

**Script** generated by TTT

Title: groh: profile1 (13.05.2014)

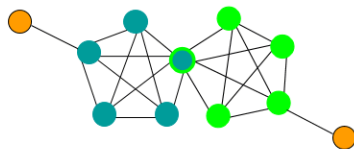
Date: Tue May 13 11:59:57 CEST 2014

Duration: 89:39 min

Pages: 83

- A subset  $U \subseteq V$  of a Graph  $(V, E)$  is a **clique** if  $G([U])$  is a **complete** graph;  $G([U])$  is the sub-graph induced by  $U$ .
- A clique is **maximal** if there is no clique  $U'$  with  $U \subset U'$  in  $G$
- A clique is a **maximum** clique if there is no clique with more vertices in  $G$

- If Graph is dense  $\rightarrow$  cliques exist:
- (Turan 1941 (see [2])): If  $|E| > |V|^2/2 \cdot (k-2)/(k-1)$  then  $G$  contains a **clique of size  $k$** .
- Usually many different maximal cliques exist in a graph; they can even overlap without being identical



- Cliques are very “strict”  $\rightarrow$  **Alternative candidates** for groups:  
Distance based structures:
  - $U$  is **N-clique** iff  $\forall u, v \in U : \text{dist}_G(u, v) \leq N$  (non-local def.!)
    - $U$  is **N-club** iff  $\text{diam}(G([U])) \leq N$
    - $U$  is **N-clan** iff  $U$  is maximal N-clique and  $\text{diam}(G([U])) \leq N$
- **Criticisms:**
  - Since dist is evaluated w.r.t. to  $G$  and not  $G([U])$  (thus N-cliques are not local structures), **N-cliques need not even be connected and can have a diameter  $\text{diam}(G([U]) > N$**

• Cliques are very “strict” → **Alternative candidates** for groups:  
**Distance based structures:**

- U is **N-clique** iff  $\forall u,v \in U : \text{dist}_G(u,v) \leq N$  (non-local def.!)
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$

• **Criticisms:**

- Since dist is evaluated w.r.t. to G and not G([U]) (thus N-cliques are not local structures), **N-cliques need not even be connected and can have a diameter  $\text{diam}(G([U]) > N$**



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$

• → **N-clan**: restrict dist-condition to paths of nodes **within** the structure: easy to find (just drop all n-cliques with diameter greater than N)

• → **N-club**: regard all induced graphs with diameter less than N: harder to find

- It can be shown / seen from the def.:
  - all N-clans are N-cliques;
  - all N-clubs are contained within N-cliques;
  - all N-clans are n-clubs
  - there are N-clubs that are not N-clans



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$

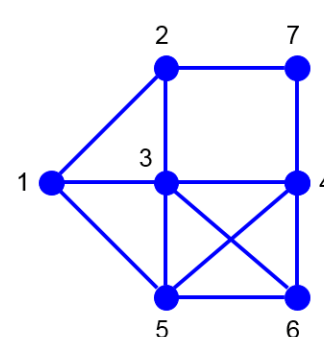
• → **N-clan**: restrict dist-condition to paths of nodes **within** the structure: easy to find (just drop all n-cliques with diameter greater than N)

• → **N-club**: regard all induced graphs with diameter less than N: harder to find

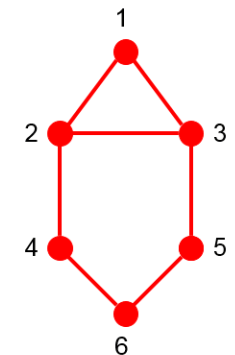
- It can be shown / seen from the def.:
  - all N-clans are N-cliques;
  - all N-clubs are contained within N-cliques;
  - all N-clans are n-clubs
  - there are N-clubs that are not N-clans



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



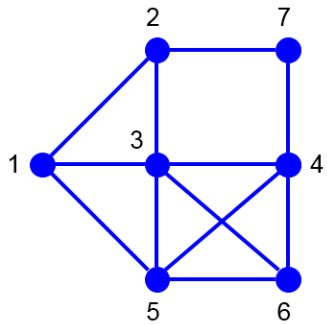
cliques: {1, 2, 3}, {1, 3, 5}, {3, 4, 5, 6}



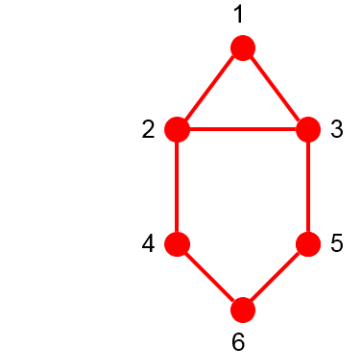
2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
 2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
 2-clan: {2, 3, 4, 5, 6}



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



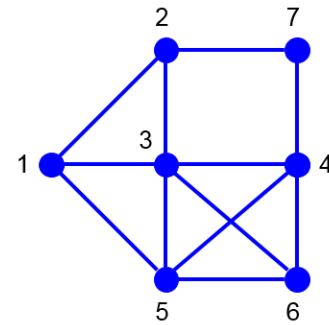
cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}



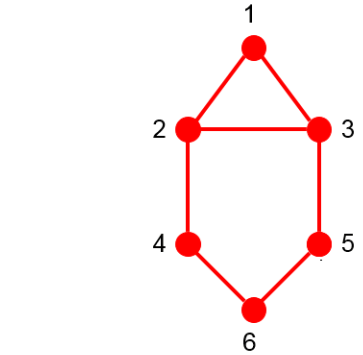
2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



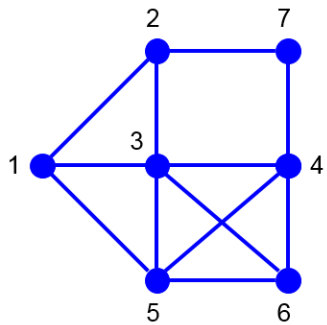
cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}



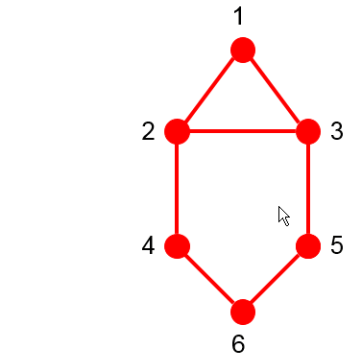
2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



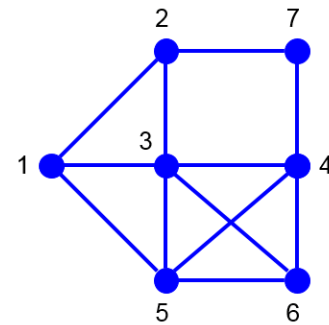
cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}



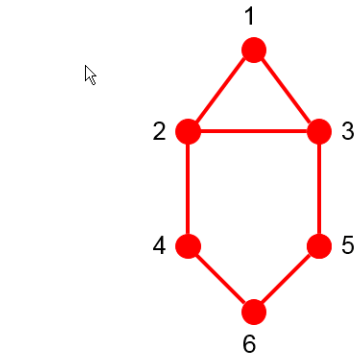
2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}



- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}

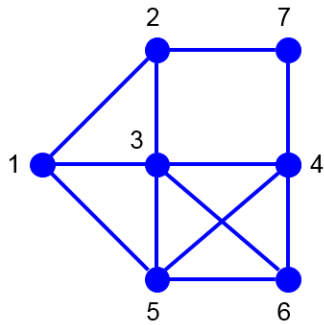


2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}

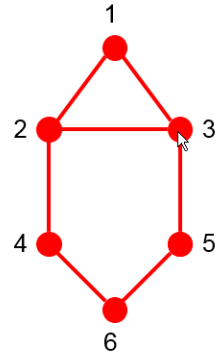


# N-Cliques, N-Clubs, N-Clans

- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



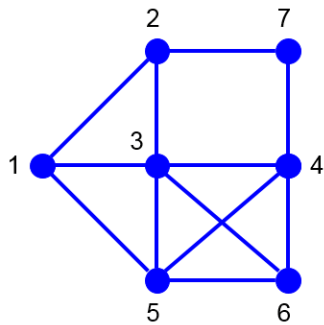
cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}



