

TEORÍA DE ALGORITMOS  
(75.29) CURSO BUCHWALD - GENENDER

# Trabajo práctico 3: Comunidades NP-Completas

6 de julio de 2025

## Integrantes

Alumno	Padrón
Alan Valdevenito	107585
Ignacio Zanoni Mutti	110884
Francisco Gutierrez	103543

## Índice

<b>1. Demostración</b>	<b>4</b>
1.1. Clustering por bajo diámetro se encuentra en NP . . . . .	5
1.2. Clustering por bajo diámetro es un problema NP-Completo . . . . .	6
1.2.1. Reduccion planteada . . . . .	6
1.2.2. Si hay K-Coloreo, entonces hay Clustering por bajo diametro . . . . .	7
1.2.3. Si hay Clustering por bajo diametro, entonces hay K-Coloreo . . . . .	7
1.2.4. Conclusion . . . . .	7
1.2.5. Ejemplo . . . . .	8
<b>2. Algoritmo de Backtracking</b>	<b>10</b>
2.1. Codigo . . . . .	10
2.2. Podas . . . . .	12
2.3. Resultados de ejecucion y mediciones de tiempo . . . . .	12
2.3.1. Resumen . . . . .	12
2.3.2. Archivo 10.3.txt con a lo sumo 2 clusters . . . . .	13
2.3.3. Archivo 10.3.txt con a lo sumo 5 clusters . . . . .	14
2.3.4. Archivo 22.3.txt con a lo sumo 3 clusters . . . . .	14
2.3.5. Archivo 22.3.txt con a lo sumo 4 clusters . . . . .	14
2.3.6. Archivo 22.3.txt con a lo sumo 10 clusters . . . . .	14
2.3.7. Archivo 22.5.txt con a lo sumo 2 clusters . . . . .	15
2.3.8. Archivo 22.5.txt con a lo sumo 7 clusters . . . . .	15
2.3.9. Archivo 30.3.txt con a lo sumo 2 clusters . . . . .	15
2.3.10. Archivo 30.3.txt con a lo sumo 6 clusters . . . . .	16
2.3.11. Archivo 30.5.txt con a lo sumo 5 clusters . . . . .	16
2.3.12. Archivo 40.5.txt con a lo sumo 3 clusters . . . . .	16
2.3.13. Archivo 45.3.txt . . . . .	17
2.3.14. Archivo 50.3.txt con a lo sumo 3 clusters . . . . .	18
2.4. Conjunto de datos generado . . . . .	19
<b>3. Algoritmo de Programación Lineal</b>	<b>22</b>
3.1. Codigo . . . . .	22
3.2. El problema . . . . .	23
3.2.1. Definición de variables . . . . .	23
3.2.2. Restricciones . . . . .	23
3.2.3. Número total de Restricciones . . . . .	24
3.3. Mediciones de tiempo y comparación con Backtracking . . . . .	24
<b>4. Algoritmo de Aproximacion</b>	<b>25</b>
4.1. Codigo . . . . .	25
4.2. Tiempos de ejecución y resultados . . . . .	26
4.3. Cota de la aproximación . . . . .	27

---

4.4. Complejidad . . . . .	27
4.5. Conjunto de datos generado . . . . .	28
<b>5. Conclusiones</b>	<b>29</b>
<b>6. Anexo</b>	<b>30</b>
6.1. Correccion 1: Demostracion . . . . .	30
6.1.1. Clustering por bajo diámetro se encuentra en NP . . . . .	30
6.2. Correccion 3: Demostracion . . . . .	33
6.2.1. Clustering por bajo diámetro es un problema NP-Completo . . . . .	33
6.3. Correccion 4: Algoritmo de Backtracking . . . . .	36
6.3.1. Grafico de tiempo de ejecucion . . . . .	36
6.3.2. Graficos de variabilidad del valor K . . . . .	37
6.4. Correccion 5: Algoritmo de Programacion Lineal . . . . .	39
6.4.1. Corrección a . . . . .	39
6.4.2. Corrección b . . . . .	39
6.4.3. Corrección c . . . . .	40
6.5. Correccion 6: Algoritmo de Aproximacion . . . . .	41
6.5.1. Comparación entre el óptimo y la aproximación . . . . .	41
6.5.2. Gráfico de tiempo de ejecución . . . . .	42
6.5.3. Gráficos de variabilidad del valor K . . . . .	42
6.5.4. Gráficos de variabilidad del valor K en grafos lista . . . . .	46
6.6. Comparacion . . . . .	48

## 1. Demostración

Problema de **Clustering por bajo diámetro**: Dado un grafo no dirigido y no pesado, un número entero  $k$  y un valor  $C$ , ¿es posible separar los vértices en a lo sumo  $k$  grupos/clusters disjuntos, de tal forma que todo vértice pertenezca a un cluster, y que la distancia máxima dentro de cada cluster sea a lo sumo  $C$ ? (Si un cluster queda vacío o con un único elemento, considerar la distancia máxima como 0).

Al calcular las distancias se tienen en cuenta tanto las aristas entre vértices dentro del cluster, como cualquier otra arista dentro del grafo.

Para realizar la demostración utilizaremos el problema de K-Coloreo el cual vimos en clase que es NP-Completo.

Problema de **K-Coloreo**: Dado un grafo y  $K$  colores diferentes, ¿es posible colorear los vértices con a lo sumo  $k$  colores de forma tal que los vértices adyacentes no compartan color?.

Para demostrar que el problema de **Clustering por bajo diámetro** es un problema NP-Completo debemos:

- Demostrar que el problema se encuentra en NP
- Reducir el problema de **K-Coloreo** al problema de **Clustering por bajo diámetro**:  
 $\text{K-Coloreo} \leq_P \text{Clustering por bajo diámetro}$

## 1.1. Clustering por bajo diámetro se encuentra en NP

Para que el problema se encuentre en NP, debe haber un verificador eficiente.

En otras palabras, debe haber un verificador que ejecute en tiempo polinomial.

Verificador: Recibe una instancia del problema y una solución.

```
1 # Al calcular las distancias se tienen en cuenta tanto las aristas entre vertices
  dentro del cluster, como cualquier otra arista dentro del grafo
2
3 def bfs_max_dist_en_grafo_completo(grafo, cluster, origen):
4     visitados = set()
5     visitados.add(origen)
6
7     distancias = {origen: 0}
8     cola = deque([origen])
9
10    max_dist = 0
11
12    while cola:
13        v = cola.popleft()
14
15        for w in grafo.adyacentes(v):
16            if w not in visitados:
17                visitados.add(w)
18                distancias[w] = distancias[v] + 1
19                cola.append(w)
20
21    # Solo consideramos distancias hacia otros vertices del cluster
22    for nodo in cluster:
23        if nodo != origen and nodo in distancias:
24            max_dist = max(max_dist, distancias[nodo])
25
26    return max_dist
27
28 def verificador(grafo, k, C, clusters):
29
30    # Validamos que la solución tenga a lo sumo k clusters
31    if len(clusters) > k: return False
32
33    clusters = [set(c) for c in clusters]
34
35    # Validamos que todo vertice pertenece a un cluster
36    for v in grafo.obtener_vertices():
37        pertenece = False
38
39        for c in clusters:
40            if v in c:
41                pertenece = True
42                break
43
44        if not pertenece: return False
45
46    # Validamos que la distancia máxima dentro de cada cluster sea a lo sumo C
47    for cluster in clusters:
48        for v in cluster:
49            max_dist = bfs_max_dist(grafo, cluster, v)
50
51            if max_dist > C: return False
52
53    return True
```

Complejidad:  $O(k \times W \times (V + E))$  con  $k$  la cantidad de clusters,  $W$  la cantidad máxima de vértices dentro de un cluster,  $V$  la cantidad de vértices y  $E$  la cantidad de aristas

- Validar que la solución tenga a lo sumo  $k$  clusters tiene una complejidad de  $O(1)$
- Validar que todo vértice pertenece a un cluster tiene una complejidad de  $O(V \times k)$  con  $V$  la cantidad de vértices y  $k$  la cantidad de clusters
- Validar que la distancia máxima dentro de cada cluster sea a lo sumo  $C$ : Para cada cluster usamos un recorrido BFS desde cada vértice dentro del cluster. El recorrido BFS tiene una complejidad de  $O(V + E)$  con  $V$  la cantidad de vértices y  $E$  la cantidad de aristas. Luego, como se ejecuta un BFS por cada cluster y por cada vértice dentro del cluster tiene una complejidad de  $O(k \times W \times (V + E))$  con  $k$  la cantidad de clusters,  $W$  la cantidad máxima de vértices dentro de un cluster,  $V$  la cantidad de vértices y  $E$  la cantidad de aristas.

Luego, se ejecuta en tiempo polinomial, lo cual quiere decir que el problema se encuentra en NP.

## 1.2. Clustering por bajo diámetro es un problema NP-Completo

Debemos demostrar que  $\mathbf{K-Coloreo} \leq_P \mathbf{Clustering por bajo diámetro}$

### 1.2.1. Reduccion planteada

¿Podemos resolver el problema de **K-Coloreo** utilizando la solución del problema de **Clustering por bajo diámetro**?

Vamos a utilizar una caja negra que resuelve el problema de **Clustering por bajo diámetro** para resolver el problema de **K-Coloreo**.

El problema de **K-Coloreo** recibe un grafo  $G'$  y un valor  $k'$  que representa la cantidad de colores.

Transformación del problema: Transformamos la entrada del problema de **K-Coloreo** en una entrada del problema **Clustering por bajo diámetro**

1. El grafo  $G'$  coincide con el grafo  $G$ .
2. El valor de  $k'$  coincide con el valor de  $k$ .
3. Definimos el valor de  $C$  en cero. Es decir que el diámetro es cero.

Esta transformación requiere una cantidad de pasos polinomiales.

A continuación, para demostrar que la reducción es correcta, debemos demostrar que:

Hay solución de **K-Coloreo** con a lo sumo  $k'$  colores, si y solo si, hay solución para el problema de **Clustering por bajo diámetro** con a lo sumo  $k$  clusters.

### 1.2.2. Si hay K-Coloreo, entonces hay Clustering por bajo diametro

Hipotesis: Hay solucion de **K-Coloreo** con a lo sumo  $k'$  colores.

Tesis: Hay solucion para el problema de **Clustering por bajo diametro** con a lo sumo  $k$  clusters.

Dada nuestra hipotesis, existen  $k'$  grupos de vertices coloreados con distinto color.

Como no hay aristas entre vertices del mismo color, el subgrafo resultante es un conjunto de vertices aislados con lo cual tienen distancia maxima 0.

Luego, cada cluster tiene diametro 0 y hay a lo sumo  $k$  clusters.

### 1.2.3. Si hay Clustering por bajo diametro, entonces hay K-Coloreo

Hipotesis: Hay solucion para el problema de **Clustering por bajo diametro** con a lo sumo  $k$  clusters.

Tesis: Hay solucion de **K-Coloreo** con a lo sumo  $k'$  colores.

Dada nuestra hipotesis, existen  $k$  clusters disjuntos de tal forma que todo vertice pertenece a un cluster y la distancia maxima dentro de cada cluster es 0.

Luego, que la distancia maxima dentro de cada cluster sea cero quiere decir que no hay aristas entre vertices del mismo cluster. Esto es porque si estuvieran conectados por una arista si distancia seria mayor o igual a 1.

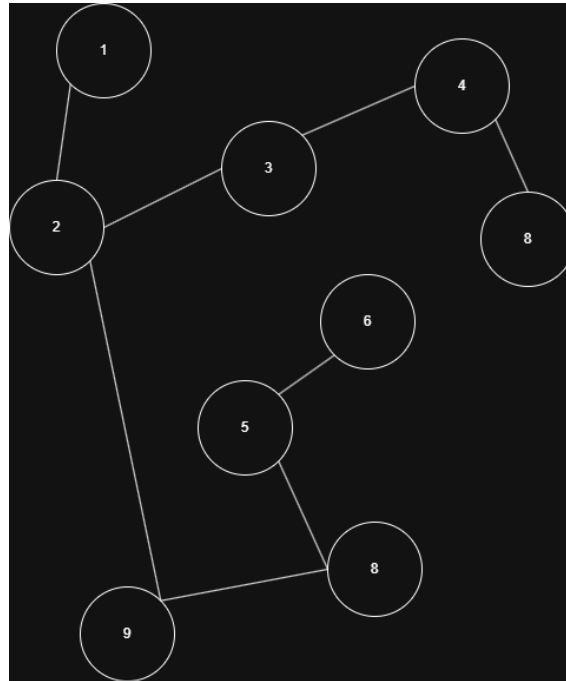
Por lo tanto podemos asignar un color distinto a cada cluster, de tal forma que es posible colorear los vertices con a lo sumo  $k'$  colores de tal forma que los vertices adyacentes no compartan color.

