

Birds of a feather talk together: user influence on language adoption

Daniel Kershaw Matthew Rowe Tassos Noulas Patrick Stacey
 Lancaster University Lancaster University, Experian Lancaster University Loughborough University
 d.kershaw1@lancaster.ac.uk matthew@matthew-rowe.com a.noulas@lancaster.ac.uk p.stacey@lboro.ac.uk

Abstract

Language is in constant flux be it from changes in meaning to the introduction of new terms. At the user level it changes by users accommodating their language in relation to whom they are in contact with. By mining diffusion's of new terms across social networks we detect the influence between users and communities. This is then used to compute the user activation threshold at which they adopt new terms dependent on their neighbours. We apply this method to four different networks from two popular on-line social networks (Reddit and Twitter). This research highlights novel results: by testing the network through random shuffles we show that the time at which a user adopts a term is dependent on the local structure, however, a large part of the influence comes from the global structure and that influence between users and communities is not significantly dependent on network structures.

1. Introduction

'I'll brb' and 'how did you vote in the Brexit' are all examples of words and phrases that are recent linguistic innovations. According to linguist David Crystal: "Although many texters enjoy breaking linguistic rules, they also know they need to be understood." [5] There is therefore a tension between the need to be expressive by using out-of-vocabulary (from here on OOV) words/phrases, and being understood. This is not solely to do with expressiveness; however, but equally relates to the (re)production of community identity. There are therefore community structures at work that enable and constrain OOV adoption; ultimately a user and community either adopt or reject a new word. Such influence is characterised by the social network a user belongs to. And yet it is by no means deterministic that in view of such social structures a user will adopt an OOV; at the level of the individual user, some people are more willing to change and be influenced than others. This can be modelled as a user specific threshold

that when breached indicated they would adopt an innovation [10], [17]. Though the challenge is how does one learn the influence exerted on a user, and the threshold at which a user adopts a new term.

The purpose of this paper is thus; to show that influence between users can be determined from mining OOV cascades, show how different users influence user language adoption to varying extents, and show how the influence between users is dependent on global and local structure of the networks. This is ultimately summarised in the following meta question: *Given the creation of out-of-vocabulary (OOV) words to what extent can their adoption be predicted, and to what extent does the structure of a social network impacts influence.*

The contribution of this work is as follows:

- **Modelling influence between user or communities through mining OOV cascades:** we show the ability to learn influence from mining historical language diffusions' from macro (between communities) and micro (between users) interactions.
- **Learn individual user adoption thresholds:** we show that individual language adoption thresholds can be learnt through the use of ROC curves.
- **Innovation adoption across word forms:** we show that users adopt innovations with little variation dependent on the innovation's POS tag.
- **Variation of influence based on network structure:** we show that network structures are highly influential on the time at which an innovation is adopted, though less influential in long run adoption.

The implications of this work are not only confined to the realms of academia with impact in both the business and security communities. In a globalised world language innovations or changes in a broader context pose a number of challenges, be it the alienation of users due to mis-communication, or individuals not understanding regionally-specific words. Further to this, changes can hinder the ability for foreign language learners to adopt a language due to the changing

meaning of words, or impeding collaborative work across cultures [4] due to different business jargon. However, understanding which users and communities have greater influence on a language will allow for foreign language teachers to pre-empt new words entering a language, or for companies to place greater emphasis on communication in the language that has the greatest influence in a given region of a network.

However there is a dark side to the Internet through the prevalence of trolling, hate crime and on-line child predators. On [Safer Internet 2016](#), Nicki Morgon the Education Secretary in the UK, stated that challenges of keeping children safe on-line included understanding the terms that children use in on-line communication “These are all terms that didn’t exist when I was young, and I suspect I’m not alone in needing them explained”.¹ However, even though they launched a website² detailing the words used and their respective meanings the true value would come from pre-empting language change and generating a dynamic list of current and future terms, allowing parents to keep one step ahead of their children.

The remainder of this paper is organised as follows: Section 2 highlights what is currently known about the topic. Section 3 introduces the data sets along with the construction of networks being explained in section 3.1 and what we class as an OOV in section 3.2. Measures and methods applied are explained in section 4, with 5 summarising the results. Section 6 outlines the final contribution of this work and future direction.

2. Related Work

Research into language and OSNs (on-line social networks) has attracted studies across a diverse set of academic disciplines, though for this section we will be focusing on user influence, information diffusion, effects of social structures and language change within OSN’s.

Language is intrinsically connected to geographical locations; leading to the ability to predict a user’s location from text [11]; though, over time as users interact and move around the landscape language innovations will ultimately diffuse. [7] showed that through the use of stochastic modelling one can infer the diffusion network of new words across the United States, results showed that language moves between cities with similar demographics and size. Limitation though could be attributed to the small number of innovation tracked

and their potential correlation to a highly mobile college demographic. However, a study concerning how language emerges over time from multiple sources [13] showed that by modelling the time series across multiple granularities of communities, one was able to pre-empt the growth and acceptance of ‘fleck’ and ‘tfw’ at a global level.

Influence is traditionally defined as “getting people to change their attitudes and behaviours” [12]. [10] modelled influence acting on user as being proportional to the number activated users in the *whole* network. A user then activates when the collective influence breaches their individual adoption/activation threshold, at this point they perceive a greater benefit than cost. However, [17] stated influence can only come from the immediate neighbourhood of a user, as for innovation such as fax machines one can only see the benefits of adoption from those they are connected with.

