

Clustering of dynamic graphs

Hauptseminar Networkvisualisation

Moritz Hamann



Abb. 1. In den Wolken: Vancouver von Cypress Mountain. Auf der ersten Seite dürfen keine Grafiken außer dieser optionalen Aufmachgrafik (Teaser) abgebildet sein.

Kurzbeschreibung—This paper presents a summary about various techniques to detect and identify densely connected nodes in a graph, so called clusters. In the first part, we introduce the concept of clusters for static graphs alongside their main properties. For dynamic graphs with time varying edge connections, these clusters may be subject to change with every time step. Therefore additional characteristics have to be introduced.

The second part describes two methods to detect, identify and track clusters in a dynamic graph. A common solution for this problem is the clustering of a static graph at each time step, and the identification of the same clusters over multiple time steps. A method is presented to track these clusters, which is independent of the underlying static graph clustering algorithm. Furthermore, we describe an extension of the k-clique percolation algorithm to dynamic graphs.

Finally, the clique percolation algorithm is applied to two different real world networks, which yields interesting results about group dynamics, with regards to the correlation of various group properties.

1 MOTIVATION

Das mathematische Konzept der Graphen ist ein essentielles Modellierungswerkzeug in der Informatik. Nicht nur lassen sich damit verschiedenste Datenstrukturen anschaulich darstellen, sondern mit ihrer Hilfe lassen sich auch jegliche Beziehungen zwischen einzelnen Objekten oder Prozessen in einem Netzwerk modellieren und untersuchen. Aus diesem Grund sind sie heutzutage nicht nur in der klassischen Informatik sowie in der Mathematik zu finden, sondern haben auch in vielen anderen Wissenschaften ihren Einzug erhalten. So werden sie genutzt um die Gruppendynamik in biologischen Netzwerken zu beschreiben, dienen als Kontrollalgorithmen für Multiagenten Systeme [?] und beschreiben Kommunikationsmuster in sozialen Netzwerken.

Um die Eigenschaften sehr großer Netzwerke analytisch untersuchen zu können, werden häufig Zufallsgraphen nach dem Model von Edgar Gilbert (nachweise?) oder Erdos-Renyi verwendet. Diese Graphen haben die Besonderheit, dass die Wahrscheinlichkeit für eine Verbindung zwischen je zwei Knoten im gesamten Netzwerk konstant ist. Dadurch entsteht ein gleichmäßiger Graph, dessen Gradverteilung binomial verteilt sind, und somit die meisten Knoten die gleiche Anzahl an Kanten haben. Mit Hilfe der Wahrscheinlichkeitstheorie, lassen sich nun die Eigenschaften dieser Graphen auch für eine sehr hohe Anzahl an Knoten bestimmen und untersuchen.

Allerdings haben Untersuchungen von realen Netzen gezeigt (nachweis), dass sich diese in den meisten Fällen von Zufallsgraphen unterscheiden. Reale Netzwerke sind häufig sogenannte Skalenfreie Netze (im Englischen 'Scale-free networks'), in denen die Anzahl der Verbindungen pro Knoten nicht binomial verteilt ist, sondern nach einem Potenzgesetz. Dadurch entsteht ein Netzwerk, in dem einzelne wenige Knoten eine große Anzahl an Verknüpfungen aufweisen, doch die Mehrzahl der Knoten weniger stark verknüpft ist. Weiterhin ist die Kantenverteilung zwischen den Knoten auch lokal sehr inhomogen, so dass sich Teilgraphen ausbilden, deren Knoten untereinander sehr stark bis komplett verknüpft sind, während sie nach Außen weniger Verbindungen aufweisen. Diese Teilgraphen werden auch 'Cluster' genannt.