### 1.2.4. Conclusion

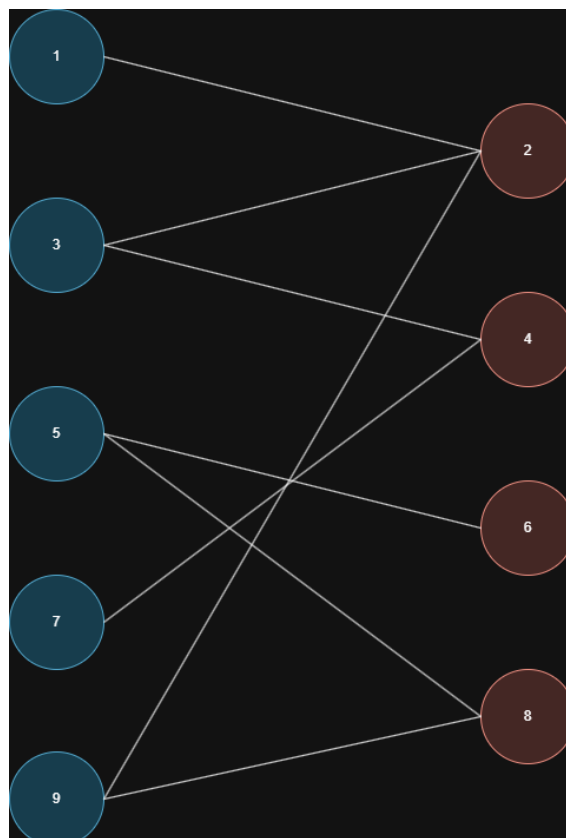
Habiendo demostrado que la reduccion es correcta, queda demostrado tambien que el problema de **Clustering por bajo diametro** es un problema NP-Completo.

### 1.2.5. Ejemplo

Consideremos el siguiente grafo  $G$ .

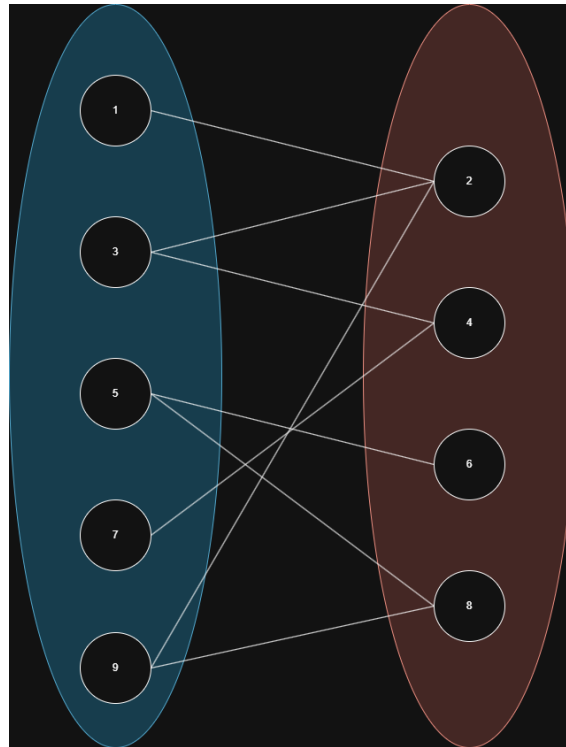


Si tenemos en cuenta el problema de **K-Coloreo** podemos ver que es posible colorear los vertices con a lo sumo  $k = 2$  colores de forma tal que los adyacentes no compartan color.





Si tenemos en cuenta el problema de **Clustering por bajo diametro** y consideramos  $C = 0$  podemos ver que es posible separar los vertices en a lo sumo  $k = 2$  clusters de tal forma que la distancia maxima dentro de cada cluster es a lo sumo  $C = 0$ . Esto ultimo quiere decir que no hay aristas dentro del grafo que unan a los vertices de un mismo cluster.



## 2. Algoritmo de Backtracking

### 2.1. Código

```
1 def bfs_distancias(grafo, origen):
2     distancias = {origen: 0}
3     visitados = set([origen])
4     cola = deque([origen])
5
6     while cola:
7         v = cola.popleft()
8         for w in grafo.adyacentes(v):
9             if w not in visitados:
10                 visitados.add(w)
11                 distancias[w] = distancias[v] + 1
12                 cola.append(w)
13     return distancias
14
15 def es_valido_extender(cluster, v, distancias_globales, mejor_resultado,
16     distancia_actual_cluster, aristas_usadas):
17     nueva_distancia = distancia_actual_cluster
18     nuevas_aristas = []
19
20     for w in cluster:
21         a1 = (v, w)
22         a2 = (w, v)
23
24         if a1 in aristas_usadas or a2 in aristas_usadas:
25             continue
26
27         d = distancias_globales[v].get(w, math.inf)
28         nueva_distancia = max(nueva_distancia, d)
29
30         # Poda: No extender si el nuevo cluster supera el mejor resultado conocido
31         if nueva_distancia >= mejor_resultado:
32             return False, distancia_actual_cluster, []
33
34         nuevas_aristas.append(a1)
35         nuevas_aristas.append(a2)
36
37     return True, nueva_distancia, nuevas_aristas
38
39 def clustering_bajo_diametro(grafo, k):
40     vertices = grafo.obtener_vertices()
41
42     # Poda: Ordenar los vertices por grado decreciente
43     vertices.sort(key=lambda v: len(grafo.adyacentes(v)), reverse=True)
44
45     clusters = [[] for _ in range(k)]
46     distancias_globales = {}
47     for v in vertices:
48         distancias_globales[v] = bfs_distancias(grafo, v)
49
50     distancias_cluster = [0 for _ in range(k)]
51
52     mejor_diametro, mejor_asignacion = backtracking(
53         grafo, vertices, k, clusters, 0, math.inf, None,
54         distancias_globales, 0, distancias_cluster, set()
55     )
56
57     return mejor_asignacion, mejor_diametro
58
59
60
61
62
63
64
```

```
65 def backtracking(grafo, vertices, k, clusters, indice, mejor_resultado,
66                 mejor_asignacion,
67                 distancias_globales, clusters_usados, distancias_cluster, aristas_usadas):
68     if indice == len(vertices):
69         peor_cluster = max(distancias_cluster[:clusters_usados])
70         if peor_cluster < mejor_resultado:
71             return peor_cluster, [list(cluster) for cluster in clusters[:
72                                     clusters_usados]]
73         return mejor_resultado, mejor_asignacion
74     actual = vertices[indice]
75
76     for i in range(clusters_usados + 1):
77
78         # Poda: No crear mas de k clusters
79         if i >= k:
80             continue
81
82         # Poda: Evitar crear clusters vac os innecesariamente
83         if len(clusters[i]) == 0 and i != clusters_usados:
84             continue
85
86         valido, nueva_distancia, nuevas_aristas = es_valido_extender(
87             clusters[i], actual, distancias_globales, mejor_resultado,
88             distancias_cluster[i], aristas_usadas
89         )
90         if not valido:
91             continue
92
93         clusters[i].append(actual)
94         vieja_distancia = distancias_cluster[i]
95         distancias_cluster[i] = nueva_distancia
96         for arista in nuevas_aristas:
97             aristas_usadas.add(arista)
98
99         nuevo_clusters_usados = clusters_usados
100         if i == clusters_usados:
101             nuevo_clusters_usados += 1
102
103         # Poda: No continuar si la peor distancia de los clusters actuales ya
104         # supera el mejor resultado
105         if max(distancias_cluster[:clusters_usados + 1]) >= mejor_resultado:
106             clusters[i].pop()
107             distancias_cluster[i] = vieja_distancia
108             for arista in nuevas_aristas:
109                 aristas_usadas.remove(arista)
110             continue
111
112         mejor_resultado, mejor_asignacion = backtracking(
113             grafo, vertices, k, clusters, indice + 1,
114             mejor_resultado, mejor_asignacion,
115             distancias_globales, nuevo_clusters_usados, distancias_cluster,
116             aristas_usadas
117         )
118
119         clusters[i].pop()
120         distancias_cluster[i] = vieja_distancia
121         for arista in nuevas_aristas:
122             aristas_usadas.remove(arista)
123
124     return mejor_resultado, mejor_asignacion
```

## 2.2. Podas

1. Poda por límite de clusters ( $i \geq k$ ): Su objetivo es evitar expandir mas clusetrs de los permitidos ( $k$ ). Si se permitiera crear mas de  $k$  clusters, no se cumpliría la restricción del problema con lo cual esta poda garantiza que no se explore ninguna asignación que rompa esta condición de tener a lo sumo  $k$  clusters.
2. Poda para evitar clusters vacíos innecesarios: Su objetivo es evitar crear nuevos clusters vacíos en posiciones que no corresponden al siguiente cluster nuevo.
3. No extender si la nueva distancia supera el mejor resultado: Cuando se intenta agregar un nuevo nodo a un cluster, se evalúa si eso haría que el diámetro del cluster supere el mejor diámetro encontrado hasta el momento. Si ya sabemos que al agregar un nodo el diámetro local se vuelve peor que el mejor encontrado, no tiene sentido continuar por ese camino.
4. No continuar si el peor cluster actual ya supera el mejor resultado: Antes de continuar con los llamados recursivos, se evalúa el peor diámetro entre los clusters ya formados. Si ese valor ya es mayor o igual al mejor resultado, no vale la pena seguir.
5. Poda de ordenamiento de vértices por grado: Su objetivo es explorar primero los nodos más difíciles o que podrían generar mas conflictos si se dejan para lo último. Es decir, aquellos con más conexiones (mayor grado).

## 2.3. Resultados de ejecución y mediciones de tiempo

### 2.3.1. Resumen

Archivo	K	Tiempo de ejecución (s)
10_3.txt	2	0.0002
10_3.txt	5	0.0003
22_3.txt	3	0.0076
22_3.txt	4	0.0010
22_3.txt	10	0.0068
22_5.txt	2	0.0010
22_5.txt	7	0.0017
30_3.txt	2	0.0057
30_3.txt	6	0.0047
30_5.txt	5	0.0173
40_5.txt	3	0.0048
45_3.txt	2	0.5131
45_3.txt	3	0.0062
45_3.txt	4	0.7859
45_3.txt	5	859.5344
45_3.txt	6	-
45_3.txt	7	-
50_3.txt	3	0.0089

Cuadro 1: Resumen de tiempos de ejecución por archivo y valor de  $K$

### 2.3.2. Archivo 10\_3.txt con a lo sumo 2 clusters

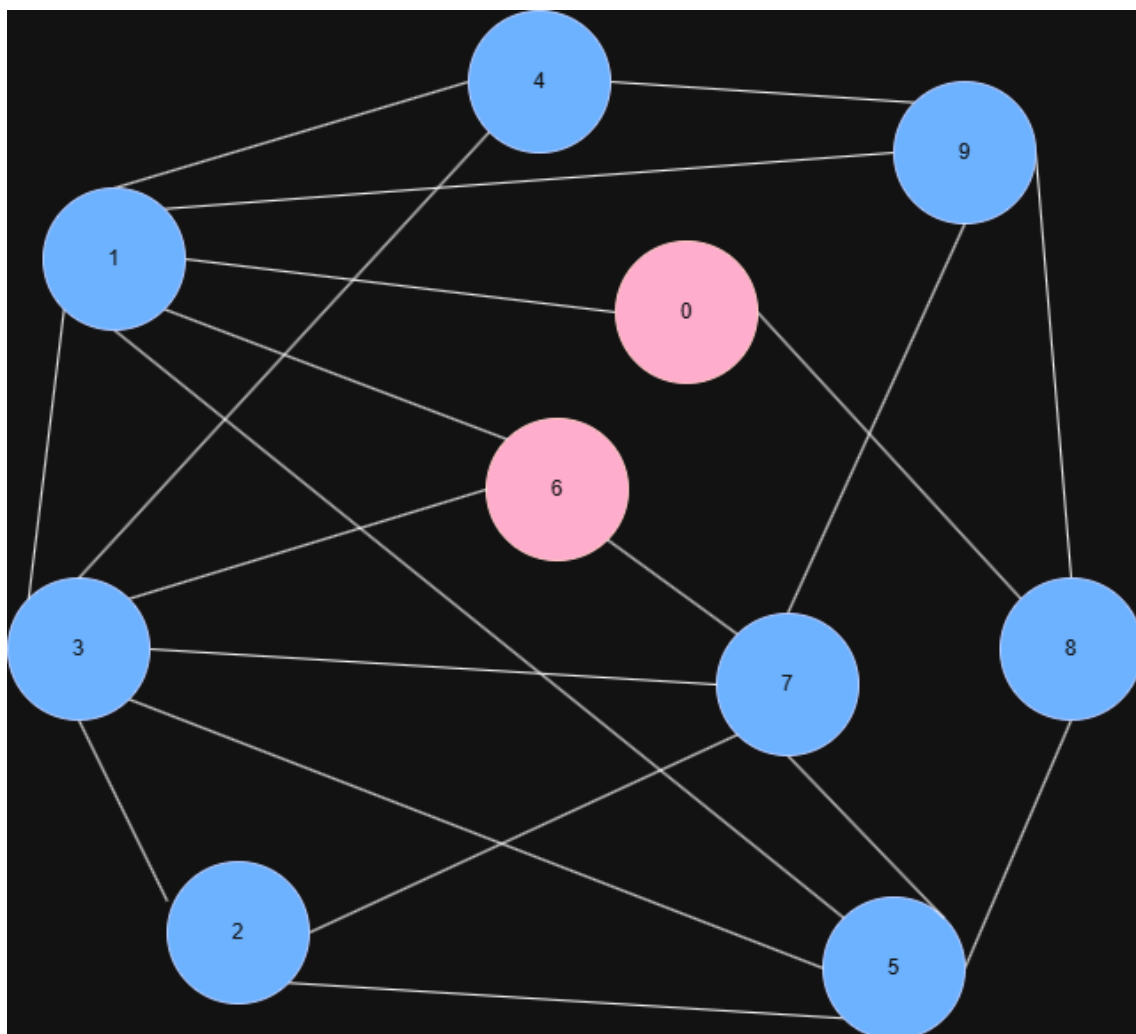
Tiempo de ejecución: 0.0002 segundos

Asignación:

1. Cluster 0 : [1, 3, 5, 7, 9, 8, 4, 2]
2. Cluster 1 : [6, 0]

Máxima distancia dentro del cluster: 2

**Grafo y clusters**



### **2.3.3. Archivo 10\_3.txt con a lo sumo 5 clusters**

Tiempo de ejecución: 0.0003 segundos

Asignación:

1. Cluster 0 : [1, 3, 5]
2. Cluster 1 : [7, 6]
3. Cluster 2 : [9, 4]
4. Cluster 3 : [8, 0]
5. Cluster 4 : [2]

Máxima distancia dentro del cluster: 1

### **2.3.4. Archivo 22\_3.txt con a lo sumo 3 clusters**

Tiempo de ejecución: 0.0076 segundos

Asignación:

1. Cluster 0 : [19, 5, 4, 3, 16, 1, 21, 2, 0, 12, 11, 8, 6]
2. Cluster 1 : [9, 17, 18, 20]
3. Cluster 2 : [7, 10, 13, 15, 14]

Máxima distancia dentro del cluster: 2

### **2.3.5. Archivo 22\_3.txt con a lo sumo 4 clusters**

Tiempo de ejecución: 0.0010 segundos

Asignación:

1. Cluster 0 : [19, 5, 4, 3, 16, 1, 21, 9, 2, 12, 11, 8, 6]
2. Cluster 1 : [7, 10, 0, 13, 14]
3. Cluster 2 : [17, 18, 20]
4. Cluster 3 : [15]

Máxima distancia dentro del cluster: 2

### **2.3.6. Archivo 22\_3.txt con a lo sumo 10 clusters**

Tiempo de ejecución: 0.0068 segundos

Asignación:

1. Cluster 0 : [19, 4, 12]
2. Cluster 1 : [5, 16, 13]
3. Cluster 2 : [3, 21]
4. Cluster 3 : [1, 8]

5. Cluster 4 : [9, 20]
6. Cluster 5 : [2, 11]
7. Cluster 6 : [7, 17, 15]
8. Cluster 7 : [10, 14]
9. Cluster 8 : [0]
10. Cluster 9 : [18, 6]

Máxima distancia dentro del cluster: 1

#### **2.3.7. Archivo 22\_5.txt con a lo sumo 2 clusters**

Tiempo de ejecución: 0.0010 segundos

Asignación:

1. Cluster 0 : [10, 11, 14, 19, 4, 21, 1, 2, 16, 18, 17, 12, 13, 8, 7, 20, 9, 0, 5, 15]
2. Cluster 1 : [6, 3]

Máxima distancia dentro del cluster: 2

#### **2.3.8. Archivo 22\_5.txt con a lo sumo 7 clusters**

Tiempo de ejecución: 0.0017 segundos

Asignación:

1. Cluster 0 : [10, 19, 21]
2. Cluster 1 : [11, 14, 1, 20]
3. Cluster 2 : [4, 18, 12, 13]
4. Cluster 3 : [2, 16, 15]
5. Cluster 4 : [17, 5, 3]
6. Cluster 5 : [8, 0]
7. Cluster 6 : [7, 9, 6]

Máxima distancia dentro del cluster: 1

#### **2.3.9. Archivo 30\_3.txt con a lo sumo 2 clusters**

Tiempo de ejecución: 0.0057 segundos

Asignación:

1. Cluster 0 : [4, 3, 13, 9, 5, 12, 26, 0, 2, 14, 23, 22, 24, 20, 16, 6, 19, 15, 21, 17, 7, 10, 28, 1, 25]
2. Cluster 1 : [18, 8, 11, 29, 27]

Máxima distancia dentro del cluster: 3

#### **2.3.10. Archivo 30\_3.txt con a lo sumo 6 clusters**

Tiempo de ejecución: 0.0047 segundos

Asignación:

1. Cluster 0 : [4, 3, 13, 9, 12, 26, 0, 23, 22, 15, 18]
2. Cluster 1 : [5, 2, 14, 20, 6, 19, 21, 1]
3. Cluster 2 : [24, 17, 7, 10, 11]
4. Cluster 3 : [16, 8, 28, 27]
5. Cluster 4 : [29]
6. Cluster 5 : [25]

Máxima distancia dentro del cluster: 2

#### **2.3.11. Archivo 30\_5.txt con a lo sumo 5 clusters**

Tiempo de ejecución: 0.0173 segundos

Asignación:

1. Cluster 0 : [5, 14, 21, 23, 24, 2, 0, 25, 4, 15, 7, 1, 17, 19, 11, 12, 16, 10, 13, 22, 26, 28, 20, 18]
2. Cluster 1 : [6, 3, 8, 9]
3. Cluster 2 : [27, 29]

Máxima distancia dentro del cluster: 2

#### **2.3.12. Archivo 40\_5.txt con a lo sumo 3 clusters**

Tiempo de ejecución: 0.0048 segundos

Asignación:

1. Cluster 0 : [0, 32, 29, 24, 36, 33, 16, 34, 25, 13, 6, 7, 22, 3, 10, 1, 2, 5, 30, 37, 38]
2. Cluster 1 : [27, 14, 28, 26, 9, 23, 17, 12, 8, 31, 35, 39]
3. Cluster 2 : [11, 21, 20, 19, 18, 15, 4]

Máxima distancia dentro del cluster: 2



### 2.3.13. Archivo 45\_3.txt

Mencionamos que particularmente para este caso de prueba, con a lo sumo 7 clusters, nuestro algoritmo no logro finalizar tras ejecutarlo por horas.

Debido a esto, probamos ejecutar este mismo caso de prueba con distinta cantidad de clusters para ver como crecia el tiempo de ejecucion.

Nuestro algoritmo finalizo su ejecucion considerando hasta a lo sumo 5 clusters.

A continuacion, presentamos los resultados obtenidos.

#### Con a lo sumo 2 clusters

Tiempo de ejecución: 0.5131 segundos

Asignación:

1. Cluster 0 : [0, 31, 2, 3, 1, 6, 26, 27, 38, 44, 17, 7, 15, 19, 40, 22, 14, 29, 30, 5, 20, 21, 37, 12, 25, 11, 42, 9, 33, 4, 13, 10, 8, 28, 41, 23, 35, 16]
2. Cluster 1 : [36, 24, 32, 39, 18, 43, 34]

Máxima distancia dentro del cluster: 3

#### Con a lo sumo 3 clusters

Tiempo de ejecución: 0.0062 segundos

Asignacion:

1. Cluster 0 : [0, 31, 2, 3, 1, 6, 26, 27, 38, 44, 17, 7, 15, 19, 40, 22, 14, 29, 30, 5, 20, 21, 37, 12, 36, 25, 11, 24, 42, 9, 33, 4, 13, 10, 8, 28, 41, 35]
2. Cluster 1 : [32, 23, 39, 43, 34]
3. Cluster 2 : [18, 16]

Máxima distancia dentro del cluster: 3

#### Con a lo sumo 4 clusters

Tiempo de ejecución: 0.7859 segundos

Asignación:

1. Cluster 0 : [0, 31, 2, 3, 1, 6, 26, 27, 38, 44, 17, 7, 15, 19, 40, 22, 14, 29, 30, 5, 20, 21, 37, 12, 36, 25, 11, 24, 42, 9, 33, 4, 13, 10, 8, 28, 41, 35]
2. Cluster 1 : [32, 23, 39, 43, 34]
3. Cluster 2 : [18, 16]

Máxima distancia dentro del cluster: 3

### Con a lo sumo 5 clusters

Tiempo de ejecución: 859.5344 segundos

Asignación:

1. Cluster 0 : [0, 31, 2, 3, 1, 6, 26, 27, 38, 44, 17, 7, 15, 19, 40, 22, 14, 29, 30, 5, 20, 21, 37, 12, 36, 25, 11, 24, 42, 9, 33, 4, 13, 10, 8, 28, 41, 35]
2. Cluster 1 : [32, 23, 39, 43, 34]
3. Cluster 2 : [18, 16]

Máxima distancia dentro del cluster: 3

### 2.3.14. Archivo 50\_3.txt con a lo sumo 3 clusters

Tiempo de ejecución: 0.0089 segundos

Asignación:

1. Cluster 0 : [31, 11, 18, 48, 3, 46, 35, 47, 14, 0, 42, 36, 10, 40, 7, 39, 28, 20, 2, 21, 8, 26, 38, 27, 45, 32, 5, 13, 34, 4, 23, 6, 44, 16, 1]
2. Cluster 1 : [9, 15, 29, 22, 12, 24, 19, 43, 41, 25]
3. Cluster 2 : [33, 37, 17, 49, 30]

Máxima distancia dentro del cluster: 3

## 2.4. Conjunto de datos generado

Archivo usado	Valor de $k$	Tiempo (s)	Distancia máxima dentro del cluster
10_1.txt	1	0.0002	9
10_1.txt	2	0.0018	4
10_1.txt	3	0.0005	3
10_1.txt	4	0.0003	2
10_1.txt	5	0.0004	1
10_1.txt	6	0.0003	1
10_1.txt	7	0.0014	1
10_1.txt	8	0.0003	1
10_1.txt	9	0.0002	1
10_1.txt	10	0.0002	0

Cuadro 2: Comparativa de resultados del algoritmo de Backtracking con distintas particiones  $k$  sobre un grafo de 10 vertices con forma de lista

- A medida que  $k$  aumenta, la distancia máxima disminuye, lo cual tiene sentido ya que al dividir más, los grupos son más pequeños.
- A partir de  $k = 5$ , la distancia máxima ya es 1 o menos.
- A  $k = 10$ , cada vértice está en su propio cluster, por lo que la distancia máxima es 0.
- El tiempo de ejecución no siempre aumenta con  $k$ , pero hay pequeñas variaciones por la cantidad de combinaciones exploradas por nuestro algoritmo .

Archivo usado	Valor de $k$	Tiempo (s)	Distancia máxima dentro del cluster
100_1.txt	1	0.0145	99
100_1.txt	50	-	-
100_1.txt	60	0.1304	1
100_1.txt	70	0.1001	1
100_1.txt	80	0.0765	1
100_1.txt	90	0.0549	1
100_1.txt	100	0.0206	0

Cuadro 3: Comparativa de resultados del algoritmo de Backtracking con distintas particiones  $k$  sobre un grafo de 100 vertices con forma de lista

- En  $k = 1$ , todo el grafo entra en un único cluster. La distancia máxima entre extremos (vertices 0 y 99) es 99.
- A partir de  $k = 60$ , el algoritmo logra una partición donde todos los clusters tienen distancia máxima menor o igual a 1.
- Para  $k = 50$ , el algoritmo no finalizó su ejecución luego de varias horas, lo que indica un punto crítico de complejidad donde las podas no son suficientes para evitar la explosión combinatoria.
- En  $k = 100$ , cada vértice está en su propio cluster, por lo tanto, la distancia máxima es 0. Esto refleja la propiedad esperada sobre que a mayor cantidad de clusters, menor será la distancia interna de cada uno, hasta alcanzar el caso donde cada vertice este en su propio cluster.
- El mayor tiempo (0.1304s) ocurre en  $k = 60$ , donde hay más combinaciones posibles que evaluar para cumplir la condición de diámetro menor o igual a 1. A partir de  $k = 70$ , el tiempo disminuye porque hay más libertad para dividir el grafo y se aplican más podas efectivas.