Though not all users are as quick to adopt an innovation, [15] showed that a user’s *innovativeness* could be defined through their time-to-adoption (*tta*) compared to other users. Thus a relatively low *tta* would mean that the user is innovative, which a high value meaning the user is conservative. However, as shown by [17] a user makes their decisions based on their neighbours with there dissension modelled as a threshold. Thus one can contrast their innovativeness at a *system* and *personal* level; a low *tta* but high threshold would mean they are only innovative relative to the system and not their personal network, were as a high *tta* and low threshold would mean the user is innovative to there personal network and the system.

Building on user threshold models proposed by [17] and [10], [9] showed that user *influence* can be modelled as a function of past action propagation’s (tagging the same photo on flickr), with the influence then decaying over time. Even though the results achieved high accuracies in predicting user actions, this was only tested on one limited data-set aiming to predict whether a user would tag an image on flickr. Additionally, the influence of a community on a signal user can be seen in [6] which showed that users adapt their language to that of the community as they join. This effect can also be used to predict when a user is going to leave a community as their language diverges away from the global language model. However, this research was only performed on a small specific ‘beer community’, where there would be a naturally higher convergence as users used more ‘technical’ terms.

Though it is not only predicting who will adopt an action but also how many people will adopt the same action. Through modelling the diffusion of memes, [19] identified that for a meme to spread it is not only the

¹Nicky Morgan: We simply can’t know everything our children are doing online - 09/02/2016

²Online teen speak - Parent Info

initial popularity of the content that is important (as stated in [16]) but also who initially uses a meme, with initial users needing to have a set of diverse topics and interest. Though [20] proposed diffusions are highly dependent on the network structure, by comparing simulated and real-world diffusion they identified the effects of *homophily* and *social influence*. However they did not differentiate between the two effect, which have been shown to auto-correlate. [1] proposed that through matched pair sampling one was able to distinguished between *homophily* and *social influence*, showing that homophily accounted for 50% of the persuasive behavioural contagion.

There has been a vast amount of work that covers language change, information diffusion and user influence; though each has its own limitations, be it from over or under sampling to limited testing on alternate communities. For the first time, in this paper, we show how users and communities influence each others' language; along with the effect of a network's structure on inter user influence. This is achieved by applying a known influence framework across multiple network abstractions.

3. Data

One of the limitations of previous research is the reliance on one social network, so in this work we draw on two distinct networks: [Reddit](#) and [Twitter](#). Even though both social networks are highly popular³⁴ they both can be conceptualised in different ways; Twitter is a personal broadcast network allowing user to express messages and emotion without necessarily getting a response. Alternatively as Reddit is content-focused and structured into self-governing *Subreddits* that have a particular topic drawing user substitutions and comments.

Twitter has been extensively used in academic research due to its relative widespread adoption and the availability of a publicly accessible API which provides up-to a 1% sample from the global fire-hose. For this study we bound the results return from Twitter to those having originated from within the UK, thus limiting the sample to Tweets that only contain GPS coordinates within the UK. Even though studies have shown that only 4% of tweets contain GPS tags⁵, the sample that was collected from September 2014 until June 2015 contained 111 million tweets. The second data source is an 18 months dump (January 2013 to June 2015) of comment posts from Reddit. This data comes from

³PewResearCenter - 6% of Online Adults are reddit Users

⁴PewResearCenter - Mobile messaging and social media 2015

⁵PewResearCenter - Location-Based Services

Table 1: Dataset Description

| | Reddit | Twitter |
|-----------------|---------------|-------------|
| Unique Words | 2,942,555 | 526,342 |
| Posts | 1,054,976,755 | 111,067,539 |
| Innovations | 2,712,629 | 373,217 |
| Days in Dataset | 880 | 283 |

Table 2: Network Description

| Network | Nodes | Edges | Communities |
|------------------|---------|-----------|-------------|
| Twitter Geo | 2,910 | 436,849 | 14 |
| Twitter Mention | 283,755 | 329,440 | 39,767 |
| Reddit Comment | 861,955 | 2,402,202 | 36,885 |
| Reddit Subreddit | 15,457 | 142,285 | 407 |

a larger data release that spans the entire existence of Reddit from conception in 2007 through to mid 2014.⁶

3.1. Networks

One of the issues when studying OSNs is the need to infer network structure; as the explicitly define relationships on Facebook (friends) and Twitter (following) are challenging to collect and guarded by the respective companies. For this work we thus aim to learn the influence between users or communities through two abstractions of networks from each data source; one representing the interaction of users (*micro*) with the second modelling interactions between communities (*macro*); this allows us to contrast different concepts of influence.

Ultimately the networks will take the form a of *directed social graph*. Where the graph (G) is defined as a quad $G = (V, E, T, W)$, containing vertices $v, u \in V$ and edges between the vertices $(v, u) \in E$; denoting an outward connection from v to u . The quad also includes the time (T) when the edge was created, while $w \in W$ denotes the weight of a given edge. Edges are only added over time and never removed, and there are also no self looping edges.