Diese Cluster spielen in vielen Anwendungsgebieten eine wichtige Rolle. Betrachtet man zum Beispiel den Graph der Freundschaftsbeziehungen in einem sozialen Netzwerk, lassen sich mithilfe von Angaben anderer Benutzer, sowie lokaler Cluster, unter anderem Rückschlüsse auf gemeinsame Interessen, Wohnorte oder Freunde der einzelnen Benutzer schließen. Diese Informationen bieten dem sozialen Marketing eine bis vor kurzem unbekannte Menge an Möglichkeiten ihre Produkte zielgerichteter und persönlicher zu vermarkten. Aber auch in dynamischen Graphen, in denen Knoten und Kanten sehr häufig wechseln können, ist es wichtig Cluster zu finden. Betrachtet man den Verbindungsgraph eines dezentralen, kabellosen Ad-Hoc Netzwerks, so lassen sich mit Hilfe von Clustering Verfahren Teilnetze

finden die geographisch Eine zusätzliche Herausforderung zur eigentlichen Clusteranalyse ist hierbei allerdings

Übersicht

Diese Arbeit gibt einen Überblick über die Eigenschaften dieser Cluster, ihr

2 GRUNDLAGEN

2.1 Formale Definition eines Clusters

In diesem Kaptiel wird versucht eine formale Definition zu geben, was ein Cluster in Graphen zu. Während die gewünschten Eigenschaften

2.1.1 Eigenschaften von Clustern

Der Artikel von S.E. Schaeffer [?] bietet eine umfassende Zusammenfassung über bisherige Clustering Verfahren, und versucht eine Definition für Cluster anhand gewünschter Eigenschaften zu geben.

Betrachtet man den Teilgraphen Ω eines kompletten Graphen Υ , so müssen mehrere Bedingungen erfüllt sein, damit Ω ein Cluster wird. Natürlich sollten alle Knoten aus Ω verbunden sein, was bedeutet das zwischen jedem Paar aus Knoten u und v mit $u, v \in \Omega$ ein Pfad existiert. Ist dies nicht der Fall so ist der gesamte Graph nicht verbunden, und das Clustering sollte auf den einzelnen Teilgraphen gesondert betrachtet werden. Weiterhin sollte der Teilgraph Ω eine hohe Kantendichte zwischen seinen Knoten aufweisen. Dies ist der Fall, wenn mehrere Pfade zwischen den Knoten aus Ω existieren, so dass jeder Pfad möglichst wenig Elemente aus $\Upsilon \setminus \Omega$ enthält.

Der Grad $d(v)$ eines Knoten v ist definiert als die Anzahl der Kanten zu anderen Knoten im Graphen Υ . Ist nun $v \in \Omega$ wobei Ω wieder ein Teilgraph von Υ ist, so lässt sich der Grad in einen externen und internen Teil unterscheiden. Dabei ist der interne Grad die Anzahl der Kanten von v zu anderen Knoten aus Ω , der externe Grad die Anzahl der Kanten von v zu allen anderen Knoten aus $\Upsilon \setminus \Omega$. Dabei gilt:

$$d_{int}(v, \Omega) = |\Gamma(v) \cap \Omega| \quad (1)$$

$$d_{ext}(v, \Omega) = |\Gamma(v) \cap (\Upsilon \setminus \Omega)| \quad (2)$$

$$d(v) = d_{int}(v, \Omega) + d_{ext}(v, \Omega) \quad (3)$$

wobei $\Gamma(v)$ die direkten Nachbarn von v sind. Eine Eigenschaft, die den Teilgraphen Ω zu einem Cluster werden lässt, ist ein hohes Verhältniss von internem zu externem Grad für alle Knoten $v \in \Omega$, d.h. die Knoten eines Cluster haben untereinander wesentlich mehr Verknüpfungen als zu den Knoten des restlichen Graphen.

Eine weiteres Kriterium für die Qualität eines Clusters ist die sogenannte interne Clusterdichte. Die allgemeine Dichte eines Graphen $\Upsilon = (V, E)$ mit der Knotenmenge V und der Kantenmenge E ist definiert als das Verhältniss der Summe aller Kanten durch die Anzahl aller möglichen Kanten in Graph:

$$\rho(\Upsilon) = \frac{|E|}{\binom{|V|}{2}} = \frac{2|E|}{|V|(|V| - 1)} \quad (4)$$

Somit lässt sich die interne Clusterdichte eines Clusters Ω definieren als

$$\rho_{int}(\Omega) = \frac{|\{\{u, v\} | u, v \in \Omega\}|}{|\Omega|(|\Omega| - 1)} \quad (5)$$

wobei $\{u, v\}$ eine Kante zwischen den Knoten u und v darstellt. Die fehlende 2 im Zähler im Vergleich zur allgemeinen Graphendichte ist dadurch zu erklären, dass $\{u, v\}$ und $\{v, u\}$ zwar die gleiche Kante darstellen, aber zwei unterschiedliche Elemente sind, wodurch die Anzahl der Kanten in $\{\{u, v\} | u, v \in \Omega\}$ verdoppelt wird.

