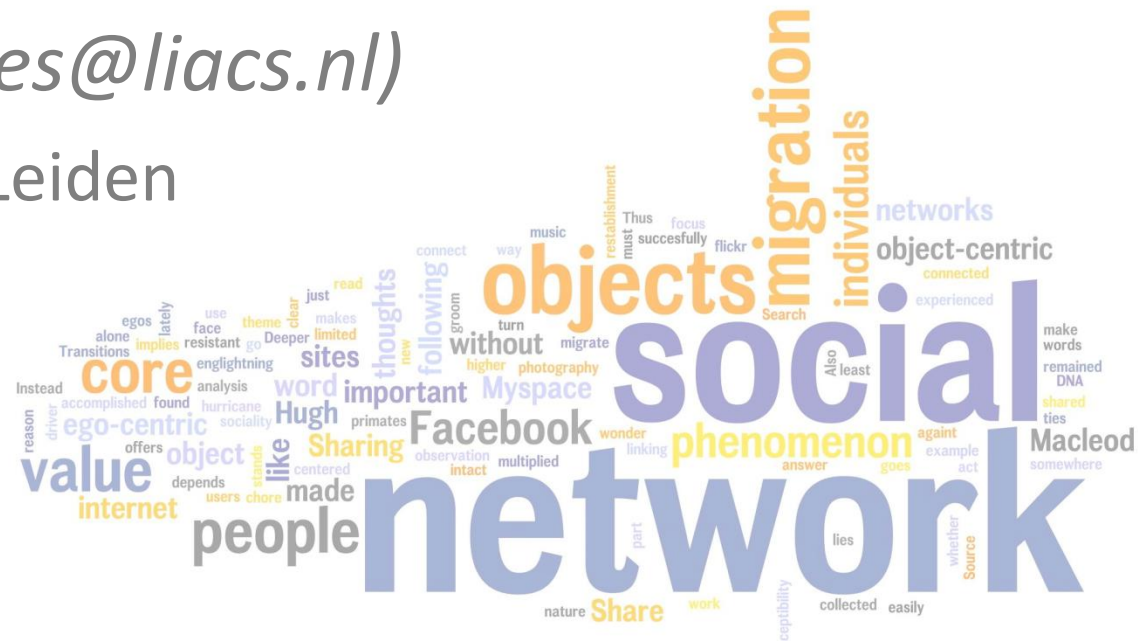




Data Mining – November 11, 2013

LIACS, Universiteit Leiden



Overview

- Social Network Analysis
- Graph Mining
- Online Social Networks
- Friendship Graph
- Semantics
- Example
- Conclusions

Social Network Analysis

- **Social Network Analysis:** the study of social networks to understand their structure and behavior.
- **Social Network:** a social structure of people, related (directly or indirectly) to each other through a common relation or interest.
- Social Networks != Social Media



Social Network Analysis

- **Social Network Analysis (SNA)**
 - Sociology
 - Algorithms
 - Data Mining
- **Social Networks**
 - Real-life (explicit)
 - Online (explicit)
 - Derived (implicit)
e-mail networks, citation networks, co-author networks, terrorist collaboration networks

SNA Research Topics

- Network analysis
 - Measuring
 - Modelling
 - Prediction
- Spread of Information
- Trust & Authority
- Community Detection
- Semantics
- Anonymity & Privacy



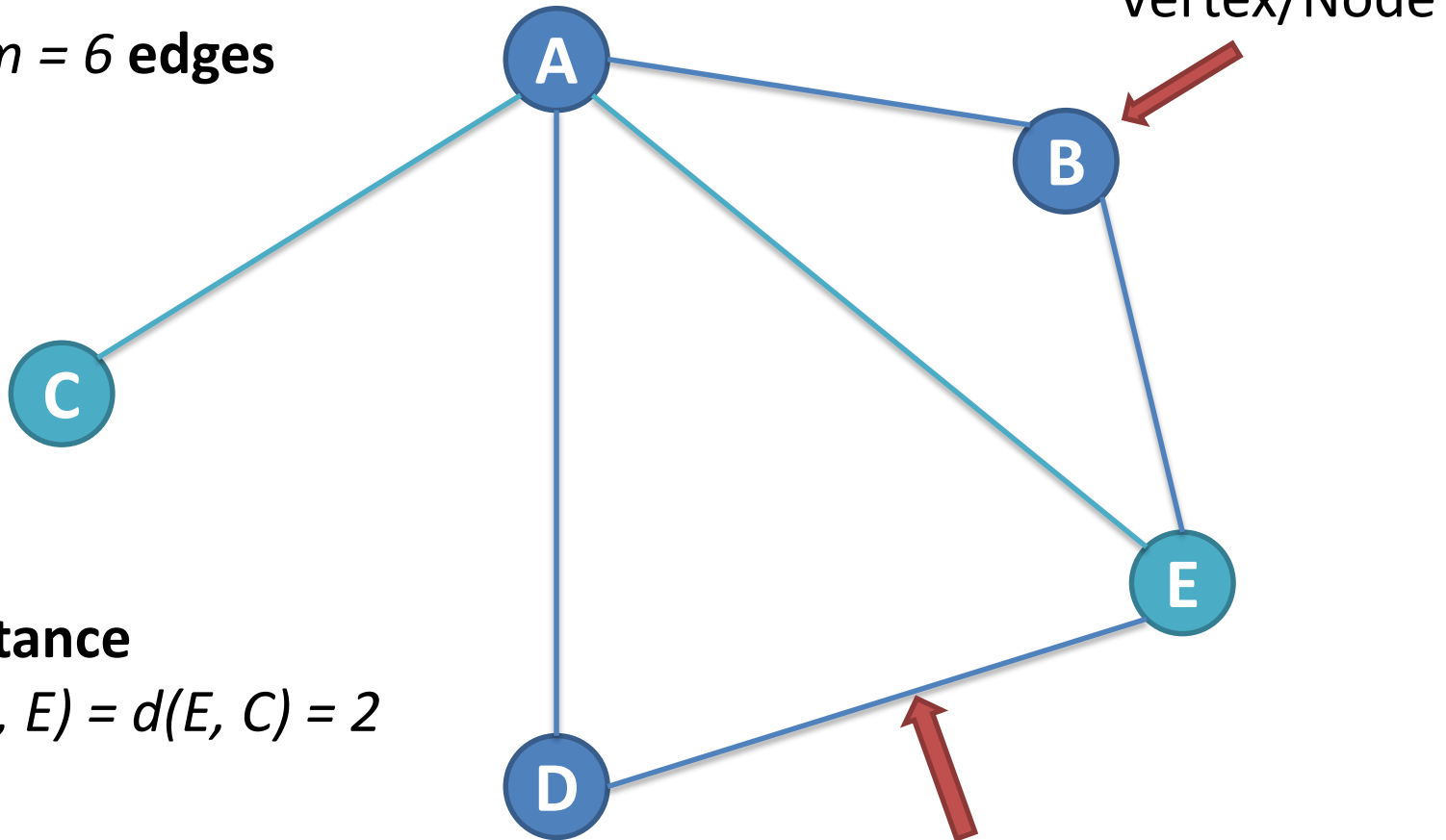
Graph Mining

Graphs

$n = 5$ nodes
 $m = 6$ edges

Distance

$$d(C, E) = d(E, C) = 2$$



Relationship/Edge/Link

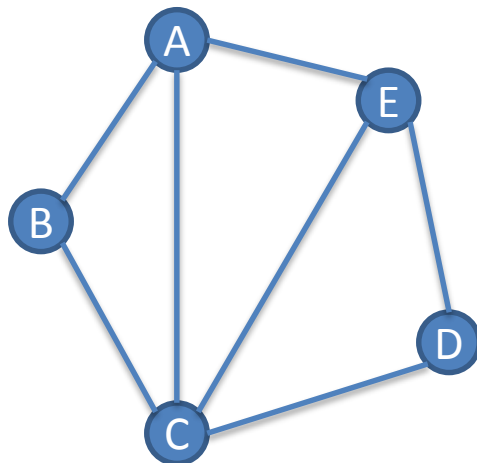
Data / Graph Mining

- **Data Mining:** looking for patterns in **Databases**

id	name	date_of_birth	email
23	B. Stinson	1983-08-15	goforbarney@me.com
42	T. Mosby	1983-08-15	mosbydesigns@gmail.com

classification, clustering, outlier detection, ...

- **Graph Mining:** looking for patterns in **Graphs**



So what are we looking for?

