

Learning from Networks

Graph Clustering

Fabio Vandin

December 4th, 2024

Graph Clustering: Definition

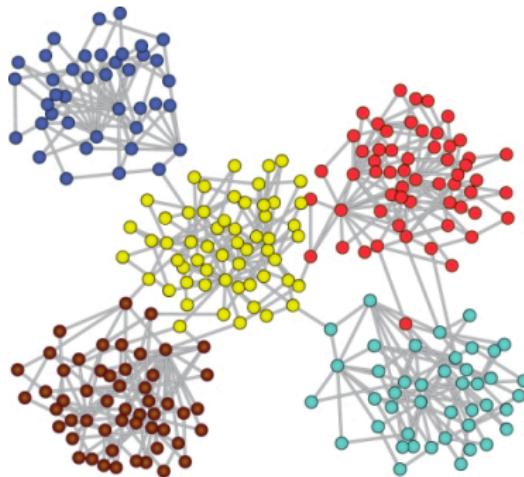
Given: graph $G = (V, E)$

Goal: partition V into clusters so that *similar vertices* are in the same cluster and *different vertices* are in different clusters.



Graph Clustering: Definition (continue)

Intuition: the similarity between vertices are represented by the edges



Given: connected graph $G = (V, E)$

Goal: partition V so that there are many edges within each cluster and few edges between clusters.

Many different formalizations based on this intuition.

Note: sometimes clusters in a graph are called *communities*

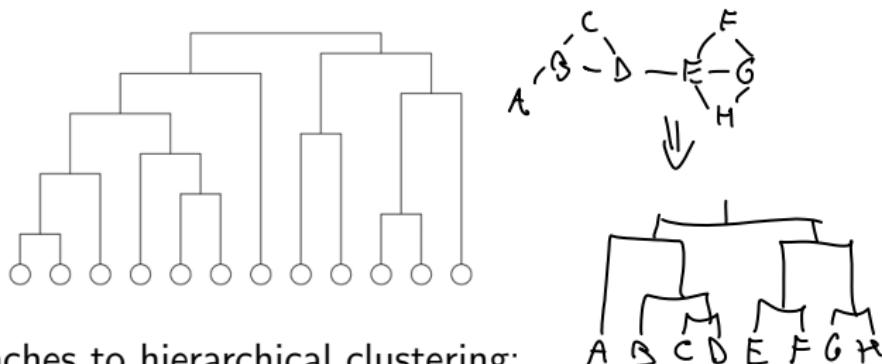
Graph Clustering: Approaches

We will see different types of approaches for clustering:

- hierarchical clustering
- cost-based clustering
- (spectral clustering)

Hierarchical Clustering

The output is a **dendrogram**, representing the clustering structure of the whole graph **G**.



Two general approaches to hierarchical clustering:

- **agglomerative approach:** start with each node in a cluster, iteratively join clusters \Rightarrow Ravasz algorithm
- **divisive approach:** start with all nodes in a cluster, iteratively split clusters \Rightarrow Girvan-Newman algorithm

Ravasz Algorithm

Algorithm AgglomerativeClustering(G)

Input: connected graph $G = (V, E)$

Output: dendrogram whose leaves are the elements of V

- 1 assign each node u to its own cluster C_u ;
- 2 for all pairs $u, v \in V, u \neq v$: compute their similarity $\text{sim}(u, v)$
- 3 repeat until all nodes are in a single cluster:
 - 4 find the pair of clusters C_1, C_2 with highest similarity $\text{sim}(C_1, C_2)$ (ties broken arbitrarily)
 - 5 merge clusters C_1, C_2 in a single cluster C'
 - 6 compute similarity between C' and all other clusters
- 7 return the corresponding dendrogram

Different variants depending on the definition of $\text{sim}(u, v)$ and the definition of $\text{sim}(C_1, C_2)$.

