

AGT Summary

June 12, 2024

Contents

1	Netzwerke und Zentralität	2
1.1	Charakterisierung der wichtigsten Ecke	2
1.2	Berechnung der Zentralitätsmaße	2
1.3	Random Walks auf Graphen	3
1.4	Eigenwert Zentralität	4
1.5	PageRank	5
2	Clustering	5
2.1	Berechnung des Clustering Koeffizienten	6
2.2	Gomory-Hu Clustering	7
2.3	Berechnung des Gomory-Hu Baums	8
2.4	ratio cuts	8
2.5	Cluster-Metrik	9
3	Streaming	10
3.1	HyperLogLog	10
3.2	Streaming mit Graphen	11
3.3	Spanners	11
3.4	sparsifier	12
3.5	Minimaler Schnitt als Streaming Algorithmus	13
3.6	Maximalgrad Streaming	14

1 Netzwerke und Zentralität

1.1 Charakterisierung der wichtigsten Ecke

Hierfür gibt es mehrere Möglichkeiten

- größter Einfluss
- wichtig für Informationsfluss

Die Wichtigkeit wird mit einem Zentralitätsmaß gemessen.

Definition 1.1.1. Zentralitätsmaße sind sehr unterschiedlich. Es muss nur erfüllt sein, dass bei einem Sterngraphen das Zentrum das größte Zentralitätsmaß erhält. Möglich sind Bewertungen nach

1. dem Maximalgrad (*degree centrality*)
2. der durchschnittlichen Entfernung zu anderen Ecken (*closeness centrality*) (bzw. der Kehrwert davon)
3. der Anzahl der Komponenten, die mit dieser Ecke verbunden sind (*betweenness centrality*). Dafür sei $\sigma_{s,t}$ die Anzahl der kürzesten $s - t$ -Wege. $\sigma_{s,t}(v)$ für $v \neq s, t$ ist dann die Anzahl der kürzesten $s - t$ -Wege, die durch v gehen. Damit gilt

$$betweenness(v) = \sum_{s,t \in V()G \setminus \{v\}} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}$$

1.2 Berechnung der Zentralitätsmaße

Wir führen nur die Berechnung der betweenness ein. Die anderen beide Maße sind sehr einfach.

Der Algorithmus zur Berechnung von $\sigma_{s,t}$ ist an Dijkstra angelehnt. Beginnend mit s wird die Anzahl der Nachbarn von s bestimmt. Anschließend die Anzahl der Knoten mit Abstand 2 usw. Um die Komplexität der Algorithmen zu bestimmen, werden im Folgenden einige Annahmen getroffen:

1. Knotenadjazenz kann in $\mathcal{O}(1)$ bestimmt werden
2. Kanteninzidenz kann in $\mathcal{O}(1)$ bestimmt werden
3. die Nachbarschaft eines Knoten wird in $\mathcal{O}(1)$ pro Knoten bestimmt
4. die zu einem Knoten inzidenten Kanten können in $\mathcal{O}(1)$ pro Kante bestimmt werden

5. alle elementaren Operationen (z.B. Kante löschen) in $\mathcal{O}(1)$.

Auf diese Weise kann man leicht sehen, dass die Laufzeit zur Berechnung von $\sigma_{s,t}$ für alle s, t in $\mathcal{O}(n \cdot m)$ implementiert werden kann. Wir nehmen nun an, $\sigma_{s,t}$ sei bekannt und wir definieren

$$\rho_s(v) = \sum_{t \neq v} \frac{\sigma_{s,t}(v)}{\sigma_{s,t}}$$

Kennt man nun alle $\rho_s(v)$, dann ist

$$betweenness(v) = \frac{1}{2} \sum_{s \neq v} \rho_s(v)$$

Lemma 1.2.1. *Sei v ein Knoten mit Distanz mindestens $d \geq 1$ zu s und sei L die Menge der Knoten mit Distanz $d+1$ zu s . Dann ist*

$$\rho_s(v) = \sum_{w \in L \cap N(v)} \frac{\sigma_{s,w}}{\sigma_{s,v}} (1 + \rho_s(w))$$

Mit dieser Überlegung lässt sich ein Algorithmus finden, der die *betweenness* jedes Knotens in $\mathcal{O}(mn)$ berechnet.