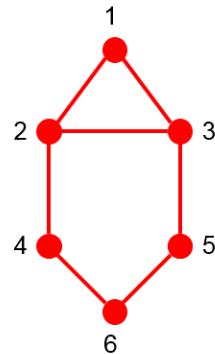
2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}

# N-Cliques, N-Clubs, N-Clans

- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



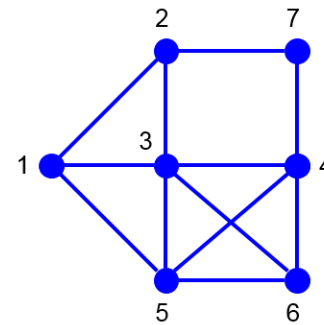
cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}



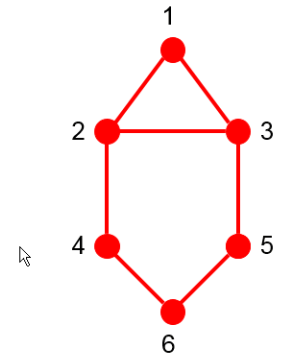
maximal 2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}

# N-Cliques, N-Clubs, N-Clans

- U is **N-clique** iff  $\forall u,v \in V : \text{dist}_G(u,v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$



cliques: {1, 2, 3}, {1, 3, 5},  
{3, 4, 5, 6}



maximal 2-cliques: {1, 2, 3, 4, 5}, {2, 3, 4, 5, 6}  
2-clubs: {1, 2, 3, 4}, {1, 2, 3, 5}, {2, 3, 4, 5, 6}  
2-clan: {2, 3, 4, 5, 6}

- U is **N-clique** iff  $\forall u, v \in V : \text{dist}_G(u, v) \leq N$
- U is **N-club** iff  $\text{diam}(G([U])) \leq N$
- U is **N-clan** iff U is maximal N-clique and  $\text{diam}(G([U])) \leq N$

- → **N-clan**: restrict dist-condition to paths of nodes **within** the structure: easy to find (just drop all n-cliques with diameter greater than N)
- → **N-club**: regard all induced graphs with diameter less than N: harder to find
- It can be shown / seen from the def.:
  - all N-clans are N-cliques;
  - all N-clubs are contained within N-cliques;
  - all N-clans are n-clubs
  - there are N-clubs that are not N-clans



## Cliques: Algorithms

- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of U if they are in E. Although up to  $\binom{|V|}{2}$  such pairs may exist: If |E| edges have been searched, search is over)
  - **Determine if clique U is maximal** (Test all vertices in V-U if they are connected to all vertices in U; again: If |E| edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing U: Assume vertices are sorted;** Test all vertices in V-U in ascending order: if they are connected to all vertices in U, add to U; again: If |E| edges have been searched, search is over)



- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of U if they are in E. Although up to  $\binom{|V|}{2}$  such pairs may exist: If |E| edges have been searched, search is over)
  - **Determine if clique U is maximal** (Test all vertices in V-U if they are connected to all vertices in U; again: If |E| edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing U: Assume vertices are sorted;** Test all vertices in V-U in ascending order: if they are connected to all vertices in U, add to U; again: If |E| edges have been searched, search is over)



## Cliques: Algorithms

- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of U if they are in E. Although up to  $\binom{|V|}{2}$  such pairs may exist: If |E| edges have been searched, search is over)
  - **Determine if clique U is maximal** (Test all vertices in V-U if they are connected to all vertices in U; again: If |E| edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing U: Assume vertices are sorted;** Test all vertices in V-U in ascending order: if they are connected to all vertices in U, add to U; again: If |E| edges have been searched, search is over)



## Cliques: Algorithms

- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of  $U$  if they are in  $E$ . Although up to  $\binom{|V|}{2}$  such pairs may exist: If  $|E|$  edges have been searched, search is over)
  - **Determine if clique  $U$  is maximal** (Test all vertices in  $V-U$  if they are connected to all vertices in  $U$ ; again: If  $|E|$  edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing  $U$ : Assume vertices are sorted**; Test all vertices in  $V-U$  in ascending order: if they are connected to all vertices in  $U$ , add to  $U$ ; again: If  $|E|$  edges have been searched, search is over)



## Cliques: Algorithms

- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of  $U$  if they are in  $E$ . Although up to  $\binom{|V|}{2}$  such pairs may exist: If  $|E|$  edges have been searched, search is over)
  - **Determine if clique  $U$  is maximal** (Test all vertices in  $V-U$  if they are connected to all vertices in  $U$ ; again: If  $|E|$  edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing  $U$ : Assume vertices are sorted**; Test all vertices in  $V-U$  in ascending order: if they are connected to all vertices in  $U$ , add to  $U$ ; again: If  $|E|$  edges have been searched, search is over)



## Cliques: Algorithms

- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of  $U$  if they are in  $E$ . Although up to  $\binom{|V|}{2}$  such pairs may exist: If  $|E|$  edges have been searched, search is over)
  - **Determine if clique  $U$  is maximal** (Test all vertices in  $V-U$  if they are connected to all vertices in  $U$ ; again: If  $|E|$  edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing  $U$ : Assume vertices are sorted**; Test all vertices in  $V-U$  in ascending order: if they are connected to all vertices in  $U$ , add to  $U$ ; again: If  $|E|$  edges have been searched, search is over)



## Cliques: Algorithms

- In connection with cliques: Algorithms with **time complexity**  $O(|E| + |V|)$  for the problems:
  - **Determine if  $U \subseteq V$  is a clique** (Test pairs of vertices of  $U$  if they are in  $E$ . Although up to  $\binom{|V|}{2}$  such pairs may exist: If  $|E|$  edges have been searched, search is over)
  - **Determine if clique  $U$  is maximal** (Test all vertices in  $V-U$  if they are connected to all vertices in  $U$ ; again: If  $|E|$  edges have been searched, search is over)
  - **Compute lexicographically smallest maximal clique containing  $U$ : Assume vertices are sorted**; Test all vertices in  $V-U$  in ascending order: if they are connected to all vertices in  $U$ , add to  $U$ ; again: If  $|E|$  edges have been searched, search is over)



## Cliques: Algorithms

- **Naive algorithm for finding the maximum clique:** Exhaustive Search: Compute all subsets of  $V$  and check for clique  $\rightarrow O(|V|^2 2^M)$  worst case
- **Is there a substantially better algorithm? Probably not:** The decision problem  $\text{CLIQUE}(G,k)$ : “Has  $G$  a clique of size at least  $k$ ” is **NP complete**. (Solving the decision problem in time  $T(|V|)$  would yield an alg for determining the maximal  $k$  in  $O(T(|V|) \log(|V|))$  via binary search.

Proof: Reduce  $\text{CLIQUE}$  on  $\text{SATISFIABILITY}$ ;

- $\rightarrow$  Unless  $P=NP$  there is **no P algorithm  $X$  to compute clique of size  $k$**  if  $G$  is guaranteed to contain such a clique.

Proof: Having  $X$  we could decide  $\text{CLIQUE}$  in  $P$  time;



## Cliques: Algorithms

- **Naive algorithm for finding the maximum clique:** Exhaustive Search: Compute all subsets of  $V$  and check for clique  $\rightarrow O(n^2 2^n)$
- **Is there a slightly better algorithm? Yes:** We can improve the exponential function’s base from 2 to approx. 1.38 (see [2])
- Can we expect to find the answer to the problem “how many cliques with exactly  $k$  nodes exist”?  $\rightarrow$  Exhaustive search:  $\Theta(|V|^k)$
- For triangles we can be better than  $\Theta(|V|^3)$ : An alg with  $O(|V|^{2.376})$  exists 😊
- For other  $k$ : An analogous technique allows for alg with  $O(|V|^{\beta(k)})$  with  $\beta(k) < k$



## Cliques: Algorithms

- **Naive algorithm for finding the maximum clique:** Exhaustive Search: Compute all subsets of  $V$  and check for clique  $\rightarrow O(n^2 2^n)$
- **Is there a slightly better algorithm? Yes:** We can improve the exponential function’s base from 2 to approx. 1.38 (see [2])
- Can we expect to find the answer to the problem “how many cliques with exactly  $k$  nodes exist”?  $\rightarrow$  Exhaustive search:  $\Theta(|V|^k)$
- For triangles we can be better than  $\Theta(|V|^3)$ : An alg with  $O(|V|^{2.376})$  exists 😊
- For other  $k$ : An analogous technique allows for alg with  $O(|V|^{\beta(k)})$  with  $\beta(k) < k$