Complexity? In general: $\Theta(|V|^2)$ computations of $\text{sim}(u, v)$ and of $\text{sim}(C_1, C_2)$

Ravasz Algorithm (continue)

Common choice for $\text{sim}(u, v)$:

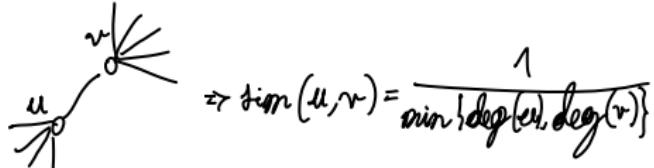
$$\text{sim}(u, v) = \frac{|\mathcal{N}(u) \cap \mathcal{N}(v)| + A_{uv}}{\min\{\deg(u), \deg(v)\} + 1 - A_{uv}}$$

where A is the adjacency matrix of G

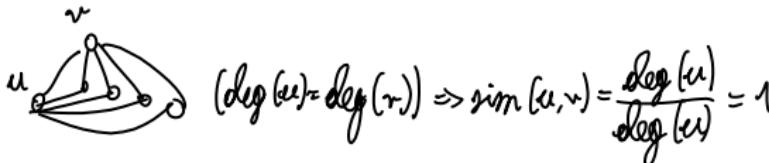
a)



b)

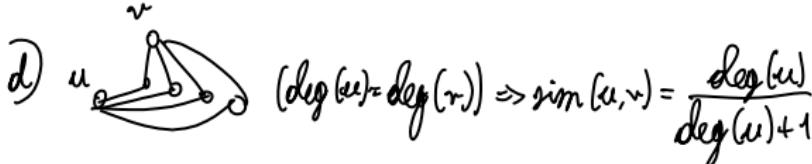


c)



$$(\deg(u) = \deg(v)) \Rightarrow \text{sim}(u, v) = \frac{\deg(u)}{\deg(u)} = 1$$

d)



$$(\deg(u) = \deg(v)) \Rightarrow \text{sim}(u, v) = \frac{\deg(u)}{\deg(u)+1}$$

Ravasz Algorithm (continue)

Common choices for $\text{sim}(C_1, C_2)$ define different types of linkage clustering:

- single linkage clustering: $\text{sim}(C_1, C_2) = \min_{u \in C_1, v \in C_2} \text{sim}(u, v)$

- average linkage clustering:

$$\text{sim}(C_1, C_2) = \frac{1}{|C_1| \cdot |C_2|} \sum_{u \in C_1, v \in C_2} \text{sim}(u, v)$$

- complete linkage clustering: $\text{sim}(C_1, C_2) = \max_{u \in C_1, v \in C_2} \text{sim}(u, v)$

Example

Girvan-Newman Algorithm

Based on the idea of iteratively removing the most *central* edge in the graph $G = (V, E)$.



Various definitions of *centrality* for edges, but the most common one is *link betweenness*.

Link betweenness

Let $\sigma_{s,t}$ be the number of shortest paths from node s to node t .

Let $\sigma_{s,t}(e)$ be the number of shortest paths from node s to node t that pass through edge e .

Definition

Given a connected graph $G = (V, E)$ and an edge $e \in E$ the link betweenness $b(e, G)$ of e in G :

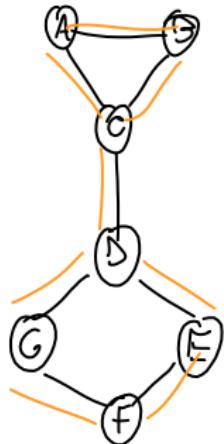
$$b(e, G) = \sum_{s,t \in V: s \neq t} \frac{\sigma_{s,t}(e)}{\sigma_{s,t}}$$

Complexity of computing $b(e, G)$ for all edges $e \in E$?

$$\Theta(|V| \cdot |E|)$$

Example

6:



e	$f(e, G)$
A, B	1
A, C	
B, C	
C, D	12 \Rightarrow (
D, E	
E, F	3. S

12 \Rightarrow (nodo in triangolo) \Rightarrow (nodo in diamante)

Girvan-Newman Algorithm (continue)

Algorithm GNClustering(G)

