Inferring Political Opinion From Social Media Data

Angus Scott

4th Year Project Report
Artificial Intelligence and Computer Science
School of Informatics
University of Edinburgh

2015

Abstract

This is an example of infthesis style. The file skeleton.tex generates this document and can be used to get a "skeleton" for your thesis. The abstract should summarise your report and fit in the space on the first page. You may, of course, use any other software to write your report, as long as you follow the same style. That means: producing a title page as given here, and including a table of contents and bibliography.

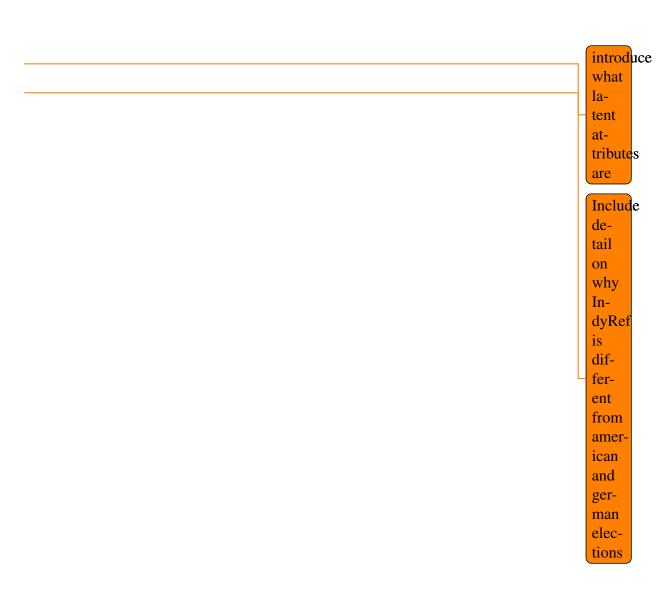
Acknowledgements

Insert Acknowledgements here.

Table of Contents

1	Introduction		3	
2 Interim Report Only			4	
3	Related Work			
	3.1	Network Structure	5	
	3.2	User Communication	6	
4	Data and Methodolgy			
	4.1	Scotland Independence Referendum Background	7	
	4.2	The Twitter Platform	7	
	4.3	Dataset	8	
	4.4	4.4 General approach to user profiling		
	4.5	4.5 Models		
		4.5.1 Naive Bayes	9	
		4.5.2 Topic Modelling	9	
		4.5.3 SVM	9	
		4.5.4 Gradient Boosted Decision Trees	9	
5	Eva	luation	10	
6	Con	iclusions and Further Work	11	

Introduction



Interim Report Only

Included todo notes for sections as I have ideas.

Covered sections: Completed sections on introduction to Twitter Platform, and Dataset. Plan to work on Related Work section: Need to outline more detail, and plan to include some heavy detail on Topic Modelling, as I'll be covering these in detail in other courses.

Plan to build models over the next couple of weeks, current schedule: By week 4 Finish annotation and build SVM and NB models for greater set of features than hashtags. Plan to use hashtags model as baseline. week 5/6 Get cross annotator scores, and get Topic Models done

Related Work

The inference of latent attributes from different sources of media has been well considered both in the context of traditional media and new media.

A natural extension to this task has been to apply similar methods to the microblogging service Twitter. Previous work includes detection of *gender*, *ethinicity*, *brand loyalty*, *regional origin* and *age*. [?] [?]

The detection of *Political Affiliation* has also been considered. Most of the research in this area has been focused on one of two political systems. First is the American system, with the 2010 Congressional Midterms [?] [?] [?] or the 2012 Presidential Elections [?]. Second is the German system, with the 2009 federal election to the national parliament [?] [?].

Approaches to the task of inference of latent attributes from twitter can be broadly split into two categories:

- 1. Use of the networked structure of the data
- 2. Analysis of user communication streams

3.1 Network Structure

It has long been known that individuals form information networks corresponding to their own political preferences [?], and this has continued into online social networks [?]. This demonstrates that people who participate in particular networks, are likely to share similar political values or opinions, which could be used to infer someones political beliefs.