Graph Mining

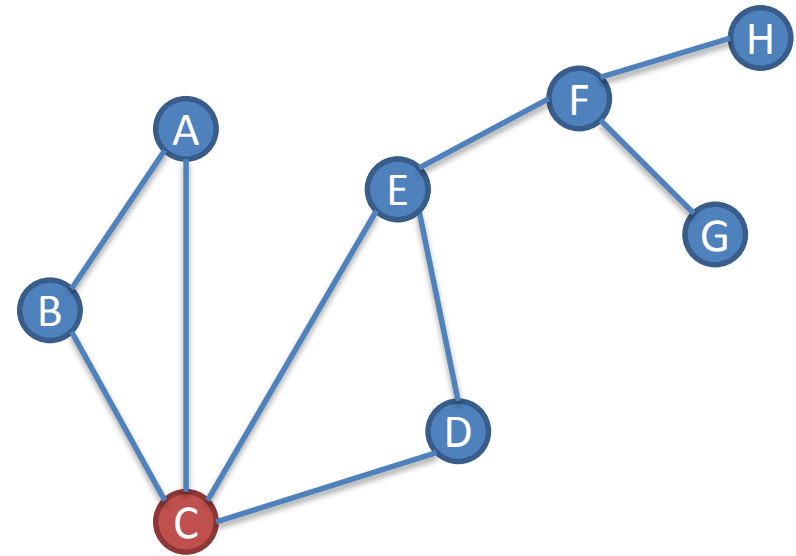
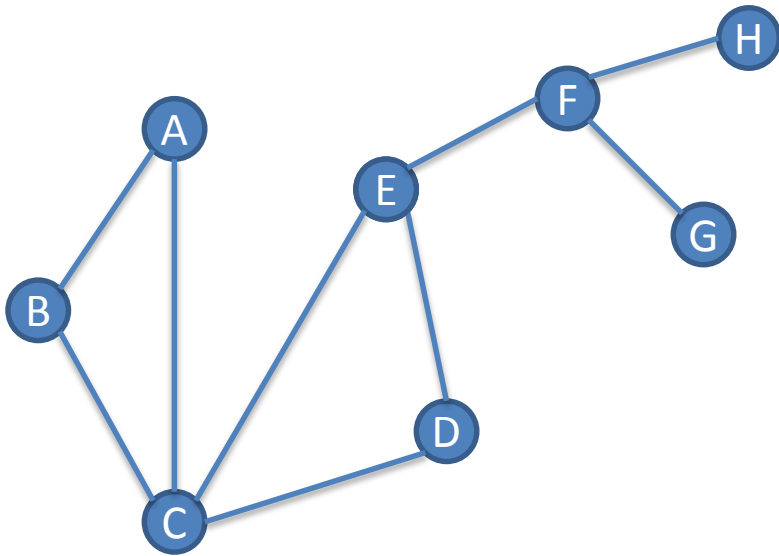
- Graph models
- Centrality measures (local)
- Graph properties (global)
- Frequent subgraphs (pattern mining)
- Detecting cliques (clustering)
- Link prediction (classification)
- Various application domains

Centrality measures

- Degree centrality
- Betweenness centrality
- Closeness centrality
- Graph centrality (eccentricity centrality)
- Eigenvector centrality
- Random walk centrality
- Hyperlink Induced Topic Search (HITS)
- PageRank

Centrality

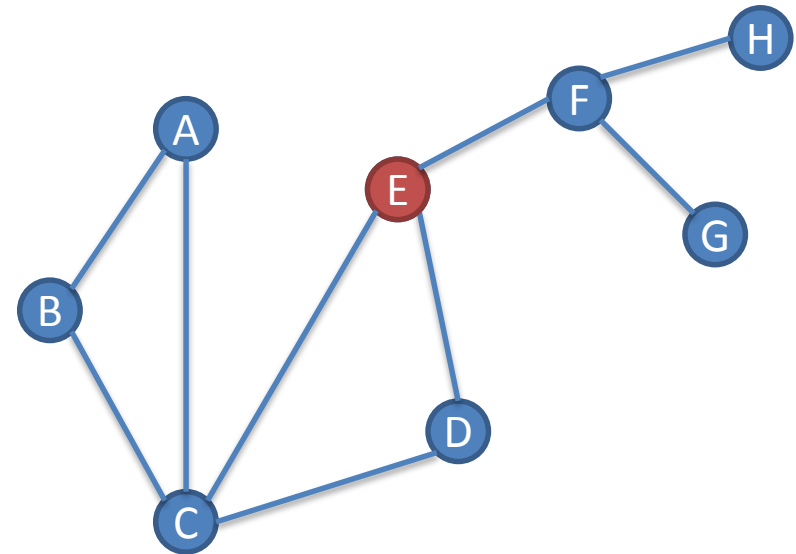
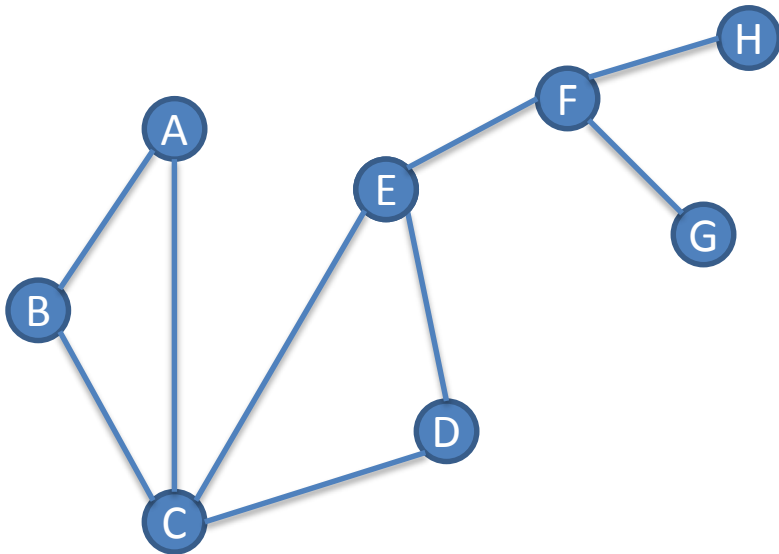
- Who has a central position in this graph?



Degree Centrality:
C has the highest degree

Centrality

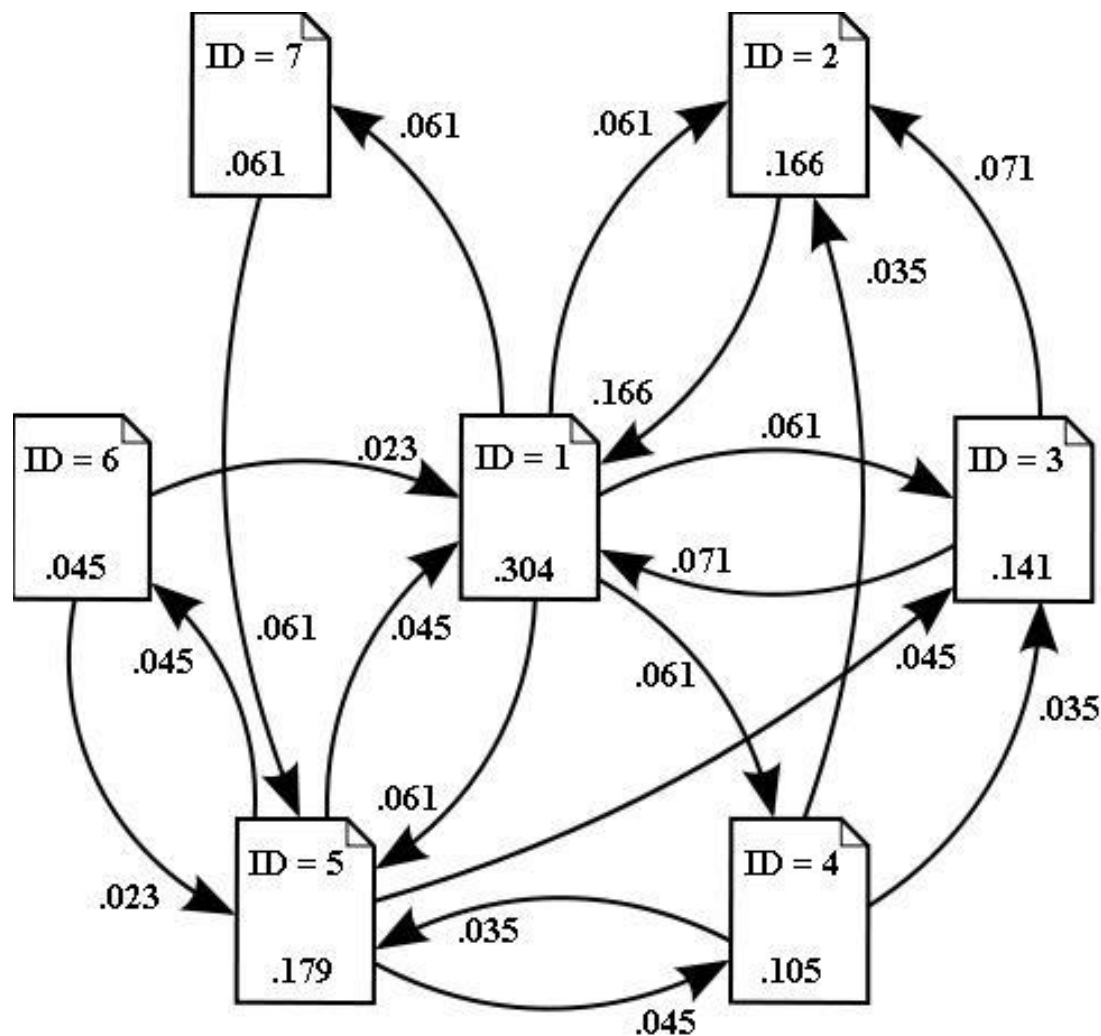
- Who has a central position in this graph?



