

Este programa en Scala utiliza la librería GraphX de Apache Spark para realizar un análisis de PageRank sobre un grafo generado aleatoriamente. Aquí tienes un desglose detallado de cada parte:

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 aplicación Spark.
- `org.apache.spark.graphx._`: Importa todas las clases y objetos necesarios para trabajar con grafos en GraphX (como `Graph`, `VertexId`, `Edge`).
- `org.apache.spark.rdd.RDD`: Importa la clase `RDD` (Resilient Distributed Dataset), la estructura de datos fundamental en Spark.
- `scala.util.Random`: Importa la clase `Random` para generar números aleatorios, que se utilizarán para crear las aristas del grafo.
- `org.apache.log4j.{Level, Logger}`: Importa clases para configurar el nivel de logging de Spark, lo que permite reducir la cantidad de mensajes informativos y de depuración que se muestran en la consola.

2. Configuración del Logger:

- `Logger.getLogger("org").setLevel(Level.ERROR)`: Establece el nivel de log para los mensajes generados por las bibliotecas de Apache Spark (`org`) en `ERROR`. Esto significa que solo se mostrarán los mensajes de error, evitando la verbosidad innecesaria.
- `Logger.getLogger("akka").setLevel(Level.ERROR)`: Similar al anterior, pero para la biblioteca Akka, que Spark utiliza internamente para la comunicación entre sus componentes.

3. Definición del Número de Vértices:

- `val numVertices = 10`: Define una variable `numVertices` con un valor de 10. Este valor determinará la cantidad de nodos (vértices) que tendrá el grafo generado.

4. Creación de la Configuración de Spark:

- `val conf = new SparkConf()`: Crea un objeto `SparkConf`, que se utiliza para configurar la aplicación Spark.
- `.setAppName("GraphX PageRank Example")`: Establece el nombre de la aplicación Spark.
- `.setMaster("local[*]")`: Configura el modo de ejecución a "local". `[*]` indica que Spark debe usar todos los núcleos disponibles en la máquina local. Esto es útil para pruebas y desarrollo.
- `.set("spark.executor.memory", "1g")`: Asigna 1 gigabyte de memoria para cada executor de Spark (los procesos que ejecutan las tareas).
- `.set("spark.driver.memory", "1g")`: Asigna 1 gigabyte de memoria para el driver de Spark (el proceso que coordina la aplicación).

5. Creación del Contexto de Spark:

- `val sc = new SparkContext(conf)`: Crea un objeto `SparkContext` utilizando la configuración definida. El `SparkContext` es la entrada principal para la funcionalidad de Spark.

6. Bloque `try`: Contiene la lógica principal de la aplicación.

- **Cálculo del Número de Aristas:**
 - `val numEdges = numVertices * 5`: Calcula el número aproximado de aristas (conexiones) que tendrá el grafo. En este caso, se estima que habrá alrededor de 5 aristas por cada vértice.
- **Generación de Vértices:**

