



Who Tweets? Deriving Demographic Information from Twitter

Dr Luke Sloan

Email: SloanLS@cardiff.ac.uk

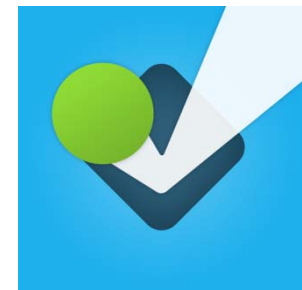
Follow: @drlukesloan

Collaborative Online Social Media ObServatory (COSMOS)

Luke Sloan, William Housley, Matthew Williams, Adam Edwards, Pete Burnap,
Jeffrey Morgan, Rob Proctor & Omer Rana

Outline

- The Digital Revolution
- Emergence of Social Media
- Reflections on Representation
- 'Naturally Occurring' Data
- COSMOS
- Deriving Demographics
- Current Work



The Digital Revolution

- Every minute...
 - 48 hours of video uploaded to YouTube
 - 204,166,667 emails are sent
 - 2,000,000 search queries on Google
 - 217 users join the mobile web
 - 571 new website are created
 - 100,000+ tweets are made
 - 3,125 photos added to Flickr
 - 684,478 pieces of content shared on Facebook

Source: <http://mashable.com/2012/06/22/data-created-every-minute/> [accessed April 2014]

It's all data!

Emergence of Social Media I

- Social Media is potentially a rich source of naturally occurring data on beliefs, attitudes, reactions and opinions
- For example, Twitter can be used for...
 - Brand tracking with sentiment (Scarfi 2012)
 - Predicting movie revenue (Asur & Huberman 2010)
 - Advance earthquake warning (Sakaki et al. 2010)
 - Predicting election results (Tumasjan et al. 2010)

Unlike traditional social science data collection, this can all be gathered for free...

Emergence of Social Media II

- The volume of data is phenomenal compared to a social survey...
 - ‘Spritzer’ at 1% approx. 3.5m tweets a day (free)
 - ‘Garden Hose’ at 10% approx 35m tweets a day (make a case)
 - ‘Fire Hose’ at 100% approx. 350m tweets a day (payment only)

Reflections on Representation

- But is Twitter representative?
 - Representative of who?
 - Population data for Twitter users (transactional)
- If Twitter data can be used to predict (*understand*) real world events then it should be representative at some level (but young people tweet, young people go to the cinema?)
- Much potential, but must be realised in a (social) scientific way
- At the heart of this is finding out who Tweets...