1.3 Random Walks auf Graphen

Wir wählen zunächst einen Startknoten v_0 bezüglich einer Wahrscheinlichkeitsverteilung $\pi^{(0)}$. Anschließend wird mit Gleichverteilung ein zufälliger Nachbar v_1 von v_0 gezogen usw. Wenn man sich nun die Frage stellt, was die Wahrscheinlichkeit ist, dass der erste gezogene Knoten (d.h. v_1) gleich dem Knoten u ist, dann entspricht das der Wahrscheinlichkeit, dass u ein Nachbar von v_0 ist mal der Wahrscheinlichkeit, dass anschließend u gezogen wird. Da der zweite Schritt gleichverteilt ist, ergibt sich

$$\pi_u(1) = \sum_{v \in N(u)} \pi_v^{(0)} \cdot \frac{1}{d(v)}$$

Das wird geschrieben als Transition Matrix mit

$$T_{uv} = \begin{cases} \frac{1}{d(v)}, & \text{wenn } uv \in E \\ 0, & \text{sonst} \end{cases}$$

Schreibt man die Wahrscheinlichkeitsverteilung $\pi^{(0)}$ einfach als Vektor, dessen Komponenten zu 1 addieren, ergibt sich

$$\pi^{(n+1)} = T\pi^{(n)}$$

für $n \geq 0$. Um zu überprüfen, ob ein bestimmter Knoten im Random Walk jemals besucht wird, muss die Grenzwertverteilung bestimmt werden

$$\pi^* = \lim_{k \rightarrow \infty} T^k \pi^{(0)}$$

Existiert π^* , dann ist π^k eine Cauchy-Folge und man kann leicht sehen, dass $T\pi^* = \pi^*$. Dann ist π^* also ein Eigenvektor zum Eigenwert 1 von T .

Satz 1.3.1 (Perron-Frobenius). *Sei $A \in \mathbb{R}^{n \times n}$ sodass $\exists k \in \mathbb{N}$ mit $A_{ij}^k > 0$ für alle $i, j \in [n]$. Dann gibt es einen eindeutigen Eigenwert λ^* mit größtem Betrag. Wenn $\lambda^* > 0$ gibt es einen positiven Eigenvektor v^* zu λ^* und alle anderen Eigenvektoren zu λ^* sind Vielfache von v^* . Ist außerdem $\lambda^* = 1$, dann konvergiert $v^{(k+1)} = Av^{(k)}$ gegen ein Vielfaches von v^* für alle positiven Startvektoren $v^{(0)} > 0$.*

Damit kann man sich überzeugen, dass T Eigenwert 1 hat und dass das der größte Eigenwert ist. Außerdem erfüllt T die Eigenschaft aus obigem Theorem, wenn G zusammenhängend und nicht bipartit ist.

1.4 Eigenwert Zentralität

Für einen Knoten v verwenden wir wieder die Matrix T und nehmen $\pi^* > 0$ als den Eigenvektor zum Eigenwert 1 mit $\|\pi^*\| = 1$. Der Eintrag π_v^* ist dann die Eigenwert Zentralität von v .

Das Problem dieses Zentralitätsbegriffs ist, dass man den Eigenwert recht leicht erraten kann. Betrachte dazu

$$\bar{\pi}_v = \frac{d(v)}{2|E|} \quad \forall v \in V$$

Es ist leicht zu sehen, dass dieser Vektor ein Eigenvektor zum Eigenwert 1 ist. Es folgt also, dass wir einen neuen Begriff haben, der aber sehr ähnlich zur *degree centrality* ist. In gerichteten Graphen ist der Begriff ein wenig hilfreicher. Der wichtigste Unterschied ist die modifizierte Matrix T mit

$$T_{vu} = \begin{cases} \frac{1}{d^+(u)}, & \text{wenn es eine Kante von } u \text{ nach } v \text{ gibt} \\ 0, & \text{sonst} \end{cases}$$

Das größere Problem sind Senken (d.h. Knoten v mit ausgehendem Grad $d^+(v) = 0$). Das kann gelöst werden, indem der Prozess neugestartet wird (d.h. eine Kante zu jedem anderen Knoten eingeführt wird).