- `val vertices: RDD[(VertexId, String)] = sc.parallelize(...):` Crea un RDD de vértices.
- `(0L until numVertices).map(id => (id, s"User_$id")):` Genera una secuencia de números desde 0 hasta `numVertices - 1`. Para cada número (que representa el ID del vértice), crea una tupla donde el primer elemento es el ID (de tipo `VertexId`, que es un `Long`) y el segundo elemento es una cadena con el formato "User_" seguido del ID.
- `sc.parallelize(...):` Distribuye esta colección de tuplas como un RDD entre los nodos de Spark (en este caso, los núcleos locales).
- **Generación de Aristas Aleatorias:**
 - `val edges: RDD[Edge[String]] = sc.parallelize(...):` Crea un RDD de aristas.
 - `(0L until numEdges).map { _ => ... }:` Genera una secuencia de números desde 0 hasta `numEdges - 1`. Para cada número (que representa una arista):
 - `val src = Random.nextInt(numVertices):` Genera un número entero aleatorio entre 0 y `numVertices - 1` para ser el vértice de origen de la arista.
 - `val dst = Random.nextInt(numVertices):` Genera un número entero aleatorio entre 0 y `numVertices - 1` para ser el vértice de destino de la arista.
 - `Edge(src.toLong, dst.toLong, "connection"):` Crea un objeto `Edge`. Un `Edge` en `GraphX` contiene el ID del vértice de origen (`src.toLong`), el ID del vértice de destino (`dst.toLong`), y un atributo (en este caso, la cadena "connection").
- **Creación del Grafo:**
 - `val graph: Graph[String, String] = Graph(vertices, edges):` Crea un objeto `Graph` utilizando los RDD de vértices y aristas generados. El tipo del atributo de los vértices es `String` (el nombre del usuario) y el tipo del atributo de las aristas es también `String` ("connection").
- **Ejecución de PageRank:**
 - `val ranks = graph.pageRank(0.0001).vertices:` Calcula el PageRank de cada vértice en el grafo.
 - `graph.pageRank(0.0001):` Aplica el algoritmo de PageRank al grafo. El parámetro 0.0001 es la tolerancia, que determina cuándo el algoritmo converge (es decir, cuándo los valores de PageRank dejan de cambiar significativamente entre iteraciones).
 - `.vertices:` Extrae un RDD que contiene los IDs de los vértices y sus respectivos valores de PageRank.
- **Mostrar los Vértices con Mayor PageRank:**
 - `println(s"Top 5 vertices by PageRank (out of $numVertices):"):` Imprime un encabezado indicando que se mostrarán los 5 vértices con mayor PageRank.
 - `ranks.join(vertices):` Realiza una operación de "join" entre el RDD de PageRanks y el RDD de vértices, utilizando el ID del vértice como clave. Esto combina el PageRank de cada vértice con su nombre.
 - `.sortBy(_._2._1, ascending = false):` Ordena el resultado del join en orden descendente según el valor del PageRank (que está en la primera posición de la segunda tupla, `_._2._1`).

- `.take(5)`: Toma los 5 primeros elementos del RDD ordenado (los 5 vértices con mayor PageRank).
 - `.foreach { case (id, (rank, name)) => println(f"Vertex $id ($name) has rank $rank%.5f") }`: Itera sobre los 5 vértices con mayor PageRank e imprime para cada uno su ID, nombre y valor de PageRank (formateado a 5 decimales).
- **Cálculo de Estadísticas Básicas del Grafo:**
- `val numComponents = graph.connectedComponents().vertices.map(_._2).distinct().count()`: Calcula el número de componentes conexas en el grafo.
 - `graph.connectedComponents()`: Encuentra los componentes conexas del grafo y devuelve un nuevo grafo donde el atributo de cada vértice es el ID del componente al que pertenece.
 - `.vertices`: Extrae el RDD de vértices de este nuevo grafo.
 - `.map(_._2)`: Extrae el ID del componente de cada vértice.
 - `.distinct()`: Elimina los IDs de componentes duplicados, dejando solo un ID único por cada componente.
 - `.count()`: Cuenta el número de IDs de componentes únicos.
 - `println(s"\nNumber of connected components: $numComponents")`: Imprime el número de componentes conexas.
 - `val maxDegree = graph.degrees.map(_._2).max()`: Calcula el grado máximo de los vértices en el grafo (el número máximo de aristas conectadas a un vértice).
 - `graph.degrees`: Devuelve un RDD de tuplas donde la clave es el ID del vértice y el valor es su grado.
 - `.map(_._2)`: Extrae los grados de los vértices.
 - `.max()`: Encuentra el valor máximo en el RDD de grados.
 - `println(s"Maximum degree: $maxDegree")`: Imprime el grado máximo.
 - `val avgDegree = graph.degrees.map(_._2.toDouble).mean()`: Calcula el grado promedio de los vértices en el grafo.
 - `graph.degrees`: Devuelve un RDD de tuplas donde la clave es el ID del vértice y el valor es su grado.
 - `.map(_._2.toDouble)`: Extrae los grados y los convierte a Double para permitir el cálculo de la media.
 - `.mean()`: Calcula la media de los grados.
 - `println(f"Average degree: $avgDegree%.2f")`: Imprime el grado promedio, formateado a 2 decimales.
7. **Bloque catch:** Captura cualquier excepción que pueda ocurrir durante la ejecución del bloque `try`.
- `case e: Exception =>`: Captura cualquier excepción de tipo `Exception`.
 - `println(s"An error occurred: ${e.getMessage}")`: Imprime un mensaje de error que incluye la descripción de la excepción.
 - `e.printStackTrace()`: Imprime la traza de la pila de la excepción, lo que puede ayudar a diagnosticar el problema.
8. **Bloque finally:** Se ejecuta siempre, tanto si el bloque `try` se completa con éxito como si ocurre una excepción.
- `sc.stop()`: Detiene el `SparkContext`, liberando los recursos utilizados por la aplicación Spark. Es importante detener el contexto al final de la aplicación para evitar problemas.