Zusätzlich zur internen Clusterdichte, existiert noch eine sogenannte externe Clusterdichte zwischen verschiedenen Clustern Ω_i eines Graphen Υ . Sie ist definiert als das Verhältniss der Summe aller Kanten

zu der Summe aller möglichen Kanten zwischen den verschiedenen Clustern Ω_i :

$$\rho_{ext}(\Upsilon | \Omega_1 \dots \Omega_k) = \frac{|\{\{v, u\} | v \in \Omega_i, u \in \Omega_j, i \neq j\}|}{|V|(|V| - 1) - \sum_{i=1}^k |\Omega_i|(|\Omega_i| - 1)} \quad (6)$$

wobei $|\Omega_i|$ die Anzahl der Knoten des Teilgraphen Ω_i darstellt. Im Allgemeinen kann man sagen, dass für ein gutes Clustering auf einem Graph Υ gelten sollte:

$$\rho_{int}(\Omega_i) > \rho(\Upsilon) > \rho_{ext}(\Upsilon | \Omega_1 \dots \Omega_k) \quad (7)$$

$$\forall i = 1 \dots k$$

Abb. 2 zeigt drei verschiedene Cluster unterschiedlicher Qualität im Vergleich. Dabei representieren die schwarz hervorgehobenen, beliebig gewählten Teilgraphen jeweils einen Cluster. Der linke Cluster weist eine sehr hohe interne Dichte auf, und hat kaum Kanten mit Knoten ausserhalb. Daher ist Qualität dieses Clusters sehr hoch. Der mittlere Cluster hat zwar die gleiche Anzahl an internen Kanten, weist aber im Gegensatz zum Linken eine wesentliche höhere Kantenzahl zu Knoten ausserhalb des Cluster auf. Zwar hat der rechte Cluster nur wenige Kanten nach aussen, allerdings ist aber die Kantendichte innerhalb des Clusters minimalst, was ihn zum schlechtesten Cluster der drei macht.

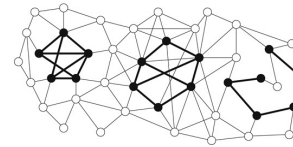


Abb. 2. Drei verschiedene Cluster unterschiedlicher Qualität [referenz]

2.1.2 k-Cliques

Einen Teilgraphen Ω mit k Knoten nennt man k -Clique, falls alle k Knoten dieses Teilgraphen direkt miteinander verbunden sind, und Ω somit vollständig ist. Die entstehende Topologie des Teilgraphen Ω ist natürlich vom Parameter k abhängig. Abb. ?? zeigt k -Cliques für verschiedene Werte von k .

Da jede k -Clique vollständig ist, ist die in Kapitel 2.1.1 definierte interne Graphdichte mit $\rho_{int}(\Omega) = 1$ maximal. Somit stellen k -Cliques theoretisch gute Kandidaten für Cluster da. Beschränkt man sich bei der Clusterfindung auf einzelne k -Cliques, werden in den meisten Fällen viele Cluster nicht gefunden, da Forderung an einen vollständig verknüpften Cluster zu restriktiv ist.

2.2 Dynamische Graphen

Ein klassischer Graph $\Omega = (V, E)$ ist eine Kombination einer Menge an Knoten V und Kanten E zu einem bestimmten Zeitpunkt t . Für viele Anwendungen ist es aber essentiell das Netzwerk über einen Zeitraum mit mehreren Zeitschritten t_i zu betrachten und untersuchen.

