



# Twitter users change word usage according to conversation-partner social identity



Nadine Tamburrini<sup>a</sup>, Marco Cinnirella<sup>b</sup>, Vincent A.A. Jansen<sup>a</sup>, John Bryden<sup>a,\*</sup>

<sup>a</sup> School of Biological Sciences, Royal Holloway University of London, Egham TW20 0EX, UK

<sup>b</sup> Department of Psychology, Royal Holloway University of London, Egham TW20 0EX, UK

## ARTICLE INFO

### Keywords:

Twitter  
Community structure  
Social identity  
Language accommodation  
Linguistic convergence  
Social network analysis

## ABSTRACT

This paper investigates how people express social identity at a large scale on a social network. We looked at communities of users on the Twitter website, and tested two established social-psychology theories that are usually performed at local scale. We found evidence of *Communication Accommodation Theory*, where community members vary their language characteristics depending on which community they are communicating with. We also found the level of linguistic variation correlated with how isolated a community was: evidence that there is *Convergence* between linked members. This demonstrates the power of methods which analyse subtle human behaviour on social networks.

© 2014 Elsevier B.V. All rights reserved.

## 1. Introduction

Social identity is that proportion of an individual's self-concept that derives from membership of a social group (Tajfel and Turner, 1979). Group affiliation has functions of enhancing cooperation (Boyd and Richerson, 2009) and allowing individuals to define others through the group they belong to, in the same way that the individual defines him or herself through the identity of their own group (Ashforth and Mael, 1989). Group members share behaviour and social norms. This shared behaviour in social groups is thought to be generated through processes on social networks such as convergence of behaviour due to social relationships (Hormuth, 1990; Ethier and Deaux, 1994).

The way we use language is strongly associated with our social identity (Scott, 2007). The convergence of behaviour, proposed by social identity theory, is often studied through the language used within social groups. This demonstrates how language is more than just a means of communication and sociolinguistic studies have shown that varieties of a language can be strongly associated with social or cultural groups (Gumperz, 1958; Labov, 1966; Carroll, 2008; Bryden et al., 2013).

By using language as a proxy for social behaviour, studies have been able to understand how expression of social identity is often strongly context dependent: people will behave differently depending on which social identity has the strongest salience in

the current situation (Hogg and Reid, 2006). Studies show how this often manifests in the accommodation of language according to the social identity of the interlocutor (Giles, 1973; Gallois et al., 2005). Individuals negotiate the social distance between themselves and the person with whom they are conversing, and are therefore in control of its creation and maintenance (Shepard et al., 2001). For example, Iwasaki and Horie (2000) reported how Thai speakers would adjust their linguistic registers when interacting with strangers. These studies look at specific groups or social situations, but we do not know whether this behaviour can be found at a large scale across many groups where these groups are allowed to freely interact with one another.

Online social networking platforms are providing us with a large scale platform to study human behaviour. With over 200 million monthly active users (Costolo, 2013), the Twitter social network is particularly useful due to its publicly accessible nature (Virk, 2011) and network size. The analysis of large networks brings with it considerable statistical power that allows for the detection of patterns that in traditional, smaller scale network studies would be undetectable. Twitter functions as a micro-blogging website, working on the premise of users sharing their opinions and thoughts in brief messages (maximum 140 characters), which are referred to as "tweets". An investigation into the reasons why people post on the Twitter website by Java et al. (2007) found that about one eighth of posts were conversational messages rendering Twitter as a prime resource for public access to naturally occurring communication (Danescu-Niculescu-Mizil et al., 2011) making this public resource an excellent place to study the expression of social identity.

\* Corresponding author. Tel.: +44 (0) 1784 414189.  
E-mail address: [john.bryden@rhul.ac.uk](mailto:john.bryden@rhul.ac.uk) (J. Bryden).

The study of how identity affects our use of language online is a growing field. There is evidence for communication accommodation between offline conversation partners (Danescu-Niculescu-Mizil et al., 2011) showing that syntax, pitch, gestures, word choice, length or form can differ according to interlocutor. Evidence for linguistic convergence online is mixed with studies finding evidence both for (e.g. Riordan et al., 2013) and against (e.g. Christopherson, 2011) the existence of convergence in online communication. The anonymity sometimes engendered in computer mediated environments can act to enhance the significance of social identity in contexts where a relevant shared group membership is salient to users (Postmes et al., 2000). Consequently, social identity can be heightened which explains why some group phenomena, such as polarisation of attitudes, and stereotyping, can seem enhanced in some online environments (e.g. Postmes et al., 2001). This is evident due to collective identity amongst communities of websites of environmental activists (Ackland and O'Neil, 2011). However, such studies of social identity in computer-mediated-communication are still in their relative infancy and this research aims to contribute to the further development of this field by looking to expressly link communication accommodation and convergence to social groups that have formed on Twitter.