Archivo usado	Valor de $k$	Tiempo (s)	Distancia máxima dentro del cluster
500_1.txt	1	0.3008	499
500_1.txt	250	-	-
500_1.txt	300	22.4992	1
500_1.txt	350	17.7464	1
500_1.txt	400	12.0435	1
500_1.txt	450	5.8424	1
500_1.txt	490	1.6157	1
500_1.txt	500	0.8165	0

Cuadro 4: Comparativa de resultados del algoritmo de Backtracking con distintas particiones  $k$  sobre un grafo de 500 vértices con forma de lista

- En  $k = 1$ , todos los vértices se agrupan en un único cluster, dando una distancia máxima de 499.
- La distancia máxima se reduce abruptamente a 1 al alcanzar  $k = 300$ , y se mantiene constante hasta  $k = 490$ .
- Para  $k = 250$ , el algoritmo no finalizó su ejecución luego de varias horas, lo que indica un punto crítico de complejidad donde las podas no son suficientes para evitar la explosión combinatoria.
- Para  $k = 500$ , cada vértice se encuentra en su propio cluster, por lo tanto, la distancia máxima es 0.
- El tiempo de ejecución disminuye a medida que  $k$  se aproxima a 500, ya que el espacio de búsqueda del algoritmo de backtracking se reduce.
- Este comportamiento refleja cómo el algoritmo explora menos combinaciones cuando los clusters se acercan al caso trivial (cada nodo en su propio grupo).

## 3. Algoritmo de Programación Lineal

### 3.1. Código

```
1 def bfs_distancias(grafo, origen):
2     distancias = {origen: 0}
3     visitados = set([origen])
4     cola = deque([origen])
5
6     while cola:
7         v = cola.popleft()
8         for w in grafo.adyacentes(v):
9             if w not in visitados:
10                 visitados.add(w)
11                 distancias[w] = distancias[v] + 1
12                 cola.append(w)
13     return distancias
14
15 def programacion_lineal(grafo, vertices, k, clusters, distancias_globales):
16     n = len(vertices)
17     k = min(k, n)
18
19     # Creo el problema y agrego las variables
20     problema = pulp.LpProblem("Clustering_Bajo_Diametro", pulp.LpMinimize)
21
22     x = pulp.LpVariable.dicts("x", [(v, c) for v in vertices for c in range(k)],
23                               cat="Binary")
24     y = pulp.LpVariable.dicts("y", [c for c in range(k)], cat="Binary")
25     D = pulp.LpVariable("D", lowBound=0, cat="Integer")
26
27     problema += D
28
29     # Agrego las restricciones
30     for v in vertices:
31         problema += pulp.lpSum(x[v, c] for c in range(k)) == 1, f"asignacion_{v}"
32
33     for c in range(k):
34         for v in vertices:
35             problema += x[v, c] <= y[c]
36
37     problema += pulp.lpSum(y[c] for c in range(k)) <= k, "limite_clusters"
38
39     for u in vertices:
40         for v in vertices:
41             if u < v:
42                 d_uv = distancias_globales[u].get(v, math.inf)
43                 if math.isfinite(d_uv):
44                     for c in range(k):
45                         problema += d_uv * (x[u, c] + x[v, c] - 1) <= D, f"
46                         diametro_{u}_{v}_{c}"
47
48     # Resuelvo el problema y reconstruyo la solucion
49     status = problema.solve()
50
51     for v in vertices:
52         for c in range(k):
53             if pulp.value(x[v, c]) > 0.5:
54                 clusters[c].append(v)
55                 break
56
57     return int(pulp.value(D)), clusters
58
59
60
61
62
63
```

```
64 def clustering_bajo_diametro(grafo, k):
65
66     vertices = grafo.obtener_vertices()
67     vertices.sort(key=lambda v: len(grafo.adyacentes(v)), reverse=True)
68
69     clusters = [[] for _ in range(k)]
70     distancias_globales = {}
71     for v in vertices:
72         distancias_globales[v] = bfs_distancias(grafo, v)
73
74     return grafo, vertices, k, clusters, distancias_globales
75
76 def clustering_bajo_diametro_pl(grafo, vertices, k, clusters, distancias_globales):
77     mejor_diametro, mejor_asignacion = programacion_lineal(grafo, vertices[:], k,
78         clusters[:], distancias_globales)
79
80     return mejor_asignacion, mejor_diametro
81
82 def main():
83     args = parsear_argumentos()
84
85     grafo = crear_grafo(args.archivo)
86
87     grafo, vertices, k, clusters, distancias_globales = clustering_bajo_diametro(
88         grafo, args.clusters)
89     return clustering_bajo_diametro_pl(grafo, vertices, k, clusters,
90         distancias_globales)
```

## 3.2. El problema

### 3.2.1. Definición de variables

Las variables utilizadas para la implementación del problema son 3:

- $X_{u,c}$ . La cual es de tipo binario y representa si el vertice  $v$  se encuentra dentro del cluster  $c$  (1 si  $v$  se encuentra en  $c$ , 0 si no).
- $Y_c$ . También de tipo binario y representa si el cluster  $c$  contiene elementos (1 si el cluster contiene elementos 0 si está vacío).
- $D$ . Por último la variable  $D$ , que es de tipo entero y tiene un valor mínimo de 0, la cual representa el diametro máximo de nuestra clauseterización.

Para la definición del problema se hizo foco justamente en esta última variable 'D'. La cual buscamos optimizar minimizándola.

### 3.2.2. Restricciones

- Asignación única de vértices, cada vértice debe pertenecer exactamente a un cluster:

$$\sum_{i=1}^k X_{i,c} = 1 \quad \forall i \in V$$

- Límite de clusters, el número de clusters no puede ser superior a  $k$ :

$$\sum_{c=1}^k y_c = 1$$

- Un cluster está activo si al menos un vértice está asignado a él:

$$x_{i,c} \leq y_c \quad \forall i \in V, \forall c \in \{1, \dots, k\}$$

- Límite de diámetro, Para cada par de vertices (i, j) asignados al mismo cluster, la distancia entre ellos debe ser menor o igual al diámetro máximo

$$d_{i,j} \cdot (x_{u,c} + x_{v,c} - 1) \leq D \quad \forall u, v \in V, \forall c \in \{1, \dots, k\}$$

### 3.2.3. Número total de Restricciones

Si analizamos el último item podemos ver que se agrega una restricción para cada par de vertices u, v en cada cluster. Por ende el número total de restricciones es equivalente a  $O(n * k)$  en notación Big O

### 3.3. Mediciones de tiempo y comparación con Backtracking

Archivo	K	T.ejecución (s)	diferencia contra backtracking
10_3.txt	2	0.0282	+0.0280
10_3.txt	5	0.0911	+0.0908
22_3.txt	3	0.3671	+0.3595
22_3.txt	4	0.5758	+0.5748
22_3.txt	10	-	-
22_5.txt	2	0.1324	+0.1314
22_5.txt	7	-	-
30_3.txt	2	1.4357	+1.4300
30_3.txt	6	0.0047	=
30_5.txt	5	1.9123	+1.895
40_5.txt	3	4.4245	+4.4197
45_3.txt	2	0.5159	+0.0028
45_3.txt	3	1.6598	+1.6918
45_3.txt	4	5.1351	+4.392
45_3.txt	5	331.7997	-527.7347
45_3.txt	6	-	-
45_3.txt	7	-	-
50_3.txt	3	2.0628	+2.0539

Cuadro 5: Comparación de tiempos de ejecución por archivo y valor de K PL vs backtracking



## 4. Algoritmo de Aproximacion

Implementación del Algoritmo de Louvain para aproximar el problema mediante la maximización de la modularización, buscando K clusters en el grafo. La modularización es un valor que se obtiene con la siguiente fórmula:

$$Q = \frac{1}{2m} \sum_i \sum_j \left( w(v_i, v_j) - \frac{k_i k_j}{2m} \right) \delta(c_i, c_j)$$

Y nos indica qué tan dividido está el grafo en comunidades o clusters.

### 4.1. Codigo

```
1 def algoritmo_louvain(grafo, K):
2     comunidades = {v: v for v in grafo.obtener_vertices()}
3     i = 0
4     while True:
5         i += 1
6         comunidades = fase1(grafo, comunidades)
7         nuevo_grafo = fase2(grafo, comunidades)
8         if len(nuevo_grafo.obtener_vertices()) == len(grafo.obtener_vertices()) or
9         len(nuevo_grafo.obtener_vertices()) <= K:
10             break
11
12     nuevas_comunidades = {}
13     for v in grafo.obtener_vertices():
14         nuevas_comunidades[v] = comunidades[v]
15
16     grafo = nuevo_grafo
17     comunidades = nuevas_comunidades
18
19     return comunidades
20
21 def fase1(grafo, comunidades):
22     cambiado = True
23     while cambiado:
24         cambiado = False
25         for nodo in grafo.obtener_vertices():
26             mejor, delta = mejor_comunidad(grafo, comunidades, nodo)
27
28             if comunidades[nodo] != mejor and delta > 0:
29                 comunidades[nodo] = mejor
30                 cambiado = True
31     return comunidades
32
33 def fase2(grafo, comunidades):
34     nuevo_grafo = Grafo(False)
35     clusters = agrupar_por_comunidad(comunidades)
36     for c in clusters:
37         nuevo_grafo.agregar_vertice(c)s
38     for v, w in obtener_aristas(grafo):
39         c1, c2 = comunidades[v], comunidades[w]
40         peso = grafo.peso_arista(v, w)
41
42         if not nuevo_grafo.pertenece_vertice(c1):
43             nuevo_grafo.agregar_vertice(c1)
44         if not nuevo_grafo.pertenece_vertice(c2):
45             nuevo_grafo.agregar_vertice(c2)
46         if nuevo_grafo.estan_unidos(c1, c2):
47             nuevo_peso = nuevo_grafo.peso_arista(c1, c2) + peso
48             nuevo_grafo.agregar_arista(c1, c2, nuevo_peso)
49         else:
50             nuevo_grafo.agregar_arista(c1, c2, peso)
51     return nuevo_grafo
52
53
```

```
54 def mejor_comunidad(grafo, comunidades, nodo):
55     original = comunidades[nodo]
56     vecinos = grafo.adyacentes(nodo)
57     comunidades_vecinas = set(comunidades[v] for v in vecinos)
58     mejor_delta = 0
59     mejor_com = original
60
61     for c in comunidades_vecinas:
62         comunidades[nodo] = c
63         nueva_mod = modularizacion(grafo, comunidades)
64         comunidades[nodo] = original
65         actual_mod = modularizacion(grafo, comunidades)
66         delta_Q = nueva_mod - actual_mod
67
68         if delta_Q > mejor_delta:
69             mejor_delta = delta_Q
70             mejor_com = c
71
72     return mejor_com, mejor_delta
73
74
75 def modularizacion(grafo, comunidades):
76     Q = 0
77     m = len(obtener_aristas(grafo))
78
79     for i in grafo.obtener_vertices():
80         for j in grafo.adyacentes(i):
81             if comunidades[i] == comunidades[j]:
82                 peso_ij = grafo.peso_arista(i, j)
83                 k_i = len(grafo.adyacentes(i))
84                 k_j = len(grafo.adyacentes(j))
85
86                 Q += peso_ij - (k_i * k_j) / (2 * m)
87
88     return Q / (2*m)
```

## 4.2. Tiempos de ejecución y resultados

Archivo	K	Tiempo de ejecución (s)
10_3.txt	2	0.0165
22_3.txt	3	0.2880
22_5.txt	2	0.2765
30_3.txt	2	0.4159
30_5.txt	5	1.0118
40_5.txt	3	1.0691
45_3.txt	7	2.1744
50_3.txt	3	2.1075

Cuadro 6: Resumen de tiempos de ejecución por archivo y valor de  $K$

### 4.3. Cota de la aproximación

Para establecer una cota, tomamos los resultados de cualquiera de las pruebas y debemos comparar el resultado óptimo con el aproximado con la siguiente fórmula:

$$\frac{A(I)}{z(I)} \leq r(A)$$

- Siendo  $r(A)$  el valor de la cota.
- $A(I)$  el resultado de la aproximación.
- $z(I)$  el resultado óptimo.

Tomando como instancia la prueba 10\_3.txt con  $K = 2$ .

- El resultado óptimo es una distancia máxima dentro del cluster de 2, es decir  $z(I) = 2$ .
- El resultado aproximado es  $A(I) = 3$ .
- Lo que resulta en  $(1, 5 \leq r(A))$

Por lo tanto estamos ante una 1,5 Aproximación.

### 4.4. Complejidad

La complejidad del algoritmo es  $O(E * \log(V))$ .

#### 4.5. Conjunto de datos generado

Archivo usado	Valor de $k$	Tiempo (s)	Distancia máxima dentro del cluster
10_1.txt	1	0.0030	5
10_1.txt	2	0.0030	7
10_1.txt	3	0.0023	7
10_1.txt	4	0.0040	7

Cuadro 7: Comparativa de resultados del algoritmo de Louvain con distintas particiones  $k$  sobre un grafo de 10 vertices con forma de lista

A partir de  $K = 2$ , la distancia máxima es 7 y no cambia, y encuentra 2 clusters o comunidades.

Archivo usado	Valor de $k$	Tiempo (s)	Distancia máxima dentro del cluster
100_1.txt	1	0.1971	90
100_1.txt	50	0.1996	90
100_1.txt	100	0.2305	90

Cuadro 8: Comparativa de resultados del algoritmo de Louvain con distintas particiones  $k$  sobre un grafo de 100 vertices con forma de lista

Sin importar el valor de  $K$ , la distancia máxima es 90 y encuentra 3 clusters o comunidades.

