



MASTER EN DATA SCIENCE

Curso Académico 2018/2019

Trabajo Fin de Master

# Divider Greedy algorithm for performing community detection in social networks

Autor : Michel Maes Bermejo

Tutor : Jesús Sánchez-Oro Calvo

# Resumen

En este proyecto abordaremos la creación de un algoritmo capaz de detectar comunidades en redes sociales, un problema en el que los algoritmos tradicionales no ofrecen soluciones rápidas. Para ello usaremos un enfoque metaheurístico basado en la división de comunidades e intentaremos mejorar sus prestaciones a lo largo de varias iteraciones. (PENDIENTE DE TERMINAR TRAS FINALIZAR LA MEMORIA)

# Índice general

<b>1. Introducción y motivación</b>	<b>1</b>
<b>2. Objetivos</b>	<b>2</b>
<b>3. Tecnologías, Herramientas y Metodologías</b>	<b>3</b>
3.1. Tecnologías . . . . .	3
3.1.1. Java . . . . .	3
3.2. Herramientas . . . . .	4
3.2.1. Control de versiones: Git . . . . .	4
3.2.2. Entorno de desarrollo: IntelliJ . . . . .	4
3.3. Metodologías . . . . .	4
<b>4. Descripción informática</b>	<b>5</b>
4.1. Requisitos . . . . .	5
4.1.1. Requisitos funcionales . . . . .	5
4.1.2. Requisitos no funcionales . . . . .	5
4.2. Diseño e implementación . . . . .	6
4.2.1. ConstRandom: Algoritmo aleatorio . . . . .	7
4.2.2. ConstDivider: Primer aproximación . . . . .	7
4.2.3. ConstDividerGreedy: Segunda aproximación . . . . .	8
4.2.4. ConstDividerGreedy + Local Search: Última versión . . . . .	9
4.2.5. Paralelización . . . . .	11
<b>5. Estudio comparativo</b>	<b>13</b>

## *ÍNDICE GENERAL*

<b>6. Conclusiones del proyecto y trabajos futuros</b>	<b>14</b>
<b>Bibliografía</b>	<b>15</b>

# Capítulo 1

## Introducción y motivación

Desde su nacimiento, Internet ha logrado conectar a las personas de forma sencilla. El ejemplo más claro son las redes sociales, que han crecido de forma contundente los últimos años. La rapidez con la que se transmite la información ha provocado que muchas personas las utilicen como medio de información de referencia frente a los tradicionales. Esto ha desencadenado un creciente interés de diferentes marcas e incluso de otros medios de aprovechar este fenómeno para su beneficio, aprovechando la estructura de red (o grafo) en la que se basa.

Desde el punto de vista de un científico de datos, resulta interesante estudiar la estructura que forman estas redes, que constituyen un ejemplo perfecto de un grafo. Una de las áreas más interesantes que utiliza como base este tipo de redes es la detección de comunidades, que cuenta con multitud de aplicaciones.

Actualmente existen multitud de algoritmos que abordan el problema de la detección de comunidades. Los más tradicionales y exactos no son viables para abordar los grandes volúmenes de datos que generan las redes sociales, por lo que se ha optado por soluciones heurísticas, no tan exactas, pero mucho más rápidas.

La motivación de este proyecto nace de la idea de crear una solución heurística diferente a las preexistentes (normalmente basadas en algún algoritmo destructivo), optando por un enfoque destructivo.

# Capítulo 2

## Objetivos

El objetivo principal de este proyecto será la creación de un algoritmo de detección de comunidades. Para ello utilizaremos un enfoque destructivo (comenzar con todos los nodos del grafo en la misma comunidad e ir haciendo particiones) contrario a un enfoque más utilizado, constructivo (comenzar con cada nodo en su propia comunidad e ir agrupandolos). Durante la construcción, usaremos la modularidad como métrica para evaluar nuestra solución. Realizaremos una serie de iteraciones en las cuales iremos mejorando el algoritmo con el fin de obtener la mejor versión del mismo.

Finalmente, compararemos nuestro algoritmo con otros preexistentes para valorar si es efectivo el enfoque elegido. En esta fase usaremos otras métricas sobre las particiones obtenidas como la conductancia o cobertura.

