

Beyond Spectral Clustering: A Comparative Study of Community Detection for Document Clustering

Kostadin Cvejosi

Fraunhofer Institute for
Intelligent Analysis and Information Systems IAIS



- embed the documents in some vector space
- define graph over the vector space
- apply the Overlap Hierarchical Algorithm

Graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$:

- \mathcal{V} - nodes; $\mathcal{V} = \{v_i\}_{i=1}^n$, where $n = |\mathcal{V}|$,
- $\mathcal{E} \subseteq (\mathcal{V} \times \mathcal{V})$ - edges

Subgraph $\mathcal{C} = (\mathcal{V}(\mathcal{C}), \mathcal{E}(\mathcal{C}))$ of a graph \mathcal{G} :

- $\mathcal{V}(\mathcal{C}) \subseteq \mathcal{V}$ - nodes; $n_c = |\mathcal{C}|$,
- $\mathcal{E}(\mathcal{C}) = \{(v_i, v_j) : v_i, v_j \in \mathcal{V}(\mathcal{C}) \wedge (v_i, v_j) \in \mathcal{E}\}$ - edges

Partition \mathcal{P} of a graph \mathcal{G} :

$$\mathcal{P} \equiv \{\mathcal{C}_1, \dots, \mathcal{C}_M\} \text{ such that } \bigcup_{i=1}^M \mathcal{V}(\mathcal{C}_i) = \mathcal{V}(\mathcal{G})$$

Overlap Hierarchical Algorithm (OH) - Notations

- **internal degree** ($k_v^{int}(\mathcal{C})$) of $v \in \mathcal{C}$ # of edges with nodes in \mathcal{C} ,
- **external degree** ($k_v^{ext}(\mathcal{C})$) of $v \in \mathcal{C}$ # of edges with nodes not in \mathcal{C} ,
- **internal(external) degree** of subgraph \mathcal{C} :

$$k_{\mathcal{C}}^{int(ext)} = \sum_{v \in \mathcal{C}} k_v^{int(ext)}(\mathcal{C})$$

- **fitness** f_C of a subgraph C

$$f_C = \frac{k_C^{int}}{(k_C^{int} + k_C^{ext})^\alpha}$$

- **node fitness** f_C^v with respect to a subgraph C :

$$f_C^v = f_{C+\{v\}}^v - f_{C-\{v\}}^v$$

Algorithm 1 Natural community of a node v_s

Require: Graph $\mathcal{G} \equiv (\mathcal{E}, \mathcal{V})$, v (seed node)

- 1: $\mathcal{C} \leftarrow \{v\}$
 - 2: **while** there exists a node $v \in \Gamma(\mathcal{C})$ such that $f_{\mathcal{C}}^v > 0$ **do**
 - 3: $u = \operatorname{argmax}_{\hat{u} \in \Gamma(\mathcal{C})} f(\hat{u})$
 - 4: $\mathcal{C} \leftarrow \mathcal{C} \cup \{u\}$
 - 5: **while** there exist $\tilde{u} \in \mathcal{C}$ such that $f_{\mathcal{C}}^{\tilde{u}} < 0$ **do**
 - 6: $\mathcal{C} \leftarrow \mathcal{C} \setminus \{\tilde{u}\}$
 - 7: **end while**
 - 8: **end while**
 - 9: **return** \mathcal{C}
-

$\Gamma(\mathcal{C}) = \{u | u \in \mathcal{V} \text{ and } u \notin \mathcal{C} \text{ and } (u, v) \in \mathcal{E}\}$ - neighbors of community \mathcal{C}

Performance Improvement of OH Algorithm

- inclusion new node in the subgraph

$$k_C^{ext} \leftarrow k_C^{ext} - k_g^{int} + k_g^{ext}$$

$$k_C^{int} \leftarrow 2((k_C^{int}/2) + k_g^{int})$$

- removing node from the subgraph

$$k_C^{ext} \leftarrow k_C^{ext} + k_t^{int} - k_t^{ext}$$

$$k_C^{int} \leftarrow 2((k_C^{int}/2) - k_t^{int})$$

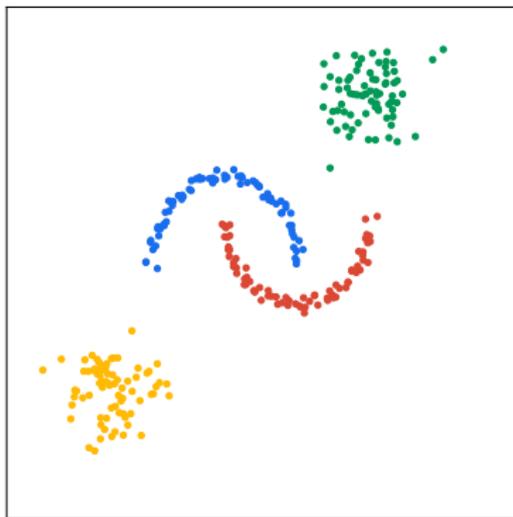
$$f(\mathcal{C}_i) = \frac{\text{vol}(\mathcal{C}_i) - W(\mathcal{C}_i, \bar{\mathcal{C}}_i)}{\text{vol}(\mathcal{C}_i)^\alpha} = \frac{1}{\text{vol}(\mathcal{C}_i)^{\alpha-1}} \left[1 - \frac{W(\mathcal{C}_i, \bar{\mathcal{C}}_i)}{\text{vol}(\mathcal{C}_i)} \right]$$

$$F_{\mathcal{P}} = \sum_i f(\mathcal{C}_i) = \sum_i \left[1 - \frac{W(\mathcal{C}_i, \bar{\mathcal{C}}_i)}{\text{vol}(\mathcal{C}_i)} \right]$$

