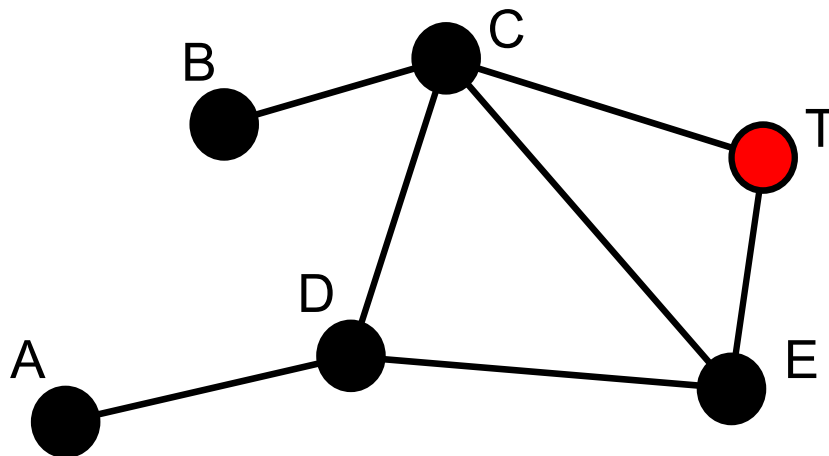


Emergent Groups: Detecting an Emergent Subgroup

- Clumpiness
- Regions
- Subgroups

Transitivity

- Proportion of triples with 3 ties as a proportion of triples with 2 or more ties
 - Aka the clustering coefficient



$$cc = 2/6 = 33\%$$

$\{C, T, E\}$ is a transitive triple, but $\{B, C, D\}$ is not. $\{A, D, T\}$ is not counted at all.

Network Regions

Network Regions

- Large “contiguous” areas
- Areas that contain cohesive subgroups
- We will cover:
 - Components
 - K-Cores

Graph Terminology

- A graph $G(V,E)$ consists of a set of nodes V and a set of lines E . Each line $e \in E$ consists of a pair of nodes (u,v)
- A graph G' is a subgraph of a graph G if every line in $E(G')$ is in $E(G)$, and every node in $V(G')$ is in $V(G)$.
- The subgraph S induced by a set of nodes consists of those nodes together with all ties among them

Components

- A subgraph S of a graph G is a component if S is maximal and connected
 - Connected means that every node can reach every other by some path (no matter how long)

Components in Digraphs

- If G is a digraph, then
 - S is a weak component if it is a component of the underlying (undirected) graph
 - i.e., we allow semi-paths rather than require true directed paths
 - S is a strong component if for all u, v in S , there is a path from u to v

Notes on Components

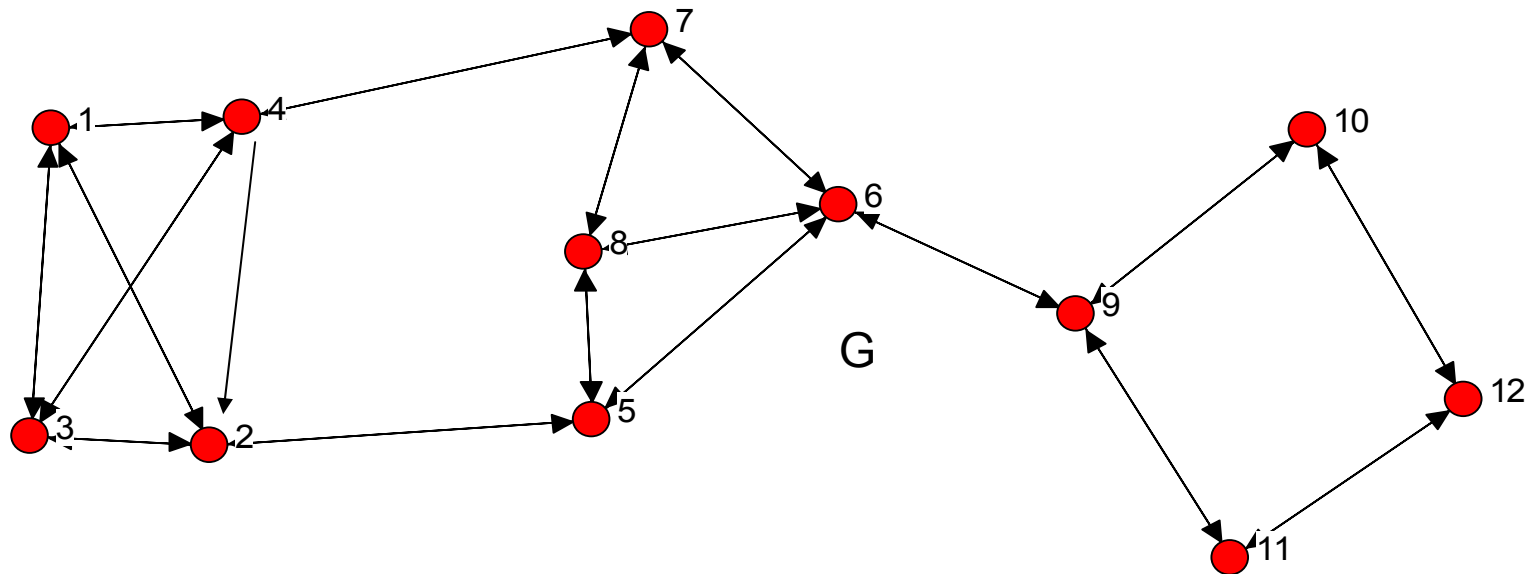
- Isolates are (very small) components
- Finding components is often first step in analysis of large graphs
 - Often analyze each component separately, or discard very small components
 - Many network measures require a connected graph, so they don't work on graphs with multiple components

Alpha Operator

- Let $\alpha(S1, S2)$ be the number of ties from members of set $S1$ to members of the set $S2$
- $\alpha(u, S)$ is number of ties node u has with members of set S
- $\alpha(S) = \alpha(S, V-S)$ is number of ties from members of set S to members of $V-S$ (i.e., all other nodes)

K-Core

- A subgraph S is a k -core if for all $u \in S$, $\alpha(u, S) \geq k$, and S is maximal



- $S=G$ is 1-core & 2-core; $S = \{1..8\}$ is 3-core
- There is no 4-core or higher

K-Core Notes

- Finds areas within which cohesive subgroups may be found
- Identifies fault lines across which cohesive subgroups do not span
- In large datasets, you can successively examine the 1-cores, the 2-cores, etc.
 - Progressively narrowing to core of network

Cohesive Subgroups

Cohesive Subgroups

- Initially conceived of as formalizations of fundamental sociological concepts
 - Primary groups
 - Emergent groups
- Now typically thought of in terms of a technique for identifying groups within networks