## Cliques: Algorithms

- **Naive algorithm for finding the maximum clique:** Exhaustive Search: Compute all subsets of  $V$  and check for clique  $\rightarrow O(n^2 2^n)$
- **Is there a slightly better algorithm? Yes:** We can improve the exponential function’s base from 2 to approx. 1.38 (see [2])
- Can we expect to find the answer to the problem “how many cliques with exactly  $k$  nodes exist”?  $\rightarrow$  Exhaustive search:  $\Theta(|V|^k)$
- For triangles we can be better than  $\Theta(|V|^3)$ : An alg with  $O(|V|^{2.376})$  exists 😊
- For other  $k$ : An analogous technique allows for alg with  $O(|V|^{\beta(k)})$  with  $\beta(k) < k$



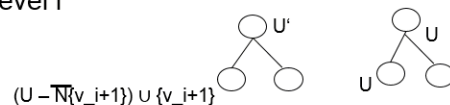
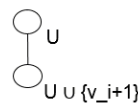
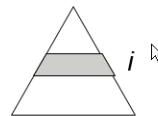
- **Naive algorithm for finding the maximum clique:** Exhaustive Search: Compute all subsets of  $V$  and check for clique  $\rightarrow O(n^2 2^n)$
- **Is there a slightly better algorithm? Yes:** We can improve the exponential function's base from 2 to approx. 1.38 (see [2])
- Can we expect to find the answer to the problem "how many cliques with exactly  $k$  nodes exist"?  $\rightarrow$  Exhaustive search:  $\Theta(|V|^k)$
- For triangles we can be better than  $\Theta(|V|^3)$ : An alg with  $O(|V|^{2.376})$  exists 😊
- For other  $k$ : An analogous technique allows for alg with  $O(|V|^{\beta(k)})$  with  $\beta(k) < k$

- **Naive algorithm for finding the maximum clique:** Exhaustive Search: Compute all subsets of  $V$  and check for clique  $\rightarrow O(n^2 2^n)$
- **Is there a slightly better algorithm? Yes:** We can improve the exponential function's base from 2 to approx. 1.38 (see [2])
- Can we expect to find the answer to the problem "how many cliques with exactly  $k$  nodes exist"?  $\rightarrow$  Exhaustive search:  $\Theta(|V|^k)$
- For triangles we can be better than  $\Theta(|V|^3)$ : An alg with  $O(|V|^{2.376})$  exists 😊
- For other  $k$ : An analogous technique allows for alg with  $O(|V|^{\beta(k)})$  with  $\beta(k) < k$



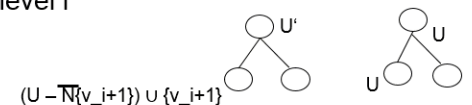
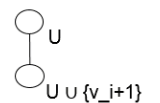
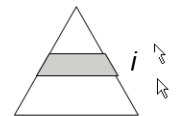
Cliques: Enumeration

- **Binary tree**,  $n$  levels, leaves only at level  $n$
- each level  $i \leftrightarrow$  vertex  $v_i$
- nodes at level  $i$ : maximal cliques in  $G[v_1, v_2, \dots, v_i]$
- level  $i+1$ : determine children of node  $U$  at level  $i$ : two cases :
  - (1)  $v_{i+1}$  adjacent to all nodes in  $U$  ( $U \subseteq N(v_{i+1})$ ):
    - $\rightarrow U \cup \{v_{i+1}\}$  is maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
    - $\rightarrow U \cup \{v_{i+1}\}$  is only child of  $U$
  - (2)  $\exists$  vertex in  $U$  not adjacent to  $v_{i+1}$  ( $U \not\subseteq N(v_{i+1})$ ):
    - $\rightarrow U$  itself is a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
    - $\rightarrow$  if  $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal it is also a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
    - $\rightarrow (U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is potentially child of many possible nodes  $\rightarrow$  define it als child of lexicographically smallest clique (node)  $U'$  at level  $i$



Cliques: Enumeration

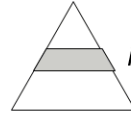
- **Binary tree**,  $n$  levels, leaves only at level  $n$
- each level  $i \leftrightarrow$  vertex  $v_i$
- nodes at level  $i$ : maximal cliques in  $G[v_1, v_2, \dots, v_i]$
- level  $i+1$ : determine children of node  $U$  at level  $i$ : two cases :
  - (1)  $v_{i+1}$  adjacent to all nodes in  $U$  ( $U \subseteq N(v_{i+1})$ ):
    - $\rightarrow U \cup \{v_{i+1}\}$  is maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
    - $\rightarrow U \cup \{v_{i+1}\}$  is only child of  $U$
  - (2)  $\exists$  vertex in  $U$  not adjacent to  $v_{i+1}$  ( $U \not\subseteq N(v_{i+1})$ ):
    - $\rightarrow U$  itself is a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
    - $\rightarrow$  if  $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal it is also a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
    - $\rightarrow (U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is potentially child of many possible nodes  $\rightarrow$  define it als child of lexicographically smallest clique (node)  $U'$  at level  $i$



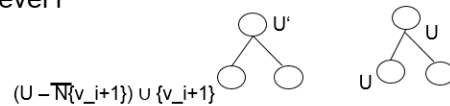
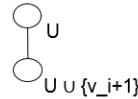


## Cliques: Enumeration

- Binary tree, n levels, leaves only at level n
- each level  $i \leftrightarrow$  vertex  $v_i$
- nodes at level  $i$ : maximal cliques in  $G[v_1, v_2, \dots, v_i]$
- level  $i+1$ : determine children of node  $U$  at level  $i$ : two cases :

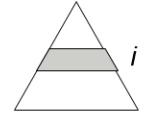


- (1)  $v_{i+1}$  adjacent to all nodes in  $U$  ( $U \subseteq N(v_{i+1})$ ):
  - $U \cup \{v_{i+1}\}$  is maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - $U \cup \{v_{i+1}\}$  is only child of  $U$
- (2)  $\exists$  vertex in  $U$  not adjacent to  $v_{i+1}$  ( $U \not\subseteq N(v_{i+1})$ ):
  - $U$  itself is a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - if  $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal it is also a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is potentially child of many possible nodes → define it als child of lexicographically smallest clique (node)  $U'$  at level  $i$

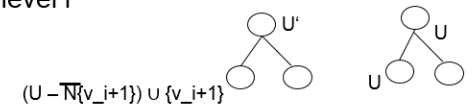
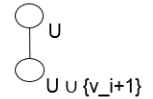


## Cliques: Enumeration

- Binary tree, n levels, leaves only at level n
- each level  $i \leftrightarrow$  vertex  $v_i$
- nodes at level  $i$ : maximal cliques in  $G[v_1, v_2, \dots, v_i]$
- level  $i+1$ : determine children of node  $U$  at level  $i$ : two cases :



- (1)  $v_{i+1}$  adjacent to all nodes in  $U$  ( $U \subseteq N(v_{i+1})$ ):
  - $U \cup \{v_{i+1}\}$  is maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - $U \cup \{v_{i+1}\}$  is only child of  $U$
- (2)  $\exists$  vertex in  $U$  not adjacent to  $v_{i+1}$  ( $U \not\subseteq N(v_{i+1})$ ):
  - $U$  itself is a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - if  $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal it is also a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is potentially child of many possible nodes → define it als child of lexicographically smallest clique (node)  $U'$  at level  $i$

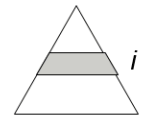


## Cliques: Algorithms

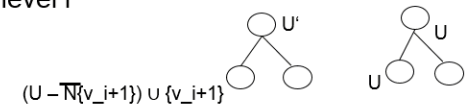
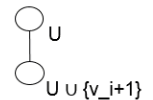
- Naive algorithm for finding the maximum clique: Exhaustive Search: Compute all subsets of  $V$  and check for clique →  $O(n^2 2^n)$
- Is there a slightly better algorithm? Yes: We can improve the exponential function's base from 2 to approx. 1.38 (see [2])
- Can we expect to find the answer to the problem "how many cliques with exactly  $k$  nodes exist"? → Exhaustive search:  $\Theta(|V|^k)$
- For triangles we can be better than  $\Theta(|V|^3)$ : An alg with  $O(|V|^{2.376})$  exists ☺
- For other  $k$ : An analogous technique allows for alg with  $O(|V|^{\beta(k)})$  with  $\beta(k) < k$