Input: connected graph $G = (V, E)$

Output: dendrogram whose leaves are the elements of V

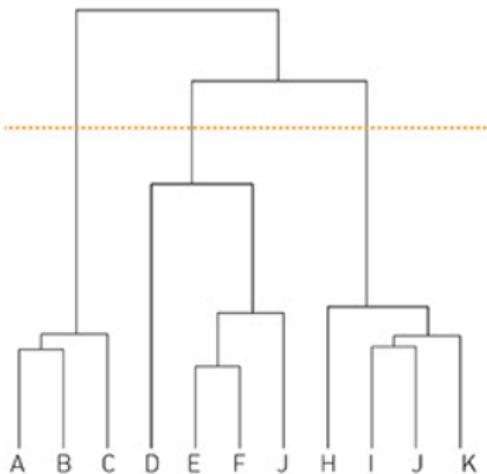
- 1 assign all nodes u to a single cluster C ;
- 2 repeat until all nodes are in different clusters: *m iterations*
 - 3 for each cluster C :
 - 4 for each edge $e \in C$: compute $b(e, C)$
 - 5 let e_{\max} the edge of maximum betweenness, and let $C(e)$ its cluster;
 - 6 remove e from $C(e)$;
 - 7 report the corresponding dendrogram

Complexity? In general: $\Theta(|E|^2|V|)$.

$$\frac{|E|^2}{m} |V|$$

Hierarchical Clustering: Getting a Clustering

The output of hierarchical clustering is a dendrogram, not a clustering. How do we obtain a clustering?

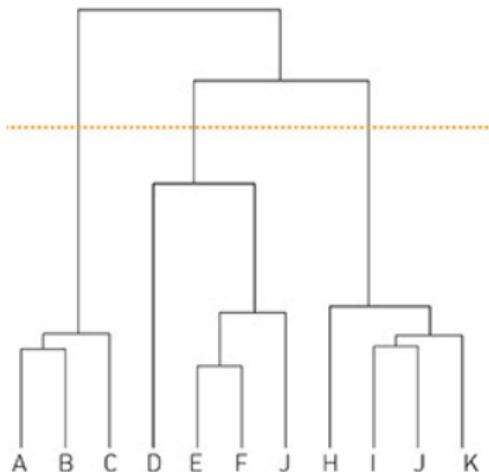


By cutting the dendrogram at a given level.

How do we select where to cut?

Hierarchical Clustering: Getting a Clustering (continue)

How do we select where to cut?



- if we know the number k of clusters we want: pick a level resulting in k clusters
- if we do not know k : define a score for clusterings, and pick the clustering from the dendrogram of maximum score

Cost-based Clustering

Common approach in clustering (not only for graphs):

- define a cost function over possible partitions of the objects
- find the partition (=clustering) of minimal cost

Modularity

Idea: a cluster should contain more edges than expected in a random graph.

Definition

Given a graph $G = (V, E)$ with $|V| = n$, $|E| = m$ the modularity $M(S)$ of a subset $S \subseteq V$ of the vertices of G is

$$M(S) = \frac{1}{2m} \sum_{u,v \in S} \left(A_{uv} - \frac{\deg(u)\deg(v)}{2m} \right)$$

Intuition: measures the difference between the number of edges within each cluster with the expected number of edges under the Chung-Lu model for random graphs.

Modularity (continue)

The modularity of a clustering of G is the sum of the modularity of each cluster.

Definition

Given a clustering $\mathcal{C} = C_1, C_2, \dots$ of graph $G = (V, E)$ with $|V| = n$, $|E| = m$, the modularity $M(\mathcal{C})$ of \mathcal{C} is:

$$\begin{aligned} M(\mathcal{C}) &= \sum_{C \in \mathcal{C}} M(C) \\ &= \frac{1}{2m} \sum_{C \in \mathcal{C}} \sum_{u, v \in C} \left(A_{uv} - \frac{\deg(u)\deg(v)}{2m} \right) \end{aligned}$$

Modularity (continue)

Proposition

Given a clustering $\mathcal{C} = C_1, C_2, \dots$ of graph $G = (V, E)$ with $|V| = n$, $|E| = m$, the modularity $M(\mathcal{C})$ of \mathcal{C} is equal to

$$M(\mathcal{C}) = \sum_{C \in \mathcal{C}} \left(\frac{|E(C)|}{m} - \left(\frac{\sum_{u \in C} \deg(u)}{2m} \right)^2 \right)$$

where $E(C)$ are the edges between nodes in cluster C :

$$E(C) = \{(u, v) \in E : u \in C, v \in C\}$$

ci serve solo sapere $|E(C)|$ e iterare su nodi (non copie di nodi)