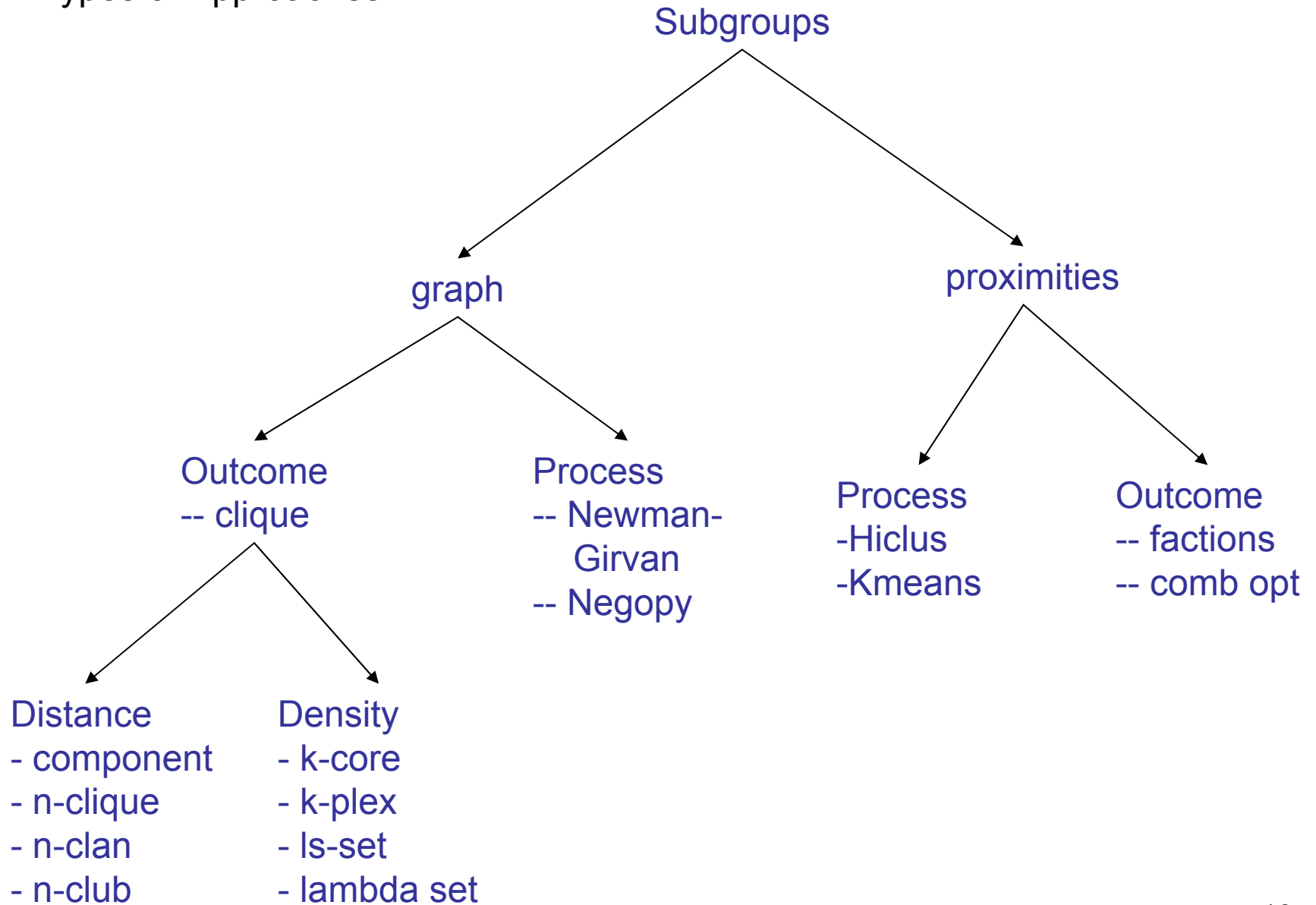
Canonical Hypothesis

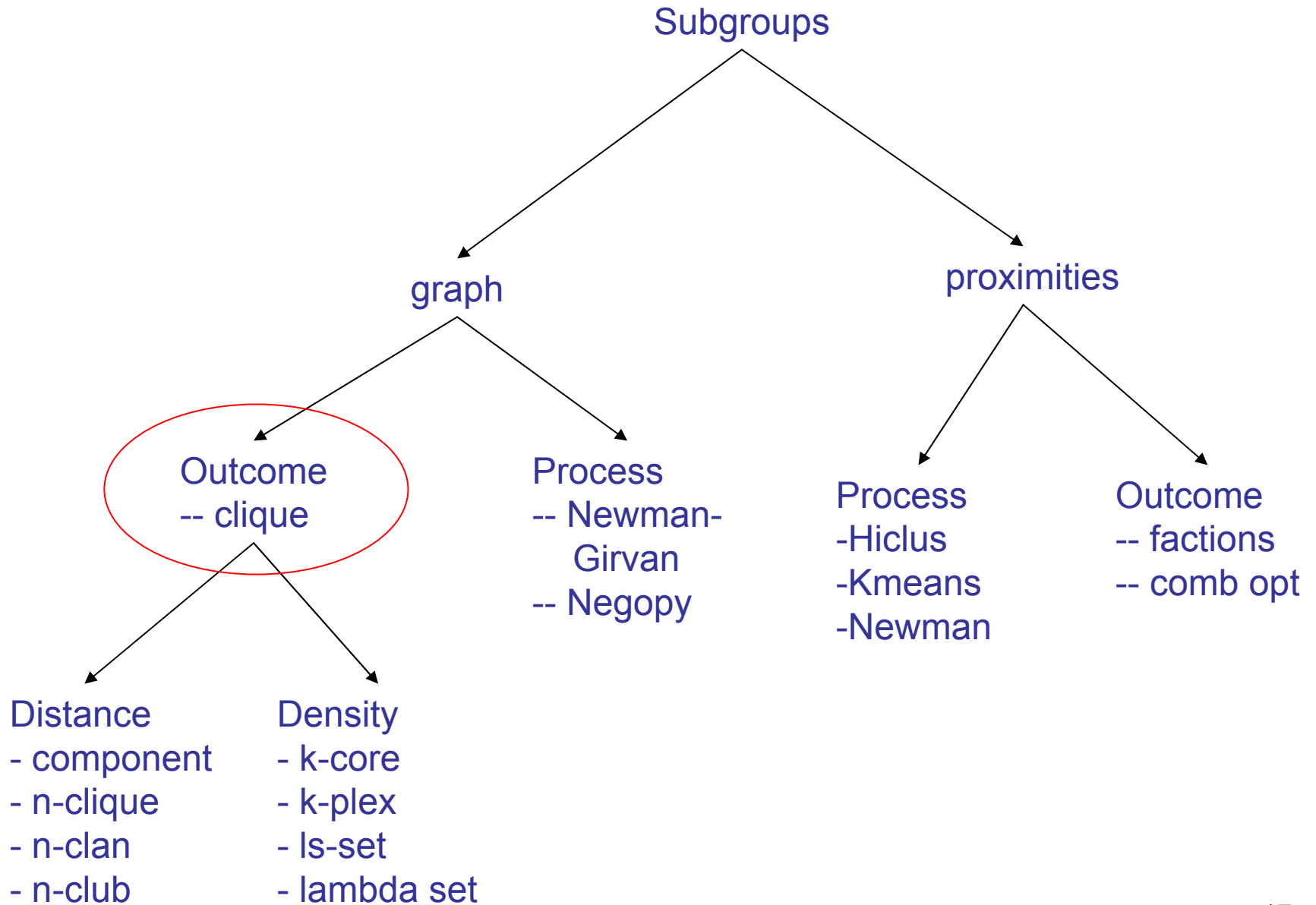
- Members of group will have similar outcomes
 - Ideas, attitudes, illnesses, behaviors
- Due to interpersonal transmission
 - transference
 - Influence / persuasion
 - Co-construction of beliefs & practices
 - As in communities of practice
- So group membership is independent var used to predict commonality of attitudes, beliefs, etc.

Typology of Subgroups

	Process	Outcome
Network / Graph theory	Newman-Girvan	Clique, n-clique, n-clan, n-club, k-plex, ls-set, lambda-set, k-core, component
Proximities / Clustering	Johnson's Hierarchical clustering; k-means; MDS	Factions, combinatorial optimization

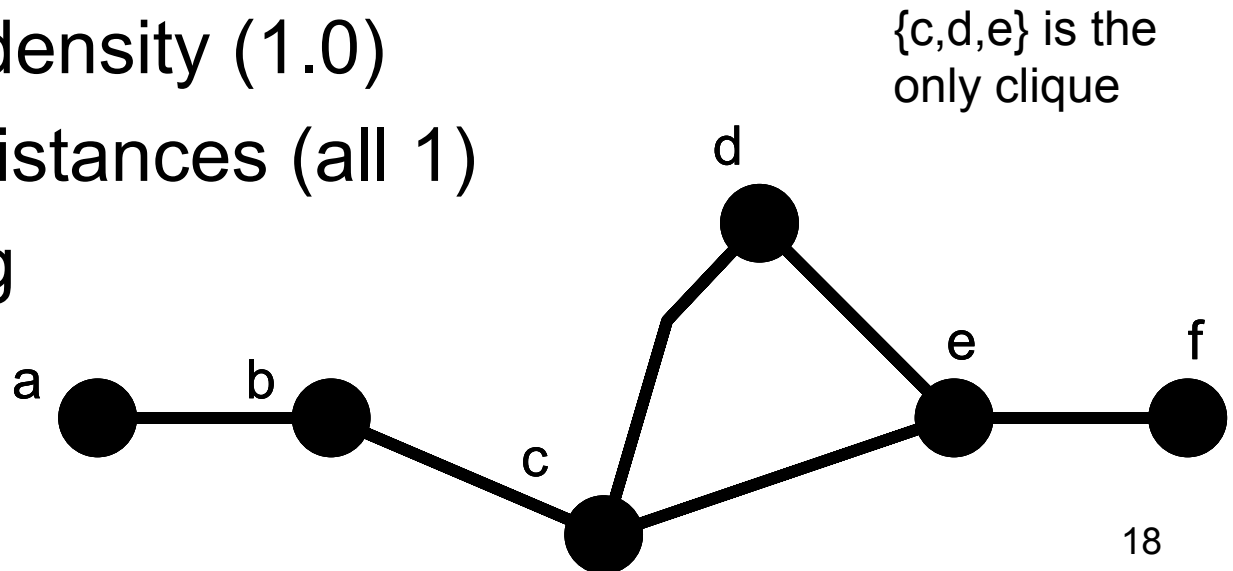
Types of Approaches





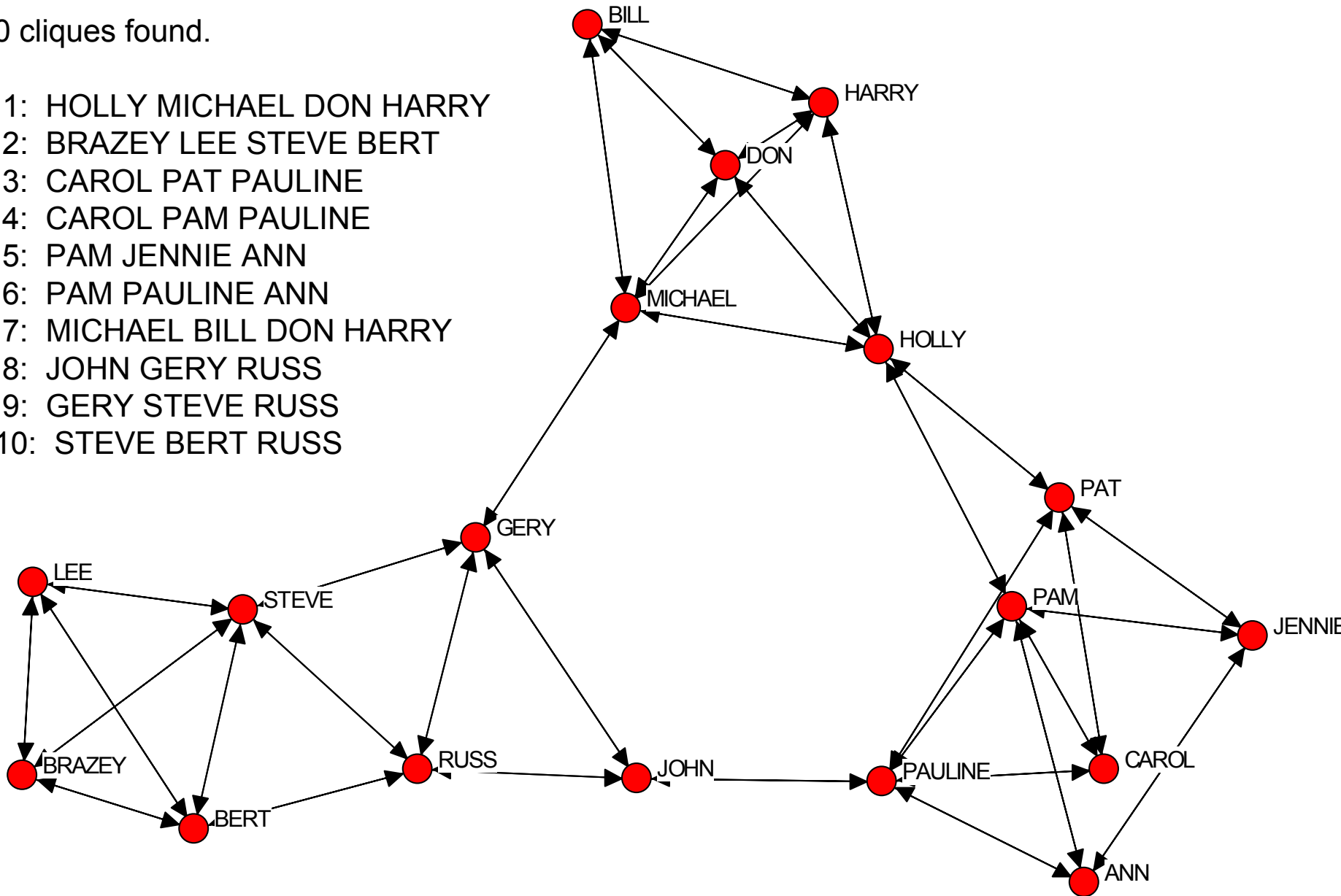
Cliques

- Definition
 - Maximal, complete subgraph
 - Set S s.t. for all u, v in S , (u, v) in E
- Properties
 - Maximum density (1.0)
 - Minimum distances (all 1)
 - overlapping
 - Strict



10 cliques found.

- 1: HOLLY MICHAEL DON HARRY
- 2: BRAZEY LEE STEVE BERT
- 3: CAROL PAT PAULINE
- 4: CAROL PAM PAULINE
- 5: PAM JENNIE ANN
- 6: PAM PAULINE ANN
- 7: MICHAEL BILL DON HARRY
- 8: JOHN GERY RUSS
- 9: GERY STEVE RUSS
- 10: STEVE BERT RUSS



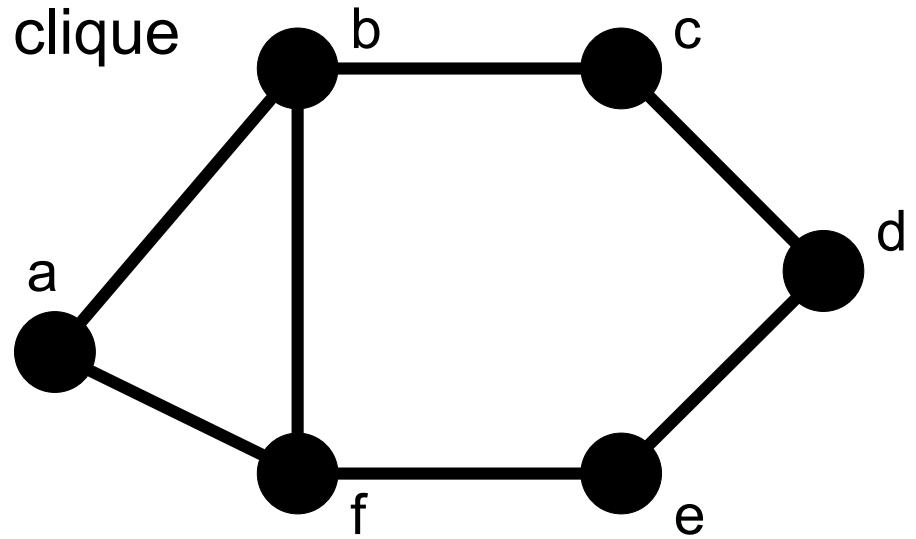
Types of Relaxations

- Distance (length of paths)
 - N-clique, n-clan, n-club
- Density (number of ties)
 - K-plex, ls-set, lambda set, k-core, component

N-cliques

- Definition
 - Maximal subset s.t. for all u, v in S , $d(u, v) \leq n$
 - Distance among members less than specified maximum
 - When $n = 1$, we have a clique

