

Heuristiken für das Entfernen von verbotenen Teilgraphen

Paul Walger

31. März 2016

Inhaltsverzeichnis

1	Abstract	3
2	Einleitung	3
2.1	Motivation	3
2.2	Anwendungsbeispiele	3
2.2.1	Soziale Netzwerke	3
2.2.2	Protein interaction networks	3
2.2.3	Bicluster Editing	3
2.3	Definitionen	3
2.3.1	Notationen und Definitionen	3
2.3.2	Problemstellung	4
2.4	Ähnliche Arbeiten	4
3	Algorithmen	4
3.1	Top-Bottom	5
3.1.1	RandomChange	5
3.1.2	Random	5
3.2	Bottom-Top	5
3.2.1	Extend	7
3.3	Grow-Reduce	7
3.4	Explored-Grow-Reduce	7
3.5	Lineare Programmierung	9
3.5.1	Lineare Optimierung	9
3.5.2	Das Model des Graphen	9
3.6	Relaxierte Lineare Programmierung	11
4	Aufbau der Test	11
4.1	Datensätze	11
4.1.1	Barabási–Albert	12
4.1.2	Erdős-Rényi	12
4.1.3	duplication divergence	12
4.1.4	newman watts strogatz	12
4.1.5	powerlaw	12
4.1.6	random	12
4.1.7	UCINetworkDataRepository	12
4.1.8	bio1	13
4.1.9	bio2	13
4.2	Optimale Lösung	13

5	Implementation	13
5.1	Repräsentation vom Graphen	13
5.2	Das Finden von induzierten Subgraphen	13
5.2.1	Vergleich VFLib, Boost und eigene Implemtation . . .	13
6	Auswertung	14
7	Vergleich mit anderen Heuristiken	14
7.1	Cluster-Editing	14
7.1.1	2K-Heuristik	14
7.1.2	Andere Heuristiken	15
7.2	Quasi-Threshold Mover	15
8	Zukünftige Forschungsmöglichkeiten	16
9	Zusammenfassung	16
10	Anhang	16
10.1	Kleine Graphen	16

1 Abstract

2 Einleitung

2.1 Motivation

2.2 Anwendungsbeispiele

2.2.1 Soziale Netzwerke

(P_4, C_4) -freie Graphen modellieren eine soziale Struktur. [11]

Ähnlich dazu sind (P_5, C_5) -freie Graphen die auch soziale Strukturen modellieren und dafür geeignet sind Gemeinschaften zu identifizieren. [12]

2.2.2 Protein interaction networks

$(2K_2, C_4, C_5)$ -freie Graphen haben gewisse Vorteile für die Untersuchung von Interaktionsnetzwerken von Proteinen [5].

2.2.3 Bicluster Editing

[7] [9]

2.3 Definitionen

2.3.1 Notationen und Definitionen

Mit Graphen sei im Folgenden stets ein ungerichteter, einfacher Graph gemeint. Wenn nicht anders angegeben ist $G = (V, E)$ ein Graph, V die Menge seiner Knoten und E die Menge seiner Kanten.

$V(G)$ ist die Menge der Knoten des Graphen G . $E(G)$ ist die Menge der Kanten des Graphen G . $N(u)$ ist die Nachbarschaft vom Knoten u . $N^*(u)$ ist die Nachbarschaft von u mit u inklusive.

Sei $G = (V, E)$ ein Graph und $S \subseteq V$ eine beliebige Knotenmenge von V . Dann ist $G[S]$ der auf S induzierte Subgraph von G mit $G[S] = (S, E \cap \{\{u, v\} \mid u \in S \wedge v \in S\})$

Sei $H = (V_H, E_H)$ und $G = (V, E)$ zwei Graphen. Ein Subgraph-Isomorphismus von H nach G ist eine Funktion $f : V_H \rightarrow V$ sodass wenn $(u, v) \in E_H$, dann auch $(f(u), f(v)) \in E$. f ist ein induzierter Subgraph-Isomorphismus, wenn es auch gilt, dass wenn $(u, v) \notin E_H$, dann auch $(f(u), f(v)) \notin E$.