3.1.1. Micro. At a micro level we model the graphs through user interactions. Within Twitter, users interact with each other in a number of ways, however the predominate form is mentioning fellow users in Tweets (through the inclusion of the '@' symbol and a username). We use this to build a user to user graph, where a relationship from user $v \rightarrow u$ is inferred if u mentions user v ; the edge time (t) is when this interaction first happens, with the weight being the total number of

⁶Data-set available on the [Internet Archive](#)

times u mentions v . Both users u and v must also exist within the data-set.

Similarly within Reddit users comment on each others' posts forming a chain of interactions, thus we define a relationship between users if user u comments on a post of user v thus forming an edge $u \rightarrow v$, the time (t) of the edge would be the first time this happened, and the weight (w) would be the number of times user u has commented on a post of v .

3.1.2. Macro. Even though users may not interact with each other does not mean that they are not exposed to each others' information by observing the network. Collectively a group of users may also exert influence over other collections of users; for this reason we cluster together content (posts and tweets) generated within the same communities (subreddits or postcodes), and generate an edge between these nodes by extracting the users traversing across the network *between* the nodes.

Language is connected to geography, with users taking language with them when they travel. Thus by generating a network based on users travelling between postcodes represents the interactions between locations in the UK, and potential the diffusion of language. Each tweet initially is assigned the postcode from which it originated, with the edges representing movement between locations. The weight of an edge ($w_{i,j}$) represents the number of users moving from $i \rightarrow j$ consecutively.

Similarly within Reddit, users interact and move around different subreddits depending on their current interests or in reaction to popular content. A similar method to that detailed above can be applied to extract the interaction between subreddits. The weight between two nodes is the number of users moving consecutively from one subreddit to the next, with the associated edge time being the first time a user first moved between the two.

3.1.3. Graph Filtering. However not all edges are significant, as a user whom has mentioned a user once is not as important as one that has been mentioned 100 times. For this reason we extract the backbone network by filtering edges that are not statistically significant using the backbone extraction algorithm [14]. Additionally, we apply a fast unfolding community detection algorithm [3] to each of the four networks to identified community of nodes.

3.2. Innovations

The premise of this work is to predict the adoption of OOVs, thus we must first classify what is and what

is not an OOV (a language innovation). For this work we will be stating that a OOV is a word that does not appear within the British National Corpus (BNC) [2]. The BNC was chosen to be the baseline for British Language as it is the most comprehensive study of British English Language in recent times; taking its sources not only from books, but also newspapers, written communication and oral discourse transcripts. This research is only inserted in the emergence and diffusion of new innovations, thus only OOVs that appeared after the first month of data collection were used. However, on initial manual inspection a large number of OOV's were only used by one user (predominantly bot accounts) thus OOVs also had to appear over 10 times and be used by more that 10 users.

The function of words in communication vary; for this reason we broke the analysis down into distinct classes of words. This is achieved by initially POS (Part-of-speech) tagging each data-set, though each innovation which is being assessed could have multiple classes. For simplicity the class assigned to each OOV is the class with the highest count. This is implemented through the use of TwitterNLP [8], due to its ability to deal with noisy social media data.

4. Method

People accommodate their language to make it similar to that of the people around them, thus this research aims to predict when people adopt new terms in response to exposure from their neighbours. We propose that as shown in [17] that each user (u) has an individual threshold σ_u of *joint influence* at which they then adopt a OOV. We therefore used the framework proposed by [9] to model influence between users as a function of previous join actions (propagations).

Goyal et al. [9] stated that influence between users $i_{u,v}$ (influence of v on u) can be learnt as a function of their previous joint actions that have propagated between the two users (from v to u). After learning the influence one can use the joint probability across all active neighbours of user u to express the current joint influence (i_u) on u to adopt action a . To predict if a user is going to adopt an action a the joint influence i_u would need to be higher than the individual threshold σ_u , if they adopt with the value less than the threshold this would be classified as a true negative.

This breaks the analysis down into two distinct stages; learning the influence between two users (section 4.1) and then predicting user adoption of terms (section 4.2).

To learn the influence between users ($i_{u,v}$ where $i_{u,v} \in [0,1]$), we first define a number of basic

measures; O_v and O_u are the number of distinct OOVs used by users v and u respectively, alternatively $O_{v|u}$ is the number of distinct OOV's v or u have used (i.e. the union of their vocabulary). The number of propagations of OOV's between users is defined as O_{v2u} , this is the number of OOV's that were first used by v and then by u , thus $t_v(o) < t_u(o)$ (with function $t_v(o)$ and $t_u(o)$ returning the time that the OOV was used by each user). Though propagation cannot occur if an edge between users has not yet been created, meaning that propagation must also fulfil the following $e_t(v, u) < t_v(o) < t_u(o)$ where $e_t(v, u)$ returns the creation time of the directed edge from v to u .

4.1. Learning Influence

We now define four measures that quantify the influence ($i_{v,u}$) of user v on u . Each measure is based on the values above aiming to quantify influence in different ways.

4.1.1. Bernoulli. We first state that influence is proportional to the fraction of OOVs that have propagated from user v to u as a fraction of *all* the innovations that v has used:

$$p_{v,u} = \frac{O_{v2u}}{O_v} \quad (1)$$

Thus, if all the OOVs that v use end up being used by u then the value would be 1.