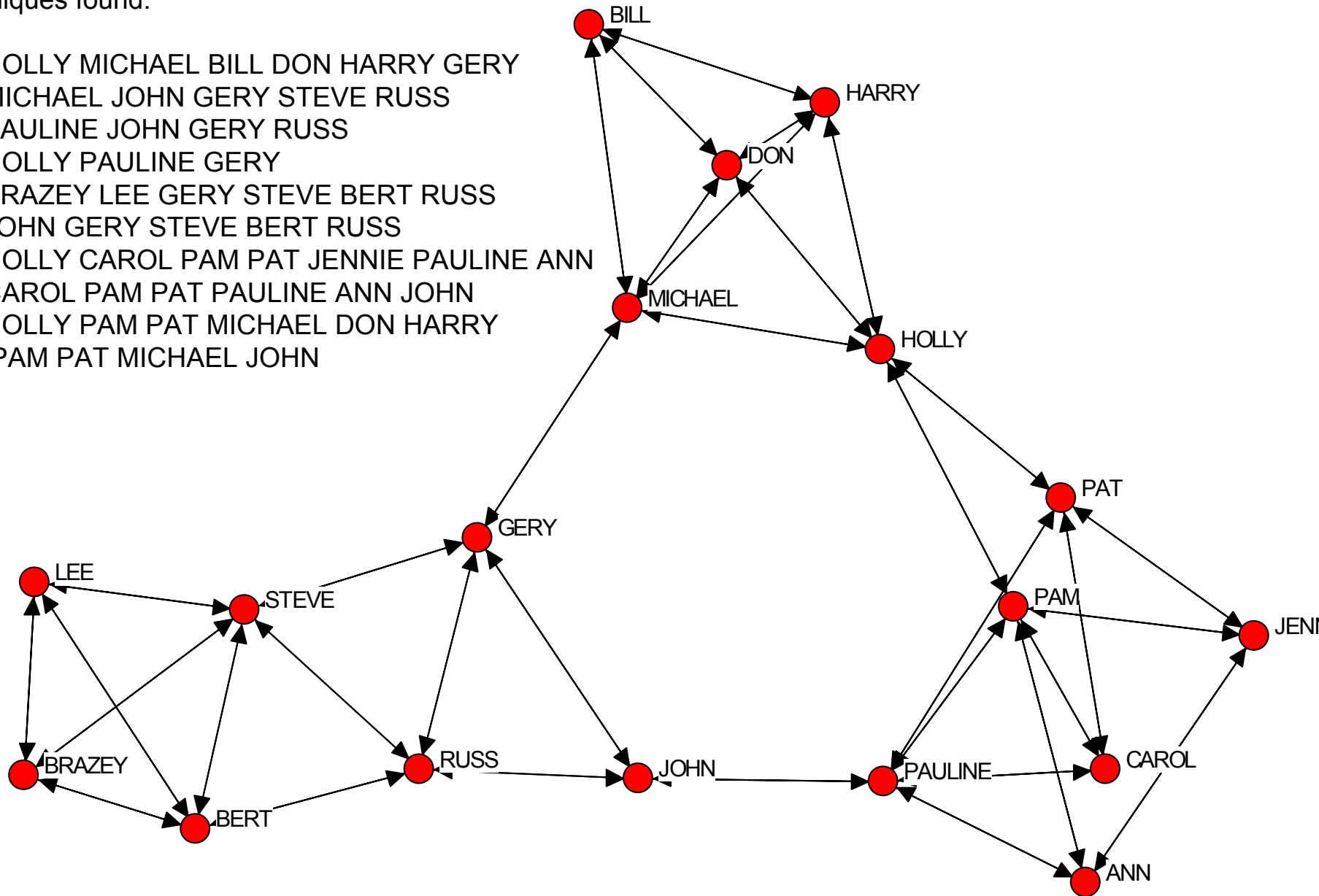
- Properties
 - Relaxes notion of clique
 - Avg distance can be greater than 1



Is $\{a, b, c, f, e\}$ a 2-clique?
yes

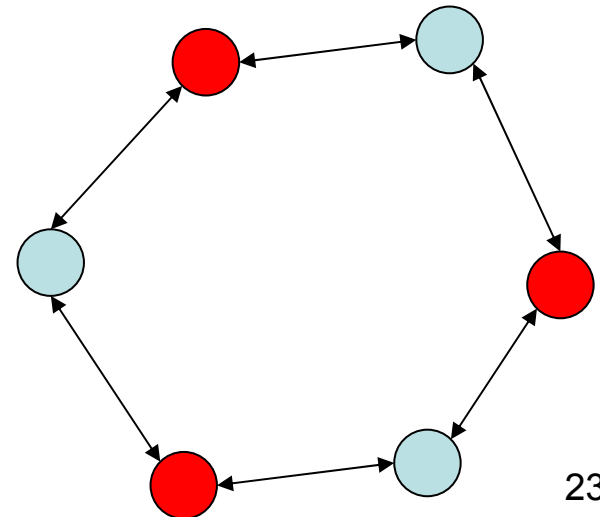
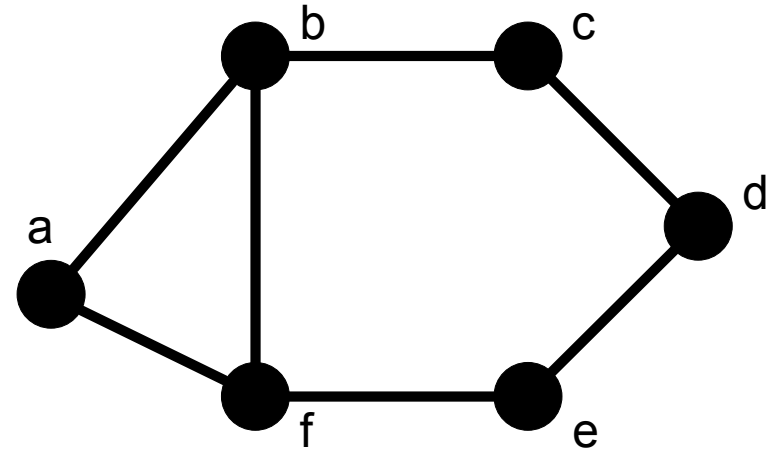
10 2-cliques found.

- 1: HOLLY MICHAEL BILL DON HARRY GERY
- 2: MICHAEL JOHN GERY STEVE RUSS
- 3: PAULINE JOHN GERY RUSS
- 4: HOLLY PAULINE GERY
- 5: BRAZEY LEE GERY STEVE BERT RUSS
- 6: JOHN GERY STEVE BERT RUSS
- 7: HOLLY CAROL PAM PAT JENNIE PAULINE ANN
- 8: CAROL PAM PAT PAULINE ANN JOHN
- 9: HOLLY PAM PAT MICHAEL DON HARRY
- 10: PAM PAT MICHAEL JOHN



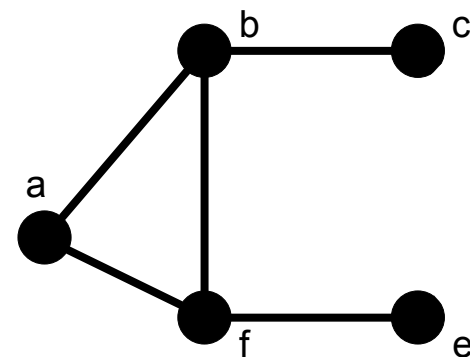
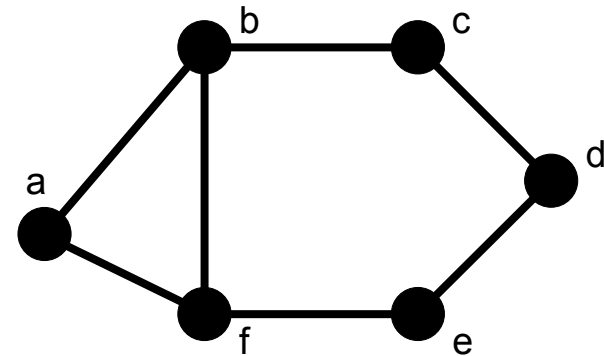
Issues with N-Cliques

- Overlapping
 - $\{a,b,c,f,e\}$ and $\{b,c,d,f,e\}$ are both 2-cliques
- Membership criterion satisfiable through non-members
- Even 2-cliques can be fairly non-cohesive
 - Red nodes belong to same 2-clique but none are adjacent



Subgraphs

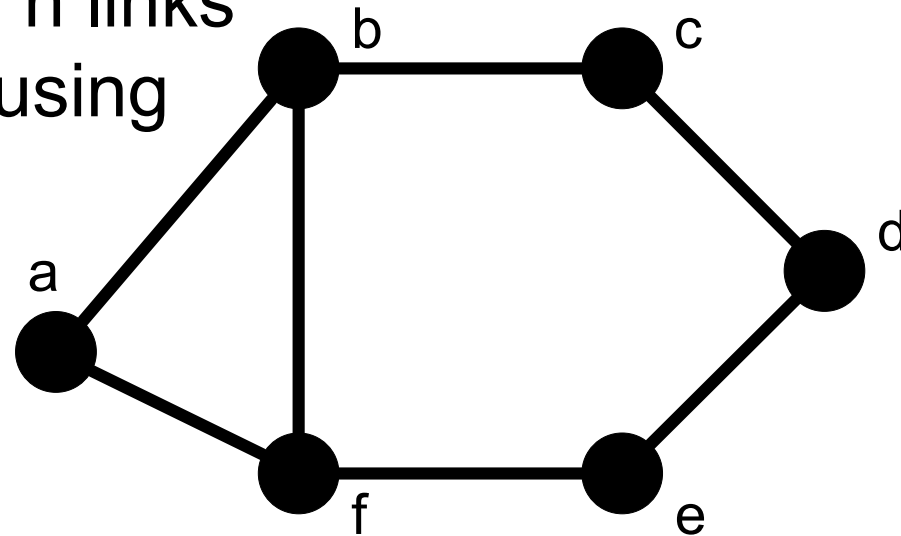
- Set of nodes
 - Is just a set of nodes
- A subgraph
 - Is set of nodes together with ties among them
- An induced subgraph
 - Subgraph defined by a set of nodes
 - Like pulling the nodes and ties out of the original graph



Subgraph induced by $\{a,b,c,f,e\}$

N-Clan

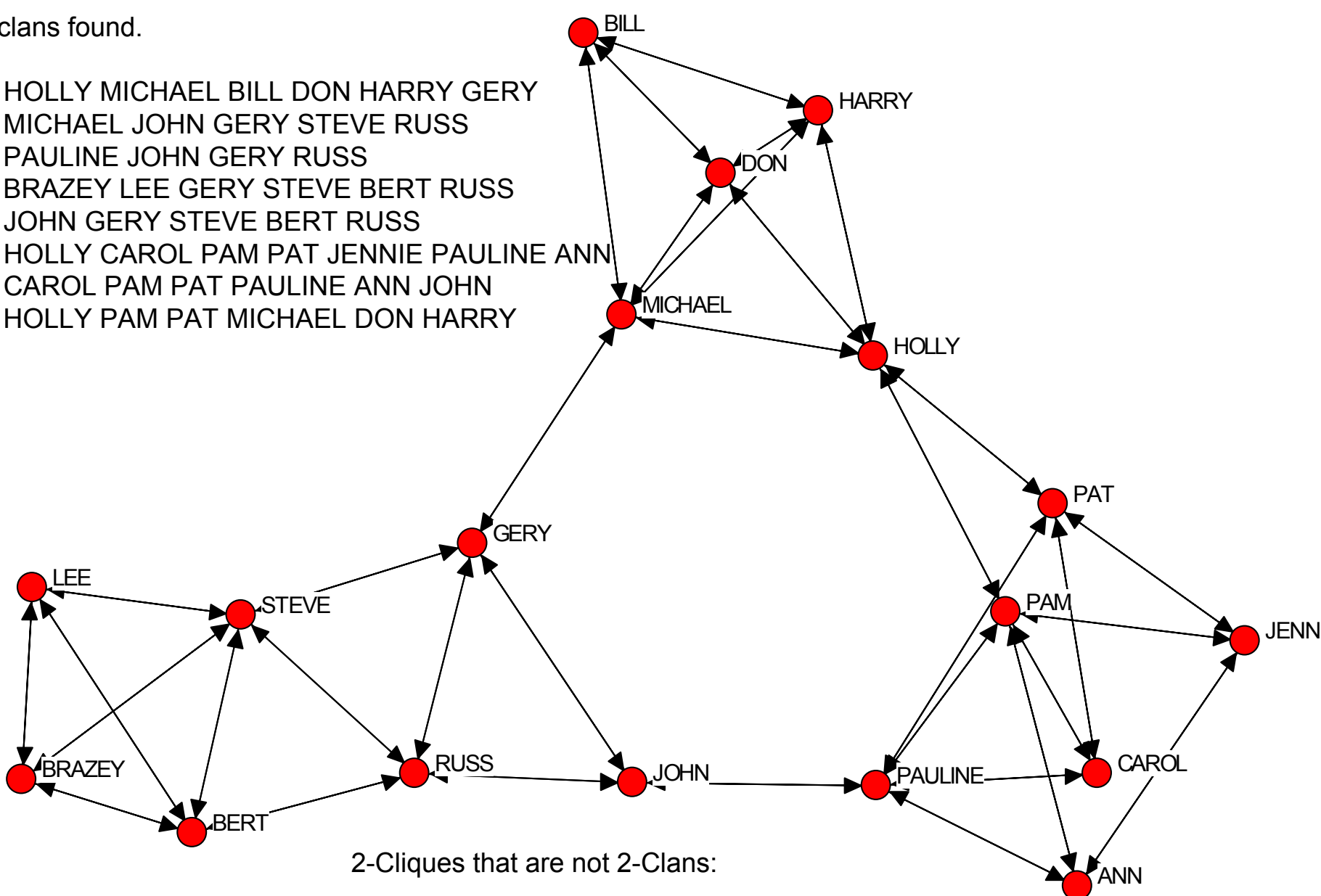
- Definition
 - An n -clique S whose diameter in the subgraph induced by S is $\leq n$
 - Members of set within n links of each other without using outsiders
- Properties
 - More cohesive than n -cliques



Is $\{a,b,c,f,e\}$ a 2-clan?
25

8 2-clans found.