# Capítulo 3

## Tecnologías, Herramientas y Metodologías

### 3.1. Tecnologías

#### 3.1.1. Java



Figura 3.1: Logo Java

Java (Figura 3.1) es un lenguaje de programación de propósito general, concurrente y orientado a objetos. Su sintaxis deriva en gran medida de C y C++. Uno de los principales atractivos de Java es su máquina virtual (JVM) que nos permite ejecutar nuestro código Java en cualquier dispositivo, independientemente de la arquitectura.

## 3.2. Herramientas

### 3.2.1. Control de versiones: Git



Figura 3.2: Logo Git

Git<sup>1</sup> (Figura 3.2) es un software de control de versiones diseñado por Linus Torvalds, pensando en la eficiencia y la confiabilidad del mantenimiento de versiones de aplicaciones cuando éstas tienen un gran número de archivos de código fuente.

### 3.2.2. Entorno de desarrollo: IntelliJ



Figura 3.3: Logo IntelliJ

IntelliJ<sup>2</sup> (Figura 3.3) es un IDE para Java desarrollado por JetBrains ideado para mejorar la productividad del programador.

## 3.3. Metodologías

Dado que pretendemos obtener la mejor versión de nuestro algoritmo, el desarrollo seguirá una metodología iterativa-incremental<sup>3</sup>, dónde en cada iteración se propondrán y evaluarán diferentes mejoras sobre el algoritmo.

---

<sup>1</sup><https://git-scm.com/>

<sup>2</sup><https://www.jetbrains.com/idea/>

<sup>3</sup><https://proyectosagiles.org/desarrollo-iterativo-incremental/>



# Capítulo 4

## Descripción informática

En este apartado se abordará la construcción del proyecto. Este proyecto consta de una colección de algoritmos de detección de comunidades, que parte desde un algoritmo de detección aleatoria, pasando por la evolución de un algoritmo destructivo hasta obtener la mejor versión del mismo aplicando una búsqueda local.

### 4.1. Requisitos

#### 4.1.1. Requisitos funcionales

Los requisitos funcionales de la aplicación son simples dado que se trata de un algoritmo:

- Debe poder leer un fichero de entrada (correspondiente al grafo a tratar).
- Dado esa entrada y un algoritmo, debe devolver los clústers resultantes al realizar la detección de comunidades.

#### 4.1.2. Requisitos no funcionales

Los requisitos no funcionales de la aplicación se reducen a la calidad de la solución ofrecida y al tiempo de ejecución:

- Maximizar la modularidad.

- Ofrecer el menor tiempo de cálculo posible<sup>1</sup>.

### Modularidad

Las soluciones de los algoritmos de detección de comunidades se construyen maximizando su valor de modularidad, una medida utilizada para medir la fuerza de división de un grafo en clústers. Esta métrica se define cómo la fracción de enlaces que caen dentro de los clústers dados menos el valor esperado que dicha fracción hubiese recibido si los enlaces se hubiesen distribuido al azar. De forma matemática:

$$Md(S, G) = \sum_{j=1}^{max(S)} (e_{jj} - a_j^2)$$

dónde  $S$  es la solución y  $G$  el grafo, siendo  $e_{jj}$  la fracción de enlaces con ambos vértices finales en el mismo grupo:

$$e_{jj} = \frac{|\{(v, u) \in E : S_v = S_u = j\}|}{|E|}$$

y  $a_j$  la fracción de enlaces con al menos un extremo en la misma comunidad:

$$a_j = \frac{|\{(v, u) \in E : S_v = j\}|}{|E|}$$

siendo en ambos caso  $E$  el conjunto de las aristas de  $G$ .

La modularidad toma valor en el intervalo  $[-.5, 1)$

## 4.2. Diseño e implementación

A continuación, se detallará la construcción incremental del algoritmo, explicando cada versión creada. Las premisas para todos los algoritmos son las siguientes:

- Se comienza con todos en la misma comunidad.
- Se asume un grafo conexo

<sup>1</sup>Existen algoritmos que ofrecen soluciones óptimas a este problema, pero no son aplicables a grafos grandes dado que requieren demasiado tiempo de computo.

**Grafo conexo**

Un grafo conexo  $\overline{G}$  es todo grafo  $G$  dónde para cada par de vertices  $(x, y)$  existe un camino  $P$ .

Para estimar cómo de buena es cada versión del algoritmo usaremos un grafo de ejemplo representativo que llamaremos *Instancia de Prueba* (214 nodos y 1954 aristas).