2.3.2 Problemstellung

2.4 Ähnliche Arbeiten

Implicit Hitting Set hilft hier leider nicht viel.[10]

Approximation von H-Free Editing für monotone graphen eigenschaften:
 $o(n^2)$ ist effizient, aber $O(n^{2-\epsilon})$ ist NP-Hard.[2]

3 Algorithmen

Die nachfolgenden beschriebenen Algorithmen basieren alle auf dem folgenden Prinzip: Suche einen validen Graphen, welcher die verbotenen Subgraphen nicht enthält. Dann gebe, die Differenz zwischen dem erstellten validen Graphen und dem Eingabegraphen. Dies wird in dem Algorithmus 1 noch einmal beschrieben. Dabei steht SOLVEALGO für einen der Algorithmen, die wir in den folgenden Abschnitten betrachten werden.

Algorithm 1 Genereller Aufbau

```
1: function SOLVE(graph, forbidden, iterations)
2:   bestGraph  $\leftarrow (\emptyset, \emptyset)$ 
3:   for i = 1 to iterations do
4:     validGraph  $\leftarrow$  SOLVEALGO(graph, forbidden)
5:     if DIFF(bestGraph, graph) < DIFF(validGraph, graph) then
6:       bestGraph  $\leftarrow$  validGraph
7:     end if
8:   end for
9:   print DIFF(graph, bestGraph)
10: end function
```

Da alle Ansätze diesen Schritte enthalten und sich nur in dem unterscheiden, wie der valide Graph gefunden wird, wird folgend nur dieser Aspekt betrachtet.

Die entwickelten Ansätze sind in 3 große Gruppen zu unterteilen. Der Top-Bottom-Ansatz nimmt den Graphen und ändert ihn solange, bis ein gültiger Graph entsteht. Der Bottom-Top-Ansatz fängt mit einem leeren oder vollen Graphen an, und ändert solange Knoten, bis man möglichst nahe an dem Eingabegraphen ist. Der Grow-Reduce-Ansatz kombiniert diese beiden Ansätze, indem es unterschiedliche Stadien gibt...

3.1 Top-Bottom

Der Top-Bottom-Ansatz nimmt den Graphen und ändert ihn solange, bis ein gültiger Graph entsteht.

3.1.1 RandomChange

Das ist der einfachste Algorithmus.

Algorithm 2 RandomChange

```
1: function RANDOMCHANGESOLVE(graph, forbidden)
2:   for Graph f  $\in$  forbidden do
3:     forbiddenSubgraph  $\leftarrow$  FINDFS(graph, f)
4:     while forbiddenSubgraph  $\neq \emptyset$  do
5:       change a random edge  $\in$  forbiddenSubgraph
6:       forbiddenSubgraph  $\leftarrow$  FINDFS(graph, f)
7:     end while
8:   end for
9:   return graph
10: end function
```

3.1.2 Random

Es ist wie RandomChange¹, aber bereits editierte Kanten werden mit einer geringeren Wahrscheinlichkeit geändert. Auch hat es für kleine Graphen ein Konvergenzkriterium. Dieses Konvergenzkriterium besteht darin, dass nach für jede Änderung, die Anzahl der verbotenen Subgraphen gezählt wird und die nur dann ausgeführt wird, wenn die Anzahl der verbotenen Subgraphen dadurch weniger wird. Der allgemeine Vorgehensweise wird im Algorithmus 3 beschrieben. Der große Unterschied zu zum RandomChange sehen wir ab Zeile 6. Dort wählen wird die Kante ausgewählt welche geändert werden soll, aber eine bereits geänderte Kante bekommt eine Wahrscheinlichkeit zugewiesen die 4-mal kleiner ist(siehe Zeile 10).