## Cliques: Enumeration

- Binary tree, n levels, leaves only at level n
- each level  $i \leftrightarrow$  vertex  $v_i$
- nodes at level  $i$ : maximal cliques in  $G[v_1, v_2, \dots, v_i]$
- level  $i+1$ : determine children of node  $U$  at level  $i$ : two cases :

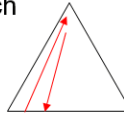


- (1)  $v_{i+1}$  adjacent to all nodes in  $U$  ( $U \subseteq N(v_{i+1})$ ):
  - $U \cup \{v_{i+1}\}$  is maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - $U \cup \{v_{i+1}\}$  is only child of  $U$
- (2)  $\exists$  vertex in  $U$  not adjacent to  $v_{i+1}$  ( $U \not\subseteq N(v_{i+1})$ ):
  - $U$  itself is a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - if  $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal it is also a maximal clique in  $G[v_1, v_2, \dots, v_i, v_{i+1}]$
  - $(U - \overline{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is potentially child of many possible nodes → define it als child of lexicographically smallest clique (node)  $U'$  at level  $i$



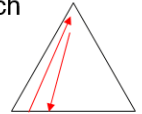
## Cliques: Enumeration

- tree constructed in principle
- Algorithm: **traverse and construct tree depth first, output leaves**
- necessary **primitives**:
  - **Parent(U,i)**: for node U at level i determine lexicographically smallest clique in  $G[v_1, v_2, \dots, v_{i-1}]$  :  $O(m+n)$
  - **LeftChild(U,i)**: either U or  $U \cup \{v_{i+1}\}$  :  $O(m+n)$
  - **RightChild(U,i)**:
    - if case (1): no right child;
    - if case (2): if  $X := (U - \mathbb{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal (takes  $O(m+n)$  time to determine) then X is right child if  $U = \text{Parent}(X, i+1)$  (takes  $O(m+n)$ )
- longest path** between two leaves passes over  $2n-1$  nodes; for each node:  $O(m+n)$  time → maximal delay:  $O(n^3)$



## Cliques: Enumeration

- tree constructed in principle
- Algorithm: **traverse and construct tree depth first, output leaves**
- necessary **primitives**:
  - **Parent(U,i)**: for node U at level i determine lexicographically smallest clique in  $G[v_1, v_2, \dots, v_{i-1}]$  :  $O(m+n)$
  - **LeftChild(U,i)**: either U or  $U \cup \{v_{i+1}\}$  :  $O(m+n)$
  - **RightChild(U,i)**:
    - if case (1): no right child;
    - if case (2): if  $X := (U - \mathbb{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal (takes  $O(m+n)$  time to determine) then X is right child if  $U = \text{Parent}(X, i+1)$  (takes  $O(m+n)$ )
- longest path** between two leaves passes over  $2n-1$  nodes; for each node:  $O(m+n)$  time → maximal delay:  $O(n^3)$



## Cliques: Enumeration

- tree constructed in principle
- Algorithm: **traverse and construct tree depth first, output leaves**
- necessary **primitives**:
  - **Parent(U,i)**: for node U at level i determine lexicographically smallest clique in  $G[v_1, v_2, \dots, v_{i-1}]$  :  $O(m+n)$
  - **LeftChild(U,i)**: either U or  $U \cup \{v_{i+1}\}$  :  $O(m+n)$
  - **RightChild(U,i)**:
    - if case (1): no right child;
    - if case (2): if  $X := (U - \mathbb{N}\{v_{i+1}\}) \cup \{v_{i+1}\}$  is maximal (takes  $O(m+n)$  time to determine) then X is right child if  $U = \text{Parent}(X, i+1)$  (takes  $O(m+n)$ )
- longest path** between two leaves passes over  $2n-1$  nodes; for each node:  $O(m+n)$  time → maximal delay:  $O(n^3)$



## Plexes and Cores

- A N-Plex is **maximal** iff it is not strictly contained in another larger N-Plex; A N-Plex is a **maximum** N-Plex iff it has a maximum number of vertices among all N-Plexes in G
- N-Plexes are closed under exclusion; N-Plexes are nested → **good candidates for social groups**
- Furthermore: If the set of nodes V of G is an N-Plex we have: (Combinatorial Proof: see [2]):
  - if  $N < (|V|+2) / 2$  then  $\text{diam}(G) \leq 2$
  - If an N-Plex has a not too large “sloppyness” N “its” diameter is very small → Socially realistic
- So far: wonderful but:



## Plexes and Cores

- But the decision problem  $\text{NPLEX}(G,k,N)$ : “Does  $G$  contain an  $N$ -Plex of size at least  $k$ ?” is **NP-complete for all  $N$** . Proof: Informal argument:  $\text{CLIQUE}(G,k)=\text{NPLEX}(G,k,1)$ ; formal: reduce  $\text{CLIQUE}(G,k)$  to  $\text{N-PLEX}$  (see [2]).
- Ok, so (instead of fixing how many edges can be missed at most): why not **demand a minimum degree** for every node in the structure?
- A subset  $U \subseteq V$  is a  **$N$ -Core** iff  $\delta(G([U])) \geq N$
- (Defs for maximal and maximum apply accordingly)
- If  $U$  is  $N$ -core,  $U$  is also a  $(|V|-N)$ -plex; Socially: If  $N$  is small compared to  $|U|$ ,  $N$ -cores are not very meaningful



## Plexes and Cores

- But the decision problem  $\text{NPLEX}(G,k,N)$ : “Does  $G$  contain an  $N$ -Plex of size at least  $k$ ?” is **NP-complete for all  $N$** . Proof: Informal argument:  $\text{CLIQUE}(G,k)=\text{NPLEX}(G,k,1)$ ; formal: reduce  $\text{CLIQUE}(G,k)$  to  $\text{N-PLEX}$  (see [2]).
- Ok, so (instead of fixing how many edges can be missed at most): why not **demand a minimum degree** for every node in the structure?
- A subset  $U \subseteq V$  is a  **$N$ -Core** iff  $\delta(G([U])) \geq N$
- (Defs for maximal and maximum apply accordingly)
- If  $U$  is  $N$ -core,  $U$  is also a  $(|V|-N)$ -plex; Socially: If  $N$  is small compared to  $|U|$ ,  $N$ -cores are not very meaningful



## Plexes and Cores

- But the decision problem  $\text{NPLEX}(G,k,N)$ : “Does  $G$  contain an  $N$ -Plex of size at least  $k$ ?” is **NP-complete for all  $N$** . Proof: Informal argument:  $\text{CLIQUE}(G,k)=\text{NPLEX}(G,k,1)$ ; formal: reduce  $\text{CLIQUE}(G,k)$  to  $\text{N-PLEX}$  (see [2]).
- Ok, so (instead of fixing how many edges can be missed at most): why not **demand a minimum degree** for every node in the structure?
- A subset  $U \subseteq V$  is a  **$N$ -Core** iff  $\delta(G([U])) \geq N$
- (Defs for maximal and maximum apply accordingly)
- If  $U$  is  $N$ -core,  $U$  is also a  $(|V|-N)$ -plex; Socially: If  $N$  is small compared to  $|U|$ ,  $N$ -cores are not very meaningful



## Plexes and Cores

- But the decision problem  $\text{NPLEX}(G,k,N)$ : “Does  $G$  contain an  $N$ -Plex of size at least  $k$ ?” is **NP-complete for all  $N$** . Proof: Informal argument:  $\text{CLIQUE}(G,k)=\text{NPLEX}(G,k,1)$ ; formal: reduce  $\text{CLIQUE}(G,k)$  to  $\text{N-PLEX}$  (see [2]).
- Ok, so (instead of fixing how many edges can be missed at most): why not **demand a minimum degree** for every node in the structure?
- A subset  $U \subseteq V$  is a  **$N$ -Core** iff  $\delta(G([U])) \geq N$
- (Defs for maximal and maximum apply accordingly)
- If  $U$  is  $N$ -core,  $U$  is also a  $(|V|-N)$ -plex; Socially: If  $N$  is small compared to  $|U|$ ,  $N$ -cores are not very meaningful



