

Challenging Social Media Analytics: Web Science Perspectives

Ramine Tinati, Olivier Philippe, Catherine Pope, Leslie Carr, Susan Halford
Web Science Institute, University of Southampton
Southampton, Hampshire
United Kingdom

{R.Tinati, op1e10, C.J.Pope, lac, Susan.Halford }@soton.ac.uk

ABSTRACT

In this paper we outline some of the challenges for social media analytics and – at the same time - challenge existing approaches to social media analysis. Specifically, we suggest that there is an unhelpful gulf between social scientific approaches, which offer rich theoretical and methodological understandings of the social; and computational approaches which offer sophisticated methods for data harvesting, interrogation and modelling. Brought together these approaches might meet the challenges facing social media analytics and produce a different order of understanding. We offer two preliminary examples of this synthesis in practice: first, we show how established computational tools might be harnessed to address theoretically grounded empirical questions about the social; and second we consider social theories might inspire the development of new methodological tools for social media analytics. In doing so, we aim to contribute to the development of interdisciplinary social media analytics with in a broader framework of Web Science.

Categories and Subject Descriptors

H.1.1 Systems and Information Theory

Keywords

Social Theory, Social Media, Twitter, Methodology, Interdisciplinarity

1. INTRODUCTION

The phenomenal growth of Web-based social media data over recent years is currently provoking enormous interest and activity from researchers across a range of disciplines. For most, if not all, the lure of these data is that they offer important insights into *the social*: that is, into the nature of interactions between individuals; the formation of, and distinction between, groups; and the shared meanings and practices – as well as the divisions and inequalities – that characterise our everyday lives. In this respect, social media data offer information at a scale hitherto unimagined in social research [38]. Furthermore, the proportionality of social media offers information (in principle at

least) on ‘whole’ populations, rather than sub-sets; the information is dynamic – captured in real time and over time; and social media provide data on what people say and do ‘in the wild’, rather than what they say they do in response to researchers’ questionnaires and interviews. Furthermore, the digital nature of the data offers unparalleled opportunities for data mining and linking [5,18]. In short, the promise is that social media data will mark a step-change in our understanding of the social world.

However, there are some considerable challenges to be faced before this promise might be realised. Specifically, these relate to the development of data sources and methodologies that will allow us to interrogate and interpret social media data in ways that address complex questions about the social. In part this is a question of data construction (harvesting and archiving). As Geoffrey Bowker now famously observed, ‘... raw data is an oxymoron’ (p.g. 184) [4]. Choices are always made about how to simplify and structure data and these choices bear implications for the kinds of questions that can be asked, and answered. However, in this paper we concentrate on the related question of methodology: that is, the overall design of research from the conceptualisation of questions, to methods and tools, to data analysis and interpretation. Specifically, we will suggest that there is currently a methodological impasse in social media research that must be overcome if we are to realise the contribution that social media data might make to understanding the social. To put it bluntly, whilst the social sciences bring the expertise to construct and interrogate social research questions, underpinned by rich theoretical and methodological traditions, they lack the repertoire of methods necessary to engage with the inherent qualities of social media data. Meanwhile, the computer sciences bring critical expertise for the interrogation of social media data, underpinned by rich computational techniques of large scale data management and modelling, but they lack the theoretical and methodological repertoire necessary to make the most of the methods in addressing complex questions about the social.

This may seem provocative but it is not intended to be so. The historical evolution of academic disciplines has produced divisions of labour that enable the growth of in-depth expertise but – as is increasingly recognised by governments, funding councils and researchers alike – this has siloed knowledge and expertise and, in doing so, limited our understanding of the world. Rather than falling into familiar routines linked to one disciplinary approach over another, our aim here is to evaluate how the combined strengths of the social and computational sciences set an agenda for social media analytics that transcends both the historical divisions and the hierarchical politics of the academy. In what follows, we begin with an outline of the conceptual and methodological framings that drive social science interest in social media, taking Twitter as one example. Next, we consider the methods developed for analysis of Twitter data in

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org

WebSci '14, June 23 - 26 2014, Bloomington, IN, USA
Copyright 2014 ACM 978-1-4503-2622-3/14/06...\$15.00.
<http://dx.doi.org/10.1145/2615569.2615690>

computational research. To illustrate the potential for synthesis across these two streams of Twitter research, we present vignettes of two new studies: one that develops a new form of Social Network Analysis to study political protest; and another that draws on social theory to develop a new computational tool for social media analysis. Finally in our conclusion we suggest some key concerns for the future development of social media analytics.

