grafo.
- **Size of each component:** Muestra el tamaño (número de vértices) de cada componente conectado. Ambos componentes en este ejemplo tienen 4 vértices.
- **Largest component:** Identifica el componente más grande en términos de número de vértices. En este caso, ambos componentes tienen el mismo tamaño, por lo que se mostrará uno de ellos (en este caso, el componente con ID 1).