- 1: HOLLY MICHAEL BILL DON HARRY GERY
- 2: MICHAEL JOHN GERY STEVE RUSS
- 3: PAULINE JOHN GERY RUSS
- 5: BRAZEY LEE GERY STEVE BERT RUSS
- 6: JOHN GERY STEVE BERT RUSS
- 7: HOLLY CAROL PAM PAT JENNIE PAULINE ANN
- 8: CAROL PAM PAT PAULINE ANN JOHN
- 9: HOLLY PAM PAT MICHAEL DON HARRY

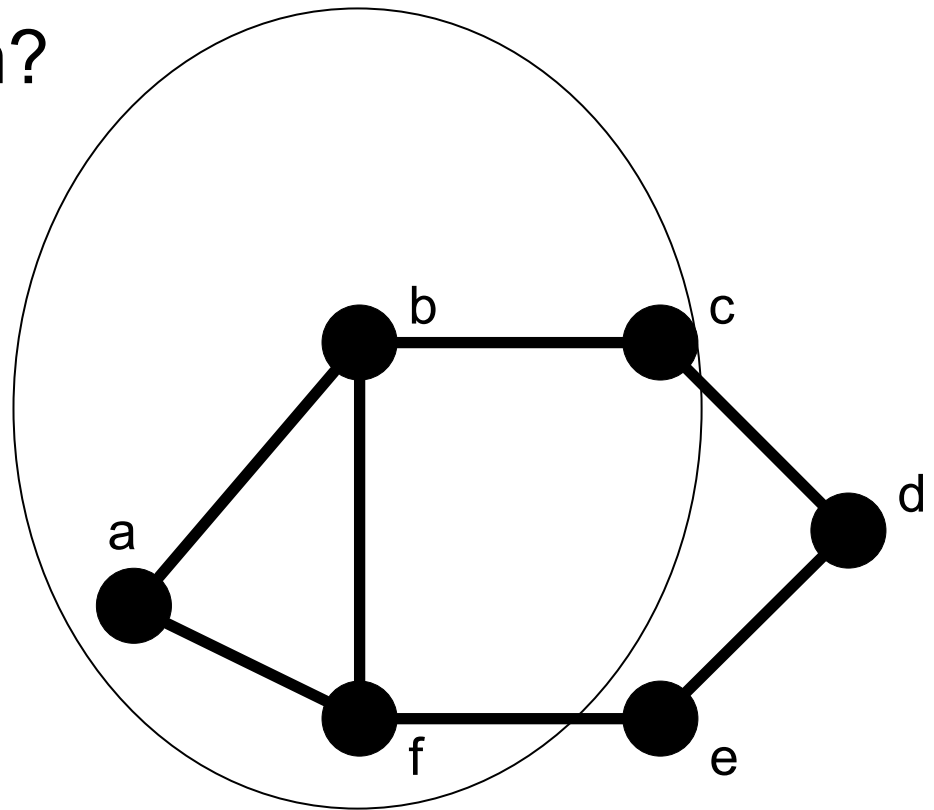


2-Cliques that are not 2-Clans:

- 4: HOLLY PAULINE GERY
- 10: PAM PAT MICHAEL JOHN

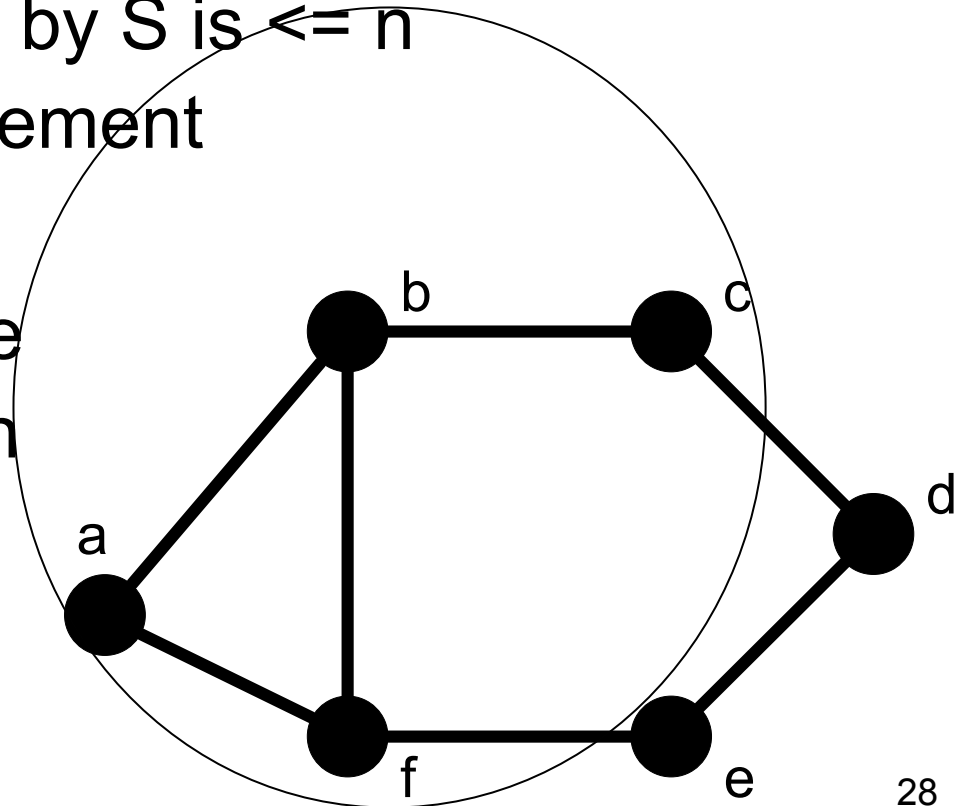
N-Clan Issues

- n-clique membership a bother
 - Is $\{a,b,c,f\}$ a 2-clan?
 - List all 2-clans
- few found in data
- overlapping

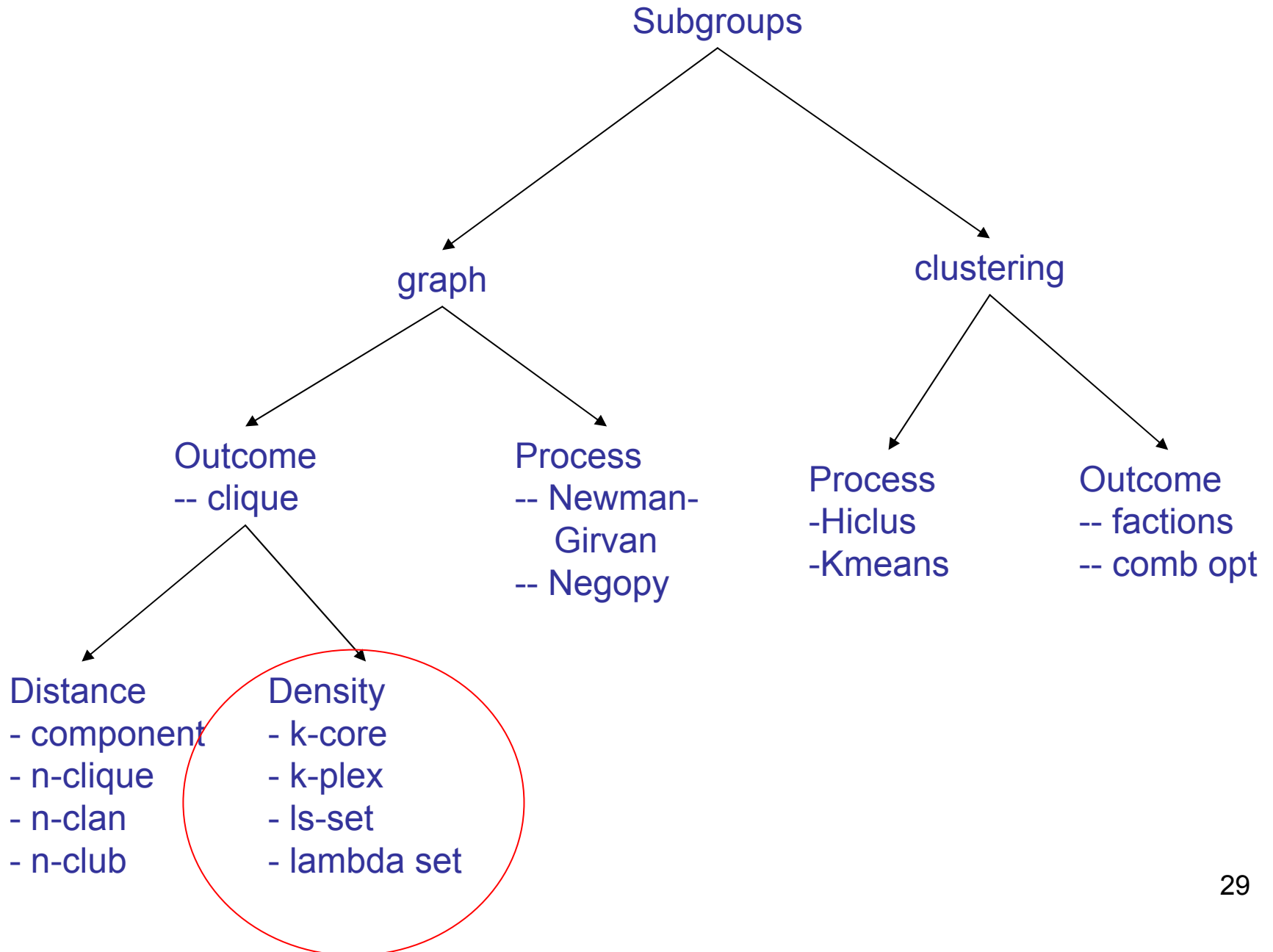


N-Club

- Definition
 - A maximal subset S whose diameter in the subgraph induced by S is $\leq n$
 - No n -clique requirement
- Properties
 - Painful to compute
 - More plentiful than n -clans
 - overlapping



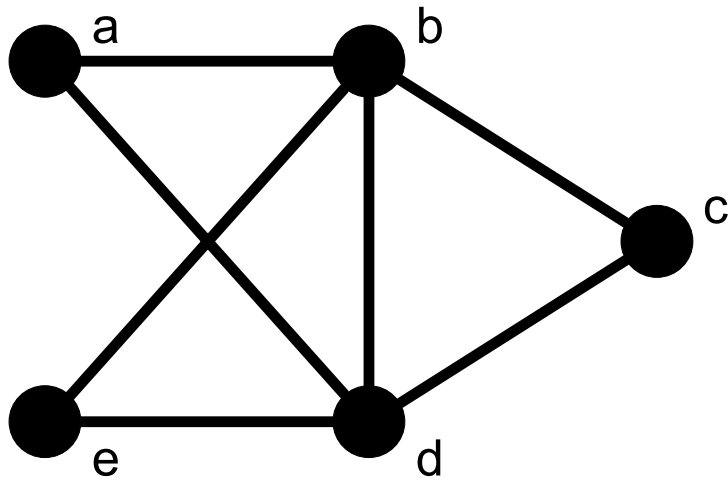
Is $\{a,b,c,f\}$ a 2-club?