3.2 Bottom-Top

Die Bottom-Top-Ansätze zeichnet sich dadurch aus, dass wir mit einem Graphen beginnen, der die selben Knoten wie der Eingabegraph hat, aber keine Kanten. Dieser Graph ist somit valide, weil er keine verbotenen Subgraphen

¹Algorithmus 2

Algorithm 3 Random

```
1: function STATERANDOM2SOLVE(graph, forbidden)
2:   for Graph  $f \in$  forbidden do
3:     forbiddenSubgraph  $\leftarrow$  FINDFS(graph,  $f$ )
4:     while forbiddenSubgraph  $\neq \emptyset$  do
5:       foundEdge  $\leftarrow \emptyset$ 
6:       while true do
7:          $e \leftarrow$  random edge from forbiddenSubgraph
8:         prob  $\leftarrow 1 / \#(E(f))$ 
9:         if  $e$  is already visited then
10:          prob  $\leftarrow$  prob / 4
11:         end if
12:         if random number from  $[0,1] >$  prob then
13:          foundEdge  $\leftarrow e$ 
14:          break
15:         end if
16:         flip  $e$  in graph
17:       end while
18:       forbiddenSubgraph  $\leftarrow$  FINDFS(graph,  $f$ )
19:     end while
20:   end for
21:   return graph
22: end function
```

enthält. Dies ist ein Vorteil gegenüber den Top-Bottom-Ansätzen, da es möglich ist immer einen validen Graphen zu haben und somit jederzeit terminieren.

3.2.1 Extend

Algorithm 4 F-Free Extend

```

1: function STATEEXTENDSOLVE(input, forbidden)
2:   graph = (V(input),  $\emptyset$ )
3:   while true do
4:     for each Edge  $e \in \text{DIFFERENCE}(\text{graph}, \text{input})$  do
5:       try to flip  $e$ , revert if it produces an invalid graph
6:     end for
7:     break if there was no change
8:   end while
9:   return graph
10: end function

```

3.3 Grow-Reduce

Der Grow-Reduce-Ansatz ist ein Art von einer greedy randomized adaptive search procedure (GRASP). GRASPH beschreiben Siehe [3] Der Grow-Reduce-Ansatz sieht wie folgt aus: Begonnen wird mit einem Graphen, der die selben Knoten wie der Eingabegraph hat, aber keine kanten. Dann wird in jeder Iteration ein Knoten und seine Umgebung hinzugefügt und durch lokale Suche werden alle neu entstandenen verbotenen Subgraphen wieder entfernt. Die ist im Algorithmus 5 zu sehen.

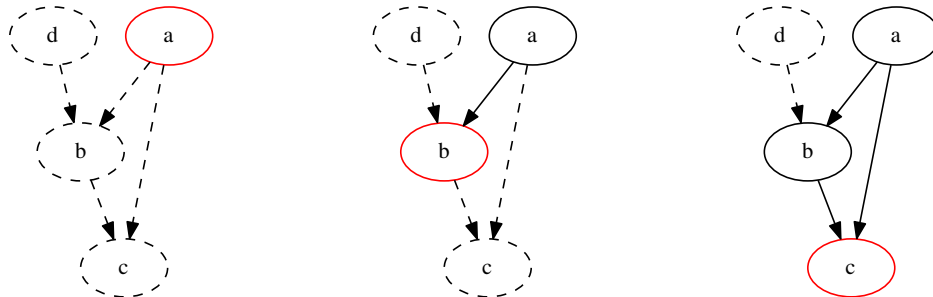
3.4 Explored-Grow-Reduce

Der Explored-Grow-Reduce-Ansatz ist dem Grow-Reduce-Ansatz ähnlich, bis auf das, es in der Grow-Phase nur die Kanten zu Knoten hinzufügt, die bereits erforscht sind.

