

UNIT-11

EXTRACTION AND MINING COMMUNITIES IN WEB SOCIAL NETWORKS

* The extraction of web community utilize web community chart A graph of communities.

* The main advantage of web community chart is existence of relevance between communities.

Notation use:

* t_1, t_2, \dots, t_n ; Currently, a month is used as the unit time.

* $w(t_k)$; The web archive time at time t_k .

* $c(t_k)$; The web community chart time t_k .

* $c(t_k), d(t_k), e(t_k), \dots$ Communities in $c(t_k)$.

Types of changes:

* Emerge

* Dissolve

* Growth and Shrink

- * Split
- * Merge.

Evolution Metrics:

- * Growth Rate
- * Stability
- * Disappearance rate
- * Merge rate
- * Split Rate.

Other Metrics:

- + Web Archives and Graphs.
- + Split and Merged Communities.
- + Emerged and Dissolved Communities.
- + Growth Rate.

Detecting Communities in Social Networks:

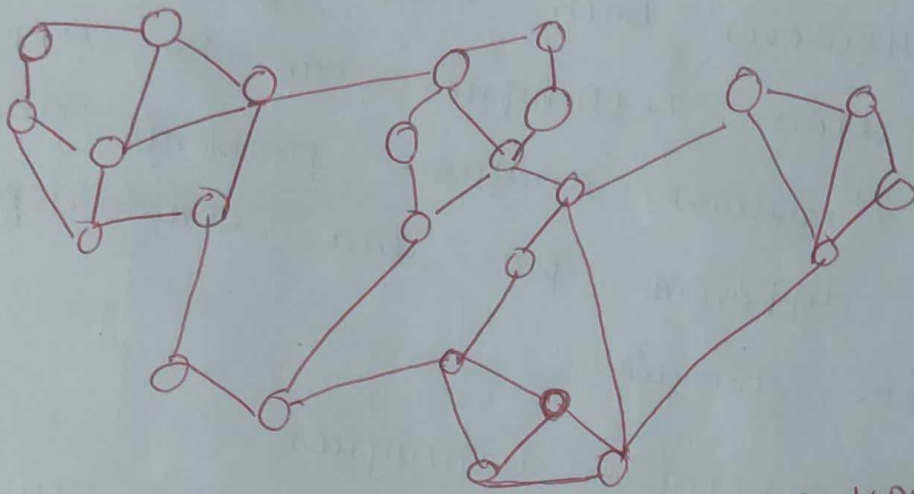
* Communities will help us understand the structure of given social networks.

* Communities are regarded as components of given social networks.

* Communities will play important roles, when we visualize large-scale

Social networks

* Relations of the communities clarify the processes of information sharing and information diffusions.



Many of your familiar with networks?

* social media sites such as Facebook, Instagram, Twitter etc.

* Communities are a property of many networks in which a particular network may have multiple communities. Such that nodes inside a community are densely connected.

Why community detection?

* Community detection can be used in machine learning to detect groups with

Similar properties.

Community Detection Vs Clustering:

* Clustering is a machine learning technique in which similar data points are grouped into the same cluster based on their attributes.

* However both clustering and community detection techniques can be applied to many network analysis problems and may raise different pros and cons depending on the domain.

Community Detection Techniques:-

* Community detection methods can be broadly categorized into two types.

* Agglomerative Methods

* Divisive Methods

Agglomerative Methods:

* Edges are added one by one to a graph which only contains nodes.

* Edges are added from the stronger edge to the weaker edge.

* Divisive methods follow the opposite of agglomerative methods.

* In these, edges are removed one by one from a complete graph.

* There can be any number of communities in a given network and they can be varying sizes.

Types of Communities:

1) URBAN

U - usually a large population so it can be noisy.

R - Residents can take buses or taxis to work.

B - Buildings are close by one another so sometimes people walk

A - Apartments are a popular place to live so the buildings are all.

N - Night life is busy because you can go to movies, theaters, or restaurants.

Definition of Community

* The word Community represents subnetwork whose edges connecting inside of it are denser than the edge.

Types of Communities: [Continuation] (2) 2 mark

- * Rural Community
- * Traditional Community
- * Solidarity Community
- * Urban Community
- * Neighbourhood Community.

— x —

Definition of Community.

* The word "Community" means a subnetwork whose edges connecting inside of it are denser than the edges connecting outside of it.

Community can be classified into the following three categories;

1. Local definitions
2. Global definitions
3. Definition Based on vertex similarity.

Local definitions:-

LAN, Local Area network

* The computers in a LAN connect to each other via TCP/IP ethernet or Wi-Fi.

* LAN is normally exclusive to an organization such as school, office, association or church.