4.1.2. Jaccard. Alternatively, influence is proportional to the number of innovations that have propagated (O_{v2u}) out of the union of all innovations used across the two users ($O_{v|u}$):

$$p_{v,u} = \frac{O_{v2u}}{O_{v|u}} \quad (2)$$

This means that if all of v 's OOV's propagate but u uses a large amount of other OOV's as well then the value will be lower than that computed in equation 1.

4.1.3. Partial Credits. However, when users adopt a new term it could be said that each of their neighbours (whom have used that OOV before) all have an equal part to play in the user adopting a new term; thus it could be said that they share equal credit:

$$credit_{v,u}(o) = \frac{1}{\sum_{w \in S} I(t_w(o) < t_u(o))} \quad (3)$$

Where S is a list of activated neighbours of v (e.g. users connected to v who have adopted the OOV before) and with the function I acting as an indicator function that

returns 1 if the neighbour w has used the OOV before u .

We then modify equations 1 and 2 to incorporate the partial credit definition (equation 3). Instead of influence being defined as the number of propagations, it is instead defined as the average credit per OOV used by v or:

$$p_{v,u} = \frac{\sum_{o \in O} credit_{v,u}(o)}{O_v} \quad (4)$$

Or the average credit used across the union of all OOV's used across users v and u :

$$p_{v,u} = \frac{\sum_{o \in O} credit_{v,u}(o)}{O_{u|v}} \quad (5)$$

4.2. Computing Joint Influence

The measures in equations 1, 2, 4 and 5 aim to quantify the influence between users, these metrics can be used to predict user adoption of new terms through computing the joint influence exerted on the user by active neighbours.

The joint probability ($i_u(S)$) can be computed by utilising the *monotonic* and *sub-modular* nature of the influence probabilities;

$$i_u(S) = 1 - \prod_{v \in S} (1 - i_{v,u}) \quad (6)$$

Where S is the set of active neighbours of node u .

Though, the influence that a user exerts might not be constant, with the influence decreasing over time after they themselves have adopted the OOV. A reduction in influence between user may be for a number of reasons, from the OOV dispersing off a users Twitter time-line, or that the users not coming in contact again.

Thus we attempt to model the decay of influence between two users as a function of the average time of propagation ($\tau_{v,u}$). The decay takes two forms, a *discrete* form where the influence stays constant for a set amount of time, or *continuous* form where influence decay happens exponentially.

First, we define the average propagation time of a OOV between two users v and u , this is defined as $\tau_{v,u}$:

$$\tau_{v,u} = \frac{\sum_{o \in O_{v,u}} (t_u(o) - t_v(o))}{O_{v2u}} \quad (7)$$

With $O_{v,u}$ being the set of OOV's that have propagated from v to u , and $t_u(o)$ being the time that u adopted o . As before O_{v2u} is the number of action propagating from v to u .

To model the decay of influence in a basic form we only allow a user to have influence over another for the length of $\tau_{v,u}$. This is to say that after a

user u is exposed to an OOV by v the influence will reduce to 0 once $\tau_{v,u}$ has elapsed, thus the influence window is $[t_v, t_v + \tau_{v,u}]$. At the point when a users influence reduces the joint probability (equation 6) can be updated with the following equation:

$$i_u(S, w) = \frac{i_u(S) - i_{w,u}}{1 - i_{w,u}} \quad (8)$$

Where S is the set of active nodes before w becomes inactive, w is the node that has become inactive. This means that the whole probability does not have to be recomputed at search step, rather just updated.

In reality, influence does not just vanish instead it diminishes over time, thus for the final variation instead of the influence being fixed for a set amount of time it instead decays exponentially:

$$i_{v,u}^t = i_{v,u}^0 e^{\frac{-(t-t_v)}{\tau_{v,u}}} \quad (9)$$

With $i_{v,u}^t$ being the influence from user v on u at time t . Thus the maximum influence would be when $t = 0$. Thus the new joint probability function is:

$$i_u^t(S) = 1 - \prod_{v \in S} (1 - i_{v,u}^t) \quad (10)$$

As stated at the beginning of this section the aim is to first learn the influence probabilities, then to use these learnt values to infer the current threshold of each user dependent on their exposure to OOVs. To achieve this we use 80% of the data to train the model with the remaining 20% used to test the ability to predict OOV adoption. As mentioned above, each user has an individual threshold σ_u which must be breached by $i_u(o)$ for the user to adopt the OOV, to infer this individual threshold we use ROC curves to determine the optimum trade-off between false and true positives, as users may have been exposed to an OOV and never adopted it; this point will be the individual threshold σ_u .

4.3. Measuring Network Effect

We want to assess to what extent the network structure affects users' adoption of new language. This is broken down into two assessments, initially measuring the effect of randomising the network (Section 4.3.1) to see the effect of time and neighbours; secondly measuring the effect of community structure (Section 4.3.2) by distinguishing between inter and intra community effects.

4.3.1. Random Network. To understand the effect of network structures we will shuffle the four networks, with the aim of randomising the edges, along with the

edge times. However social networks are defined by their degree distribution, thus even though we shuffle the edges we aim to maintain the degree distribution. As proposed in [18] instead of shuffling the edges we iterate over each edge, randomly selecting an alternate edge and swapping the source of each edge; thus maintaining the degree distribution. With these four new graphs we then *relearn* the influence measures across the new networks.