1.5 PageRank

PageRank ist der Suchalgorithmus von Google. Er funktioniert in den folgenden Schritte, die sehr ähnlich zum Eigenwertzentralität sind

1. Wähle unter Gleichverteilung einen Startknoten
2. Mit Wahrscheinlichkeit $1 - \alpha$ (α konstant) wähle einen Nachbarn und gehe dorthin.
3. Mit Wahrscheinlichkeit α wähle einen neuen Startknoten.

Die Transitionsmatrix definiert nun eine andere Matrix

$$P = (1 - \alpha)T + \frac{\alpha}{n}J$$

wobei J die Matrix mit nur 1 Einträgen ist. Es ergibt sich der Prozess

$$\pi^{(k+1)} = P\pi^{(k)}$$

Da P positiv ist, ergibt Theorem 1.3.1 die Existenz der Grenzwertverteilung p_v^* . Es gilt $\text{PageRank}(v) = \pi_v^*$. Da der erste Teil von P *sparse* ist, kann die Iteration relativ effizient durchgeführt werden.

2 Clustering

Wie führen zunächst den Begriff des Clustering-Koeffizienten ein. Sei $v \in V$. dann ist

$$C(v) = \frac{|E_G[N(v)]|}{\binom{|N(v)|}{2}}$$

Der durchschnittliche Clustering-Koeffizient ist dann

$$C(G) = \frac{1}{|V|} \sum_{v \in V} C(v)$$

Ein Zufallsgraph mit Kantenwahrscheinlichkeit p hat im Erwartungswert eine Kantendichte $\frac{|E|}{\binom{n}{2}}$ von ungefähr p . Der Clustering-Koeffizient ist ebenso ungefähr p .

2.1 Berechnung des Clustering Koeffizienten

Es ist leicht zu sehen, dass man den Clustering Koeffizienten eines einzelnen Knotens v in $\mathcal{O}(d(v)^2)$ berechnen kann. Um den durchschnittlichen Wert zu bestimmen genügt daher eine Laufzeit von $\mathcal{O}(\sum_{v \in V(G)} d(v)^2)$. Die Summe lässt sich nach oben abschätzen durch $2mn$ wodurch die Laufzeit bei $\mathcal{O}(2mn)$ liegt.

Ist $d(v)$ klein, so ist der Algorithmus sehr effizient, aber ist $d(v) \gg \sqrt{m}$ so lässt sich eine Verbesserung erzielen, indem für jede Kante uw überprüft wird, ob u, v, w ein Dreieck bilden. Kombiniert man diese beiden Überlegungen zu einem Algorithmus mittels einer Fallunterscheidung, so erhält man einen Algorithmus zur Berechnung des durchschnittlichen Clustering Koeffizienten mit einer Laufzeit von $\mathcal{O}(m^{\frac{3}{2}})$.

Es gibt außerdem einen randomisierten Ansatz für die Schätzung des durchschnittlichen Clustering Koeffizienten auf Graphen mit Minimalgrad mindestens 2. Hierfür wird zunächst eine Konstante $k \in \mathbb{N}$ festgelegt. Anschließend werden nacheinander k Knoten v_1, \dots, v_k zufällig gezogen und aus $N(v_i)$ werden jeweils zwei Nachbarn u_i, w_i zufällig gezogen. Es wird gezählt, wie viele dieser Nachbarn der k Knoten mit v_i ein Dreieck aufspannen und diese Anzahl anschließend durch k geteilt.

Satz 2.1.1. *Sei $\varepsilon > 0, \delta > 0$ und $k = \lceil \ln(\frac{2}{\delta}) / (2\varepsilon^2) \rceil$. Dann hat der Algorithmus eine Laufzeit von $\mathcal{O}(\ln(\frac{1}{\delta}) / \varepsilon^2 \cdot \ln n)$ und mit Wahrscheinlichkeit mindestens $1 - \delta$ unterscheidet sich der berechnete Wert um maximal ε vom tatsächlichen Wert.*

Sei G ein Graph mit $V(G) = U \dot{\cup} W$. Wir schreiben

$$\partial_G = \{uw \in E(G) : u \in U, w \in W\}$$

Für $A \subset V$ ist $\partial_G(A)$ die Anzahl der Kanten zwischen A und $v \setminus A := B$. Ist w eine Funktion, die jeder Kante ein Gewicht zuordnet, definiere

$$w(F) = \sum_{e \in F} w(e)$$

für alle $F \subset E$.