Global definition;

Social networking platform; Facebook,

LinkedIn, Myspace, Orkut.

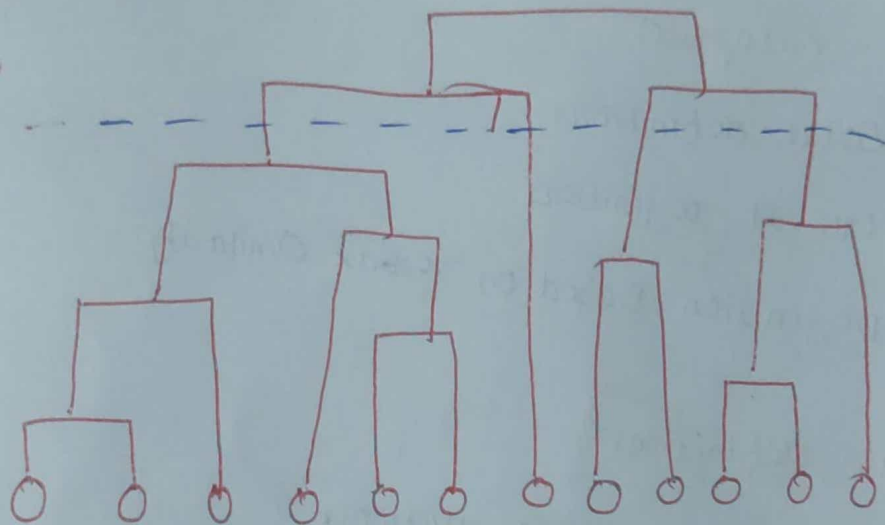
Microblogs; Twitter, Google Buzz.

Photo and video sharing; Flickr, YouTube.

Definitions Based on vertex similarity.

- * Communities are groups of vertices which are similar to each other.
- * Some quantitative criterion is employed to evaluate the similarity between each pair of vertices.

Ex:



- * This structure is called dendrogram, and highly similar vertices are connected in the lower part of the dendrogram.

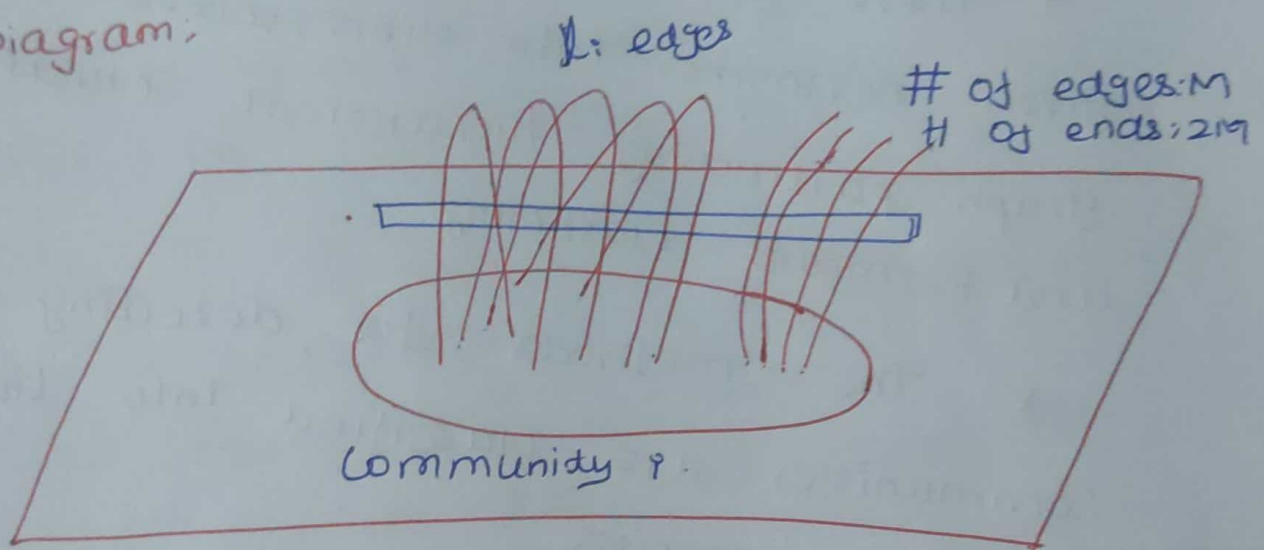
- * Subtrees obtained by cutting the dendrogram with horizontal line correspond to communities.

Evaluating Communities:-

- * Real Community Structure.
- * Quality function for evaluating how good partition is needed.

- * The most popular quality function is modularity of Newman and Girvan.
- * The fraction of edges of the network inside the community.
- * Second term represents the expected fraction of edges.

