

Linking survey and social media data in Understanding Society

Tarek Al BaghalISER University of Essex









Acknowledgments

Curtis Jessop – NatCen Social Research

Luke Sloan – University of Cardiff

Alexander Wenz – University of Mannheim

Social Media (in the UK)

- 2011: 45% access Internet to use social media
- 2020: 70% access Internet to use social media
 - 97% of 16-24; 91% of 25-34; 90% of 35-44
- 79% say they <u>use</u> Facebook account
- 47% say they <u>use</u> Twitter account
 - Likely skewed online panel data
- 33 million monthly-active FB users (2013)
- 15 million monthly-active Twitter users (2013)

What are we trying, and why?

 Link survey participants' answers to publicly available information from their Twitter accounts

 Survey data benefits from real-time, 'natural' behavioural and attitudinal data

Between wave measures(?)

What are we trying, and why?

Adds the 'who' to Twitter data – creates a sample frame

Analysis of different groups on Twitter

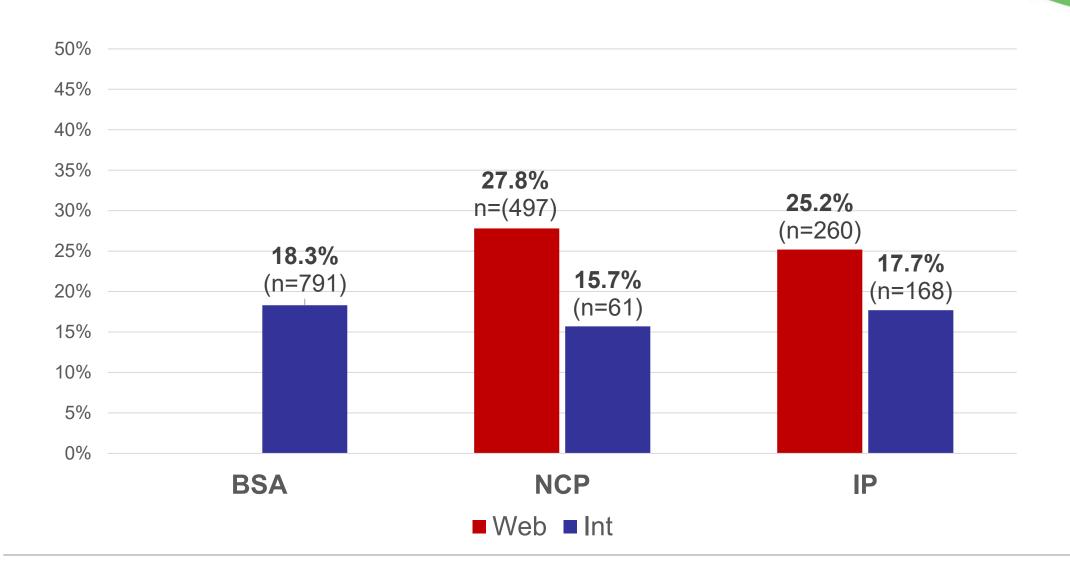
Find acceptable archiving methods

Complement, not contrast!