2. TWITTER AND THE SOCIAL

Social media platforms like Twitter provide a digital trace of human expression, action and interaction of interest to social scientists across the spectrum from psychology to political science, geography and sociology. These digital traces provide us with data at a scale rarely encountered in the social sciences and of a nature that is tremendously time consuming and often difficult to come by: as Latour [28] suggests ‘... it is as if the inner workings of private worlds have been pried open’ (p.2). Broadly speaking, we can identify two distinct types of social science interest in Twitter.

First, social media offer social scientists new data on the subjects that are already well established topics of research. For example, Geographers can use geo-tagging to learn more about the spatialities of social ties [41]; Political Scientists can follow unfolding political protests online [40,43] and the exchange of information between communities of different languages [7], and Sociologists have new data with which to explore identity [19,30,31]. *Second, social media are sometimes seen as part of a paradigmatic shift in the nature of society itself* linked to the emergence of the ‘information age’ [2], ‘network society’ [8,9,10] and the ‘mobilities’ turn [46,47] in social theory. This turns the process of social research on its head. Instead of starting with categories or concepts assumed to define the social and seeking to trace their iteration in the empirical world, the point becomes to trace the *emergence of the social in the dynamic flows of people, objects, images and information* (p.g. 190) [46].

In both cases, social scientists’ experience in researching these areas raises challenges for social media analytics. Contrary to first appearances, it is no simple matter to link social scientific understandings of the social to social media data. Not least, there are sophisticated and competing approaches to *theorising core concepts* - friendship, influence and identity, for instance – that have rich histories in the social sciences and cannot be taken as self-evident. Think of identity for example. This subject evokes enormous contention both between and within disciplines: is identity innate, contextual or discursive? Is it static or dynamic? Are identities coherent or fragmented? The answers to these questions are linked to wider epistemological positions with enormous consequences for the way that research questions are framed, the methods chosen and interpret findings. Linked to this, social scientists have developed an *extensive repertoire of research methods* with which to pursue these complex concepts. Whilst quantitative modelling of large data sets might allow us to answer some questions; others will require in-depth interviews, visual methods of data collection, focus groups or oral histories. We know that different methods will produce different types of data, and different insights. Ticking a box in a questionnaire is not the same as articulating complex emotions in an interview; or recording a visual diary over a longer period of time. Similarly, methods of analysis will – of course – shape the findings. The point is that *theory, methods and interpretation are interwoven* and we must attend to the implications of this for social media research.

However, whilst this substantive and methodological expertise is key to analysing social media data we suggest that, to date, the scope for social scientific research using social media data has been limited by their methodological repertoire. Specifically, that social scientists have approached Big Data with methods that cannot explore many of the particular qualities that make it so appealing to use *viz.* the scale, proportionality, dynamism and relationality described above. Rather, Big Data has commonly been approached with *small scale* content analysis – looking at small numbers of users – or larger scale *random or purposive samples* of tweets. Rendering Twitter data manageable in this way overrides its nature as ‘big’ data, by-passing the scale of the data for its availability or imposing an external structure by sampling users or tweets according to a priori criteria, external to the data themselves. Furthermore, most previous social science studies are snapshots, categorising content and user-types rather than following the data as it emerges dynamically or exploring the nature of online social networks.

3. COMPUTATIONAL APPROACHES TO TWITTER

Meanwhile, the computer sciences’ interest in Twitter begins from quite a different starting point. In particular, there is interest in understanding Twitter at the macro level. These studies aim to explore *the network as a whole*, using computational social network analysis (SNA); well-documented techniques for analysing network graphs which have often been applied to other – similar – large-scale network sources (e.g. the Web graph). These studies use SNA to describe aspects of the Twitter network (friends, followers, retweets, mentions) reveal characteristics including size [42], connectivity [50] and its small-world features [24]. These techniques allow comparison with other Web phenomenon [16] and indeed, other network structures (cancer cells or neural networks [14]). In short, these methods have been driven by questions about the mathematical structure of nodes and edges and by data modelling rather than analyses of the specifically social nature of Twitter.