Definition 2.1.2. Die *expansion* von A ist

$$\frac{w(\partial_G(A))}{\min\{|A|, |B|\}}$$

Der *ratio cut* von A ist

$$\frac{w(\partial_G(A))}{|A| |B|}$$

2.2 Gomory-Hu Clustering

Definition 2.2.1. Sei G ein Graph und w die Kantengewichte. Der Gomory-Hu-Baum T für G ist ein Graph mit

- $V(T) = V(G)$
- beim Löschen einer Kante uv aus dem Baum entstehen zwei Komponenten $T_{uv}(u)$ und $T_{uv}(v)$. Für alle $uv \in E(T)$ soll gelten

$$w(\partial_G(T_{uv}(u))) = \min_{X \subseteq V(G): v \notin X \ni u} w(\partial_G(X))$$

Wir definieren darauf aufbauend

$$\lambda(s, t) = \min_{X \subseteq V(G): t \notin X \ni s} w(\partial_G(X))$$

Lemma 2.2.2. Sei G ein Graph mit Kantengewichten w und T ein Gomory-Hu-Baum für G, w . Seien $s, t \in V(G), s \neq t$, sei P der st -Weg in T und $uv \in E(P)$ mit

$$\min_{ab \in E(P)} w(\partial_G(T_{ab}(a))) = w(\partial_G(T_{uv}(u)))$$

Dann ist $w(\partial_G T_{uv}(u)) = \lambda(s, t)$.

Satz 2.2.3. Für alle G, w existiert ein Gomory-Hu-Baum. Ein solcher Baum kann in $\mathcal{O}(n\tau)$ berechnet werden. τ ist dabei die Laufzeit um einen gewichtsminimalen $s - t$ Schnitt für beliebige s, t zu finden.

Mit diesem Konzept kann nun ein Clustering in den folgenden Schritten gefunden werden

1. füge einen universellen Knoten t mit Kantengewichten α zurück
2. berechne den G-H-Baum T
3. gib die Komponenten von $T - t$ als Cluster zurück

Lemma 2.2.4. Wir arbeiten auf einem Graphen G mit Gewichten w . Mit dem G-H-clustering erreichen wir ein Cluster $C \subseteq V(G)$. Dann gilt

$$\frac{w(\partial_G C)}{|V \setminus C|} \leq \alpha$$

Lemma 2.2.5. *Teilt man in der Situation von oben das Cluster C noch weiter in Q, P , dann gilt*

$$\frac{w(\partial_G(P, Q))}{\min\{|P|, |Q|\}}$$

2.3 Berechnung des Gomory-Hu Baums

Lemma 2.3.1. *Die Gewichte eines Cuts sind submodular. Für alle $U, W \subseteq V$ gilt*

$$w(\partial U) + w(\partial W) \geq w(\partial(W \cap U)) + w(\partial(U \cup W))$$

Lemma 2.3.2. *Seien $s, t \in V(G)$ und sei für $X \subseteq V$ $\partial_G X$ ein minimaler $s-t$ -Cut. Nun wird X zu einem neuen Knoten v_x kontrahiert. Wobei mehrfache Kanten als eine Kante mit der Summe der Gewichte eingeführt wird. Sei $p, q \in V \setminus X$ und $U \subseteq V \setminus X$ sodass $\partial_{G/X}(U \cup v_x)$ ein minimaler $p-q$ -Cut von G/X ist. Dann ist $\partial_G(U \cup X)$ ein minimaler $p-q$ -Cut in G .*

Definition 2.3.3 (Teil GH Baum). Sei $R \subseteq V(G)$. Dann ist ein Baum $T = (R, F)$ mit einer Partition $(C_r)_{r \in R}$ von R ein GH Baum für R , wenn

$$\forall uv \in F : \partial_G \left(\bigcup_{r \in V(T_{uv}(u))} C_r \right)$$

Ist $R = V(G)$, dann entspricht der GH-Baum für R dem GH-Baum für G .

Mit dieser Überlegung lässt sich der GH-Baum von G rekursiv aufbauen.