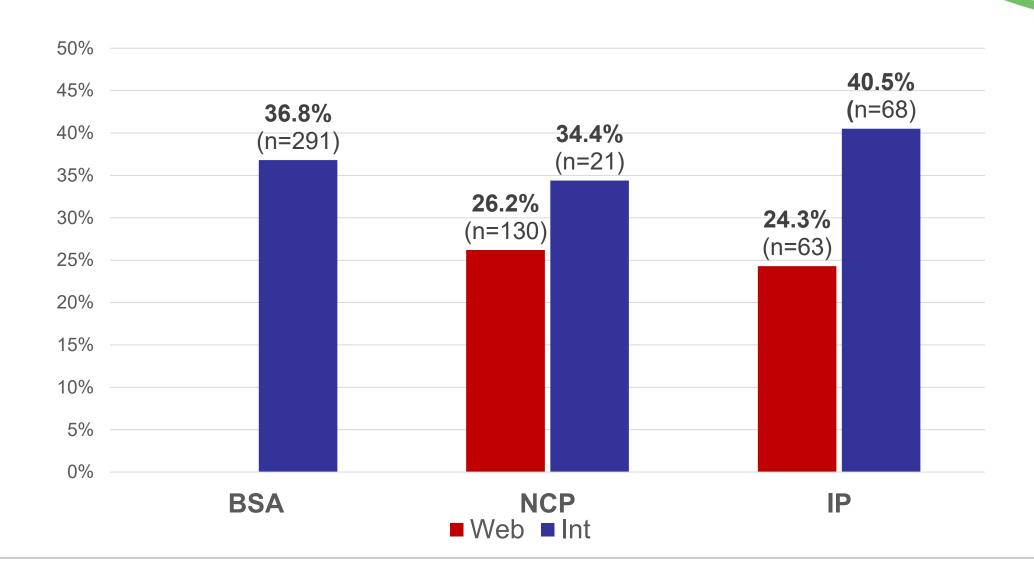
Consent in 3 UK Studies

- British Social Attitudes (2015)
 - Cross-sectional, F2F
- NatCen Panel (July 2017)
 - Longitudinal, BSA 2015,2016 samples
 - Sequential MM Web CATI
- Innovation Panel (2017)
 - Longitudinal, annual since 2008
 - CAPI, & Sequential MM Web CAPI

How many Twitter Users?



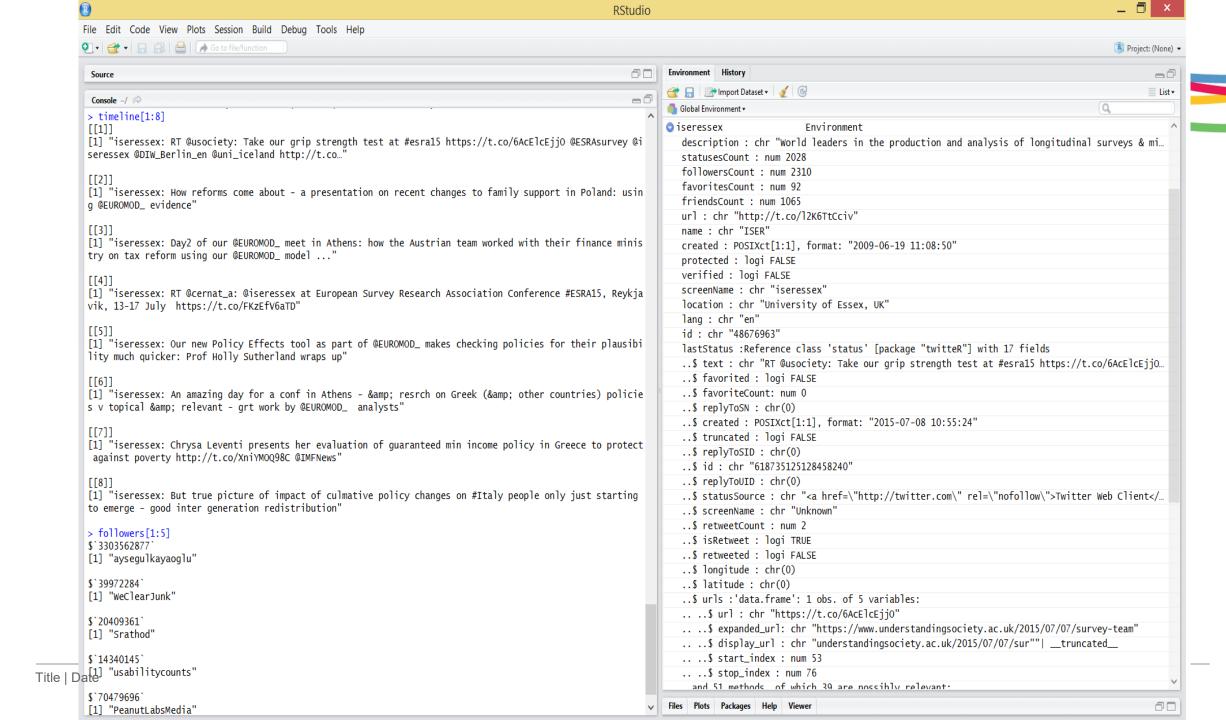
Consent Among Users



Significant Characteristics

- British Social Attitudes (2015)
 - Age: Younger Older
- NatCen Panel (July 2017)
 - Age: Younger Older
 - Sex: Male Female Female
- Innovation Panel (2017)
 - Mode: F2F
 Web





Programs for Coding the Data

- Tidyverse, tm (text mining), sentiment in R
- Linguistic Inquiry and Word Count (LIWC)
 - Extensive built-in dictionary
 - Text mines, outputting 80 variables.
 - Linguistic
 - Psychological
 - Social and biological processes
 - Beliefs
 - Socio-economic issues

Comparing Data Sources

- Can compare at micro-level
- Measures may show similarity, dissimilarity

E.g, partner status

- Are surveys gold standard in this case?
- Who is discordant?