## Plexes and Cores

- If  $U_1$  and  $U_2$  are N-cores,  $U_1 \cup U_2$  is also an N-core; Socially: not very desirable; (math:  $\rightarrow$  maximum N-core is unique)
- If  $U_1$  and  $U_2$  are connected maximal N-cores  $\rightarrow U_1$  and  $U_2$  are disjoint; Socially: not very desirable
- N-cores are **not closed under exclusion** (Example: cycle is 2-core but subset not) and are generally **not nested**
- N-cores **need not be connected**
- $\rightarrow$  N-cores do not seem to be good candidates for social groups
- But: We can easily compute the maximum N-core (see [2])

## Plexes and Cores

- If  $U_1$  and  $U_2$  are N-cores,  $U_1 \cup U_2$  is also an N-core; Socially: not very desirable; (math:  $\rightarrow$  maximum N-core is unique)
- If  $U_1$  and  $U_2$  are connected maximal N-cores  $\rightarrow U_1$  and  $U_2$  are disjoint; Socially: not very desirable
- N-cores are **not closed under exclusion** (Example: cycle is 2-core but subset not) and are generally **not nested**
- N-cores **need not be connected**
- $\rightarrow$  N-cores do not seem to be good candidates for social groups
- But: We can easily compute the maximum N-core (see [2])



## LS-Sets and Lambda Sets

- We will now look at some **non-local** concepts:
- **LS-sets** (Luccio-Sami-Sets): Formalize the paradigm „Intra-cluster coherence, inter-cluster decoherence“: An LS-set is a „*network region where internal ties are more significant than external ties*“ [2] (extreme version of paradigm: „strong alliance“: complete component (disconnected clique))
- A subset  $U \subseteq V$  is a **LS-set** iff all proper subsets of  $U$  share more ties with the network outside than  $U$  does:  
$$U' \subset U \rightarrow |E(U', V-U)| > |E(U, V-U)|$$
- LS-sets have some **interesting properties** (e.g. no trivial overlaps: if  $u_1 \cap u_2 \neq \emptyset \rightarrow u_1 \subseteq u_2$  or  $u_2 \subseteq u_1$ ; min degree of LS set is at least half the number of outgoing edges) and may be good candidates for social groups (overlap property can be a counter argument)
- LS-sets can be computed with **reasonable complexity** ([2])

## LS-Sets and Lambda Sets

- **Lambda Sets** are another related concept: In a lambda set, members are connected to other members by more (edge disjoint) paths than to outside nodes
- Let  $\lambda(u,v)$  denote the number of edge-disjoint paths between nodes  $u$  and  $v$ ;  
A subset  $U \subseteq V$  is a Lambda set iff  
$$\min_{u,v \in U} \lambda(u,v) > \max_{u \in U, v \in V-U} \lambda(u,v)$$
- Lambda sets can be computed in P time [2]



- **Lambda Sets** are another related concept: In a lambda set, members are connected to other members by more (edge disjoint) paths than to outside nodes
- Let  $\lambda(u,v)$  denote the number of edge-disjoint paths between nodes  $u$  and  $v$ ;  
A subset  $U \subseteq V$  is a Lambda set iff

$$\min_{u,v \in U} \lambda(u,v) > \max_{u \in U, v \in V-U} \lambda(u,v)$$

- Lambda sets can be computed in P time [2]



### Graph Clustering

- Usually: **Clustering: Unsupervised Classification** (“partition a set of patterns  $\{v_i\}$  into subsets (classes)  $\{C_j\}$  without training data”)
- Contrast: **Supervised Classifier**: Train the parameters of a **classifier system** (e.g. the weights of a neural net, the probabilities of a naive Bayes classifier, the support weights of a support vector machine , etc.) with a pre-classified set of “training instances”  $\{(v_{train_i}, C(v_{train_i}))\}$
- **AI: Application** of Supervised Classifiers manifold; e.g.: text-classification (Classification (supervised or unsupervised): one of the most elementary intelligent “actions”
- Usually: Clustering: In **metric spaces**: example: K-Means. We will come back to this in next lecture (Clustering of Profile Elements)



- **Lambda Sets** are another related concept: In a lambda set, members are connected to other members by more (edge disjoint) paths than to outside nodes
- Let  $\lambda(u,v)$  denote the number of edge-disjoint paths between nodes  $u$  and  $v$ ;  
A subset  $U \subseteq V$  is a Lambda set iff

$$\min_{u,v \in U} \lambda(u,v) > \max_{u \in U, v \in V-U} \lambda(u,v)$$

- Lambda sets can be computed in P time [2]



### Graph Clustering

- Usually: **Clustering: Unsupervised Classification** (“partition a set of patterns  $\{v_i\}$  into subsets (classes)  $\{C_j\}$  without training data”)
- Contrast: **Supervised Classifier**: Train the parameters of a **classifier system** (e.g. the weights of a neural net, the probabilities of a naive Bayes classifier, the support weights of a support vector machine , etc.) with a pre-classified set of “training instances”  $\{(v_{train_i}, C(v_{train_i}))\}$
- **AI: Application** of Supervised Classifiers manifold; e.g.: text-classification (Classification (supervised or unsupervised): one of the most elementary intelligent “actions”
- Usually: Clustering: In **metric spaces**: example: K-Means. We will come back to this in next lecture (Clustering of Profile Elements)



## Graph Clustering

- Usually: **Clustering: Unsupervised Classification** (“partition a set of patterns  $\{v_i\}$  into subsets (classes)  $\{C_j\}$  without training data”)
- Contrast: **Supervised Classifier**: Train the parameters of a **classifier system** (e.g. the weights of a neural net, the probabilities of a naive Bayes classifier, the support weights of a support vector machine , etc.) with a pre-classified set of “training instances”  $\{(v_{train_i}, C(v_{train_i}))\}$
- **AI: Application** of Supervised Classifiers manifold; e.g.: text-classification (Classification (supervised or unsupervised): one of the most elementary intelligent “actions”
- Usually: Clustering: In **metric spaces**: example: K-Means. We will come back to this in next lecture (Clustering of Profile Elements)



## Graph Clustering

- Given directed, weighted graph  $G=(V,E,w)$ ; A **graph clustering**  $\mathbf{C}=\{C_1, C_2, \dots, C_k\}$  is a **partition of  $V$**  into non-empty subsets  $C_i$ ;
- **Notations:**
  - $E(C_i, C_j)$ : Set of edges in  $G$  from  $C_i$  to  $C_j$ ;
  - $E(\mathbf{C}) = \bigcup_{i=1, \dots, k} E(C_i)$ : Set of intra-cluster edges;
  - $\overline{E(\mathbf{C})} = E \setminus E(\mathbf{C})$ : Set of inter-cluster edges;
  - $m(\mathbf{C}) = |E(\mathbf{C})|$ ;  $\overline{m}(\mathbf{C}) = |\overline{E(\mathbf{C})}|$ ;
  - $G([C_i])$ : subgraph induced by  $C_i$ ;
  - $\mathbf{C}$  with  $k=1$ : *1-clustering*;  $\mathbf{C}$  with  $k=|V|$ : *singletons*; (both: *Trivial clustering*);
  - $\mathbf{C}$  with  $k=2$ : *cut*;
  - $\mathbf{A}(G)$ : Set of all possible clusterings on  $G$ ;



## Graph Clustering

- Usually: **Clustering: Unsupervised Classification** (“partition a set of patterns  $\{v_i\}$  into subsets (classes)  $\{C_j\}$  without training data”)
- Contrast: **Supervised Classifier**: Train the parameters of a **classifier system** (e.g. the weights of a neural net, the probabilities of a naive Bayes classifier, the support weights of a support vector machine , etc.) with a pre-classified set of “training instances”  $\{(v_{train_i}, C(v_{train_i}))\}$
- **AI: Application** of Supervised Classifiers manifold; e.g.: text-classification (Classification (supervised or unsupervised): one of the most elementary intelligent “actions”
- Usually: Clustering: In **metric spaces**: example: K-Means. We will come back to this in next lecture (Clustering of Profile Elements)