2.4 ratio cuts

Die Idee ist, dass für ein gegebenes $k \in \mathbb{N}$ eine Partition C_1, \dots, C_k berechnet wird, sodass

$$\sum_{i=1}^k w(\partial C_i)$$

minimal ist. Damit nicht nur isolierte Knoten geclustert werden, soll der *ratio cut* minimiert werden:

$$\min_{C_1, \dots, C_k} \sum_{i=1}^k \frac{w(\partial C_i)}{|C_i|}$$

Da dieses Problem aber NP-schwer ist, soll stattdessen eine Annäherung gefunden werden.

Definition 2.4.1. Wie immer ist $G = (V, E)$ und w eine Gewichtsfunktion auf den

Kanten. Die Laplace Matrix $L \in \mathbb{R}^{V \times V}$ von G ist

$$L_{u,v} = \begin{cases} -w_{uv}, & \text{if } uv \in E \\ \sum_{e \in \delta(u)} w(e), & \text{if } u = v \\ 0, & \text{sonst} \end{cases}$$

Man kann schreiben $L = B^T \cdot D \cdot B$ wobei $D \in \mathbb{R}^{E \times E}$ im Feld (e, e) das Gewicht $w(e)$ hat und 0 sonst. B ist aus $\mathbb{R}^{E \times V}$. Für B geben wir allen Kanten aus G eine Richtung vor. In der Spalte $e = (u, v)$ steht 1 in Spalte v , -1 in Spalte u , wenn (v, u) eine Kante ist und 0 sonst.

Wir können nun $D^{1/2}$ definieren als jede Matrix die Wurzel genau die Einträge der Matrix D hat. Diese Definition ist daher sinnvoll, da D eine Diagonalmatrix ist. Es folgt $L = (D^{1/2}B)^T(D^{1/2}B)$ und

$$x^T L x = \|D^{1/2} B x\|_2^2 = \sum_{uv \in E} w_{uv} (x_u - x_v)^2$$

Lemma 2.4.2. *Seien G und w wie immer. Sei L die Laplace Matrix von G . Dann*

1. *L ist symmetrisch und positiv semi-definit*
2. *kleinster Eigenwert ist 0*
3. *ist \mathcal{C} die Menge der Komponenten von G . Dann ist $\chi_C, C \in \mathcal{C}$ orthogonale Eigenbasis des Eigenwerts 0.*

Zurück zu ratio cuts. Sei $A \subsetneq V$ mit $A \neq \emptyset$. Wähle

$$\mathbb{R}^V \ni z = \sqrt{\frac{|A|}{|V|}} \chi_A - \sqrt{\frac{|A|}{|V|}} \chi_{\bar{A}}$$

woraus folgt

$$z^T L z = |V| \sum_{uv \in \partial A} w_{uv} \frac{1}{|A|} + \frac{1}{|\bar{A}|} \sum_{j=1, j \neq i}^k |\partial(C_i, C_j)|$$

In der Theorie ist dieses Verfahren sehr interessant, allerdings normalerweise nicht praktisch umsetzbar. Deswegen wird das nicht weiter verfolgt.

2.5 Cluster-Metrik

Wenn die verschiedenen Algorithmen miteinander verglichen werden, ist ein Test-Graph nötig. Solche Testfälle werden meist so aufgebaut, dass ein objektiv bestes Clustering ex-

istiert (*ground truth*). Dann muss verglichen werden, welcher Algorithmus ein Clustering produziert, dass am nächsten an dieses Clustering heranreicht. Dafür sind Clustering-Metriken nötig.

Rand Distanz: Seien A, B zwei Clusterings. Wir nennen $u, v \in V$ eine Unstimmigkeit, wenn $\exists a \in A : u, v \in a$ aber $\nexists b \in B : u, v \in b$ oder umgekehrt. $d(A, B)$ als die Anzahl der Unstimmigkeiten ist eine Metrik. Diese Metrik ist leicht zu berechnen, denn

$$d(A, B) = \sum_{a \in A} \binom{|a|}{2} + \sum_{b \in B} \binom{|b|}{2} - 2 \sum_{a \in A, b \in B} \binom{|a \cap b|}{2}$$