Archivo usado	Valor de $k$	Tiempo (s)	Distancia máxima dentro del cluster
500_1.txt	1	10.3451	490
500_1.txt	250	10.6549	490
500_1.txt	500	10.6815	490

Cuadro 9: Comparativa de resultados del algoritmo de Louvain con distintas particiones  $k$  sobre un grafo de 500 vertices con forma de lista

Ocurre lo mismo que en las pruebas anteriores. Sin importar el valor de  $K$ , la distancia máxima es 490 y encuentra 3 clusters o comunidades.

## 5. Conclusiones

### Algoritmo de Backtracking

El archivo de prueba `45_3.txt` representa un caso particularmente ilustrativo para estudiar el comportamiento del algoritmo de *backtracking* ante variaciones del parámetro  $K$ . En este caso, observamos una tendencia creciente muy marcada en los tiempos de ejecución:

1. Para  $K = 2$ , el algoritmo resolvió el problema en aproximadamente **0.5 segundos**.
2. Para  $K = 3$ , el algoritmo resolvió el problema en aproximadamente **0.0062 segundos**.  
Notamos que hubo una reducción en el tiempo de ejecución.
3. Para  $K = 4$ , el tiempo aumentó nuevamente a aproximadamente **0.79 segundos**.
4. Para  $K = 5$ , se produjo una explosión en el tiempo de ejecución, alcanzando más de **859 segundos** (casi 15 minutos).
5. A partir de  $K = 6$ , el algoritmo no logró finalizar su ejecución tras varias horas.

Este comportamiento pone en evidencia la sensibilidad del algoritmo ante variaciones en el valor de  $K$ , ya que pequeños incrementos pueden generar un crecimiento exponencial en el espacio de soluciones a explorar.

## 6. Anexo

### 6.1. Correccion 1: Demostracion

#### 6.1.1. Clustering por bajo diámetro se encuentra en NP

Para que el problema se encuentre en NP, debe haber un verificador eficiente.

En otras palabras, debe haber un verificador que ejecute en tiempo polinomial.

Verificador: Recibe una instancia del problema y una solucion.

```
1 # Al calcular las distancias se tienen en cuenta tanto las aristas entre vertices
   dentro del cluster, como cualquier otra arista dentro del grafo
2
3 def bfs_max_dist_en_grafo_completo(grafo, cluster, origen):
4     visitados = set()
5     visitados.add(origen)
6
7     distancias = {origen: 0}
8     cola = deque([origen])
9
10    max_dist = 0
11
12    while cola:
13        v = cola.popleft()
14
15        for w in grafo.adyacentes(v):
16            if w not in visitados:
17                visitados.add(w)
18                distancias[w] = distancias[v] + 1
19                cola.append(w)
20
21    # Solo consideramos distancias hacia otros vertices del cluster
22    for nodo in cluster:
23        if nodo != origen and nodo in distancias:
24            max_dist = max(max_dist, distancias[nodo])
25
26    return max_dist
```

```
1 def verificador(grafo, k, C, clusters):
2
3     # Validamos que la solución tenga a lo sumo k clusters
4     if len(clusters) > k: return False
5
6     clusters = [set(c) for c in clusters]
7
8     # Validamos que los clusters sean disjuntos
9     vertices_en_clusters = set()
10
11     for cluster in clusters:
12         for v in cluster:
13
14             if v in vertices_en_clusters:
15                 return False
16
17             vertices_en_clusters.add(v)
18
19     # Validamos que se cubran todos los vertices del grafo
20     if set(grafo.obtener_vertices()) != vertices_en_clusters:
21         return False
22
23     # Validamos que la distancia máxima dentro de cada cluster sea a lo sumo C
24     for cluster in clusters:
25         for v in cluster:
26             max_dist = bfs_max_dist(grafo, cluster, v)
27
28             if max_dist > C: return False
29
30     return True
```

Complejidad: La complejidad total resulta ser  $O(Vx(V + E))$  con  $V$  la cantidad de vertices y  $E$  la cantidad de aristas.

1. Validar que la solución tenga a lo sumo  $k$  clusters tiene una complejidad de  $O(1)$ .
2. Validar que los clusters sean disjuntos tiene una complejidad de  $O(V)$  con  $V$  la cantidad de vertices ya que iteramos cada cluster y por cada cluster iteramos todos sus vertices, lo cual implica que, al ser disjuntos, cada vértice se recorre exactamente una vez. Dado que utilizamos un conjunto (set), las operaciones de agregar y consultar si un vertex se encuentra en dicho conjunto son  $O(1)$ .
3. Validar que se cubran todos los vertices del grafo tiene una complejidad de  $O(V)$  con  $V$  la cantidad de vertices ya que se compara el conjunto de vertices del grafo, con el conjunto de vertices en clusters.
4. Validar que la distancia máxima dentro de cada cluster sea a lo sumo  $C$  tiene una complejidad de  $O(Vx(V + E))$  con  $V$  la cantidad de vertices y  $E$  la cantidad de aristas ya que iteramos cada cluster y por cada cluster iteramos todos sus vertices, donde para cada vertex se hace un BFS sobre el grafo completo. Recordamos que la complejidad de un recorrido BFS es  $O(V + E)$  y mencionamos que como los clusters son disjuntos, al recorrer cada cluster y recorrer sus vertices estamos recorriendo todos los vertices del grafo una única vez, con lo cual es por esto que por cada vertex del grafo hacemos un recorrido BFS.

Anteriormente mencionamos tanto a  $k$  como a  $W$  en la complejidad, donde  $W$  era la cantidad máxima de vertices en un cluster. Sin embargo, dado que los clusters son disjuntos y cubren todos los vértices del grafo, la suma total de vértices en todos los clusters es exactamente  $V$ , con lo cual al validar que la distancia máxima dentro de cada cluster sea a lo sumo  $C$  se ejecutan  $V$  recorridos BFS. Es por esto que en el nuevo análisis de la complejidad no son mencionados.

Luego, se ejecuta en tiempo polinomial, lo cual quiere decir que el problema se encuentra en NP.



## 6.2. Correccion 3: Demostracion

### 6.2.1. Clustering por bajo diámetro es un problema NP-Completo

Para demostrar que el problema de **Clustering por bajo diámetro** es NP-Completo, debemos demostrar que existe una reducción polinomial desde algún problema NP-Completo. En nuestro caso, elegimos realizar la reducción utilizando el problema de **K-Coloreo**.

Luego, debemos demostrar que **K-Coloreo**  $\leq_P$  **Clustering por bajo diámetro**

#### 6.4.1.1. Reduccion planteada

¿Podemos resolver el problema de **K-Coloreo** utilizando la solucion del problema de **Clustering por bajo diametro**?

Vamos a utilizar una caja negra que resuelve el problema de **Clustering por bajo diametro** para resolver el problema de **K-Coloreo**.

El problema de **K-Coloreo** recibe un grafo  $G$  y un valor  $k$  que representa la cantidad de colores.

Transformacion del problema: Transformamos la entrada del problema de **K-Coloreo** en una entrada del problema **Clustering por bajo diametro**

1. A partir del grafo  $G$  construimos el grafo complemento  $G'$  donde cada arista  $(u, v) \in G'$  si y solo si  $(u, v) \notin G$
2. El valor de  $k$  del problema de **K-Coloreo** coincide con el valor de  $k'$  del problema de **Clustering por bajo diametro**
3. Definimos el valor de  $C = 1$  de tal forma que la distancia maxima permitida dentro de cada cluster es 1 teniendo en cuenta que al calcular las distancias se tienen en cuenta tanto las aristas entre vertices dentro del cluster, como cualquier otra arista dentro del grafo

Esta transformacion requiere una cantidad de pasos polinomiales.

A continuacion, para demostrar que la reduccion es correcta, debemos demostrar que:

Hay solucion de **K-Coloreo** con a lo sumo  $k$  colores en  $G$ , si y solo si, hay solucion para el problema de **Clustering por bajo diametro** con a lo sumo  $k'$  clusters y  $C = 1$  en  $G'$ .

#### 6.4.1.2. Si hay K-Coloreo, entonces hay Clustering por bajo diametro

Hipotesis: Hay solucion de **K-Coloreo** con a lo sumo  $k$  colores en  $G$

Tesis: Hay solucion para el problema de **Clustering por bajo diametro** con a lo sumo  $k'$  clusters en  $G'$

Dada nuestra hipótesis, existen  $k$  grupos de vértices coloreados con distinto color, tales que no hay aristas entre vértices del mismo color en  $G$ .

Luego, en el grafo complemento  $G'$ , todos los vértices que comparten color en  $G$  están conectados entre sí, ya que no había aristas entre ellos en  $G$ . Esto significa que el conjunto de vértices con el mismo color en  $G$  forma un **clique** en  $G'$ .

Dado que en un clique toda pareja de vértices está conectada directamente, la distancia entre ellos en el grafo complemento  $G'$  es 1.

Por lo tanto, cada conjunto de un mismo color puede utilizarse como un cluster válido en el problema de **Clustering por bajo diametro**, con distancia máxima igual a 1.

#### 6.4.1.3. Si hay Clustering por bajo diametro, entonces hay K-Coloreo

Hipotesis: Hay solución para el problema de **Clustering por bajo diametro** sobre el grafo complemento  $G'$  con a lo sumo  $k'$  clusters y  $C = 1$

Tesis: Hay solucion de **K-Coloreo** con a lo sumo  $k$  colores para el grafo original  $G$

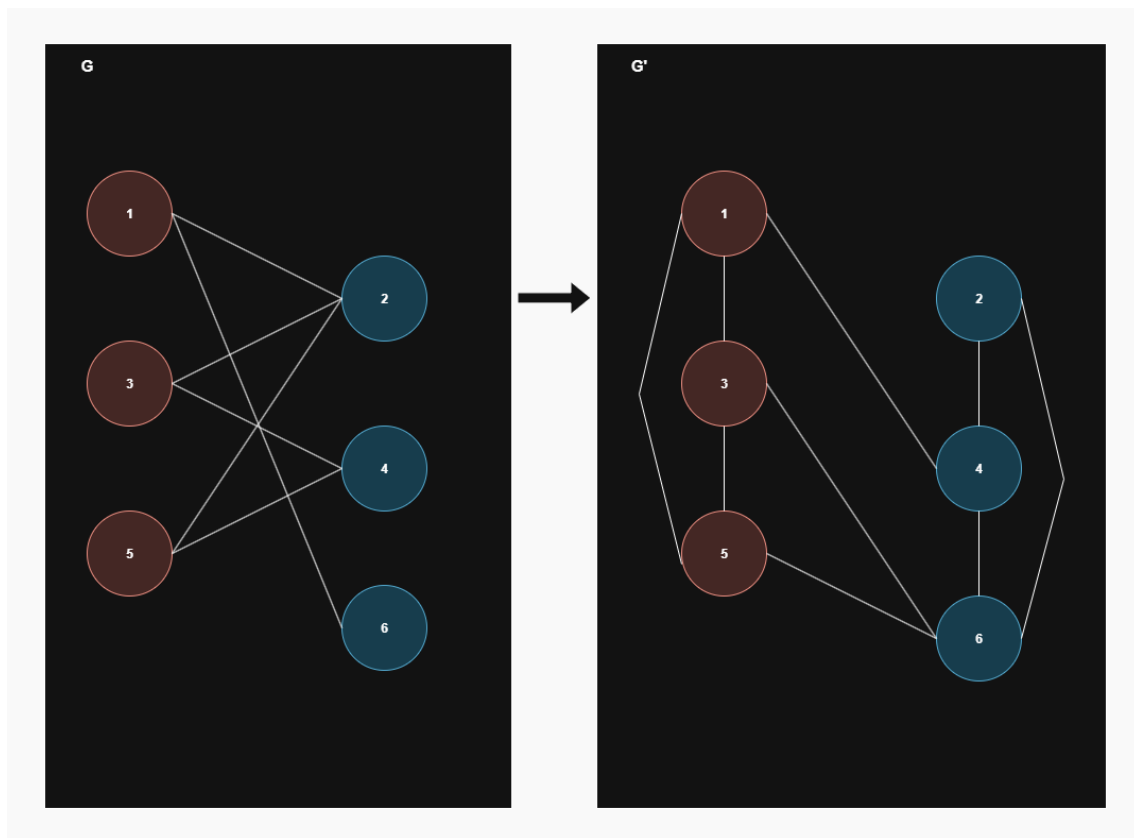
Dada nuestra hipótesis, existen  $k$  clusters disjuntos que cubren todos los vertices del grafo, donde en cada cluster la distancia maxima entre pares de vertices (medida en el grafo completo  $G'$ ) es a lo sumo 1.

Eso implica que, en cada cluster, todos los vertices estan conectados entre si en  $G'$ , es decir, forman un clique.