Betweenness Centrality:

E is part of the largest
number of shortest paths

Google PageRank



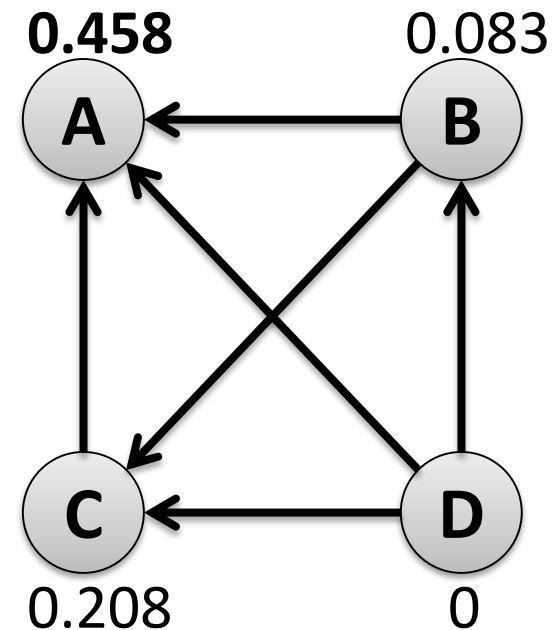
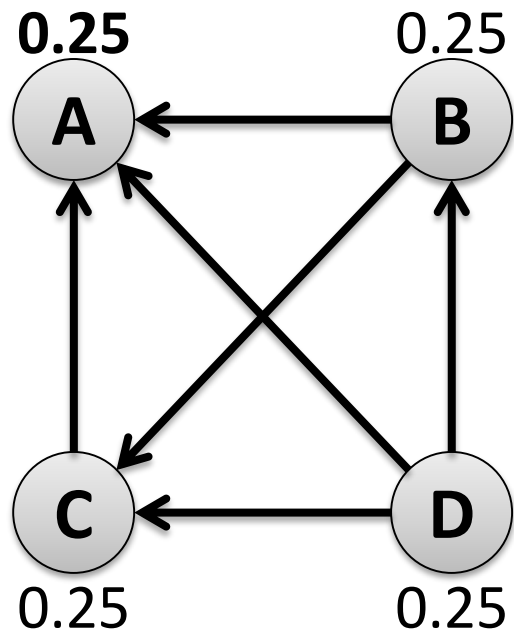
Google PageRank

- 4 webpages A , B , C and D ($N = 4$)
- Initially: $PR(A) = PR(B) = PR(C) = PR(D) = 1/n$
- $L(A)$ is the outdegree of page A
- Now if B , C and D each link to A , the simple PageRank $PR(A)$ of a page A is equal to:

$$PR(A) = \frac{PR(B)}{L(B)} + \frac{PR(C)}{L(C)} + \frac{PR(D)}{L(D)}$$

Google PageRank

$$PR(A) = \frac{PR(B)}{2} + \frac{PR(C)}{1} + \frac{PR(D)}{3}$$



Google PageRank

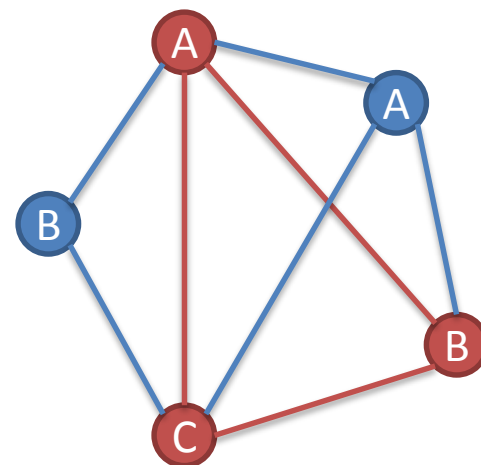
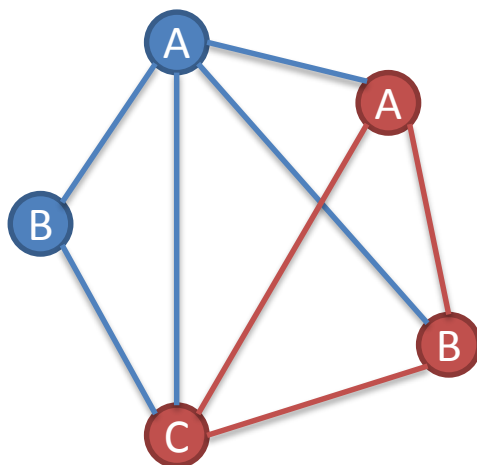
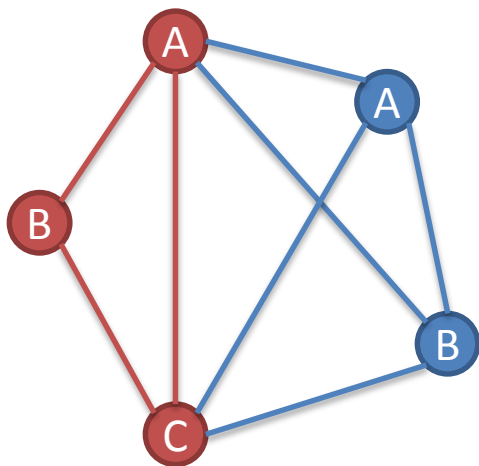
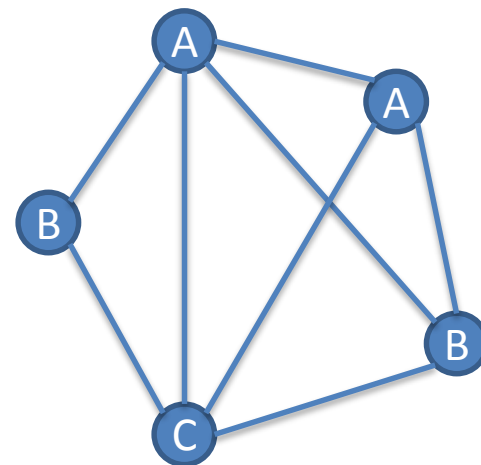
- **PageRank** as suggested by Larry Page in 1999
- N = number of pages, p_i and p_j are pages
- $M(p_i)$ is the set of pages linking to p_i
- $L(p_j)$ is the outdegree of p_j
- $d = 0.85$, 85% chance to follow a link, 15% chance to jump to a random page (random surfer)

- $t = 0$
$$PR(p_i; 0) = \frac{1}{N}.$$

- $t = t + 1$
$$PR(p_i; t + 1) = \frac{1 - d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j; t)}{L(p_j)}$$