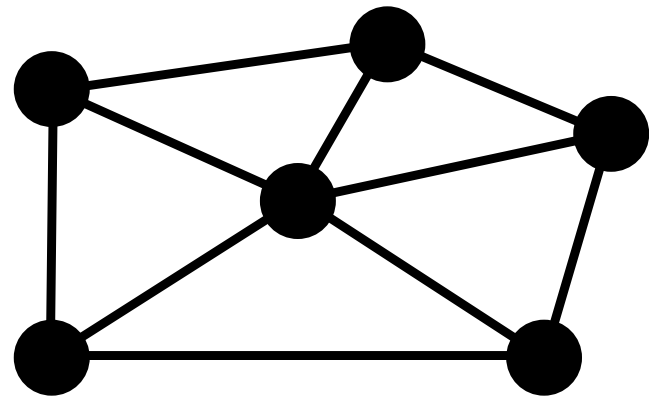
K-Plexes

- Definition
 - A k-plex is a [maximal] subset S s.t. for all u in S , $\alpha(u, S) \geq |S| - k$, where $|S|$ is size of set S
- Properties
 - Subsets of k-plexes are k-plexes
 - Limited diameter (i.e., get distance as freebie)
 - If $k < (|S| + 2)/2$ then diameter ≤ 2
 - Very numerous & overlapping
 - Sometimes better match to intuition than distance relaxations

K-Plex



Is $\{a, b, d, e\}$ a 2-plex?
Is $\{a, b, c, d, e\}$ a 2-plex?
Is $\{a, b, d\}$ a 2-plex?



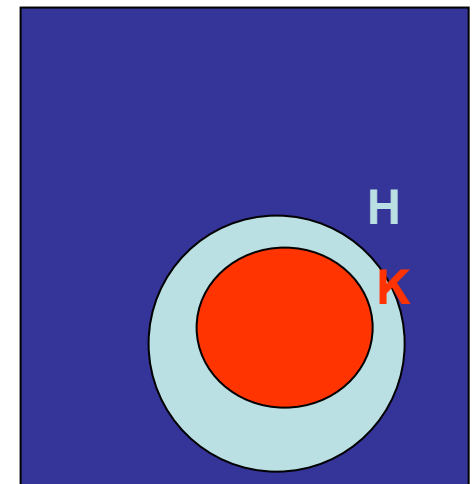
Is the graph as a whole a 2-plex?
Is it a 3-plex?

LS-Sets

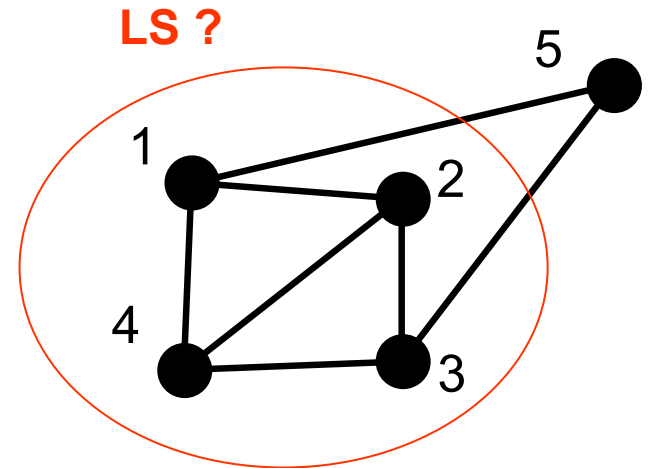
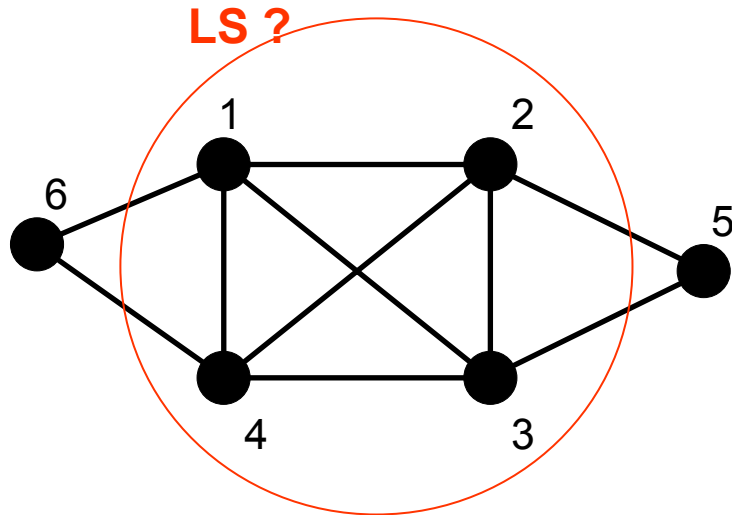
- Definition
 - Given a graph $G(V,E)$, let H be a subset of V , and let K be any proper subset of H
 - H is **LS** if $\alpha(K, H-K) > \alpha(K, V-H)$ for all K
 - All subsets of the LS set are more connected to other LS members than outsiders of LS set

or...

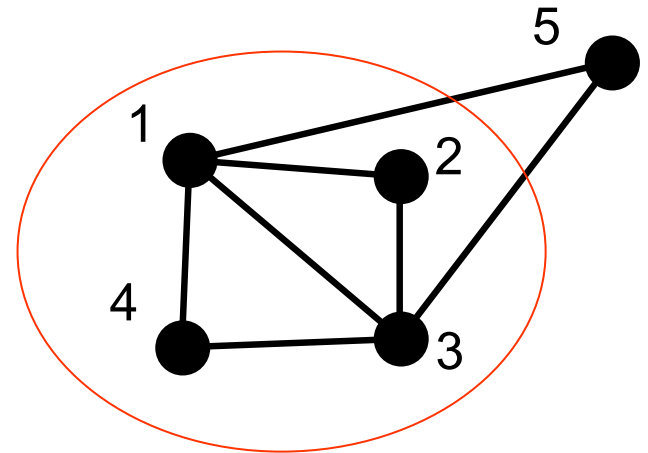
- H is **LS** if $\alpha(K) > \alpha(H)$
 - Subsets better off joining LS set
 - This one's usually easier to compute



LS-Sets



- H is LS if $\alpha(K, H-K) > \alpha(K, V-H)$
 - Use when K is large
- or ...
- H is LS if $\alpha(K) > \alpha(H)$
 - Use when K is small



LS-Sets

- Properties – very cohesive
 - Wholly nested or disjoint: no partial overlaps
 - More ties within than between (doesn't just consider density inside density)
 - Contain no minimum weight cutsets (lie on either side of “fault lines”)
 - Multiple edge-independent paths within
 - High edge-connectivity

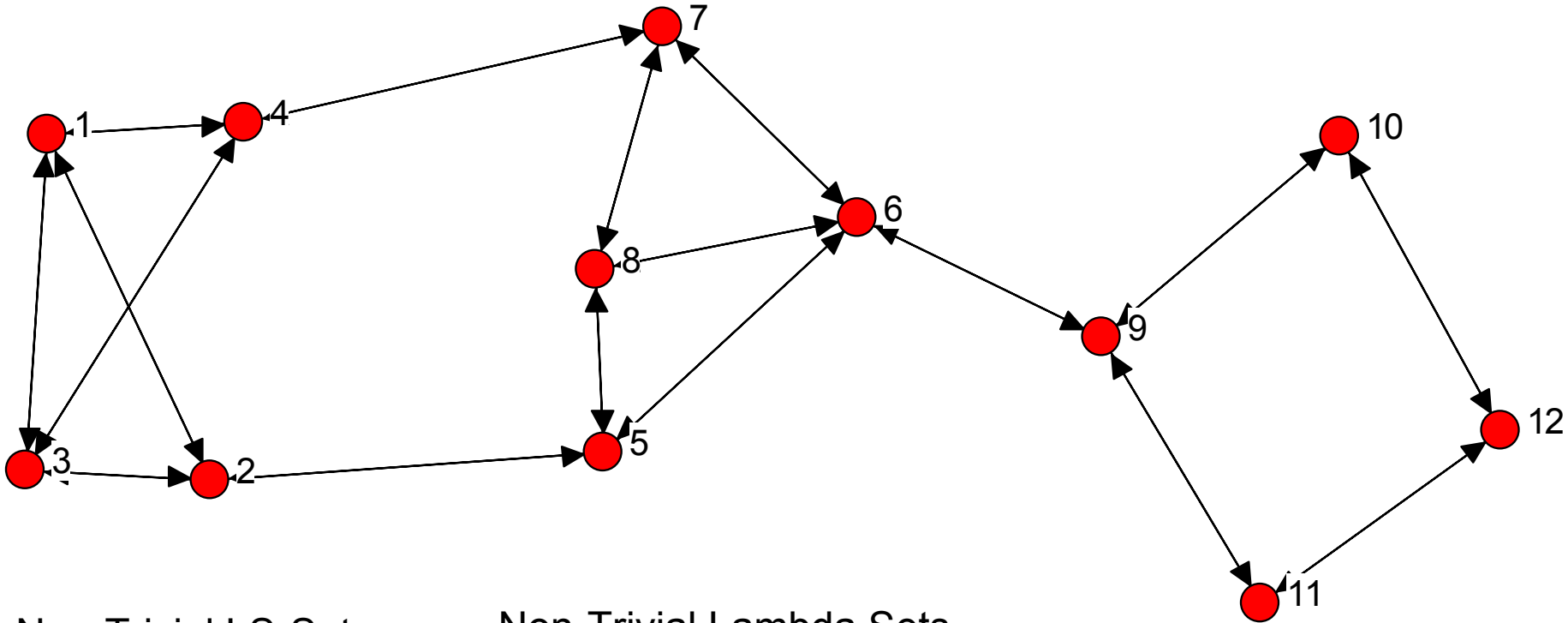
Lambda Operator

- Let $\lambda(u,v)$ be the number of edge-independent paths from node u to node v
- $\lambda(u,v)$ is also the minimum number of ties that must be removed from the network in order to disconnect u and v

Lambda Sets

- Definition
 - A set of nodes S is a lambda set if for all a, b, c in S and d not in S , $\lambda(a, b) > \lambda(c, d)$
 - More independent paths to other group members than to outsiders
- Properties
 - Robust
 - very difficult to disconnect even with intelligent attack
 - Mutually exclusive or wholly inclusive
 - No partially overlapping groups
 - Pure – like n-clubs, defined on a single attribute