Das Konzept der *dynamischen Graphen*[Quelle] stellt die zeitliche Veränderung eines Graphen als geordnete Folge von statischen Teilgraphen für jeden Zeitpunkt $t = 1, \dots, n$ da:

$$\Upsilon = \{\Omega_1 = (V_1, E_1), \Omega_2 = (V_2, E_2), \dots, \Omega_n = (V_n, E_n)\} \quad (8)$$

dabei entspricht Ω_i der Konfiguration des dynamischen Graphen Υ zum Zeitpunkt i .

3 CLUSTERING IN STATISCHEN GRAPHEN

Dieses Kaptiel gibt eine kurze Übersicht

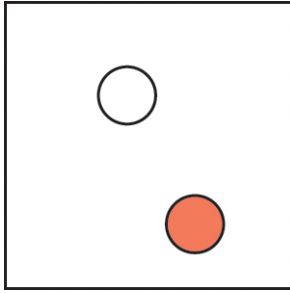


Abb. 3. Beispielillustration.

3.1 Modularitätsoptimierung

3.2 Clique Percolation Method (CPM)

4 CLUSTERING IN DYNAMIC GRAPHS

4.1 Extension of clique percolation

4.2 Time step Clustering

5 MAYBE VISUALIZATION OF DYNAMIC GRAPHS

6 GROUP EVOLUTION (RESULTS)

7 CONCLUSION

7.1 Abbildungen und Tabellen

Alle Abbildungen (siehe Abb. 3) und Tabellen (Tabelle 1) sollten zentriert sein (`\centering`). Abbildungen über beide Textspalten (Abb. 4) können mit `\begin{figure*}...\end{figure*}` eingefügt werden.

7.2 Referenzen

Literaturangaben wie beispielsweise Levoy [?] werden mit Hilfe von BibTeX erzeugt. Dazu werden die Referenzen in die Literaturliste (hier *literatur.bib*) eingetragen und entsprechend mit `\cite` referenziert.

7.3 L^AT_EX-Übersetzung

Die L^AT_EX-Datei kann mit *latex* oder *pdflatex* übersetzt werden. Dabei ist zu beachten, dass für die Übersetzung mit *latex* die Grafiken in Postscript (eps) vorliegen, für *pdflatex* entsprechend als jpg, png oder pdf. Der Ablauf ist dabei der folgende:

1. `pdflatex <quelldatei.tex>`
2. `bibtex <quelldatei>`
3. `pdflatex <quelldatei.tex>` (evtl. mehrfach)

Alternativ kann auch das mitgelieferte Makefile verwendet werden.

Tabelle 1. Vis Paper Acceptance Rate

Year	Submitted	Accepted	Accepted (%)
1994	91	41	45.1
1995	102	41	40.2
1996	101	43	42.6
1997	117	44	37.6
1998	118	50	42.4
1999	129	47	36.4
2000	151	52	34.4
2001	152	51	33.6
2002	172	58	33.7
2003	192	63	32.8
2004	167	46	27.6
2005	268	88	32.8
2006	228	63	27.6

$$\sum_{j=1}^z j = \frac{z(z+1)}{2} \quad (9)$$

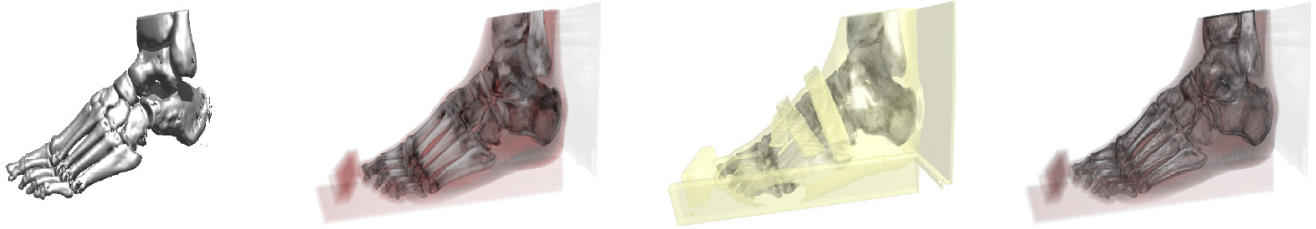


Abb. 4. Illustration über beide Textspalten hinweg. Auch Illustrationen müssen entsprechend den Quellen gekennzeichnet werden. [?]