New measures

- Can create new measures, e.g.
 - Social networks, strength of ties
 - Attitudes, "likes"

Change on more continuous dimension

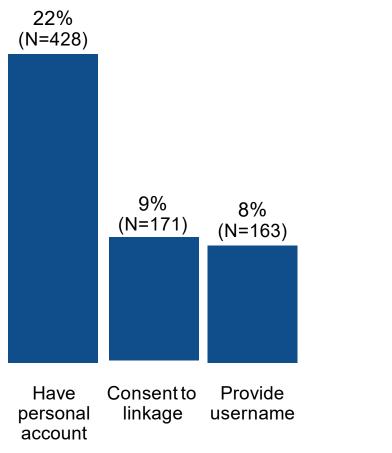
Reduction of respondent burden?

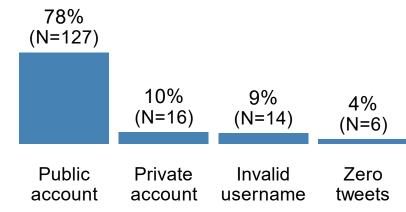
Relate to captured survey measures

Impact of Data Quantity

- What amount of Twitter data can be collected from respondents in a longitudinal survey?
- Amount can impact capture of signal in the noise
- Increase in variance, reduction in information
- Is there potential bias in substantive analyses?

Respondent linkage IP





Total Respondents: N=1,945.

Amount of Twitter Data Available

	Median	Mean	SD	Min	Max
Scraped tweets	304	933.63	1157.60	1	3199
Total tweets	306	2255.32	6057.36	1	36451
Followers	71	260.25	568.95	1	3734
Accounts followed	182	350.95	567.54	0	3912

Amount and respondent characteristics

Regression of total number of tweets (log)

- Female
- A-level or professional degree
- Number of Twitter followers
- Number of Twitter accounts followed ↔
- Frequency of Internet use ↔
- Age ↔
- Ethnicity ↔
- Marital status ↔
- HH income ↔
- Employment status ↔

Potential Bias

Relationship between

- Survey-based measure of general mental well-being
- Amount of tweets with positive/negative sentiment

Question from General Health Questionnaire (GHQ)

 "Have you recently been feeling reasonably happy, all things considered?"

More so than usual ⇒ Happy About the same as usual

Less so than usual

Much less than usual

⇒ Unhappy

Sentiment Analysis of Tweets

Words coded +/- based on Bing lexicon (Liu 2015)

- Sentiment score calculated per tweet:
 - sentiment = words_{positive} words_{negative}

Sum of +/- tweets calculated per respondent

Impact of Amount on Outcomes

Regression of unhappy response

- Number of negative tweets
- Number of followers 1
- Female
- Employed 1
- Number of positive tweets ↔
- Number of accounts followed ↔
- Age ↔
- Education ↔
- Ethnicity ↔
- Marital status ↔
- HH income ↔

More practical issues

Twitter allows for limited time-frames in aggregate

- Limited number of tweets.
 - Last 3200 from user, 5000 w/ a given keyword
- No control on content

The obsolescence of Twitter (or others)?

Secure access to linked data

Existing process to access to identifiable survey data

- Quasi-anonymisation & cut-down datasets
- Consideration of justification for research
- Training/accreditation of researchers
- Documentation of access
- Access to raw data in a secure environment

Offline access (if possible)

Not able to take data away (without review)

Archiving and Sharing

Archiving and sharing of data is important:

- Replication of results
- Maximise value of data

Particular issues:

- Who is responsible for maintaining the data?
- Deleted Tweets/withdrawn consent
 Multiple consent requests in longitudinal survey?
- Legal issues of sharing Twitter datasets

Future Plans

- Long-term: Archive, Expand to main USoc
- Improve Survey Measures
 - New Measures and Error
 - Item/Unit Nonresponse
- Improve Twitter Measures
 - E.g. Improve Prediction
- Use Knowledge for Expanded Study
 - Focus on Consent
 - And target users, specific topic survey