Lambda Sets

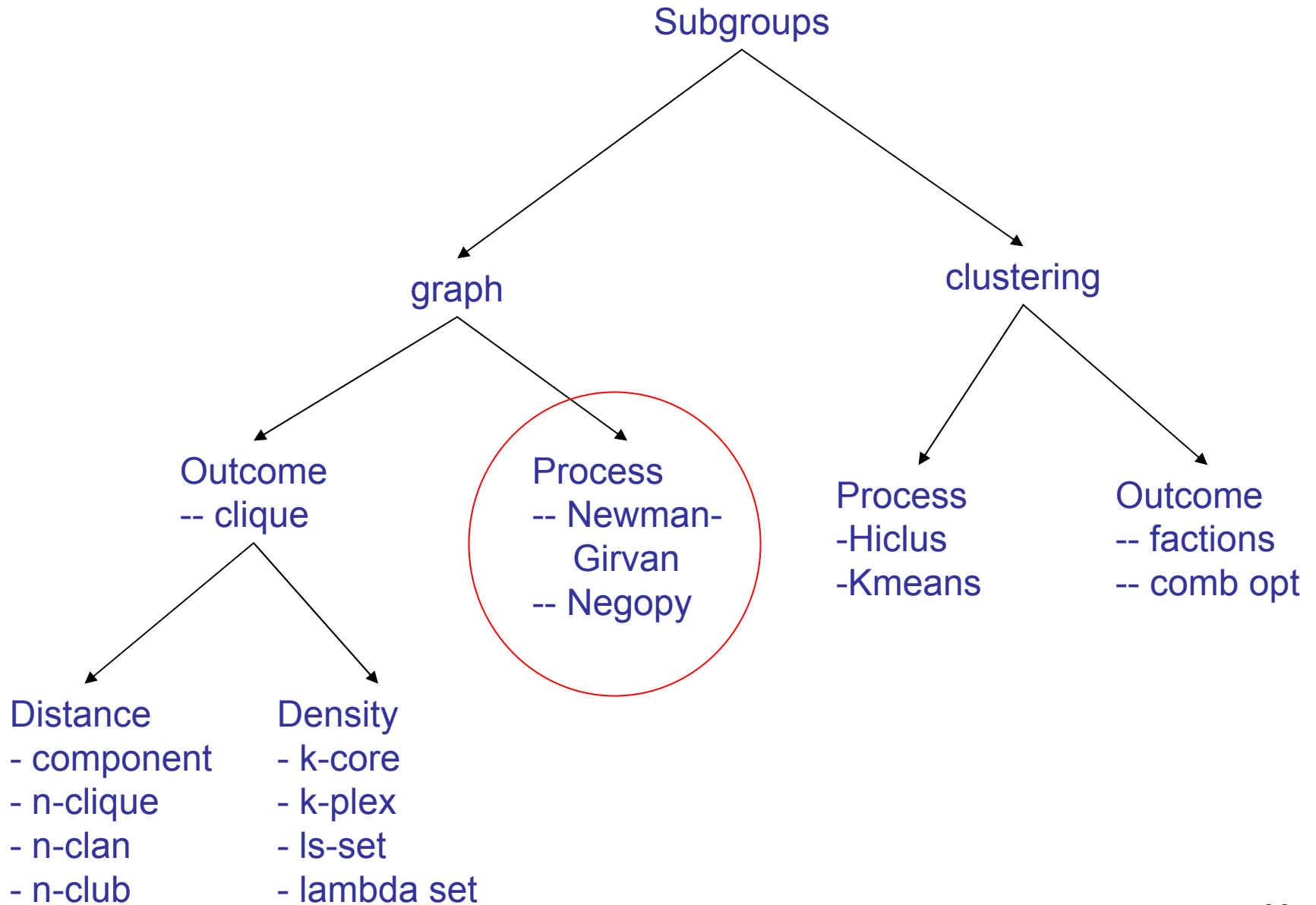


Non-Trivial LS-Sets

$\{1,2,3,4\}$
 $\{1,2,3,4,5,6,7,8\}$
 $\{9,10,11,12\}$

Non-Trivial Lambda Sets

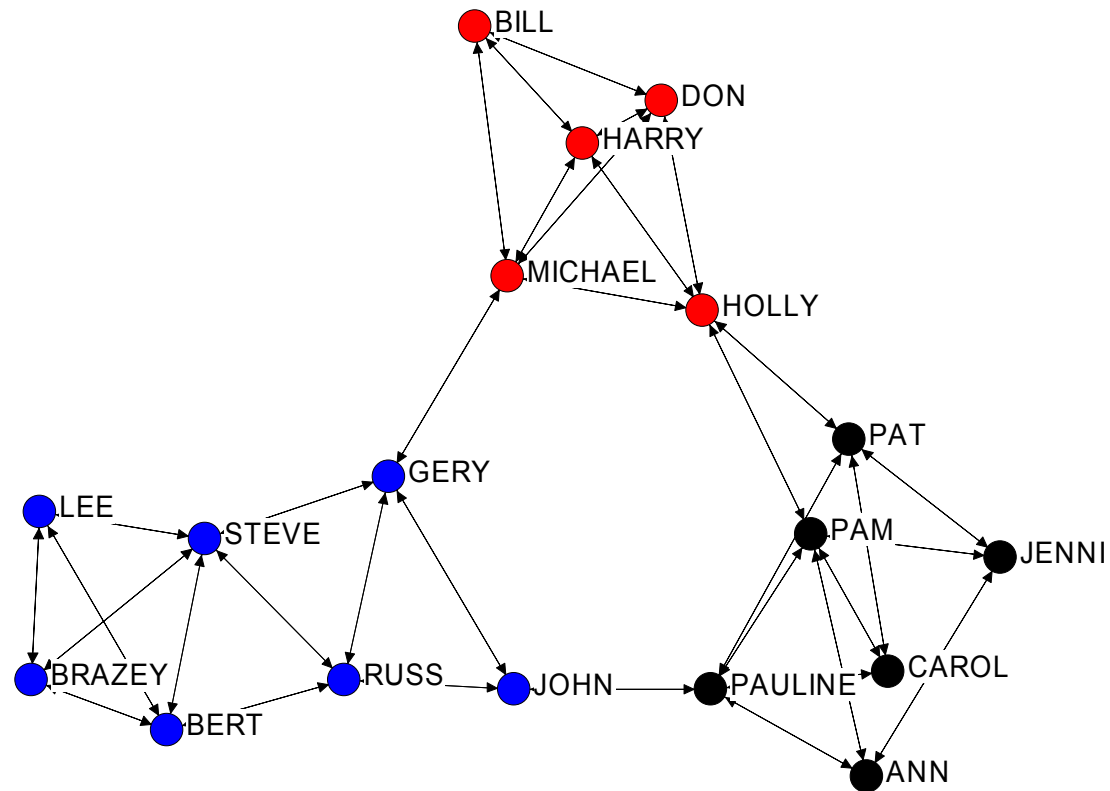
$\{1,2,3,4\}$
 $\{1,2,3,4,5,6,7,8\}$
 $\{9,10,11,12\}$
 $\{5,6,7,8\}$



Newman-Girvan

- Successively deleting the tie with the most edge betweenness, and identifying components, then recalculating betweenness
- Yields a hierarchical clustering

Newman-Girvan



Proximities / Clustering and Scaling Methods

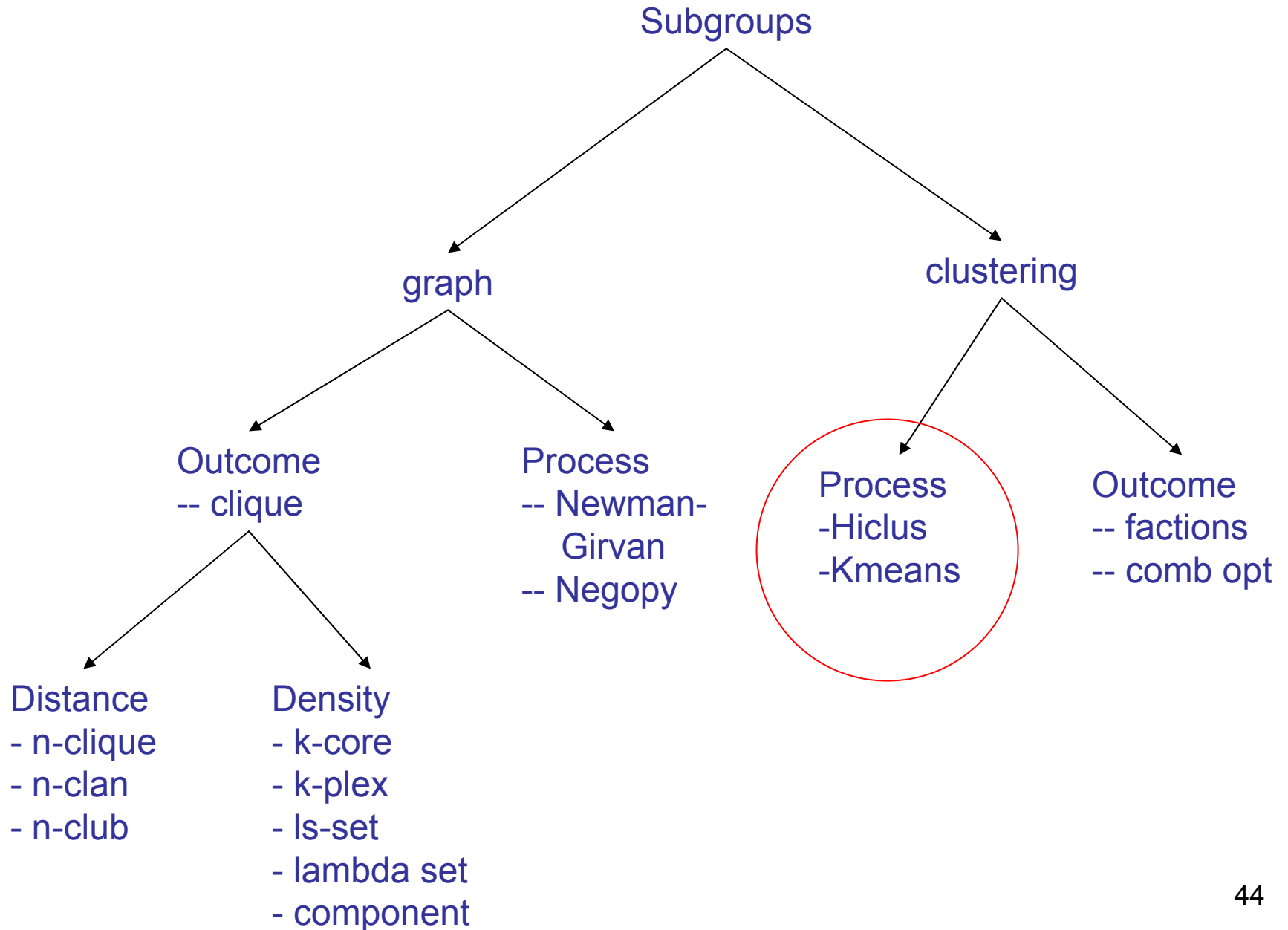
- First compute dyadic cohesion matrix
 - E.g. geodesic distance
- Then cluster or scale
 - Two major kinds of clustering routines
 - Process-defined
 - Outcome-defined
- Typical result is a partition

Partitions

- Partition P is just an assignment of nodes to classes
 - $P(i)$ gives the class of node i
 - Every node assigned to one & only one class
- A partition P is nested in partition M if for all nodes i and j , $P(i)=P(j)$ implies $M(i)=M(j)$
- Trivial partitions
 - Identity: $P(i) = i$ for all i
 - Complete: $P(i) = 1$ for all i

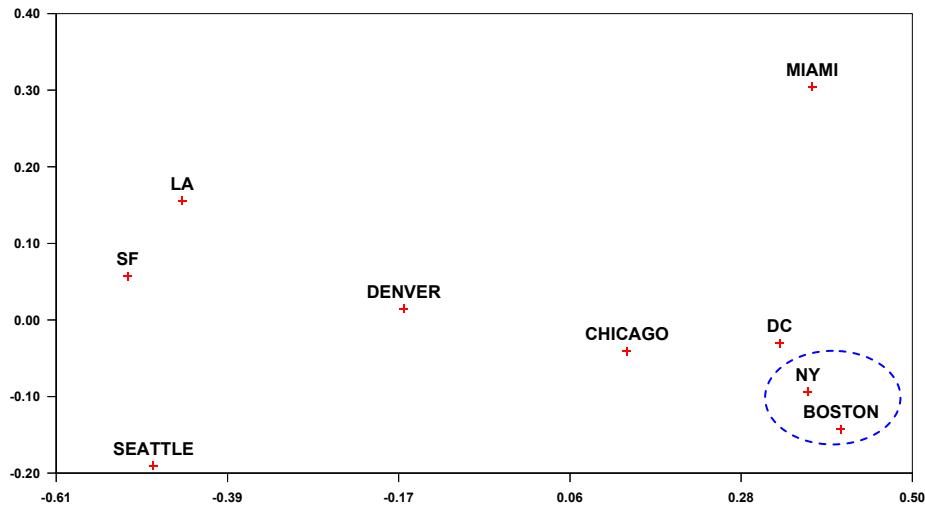
Process-Defined Clustering

- Heuristic definitions
 - Multivariate methods
 - Johnson's hierarchical
 - Wards
 - K-means
 - Graph-theoretic / Network methods
 - Newman-Girvan
- Sometimes specify number of groups a priori, sometimes not



Johnson's Hierarchical Clustering

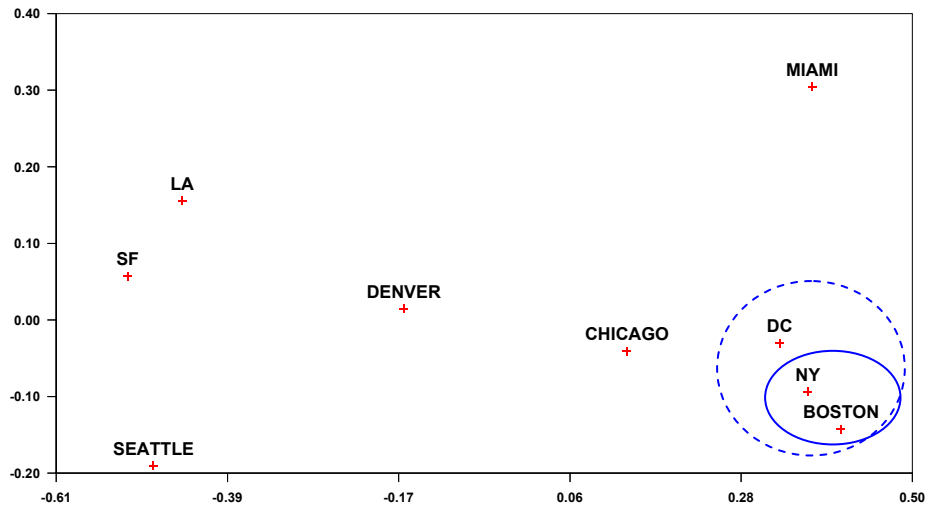
- Output is a set of nested partitions, starting with identity partition and ending with the complete partition
- Different flavors based on how distance from a point to a cluster is defined
 - Single linkage; connectedness; minimum
 - Complete linkage; diameter; maximum
 - Average, median, etc.