Dieser Unterschied wird in der Abbildung 1 verdeutlicht, wo nur die Grow Schritte visualisiert wurden, ohne die Reduce-Phase, um es zu vereinfachen. In Abbildung 1a ist der Anfangstatus zu sehen. Die gestrichelten Kanten, sind Kanten die Graphen vorhanden sind, aber noch nicht hinzugefügt worden sind und es wird der Knoten a hinzugefügt, weil aber keine anderen Knoten bisher hinzugefügt worden sind, wird werden keine Kanten hinzugefügt. In

Algorithm 5 GrowReduce

```
1: function GROWREDUCESOLVE(input, forbidden)
2:   graph  $\leftarrow$  (V(input),  $\emptyset$ )
3:   nodes  $\leftarrow$  RANDOMORDER(V(input))
4:   for node  $\in$  nodes do
5:     for neighbor  $\in$  N(node) do ▷ Grow Phase
6:       Add Edge (node, neighbor) to graph
7:     end for
8:     for f  $\in$  forbidden do ▷ Reduce Phase
9:       forbiddenSubgraph  $\leftarrow$  FINDFS(graph, f)
10:      while forbiddenSubgraph  $\neq \emptyset$  do
11:        edge  $\leftarrow$  random Edge from forbiddenSubgraph
12:        count  $\leftarrow$  #(FINDALLFS(graph, f))
13:        flip edge in graph
14:        countAfter  $\leftarrow$  #(FINDALLFS(graph, f))
15:        if countAfter  $\geq$  count then
16:          flip edge in graph
17:        end if
18:        forbiddenSubgraph  $\leftarrow$  FINDFS(graph, f)
19:      end while
20:    end for
21:  end for
22:  return graph
23: end function
```



(a) Knoten a wird hinzugefügt (b) Knoten b wird hinzugefügt (c) Knoten c wird hinzugefügt

Abbildung 1: Beispielweise Grow-Phase

Abbildung 1b wird der Knoten b hinzugefügt und weil a auch schon hinzugefügt wurde, wird auch die Kante (a, b) hinzugefügt. Aber weder (a, c) noch (a, b) werden hinzugefügt, weil c noch nicht erforscht wurde. In Abbildung 1c wird der Knoten c hinzugefügt und somit auch die Kanten (a, b) und (a, c) .

3.5 Lineare Programmierung

3.5.1 Lineare Optimierung

Bei der linearen Optimierung wird eine lineare Zielfunktion minimiert bzw. maximiert, wobei sie durch lineare Gleichungen und Ungleichungen beschränkt ist.

3.5.2 Das Model des Graphen

Wir nutzen binäre Variablen e_{uv} , wobei $u, v \in V$ sind und $u < v$ gilt. Dabei ist $e_{uv} = 1$ genau dann wenn, die kante u, v ein Teil des Lösungsgraphen ist.

Wir minimieren

$$\sum_{u,v \in V} \begin{cases} e_{u,v} & \{u, v\} \in E \\ -e_{u,v} & \{u, v\} \notin E \end{cases}$$

Da alle möglichen Bedingungen hinzuzufügen, welche bei alle verbotenen Subgraphen ausschließen wpürden, viel zum umfangreich wäre, werden die Bedingungen iterative dort hinzu gefügt, wo es einen verbotenen Teilgraphen

Algorithm 6 ExploredGrowReduce

```
1: function EXPLOREDGROWREDUCESOLVE(input, forbidden)
2:   graph  $\leftarrow$  (V(input),  $\emptyset$ )
3:   nodes  $\leftarrow$  RANDOMORDER(V(input))
4:   explored  $\leftarrow \emptyset$ 
5:   for node  $\in$  nodes do
6:     for neighbor  $\in$  N(node) do  $\triangleright$  Grow Phase
7:       if neighbor  $\in$  explored then
8:         Add Edge (node, neighbor) to graph
9:       end if
10:    end for
11:    explored  $\leftarrow$  explored  $\cup$  {node}
12:    for f  $\in$  forbidden do  $\triangleright$  Reduce Phase
13:      forbiddenSubgraph  $\leftarrow$  FINDFS(graph, f)
14:      while forbiddenSubgraph  $\neq \emptyset$  do
15:        edge  $\leftarrow$  random Edge from forbiddenSubgraph
16:        count  $\leftarrow$  #(FINDALLFS(graph, f))
17:        flip edge in graph
18:        countAfter  $\leftarrow$  #(FINDALLFS(graph, f))
19:        if countAfter  $\geq$  count then
20:          flip edge in graph
21:        end if
22:        forbiddenSubgraph  $\leftarrow$  FINDFS(graph, f)
23:      end while
24:    end for
25:  end for
26:  return graph
27: end function
```