In order to identify online groups, we look to the study of complex networks. In this field, the term communities is used to denote parts of the network that are more strongly linked within themselves than to the rest of the network, a phenomenon that has been observed in many human social networks (Porter et al., 2009). In this sense, communities are an emergent property of network structure. Much work has gone into developing methods to detect such groups from topological analysis (Fortunato, 2010), and the extent to which this is possible has been termed modularity (Newman, 2006). The communities found in this way are usually associated with groups of friends or acquaintances, or similarity in traits (Porter et al., 2009; Bryden et al., 2011; Traud et al., 2012) and have also been shown to share language features (Bryden et al., 2013). We hypothesise that communities found in online networks will share social identity and consequently we expect to find that they demonstrate communication accommodation and convergence.

In this study we focus on a specific aspect of behaviour that is strongly associated with social identity, asking whether individuals will shift their linguistic behaviour according to which social group they are messaging. The data of online communities that we used came from a previous study of the Twitter website (Bryden et al., 2013). We tested for communication accommodation by looking to see if users varied specific language characteristics according to whether they had sent conversational messages to members of the same community or to members from other communities. We tested for convergence by looking to see whether this level of language variation for a community correlated with how strongly linked a community is within itself.

## 2. Methods

The data upon which we did our tests was a network of 189,000 Twitter users. To identify users to download we used a snowball-sample where, for each user sampled, all their tweets which mentioned other users (using the '@' symbol) were recorded and any new users referenced added to a list of users from which the next user to be sampled was picked. Starting from a random user, conversational tweets, time-stamped between January 2007 to November 2009 were sampled from the Twitter website during December 2009, yielding over 200 million messages. The network was formed of bidirectional links, where both nodes had sent at least one message to one another, and weighted by the number of

tweets sent between the two users linked. We ignored messages that were copies of other messages (so-called retweets, which are identified by a case-insensitive search for the text 'RT'). In total the network had 75 million messages (tweets) directed from users of the network to one another.

The network was partitioned into communities using a modularity maximisation algorithm (Blondel et al., 2008) and a partition of the network was found where 91% of the tweets were sent by users to other users within the same community. For each community, characteristic words were generated that were used more commonly in that group than the global average (see Supplementary information for characteristic words). These allowed us to identify English speaking groups and also qualitatively summarise shared characteristics of each group. For more information on how characteristic words were generated, and an argument that the network sampled was representative of the complete Twitter network, see Bryden et al. (2013).

To investigate changes in language characteristics, we divided messages into two collections: internal messages that were sent to other members of the same group, and external messages that were sent to members of different communities. For each group, we made sure that both collections were of the same size by discarding messages at random from the larger collection. The difference in word usage between the samples from the two classes was calculated.

To calculate differences between word usage between the two samples we used text similarity measures. We used two different text measures (Gomaa and Fahmy, 2013) to confirm that the result was not an artefact generated by one of the measures. For a word  $w$  we define numbers of usages of  $w$  in the internal and external samples as  $\lambda_i(w)$  and  $\lambda_e(w)$  respectively. The first measure was the Euclidean distance between relative word usage frequencies for each collection, given by,

$$\left[ \sum_w \left( \frac{\lambda_i(w)}{\sum_v \lambda_i(v)} - \frac{\lambda_e(w)}{\sum_v \lambda_e(v)} \right)^2 \right]^{1/2} \quad (1)$$

The second measure was the quantitative version of the Jaccard distance measure (Gallagher, 1999) which is one minus the multiset intersection of the two samples divided by the multiset union. This is given by,

$$1 - \frac{\sum_w \min[\lambda_i(w), \lambda_e(w)]}{\sum_v \max[\lambda_i(v), \lambda_e(v)]} \quad (2)$$

To look at other linguistic features that can be indicative of changes in linguistic style (see e.g. Bryden et al., 2013; Wagner et al., 2013), we also calculated differences between word-ending frequencies (using both Euclidean and Jaccard distances) and apostrophe frequencies. Differences between apostrophe frequencies were calculated by calculating the frequency of apostrophes per word used by each of the two collections and then calculating the absolute difference between these two values.

## 3. Results

The partitioning of the sample network of Twitter users yielded 414 groups, with 42 groups having more than 250 users. A variety of languages were found with different groups using different languages. To eliminate the effects of a user simply changing between different languages depending on which group they were speaking to, we did the study on the 24 groups (of a size greater than 250 users) that used the English language which were selected in a previous study ((Bryden et al., 2013), and see methods).

With these English-speaking groups, we formed collections of internal and external messages for each group, and then measured

Download English Version:

<https://daneshyari.com/en/article/1129157>

Download Persian Version:

<https://daneshyari.com/article/1129157>

[Daneshyari.com](https://daneshyari.com)