Diagram:



$$Q_{ii} = \frac{L_i}{M} : \text{Fraction of edges in community } i.$$

$$Q_i^2 = \left[\frac{d_i}{2M} \right]^2 : \text{fraction of edges when connected randomly.}$$

- * A subnetwork is a community if the number of edges inside is larger than the expected in Modularity's null model.

- * The Modularity of the whole network taken as a single community, is zero
- * Modularity is always smaller than one, and it can be negative as well.

Methods for community detection and mining

- * There are naive methods for dividing given networks into subnetworks, such as graph partitioning, hierarchical clustering and K-means clustering.

- * The methods for detecting communities are classified into the following categories

- 1) Divisive algorithms.
- 2) Modularity optimization.
- 3) Spectral algorithms.
- 4) Other algorithms.

Divisive algorithms:

- * A simple way to identify communities.

- * The steps of the algorithm are as follows.

- 1) Computation of the centrality of all edges.
- 2) Removal of edge with largest centrality
- 3) Recalculation of centralities on the running network
- 4) Iteration of the cycle from step (2)

* Edge betweenness is the number of shortest paths between all vertex pairs that run along the edge.

Modularity optimization

* Modularity is a quality function for evaluating partitions.

* This is the main idea for Modularity optimization.

* It has been proved that modularity optimization is an NP hard problem.

* Famous algorithms for modularity optimization is CNM (Customer network management) algorithm.

* Another examples of the algorithms are greedy algorithm and Simulated Annealing

Spectral algorithms:

* Spectral algorithms are cut given network into pieces so that the number of edges to be cut will be minimized.

* One of the basic algorithms is spectral graph bipartitioning.

* The Laplacian matrix L of a network

* All eigenvalues of L are real and non-negative.

Other algorithms:-

* There are many other algorithms for detecting communities, such as the methods focusing on random walk, and the ones searching for overlapping cliques.

Applications of community mining algorithms:

* Network Reduction.

* Discovering scientific collaboration groups from social networks.

* Mining communities from distributed and dynamic networks.

Tools for detecting communities Social network infrastructure and communities

* Several tools have been developed for detecting communities.

classified into the following categories-

1. Detecting communities from large scale network.
2. Interactively analyzing communities from small networks.

* Graph modelling language (GML) is one of the formats for representing networks.

• The common fast greedy community ps for maximizing modularity greedily, and its results are stored in variable gr.

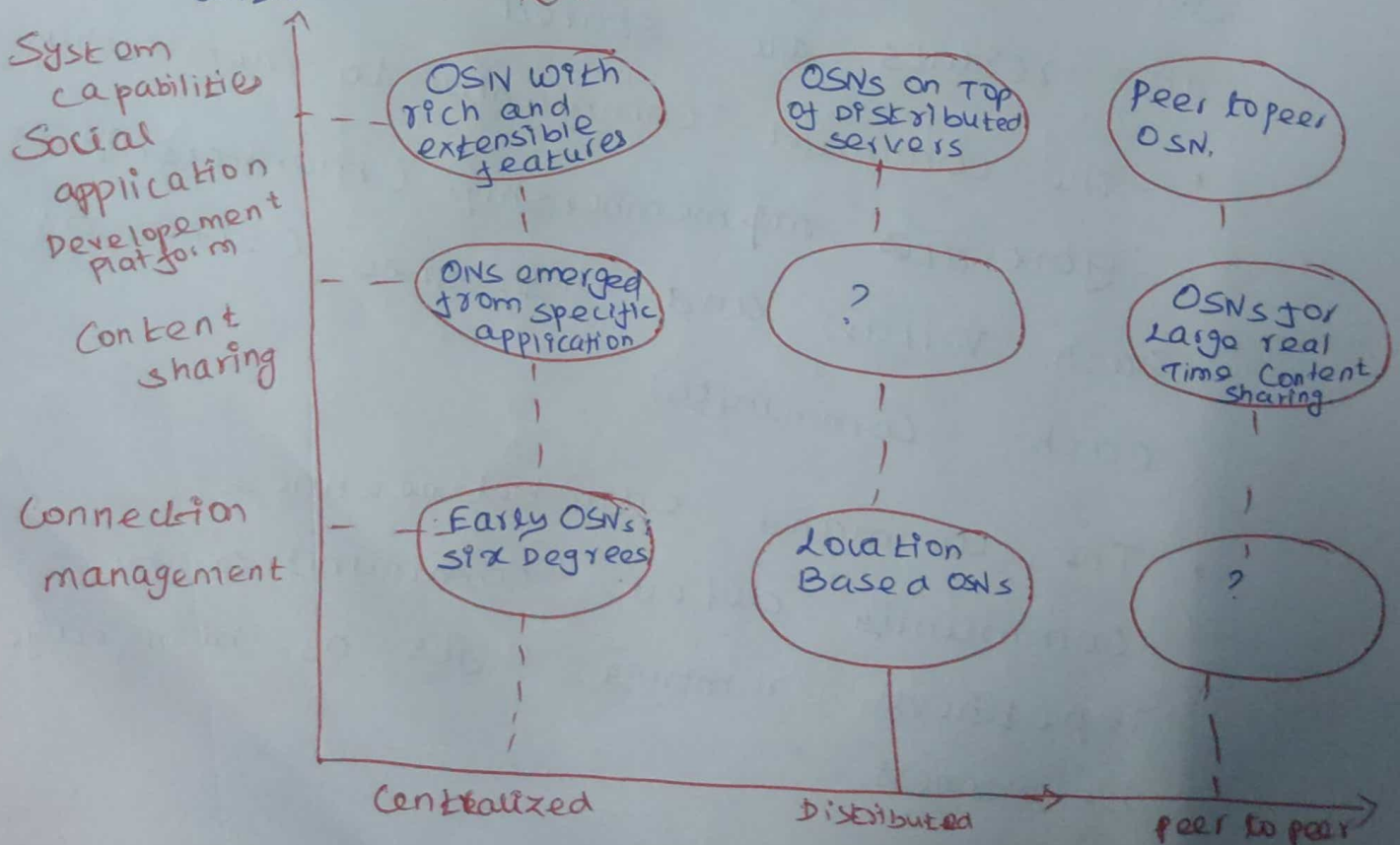
• The command community to membership generates $m \& membership$ (membership for each vertex) and $m \& size$ (size of each community).