**4.2.1. ConstRandom: Algoritmo aleatorio**

En esta primera aproximación (Algoritmo 1) que nos sirve como punto de partida y para posterior comparación, generamos los clústers de manera aleatoria; se crea un número aleatorio de clústers y se distribuyen los nodos de forma aleatorio en ellos.

**Algorithm 1** ConsRandom algorithm

---

```

1: procedure CONSTRANDOM( $G$ )                                ▷ The graph 'G'
2:    $S \leftarrow \text{EmptySolution}(G)$ 
3:    $\text{numClusters} \leftarrow \text{Random}(N(g))$                     ▷ N returns n° of graph nodes
4:   for  $i = 0$  to  $\text{numClusters}$  do
5:      $\text{createEmptyCluster}(S)$ 
6:   for  $i = 0$  to  $N(S)$  do
7:      $\text{rnd} \leftarrow \text{Random}(\text{numClusters})$ 
8:      $\text{AssignToCluster}(S, i, \text{rnd})$                         ▷ Assign node  $i$  to cluster  $\text{rnd}$  in solution  $S$ 
9:   return  $S$ 

```

---

**Modularidad para Instancia de Prueba: -0.007**

---

**4.2.2. ConstDivider: Primer aproximación**

Este algoritmo supone el primer intento de aplicar un enfoque divisor a la detección de comunidades. Los pasos del algoritmo a seguir son los siguiente (Ver junto al Algoritmo 2):

1. Comenzamos con una sola comunidad que contiene todos los nodos.
2. Se busca el nodo peor conectado (menor número de aristas) ó el primero que encuentre una arista.
3. Se crea un nuevo clúster con ese nodo y se añaden los adyacentes en la medida en la que aumente la modularidad. Se descartan los adyacentes de los nodos los cuales no han mejorado la modularidad. De esta manera, se crearán dos clústers (A y B).

4. Repetimos los pasos 1-3 con el clúster B hasta que no se generen nuevos clusters.

---

**Algorithm 2** ConstDivider algorithm
 

---

```

1: procedure CONSTDIVIDER( $G$ )                                ▷ The graph 'G'
2:    $S \leftarrow \text{EmptySolution}(G)$ 
3:    $S' \leftarrow S$ 
4:    $c \leftarrow 0$                                             ▷ Current cluster index
5:   while  $c < \text{Size}(S')$  do
6:      $K \leftarrow \text{EmptySet}()$                                 ▷ Checked Nodes
7:      $S'' \leftarrow S'$ 
8:      $\text{createEmptyCluster}(S'')$ 
9:      $wcn \leftarrow -\text{GetWorstConnectedNode}(S'')$ 
10:     $\text{Move}(wcn, c + 1, S'')$                                 ▷ Move node  $wcn$  to cluster  $c + 1$  in solution  $S''$ 
11:     $A \leftarrow \text{Adjacents}(G, wcn)$ 
12:    while  $\text{Size}(A) > 0$  do
13:       $i \leftarrow -\text{Poll}(A)$ 
14:       $K \leftarrow -\text{Append}(K, i)$ 
15:       $\text{Move}(i, c + 1, S'')$ 
16:      if  $\text{Md}(S'', G) > \text{Md}(S', G)$  then
17:         $S' \leftarrow S''$ 
18:      for  $j \in \text{Adjacents}(i)$  do
19:        if  $j \notin K$  &  $j \notin \text{Cluster}(S'', c)$  then
20:           $A \leftarrow -\text{Remove}(A, \text{node})$ 
21:           $K \leftarrow -\text{Append}(K, \text{node})$ 
22:      if  $\text{Md}(S', G) > \text{Md}(S, G)$  then
23:         $S \leftarrow S'$ 
24:      else
25:         $c \leftarrow c + 1$ 
26:  return  $S$ 

```

---

**Modularidad para Instancia de Prueba: 0.096**

---

### 4.2.3. ConstDividerGreedy: Segunda aproximación

Podemos observar como la mejora en cuanto a modularidad no es destacable (recordemos que la modularidad toma un valor entre -0.5 y 1.0). Es necesario utilizar otro enfoque. El problema de este algoritmo es que una mala elección del primer nodo puede provocar que no se formen correctamente las comunidades. Aleatorizar esta elección de este primer nodo nos abre la puerta a realizar múltiples iteraciones comenzando desde distintos nodos. Los pasos a seguir en este algoritmo serían (Ver junto al Algoritmo 3):

