

# Chapter 4

## Community Discovery: From Web Pages to Social Networks

**Damien Leprovost**

*University of Bourgogne, France*

**Lylia Abrouk**

*University of Bourgogne, France*

**David Gross-Amblard**

*University of Bourgogne, France*

### ABSTRACT

*This chapter presents a state of the art of research on the discovery of Web communities, in a general sense. For this purpose, the authors discuss various notions of communities and their related assumptions: hypertextual communities, tag communities and semantic-based communities.*

### COMMUNITY DISCOVERY: FROM WEB PAGES TO SOCIAL NETWORKS

During the past ten years, the Web has turned into an open collaborative system, often called Web 2.0, where anyone can provide information to others using publishing tools such as forums, blogs and wikis. The Web 2.0 also contains social web sites like Myspace, Facebook or Flickr, where people can annotate third-party information with tags. Various kinds of people use these collaborative

systems, ranging from simple visitors to experts of a discussion topic.

From the huge amount of the resulting web pages one can observe emerging structures, forming communities of topics. Similarly, social networks enable to build communities of people, according to their friendship or common interests. A natural challenge for the next decade is to discover and exploit these communities.

This survey presents a state of the art of Web communities emergence, and is organized as follows. The first section discusses the concept of communities and their related assumptions, and presents our analysis method. The subsequent

DOI: 10.4018/978-1-61350-513-7.ch004

sections support the resulting classification: hypertextual communities, tag-based communities, social networks and semantic communities. The last section concludes.

### Communities and Social Networks

What is a community? Indeed, community is an ambiguous term with over 120 definitions noted by Poplin (1979). According to the Wikipedia entry<sup>1</sup> on communities, “*A community is a group of interacting organisms sharing a populated environment. In human communities, intent, belief, resources, preferences, needs, risks, and a number of other conditions may be present and common, affecting the identity of the participants and their degree of cohesiveness*”. In this survey, we consider that a community is a virtual group of interacting entities (a web page, a user, etc.) sharing something or having something in common, each entity being identified in some way (by an URL or a UID). From there, we can distinguish between web pages communities and social networks. More specifically, a social network is a social structure made up of individuals (or organizations) called nodes, which are connected by one or more specific types of interdependency, such as friendship, kinship, common interest, financial exchange, emotional relationships, knowledge, and so on. Hence, a social network is a specific type of community, where relations are explicit, declared (or assumed) by its members.

### First Movement: From Implicit to Explicit Communities

As illustrated by the social network example, there is a historical and natural movement on the Internet from implicit communities to explicit ones. This operation of community discovery is the hard goal of many organizations. We can illustrate this proposition by two scenarios:

- **Web pages community discovery:** due to its openness, the Web allowed a huge number of users for producing content, according to a agreed format, HTML. Hence users or organizations first produced content, regardless of its classification or accessibility. Since group of web pages addressing the same topics were hard to find, a first step toward expliciting communities was registration: that is the effort of declaring a web site to an authority, like DMOZ<sup>2</sup>, Lycos<sup>3</sup>, or Yahoo! Directory<sup>4</sup>. A second step was the use of automatic content analysis to identify topics of interest. Finally, a great success was obtained by combining this content analysis with the examination of explicit web links between pages (Kleinberg, 1998; Brin & Page, 1998).
- **Social community discovery:** as a communication artifact, the Internet naturally hosts social communications. Being part of a thematic mailing list is an example of an explicit community. Similarly, several systems ask users to fulfill detailed profiles to find their matching partners. These approaches, not covered in the present survey, suffer from the hindrance of profile typing, often neglected by users. More recently, registering on a social network, applying for a group and making its friendship relations explicit are other examples. But systems like Facebook<sup>5</sup> are actively analyzing this explicit social graph to discover and suggest new friends to a given user. This way, the social network is expanding from a hidden, implicit state to an explicit one.

The evolution of these notions of communities therefore illustrates a movement from implicit communities to explicit ones. It is not a transformation, but rather the continuous inclusion of unknown or not understood preexisting elements.

7 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

[www.igi-global.com/chapter/community-discovery-web-pages-social/61511](http://www.igi-global.com/chapter/community-discovery-web-pages-social/61511)

## Related Content

---

### Systems Biology-Based Approaches Applied to Vaccine Development

Patricio A. Manque and Ute Woehlbier (2013). *Data Mining: Concepts, Methodologies, Tools, and Applications* (pp. 1131-1148).

[www.irma-international.org/chapter/systems-biology-based-approaches-applied/73488](http://www.irma-international.org/chapter/systems-biology-based-approaches-applied/73488)

### Enabling Efficient Service Distribution using Process Model Transformations

Ramón Alcarria, Diego Martín, Tomás Robles and Álvaro Sánchez-Picot (2016). *International Journal of Data Warehousing and Mining* (pp. 1-19).

[www.irma-international.org/article/enabling-efficient-service-distribution-using-process-model-transformations/143712](http://www.irma-international.org/article/enabling-efficient-service-distribution-using-process-model-transformations/143712)

### Schema Evolution in Multiversion Data Warehouses

Waqas Ahmed, Esteban Zimányi, Alejandro A. Vaisman and Robert Wrembel (2021). *International Journal of Data Warehousing and Mining* (pp. 1-28).

[www.irma-international.org/article/schema-evolution-in-multiversion-data-warehouses/290268](http://www.irma-international.org/article/schema-evolution-in-multiversion-data-warehouses/290268)

### A Promising Direction towards Automatic Construction of Relevance Measures

Lucianne Varn and Kourosh Neshatian (2014). *Biologically-Inspired Techniques for Knowledge Discovery and Data Mining* (pp. 221-235).

[www.irma-international.org/chapter/a-promising-direction-towards-automatic-construction-of-relevance-measures/110461](http://www.irma-international.org/chapter/a-promising-direction-towards-automatic-construction-of-relevance-measures/110461)

### On the Use of Bayesian Network in Crime Suspect Modelling and Legal Decision Support

O. E. Isafiade, A. B. Bagula and S. Berman (2016). *Data Mining Trends and Applications in Criminal Science and Investigations* (pp. 143-168).

[www.irma-international.org/chapter/on-the-use-of-bayesian-network-in-crime-suspect-modelling-and-legal-decision-support/157458](http://www.irma-international.org/chapter/on-the-use-of-bayesian-network-in-crime-suspect-modelling-and-legal-decision-support/157458)