	M	S		B		C	D
	I	E	L	O	N	D	H
	A	A	F	A	S	C	I
Level	4	6	7	8	1	2	3
206	-	-	-	-	XXX	-	-
233	-	-	-	-	XXXXX	-	-
379	-	-	XXX	XXXXX	-	-	-
671	-	-	XXX	XXXXXXXX	-	-	-
808	-	XXXXX	XXXXXXXX	-	-	-	-
996	-	XXXXX	XXXXXXXXXX	-	-	-	-
1059	-	XXXXXXXXXXXXXXXXXX	-	-	-	-	-
1075	-	XXXXXXXXXXXXXXXXXX	-	-	-	-	-

	BOS	NY	DC	MIA	CHI	SEA	SF	LA	DEN
BOS	0	206	429	1504	963	2976	3095	2979	1949
NY	206	0	233	1308	802	2815	2934	2786	1771
DC	429	233	0	1075	671	2684	2799	2631	1616
MIA	1504	1308	1075	0	1329	3273	3053	2687	2037
CHI	963	802	671	1329	0	2013	2142	2054	996
SEA	2976	2815	2684	3273	2013	0	808	1131	1307
SF	3095	2934	2799	3053	2142	808	0	379	1235
LA	2979	2786	2631	2687	2054	1131	379	0	1059
DEN	1949	1771	1616	2037	996	1307	1235	1059	0

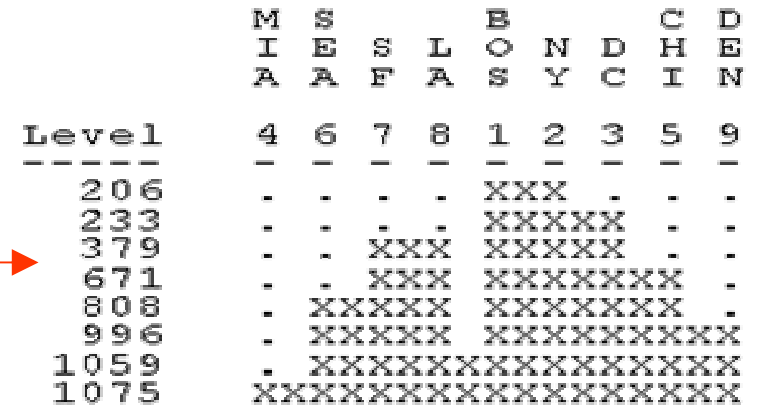
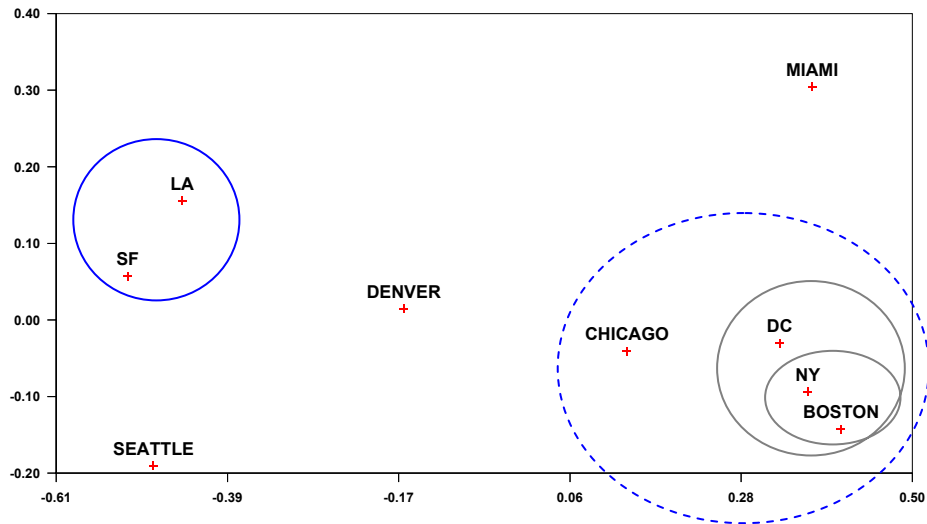
Closest distance is NY-BOS = 206, so merge these.



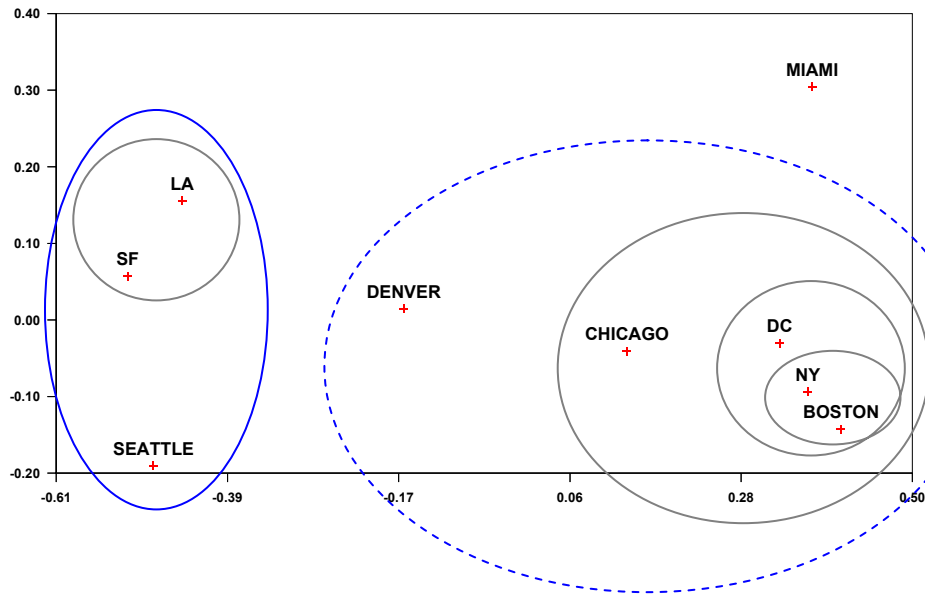
	M	S	L	B	O	N	D	H	E
	A	A	F	A	S	Y	C	I	N
Level	4	6	7	8	1	2	3	5	9
206	-	-	-	-	XXX	-	-	-	-
233	-	-	-	-	XXXXX	-	-	-	-
379	-	-	XXX	XXXXX	-	-	-	-	-
671	-	-	XXX	XXXXXXXX	-	-	-	-	-
808	-	XXXXX	XXXXXXXX	-	-	-	-	-	-
996	-	XXXXX	XXXXXXXX	-	-	-	-	-	-
1059	-	XXXXXXXXXXXXXXXXXXXX	-	-	-	-	-	-	-
1075	XXXXXXXXXXXXXXXXXXXX	-	-	-	-	-	-	-	-

	BOS N Y	DC	MIA	CHI	SEA	SF	LA	DEN
BOS/ NY	0	223	1308	802	2815	2934	2786	1771
DC	223	0	1075	671	2684	2799	2631	1616
MIA	1308	1075	0	1329	3273	3053	2687	2037
CHI	802	671	1329	0	2013	2142	2054	996
SEA	2815	2684	3273	2013	0	808	1131	1307
SF	2934	2799	3053	2142	808	0	379	1235
LA	2786	2631	2687	2054	1131	379	0	1059
DEN	1771	1616	2037	996	1307	1235	1059	0

Closest pair
is DC to
BOSNY
combo @
223. So
merge these.



	BOS/ NY/DC	MIA	CHI	SEA	SF/LA	DEN
BOS/NY/DC	0	1075	671	2684	2631	1616
MIA	1075	0	1329	3273	2687	2037
CHI	671	1329	0	2013	2054	996
SEA	2684	3273	2013	0	808	1307
SF/LA	2631	2687	2054	808	0	1059
DEN	1616	2037	996	1307	1059	0



Level

```

-----
206
233
379
671
808
996
1059
1075

```

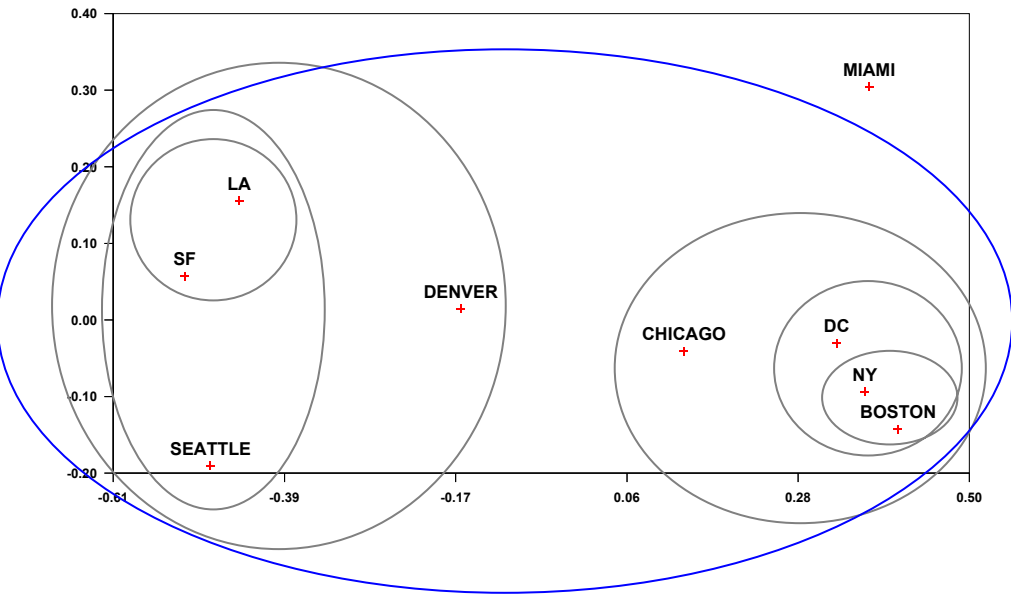
```

M S B C D
I E S L O N D H E
A A F A S Y C I N

4 6 7 8 1 2 3 5 9
- - - - -
206 . . . . XXX . . .
233 . . . . XXXXX . .
379 . . XXX XXXXX . .
671 . . XXX XXXXXXXX .
808 . XXXXX XXXXXXXX .
996 . XXXXX XXXXXXXXXXXX
1059 . XXXXXXXXXXXXXXXXXXXX
1075 XXXXXXXXXXXXXXXXXXXX

```

	BOS/ NY/D C/C HI	MIA	SF/L A/SE A	DEN
BOS/NY/DC/ CHI	0	1075	2013	996
MIA	1075	0	2687	2037
SF/LA/SEA	2054	2687	0	1059
DEN	996	2037	1059	0



	M	S		B		C	D
	I	E	S	L	O	N	H
	A	A	F	A	S	Y	C
Level	4	6	7	8	1	2	3
---	---	---	---	---	---	---	---
206	XXX	.	.
233	XXXXX	.	.
379	.	.	XXX	XXXXX	.	.	.
671	.	.	XXX	XXXXXXXX	.	.	.
808	.	XXXXX	XXXXXXXX
996	.	XXXXX	XXXXXXXXXX
1059	.	XXXXXXXXXXXXXXXXXX
1075	XXXXXXXXXXXXXXXXXX