- Se elige el primer nodo de forma aleatoria, añadiendo, como antes, los nodos adyacentes si aumenta la modularidad.
- Se crea una lista de candidatos, que contiene a todos los vertices ordenados por el número de aristas que vayan a otro clúster de mayor a menor.
- A partir de una función voraz (Ver Función 4.2.3), obtenemos un umbral  $mu$ , que restringe la lista a los vértices con un número de vertices mayor que  $mu$ .
- El siguiente vértice se selecciona al azar de esta lista restringida.
- Se repite este proceso hasta que no mejora la modularidad.

#### Función Voraz

Dado  $g_{max}$  y  $g_{min}$  como los valores mayor y menos de nuestra función voraz (el número de aristas) y  $\alpha$  un parámetro en el rango  $[0,1]$  (que indicará cómo de voraz es nuestro algoritmo, siendo 1 completamente aleatorio y 0 totalmente voraz), el umbral  $mu$  se calcula siguiendo la siguiente función:

$$mu = g_{max} - \alpha * (g_{max} - g_{min})$$

Para probar el algoritmo se han utilizado los siguientes valores de  $\alpha$ : 0.0, 0.25, 0.5, 0.75, 1.

#### 4.2.4. ConstDividerGreedy + Local Search: Última versión

La modularidad obtenida tras aplicar un enfoque voraz mejora de forma considerable. A pesar de ello, intentaremos aumentarla aún más partiendo de un heurístico no contemplado hasta ahora, la búsqueda local. Esta búsqueda local utilizará la solución de cada iteración del algoritmo anterior (ConstDividerGreedy) e intentará mejorarla cambiando algunos nodos de clúster. La búsqueda local implementada se detalla a continuación (Ver junto al Algoritmo 4):

- Para cada clúster, se crea una lista de candidatos a moverse. Esta lista de candidatos estará ordenada por el % de aristas que caen fuera del clúster.
- Cada nodo de esta lista de candidatos se intenta mover a otro clúster (incluido a un nuevo clúster). Si la modularidad mejora, se actualiza la solución, pero se siguen explorando otros movimientos.

**Algorithm 3** ConstDividerGreedy algorithm

---

```

1: procedure CONSTDIVIDERGREEDY( $G, \alpha$ ) ▷ The graph 'G' and param ' $\alpha$ '
2:    $S \leftarrow \text{EmptySolution}(G)$ 
3:    $S' \leftarrow S$ 
4:    $c \leftarrow 0$  ▷ Current cluster index
5:    $b \leftarrow \text{GetRandom}(S', c)$  ▷ Base node
6:    $\text{refactorIters} \leftarrow 10$ 
7:   while  $c < \text{Size}(S') \ \& \ b \neq -1 \ \& \ \text{refactorIters} > 0$  do
8:      $K \leftarrow \text{EmptySet}()$  ▷ Checked Nodes
9:      $S'' \leftarrow S'$ 
10:     $\text{createEmptyCluster}(S'')$ 
11:     $\text{candidates} \leftarrow \text{RestrictedList}(G, c, \alpha)$ 
12:    if  $\text{Empty}(\text{candidates})$  then
13:       $b \leftarrow \text{GetRandom}(\text{candidates})$ 
14:    else
15:       $b \leftarrow \text{GetRandom}(S', c)$ 
16:    if  $b = -1$  then
17:       $c \leftarrow 0$ 
18:       $\text{refactorIters} \leftarrow \text{refactorIters} - 1$ 
19:      continue
20:     $\text{Move}(b, c + 1, S'')$  ▷ Move node  $b$  to cluster  $c + 1$  in solution  $S''$ 
21:     $A \leftarrow \text{Adjacents}(G, wcn)$ 
22:    while  $\text{Size}(A) > 0$  do
23:       $i \leftarrow -\text{Poll}(A)$ 
24:       $K \leftarrow -\text{Append}(K, i)$ 
25:       $\text{Move}(i, c + 1, S'')$ 
26:      if  $\text{Md}(S'', G) > \text{Md}(S', G)$  then
27:         $S' \leftarrow -S''$ 
28:      for  $j \in \text{Adjacents}(i)$  do
29:        if  $j \notin K \ \& \ j \notin \text{Cluster}(S'', c)$  then
30:           $A \leftarrow -\text{Remove}(A, \text{node})$ 
31:           $K \leftarrow -\text{Append}(K, \text{node})$ 
32:      if  $\text{Md}(S', G) > \text{Md}(S, G)$  then
33:         $S \leftarrow S'$ 
34:         $c \leftarrow c + 1$ 
35:  return  $S$ 