4.3.2. Community Influence. In social graphs users and communities cluster together, thus we aim to see if influence between nodes is greater *internally* or *externally* to a community. As each network has been classified into distinct communities with [3], this means that intra community edges (E_{\hookleftarrow}) are ones that cross community boundaries, whereas inter edges (E_{\hookrightarrow}) are edges within the same community.

To compare the influence that exists internal and external to communities we compute the average intra influence (\bar{i}_{\hookleftarrow}) and average inter influence ($\bar{i}_{\hookrightarrow}$) across the network.

$$\bar{i}_{\hookrightarrow} = \frac{\sum_{e \in E_{\hookrightarrow}} w(e)}{|E_{\hookrightarrow}|} \quad (11)$$

Where $w(e)$ returns the influence of the given edge (e), this is divided by the number of internal edges ($|E_{\hookrightarrow}|$)

Similarly computing the intra community influence sums the influence of all external edges and then divide by the number of external edges.

$$\bar{i}_{\hookleftarrow} = \frac{\sum_{e \in E_{\hookleftarrow}} w(e)}{|E_{\hookleftarrow}|} \quad (12)$$

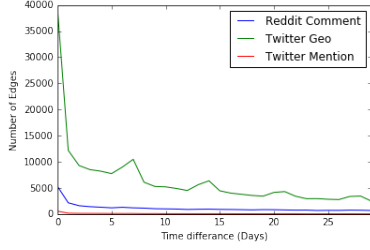
If the network structure has limited effect on the influence internal or external to a community, we would expect there to be limited difference between \bar{i}_{\hookleftarrow} and $\bar{i}_{\hookrightarrow}$. Whereas a larger inter value would indicate that the community structure is having an effect on the distribution of influence, with influence being affected by concepts such as structural trapping [20].

5. Results and Discussion

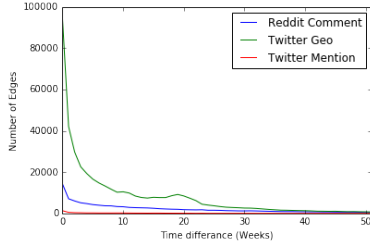
The following section outlines the main findings and results from the experiments that have been performed. The results have been split into two separate sections; section 5.1 discusses the results of learning influence and using it to predict OOV usage, were as the effect of network structure is discussed in section 5.2.

5.1. OOV Prediction

Figure 1 shows the diffusion time of OOVs across two separate granularities (Day fig. 1a, Week fig. 1b),

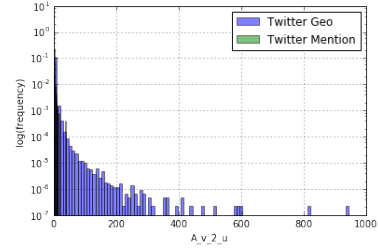


(a) Day

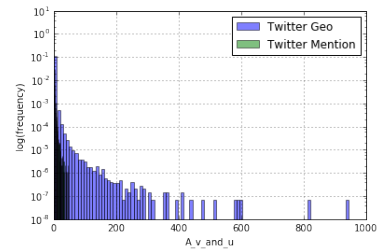


(b) Week

Figure 1: OOV diffusion frequencies



(a) O_{v2u}



(b) $O_{v\&u}$

Figure 2: Twitter-Geo and Mention network

the diffusion process follows a power law distribution, with the majority of the diffusion happening in within the first day of exposure.

To assess the accuracy of using the influence measures in learning the activation threshold of a user (see section 3) we only focus on users that have been exposed to an innovation and not users who used it without being exposed (this is due to want to predict adoption based on exposure, if there is no exposure then this would not be able to be predicted). As stated prior this prediction challenge is a binary classification, with the adoption being predicted if the joint probability is greater than the user activation threshold ($i_u(o) > \sigma_u$).

To learn each individual users activation threshold σ_u we use ROC curves; these represent the trade off between the true positive (TPR) and false positive rate (FPR) through varying the user threshold, with the aim is to find the threshold at which a user has the largest amount of True Positive and True Negatives. As can be seen in Figure 6 the results across all data set are varied, with the static models proposed in section 4.1 being able to predict with AUC highs of 0.92, though there is no discernible difference between the four variations of modelling influence. The introduction of decay functions appears to have reduced the accuracy across all models, resulting in some performing with an accuracy less than a random model ($AUC = 0.5$). This can be seen to the greatest extent in the Twitter mention network; this network is sparsely connected in comparison to other network, (Table 2), resulting in the number propagation's (O_{v2u}) being significantly

smaller than that of say the Twitter Geo network (see Figure 2). However, the overall reduction in accuracy could be due to an external unobserved process that effect language adoption, this could be in line with [10] stating that influence/exposure comes from not only the local connections but the community as a whole.

As stated in section 3 OOV's can be classified into different function sets (POS tags), with table 3 shows the AUC values for the two none credit models (equations 1 and 2) in each of the four networks. Across the board that values in table 3 show high AUC, within noticeable improvement in the Twitter Comment network. However, the majority of values are still less accurate than a random base line ($AUC < 0.5$). However abbreviations (G), verbs (V) appear to perform the greatest, potentially indicating that open class words used to describe events are more likely to be used in conversational discourse.