gibt. Dann wird der Problem gelöst und die Änderungen auf den Graphen übertragen. Dann wird wieder nach allen verbotenen Subgraphen gesucht. Dies wird solange wiederholt bis es keine mehr gibt. Nun ist die minimale Anzahl von Änderungen gefunden. Dieses Vorgehen ist im Algorithmus 7 zu sehen.

Algorithm 7 F-Free BLP

```

1: function SOLVEBLP(graph, algo:blp)
2:   for graph  $f \in$  forbidden do
3:     while FINDFS(graph,  $f$ )  $\neq \emptyset$  do
4:       for each graph  $M \in$  FINDEVORBOTENESUBGRAPHEN(graph,
5:          $f$ ) do
6:         constraint = 0
7:         for each  $\{u, v\} \in$  kanten( $M$ ) do
8:           if  $\{u, v\} \in$  kanten(graph) then
9:             constraint += 1 -  $e_{uv}$ 
10:          else
11:            constraint +=  $e_{uv}$ 
12:          end if
13:        end for
14:        addConstraint(constraint)
15:      end for
16:      graph = lpSolve(graph)
17:    end while
18:  end for
19:  return (graph)
20: end function

```

3.6 Relaxierte Lineare Programmierung

4 Aufbau der Test

4.1 Datensätze

Folgende Datensätze wurden verwendet. Da verschiedene Mengen von verbotenen Teilgraphen monotonen Grapheneigenschaften zugeordnet werden können und jeder Datensatz von Graphen und jede Methode zufällige Graphen zu erzeugen, charakteristische Eigenschaften hat, ist es notwendig verschiedene Datensätze zu verwenden und verschiedenen Methoden zur Erzeugung

von zufälligen Graphen. Insgesamt wurde 6 verschiedene Methoden zur Erzeugung von zufälligen Graphen verwendet und 3 Datensätze. Wir brachten zuerst die zufälligen Graphen.

4.1.1 Barabási–Albert

Für den Datensatz `barabasi_albert` wurde das Barabási–Albert Model verwendet, welches ein zufälliges skalenfreies Netz erzeugt.[1] Skalenfrei bedeutet hier, dass die Knotengrad einer Potenzverteilung folgt. Es gibt also viel mehr Knoten die einen geringeren Grad haben als Knoten mit einem hohen Anzahl von Nachbarn. Es wurden 56 Graphen generiert mit Knotenanzahl zwischen 10 und 150, wobei der Parameter m , welcher die Anzahl der Kanten definiert, die zu bereits bestehenden Knoten erstellt werden, zwischen 1 und 8 war.

4.1.2 Erdős-Rényi

Für den Datensatz `binomial` wurde das Erdős-Rényi Model verwendet[8], wo jede Kante eine fixe Wahrscheinlichkeit hat zu existieren oder nicht zu existieren. Es wurden dabei 54 Graphen generiert, mit einer Knotenanzahl zwischen 10 und 100 und den folgenden Wahrscheinlichkeiten: $\frac{1}{10}$, $\frac{2}{10}$, $\frac{5}{20}$, $\frac{4}{10}$, $\frac{5}{10}$, $\frac{8}{10}$.