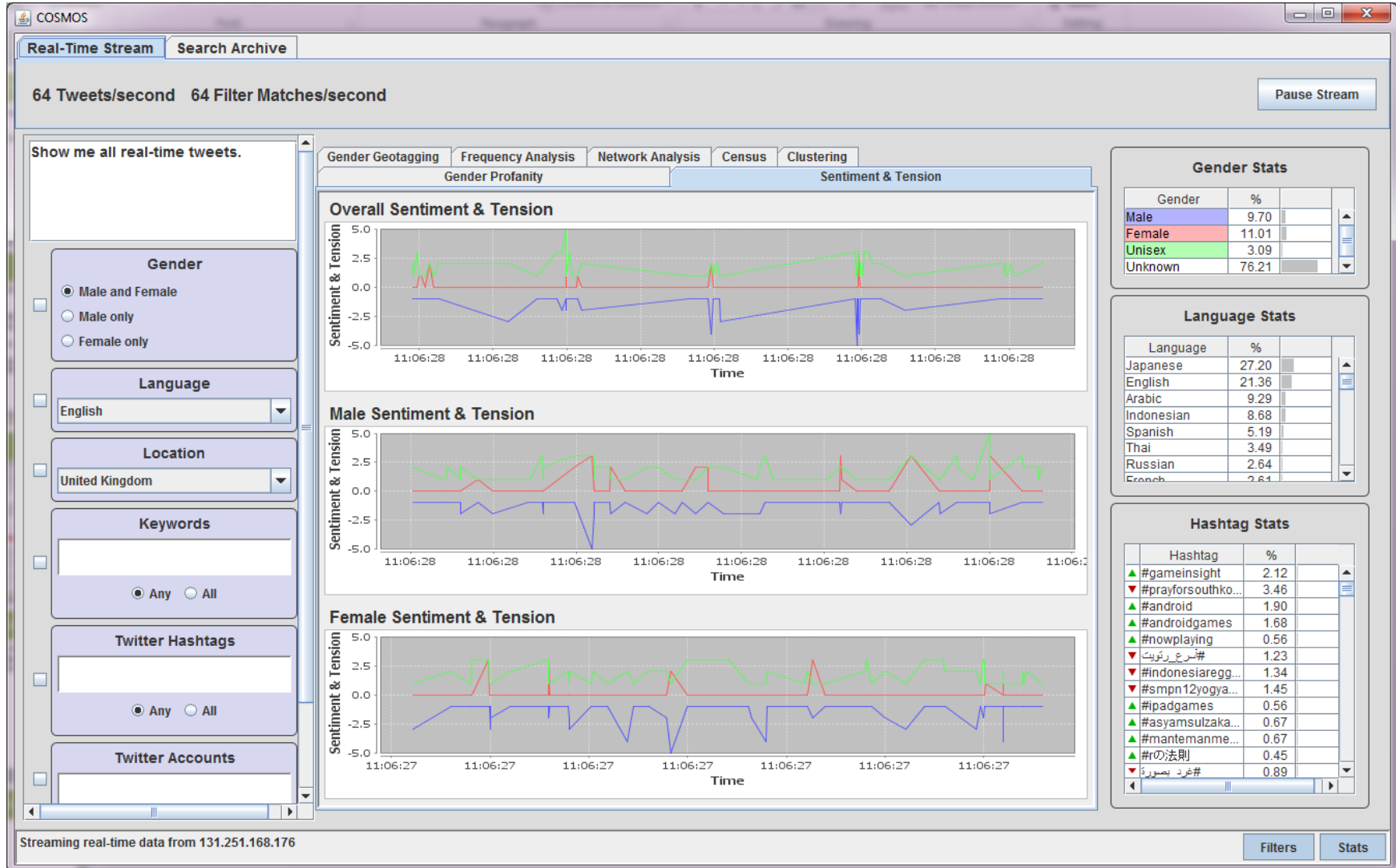
'Naturally Occurring' Data

- User generated content can be 'data light' (Mislove et al. 2011, Gayo-Avello 2012)
- Facebook is different because it stores baseline demographic information (Schwartz et al. 2013)
- Twitter has signatures, but nothing systematic
- When the data is not available we develop proxies, so why not for Twitter?

COSMOS I

- The Collaborative Online Social Media ObServatory (COSMOS) funded by...
- Develop a platform for data interoperability between:
 - naturally occurring data (e.g. social media)
 - curated data (e.g. social surveys)
 - administrative data (e.g. A&E admissions)
- Social media augments traditional social scientific investigation – it is not a surrogate! (Edwards et al. 2013)
- A key programme of work within COSMOS is on the derivation of demographic data from ‘signatures’ (aka *exhaust fumes*)

COSMOS III



COSMOS IV

COSMOS

Real-Time Stream Search Archive

46 Tweets/second 46 Filter Matches/second Pause Stream

Show me all real-time tweets.

Gender Geotagging Frequency Analysis Network Analysis **Census** Clustering

Gender Profanity Sentiment & Tension

Newham: 170,268

48% Unemployment **Ethnicity**

White British 82,390 33.8%
 White Irish 3,231 1.3%
 Other White 10,509 4.3%
 Mixed White & Black Carri 2,986 1.2%
 Mixed White & Black Afric 1,657 0.7%
 Mixed White & Asian 1,652 0.7%
 Mixed Other White 1,953 0.8%
 Asian British Indian 29,597 12.1%
 Asian British Pakistani 20,644 8.5%

Crime

Click a polygon or marker on the map for neighbourhood crime data

Newham 00BB (541723, 182438) (51.522571, 0.043046)

Colour Boroughs by Unemployment Show Tweet Markers Colour Tweet Markers by Gender