Proof

$$\text{Diam: } M(C) = \frac{1}{2m} \sum_{\text{cel}} \sum_{u, v \in C} \left(A_{uv} - \frac{\deg(u)\deg(v)}{2m} \right) =$$

$$\frac{1}{2m} \sum_{\text{cel}} \sum_{u, v \in C} A_{uv} = \frac{1}{2m} \sum_{\text{cel}} |E(C)| = \frac{1}{m} \sum_{\text{cel}} |E(C)|$$

$$\frac{1}{2m} \sum_{\text{cel}} \sum_{u, v \in C} \frac{\deg(u)\deg(v)}{2m} = \frac{1}{(2m)^2} \sum_{\text{cel}} \sum_{u, v \in C} \deg(u)\deg(v) =$$

$$= \frac{1}{(2m)^2} \sum_{\text{cel}} \left(\sum_{u \in C} \deg(u) \left(\sum_{v \in C} \deg(v) \right) \right) =$$

$$= \frac{1}{(2m)^2} \sum_{\text{cel}} \left(\left(\sum_{u \in C} \deg(u) \right) \left(\sum_{v \in C} \deg(v) \right) \right) = \frac{1}{(2m)^2} \sum_{\text{cel}} \left(\sum_{u \in C} \deg(u) \right)^2 =$$

$$= \sum_{\text{cel}} \left(\frac{\sum_{u \in C} \deg(u)}{2m} \right)^2$$

= QED ■

Example

Modularity-Based Clustering

Input: graph $G = (V, E)$

Goal: find the clustering $\mathcal{C} = C_1, C_2, \dots$ that maximizes the modularity

$$M(\mathcal{C}) = \sum_{C \in \mathcal{C}} \left(\frac{|E(C)|}{m} - \left(\frac{\sum_{u \in C} \deg(u)}{2m} \right)^2 \right)$$

Equivalent formulation: since the cost of clustering \mathcal{C} is $-M(\mathcal{C})$, the following formulation is equivalent:

Input: graph $G = (V, E)$

Goal: find the clustering $\mathcal{C} = C_1, C_2, \dots$ of minimum cost $-M(\mathcal{C})$.

Modularity-Based Clustering: Complexity

Informal: finding a clustering of maximum modularity is hard!

Problem **Modularity Clustering Problem**)

Given a graph G and a value K is there a clustering C of G such that $M(C) \geq K$?

Proposition

The Modularity Clustering Problem is NP-complete.

So? (Greedy) agglomerative algorithm

Modularity-Based Clustering: Greedy Agglomerative Approach

Algorithm GreedyModularityClustering(G)

Input: connected graph $G = (V, E)$

Output: clustering of the elements of V

- 1 $\mathcal{C}_1 \leftarrow$ clustering where each node u is assigned to its own cluster C_u ; $i \leftarrow 1$;
- 2 repeat until all nodes are in a single cluster:
 - 3 for each pair of clusters C_1, C_2 such that there exists one edge between C_1 and C_2 : compute
$$\delta(C_i, C_1, C_2) = M(C_i - C_1 - C_2 + (C_1 \cup C_2)) - M(C_i);$$
 - 4 find C', C'' that maximize $\delta(C_i, C', C'')$
 - 5 $\mathcal{C}_{i+1} \leftarrow \mathcal{C}_i - C' - C'' + (C' \cup C'')$; $i \leftarrow i + 1$;
- 6 return the clustering \mathcal{C}^* , across iterations, of maximum modularity: $\mathcal{C}^* = \arg \max_{\mathcal{C}_i, i=1,2,\dots} M(\mathcal{C}_i)$

Complexity? In general: $O(|E| \cdot |V|)$ computations of
 $\delta(C_i, C_1, C_2)$

Modularity-Based Clustering: Efficient Computation

Proposition

Let $E(C_1, C_2)$ be the edges between cluster C_1 and cluster C_2 :
 $E(C_1, C_2) = \{(u, v) \in E : u \in C_1, v \in C_2\}$. Then

$$\delta(C_i, C_1, C_2) = \frac{|E(C_1, C_2)|}{m} - \frac{(\sum_{u \in C_1} \deg(u)) (\sum_{v \in C_2} \deg(v))}{2m^2}$$