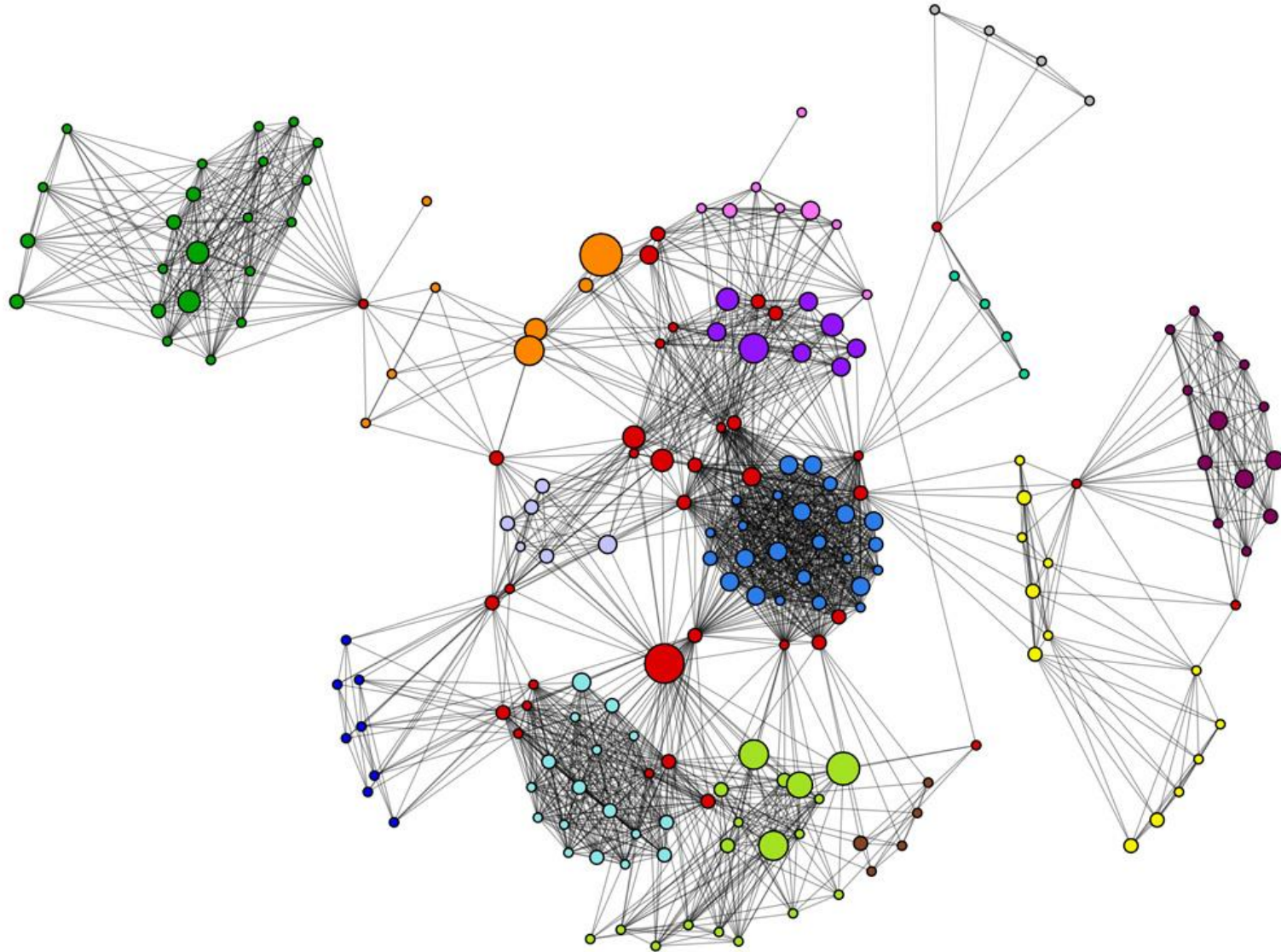
Frequent Subgraphs

- What pattern occurs frequently in this graph?



Frequent Subgraph: A-B-C

Clustering





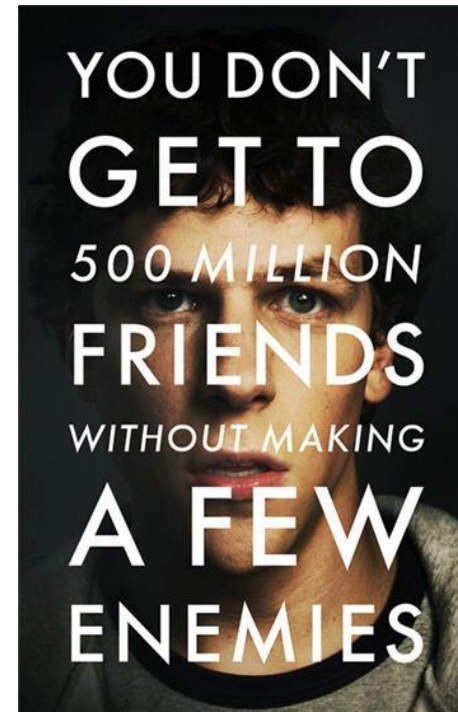


Universiteit Leiden

Online Social Networks

History

- 1997: SixDegrees.com
- 2000: Friendster
- 2003: LinkedIn & MySpace
- 2004: Hyves
- 2005: Facebook
- 2006: Twitter
-
- 2010: The Social Network (movie)

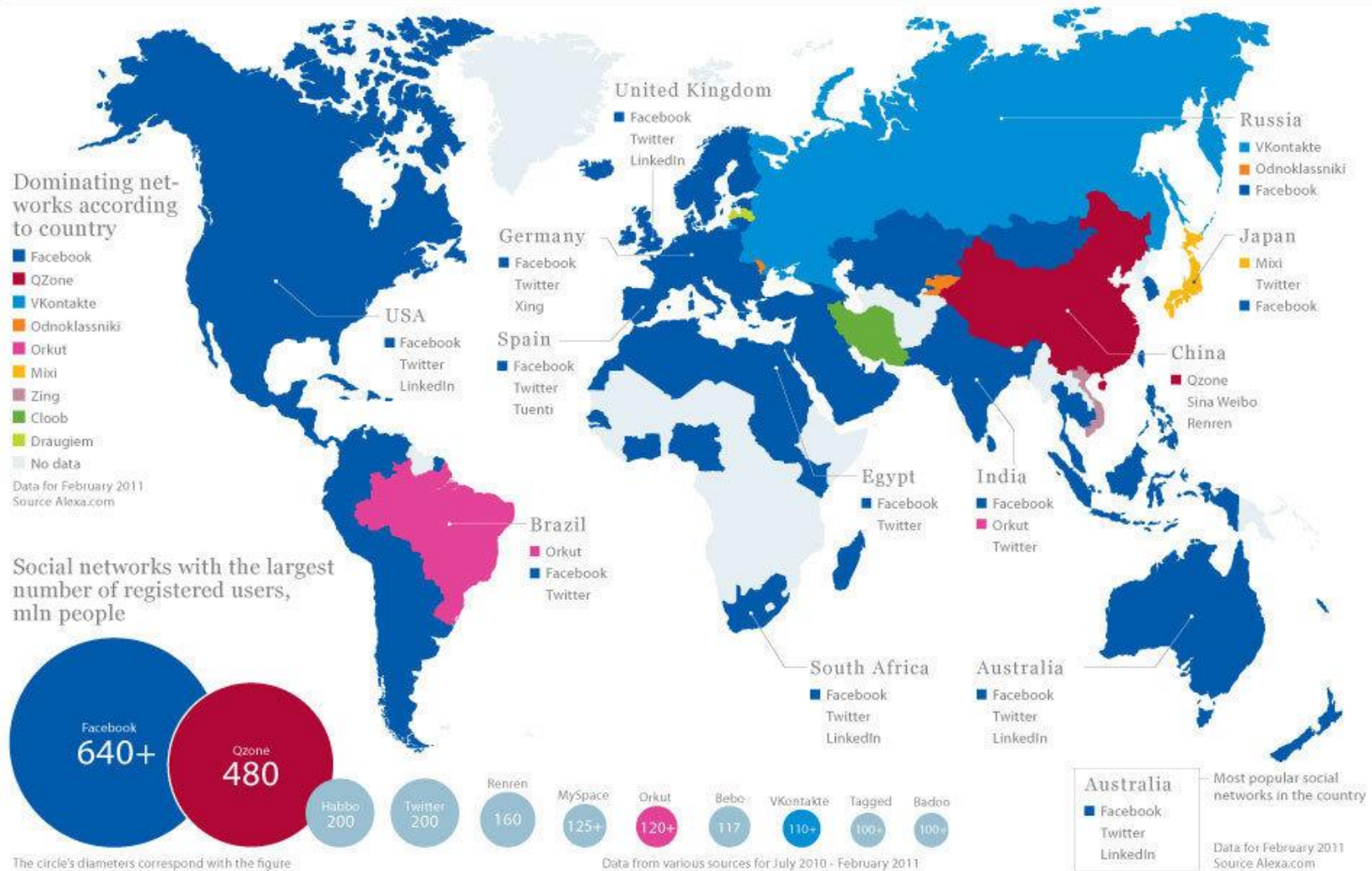


Online Social Networks

- **User** (node) has a profile
- Profiles have **attributes** (labels/annotations)
- **Explicit links** (edges)
 - Social Links / Friendship links
 - User groups
- **Implicit links**
 - Social messaging
 - Common attributes
- Directed vs. undirected links