Gender Stats

Gender	%
Male	9.53
Female	11.25
Unisex	3.25
Unknown	75.97

Language Stats

Language	%
Japanese	26.91
English	22.15
Indonesian	8.98
Arabic	8.96
Spanish	4.85
Thai	3.15
Russian	2.56
French	2.46

Hashtag Stats

Hashtag	%
#prayforsouthko...	2.71
#برع رتويت	1.20
#gameinsight	1.68
#smpn12yogya...	1.50
#android	1.50
#indonesiaregg...	1.32
#androidgames	1.08
#tonygrastafara	1.02
#ipad	0.66
#السعودية	0.60
#برتويت	0.48
#asyamsulzaka...	0.60
#الكرت	0.42

Streaming real-time data from 131.251.168.176

Filters Stats

Deriving Demographics

- Development of tools embedded within the COSMOS platform to identify signatures of demographic characteristics (Sloan et al. 2013)
- Location
- Gender
- Language

Deriving Demographics: Location I

- Three primary sources of location:
 - User profile information
 - Content of tweets (inc. ‘mundane geography’)
 - Geo-tagged tweets
- Geo-tagged tweets are the gold standard
- Allows us to locate people at the time they tweeted in existing geographies (output area level!)
- RQ: do people tweet about crime in high crime areas?