Luego, en el grafo original  $G$ , esos mismos vertices no estan conectados por ninguna arista.

Por lo tanto, podemos asignar un color distinto a cada cluster, y obtenemos una asignación de colores tal que no hay vertices adyacentes con el mismo color.

#### 6.4.1.4. Ejemplo



#### 6.4.1.5. Conclusion

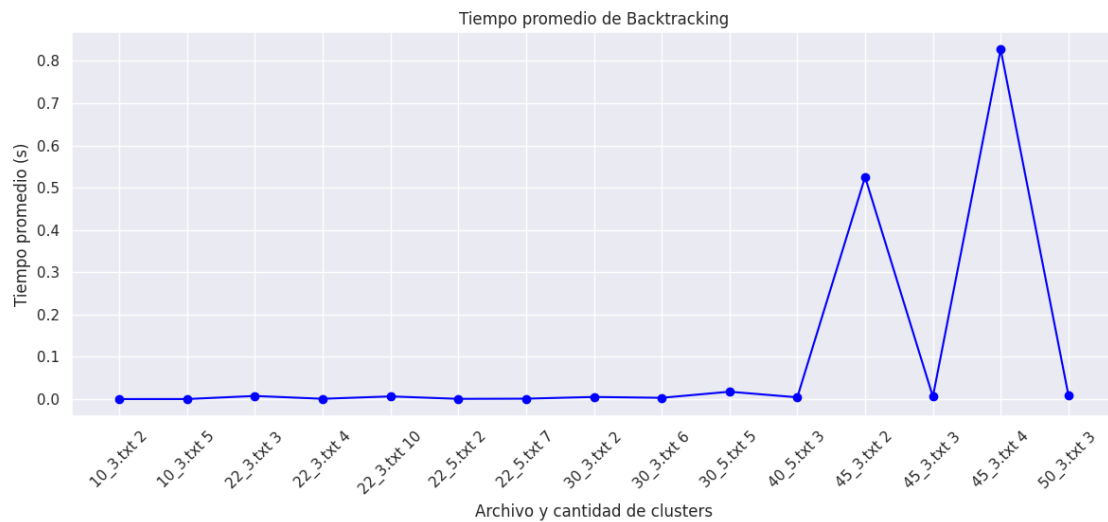
Habiendo demostrado que la reduccion es correcta, queda demostrado tambien que el problema de **Clustering por bajo diametro** es un problema NP-Completo.

## 6.3. Correccion 4: Algoritmo de Backtracking

### 6.3.1. Grafico de tiempo de ejecucion

Se realizaron mediciones de tiempo con el objetivo de comparar los tiempos de ejecucion para cada grafo y valor de  $k$  de los casos proporcionados por la catedra.

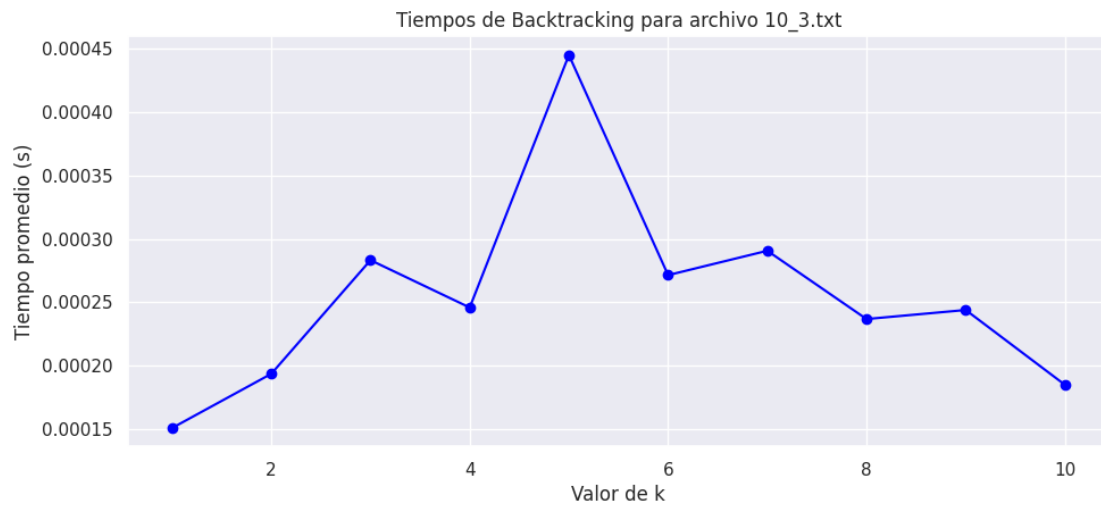
Para ello, se generaron gráficos y se utilizó el codigo proporcionado por la catedra para ejecutar el algoritmo de tal forma que tomamos más de una medición de la misma muestra y nos quedamos con el promedio para reducir el ruido en la medición. Esto utilizando el archivo `util.py` que nos proporciona la catedra y colocando la constante `RUNS_PER_SIZE` con valor 30.



Podemos ver que el caso `45_3.txt` con  $k = 2$  y el caso `45_3.txt` con  $k = 4$  son los casos que mas tiempo tardan.

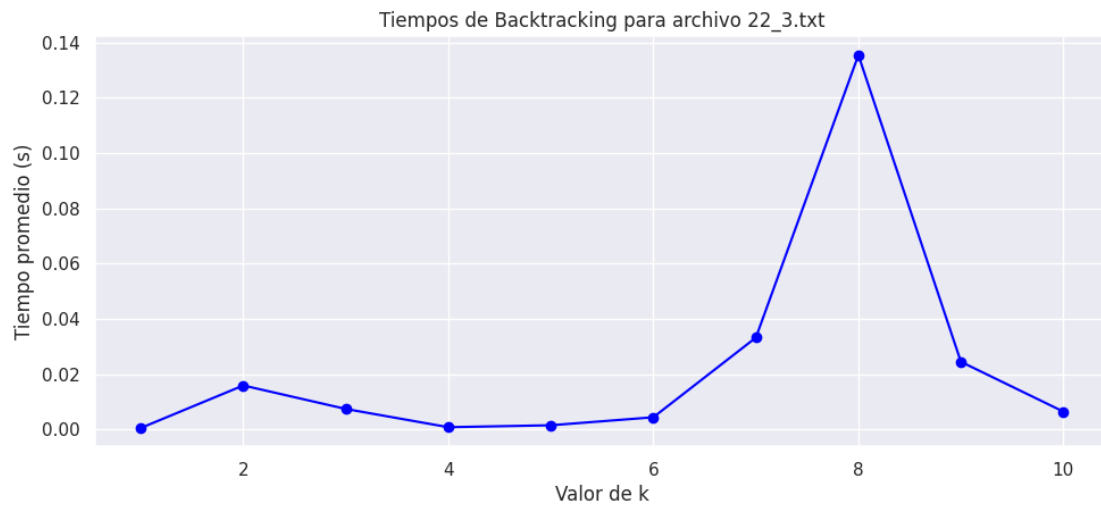
### 6.3.2. Graficos de variabilidad del valor K

A continuacion, mostramos un grafico de los tiempos de ejecucion para el caso 10\_3.txt variando el valor de k.



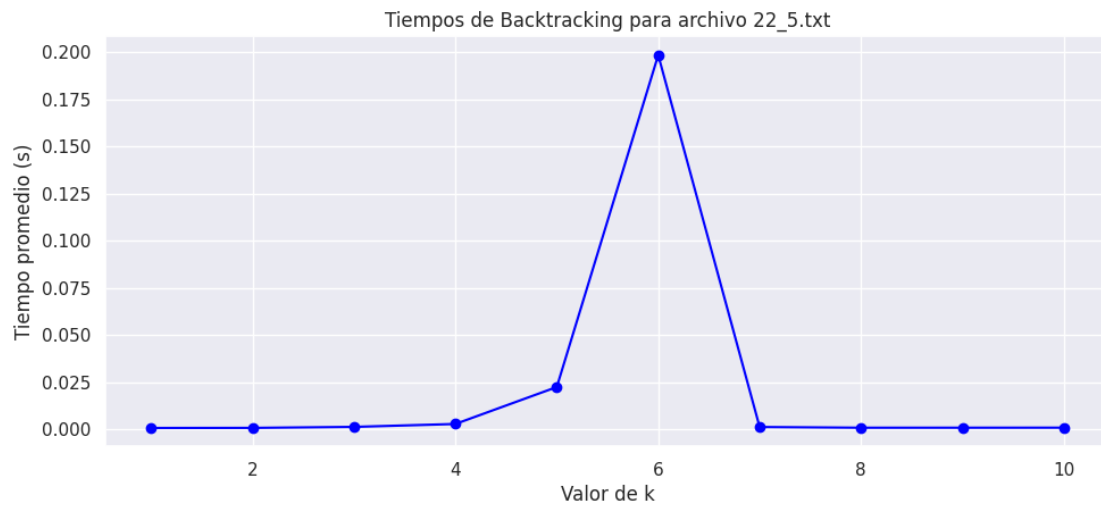
Podemos ver que para  $k = 5$  se alcanza el tiempo mas alto de ejecucion.

A continuacion, mostramos un grafico de los tiempos de ejecucion para el caso 22\_3.txt variando el valor de k.



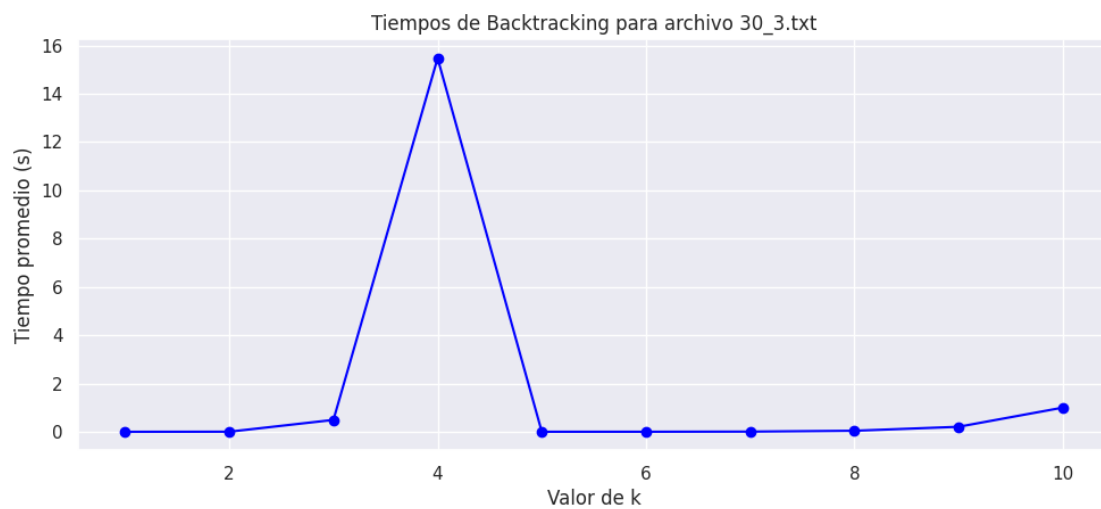
Podemos ver que para  $k = 8$  se alcanza el tiempo mas alto de ejecucion.

A continuacion, mostramos un grafico de los tiempos de ejecucion para el caso *22\_5.txt* variando el valor de  $k$ .



Podemos ver que para  $k = 6$  se alcanza el tiempo mas alto de ejecucion.

A continuacion, mostramos un grafico de los tiempos de ejecucion para el caso *30\_3.txt* variando el valor de  $k$ .



Podemos ver que para  $k = 4$  se alcanza el tiempo mas alto de ejecucion.

## 6.4. Corrección 5: Algoritmo de Programación Lineal

En PL:

- (a) Ponen variables con índices y no dicen los rangos de los índices en sí.
- (b) La inecuación  $\sum y_c = 1$  no debería ser  $\sum y_c \leq k$ ?
- (c) Nuevamente, no hay gráficos, difícil comparar con BT.

### 6.4.1. Corrección a

Las variables y sus rangos son:

- $X_{v,c}$ . Es de tipo binario y representa si el vértice  $v$  se encuentra dentro del cluster  $c$  (1 si  $v$  se encuentra en  $c$ , 0 si no). Su rango es de 1 a  $V$  (cantidad de vértices) para  $v$  y de 1 a  $k$  (cantidad de clusters) para  $c$
- $Y_c$ . También de tipo binario y representa si el cluster  $c$  contiene elementos (1 si el cluster contiene elementos 0 si está vacío). Su rango es  $[1, k]$
- $D$ . Por último la variable  $D$ , que es de tipo entero y tiene un valor mínimo de 0, la cual representa el diámetro máximo de nuestra clausurización.