There are four main types of networks that have previously been exploited to infer political persuasion:

- Follower -
- Mention -

include links to papers that have per-ofrmed tasks in different me-

dia

- Retweet -
- Hyperlink -

3.2 User Communication

Sentiment Analysis Topic Modelling Support Vector Machines Decision Trees

Data and Methodolgy

4.1 Scotland Independence Referendum Background

4.2 The Twitter Platform

Twitter is a popular microblogging and social network site that allows users to send and read short 140 character messages commonly known as *tweets*. Tweets as used in this report include the 140 character message, along with the associated metadata. This includes features such as: creation dates, user follower counts, coordinates and language.

As well as tweeting¹ to an audience of *followers*, Twitter users can interact publicly through two main approaches: *retweets* and *mentions*. Retweets often indicate agreement [?], allowing users to rebroadcast content from other users, having the effect of spreading tweets to a larger audience [?]. Alternatively, mentions work by allowing someone to refer to a particular Twitter user by including their @username in a tweet, creating a public dialogue between the referrer and referee.

Hashtags are another important feature of Twitter, allowing users to tag tweets according to topic and their intended audience. For example #IndyRef was often used to reference the topic of the Independence Referendum, or #bettertogether which generally indicated that someone was No Voter or was addressing No Voters.

Twitter has several benefits over other social networking sites [?] [?] including:

- 1. Twitter users retweet notable events and participate in the spread of realtime news.
- 2. The 140 character constraint forces tweets to be concise and to the point.
- 3. Users are highly reactionary and discussed events tend to have happened in the recent past.

¹Tweeting - The act of broadcasting a tweet

Topic Indicators:	<pre>#indyref, #scotland,</pre>
	<pre>#scotdecides, #scotlanddecides,</pre>
	#independence
Yes Campaign:	<pre>#voteyes, #yesscotland,</pre>
	#yesscot, #youyesyet
No Campaign:	<pre>#voteno, #nothanks,</pre>
	<pre>#bettertogether,</pre>
	#letsstaytogether
Political Parties/Politicians:	#snp, #labour, #conservative,
	<pre>#tory, #libdem, #salmond,</pre>
	<pre>#alexsalmond, #davidcameron,</pre>
	#gordonbrown, #alistairdarling

Table 4.1: Set of hashtags related to the campaign

4. Tweets are media rich, and include content like video, image and hyperlinks along with text.

4.3 Dataset

The dataset was collected following the referendum, and is composed of X tweets from Y users, spanning from the 1st February 2014, to 17th September 2014, the day prior to the Independence Referendum. The corpus also includes set of user classifications obtained through annotation of users tweets from the 18th of September.

The first steps in building the corpus, was to obtain tweets from the 18th of September 2014. This first set of tweets were obtained from a Twitter Mining system developed by Sasa Petrovic and Miles Osborne as part of their efforts to create the Edinburgh Twitter Corpus. From this I was given a collection of tweets that matched a set of hashtags commonly used by users involved in the Yes and No campaigns.

Hashtags, and their associated campaign are given in Table 4.1, which were obtained using the hashtagify.me tool. It is worth noting that this is a set of non-event specific hashtags. Whilst hashtags such as #PatronisingBTlady [?] are politically divisive, they have been excluded as they have a short lifespan as discussed in Section 4.2, and are therefore unlikely to be used on the final day of the referendum.

These tweets are then used as evidence for the annotators to classify a user as Y (Voted Yes in the referendum), N (Voted No in the referendum) or U (Unknown or Undeterminable). We make the assumption that tweets on the final day are the most emblematic of how a user voted, as the opportunity to flip their opinion is minimised and Twitter activity surrounding the event was at it's peak. This is done as there are both technical and ethical issues about obtaining classifications more directly.

After providing a class for each user, we then collect all their publicly available tweets, available from 1st February to 17th September, and their associated meta data using

Add number of users here

Include graph showing strength of hash-tag, and how many tweets it appears in

the Twitter API ².

4.4 General approach to user profiling

4.5 Models

4.5.1 Naive Bayes

4.5.2 Topic Modelling

4.5.3 SVM

4.5.3.1 One class SVMs

4.5.3.2 Alternative SVMs, standard?

4.5.4 Gradient Boosted Decision Trees

Detail what can be measured in Twitter and how we could use it

I plan on implementing, may be moved to eval

section

Models

²Due to a limitation of the Twitter API, you can only collect from the last 3200 published tweets.

Evaluation

Chapter 6 Conclusions and Further Work