## ABSTRACT

This book offers a comprehensive overview of the various concepts and research issues about blogs or weblogs. It introduces techniques and approaches, tools and applications, and evaluation methodologies with examples and case studies. Blogs allow people to express their thoughts, voice their opinions, and share their experiences and ideas. Blogs also facilitate interactions among individuals creating a network with unique characteristics. Through the interactions individuals experience a sense of community. We elaborate on approaches that extract communities and cluster blogs based on information of the bloggers. Open standards and low barrier to publication in Blogosphere have transformed information consumers to producers, generating an overwhelming amount of ever-increasing knowledge about the members, their environment and symbiosis. We elaborate on approaches that sift through humongous blog data sources to identify influential and trustworthy bloggers leveraging content and network information. Spam blogs or *splogs* is an increasing concern in Blogosphere, which is discussed in detail with the approaches leveraging supervised machine learning algorithms and interaction patterns. We elaborate on data collection procedures, provide resources for blog data repositories, mention various visualization and analysis tools in Blogosphere, and explain conventional and novel evaluation methodologies, to help perform research in the Blogosphere.

The book is supported by additional material, including lecture slides as well as the complete set of figures used in the book, and the reader is encouraged to visit the book website for the latest information:

http://tinyurl.com/mcp-agarwal

## KEYWORDS

blogosphere, weblogs, blogs, blog model, power law distribution, scale free networks, degree distribution, clustering coefficient, centrality measures, clustering, community discovery, influence, diffusion, trust, propagation, spam blogs, splogs, data collection, blog crawling, performance evaluation

*To my parents, Sushma and Umesh Chand Agarwal…–NA*
*To my parents, wife, and sons…–HL*

*…with much love and gratitude for everything.*

# Contents