## Graph Clustering

- Given directed, weighted graph  $G=(V,E,w)$ ; A **graph clustering**  $\mathbf{C}=\{C_1, C_2, \dots, C_k\}$  is a **partition of  $V$**  into non-empty subsets  $C_i$ ;
- **Notations:**
  - $E(C_i, C_j)$ : Set of edges in  $G$  from  $C_i$  to  $C_j$ ;
  - $E(\mathbf{C}) = \bigcup_{i=1, \dots, k} E(C_i)$ : Set of intra-cluster edges;
  - $\overline{E(\mathbf{C})} = E \setminus E(\mathbf{C})$ : Set of inter-cluster edges;
  - $m(\mathbf{C}) = |E(\mathbf{C})|$ ;  $\overline{m}(\mathbf{C}) = |\overline{E(\mathbf{C})}|$ ;
  - $G([C_i])$ : subgraph induced by  $C_i$ ;
  - $\mathbf{C}$  with  $k=1$ : *1-clustering*;  $\mathbf{C}$  with  $k=|V|$ : *singletons*; (both: *Trivial clustering*);
  - $\mathbf{C}$  with  $k=2$ : *cut*;
  - $\mathbf{A}(G)$ : Set of all possible clusterings on  $G$ ;



Given directed, weighted graph  $G=(V,E,w)$ ; A **graph clustering**  $\mathbf{C}=\{C_1, C_2, \dots, C_k\}$  is a **partition of  $V$**  into non-empty subsets  $C_i$ ;

• **Notations:**

- $E(C_i, C_j)$ : Set of edges in  $G$  from  $C_i$  to  $C_j$ ;
- $E(\mathbf{C}) = \bigcup_{i=1, \dots, k} E(C_i)$ : Set of intra-cluster edges;
- $\overline{E(\mathbf{C})} = E \setminus E(\mathbf{C})$ : Set of inter-cluster edges;
- $m(\mathbf{C}) = |E(\mathbf{C})|$ ;  $\overline{m}(\mathbf{C}) = |\overline{E(\mathbf{C})}|$ ;
- $G([C_i])$ : subgraph induced by  $C_i$ ;
- $\mathbf{C}$  with  $k=1$ : *1-clustering*;  $\mathbf{C}$  with  $k=|V|$ : *singletons*;  
(both: *Trivial clustering*);
- $\mathbf{C}$  with  $k=2$ : *cut*;
- $\mathbf{A}(G)$ : Set of all possible clusterings on  $G$ ;

Given directed, weighted graph  $G=(V,E,w)$ ; A **graph clustering**  $\mathbf{C}=\{C_1, C_2, \dots, C_k\}$  is a **partition of  $V$**  into non-empty subsets  $C_i$ ;

• **Notations:**

- $E(C_i, C_j)$ : Set of edges in  $G$  from  $C_i$  to  $C_j$ ;
- $E(\mathbf{C}) = \bigcup_{i=1, \dots, k} E(C_i)$ : Set of intra-cluster edges;
- $\overline{E(\mathbf{C})} = E \setminus E(\mathbf{C})$ : Set of inter-cluster edges;
- $m(\mathbf{C}) = |E(\mathbf{C})|$ ;  $\overline{m}(\mathbf{C}) = |\overline{E(\mathbf{C})}|$ ;
- $G([C_i])$ : subgraph induced by  $C_i$ ;
- $\mathbf{C}$  with  $k=1$ : *1-clustering*;  $\mathbf{C}$  with  $k=|V|$ : *singletons*;  
(both: *Trivial clustering*);
- $\mathbf{C}$  with  $k=2$ : *cut*;
- $\mathbf{A}(G)$ : Set of all possible clusterings on  $G$ ;



Given directed, weighted graph  $G=(V,E,w)$ ; A **graph clustering**  $\mathbf{C}=\{C_1, C_2, \dots, C_k\}$  is a **partition of  $V$**  into non-empty subsets  $C_i$ ;

• **Notations:**

- $E(C_i, C_j)$ : Set of edges in  $G$  from  $C_i$  to  $C_j$ ;
- $E(\mathbf{C}) = \bigcup_{i=1, \dots, k} E(C_i)$ : Set of intra-cluster edges;
- $\overline{E(\mathbf{C})} = E \setminus E(\mathbf{C})$ : Set of inter-cluster edges;
- $m(\mathbf{C}) = |E(\mathbf{C})|$ ;  $\overline{m}(\mathbf{C}) = |\overline{E(\mathbf{C})}|$ ;
- $G([C_i])$ : subgraph induced by  $C_i$ ;
- $\mathbf{C}$  with  $k=1$ : *1-clustering*;  $\mathbf{C}$  with  $k=|V|$ : *singletons*;  
(both: *Trivial clustering*);
- $\mathbf{C}$  with  $k=2$ : *cut*;
- $\mathbf{A}(G)$ : Set of all possible clusterings on  $G$ ;



• **Notations (continued):**

- Let  $\mathbf{C}_1=\{C_1, C_2, \dots, C_k\}$  and  $\mathbf{C}_2=\{C'_1, C'_2, \dots, C'_l\}$ ;  
 $\mathbf{C}_1 \leq \mathbf{C}_2 \leftrightarrow \forall i \exists j : C_i \subseteq C'_j$ ;  
 $\mathbf{C}_1$ : *refinement of  $\mathbf{C}_2$* ;  $\mathbf{C}_2$ : *coarsening of  $\mathbf{C}_1$* ;
- *Chain* (comparable set) of clusterings: *hierarchy*;
- Hierarchy is *total*  $\leftrightarrow$  Both trivial clusterings are contained;
- Hierarchy that contains one clustering for each  $\{1, 2, \dots, |V|\}$ :  
*complete*;
- $S(V)$ : (*Proper*) *Cut function*: Cutting the node-set in two (non-empty) subsets:  $S(V)$  and  $V \setminus S(V)$ ;



• **Notations (continued):**

- Let  $\mathbf{C}_1 = \{C_1, C_2, \dots, C_k\}$  and  $\mathbf{C}_2 = \{C'_1, C'_2, \dots, C'_l\}$ :  
 $\mathbf{C}_1 \leq \mathbf{C}_2 \leftrightarrow \forall i \exists j : C_i \subseteq C'_j$ ;  
 $\mathbf{C}_1$ : *refinement of  $\mathbf{C}_2$* ;  $\mathbf{C}_2$ : *coarsening of  $\mathbf{C}_1$* ;
- *Chain* (comparable set) of clusterings: *hierarchy*;
- Hierarchy is *total*  $\leftrightarrow$  Both trivial clusterings are contained;
- Hierarchy that contains one clustering for each  $\{1, 2, \dots, |V|\}$ : *complete*;
- $S(V)$ : (*Proper*) *Cut function*: Cutting the node-set in two (non-empty) subsets:  $S(V)$  and  $V \setminus S(V)$ ;

- **Quality measure**: Objective function  $\mathbf{A}(G) \rightarrow \mathbb{R}$  that formalizes the clustering paradigm in a special way
- $G = (V, E, w)$ : **Weight function**  $w: E \rightarrow \mathbb{R}^+$  is interpreted as “**similarity**” (higher weights correspond to more intense tie); also possible: negative weights = dissimilarity; or  $w: E \rightarrow [0, 1]$  or  $w: E \rightarrow [-1, 1]$  etc.
- **Distinguish** between no edge and edge with weight zero;
- **Notation**:  $w(E) = \sum_{e \in E} w(e)$



Quality Measures for Clusterings

- **Quality measure**: Objective function  $\mathbf{A}(G) \rightarrow \mathbb{R}$  that formalizes the clustering paradigm in a special way
- $G = (V, E, w)$ : **Weight function**  $w: E \rightarrow \mathbb{R}^+$  is interpreted as “**similarity**” (higher weights correspond to more intense tie); also possible: negative weights = dissimilarity; or  $w: E \rightarrow [0, 1]$  or  $w: E \rightarrow [-1, 1]$  etc.
- **Distinguish** between no edge and edge with weight zero;
- **Notation**:  $w(E) = \sum_{e \in E} w(e)$



6: Graph Clustering

- General **framework** for a quality index of a clustering:

$$index(\mathbf{C}) = \frac{f(\mathbf{C}) + g(\mathbf{C})}{\max\{f(\mathbf{C}') + g(\mathbf{C}') : \mathbf{C}' \in \mathbf{A}(G)\}}$$