### 6.4.2. Corrección b

Se corrigieron principalmente dos restricciones

- Asignación única de vértices, cada vértice debe pertenecer exactamente a un cluster:

$$\sum_{c=1}^k X_{v,c} = 1 \quad \forall v \in V$$

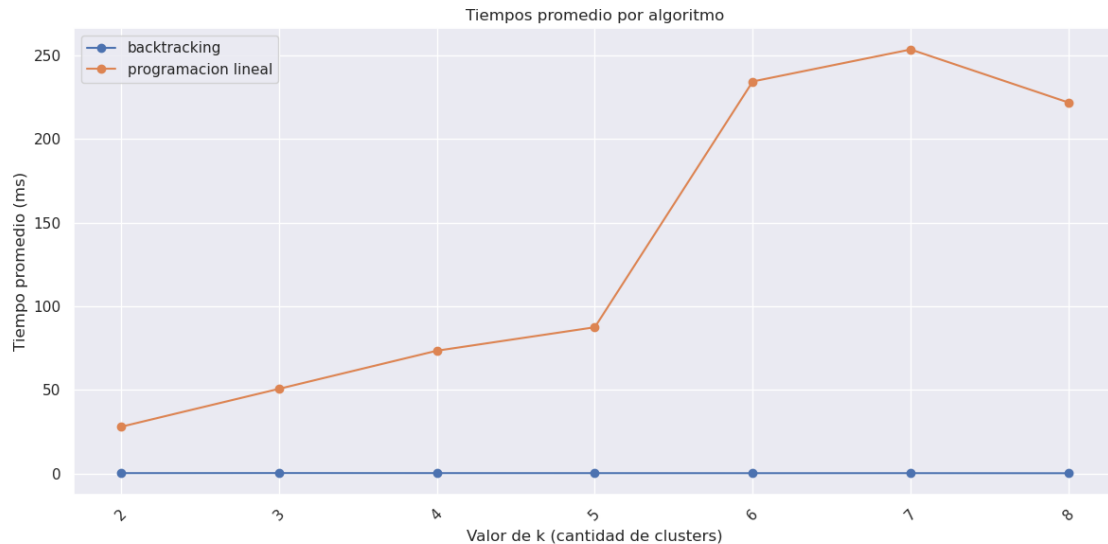
Aquí estaban mal definidos los subíndices ya que hay que iterar sobre los clusters para cada vértice

- Límite de clusters, el número de clusters 'activos' no puede ser superior a  $k$ :

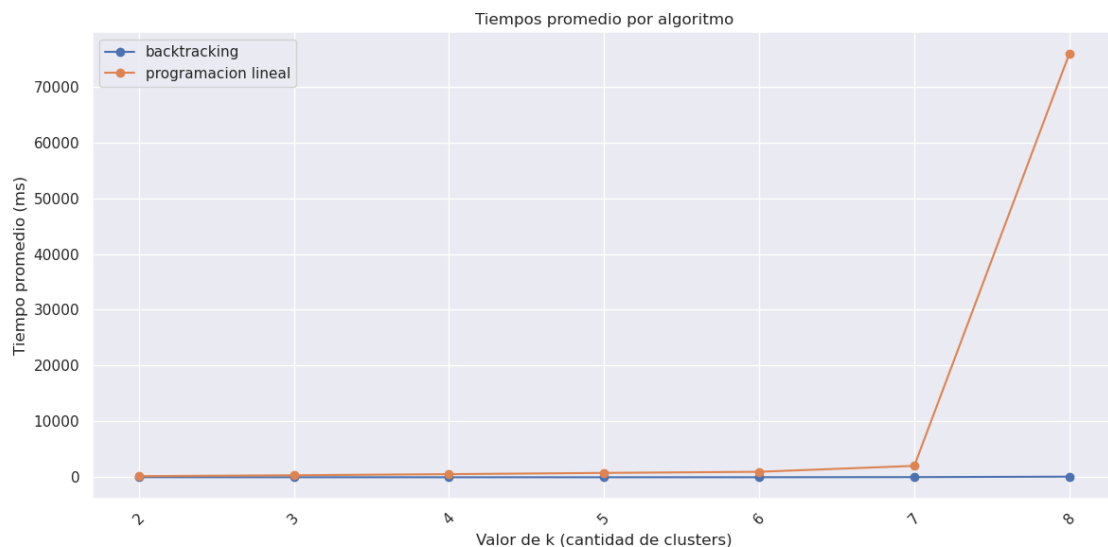
$$\sum_{c=1}^k Y_c \leq K$$

### 6.4.3. Corrección c

En la siguiente figura podemos ver la comparación entre los tiempos promedios de ejecución de los algoritmos de programación lineal y backtracking para el archivo *10\_3.txt*. Puede notarse claramente como a medida que aumenta la cantidad de clusters, el tiempo de ejecución del algoritmo de programación lineal se ve seriamente perjudicado respecto al de backtracking. También podemos ver como llega a un máximo en 7 clusters y reduce el tiempo de ejecución para 8. Esto se debe a que el archivo del grafo analizado cuenta unicamente con 10 vertices.



Para la próxima figura se representa el tiempo de ejecución entre ambos algoritmos utilizando el grafo en el archivo *22\_3.txt*. Acá podemos ver, incluso con mayor grado, como se expande el tiempo empleado por el algoritmo de programación lineal. Vemos como hay una diferencia de mas de un minuto en la ejecución aumentando únicamente en uno la cantidad de clusters 7-8. Mientras que el algoritmo de backtracking parece mantenerse en una zona aceptable, la complejidad del algoritmo de programación lineal se dispara.





## 6.5. Corrección 6: Algoritmo de Aproximación

### 6.5.1. Comparación entre el óptimo y la aproximación

A continuación se muestra una tabla con varias pruebas, comparando el diámetro máximo de los clusters generados por el algoritmo exacto (backtracking) con el resultado del algoritmo de Louvain (greedy). También se calcula la razón entre el valor aproximado y el óptimo, así como el error porcentual:

Prueba	Diámetro máximo (backtracking)	Diámetro máximo (greedy)	$A(I)/z(I)$	Error porcentual
10_3_2.txt	2	3	1,5	50 %
10_3_5.txt	1	3	3	200 %
22_3_3.txt	2	3	1,5	50 %
22_3_4.txt	2	3	1,5	50 %
22_3_10.txt	1	3	3	200 %
22_5_2.txt	2	3	1,5	50 %
22_5_7.txt	1	3	3	200 %
30_3_2.txt	3	4	1,33	33 %
30_3_6.txt	2	4	2	100 %
30_5_5.txt	2	3	1,5	50 %
40_5_3.txt	2	3	1,5	50 %
45_3_7.txt	3	4	1,33	33 %
50_3_3.txt	3	4	1,33	33 %

Viendo los resultados de la relación  $A(I)/z(I)$ , podemos establecer como una cota máxima:

$$r(A) := \frac{A(I)}{z(I)} \leq 3.$$

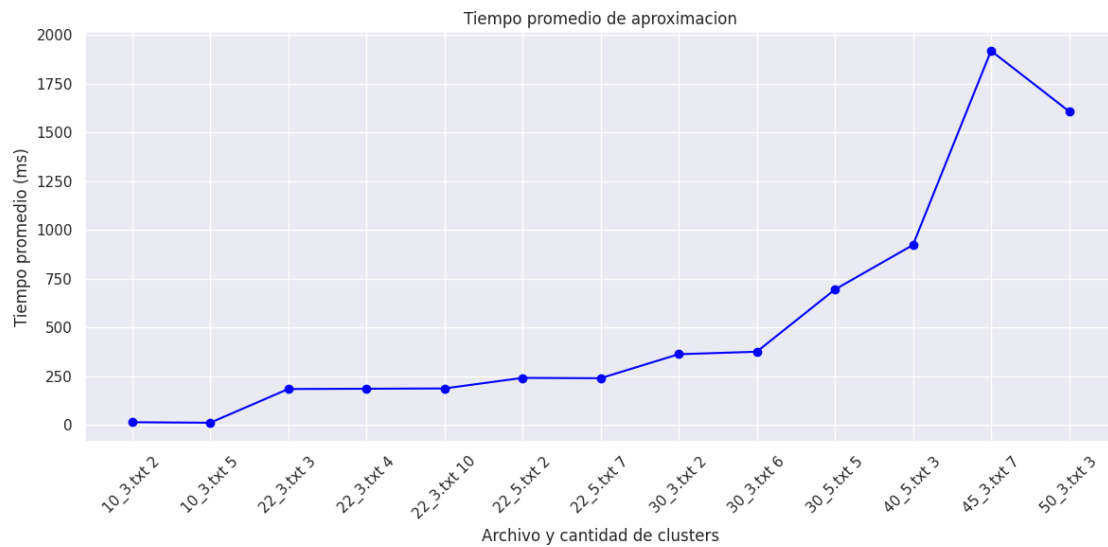
Esto nos da como resultado que, en el peor caso evaluado, el algoritmo de Louvain proporciona una 3-Aproximación respecto al diámetro óptimo.

Las demás cotas encontradas fueron:

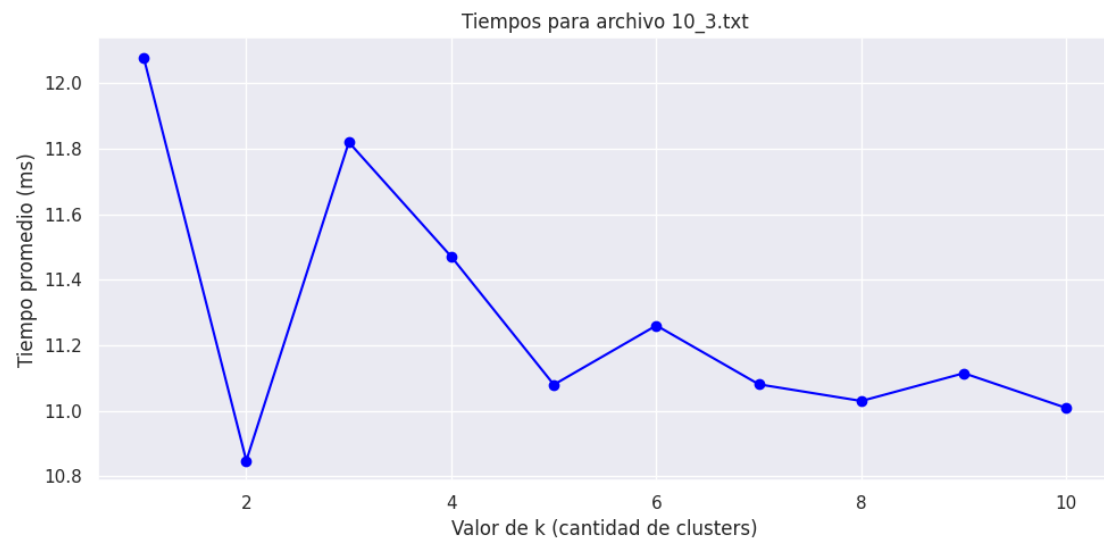
- $A(I)/z(I) \leq 1,33$ , para los casos con un error de 33 %, es decir una  $\frac{1}{3}$ -Aproximación.
- $A(I)/z(I) \leq 1,5$ , con un error de 50 %, es decir una  $\frac{1}{2}$ -Aproximación.
- $A(I)/z(I) \leq 2$ , con un error de 100 %, es decir una 2-Aproximación.