Unlike the static and discrete time models there is only one global max joint probability $i_u(o)$ when using the continuous time model (equation 9). This is the point at which there is the greatest amount of influence on the given user to adopt a term. To assess if the time of adoption is near the time at which there was the greatest amount of influence we compute the difference between the time that a user adopted a term and the time at which they had the greatest influence. This can be seen in figure 3, values < 0 indicate that they used the OOV before there global max, were as values > 0 indicate they used the term after there global max. There is a distinctive spike around 0 indicated that

Table 3: AUC values for each given POS tag

| Tag | Reddit Comment | | Reddit Traversal | | Twitter Comment | | Twitter Geo | |
|-----|----------------|----------|------------------|----------|-----------------|-----------------|-------------|----------|
| | Bernoulli | Jaccard | Bernoulli | Jaccard | Bernoulli | Jaccard | Bernoulli | Jaccard |
| ! | 0.945814 | 0.932778 | 0.919580 | 0.919668 | 0.596826 | 0.391674 | 0.821373 | 0.754288 |
| # | 0.988482 | 0.985960 | 0.992614 | 0.986222 | - | - | - | - |
| A | 0.974176 | 0.966811 | 0.941658 | 0.933533 | 0.372159 | 0.390327 | 0.828381 | 0.781653 |
| D | 0.916667 | 0.655914 | - | - | - | - | 0.710258 | 0.639674 |
| E | - | - | 0.997986 | 0.994964 | - | - | - | - |
| G | 0.949570 | 0.946256 | 0.956844 | 0.954782 | 0.815499 | 0.712872 | 0.890735 | 0.852988 |
| L | 0.961226 | 0.905481 | 0.970527 | 0.971248 | - | - | 0.769028 | 0.687771 |
| N | 0.919589 | 0.914614 | 0.923128 | 0.919163 | 0.606992 | 0.544365 | 0.847362 | 0.806809 |
| O | 0.945832 | 0.943943 | 0.902753 | 0.901661 | - | - | 0.801335 | 0.767615 |
| P | 0.797396 | 0.812736 | - | - | - | - | 0.836548 | 0.822249 |
| R | 0.960906 | 0.959497 | 0.936667 | 0.926402 | 0.426421 | 0.436959 | 0.817372 | 0.758002 |
| V | 0.950552 | 0.945546 | 0.943985 | 0.941717 | 0.635784 | 0.637408 | 0.819181 | 0.767474 |

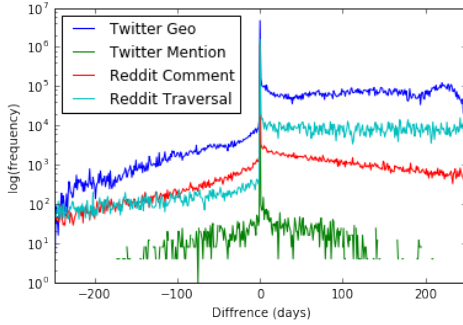


Figure 3: Time difference from global maximum of user adoption

the value of global max is on the same day they use the OOV, however long tails that reduce in frequency that further from the centre indicate that influence may decay faster than that modelled.

5.2. Network Structure

However, to what extent does the network effect influence between users or communities. Figure 4 plots the distributions of influence internal and external to each community. The aim was to see a statistically significant difference, however, the results are inconclusive showing that there are no significant difference. Though, internal to a community the spread of influence is greater, with the majority of influence being higher across the four networks internal to a network. Though, due to the sparsity of the Twitter mention network resulted in the greatest ranges of values, though with the same mean and higher external influence. Thus influence may not be effected by structural trapping as proposed in [20].

Though to what extent does the influence depend

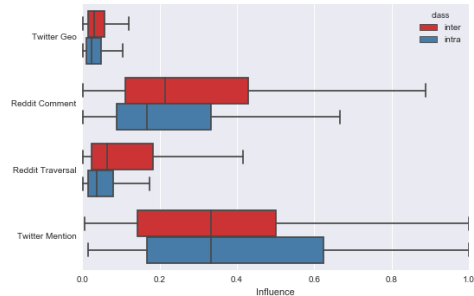


Figure 4: Influence distributions for internal (inter) and external (intra) to a community.

on the network structure, Figure 5 shows the results from shuffling the Reddit comment network. As one can see for the static time models there is roughly a 75% reduction in *AUC*, however when looking at the time dependent models the accuracy falls below that of the random base line model with values as low as 0.07. When inspecting the values from the other networks, similar patterns can be seen where the static time model has a reduction in accuracy with the time dependent models reducing to a larger extent. As [10] stated that influence might not only come from your local acquaintances, but from the ether across the whole network, however when assessing the decay models it could be that a users adoption of an OOV is dependent on the network, though the time at which the OOV is adopted is dependent on the local network.