4.1.3 duplication divergence

Anzahl: 54

4.1.4 newman watts strogatz

Anzahl: 144

4.1.5 powerlaw

Anzahl: 30

4.1.6 random

Anzahl: 40ffl

4.1.7 UCINetworkDataRepository

Anzahl: 9

4.1.8 bio1

Anzahl: 147 **Was ist die Quelle für diese Daten**

4.1.9 bio2

Anzahl: 360 **Was ist die Quelle für diese Daten**

4.2 Optimale Lösung

Um die Qualität der Lösung eines heuristischen Ansatzes bewerten zu können, ist es sehr gut die optimale Lösung zu wissen. Es gibt verschiedene Ansatz wie das Problem zu lösen sein, wir haben uns jedoch für die lineare Optimierung entschieden.

5 Implementation

5.1 Repräsentation vom Graphen

Die Graphen werden in einer Adjazenzmatrix gespeichert.

5.2 Das Finden von induzierten Subgraphen

[13] Wie verwenden einen VF Algorithmus für `FINDFS(graph,forbidden)`. Dieser gibt eine Menge von Subgraphen zurück.

5.2.1 Vergleich VFLib, Boost und eigene Implementation

•	find all p3s	count all p3s	has a p3
Spezial	0.73s	0.04s	0.00016s
VFLib	1.73s	0.87s	0.01236s
Boost	3.04s	1.68	0.00102s

Bei VFLib ist der Graph immutable und bei der Suche nach einem Subgraphen müssen wir jedes Mal den Graphen neu erstellen.

6 Auswertung

7 Vergleich mit anderen Heuristiken

7.1 Cluster-Editing

7.1.1 2K-Heuristik

Die 2K-Heuristik, basiert auf einem Kernel für das Cluster-Editing-Problem, welches maximal $2K$ Knoten liefert [6]. Wenn man dort eine Bedingung für die ?. Reduktionsregel abschwächt abschwächt, kommt eine sehr gute Heuristik für das Cluster-Editing-Problem heraus. Dabei wird die Bedingung mit jedem Durchlauf abgeschwächt.

Algorithm 8 2K Heuristik

```
1: function SOLVE2K( $g ::$  Gewichteter Graph)
2:    $a = 1,0$ 
3:   while graph hat einen P3 do
4:     for each knoten  $u \in g$  do
5:       if  $2 \cdot a \cdot \text{costClique}(g, u) + a \cdot \text{costCut}(g, u) < \#(N(u))$  then
6:         for each  $\{a, b\}$  mit  $a \in N(u), b \in N(u) \wedge a \neq b$  do
7:           merge( $a, b$ )
8:         end for
9:       end if
10:    end for
11:     $a = 0,99 \cdot a - 0,01$ 
12:  end while
13:  return graph
14: end function
15: function COSTCLIQUE(graph  $::$  Gewichteter Graph,  $u ::$  Kante)
16:   $\text{cost} = 0$ 
17:  for each  $\{a, b\}$  mit  $a \in N^*(u), b \in N^*(u) \wedge \{a, b\} \notin \text{graph}$  do
18:     $\text{cost} += |w(\{a, b\})|$ 
19:  end for
20:  return cost
21: end function
22: function COSTCUT(graph  $::$  Gewichteter Graph,  $u ::$  Kante)
23:   $\text{cost} = 0$ 
24:  for each  $\{a, b\}$  mit  $a \in N^*(u), b \notin N^*(u) \wedge \{a, b\} \in \text{graph}$  do
25:     $\text{cost} += w(\{a, b\})$ 
26:  end for
27:  return cost
28: end function
```

7.1.2 Andere Heuristiken

[3] Effiziente Algorithmen

GRASP Heuristik ILS Heuristik

7.2 Quasi-Threshold Mover

In [4] wurde ein neuer schneller und auch für große Graphen geeigneter Algorithmus entwickelt für das Quasi-Threshold Editing Problem. Quasi-