Entropie: Sei V der Größe n und seien A, B, C drei Clusterings. Diese werden in folgender Weise als Zufallsvariablen modelliert: wähle ein $v \in V$ mit gleichmäßiger Wahrscheinlichkeit. Dann ist

$$X_A : V \rightarrow A, v \mapsto a \ni v$$

eine Zufallsvariable wobei

$$\mathbb{P}[X_A = a] = \frac{|a|}{n}$$

Es sei dann weiter

$$H(A) = H(X_A) = - \sum_{a \in A} \frac{|a|}{n} \log_2 \left(\frac{|a|}{n} \right)$$

Die Entropie des Clusterings. Weiter ist

$$H(A, B) = - \sum_{a \in A, b \in B} \frac{|a \cap b|}{n} \log_2 \left(\frac{|a \cap b|}{n} \right)$$

und $VI(A, B) = 2H(A, B) - H(A) - H(B)$ die *variation of information*. Diese ist dann eine Metrik.

3 Streaming

3.1 HyperLogLog

Es wird eine Reihe an Zahlen eingelesen, die nicht in den Speicher passt. Wie kann die Anzahl der paarweise verschiedenen Zahlen festgestellt werden?

Die Zahlenreihe a_1, \dots, a_n wird beim Einlesen in Binärdarstellung umgewandelt. Die Binärzahlen werden anschließend mit einer zufälligen aber deterministischen Transformation in eine andere Binärzahl konvertiert. Bei jeder transformierten Zahl werden die letzten Stellen betrachtet und die Anzahl der 0en am Ende gezählt. Die maximale Zahl R von 0en am Ende einer Zahl wird gespeichert. Ausgegeben wird 2^R .

Satz 3.1.1. Sei F_0 die Anzahl der paarweise verschiedenen Elemente eines Streams von m Zahlen und sei Y der Output des Algorithmus ($Y = 2^R$). Jedes a des Streams liegt in $\{1, \dots, n\}$. Es wird $\mathcal{O}(\log n)$ Platz gebraucht und für alle Integer $c > 2$ ist

$$\mathbb{P} \left[\frac{1}{c} \leq \frac{Y}{F_0} \leq c \right] \geq 1 - \frac{3}{c}$$

3.2 Streaming mit Graphen

Die Knoten des Graphens sind initial im Speicher. Dieser ist nur $\mathcal{O}(n \cdot \log(n)^p)$ für ein p . Das heißt die Kanten (ungefähr n^2) können nicht gespeichert werden. Der Stream besteht nun aus den Kanten, also e_1, e_2, \dots, e_m .

3.3 Spanners

Wir wollen für einen sehr großen Graphen G eine komprimierte Version dieses Graphens konstruieren, die die jeweiligen Abstände innerhalb des Graphens bestmöglich erhält.

Definition 3.3.1. Gegeben ein Graph G , ein α -spanner ist ein Subgraph H s.d.

$$d_H(u, v) \leq \alpha d_G(u, v) \quad \forall u, v$$

Wir wollen einen $(2k+1)$ -spanner berechnen mit $\mathcal{O}(m \log n)$ Zeit. Der Algorithmus wird randomisiert laufen, wobei die erwartete Größe des Graphen $\mathcal{O}(kn^{1+\frac{1}{k}})$ ist. Der folgende Algorithmus dient als Setup.

1. initialisiere $V_0 = V, V_1, \dots, V_{k-1} = \emptyset$
2. for $i = 1$ to $k - 1$
3. for $x \in V_{i-1}$
4. mit Wahrscheinlichkeit $n^{-\frac{1}{k}}$ ergänze x zu V_i und setze $B_{i,x} = \{x\}$ und $h(x) = i$
5. endfor
6. endfor

$B_{i,x}$ sind bags um das Zentrum x und $h(x)$ ist das größte i , in dem x noch vorkommt.

1. verwende das Setup von oben
2. setze $H = (V, \emptyset)$ und $l(v) = h(v)$ für alle $v \in V$

3. for $uv = \text{next}(S)$
4. stelle sicher, dass $l(u) \leq l(v)$ durch Tauschen
5. Finde das Zentrum x des bags mit $v \in B_{l(u),x}$
6. if $u \notin B_{l(u),x}$ und es keine Kante uw in H gibt mit $w \in B_{l(u),x}$ ergänze uv zu H
7. if $l(u) < l(v)$
8. ergänze u zu $B_{l(u)+1,x}, \dots, B_{h(x),x}$
9. setze $l(u) = h(x)$
10. endif
11. endfor

Lemma 3.3.2. Die erwartete Anzahl an bags $B_{k-1,x}$ im Level $k-1$ ist $n^{\frac{1}{k}}$.