Comparing the results on both the micro (user interactions) and macro (community interactions) networks one can see that there is a higher accuracy in predicting when a community rather than an individual will adopt a new term. This could indicate that the collective influence is greater than that of the individual, showing that it is not an individual that affects change but rather a

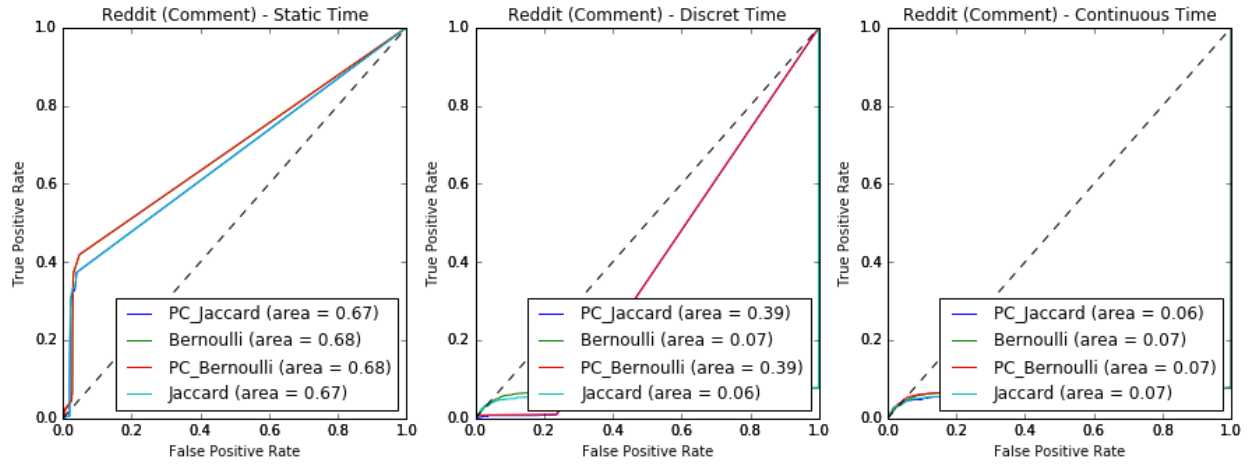


Figure 5: ROC curves for learnt on the shuffled version of the Reddit Comment network

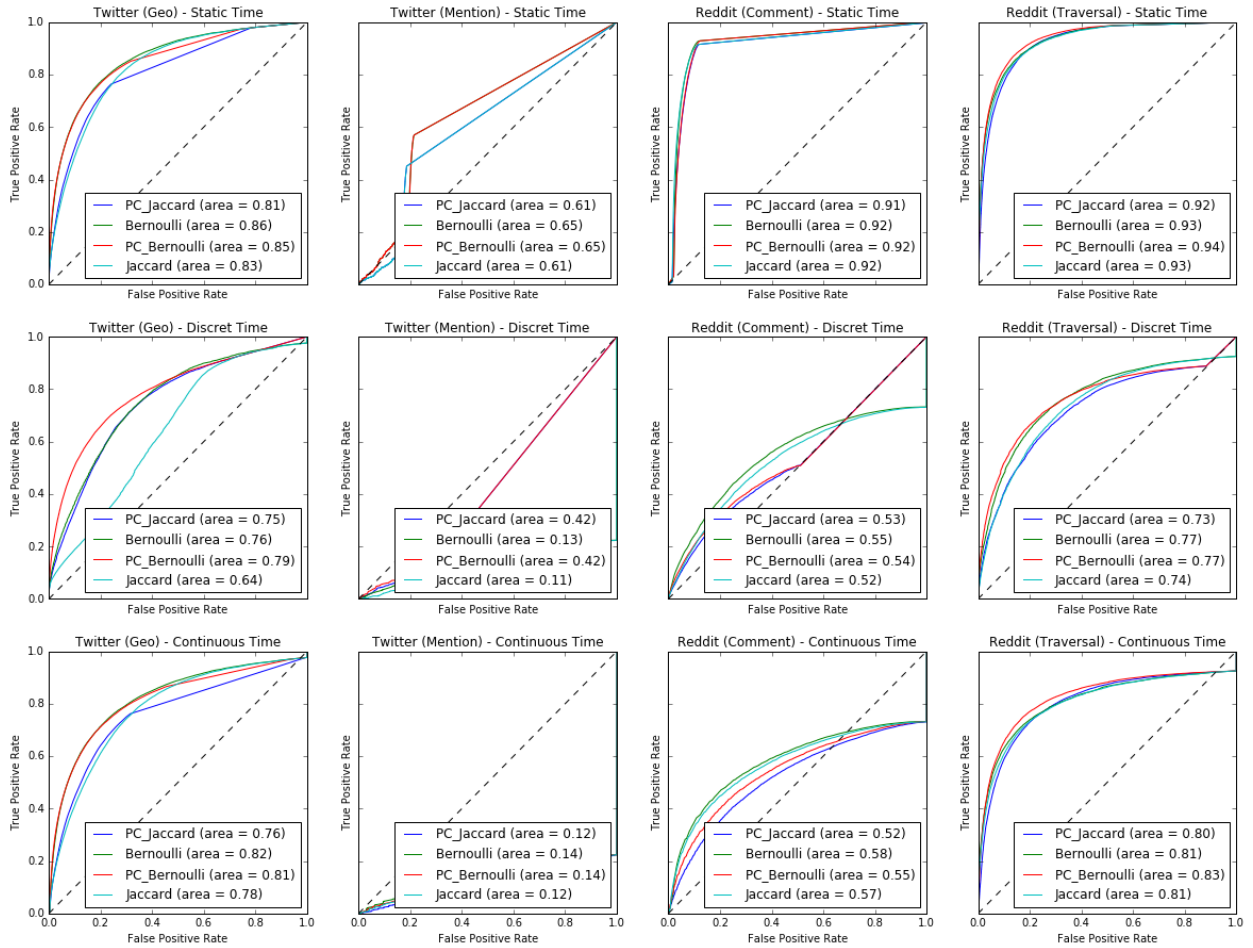


Figure 6: ROC curves for each of the four network, using the three different influence models.

collective. However for research purposes we classified an adoption as one when of term is only adopted once in a community. In reality for a community to ‘adopt’ a given term, it would more likely be when the majority of users in the community use the term, or when it has been used more than a certain number of times.