2011

The world map of social networks





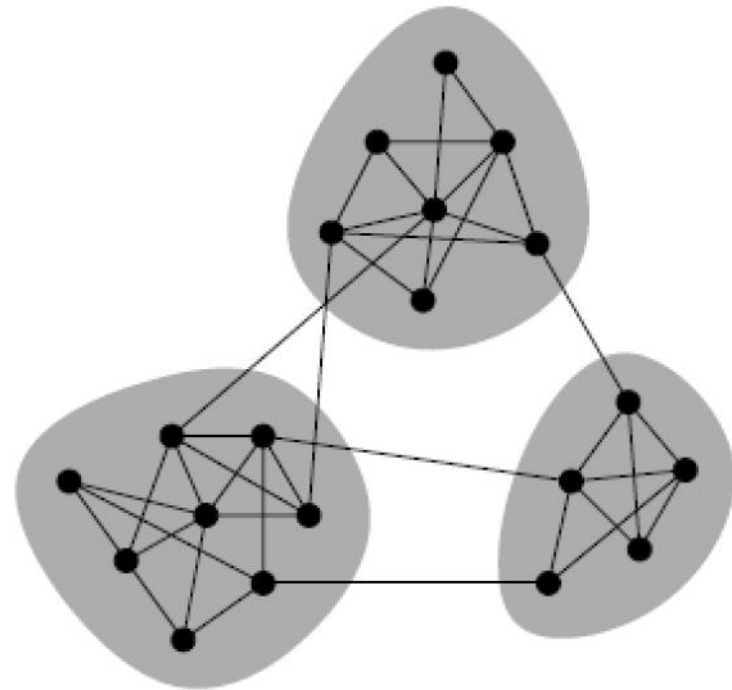
Example: Facebook

- More than **1 billion** active users
- Average user has 130 friends
- Estimated **100 billion** social links
- Over 600 million interactive objects (pages, groups and events)
- More than 45 billion pieces of content (web links, news stories, blog posts, notes, photo albums, etc.) shared each month



OSNA Research Topics

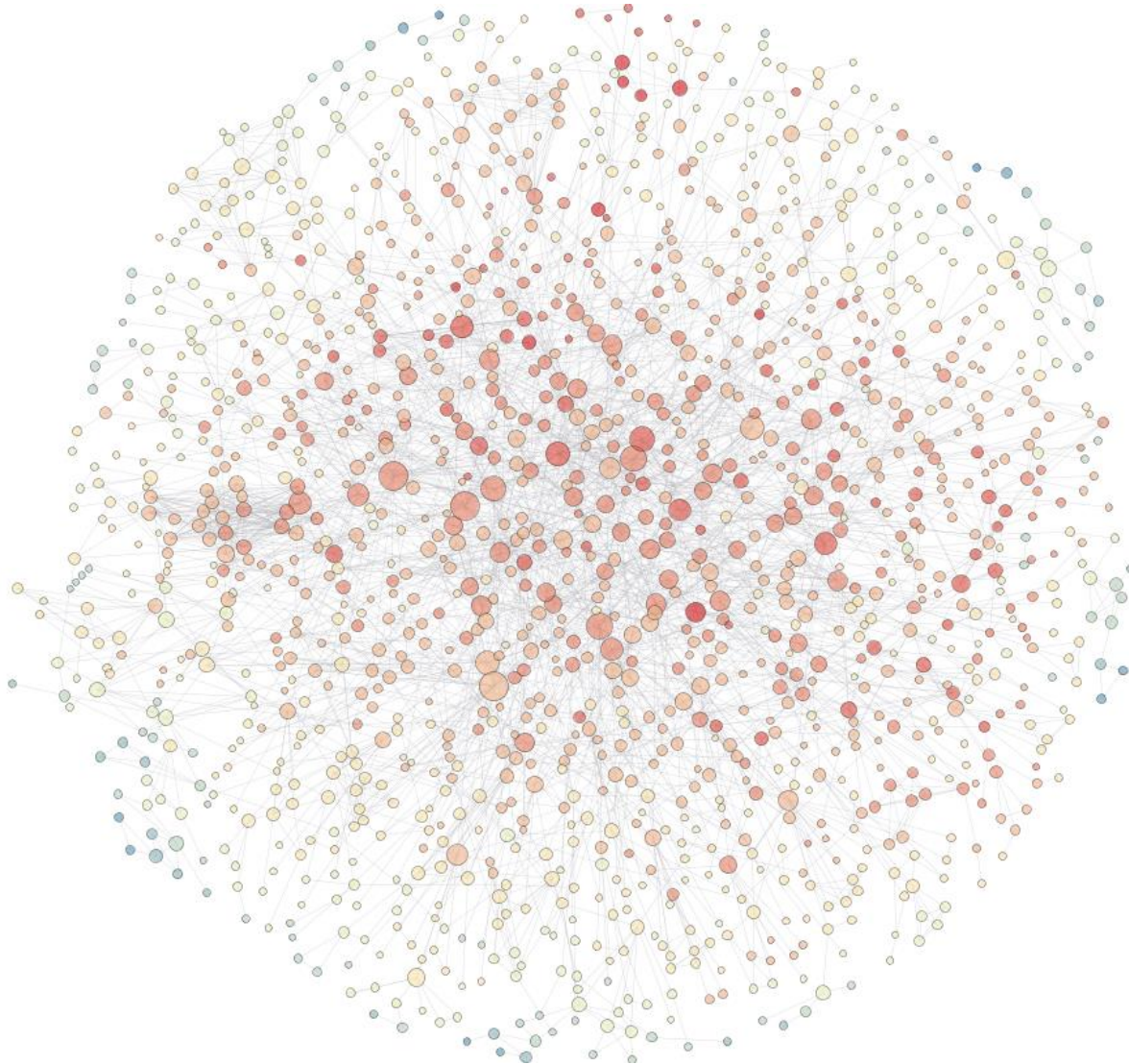
- User behavior
- Privacy & Anonymity
- Trust & Authorities
- Diffusion of information
- Sampling & Crawling
- Community Detection
- **Friendship Graph**





Friendship Graph

Friendship Graph



Friendship Graph Analysis

- **Static analysis**

- Densely connected core
- Fringe of low-degree nodes
- Few isolated communities & singletons