Beyond this, we see two broad strands of Twitter research in the computational sciences. First, there is growing interest in the *textual data ‘inside’ the Twitter network*. Named Entity Recognition (NER), is used in order to detect and extract vocabulary within tweets. This research is driven largely by the technical challenges involved: for example, how to apply NER to large streaming data sources [29]; or improve the reliability of data extraction from large volumes of unstructured textual data [36]. Way beyond anything achievable through manual processes of analysis [33,37] machine learning techniques can be used to identify events [48], to measure topic frequency and ‘popularity’, and model local and global trends [13]. These findings can be extended with sentiment analysis using Natural Language Processing (NLP) [3,20,32] (usually) to produce a quantifiable value representing positive or negative sentiment e.g. regarding political opinion and mobilization [25,40], health issues [17], and well-being [12]. These techniques have also been applied to more technical challenges such as the detection and filtering of ‘spam’ from large streams of text information [6,21,51]. Meanwhile a second strand of research pays attention to *community and user identification* [15,34] detecting user latent attributes [23,35] and user influence [1,11] and *community formation* around specific topics or events. Here complex algorithms and taxonomic models are used to identify and classify individual behavior behave within

a network, and to model and predict future behaviour patterns [22].

This brief review of computational research on Twitter demonstrates the successful application of computer science techniques in the fields of data mining and NLP to social media analytics. This facilitates increasingly fast and reliable data extraction and interrogation at a scale unachievable with social science methods. Not least, these techniques allow us to engage with the particular qualities that make social media data so appealing to social scientists, particularly those concerned with networks, mobilities and flow viz. the proportionality, temporality and dynamism of social practice 'in the wild'. However whilst this computational expertise is key to developing social media analytics, we suggest that the scope for computational research has been limited by its a-theoretical and largely technical and/or mathematical orientation. Analysing Twitter data in this way overrides its nature as 'social' data, by-passing the theoretical and methodological complexity of the data for its scale. In and of itself this may be unproblematic, depending on the questions being asked. So long as these are technical or mathematical this is entirely appropriate. However if our intention is to explore the 'social' in social media it is more troublesome.

4. WEB SCIENCE TWITTER METHODS

From the brief review above, we are only too aware of the challenges facing interdisciplinary social media analytics. At the same time, there are clearly many ways forward. In what follows we explore two examples, in this case drawing together social theory and computational techniques to achieve a richer and more insightful analysis of Twitter data.

4.1 Using SNA to Trace Information Flows and Emergent Network Roles

Flow 140 (described in detail in [44,45]) is a new network analytics platform built on the well-established techniques and metrics developed in social network analysis (SNA) studies, adjusted to explore the emergence of information flows and network roles over time. Following the sociology of networks, mobilities and flows *Flow 140* is distinguished from conventional SNA in three key ways. First, rather than providing a snapshot of the final network structure, *Flow 140* provides a dynamic mapping of the conversations and flows of information to demonstrate process: that is, how the social emerges over time. Second, and linked to the previous point, unlike traditional approaches to SNA which search for a set of *a priori* characteristics related to the structure and connectivity of a network, *Flow140* attends to the roles that emerge as the network grows over time interactions and activities of the individuals involved [27]. As such it provides a method to follow the digital traces of the social as it evolves [26]. Third, and finally, *Flow 140* transcends the distinction between macro and micro analysis, enabling both large scale data capture of the network as a whole, and associated analysis of network metrics; and in-depth qualitative analysis of the content of individual tweets. We can see not only how information flows, but what information flows; which users are connected in what ways and the roles that emerge in the process of this information flow and network formation.

Using *Flow 140* for a case study of the use of Twitter in political protest [45] revealed which users were key to the generation and flow of information and the different types of roles that were involved. These stretch beyond quantitative measures of re-tweets to include 'amplifiers' and 'aggregators' who – whilst not

necessarily highly retweeted themselves play an important role in the diffusion of information and in building connections between discrete networks. We can also see how quickly particular pieces of information flowed, through which parts of the network and that some limited pieces of information came to dominate the network over time. .