6. Conclusion and Future work

The novel contribution in this paper lies in the application of the framework proposed by [9] to the field of OOV adoption prediction, across multiple large-scale social networks. By drawing on [10], [17] the results show a high accuracy in predicting language adoption when learning the individual user activation threshold. We also show that there is a potential underlying background process affecting the adoption of language, this was shown in the decreased accuracy when modelling decays in influence, and by testing the model on random networks. Yet, our results also show that there is little dependency on the community structure in the propagation of influence, potentially showing that in an on-line world where users can move about freely there is less constant membership with one community, as users change their membership frequently.

Limitations for this work mainly from the available data not being able to mode true inter user relationship. As constructing a network from consecutive user movements/posts (*macro*) does not captured the true exposure users experience to OOV’s, as they may not interact with a community in which they saw the OOV first. Finally language change happens over an extended period of time, however this research does not separate between ‘bursty’ words (such as new product names) and word the grow over an extended period of time, thus blurring the definition of ‘meme’/‘language change’.

Future work will look in greater detail the effects of community structures and position of users within a network when predicting the final size of a OOV diffusion rather than if a user will adopt an innovation. This will look at two core concepts, strength of weak ties and access to structural holes, aiming to quantify their effect on the final diffusion size of an OOV. The belief is that users who have access to structural holes will cause large diffusions of OOV, whereas an OOV may become trapped when used in communities with strong internal bonding.

7. Data access statements

All data and code created during this research are openly available from Lancaster University data archive at <http://dx.doi.org/10.17635/lancaster/researchdata/99>.

8. Acknowledgements

This work is funded by the Digital Economy programme (RCUK Grant EP/G037582/1), which supports the High-Wire Centre for Doctoral Training (<http://highwire.lancaster.ac.uk>).

References

- [1] S. Aral, L. Muchnik, and A. Sundararajan. Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks. *PNAS*, 106(51):21544–21549, 2009.
- [2] G. Aston and L. Burnard. *The BNC handbook: exploring the British National Corpus with SARA*. Capstone, 1998.
- [3] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, physics.soc-ph(10), Oct. 2008.
- [4] E. Carmel and R. Agarwal. Tactical Approaches for Alleviating Distance in Global Software Development. *IEEE Software*, 18(2), Mar. 2001.
- [5] D. Crystal. *Txtng: The Gr8 Db8*. OUP Oxford, July 2009.
- [6] C. Danescu-Niculescu-Mizil, R. West, D. Jurafsky, J. Leskovec, and C. Potts. *No country for old members: user lifecycle and linguistic change in online communities*. WWW, May 2013.
- [7] J. Eisenstein, B. O’Connor, N. A. Smith, and E. P. Xing. Mapping the geographical diffusion of new words. *arXiv.org*, page 5268, Oct. 2012.
- [8] K. Gimpel, N. Schneider, B. O’Connor, D. Das, D. Mills, J. Eisenstein, M. Heilman, D. Yogatama, J. Flanigan, and N. A. Smith. Part-of-speech tagging for Twitter: annotation, features, and experiments. In *HLT ’11*. Association for Computational Linguistics, June 2011.
- [9] A. Goyal, F. Bonchi, and L. V. S. Lakshmanan. Learning influence probabilities in social networks. In *WSDM ’10*. ACM Request Permissions, Feb. 2010.
- [10] M. Granovetter. Threshold Models of Collective Behavior. *American journal of sociology*, 83(6):1420–1443, May 1978.
- [11] B. Han, P. Cook, and T. Baldwin. Geolocation prediction in social media data by finding location indicative words. In *COLING 2012*, pages 1045–1062. University of Melbourne, Parkville, Australia, Dec. 2012.
- [12] E. Katz and P. F. Lazarsfeld. *Personal Influence, the Part Played by People in the Flow of Mass Communications*. Transaction Publishers, 1955.
- [13] D. Kershaw, P. Stacey, and M. Rowe. Towards modelling language innovation acceptance in online social networks. *WSDM’16*, 2015.
- [14] R. S. Olson and Z. P. Neal. Navigating the massive world of reddit: Using backbone networks to map user interests in social media. *arXiv.org*, page 3387, Dec. 2013.
- [15] E. M. Rogers. *Diffusion of Innovations, 5th Edition*. Simon and Schuster, Aug. 2003.
- [16] G. Szabo and B. A. Huberman. Predicting the popularity of online content. *Communications of the ACM*, 53(8), Aug. 2010.
- [17] T. W. Valente. Social network thresholds in the diffusion of innovations. 1996.
- [18] F. Viger and M. Latapy. Efficient and simple generation of random simple connected graphs with prescribed degree sequence. *Journal of Complex Networks*, 4(1):15–37, Feb. 2016.
- [19] L. Weng and F. Menczer. Topicality and Social Impact: Diverse Messages but Focused Messengers. *arXiv.org*, (2):e0118410, Feb. 2014.
- [20] L. Weng, F. Menczer, and Y.-Y. Ahn. Virality Prediction and Community Structure in Social Networks. *arXiv.org*, June 2013.