# Chapter 7.10
# Anomaly Detection for Inferring Social Structure

**Lisa Friedland**
*University of Massachusetts Amherst, USA*

## INTRODUCTION

In traditional data analysis, data points lie in a Cartesian space, and an analyst asks certain questions: (1) What distribution can I fit to the data? (2) Which points are outliers? (3) Are there distinct clusters or substructure? Today, data mining treats richer and richer types of data. Social networks encode information about people and their communities; relational data sets incorporate multiple types of entities and links; and temporal information describes the dynamics of these systems. With such semantically complex data sets, a greater variety of patterns can be described and views constructed of the data.

This article describes a specific social structure that may be present in such data sources and presents a framework for detecting it. The goal is to identify *tribes*, or small groups of individuals that intentionally coordinate their behavior—individuals with enough in common that they are unlikely to be acting independently.

While this task can only be conceived of in a domain of interacting entities, the solution techniques return to the traditional data analysis questions. In

order to find hidden structure (3), we use an anomaly detection approach: develop a model to describe the data (1), then identify outliers (2).

## BACKGROUND

This article refers throughout to the case study by Friedland and Jensen (2007) that introduced the tribes task. The National Association of Securities Dealers (NASD) regulates the securities industry in the United States. (Since the time of the study, NASD has been renamed the Financial Industry Regulatory Authority.) NASD monitors over 5000 securities firms, overseeing their approximately 170,000 branch offices and 600,000 employees that sell securities to the public. One of NASD's primary activities is to predict and prevent fraud among these employees, called registered representatives, or *reps*. Equipped with data about the reps' past employments, education, and "disclosable events," it must focus its investigatory resources on those reps most likely to engage in risky behavior. Publications by Neville et al. (2005) and Fast et al. (2007) describe the broader fraud detection problem within this data set.

NASD investigators suspect that fraud risk depends on the social structure among reps and their employers. In particular, some cases of fraud appear to be committed by what we have termed *tribes*—groups of reps that move from job to job together over time. They hypothesized such coordinated movement among jobs could be predictive of future risk. To test this theory, we developed an algorithm to detect tribe behavior. The algorithm takes as input the employment dates of each rep at each branch office, and outputs small groups of reps who have been co-workers to a striking, or anomalous, extent.

This task draws upon several themes from data mining and machine learning:

**Inferring latent structure in data.** The data we observe may be a poor view of a system's underlying processes. It is often useful to reason about objects or categories we believe exist in real life, but that are not explicitly represented in the data. The hidden structures can be inferred (to the best of our ability) as a means to further analyses, or as an end in themselves. To do this, typically one assumes an underlying model of the full system. Then, a method such as the expectation-maximization algorithm recovers the best match between the observed data and the hypothesized unobserved structures. This type of approach is ubiquitous, appearing for instance in mixture models and clustering (MacKay, 2003), and applied to document and topic models (Hofmann, 1999; Steyvers, et al. 2004).

In relational domains, the latent structure most commonly searched for is clusters. Clusters (in graphs) can be described as groups of nodes densely connected by edges. Relational clustering algorithms hypothesize the existence of this underlying structure, then partition the data so as best to reflect the such groups (Newman, 2004; Kubica et al., 2002; Neville & Jensen, 2005). Such methods have analyzed community structures within, for instance, a dolphin social network (Lusseau & Newman, 2004) and within a company using its network of emails (Tyler et al., 2003).

Other variations assume some alternative underlying structure. Gibson et al. (1998) use notions of hubs and authorities to reveal communities on the web, while a recent algorithm by Xu et al. (2007) segments data into three types—clusters, outliers, and hub nodes.

For datasets with links that change over time, a variety of algorithms have been developed to infer structure. Two projects are similar to tribe detection in that they search for specific scenarios of malicious activity, albeit in communication logs: Gerdes et al. (2006) look for evidence of chains of command, while Magdon-Ismail et al. (2003) look for hidden groups sending messages via a public forum.

For the tribes task, the underlying assumption is that most individuals act independently in choosing employments and transferring among jobs, but that certain small groups make their decisions jointly. These tribes consist of members who have worked together unusually much in some way. Identifying these unusual groups is an instance of anomaly detection.

**Anomaly detection.** Anomalies, or outliers, are examples that do not fit a model. In the literature, the term anomaly detection often refers to intrusion detection systems. Commonly, any deviations from normal computer usage patterns, patterns which are perhaps learned from the data as by Teng and Chen (1990), are viewed as signs of potential attacks or security breaches. More generally for anomaly detection, Eskin (2000) presents a mixture model framework in which, given a model (with unknown parameters) describing normal elements, a data set can be partitioned into normal versus anomalous elements. When the goal is fraud detection, anomaly detection approaches are often effective because, unlike supervised learning, they can highlight both rare patterns plus scenarios not seen in training data. Bolton and Hand (2002) review a number of applications and issues in this area.

## Related Content

Building a Virtual Environment for Diabetes Self-Management Education and Support
Constance Johnson, Kevin Feenan, Glenn Setliff, Katherine Pereira, Nancy Hassell, Henry F. Beresford, Shelly Epps, Janet Nicollerat, William Tatum, Mark Feinglosand Allison Vorderstrasse (2013). *International Journal of Virtual Communities and Social Networking (pp. 68-80).*
www.irma-international.org/article/building-a-virtual-environment-for-diabetes-self-management-education-and-support/111359

Retrieval of Personal Public Data on Social Networks: The Risks for Privacy
Francesca Carmagnola, Francesco Osborneand Ilaria Torre (2013). *Social Network Engineering for Secure Web Data and Services (pp. 137-160).*
www.irma-international.org/chapter/retrieval-personal-public-data-social/75891

Careful What You Say: Media Control in Putin's Russia – Implications for Online Content
Katherine Ognyanova (2010). *International Journal of E-Politics (pp. 1-15).*
www.irma-international.org/article/careful-you-say/43597

Trust Modeling and Management: From Social Trust to Digital Trust
Zheng Yanand Silke Holtmanns (2013). *Examining the Concepts, Issues, and Implications of Internet Trolling (pp. 279-303).*
www.irma-international.org/chapter/trust-modeling-management/74120

The Pervasive and the Digital: Immersive Worlds in Blast Theory's 'A Machine to See With' and Dennis Del Favero's 'Scenario'
Daniel Paul O'Brien (2017). *International Journal of E-Politics (pp. 30-41).*
www.irma-international.org/article/the-pervasive-and-the-digital/186962