	BOS/ NY/D C/CH I/DE N/SF/ LA/S EA	MIA
BOS/NY/DC/CHI/DEN/SF/L A/SEA	0	1075
MIA	1075	0

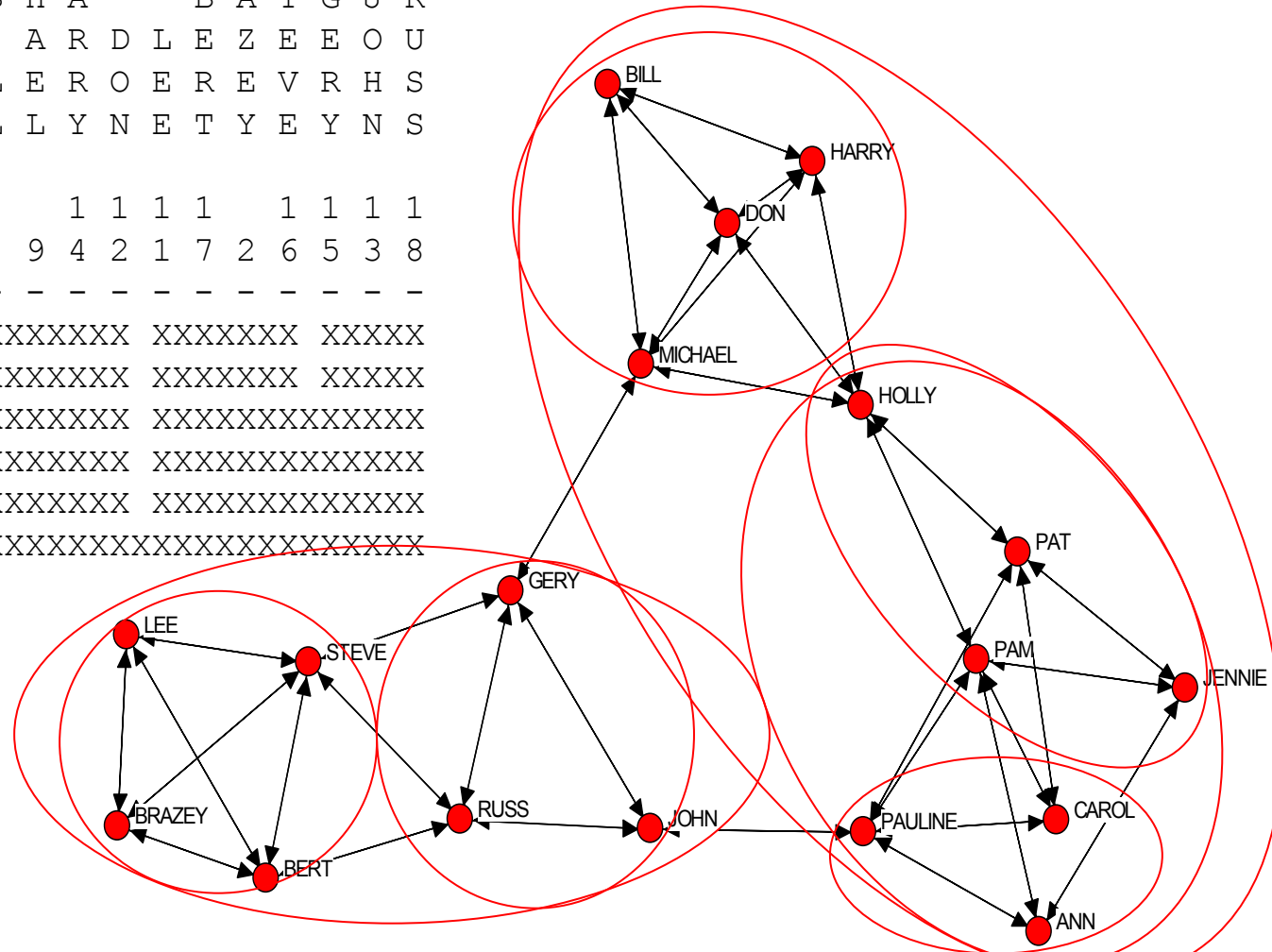
Geodesic Distances

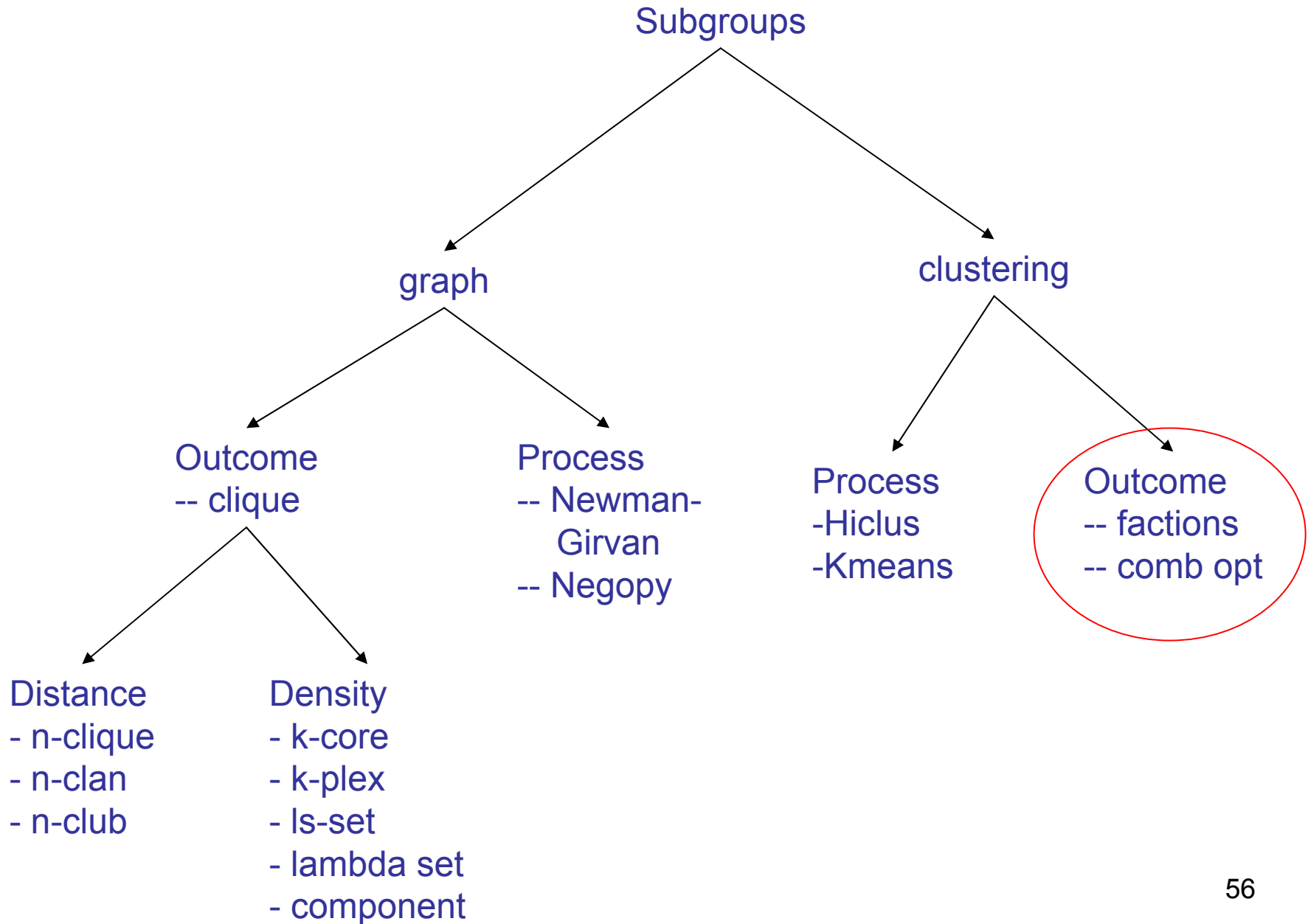
		1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8
		H	B	C	P	P	J	P	A	M	B	L	D	J	H	G	S	B	R
		-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
1	HOLLY	0	4	2	1	1	2	2	2	1	2	4	1	3	1	2	3	4	3
2	BRAZEY	4	0	5	5	5	6	4	5	3	4	1	4	3	4	2	1	1	2
3	CAROL	2	5	0	1	1	2	1	2	3	4	5	3	2	3	3	4	4	3
4	PAM	1	5	1	0	2	1	1	1	2	3	5	2	2	2	3	4	4	3
5	PAT	1	5	1	2	0	1	1	2	2	3	5	2	2	2	3	4	4	3
6	JENNIE	2	6	2	1	1	0	2	1	3	4	6	3	3	3	4	5	5	4
7	PAULINE	2	4	1	1	1	2	0	1	3	4	4	3	1	3	2	3	3	2
8	ANN	2	5	2	1	2	1	1	0	3	4	5	3	2	3	3	4	4	3
9	MICHAEL	1	3	3	2	2	3	3	3	0	1	3	1	2	1	1	2	3	2
10	BILL	2	4	4	3	3	4	4	4	1	0	4	1	3	1	2	3	4	3
11	LEE	4	1	5	5	5	6	4	5	3	4	0	4	3	4	2	1	1	2
12	DON	1	4	3	2	2	3	3	3	1	1	4	0	3	1	2	3	4	3
13	JOHN	3	3	2	2	2	3	1	2	2	3	3	3	0	3	1	2	2	1
14	HARRY	1	4	3	2	2	3	3	3	1	1	4	1	3	0	2	3	4	3
15	GERY	2	2	3	3	3	4	2	3	1	2	2	2	1	2	0	1	2	1
16	STEVE	3	1	4	4	4	5	3	4	2	3	1	3	2	3	1	0	1	1
17	BERT	4	1	4	4	4	5	3	4	3	4	1	4	2	4	2	1	0	1
18	RUSS	3	2	3	3	3	4	2	3	2	3	2	3	1	3	1	1	1	0

Hierarchical Clustering

	P							M											
	A		J					I						B					
	C	U		H	E			C	H					R	S				
	A	L		O	N			B	H	A				B	A	T	G	J	R
P	R	I	P	L	N	A	I	A	R	D	L	E	Z	E	E	O	U		
A	O	N	A	L	I	N	L	E	R	O	E	R	E	V	R	H	S		
T	L	E	M	Y	E	N	L	L	Y	N	E	T	Y	E	Y	N	S		

									1		1	1	1	1		1	1	1	1
Level	5	3	7	4	1	6	8	0	9	4	2	1	7	2	6	5	3	8	
-----	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
1.000	XXXXX	XXX	XXX					XXXXXXXX		XXXXXXXX		XXXXXXXX		XXXXX					
1.333	XXXXX	XXXXXXXX						XXXXXXXX		XXXXXXXX		XXXXXXXX		XXXXX					
1.457	XXXXX	XXXXXXXX						XXXXXXXX		XXXXXXXXXXXXXXXXXX									
1.481	XXXXXXXXXXXXXXXXXX							XXXXXXXX		XXXXXXXXXXXXXXXXXX									
2.723	XXXXXXXXXXXXXXXXXXXXXXXXXXXX							XXXXXXXXXXXXXXXXXX		XXXXXXXXXXXXXXXXXX									
3.142	XXXXXXXXXXXXXXXXXXXXXXXXXXXX							XXXXXXXXXXXXXXXXXX		XXXXXXXXXXXXXXXXXX									





Factions

- Outcome-Defined Clustering
- Input is proximity matrix X
 - Could be similarities or distances
- Assign items to clusters such that
 - For similarities, maximize similarities within cluster while minimizing similarities between clusters
 - For distances, minimize distance within cluster while maximizing distances between clusters
- Optimize explicit fitness function
 - Correlation with idealized image matrix
- Typically choose # of groups *a priori*

Factions

		5	6	3	4	7	8	1	1	1	1	1	1	1	1	1	1	1	1
		P	J	C	P	P	A	R	B	J	S	L	B	G	B	H	D	H	M
5	PAT		1	1		1											1		
6	JENNIE		1			1	1												
3	CAROL		1			1	1												
4	PAM			1	1		1	1									1		
7	PAULINE		1		1	1	1			1									
8	ANN			1		1	1												
18	RUSS									1	1		1	1					
2	BRAZEY										1	1	1						
13	JOHN					1		1						1					
16	STEVE							1	1			1	1	1					
11	LEE								1		1		1						
17	BERT							1	1		1	1							
15	GERY							1		1	1							1	
10	BILL															1	1		1
14	HARRY															1		1	1
12	DON															1	1		1
1	HOLLY		1			1											1	1	
9	MICHAEL												1		1	1	1	1	