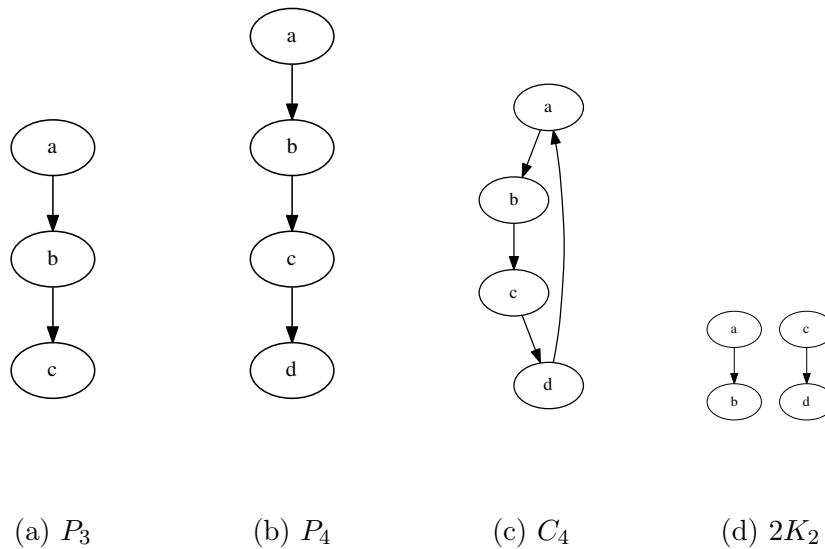


Abbildung 2: Einige Graphen

Threshold Graphen, auch bekannt als trivial perfekte Graphen lassen sich auch als (P_4, C_4) - freie Graphen charakterisieren.

Vergleich mit mit meinem Algorithmus.

8 Zukünftige Forschungsmöglichkeiten

9 Zusammenfassung

10 Anhang

10.1 Kleine Graphen

Literatur

- [1] Reka Albert and Albert-Laszlo Barabasi. Statistical mechanics of complex networks. *Reviews of Modern Physics*, 74:47, 2002.
- [2] Noga Alon and Uri Stav. Hardness of edge-modification problems. *Theoretical Computer Science*, 410(47):4920–4927, 2009.

- [3] Lucas Bastos, Luiz Satoru Ochi, Fábio Protti, Anand Subramanian, Ivan César Martins, and Rian Gabriel S. Pinheiro. Efficient algorithms for cluster editing. *Journal of Combinatorial Optimization*, 31(1):347–371, 2014.
- [4] Ulrik Brandes, Michael Hamann, Ben Strasser, and Dorothea Wagner. Fast quasi-threshold editing. *CoRR*, abs/1504.07379, 2015.
- [5] Sharon Bruckner, Falk Hüffner, and Christian Komusiewicz. A graph modification approach for finding core-periphery structures in protein interaction networks. *Algorithms for Molecular Biology*, 10:16, 2015.
- [6] Jianer Chen and Jie Meng. A 2k kernel for the cluster editing problem. *J. Comput. Syst. Sci.*, 78(1):211–220, 2012.
- [7] Gilberto F de Sousa Filho, F Lucidio dos Anjos, Luiz Satoru Ochi, Fábio Protti, Rio Tinto-PB-Brazil, and Joao Pessoa-PB-Brazil. Metaheuristic grasp for the bicluster editing problem. 2012.
- [8] Edgar N Gilbert. Random graphs. *The Annals of Mathematical Statistics*, 30(4):1141–1144, 1959.
- [9] Sara C. Madeira and Arlindo L. Oliveira. Biclustering algorithms for biological data analysis: A survey. *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, 1(1):24–45, January 2004.
- [10] Erick Moreno-Centeno and Richard M. Karp. The implicit hitting set approach to solve combinatorial optimization problems with an application to multigenome alignment. *Operations Research*, 61(2):453–468, 2013.
- [11] James Nastos and Yong Gao. Familial groups in social networks. *Social Networks*, 35(3):439–450, 2013.
- [12] Philipp Schoch. Editing to (p5, c5)-free graphs - a model for community detection? Bachelor’s thesis (Studienarbeit), Karlsruher Institut für Technologie, October 2015.
- [13] J. R. Ullmann. An algorithm for subgraph isomorphism. *Journal of the Association for Computing Machinery*, 23(1):31–42, 1976.