In theoretical terms, *Flow 140* traces the emergence of the social in Twitter activities. Furthermore, by allowing in-depth analysis of the tweet contents *Flow 140* drew attention to the importance of a wider eco-system of interactions with other socio-technical systems such as YouTube, Blogs, and photo sharing sites (and here there are promising connections to computational research making the same point more generally.

Following these links offered a richer understanding of the emerging activities and – critically – how these were connected to activities off-line. In this sense then *Flow 140* extends 'network analysis' beyond the mathematical structure of nodes and edges within Twitter platform – although these are helpful metrics. Instead, this network analysis demands attention to the connections and disconnections online and offline, across diverse fields of action. In this sense, the term 'network' refers not just to a social media network in and of itself but to the wider network in which this might play a part.

4.2 Using Social Theory to Develop New Methods for Twitter Analytics

Our second example takes social theory as its starting point – specifically, theories of social action that emphasise the emergent nature of social outcomes in the flow of everyday action [27,49] - and considers what kinds of methods would be necessary to explore questions of the social from this starting point. From a sociological perspective, the point is not to study individual, discrete actions in and of themselves but rather to understand the contexts and processes that shape these and – in turn – how these actions (re)produce the social world. Considering Twitter, we might ask: why do people tweet, why do people follow particular individuals, or what is the relationship between tweets and followers? However, by asking these questions we must confront methods - both for collecting and for analysing data.

Whilst the dominant paradigm in computational methodologies for social media generates quantitative descriptions of large data for *modelling* and *prediction* [26,39] this does not help us to explore these interactive relationships on Twitter or how they evolve over time. Rather, from this perspective, the tweet is simply a unit of data. It is only *after* collection, during the analysis, that meaning is imputed: the tweets, filtered intentionally or not, (e.g. by technical limitations or sampling techniques), are conceived as raw data, meaningless until clustering or network analysis is applied to make sense of them. This reduces interaction on Twitter to the tweet alone, rather than to the broader range of contexts and relationships in play. In comparison with off-line methods, this is analogous to reducing our understanding of complex social relationships to tick boxes in a survey asking about very particular and actions, rather than asking more in-depth questions about underlying processes and meanings or the wider contexts of action.

Developing an approach where a tweet is meaningful in its context of production rather than during the analysis forces us to rethink the method of collection [27]. This led us to the following principles for our study: (1) Define the population of interest theoretically, rather than solely by reference to technical