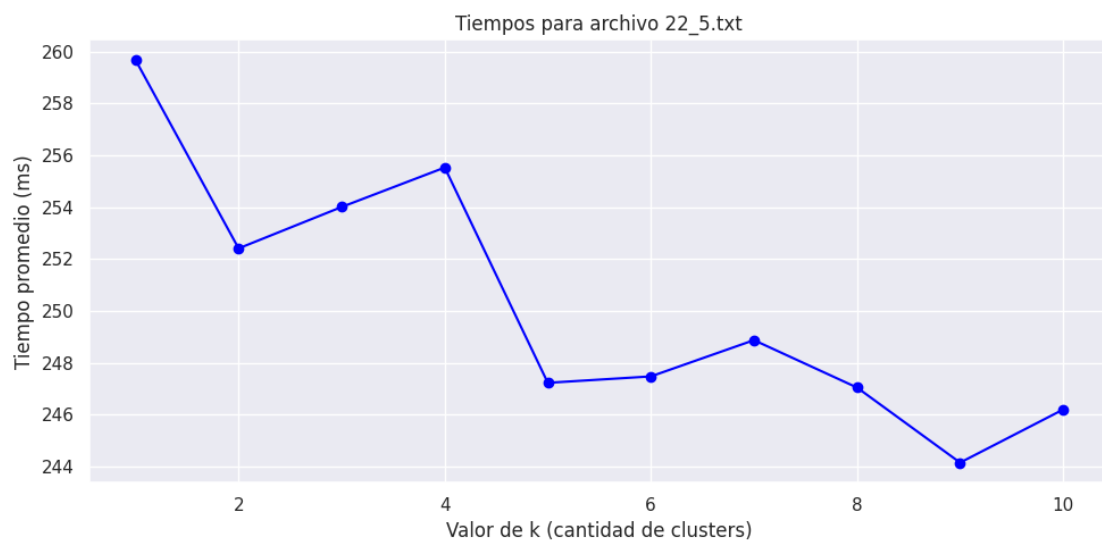
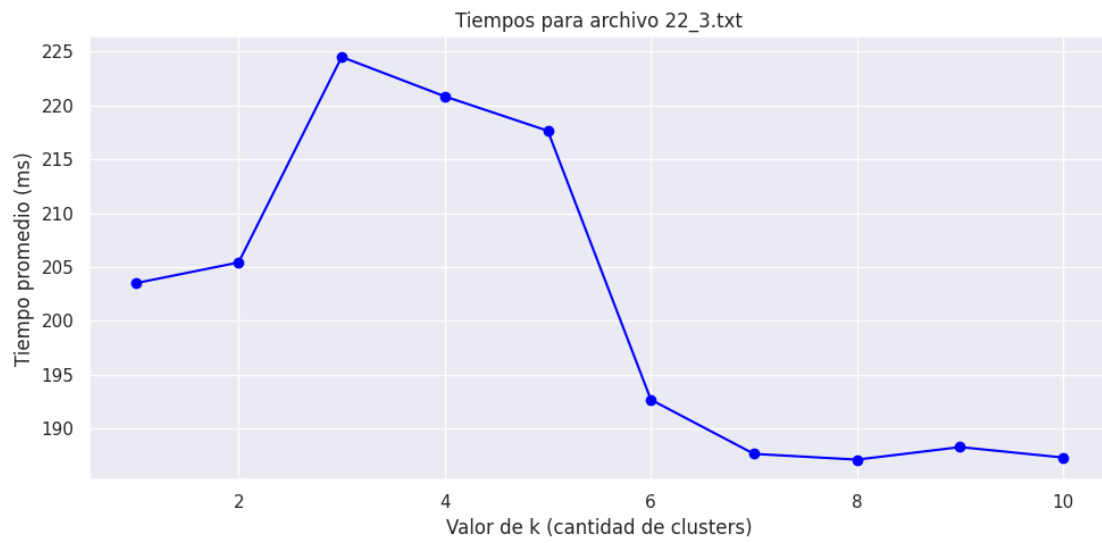
En ninguna de las instancias evaluadas, el algoritmo de Louvain logra alcanzar el resultado óptimo. Sin embargo, se mantiene dentro de un margen razonable de error considerando que se trata de un enfoque heurístico basado en maximización de modularidad y no en minimización directa del diámetro.

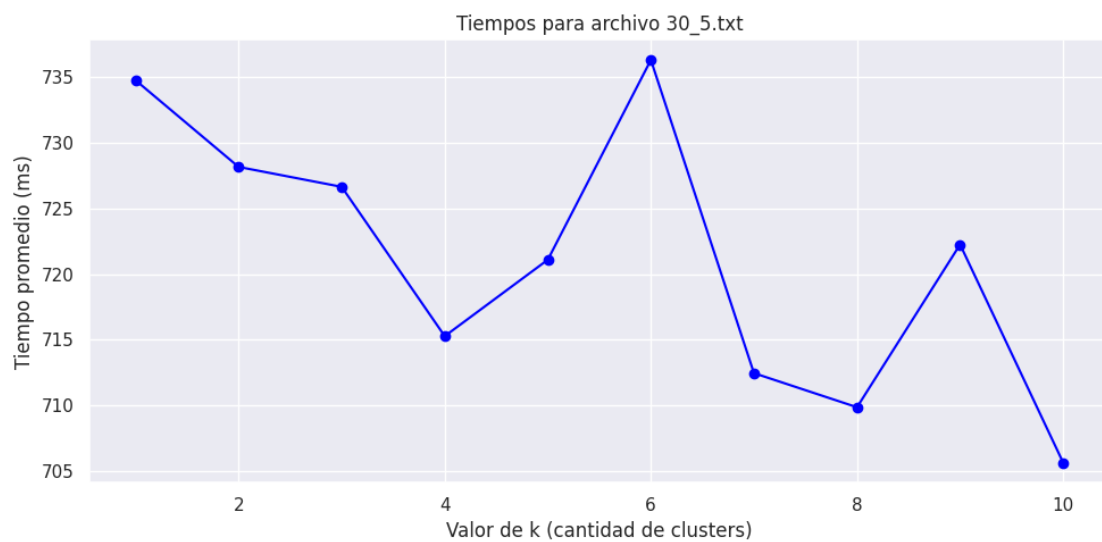
### 6.5.2. Gráfico de tiempo de ejecución

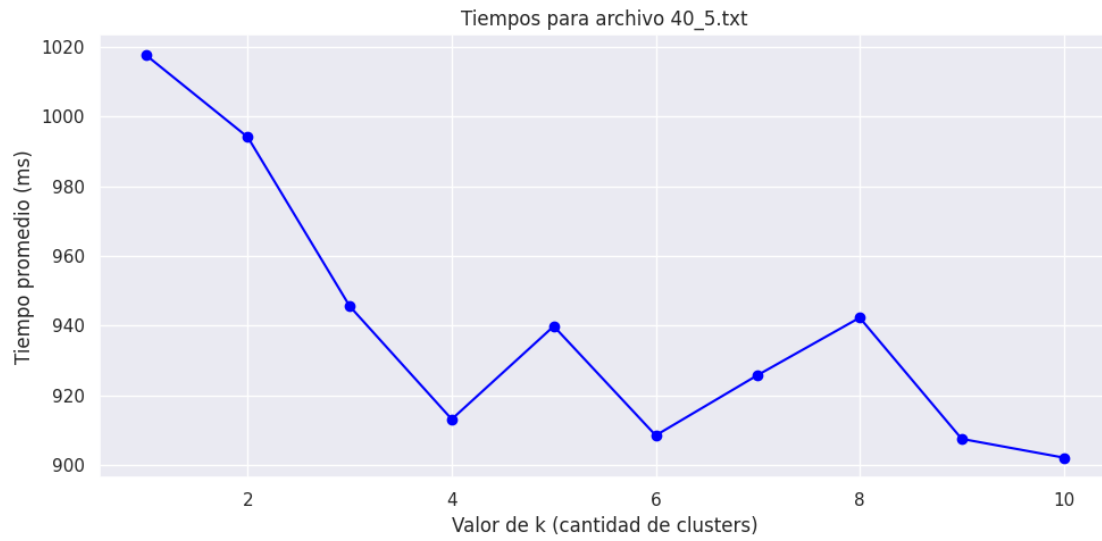


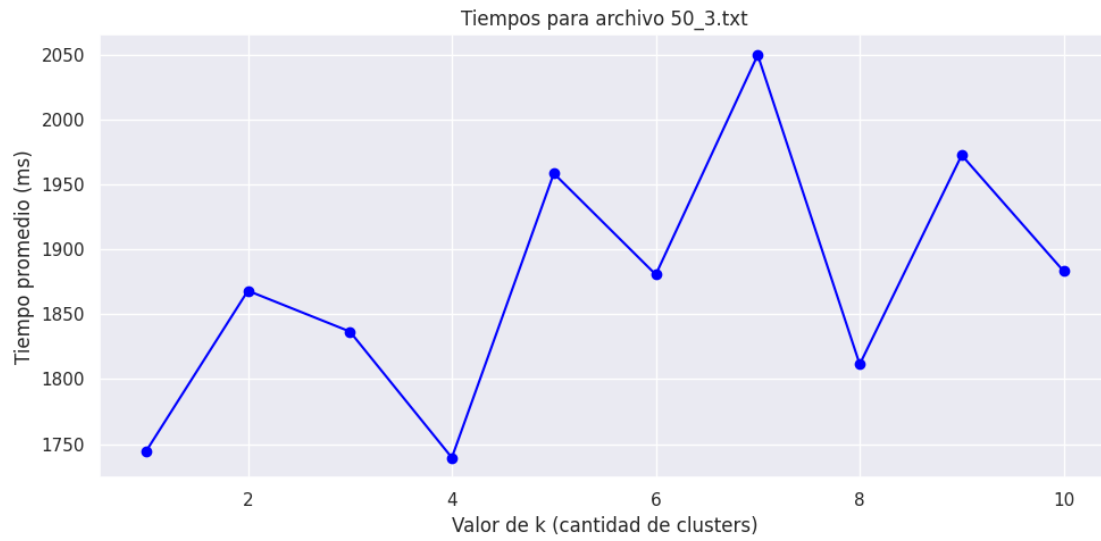
### 6.5.3. Gráficos de variabilidad del valor K





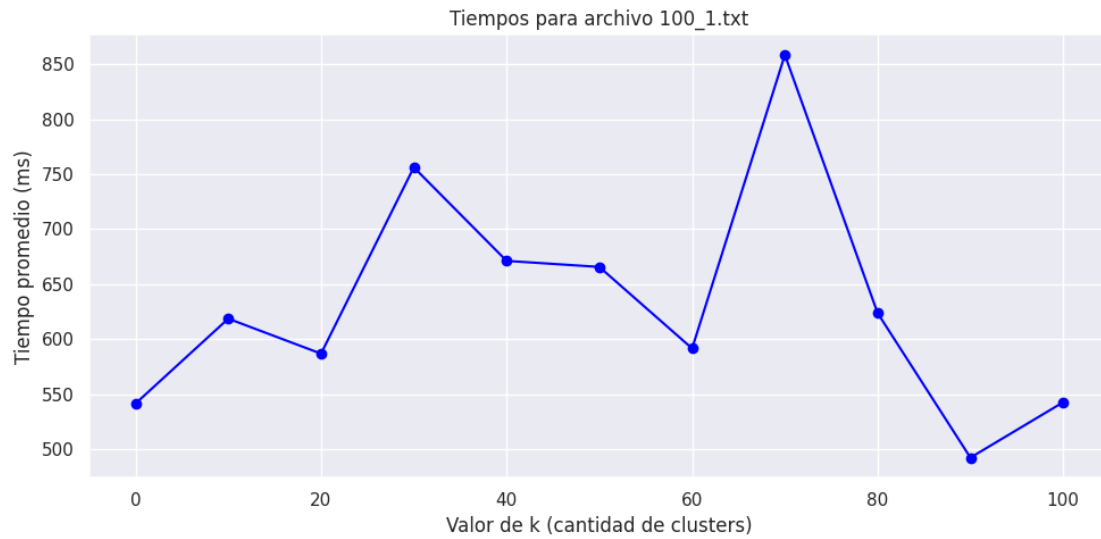






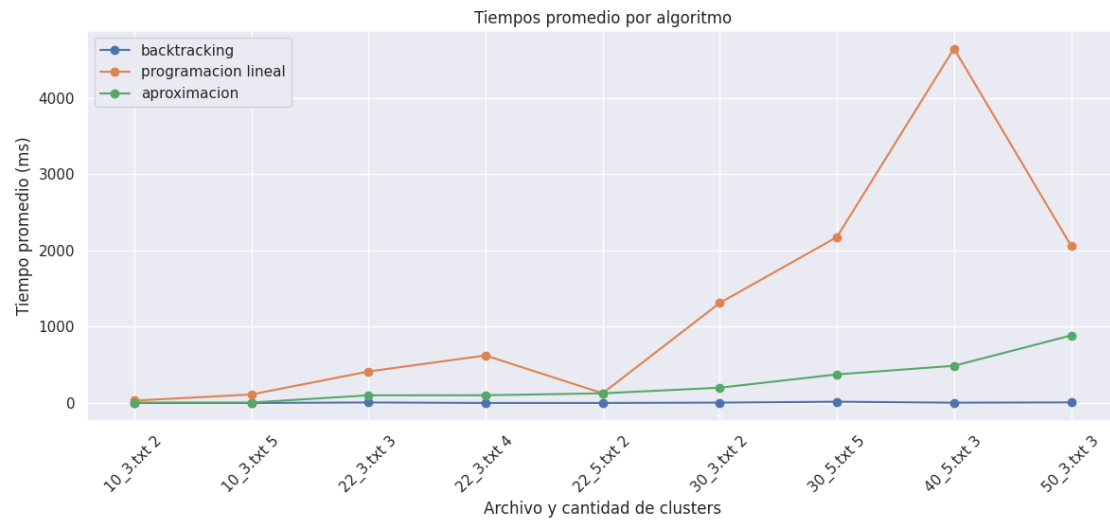
#### 6.5.4. Gráficos de variabilidad del valor K en grafos lista





Se puede notar que a medida que encuentra cierto valor de  $K$  que es el máximo posible, el tiempo de ejecución va bajando. Esto se da por cómo funciona el Algoritmo de Louvain, que como busca comunidades, cuando encuentra la mejor, el grafo ya no cambia. Y esto pasa en los ejemplos cuando se da un cierto valor de  $K$  y desde ese valor en adelante, el resultado es el mismo, es decir, las mismas comunidades que ya son las que cumplen la máxima modularización, que sería la distancia máxima dentro del cluster.

## 6.6. Comparacion



Podemos notar que el algoritmo de backtracking es el más rápido de los tres en los valores pedidos para las pruebas, quitando el caso 45.3.txt con  $K = 7$ , que es donde más afecta la curva exponencial de la complejidad.