Lemma 3.3.3. Ist $u \in B_{i,x}$ dann gibt es einen u, x -Pfad der Länge höchstens i .

Lemma 3.3.4. H ist ein $(2k-1)$ -spanner von G .

Lemma 3.3.5. $\mathbb{E}[|E(H)|] \leq kn^{1+\frac{1}{k}}$

Aus einer Vermutung von Erdős folgt, dass ein $(2k-1)$ -spanner $\Omega(n^{1+\frac{1}{k}})$ Kanten braucht.

Satz 3.3.6. Für alle Primpotenzen q und $k \in \{2, 3\}$ gibt es einen bipartiten Graphen $D_k(q)$ mit $n = 2q^k$ Knoten, $(\frac{n}{2})^{1+\frac{1}{k}}$ Kanten und Umfang mindestens $2k+2$.

3.4 sparsifier

Sei G ein Graph, $\varepsilon > 0$. Suche $H \subset G$, $w : E(H) \rightarrow \mathbb{R}_+$ sodass $\forall \emptyset \neq A \subsetneq V(G)$ gilt

$$(1 - \varepsilon) |\partial_G(A)| \leq w(\partial_H(A)) \leq (1 + \varepsilon) |\partial_G(A)|$$

Ohne das Streaming Setting ist das Problem sehr einfach:

1. Berechne kleinsten cut λ
2. sei $H = (V, \emptyset)$ und

$$p = 9 \ln\left(\frac{n}{\varepsilon^2 \lambda}\right)$$

3. $\forall e \in E(G)$ ist e in $E(H)$ mit Wahrscheinlichkeit p

4. Gewichte überall $\frac{1}{p}$

Satz 3.4.1. *Mit Wahrscheinlichkeit $1 - \mathcal{O}(\frac{1}{n})$ ist H ein sparsifier.*

Lemma 3.4.2. *Sei $\lambda \geq 1$ die kleinste Schnittgröße und $\alpha > 0$. Dann gilt: die Anzahl der Schnitte der Kardinalität $\leq \alpha\lambda$ ist begrenzt durch $n^{2\alpha}$.*

Satz 3.4.3 (Chernoff-Ungleichung).

$$\mathbb{P}[|X - \mathbb{E}[X]| > \varepsilon \mathbb{E}[X]] \leq 2e^{-\frac{\varepsilon^2 \mathbb{E}[X]}{3}}$$

Definition 3.4.4. Sei G ein Graph. Ein Teilgraph H ist ein k -Skelett, wenn $\forall C \subseteq E(G)$ mit $|C| < k$ gilt, dass C ein Schnitt in G ist genau dann, wenn C ein Schnitt in H ist.

Ohne Streaming kann ein k -Skelett mit folgendem Algorithmus gefunden werden:

1. for $i = 1, \dots, k$
2. Berechne F_i einen aufspannenden Wald von $G - \bigcup_{j=1}^{i-1} E(F_j)$
3. endfor
4. $H = \bigcup_{i=1}^k F_i$

Als Streaming-Algorithmus funktioniert das Konzept genau gleich.

1. for $e = \text{next}(S)$:
2. sei i der minimale Index, sodass $F_i \cup e$ keinen Kreis enthält. Setze $F_i = F_i \cup e$.
3. endfor

3.5 Minimaler Schnitt als Streaming Algorithmus

Wir wollen den minimalen Schnitt approximieren. Mittels Karger kann ein ε -sparsifier bestimmt werden, aber dafür ist a priori der Wert λ , also der minimale Schnitt nötig. Es wird also stattdessen $p = \frac{1}{2^r}$ für $r \in \{1, \dots, l\}$ sodass $2^l \approx \log n$. Es werden aus dem Stream die gewichteten Graphen $G, 1/p$ bestimmt. Diese sind aber immer noch zu groß für den Speicher. Deshalb wird aus $G, 2^r$ direkt mittels Streaming ein k -Skelett H_r bestimmt (für $k \approx \log n$). Anschließend wird jenes i bestimmt, sodass $\text{mincut}(H_i) < k$ und geben dann $2^i \cdot \text{mincut}(H_i)$ zurück.