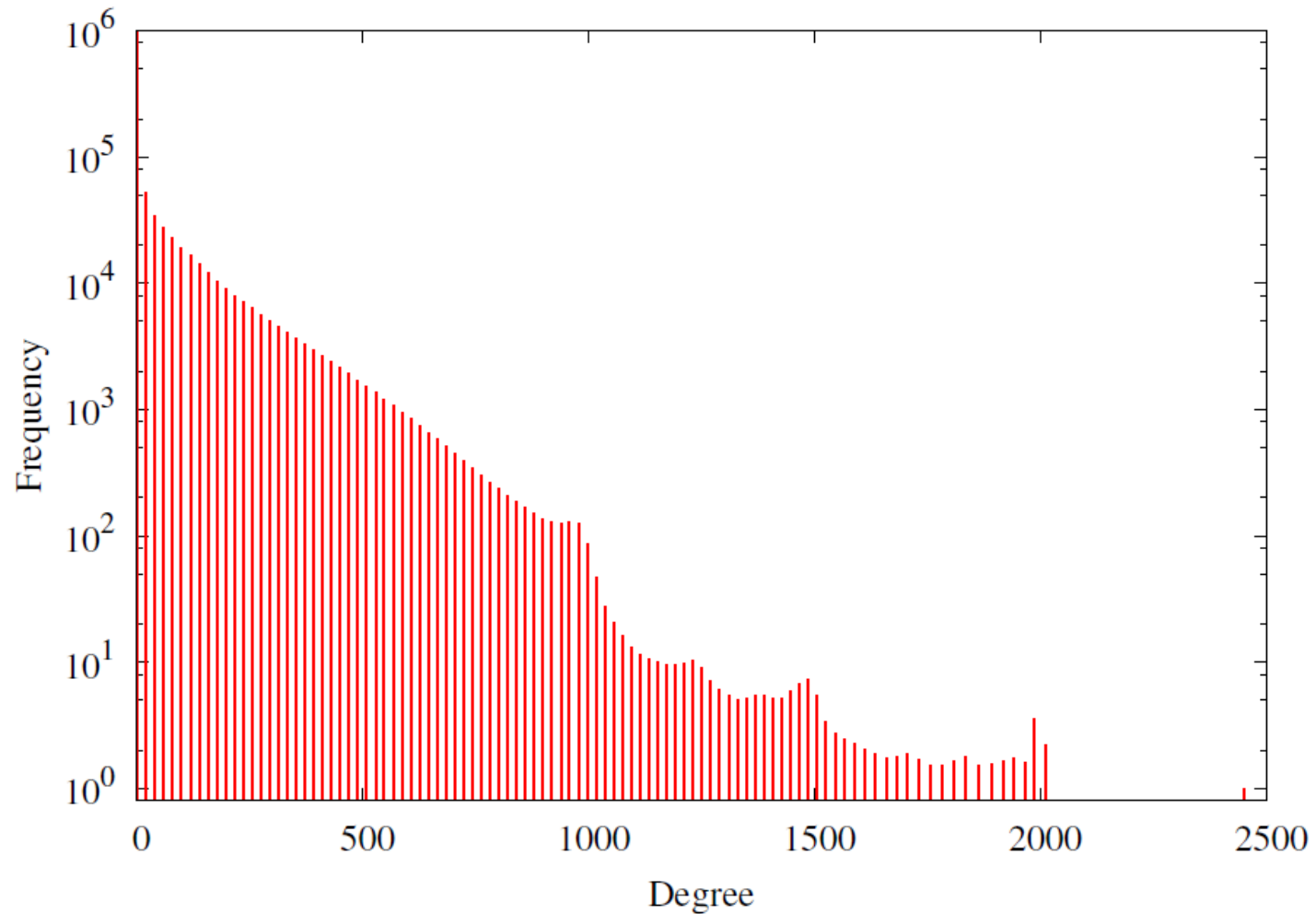
- **Static properties**

- Node **degree distribution**, average distance, diameter
- Edge/node ratio, level of symmetry
- Number of cliques, k -cliques, etc.
- **Small world phenomenon**

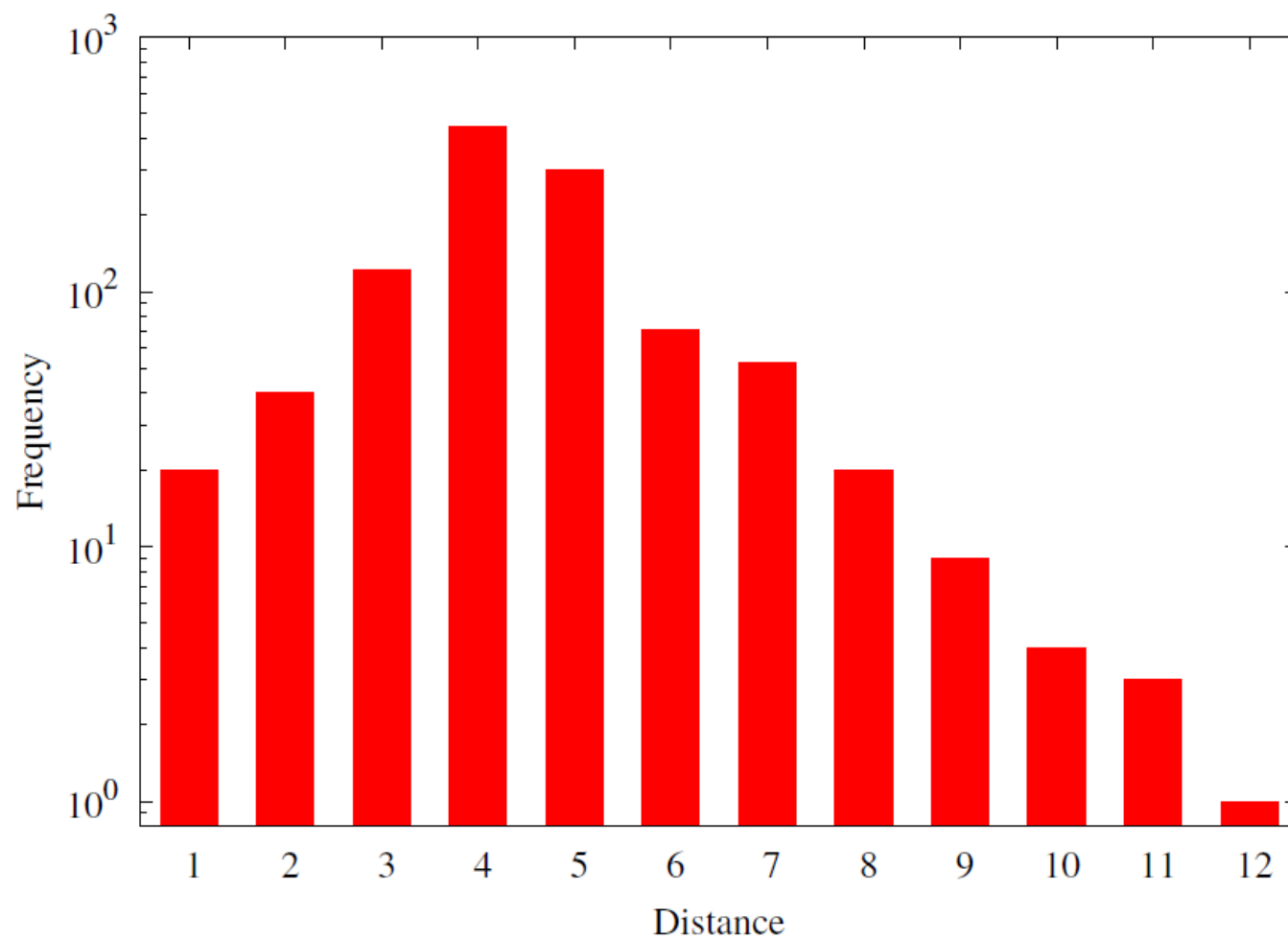
Static properties

Dataset	Nodes	Links	Average Degree	Average Distance	Δ
ASTROPHYS [12]	17,903	396K	21	4.15	14
ENRON [9]	33,696	362K	10	4.07	13
FLICKR [15]	1,624,992	30.9M	18	5.38	24
HYVES	8,057,981	871M	112	4.75	25
LIVEJOURNAL [15]	5,189,809	97.4M	19	5.48	23
ORKUT [15]	3,072,441	234M	76	4.16	10
SKITTER [11]	1,696,415	22.2M	13	5.08	31
YOUTUBE [15]	1,134,890	5.98M	5.3	5.32	24
WEB [13]	855,802	8.64M	10	6.30	24
WIKIPEDIA [5]	2,213,236	23.5M	11	4.81	18