- $f: \mathbf{A}(G) \rightarrow \mathbb{R}^+$  measures **intra cluster density** (coherence);  
 $g: \mathbf{A}(G) \rightarrow \mathbb{R}^+$  measures **inter cluster sparseness** (decoherence);





- General **framework** for a quality index of a clustering:

$$\text{index}(\mathbf{C}) = \frac{f(\mathbf{C}) + g(\mathbf{C})}{\max\{f(\mathbf{C}') + g(\mathbf{C}') : \mathbf{C}' \in \mathbf{A}(G)\}}$$

- $f: \mathbf{A}(G) \rightarrow \mathbb{R}^+$  measures **intra cluster density** (coherence);
- $g: \mathbf{A}(G) \rightarrow \mathbb{R}^+$  measures **inter cluster sparseness** (decoherence);

- General **framework** for a quality index of a clustering:

$$\text{index}(\mathbf{C}) = \frac{f(\mathbf{C}) + g(\mathbf{C})}{\max\{f(\mathbf{C}') + g(\mathbf{C}') : \mathbf{C}' \in \mathbf{A}(G)\}}$$

- $f: \mathbf{A}(G) \rightarrow \mathbb{R}^+$  measures **intra cluster density** (coherence);
- $g: \mathbf{A}(G) \rightarrow \mathbb{R}^+$  measures **inter cluster sparseness** (decoherence);



### Coverage

- First quality measure: **Coverage**

$$\gamma(\mathbf{C}) = \frac{w(E(\mathbf{C}))}{w(E)} = \frac{\sum_{e \in E(\mathbf{C})} w(e)}{\sum_{e \in E} w(e)}$$

- Thus:  $f = w(E(\mathbf{C}))$  and  $g = 0$ ;  $\rightarrow$  only accumulated intra cluster density is measured

- **Maximum value 1** achieved for  $\mathbf{C} = \{V\}$  (1-clustering)
- A clustering has coverage  $\gamma(\mathbf{C}) = 1$  iff  $\overline{E(\mathbf{C})} = \emptyset$  (clustering is composed of connected components of  $G$ ) or  $w(E(\mathbf{C})) = 0$
- $\mathbf{C}$  with  $k > 1$  can be transformed into  $\mathbf{C}'$  with  $k' < k$  and  $\gamma(\mathbf{C}) \leq \gamma(\mathbf{C}')$  by **merging two clusters** in  $\mathbf{C} \rightarrow$  optimal non-trivial  $\mathbf{C}$  is a minimum cut
- $\rightarrow$  „Monotonic“ behavior of coverage  $\rightarrow$  coverage is not a good sole quality index



### Coverage

- First quality measure: **Coverage**

$$\gamma(\mathbf{C}) = \frac{w(E(\mathbf{C}))}{w(E)} = \frac{\sum_{e \in E(\mathbf{C})} w(e)}{\sum_{e \in E} w(e)}$$

- Thus:  $f = w(E(\mathbf{C}))$  and  $g = 0$ ;  $\rightarrow$  only accumulated intra cluster density is measured

- **Maximum value 1** achieved for  $\mathbf{C} = \{V\}$  (1-clustering)
- A clustering has coverage  $\gamma(\mathbf{C}) = 1$  iff  $\overline{E(\mathbf{C})} = \emptyset$  (clustering is composed of connected components of  $G$ ) or  $w(E(\mathbf{C})) = 0$
- $\mathbf{C}$  with  $k > 1$  can be transformed into  $\mathbf{C}'$  with  $k' < k$  and  $\gamma(\mathbf{C}) \leq \gamma(\mathbf{C}')$  by **merging two clusters** in  $\mathbf{C} \rightarrow$  optimal non-trivial  $\mathbf{C}$  is a minimum cut
- $\rightarrow$  „Monotonic“ behavior of coverage  $\rightarrow$  coverage is not a good sole quality index



## Conductance

- **Clustering paradigm reformulated:** Clusters should be **well connected** (many edges need to be removed to make it unconnected); few inter cluster edges (ideally none)
- **Conductance: Measure for bottlenecks** (Bottleneck: Cut that separates  $V$  into roughly same size halves and “cuts across” relatively few edges)
- Let  $\mathbf{C}=\{C_1, V \setminus C_1\}$  be a cut. Conductance  $\phi$  of  $\mathbf{C}$  is defined as

$$\phi(\mathbf{C}) = \begin{cases} 1 & \text{if } C_1 \in \{\emptyset, V\} \\ 0 & \text{if } C_1 \notin \{\emptyset, V\}, w(\overline{E(\mathbf{C})})=0 \\ \frac{w(\overline{E(\mathbf{C})})}{\min(\sum_{e \in E(C_1, V)} w(e), \sum_{e \in E(V \setminus C_1, V)} w(e))} & \text{otherwise} \end{cases}$$

the smaller  $\phi(\mathbf{C})$ , the more „bottlenecky“ is  $\mathbf{C}$

## Conductance

- **Clustering paradigm reformulated:** Clusters should be **well connected** (many edges need to be removed to make it unconnected); few inter cluster edges (ideally none)
- **Conductance: Measure for bottlenecks** (Bottleneck: Cut that separates  $V$  into roughly same size halves and “cuts across” relatively few edges)
- Let  $\mathbf{C}=\{C_1, V \setminus C_1\}$  be a cut. Conductance  $\phi$  of  $\mathbf{C}$  is defined as

$$\phi(\mathbf{C}) = \begin{cases} 1 & \text{if } C_1 \in \{\emptyset, V\} \\ 0 & \text{if } C_1 \notin \{\emptyset, V\}, w(\overline{E(\mathbf{C})})=0 \\ \frac{w(\overline{E(\mathbf{C})})}{\min(\sum_{e \in E(C_1, V)} w(e), \sum_{e \in E(V \setminus C_1, V)} w(e))} & \text{otherwise} \end{cases}$$

the smaller  $\phi(\mathbf{C})$ , the more „bottlenecky“ is  $\mathbf{C}$

## Conductance

- **Clustering paradigm reformulated:** Clusters should be **well connected** (many edges need to be removed to make it unconnected); few inter cluster edges (ideally none)
- **Conductance: Measure for bottlenecks** (Bottleneck: Cut that separates  $V$  into roughly same size halves and “cuts across” relatively few edges)
- Let  $\mathbf{C}=\{C_1, V \setminus C_1\}$  be a cut. Conductance  $\phi$  of  $\mathbf{C}$  is defined as

$$\phi(\mathbf{C}) = \begin{cases} 1 & \text{if } C_1 \in \{\emptyset, V\} \\ 0 & \text{if } C_1 \notin \{\emptyset, V\}, w(\overline{E(\mathbf{C})})=0 \\ \frac{w(\overline{E(\mathbf{C})})}{\min(\sum_{e \in E(C_1, V)} w(e), \sum_{e \in E(V \setminus C_1, V)} w(e))} & \text{otherwise} \end{cases}$$

the smaller  $\phi(\mathbf{C})$ , the more „bottlenecky“ is  $\mathbf{C}$

## Conductance

- **Clustering paradigm reformulated:** Clusters should be **well connected** (many edges need to be removed to make it unconnected); few inter cluster edges (ideally none)
- **Conductance: Measure for bottlenecks** (Bottleneck: Cut that separates  $V$  into roughly same size halves and “cuts across” relatively few edges)
- Let  $\mathbf{C}=\{C_1, V \setminus C_1\}$  be a cut. Conductance  $\phi$  of  $\mathbf{C}$  is defined as

$$\phi(\mathbf{C}) = \begin{cases} 1 & \text{if } C_1 \in \{\emptyset, V\} \\ 0 & \text{if } C_1 \notin \{\emptyset, V\}, w(\overline{E(\mathbf{C})})=0 \\ \frac{w(\overline{E(\mathbf{C})})}{\min(\sum_{e \in E(C_1, V)} w(e), \sum_{e \in E(V \setminus C_1, V)} w(e))} & \text{otherwise} \end{cases}$$

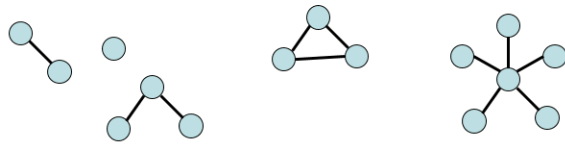
the smaller  $\phi(\mathbf{C})$ , the more „bottlenecky“ is  $\mathbf{C}$

## Conductance

- Conductance  $\phi$  of  $G$  is defined as

$$\phi(G) = \min_{C_1 \subseteq V} \phi(C_1, V \setminus C_1)$$

- Small conductance  $\leftrightarrow$  "good cut possible"
- All unconnected graphs have conductance 0
- Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star.