Satz 3.5.1. *Mit hoher Wahrscheinlichkeit gibt der Algorithmus einen Wert α zurück, sodass*

$$(1 - \varepsilon) \text{mincut}(G) \leq \alpha \leq (1 + \varepsilon) \text{mincut}(G)$$

3.6 Maximalgrad Streaming

Das Problem, den Maximalgrad eines Graphen mit Streaming zu bestimmen ist mit Speicherplatz $\mathcal{O}(n)$ leicht. Das ist aber nicht nötig, wie wir im Folgenden sehen werden.

Beispiel 3.6.1 (Zählen mit Hashing). *Wir wollen in einem allgemeinen Stream die häufigsten Elemente bestimmen.*

Die Idee ist, die Elemente des Streams zu hashen und auf einem kleinen Hash-Array die Häufigkeiten zu zählen. Das Problem ist nun, die Kollisionen zu minimieren. Eine Möglichkeit ist, das Hash-Array als Hash-Matrix aufzubauen, wobei jede Zeile einem Hash-Array entspricht. Exakte Kollisionen sind dann wesentlich seltener und um die Häufigkeiten aus der Hash-Matrix zu bestimmen wird schließlich der minimale Wert der Hash-Stellen eines Elements bestimmt.

Die Hashfunktionen h_i sollten der Form

$$h : a \mapsto ((\alpha_i a + \beta_i) \bmod p) \bmod l$$

wobei α_i, β_i gleichverteilt in $\{0, 1, \dots, p-1\}$ für eine Primzahl $p \leq n$ ist und l die Größe der Hashtabelle.

Satz 3.6.2. *Sei k die Anzahl der Hashfunktionen. Sei \hat{f}_a die bestimmte Anzahl an Vorkommen des Elements a in einem Stream der Länge m . Es gilt $\hat{f}_a \geq f_a$ die tatsächliche Anzahl an Vorkommen. Außerdem gilt,*

$$\mathbb{P}[\hat{f}_a \leq f_a + \varepsilon m] \geq 1 - \left(\frac{1}{\varepsilon l}\right)^k$$

Beispiel 3.6.3 (heavy hitters). *Für eine Grenze T wollen wir nun bestimmen, welche Elemente mindestens T mal vorkommen. Mit der Vorüberlegung ist das recht leicht, denn wir können den Stream durchgehen und überprüfen ob $\hat{f}_a \geq T$ ist. Wenn schon ergänzen wir a zu unserer Liste. Auf genau diese Weise können die Knoten mit größtem Grad eines Graphen bestimmt werden.*

Wenn wir nun einen Stream $A = (a_1, \dots)$ mit $a_i \in \{0, 1, \dots, n-1\}$ haben, so können wir a schreiben als $e_{a+1} \in \mathbb{R}^n$. Sei dann $f = \sum_{a \in A} e_{a+1}$. Sind dann c_r die Zeilen der Hash-Matrix C , so gilt der folgende

Satz 3.6.4. $\exists L$ sodass

$$\begin{pmatrix} c_1^T \\ c_2^T \\ \vdots \\ c_k^T \end{pmatrix} = Lf$$

Lemma 3.6.5 (Johnson-Lindenstrauss). *(Das soll nur eine sehr grobe Darstellung sein). Es gibt für Daten einer sehr hohen Dimension eine passende Zufallstransformation L , die mit hoher Wahrscheinlichkeit die Abstände zwischen Originaldaten beibehält, wenn es die Daten auf eine wesentlich niedriger dimensionale Bildmenge abbildet.*

3.7 Dynamisches Streaming

In diesem Fall kann sich der Graph innerhalb des Streams verändern. D.h. der Stream S besteht aus Elementen

$$S = ((e_1, \sigma_1 = \pm 1), (e_2, \sigma_2 = \pm 1), \dots)$$

wobei $+1$ bedeutet, dass die Kante hinzugefügt wird und -1 , dass die Kante gelöscht werden soll. Anstelle von ± 1 kann auch ein Gewicht w_i übergeben werden, wobei $w_i = 0$ dann bedeuten soll, dass die Kante entfernt wird. Gegeben eine Position t im Stream sei $G_t = (V, \{e_i : i \leq t \wedge \sum_{e_j=e_i} \sigma_j = 1\})$.