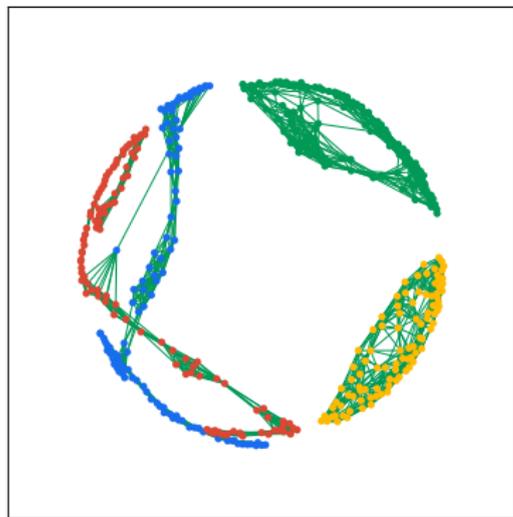
$$\min\{-F_{\mathcal{P}}\} = \min[2\text{cut}(\mathcal{P}) - M] = \min\{\text{cut}(\mathcal{P})\}$$

$$\text{cut}(\mathcal{P}) = \sum_i \frac{W(\mathcal{C}_i, \bar{\mathcal{C}}_i)}{\text{vol}(\mathcal{C}_i)}$$

Synthetic Data

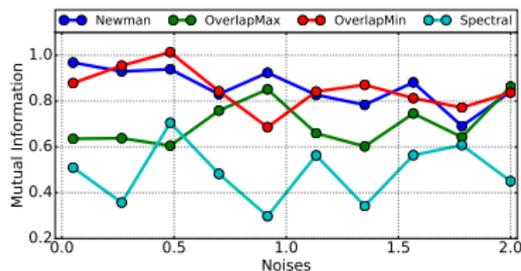


(a) spatial data

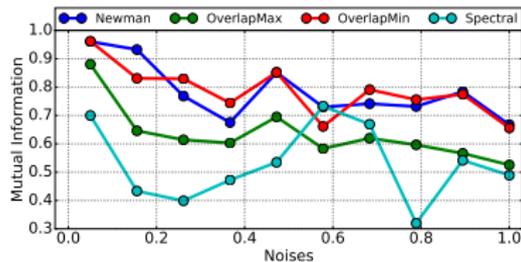


(b) associated k -NN graph

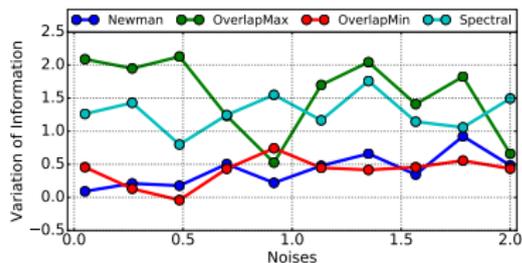
Synthetic Data



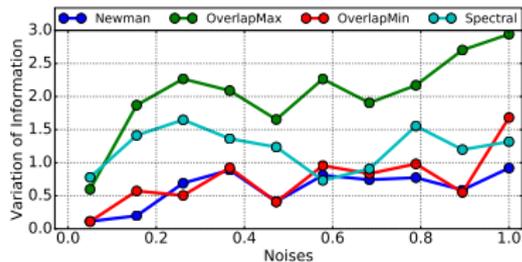
(a) Gaussian Cluster Dispersion



(c) Moon Cluster Dispersion



(b) Gaussian Cluster Dispersion



(d) Moon Cluster Dispersion

Clustering of News Articles

Com. 1: Barack Obama	
<p>"The Obama administration rules out raising medicaid eligibility age cat food stocks plummet", "This racist snowbilly was almost our vice president palin tweets that Obama is using shuck jive", "Photos and quotes from president Barack Obamas speech in cleveland ohio", "President Obama weighs in on Michigan's right to work for less attack on unions"</p>	
Subcom. 1.1: Republicans vs Obama	Subcom. 1.2: Obama Social Events
<p>"President Obama comes out of the gate powerfully against childish republicans mccain graham", "This racist snowbilly was almost our vice president palin tweets that Obama is using shuck jive", "The republican party is done at the national level"</p>	<p>"First official Obama rally of 2012 in columbus ohio", "Exclusive interactive panoramic images from the last night of the democratic national convention", "There is nothing more powerful than ordinary citizens coming together for a just cause", "Michelle Obamas 2012 dnc speech with photos and transcript"</p>

Clustering of News Articles

Com. 2: Obamacare	
<p>"Obamacare signups Friday 970k or 1.3m private enrollments 3.3m total", "It's not your freedom you're worried about theirs", "Debunking the republican lie that health insurance costs have skyrocketed under Obama", "Obamacare enrollment picks up as more Americans are getting covered", "Lowering budget with health insurance", "Republican Mike Shirkey working hard to solve an Obamacare problem that literally does not exist"</p>	
Subcom. 2.1: Republicans	Subcom. 2.2: Insurance
<p>"The president almost becomes his own anger translator over GOP Obamacare sabotage", "Republicans have now entered the hypocrisy zone on the Obamacare", "GOP Obamacare replacement shows conservatives have lost the health care battle"</p>	<p>"Funding for Michigan health centers will help the uninsured find coverage", "Obamacare day one my hideous experience", "More on GOP bill to require insurers to tell you how Obamacare is raising your rates the house republicans respond", "RNC fundraising email claims health insurance costs under Obama have gone up eight times higher than they have"</p>

Conclusion

- introduce new method of document clustering
- improve the performance of the OH algorithm
- equivalence to Ncut

- introduce new method of document clustering
- improve the performance of the OH algorithm
- equivalence to Ncut

Thank You