Degree Distribution



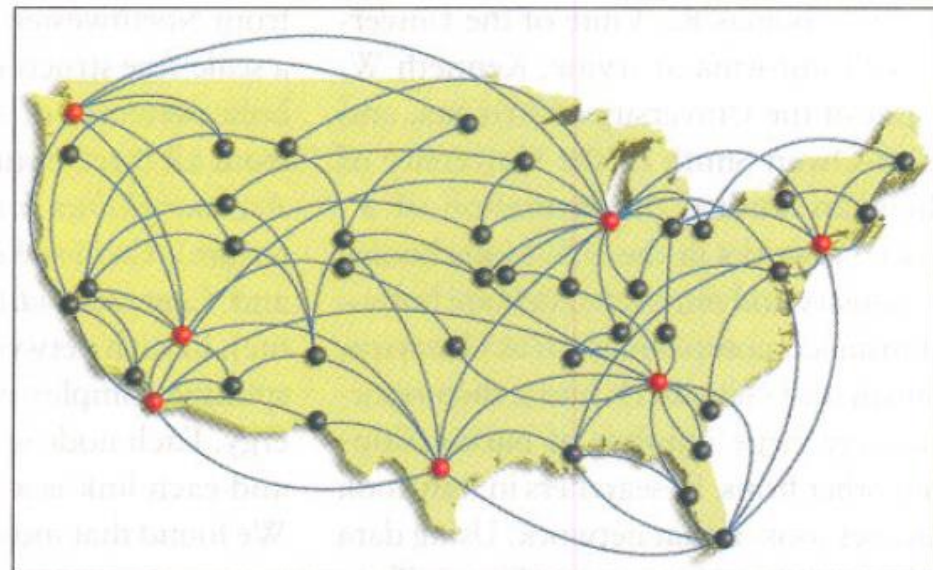
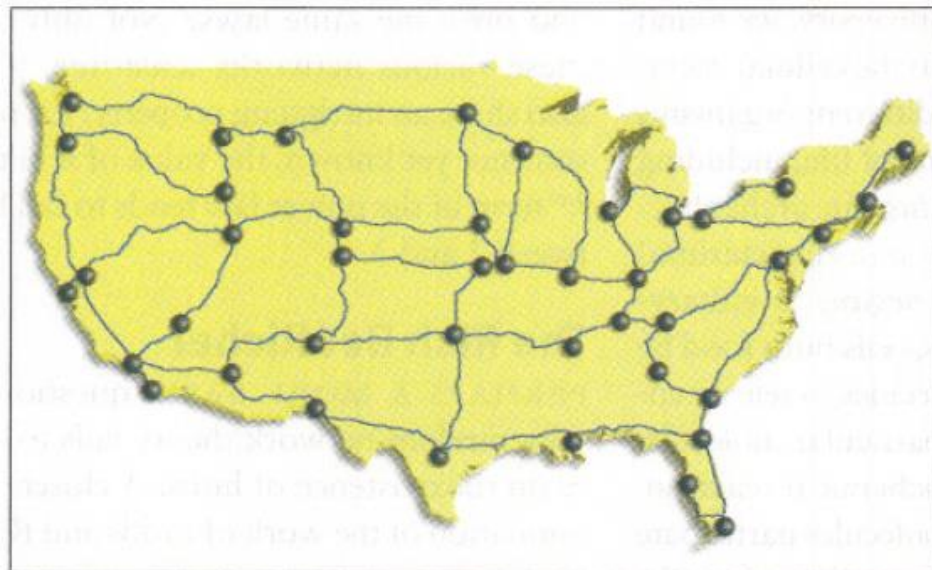
Distance Distribution



Small World Networks

- Class of networks with certain properties:
 - Sparse graphs
 - Highly connected
 - Short average node-to-node distance: $d \sim \log(n)$
 - Fat tailed power law node degree distribution
 - Densely connected core with many (near-)cliques
 - Existence of hubs: nodes with a very high degree
 - Fringe of low(er)-degree nodes

Regular vs. Small World



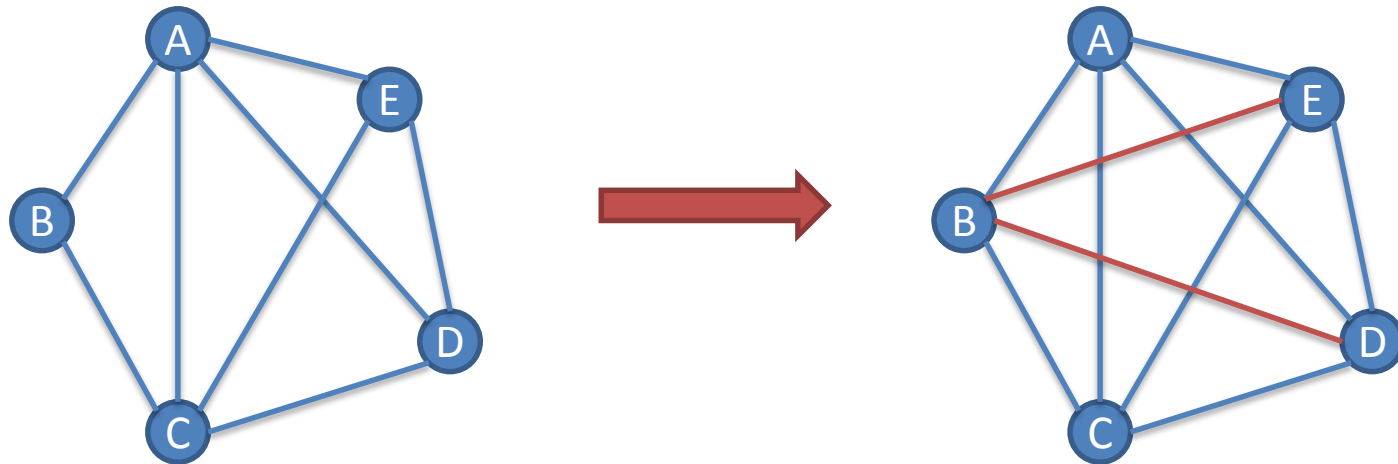
Small World Networks

- Other examples of small world networks
 - Web graphs
 - Gene networks
 - E-mail networks
 - Telephone call graphs
 - Information networks
 - Internet topology networks
 - Scientific co-authorship networks
 - Corporate networks (interlocks or ownerships)

Friendship Graph Analysis

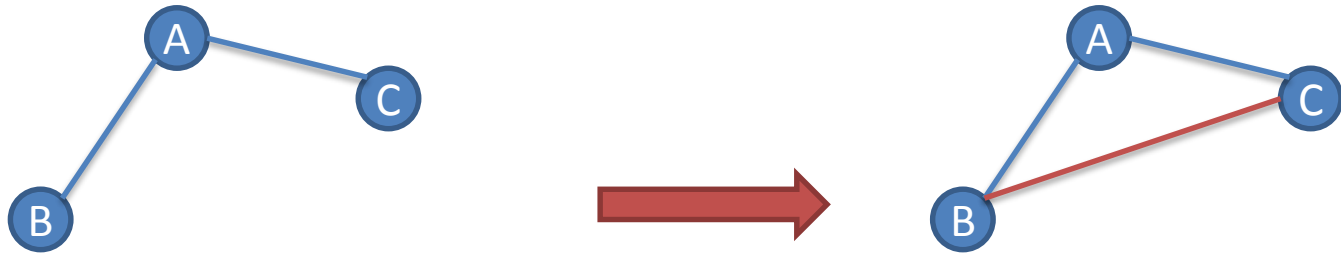
- Static analysis
- **Dynamic analysis**
 - Network evolution
 - Network modelling
 - Network growth
 - **Link Prediction**
 - Triadic Closure
 - Preferential Attachment
 - **Semantic Link Prediction**

Link Prediction



- Two principles: preferential attachment and triadic closure

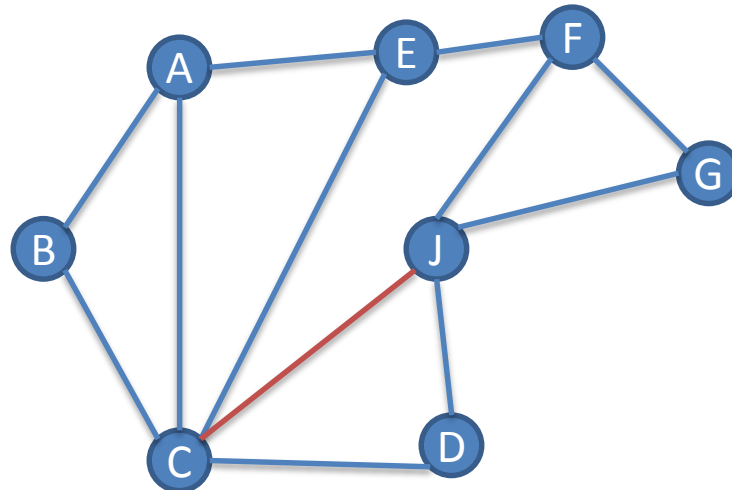
Triadic Closure



- Two principles: preferential attachment and triadic closure

Preferential Attachment

- Nodes with a large degree acquire new links at a faster rate.