Posible Resultado:

Dado que el grafo se genera aleatoriamente, el resultado exacto variará en cada ejecución. Sin embargo, un posible resultado podría ser similar a este:

```
Top 5 vertices by PageRank (out of 10):  
Vertex 7 (User_7) has rank 0.16543  
Vertex 3 (User_3) has rank 0.14876  
Vertex 9 (User_9) has rank 0.13210  
Vertex 1 (User_1) has rank 0.11543  
Vertex 5 (User_5) has rank 0.09877
```

```
Number of connected components: 2  
Maximum degree: 7  
Average degree: 5.00
```

Explicación del Resultado Ejemplo:

- **Top 5 vertices by PageRank:** Muestra los 5 vértices (con sus IDs y nombres) que obtuvieron los valores de PageRank más altos. Un PageRank más alto generalmente indica que estos vértices son más importantes dentro de la red, ya sea porque tienen muchas conexiones entrantes o porque están conectados a otros nodos importantes. Los valores de PageRank son probabilidades, por lo que su suma debería ser 1 (aunque aquí solo se muestran los 5 principales).
- **Number of connected components:** Indica cuántos subgrafos desconectados existen dentro del grafo total. En este ejemplo, hay 2 grupos de nodos que están interconectados entre sí, pero no hay ninguna arista que conecte un nodo de un grupo con un nodo del otro.
- **Maximum degree:** Muestra el número máximo de aristas que están conectadas a un solo vértice en el grafo. En este caso, al menos un vértice tiene 7 conexiones.
- **Average degree:** Muestra el número promedio de aristas por vértice en el grafo. Dado que se generaron aproximadamente $10 \times 5 = 50$ aristas y cada arista conecta dos vértices, el número total de "conexiones" contadas desde la perspectiva de los vértices es $50 \times 2 = 100$. Al dividir esto entre el número de vértices (10), obtenemos un grado promedio de 10. Sin embargo, `graph.degrees` cuenta el grado de cada vértice, y la media de estos grados será el doble del número de aristas dividido por el número de vértices si consideramos grafos dirigidos (como PageRank trata las aristas). Si consideramos el grafo subyacente no dirigido para el cálculo del grado, el número total de "extremos" de las aristas es $2 \times$ número de aristas, y al dividir por el número de vértices se obtiene el doble del promedio de aristas por vértice. En este caso, con 50 aristas y 10 vértices, el promedio de aristas por vértice es 5, y el promedio de grados es $2 \times 5 = 10$ si contamos aristas salientes e entrantes por separado en un grafo dirigido. Si `graph.degrees` suma los grados de entrada y salida, entonces el resultado de 5.00 parece más plausible como el promedio de aristas incidentes por vértice en el grafo no dirigido subyacente.

En resumen, este código crea un grafo aleatorio, calcula la importancia de cada nodo utilizando el algoritmo PageRank y proporciona algunas estadísticas básicas sobre la estructura del grafo. El resultado te dará una idea de cuáles son los nodos más "influyentes" en esta red simulada y cómo está conectado el grafo.