```

---

**Modularidad para Instancia de Prueba (1000 iteraciones y  $\alpha = 0.0$ ): 0.251**

---

Para probar el algoritmo se han utilizado los siguientes valores de  $\alpha$ : 0.0, 0.25, 0.5, 0.75, 1.

---

**Algorithm 4** LocalSearch algorithm
 

---

```

1: procedure LOCALSEARCH( $S, G$ )                                ▷ The Solution 'S' and Graph 'G'
2:    $i \leftarrow 0$                                               ▷ Current cluster index
3:   while  $i < \text{Size}(S)$  do
4:      $CL \leftarrow \text{CandidateList}(S, i)$ 
5:     while  $\text{Size}(CL) > 0$  do
6:        $c \leftarrow \text{Poll}(CL)$                                 ▷ Current candidate
7:        $S' \leftarrow S$ 
8:        $\text{createEmptyCluster}(S')$ 
9:        $j \leftarrow 0$                                         ▷ Target cluster
10:      while  $j < \text{Size}(S)$  do
11:        if  $j = i$  then
12:           $j \leftarrow j + 1$ 
13:          continue
14:           $\text{Move}(c, j, S')$                                 ▷ Move node  $c$  to cluster  $j$  in solution  $S'$ 
15:          if  $\text{Md}(S', G) > \text{Md}(S, G)$  then
16:             $S \leftarrow S'$ 
17:             $\text{Move}(c, i, S')$                                 ▷ Return node  $c$  to cluster  $i$  in solution  $S'$ 
18:             $j \leftarrow j + 1$ 
19:           $i \leftarrow i + 1$ 
20:   return  $S$ 

```

---

**Modularidad para Instancia de Prueba (1000 iteraciones y  $\alpha = 0.25$ ): 0.3295**

---

Observamos que para  $\alpha = 0.25$  obtenemos un resultado notablemente mejor que con el anterior algoritmo y consideramos *ConstDividerGreedy + Local Search* la última versión de nuestro algoritmo.

#### 4.2.5. Paralelización

Una vez obtenida una modularidad aceptable, terminaremos el desarrollo procurando optimizar el uso de los recursos de la máquina dónde se ejecute el algoritmo. Se hará uso de la librería estándar de java para el manejo de hilos, *java.util.concurrent*, concretamente de los Ejecutores (*Executors*), que nos permiten gestionar de forma sencilla y óptima la casuística de la programación concurrente. Las tareas a paralelizar de nuestro programa serán las iteraciones del algoritmo.

Para demostrar la eficacia de usar paralelización en este algoritmo, usaremos de ejemplo la ejecución del algoritmo más lento: *ConstDividerGreedy + LocalSearch*. La siguiente prueba

se ha realizado en una máquina Ubuntu 16.04 con 16GB de RAM y 4 procesadores (sin *hyper-threading*).

Tiempos de ejecución (en milisegundos) para 1000 iteraciones:

- Secuencial: 446.187 ms
- Paralelo (4 hilos): 132.753 ms

Podemos observar que se reduce de forma considerable (hasta un 70 %) el tiempo de ejecución, atendiendo a nuestro requisito no funcional de ofrecer el menor tiempo de ejecución posible.



# Capítulo 5

## Estudio comparativo

El estudio comparativo que se llevará a cabo a continuación constará de dos secciones:

- La primera usará un conjunto de grafos de '*entrenamiento*' que servirá para comparar los actuales algoritmos junto a los posibles valores de sus parámetros en términos de modularidad.
- La segunda usará un conjunto de grafos de *test* (distintos al primero) que servirán para comparar la mejor versión de nuestro algoritmo con otras soluciones pre-existentes en términos de conductancia y cobertura.

## **Capítulo 6**

# **Conclusiones del proyecto y trabajos futuros**

Prueba cita [1].

# Bibliografía

- [1] J. Sánchez-Oro and A. Duarte. Iterated greedy algorithm for performing community detection in social networks. 2018.