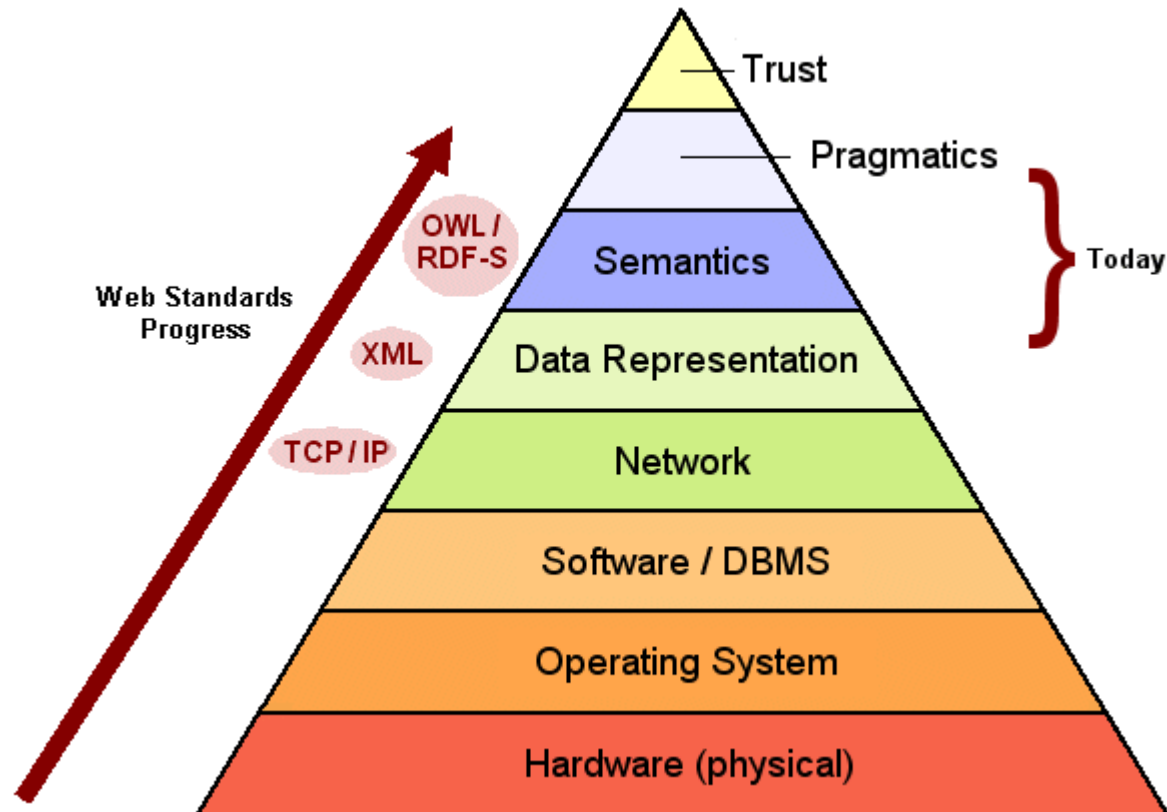


Universiteit Leiden

Semantics

"Call for Semantics"

- Structured data volumes grow rapidly
- "Semantic Web"

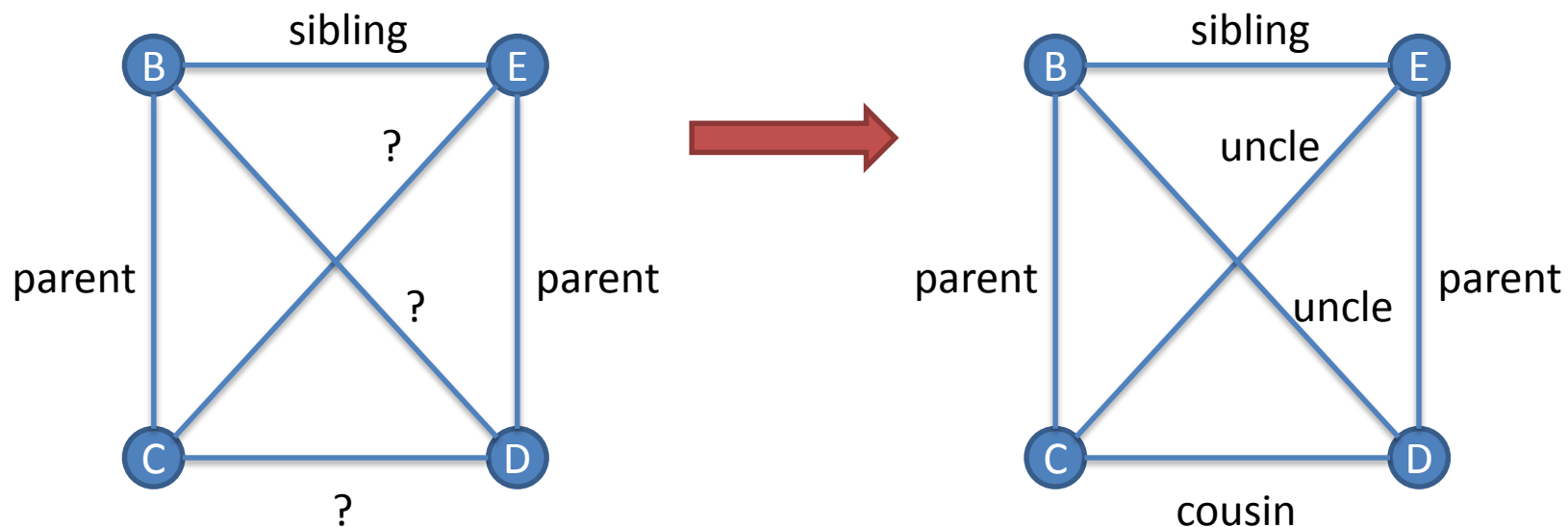


"Call for Semantics"

- What is the meaning of a link?
- What type of relation is defined by a link?
- Are there any **wrong** links?
- What is the **strength** of a link?
- Are link descriptions missing?

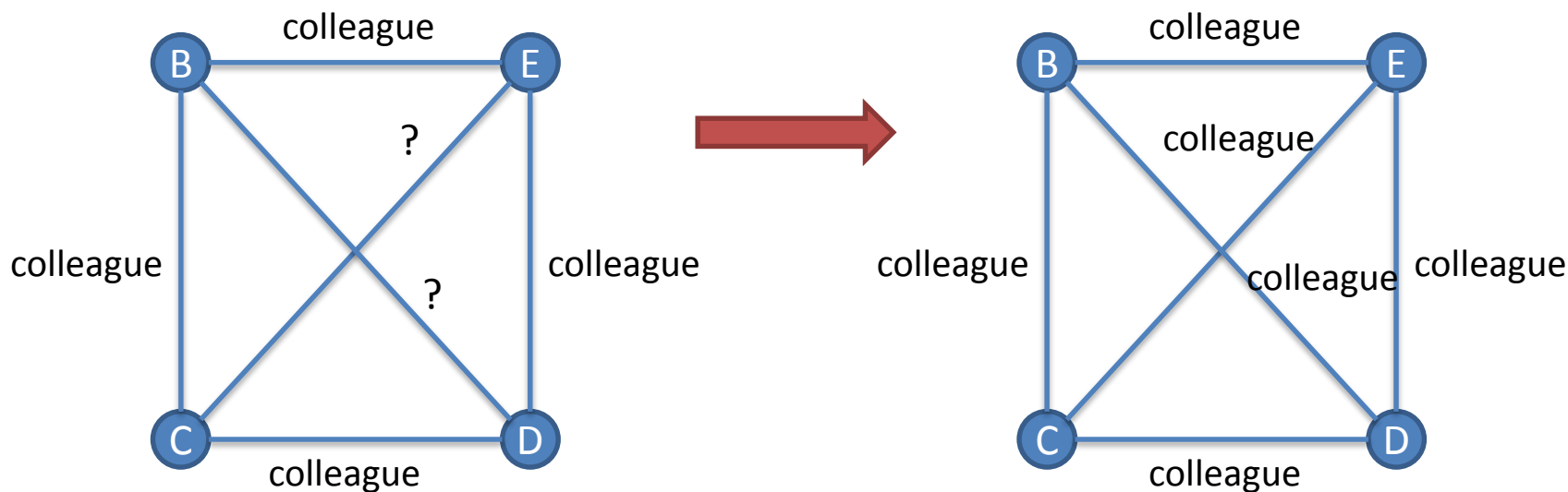
Semantic Link Prediction

- Deterministic



Semantic Link Prediction

- Predictive (learning, classification, etc.)



Conclusions

- Online social networks are an **excellent domain of study** for data (or graph-) miners
- Social Network Analysis is important for **many areas of research**, not only computer science
- **Semantics** within large networks are becoming increasingly more important
- Challenges may be found in the **temporal** (dynamic) **analysis** of social networks