• The command edge.betweenness: community detect communities by repeatedly removing edges of high edge betweenness.

Decentralized online social networks

- * Online social network (OSN) is an online platform that
- provides services for a user to build a public profile and to explicitly declare the connection between his or her profile with those of the other users.
 - enables a user to share information and content with the chosen users or public.
 - supports the development and usage of social applications.
 - user can interact and collaborate

With both friends and strangers.



* OSNs with rich and extensive features.

- Facebook

- MySpace.

* OSN emerged from specific applications:

- Flickr

- YouTube

- Twitter.

* Early OSNs: six degrees,

- Friendster

- LinkedIn

* OSNs on top of distributed trusted servers:
young et al. (2009, W3C).

* Location Based OSNs:-

- Dogdeball.

* peer-to-peer OSNs:-

- academic initiatives.

* OSNs for large real-time content sharing.

* Wuia.

* Tribler.

challenges for DASN: [DISK Operating System].

* Storage

* updates

* topology

* search, addressing.

- * Openness to New Applications
- * security
- * Robustness
- * Limited peers
- * Locality

General purpose DOSNs. [Disk Operating System].

Graphical User Interface

Basic applications

Third-party applications

Application programming interfaces (API)

Distributed search

Messaging

publish/subscribe

Application management

shared space

User account

Trust

Access and control security

...

Social network support

Distributed or P2P Storage Systems

Distributed or P2P overlay management

Physical communication network.

* The Reference Architecture consists of six layers.

* Lower layer of this architecture is the Physical Communication network.

* The distributed or P2P overlay management provides core functionalities to manage resources in the supporting infrastructure of the system.

* on the top of this overlay is the Decentralized data management layer,

* which implements functionalities of a distributed or peer-to-peer information system to query, insert, and update various persistent object to the systems.

* the social networking layer implements all basic functionalities.

* social networking layer exposes and implements an application programming interface. To support development of new applications

* The top layer of the architecture includes the user interface to the system and various applications.

* The of & Top of the development platform provided by the DOSN.

Mult-Relational Characterization of dynamic social network communities:

* The characterization of communities in online social media is presented using computational approaches grounded on the observation from social science.

* Motivation: human community as meaning-making eco-system:- The semantic is an emergent artifact of human activity.

* Human activity is mostly social, and social networks.

Motivating applications:

* Real human network communities.

Example Applications include:

* Context-sensitive information search and recommendation.

* Content organization, tracking and monitoring.

Data characteristics and challenges:-

* Large volumes of social media data are being generated from various social media platforms including blogs, Facebook, Twitter, Digg, Flickr

The key characteristics of online social media data include:-

- * Voluminous.
- * Dynamic
- * Context-rich.

Approaches to three problems:-

* Mutual awareness:- It is a bidirectional relationship indicating how well a pair of bloggers is aware of each other, as a

* Facet Net: The community structure at a given timestep is determined both by the observed networked data.

* Metafac: Metafac is the first graph-based tensor factorization framework for analyzing the dynamics of heterogeneous social networks.