Deriving Demographics: Location II



Source: Sloan et al. 2013

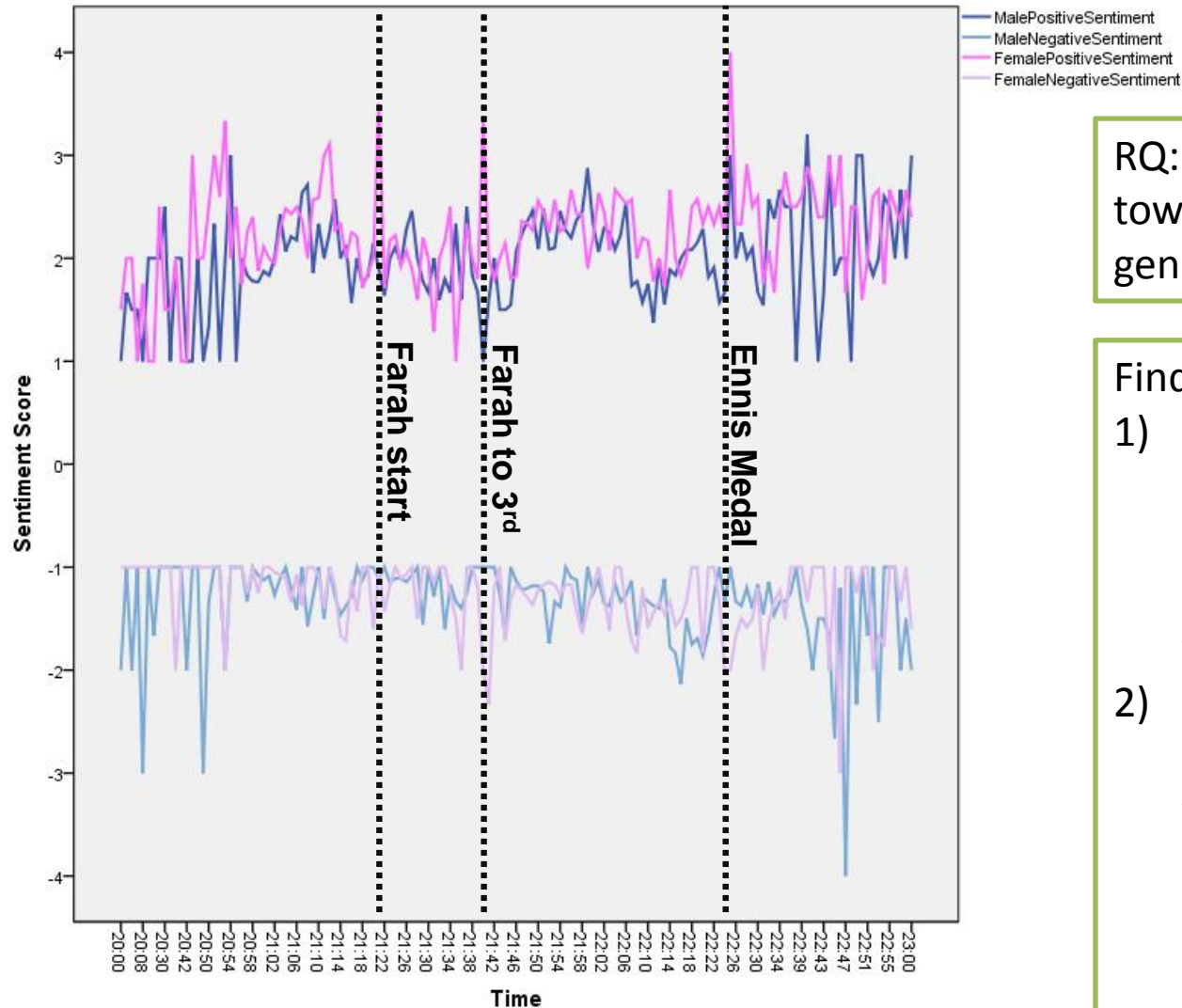
Deriving Demographics: Language

- Two methods of identifying language:
 - The language of the Twitter interface
 - The language of the Tweet (Java library)
- Detecting language is important for efficiency of other analytical tools (e.g. sentiment analysis)
- 40% of content is in English
- RQ: do spatial patterns of language use recorded on the Census correspond with Twitter maps?
- Note that 'hard to reach groups' may use Twitter...

Deriving Demographics: Gender I

- Use the name field of the Twitter profile
- Clean the data to extract a first name and compare against a large database of first names
- Important to categorise 'unisex' and 'unknown'
- Of those we could identify: 48.8% male and 51.2% female... exactly the same as the 2011 Census

Deriving Demographics: Gender II



RQ: How does sentiment towards Team GB differ by gender?

Findings:

- 1) Sentiment peaks reflect real world events (relationship between social media and real world)
- 2) Sentiment differs between men and women (difference is so pronounced that gender detection method appears to work)

Current Work

- Identifying age from signature data:
 - Preliminary analysis suggests usable age data for 0.35% of Twitter users
 - Note that 0.35% of 645m is 2.25m (approx 40% of which is English language)
- Identify occupation from signature data:
 - Linked to SOC2010 codes
 - Enables allocation into NS-SEC groups

Discussion

References

Asur & Huberman (2010) Predicting the Future with Social Media, *Social Computing Lab*, HP Labs in Palo Alto

Edwards et al. (2013) Computational social science and methodological innovation: surrogacy, augmentation or reorientation?, *International Journal of Social Research Methods*, 16:3

Gayo-Avello (2012) I wanted to Predict Elections with Twitter and all I got was this Lousy Paper: A Balanced Survey on Election Prediction using Twitter Data, *Department of Computer Science*, University of Oviedo Spain

Mislove et al. (2011) Understanding the demographics of Twitter users, *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*

Sakaki et al. (2010) Earthquake shakes Twitter users: real-time event detection by social sensors, *presented at WWW 2010*, Raleigh, NC USA

Scarfi (2012) Social Media and the Big Data Explosion, *Forbes* (www.forbes.com/sites/onmarketing/2012/06/28/social-media-and-the-big-data-explosion/)

Schwartz et al. (2013) Personality, Gender, and Age in the Language of Social Media: The Open-Vocabulary Approach, *PLOS ONE*, 8:9 (DOI: 10.1371/journal.pone.0073791)

Sloan et al. (2013) Knowing the Tweeters: Deriving Sociologically Relevant Demographics from Twitter, *Sociological Research Online*, 18:3 (<http://www.socresonline.org.uk/18/3/7.html>)

Tumasjan et al. (2010) Predicting elections with Twitter: what 140 characters reveal about political sentiment, *Proceedings of the Fourth International AAAI Conference on Weblogs and Social Media*



Who Tweets?

Dr Luke Sloan

Email: SloanLS@cardiff.ac.uk

Follow: @drlukesloan

Collaborative Online Social Media ObServatory (COSMOS)

Luke Sloan, William Housley, Matthew Williams, Adam Edwards, Pete Burnap,
Jeffrey Morgan, Rob Proctor & Omer Rana