

Algorithmische Methoden der Netzwerkanalyse

Marco Gaertler

27. Januar, 2009

1 Clustering

2 Qualitätsmaße

Satz

Das Problem zu einem gegebenen Graphen $G = (V, E)$ zu entscheiden, ob es eine Kantenreduktionsmenge gibt, die höchstens K Elemente enthält (und den Graphen in **mindestens drei** Cluster aufspaltet), ist \mathcal{NP} -vollständig.

Anmerkung: Es gibt eine Konstante $\varepsilon > 0$, so dass es \mathcal{NP} -schwer ist das obige, allgemeine Entscheidungsproblem mit einem Faktor von $1 + \varepsilon$ zu approximieren.

Berechnung mittels ILP

Satz

Das folgende ILP berechnet die Editierdistanz eines beliebigen Graphen $G = (V, E)$, wobei n^2 binäre Variablen $X_{u,v}$ und $n + n^2 + n^3$ viele Nebenbedingungen benötigt werden.

$$\text{obj : min } \sum_{u \in V} \sum_{v \in V} \left((1 - E_{u,v}) \cdot X_{u,v} + E_{u,v} \cdot (1 - X_{u,v}) \right)$$

$$\forall v \in V : X_{v,v} = 1$$

$$\forall u, v \in V : X_{u,v} - X_{v,u} = 0$$

$$\forall u, v, w \in V : X_{u,v} + X_{v,w} - 2 \cdot X_{u,w} \leq 1$$

1 Clustering

2 Qualitätsmaße

Bewertungsfunktion

6/12

Definition

Eine Bewertungsfunktion index ordnet jeder Clustering eines Graphens eine nicht-negative reelle Zahl zu (typischer in $[0, 1]$), die die Qualität der Clustering bezüglich eines Paradigmas beschreibt.

Verfeinerung

Eine Bewertungsfunktion index für das Paradigma *intra-cluster density vs. inter-cluster sparsity* hat oft die Form:

$$\text{index}(\mathcal{C}) = \frac{f(\mathcal{C}) + g(\mathcal{C})}{\max_{\mathcal{C}'} (f(\mathcal{C}') + g(\mathcal{C}'))} ,$$

wobei f die Dichte der Cluster und g die Düntheit der Verbindungen mißt.

Interpretation

7/12

- Bewertungsfunktionen stellen eine Formalisierung der umgangssprachlich Beschreibung eines Paradigmas da
- je größer der Wert einer Bewertungsfunktion desto besser ist die Clusterung bzgl. des Paradigmas
- die Bestimmung des Nenners, also die Berechnung der maximalen Qualität, ist oft schwer
- $\max_{C'}(f(C') + g(C'))$ wird oft durch eine obere Schranke M ersetzt

Distanz/Ähnlichkeit zum Clustergraphen

8/12

Idee: Eine Clusterung ist dann gut, wenn es weniger Kantenkorrekturen bedarf, um den Clustergraphen zu erzeugen, der die Cluster als Zusammenhangskomponenten hat.

Distanz/Ähnlichkeit zum Clustergraphen

8/12

Idee: Eine Clusterung ist dann gut, wenn es weniger Kantenkorrekturen bedarf, um den Clustergraphen zu erzeugen, der die Cluster als Zusammenhangskomponenten hat.

$$f(C) = \sum_{C \in \mathcal{C}} |E(C)|$$

$$g(C) = \frac{1}{2} \sum_{C \in \mathcal{C}} \sum_{\substack{u \in C, \\ v \in V \setminus C}} [\{u, v\} \notin E]$$

$$M = \binom{n}{2}$$

Der zugehörige Index wird als *Performance* $\text{perf}(C)$ bezeichnet.

Einschränkung auf die Dichte

9/12

Einschränkung: Messe nur die Dichte der Cluster und nicht die Düntheit dazwischen.

$$f(\mathcal{C}) = \sum_{C \in \mathcal{C}} |E(C)| \quad g(\mathcal{C}) = 0 \quad \max = |E|$$

Der zugehörige Index wird als *Coverage* $\text{cov}(\mathcal{C})$ bezeichnet.

Einschränkung auf die Dichte

9/12

Einschränkung: Messe nur die Dichte der Cluster und nicht die Düntheit dazwischen.

$$f(\mathcal{C}) = \sum_{C \in \mathcal{C}} |E(C)| \quad g(\mathcal{C}) = 0 \quad \max = |E|$$

Der zugehörige Index wird als *Coverage* $\text{cov}(\mathcal{C})$ bezeichnet.

Lemma

Eine Clusterung \mathcal{C} eines Graphens G hat genau dann eine Coverage von 1, wenn die Cluster aus der disjunkten Vereinigung von Zusammenhangskomponenten bestehen.

Schnitte als Maß für Dichte

10/12

Definition

Sei $G = (V, E)$ ein Graph und $C' = (C_1, C_2)$ ein beliebiger Schnitt. Die *Conductance* (*Leitfähigkeit*) $\varphi(C')$ von C' ist definiert als:

$$\varphi(C') := \begin{cases} 1 & , \text{ falls } C_1 \in \{\emptyset, V\} \\ 0 & , \text{ falls } C_1 \notin \{\emptyset, V\}, \overline{E(C')} = 0 \\ \frac{|E(C')|}{|E(C')| + \min(|E(C_1)|, |E(C_2)|)} & , \text{ sonst} \end{cases} .$$

Die *Conductance* $\varphi(G)$ von G ist definiert als

$$\varphi(G) := \min_{C \subseteq V} \varphi((C, V \setminus C)) .$$

Conductance-Qualitätsmaß

11/12

Definition

Sei $G = (V, E)$ ein Graph und $\mathcal{C} = (C_1, \dots, C_k)$ eine Clustering von G . Der Index *Intra-Cluster Conductance* ist definiert durch:

$$f(\mathcal{C}) = \min_{1 \leq i \leq k} \varphi(G[C_i]) \quad g \equiv 0 \quad M = 1 .$$

und der Index *Inter-Cluster Conductance* ist gegeben durch:

$$M = 1 \quad f \equiv 0 \quad g = \begin{cases} 1 & , \text{ falls } \mathcal{C} = \{V\} \\ 1 - \max_{1 \leq i \leq k} \varphi((C_i, V \setminus C_i)) & , \text{ sonst} \end{cases} .$$

Definition

Sei $G = (V, E)$ ein Graph und $\mathcal{C} = (C_1, \dots, C_k)$ eine Clustering von G . Das Maß *Modularity* ist definiert durch:

$$\text{mod}(\mathcal{C}) := \sum_{i=1}^k \left(\frac{|E(C_i)|}{|E|} - \frac{1}{4|E|^2} \left(\sum_{v \in C_i} \text{deg}(v) \right)^2 \right)$$

Definition

Sei $G = (V, E)$ ein Graph und $\mathcal{C} = (C_1, \dots, C_k)$ eine Clustering von G . Das Maß *Modularity* ist definiert durch:

$$\text{mod}(\mathcal{C}) := \sum_{i=1}^k \left(\frac{|E(C_i)|}{|E|} - \frac{1}{4|E|^2} \left(\sum_{v \in C_i} \text{deg}(v) \right)^2 \right)$$

Es gilt: $\text{mod}(\mathcal{C}) \in [-0.5, 1]$.