- epidemics using search engine query data. *Nature* 457, 7232 (2009), 1012–4.
- [18] Halford, S., Pope, C., and Weal, M. Digital Futures? Sociological Challenges and Opportunities in the Emergent Semantic Web. *Sociology* 47, 1 (2012), 173–189.
- [19] Hargittai, E. and Litt, E. The tweet smell of celebrity success: Explaining variation in Twitter adoption among a diverse group of young adults. *New Media & Society* 13, 2011, 824–842.
- [20] Hu, X., Tang, L., Tang, J., and Liu, H. Exploiting social relations for sentiment analysis in microblogging. *Proceedings of the sixth ACM international conference on Web search and data mining - WSDM '13*, (2013), 537.
- [21] Hurlock, J. and Wilson, M.L. Searching Twitter: Separating the Tweet from the Chaff. *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, AAAI (2011), 161–168.
- [22] Kairam, S., Wang, D., and Leskovec, J. The life and death of online groups: Predicting group growth and longevity. *WSDM'12*, (2012).
- [23] Kim, D., Jo, Y., Moon, I.-C., and Oh, A. Analysis of Twitter Lists as a Potential Source for Discovering Latent Characteristics of Users. *CHI 2010 Workshop on Microblogging: What and How Can We Learn From It?*, (2010).
- [24] Kwak, H., Lee, C., Park, H., and Moon, S. What is Twitter, a social network or a news media? *19th international conference on the World Wide Web*, (2010), 591–600.
- [25] Larsson, a. O. and Moe, H. Studying political microblogging: Twitter users in the 2010 Swedish election campaign. *New Media & Society* 14, 5 (2011), 729–747.
- [26] Latour, B., Jensen, P., Venturini, T., Grauwlin, S., and Boullier, D. The Whole is Always Smaller Than Its Parts - A Digital Test of Gabriel Tarde's Monads. *British Journal of Sociology*, (2011), 1–21.
- [27] Latour, B. *Reassembling the Social: An Introduction to Actor-Network-Theory by Bruno Latour*. Oxford University Press, 2005.
- [28] Latour, B. Beware, your imagination leaves digital traces. *Times Higher Literary Supplement*, April (2007).
- [29] Li, C., Weng, J., He, Q., Yao, Y., and Datta, A. TwiNER: named entity recognition in targeted twitter stream. *SIGIR '12 Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, (2012), 721–730.
- [30] Marwick, A. and boyd, danah. To See and Be Seen: Celebrity Practice on Twitter. *Convergence: The International Journal of Research into New Media Technologies* 17, 2 (2011), 139.
- [31] Murthy, D. Towards a Sociological Understanding of Social Media: Theorizing Twitter. *Sociology* 46, 6 (2012), 1059–1073.
- [32] Park, J. and Lee, W. Revolution 2.0 in Tunisia and Egypt: Reactions and sentiments in the online world. *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, AAAI (2011).
- [33] Paul, M.J. and Dredze, M. You Are What You Tweet: Analyzing Twitter for Public Health. *Artificial Intelligence*, (2011), 265–272.
- [34] Pennacchiotti, M. and Popescu, A. A Machine Learning Approach to Twitter User Classification. *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, (2011), 281–288.
- [35] Rao, D., Yarowsky, D., Shreevats, A., and Gupta, M. Classifying Latent User Attributes in Twitter. *Proceedings of the 2nd international workshop on Search and mining user-generated*, ACM Press (2010).
- [36] Ritter, A., Clark, S., and Etzioni, O. Named entity recognition in tweets: an experimental study. *Conference on Empirical Methods*, (2011), 1524–1534.
- [37] Sadilek, A. and Kautz, H. Modeling spread of disease from social interactions. *Proceedings of the Sixth International AAAI Conference on Weblogs and Social Media*, (2012), 322–329.
- [38] Savage, M. and Burrows, R. The Coming Crisis of Empirical Sociology. *Sociology The Journal Of The British Sociological Association* 41, 5 (2007), 885–899.
- [39] Savage, M. and Burrows, R. Some Further Reflections on the Coming Crisis of Empirical Sociology. *Sociology The Journal Of The British Sociological Association* 43, 4 (2009), 762–772.
- [40] Segerberg, A. and Bennett, W.L. Social Media and the Organization of Collective Action: Using Twitter to Explore the Ecologies of Two Climate Change Protests. *The Communication Review* 14, 3 (2011), 197–215.
- [41] Takhteyev, Y., Gruz, A., and Wellman, B. Geography of Twitter networks. *Social Networks* 34, 1 (2011), 1–25.
- [42] Teutle, A.R.M. Twitter: Network properties analysis. *2010 20th International Conference on Electronics Communications and Computers (CONIELECOMP)*, (2010), 180–186.
- [43] Theocharis, Y. Young people, political participation and online postmaterialism in Greece. *New Media & Society* 13, 2 (2011), 203–223.
- [44] Tinati, R., Carr, L., and Hall, W. Identifying communicator roles in twitter. *Proceedings of the 21st international conference companion on World Wide Web*, (2012), 1161–1168.
- [45] Tinati, R., Halford, S., Carr, L., and Pope, C. Big Data: Methodological Challenges and Approaches for Sociological Analysis. *Sociology*, (2014).
- [46] Urry, J. Mobile sociology. *British Journal of Sociology* 51, 1 (2000), 185–203.
- [47] Urry, J. *Mobilities*. Wiley, 2007.
- [48] Weng, J. and Lee, B. Event Detection in Twitter. *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*, AAAI (2011).
- [49] Wittel, A. Toward a Network Sociality. *Theory, Culture & Society* 18, 2001, 51–76.
- [50] Wu, S., Hofman, J.M., Watts, D.J., and Mason, W.A. Who Says What to Whom on Twitter. *Proceedings of the World Wide Web 2011*, (2011).
- [51] Yardi, S., Romero, D., Schoenebeck, G., and Boyd, D. Detecting spam in a twitter network. *First Monday* 15, 1 (2010).