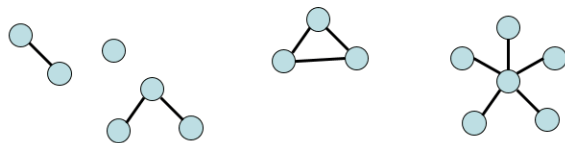


## Conductance

- Conductance  $\phi$  of  $G$  is defined as

$$\phi(G) = \min_{C_1 \subseteq V} \phi(C_1, V \setminus C_1)$$

- Small conductance  $\leftrightarrow$  "good cut possible"
- All unconnected graphs have conductance 0
- Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star.

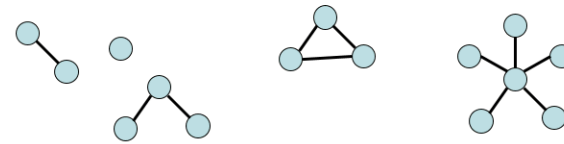


## Conductance

- Conductance  $\phi$  of  $G$  is defined as

$$\phi(G) = \min_{C_1 \subseteq V} \phi(C_1, V \setminus C_1)$$

- Small conductance  $\leftrightarrow$  "good cut possible"
- All unconnected graphs have conductance 0
- Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star.

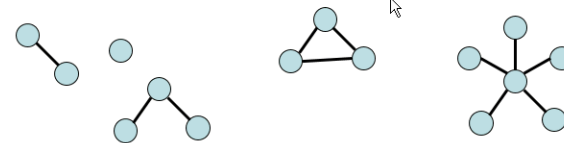


## Conductance

- Conductance  $\phi$  of  $G$  is defined as

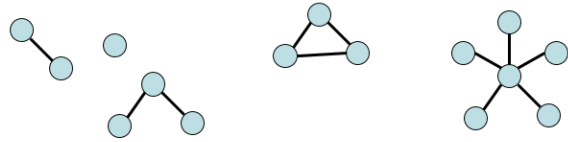
$$\phi(G) = \min_{C_1 \subseteq V} \phi(C_1, V \setminus C_1)$$

- Small conductance  $\leftrightarrow$  "good cut possible"
- All unconnected graphs have conductance 0
- Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star.



## Conductance

• **Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star. (Proof: see [1])



• **Proof** ←  $\sum_{e \in E(C_{-1}, V)} w(e) = w(E(C_{-1})) + w(\overline{E(C)}) \rightarrow$

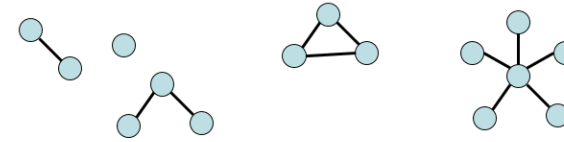
$$\frac{w(\overline{E(C)})}{\min(\sum_{e \in E(C_{-1}, V)} w(e), \sum_{e \in E(V \setminus C_{-1}, V)} w(e))} =$$

$$\frac{w(\overline{E(C)})}{w(\overline{E(C)}) + \min(w(E(C_{-1})), w(E(V \setminus C_{-1})))} = 1$$

=0 if star or at most 3 nodes

## Conductance

• **Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star. (Proof: see [1])



• **Proof** ←  $\sum_{e \in E(C_{-1}, V)} w(e) = w(E(C_{-1})) + w(\overline{E(C)}) \rightarrow$

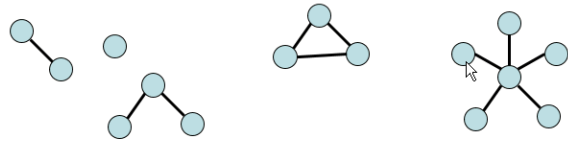
$$\frac{w(\overline{E(C)})}{\min(\sum_{e \in E(C_{-1}, V)} w(e), \sum_{e \in E(V \setminus C_{-1}, V)} w(e))} =$$

$$\frac{w(\overline{E(C)})}{w(\overline{E(C)}) + \min(w(E(C_{-1})), w(E(V \setminus C_{-1})))} = 1$$

=0 if star or at most 3 nodes

## Conductance

• **Theorem:** If  $G$  is undirected and positively weighted,  $G$  has maximum conductance  $\phi(G)=1$  iff  $G$  is connected and has at most three nodes or is a star. (Proof: see [1])



• **Proof** ←  $\sum_{e \in E(C_{-1}, V)} w(e) = w(E(C_{-1})) + w(\overline{E(C)}) \rightarrow$

$$\frac{w(\overline{E(C)})}{\min(\sum_{e \in E(C_{-1}, V)} w(e), \sum_{e \in E(V \setminus C_{-1}, V)} w(e))} =$$

$$\frac{w(\overline{E(C)})}{w(\overline{E(C)}) + \min(w(E(C_{-1})), w(E(V \setminus C_{-1})))} = 1$$

=0 if star or at most 3 nodes

## Conductance

• Conductance appears to be a **good element for a quality measure**, but: calculating it is **NP hard** ⊗ but: it can be approximated with guarantee  $O((\log |V|)^{1/2})$  (see [1]) (This means: approximation \*  $O((\log |V|)^{1/2})$  = true value)

• The conductance of a complete graph is asymptotically 0.5: Let  $n$  be an integer:

$$\phi(K_n) = \begin{cases} 0.5 \frac{n}{n-1} & \text{if } n \text{ is even} \\ 0.5 + \frac{1}{n-1} & \text{if } n \text{ is odd} \end{cases}$$

## Conductance

- Conductance appears to be a **good element for a quality measure**, but: calculating it is **NP hard** ⊗ but: it can be approximated with guarantee  $O((\log |V|)^{1/2})$  (see [1]) (This means: approximation \*  $O((\log |V|)^{1/2})$  = true value)
- The conductance of a complete graph is asymptotically 0.5: Let  $n$  be an integer:

$$\varphi(K_n) = \begin{cases} 0.5 \frac{n}{n-1} & \text{if } n \text{ is even} \\ 0.5 + \frac{1}{n-1} & \text{if } n \text{ is odd} \end{cases}$$



## Conductance

- With conductance we can define two appropriate **quality measures** for clusterings:
- **First measure:**  $g=0$  and  $f(\mathbf{C}) = \min_{1 \leq i \leq k} \varphi(G[\mathbf{C}_i])$
- **If the measure is small:** At least one of the clusters (more precisely: the induced subgraph) contains at least one bottleneck → This cluster is **too coarse** → Use minimum conductance cut to cut this cluster in “halves”
- **From theorem before:** Only clusterings where the clusters induce subgraphs that are stars or have size at most three have  $f=1$  ( $f$  is called *intra cluster conductance*)



## Conductance

- With conductance we can define two appropriate **quality measures** for clusterings:
- **First measure:**  $g=0$  and  $f(\mathbf{C}) = \min_{1 \leq i \leq k} \varphi(G[\mathbf{C}_i])$
- **If the measure is small:** At least one of the clusters (more precisely: the induced subgraph) contains at least one bottleneck → This cluster is **too coarse** → Use minimum conductance cut to cut this cluster in “halves”
- **From theorem before:** Only clusterings where the clusters induce subgraphs that are stars or have size at most three have  $f=1$  ( $f$  is called *intra cluster conductance*)



## Conductance

- With conductance we can define two appropriate **quality measures** for clusterings:
- **First measure:**  $g=0$  and  $f(\mathbf{C}) = \min_{1 \leq i \leq k} \varphi(G[\mathbf{C}_i])$
- **If the measure is small:** At least one of the clusters (more precisely: the induced subgraph) contains at least one bottleneck → This cluster is **too coarse** → Use minimum conductance cut to cut this cluster in “halves”
- **From theorem before:** Only clusterings where the clusters induce subgraphs that are stars or have size at most three have  $f=1$  ( $f$  is called *intra cluster conductance*)