Dim: per definizione:

$$\begin{aligned}\delta(C_i, C_1, C_2) &= M(\overbrace{C_i \setminus C_1 \setminus C_2 \cup (C_1 \cup C_2)}^{C_i}) - M(C_i) = \\ &= \left(\sum_{c \in C_i} M(c) \right) - \left(\sum_{c \in C_i} M(c) \right) = \text{(cluster che non sono } C_1, C_2 \text{ si cancellano)} \\ &= M(C_1 \cup C_2) - M(C_1) - M(C_2) =\end{aligned}$$

$$= \frac{|\mathbb{E}(C_1 \cup C_2)|}{m} - \left(\frac{\sum_{v \in C_1 \cup C_2} \deg(v)}{2m} \right)^2 - \left(\frac{|\mathbb{E}(C_1)|}{m} - \left(\frac{\sum_{v \in C_1} \deg(v)}{2m} \right)^2 \right) -$$

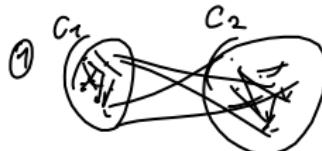
$$- \left(\frac{|\mathbb{E}(C_2)|}{m} - \left(\frac{\sum_{v \in C_2} \deg(v)}{2m} \right)^2 \right) =$$

(1)

$$= \frac{|\mathbb{E}(C_1 \cup C_2)| - |\mathbb{E}(C_1)| - |\mathbb{E}(C_2)|}{m} - \left\{ \left(\frac{\sum_{v \in C_1} \deg(v) + \sum_{v \in C_2} \deg(v)}{2m} \right)^2 + \right.$$

$$\left. + \left(\frac{\sum_{v \in C_1} \deg(v)}{2m} \right)^2 + \left(\frac{\sum_{v \in C_2} \deg(v)}{2m} \right)^2 \right\} =$$

(2)



$$\Rightarrow = \frac{|\mathbb{E}(C_1, C_2)|}{m}$$

$$(2) - (a+b)^2 + a^2 + b^2 = -2ab$$

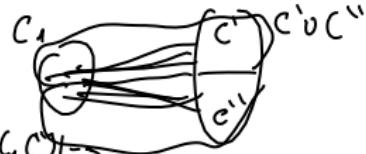
$$= \frac{|\mathbb{E}(C_1, C_2)|}{m} - \frac{2(\sum_{v \in C_1} \deg(v))(\sum_{v \in C_2} \deg(v))}{2 \cdot m^2}$$

Modularity-Based Clustering: Efficient Computation (continue)

Proposition

In every iteration of the repeat-until loop, the values $|E(C_1, C_2)|$ for all $C_1, C_2 \in \mathcal{C}$ and $\sum_{u \in C} \deg(u)$ for all $C \in \mathcal{C}$ can be efficiently updated in total time $O(|E|)$.

- Dim:
- inizio: $\forall u \in V$, se cluster $C_u \Rightarrow |E(C_1, C_2)|$ dipende dalla presenza singolo bordo $\Rightarrow \sum_{u \in C} \deg(u)$ calcolabile in tempo totale $\Theta(\sum_{u \in V} \deg(u)) = \Theta(m)$
 - suppongo valori calcolati questi prima di inizio it. \Rightarrow qui, C', C'' sostituiti da $C' \cup C'' \Rightarrow \forall$ paio C_1, C_2 che non includono $C' \cup C''$, no calcoli \Rightarrow vendo coppia $C_1, C' \cup C''$:
 $E(C_1, C' \cup C'') = E(C_1, C') \cup E(C_1, C'')$,
 $E(C_2, C') \cap E(C_2, C'') = \emptyset \Rightarrow |E(C_1, C' \cup C'')| = |E(C_1, C')| + |E(C_1, C'')| \Rightarrow$
 $\Rightarrow \leq |V|$ copie di questo tipo \Rightarrow computation time $O(|V|) \in O(|E|)$



\forall cluster $C \neq C_i \cup C'$, $\sum_{u \in C} \deg(u)$ già calcolato

dato $C \cup C'$: $\sum_{u \in C \cup C'} \deg(u) = \sum_{u \in C} \deg(u) + \sum_{u \in C'} \deg(u) \Rightarrow$
 \Rightarrow valori già disponibili \Rightarrow tempo $\Theta(1)$ ■

Example

