# Social Media Data

DSIER [/dɪˈzaɪər/]

Julian Hinz

**Bielefeld University** 





SOURCES ENDS ANDRISE, COMMAN ADRESING RESULTS, AND ANNOUNCEMENTS, CHALLE TICHARA, CCIMI ADMIONT SOCIAL IDDA UBS ANN ANT ERREENTUNDLE INNYRAULA DIREE ADDE INNYRAULE DAN RESINGEN AND O RIME REGORD BARD ON HE UNIEND REGORD MARKING COMPANABLEMIST SOURCE LAX, AND MENDAMING AND ANDRES AND ANNOUNCEMENTATIONALE INNYRAULA SIGNIFICANT SOURCE DAN RESINGEN AND ONNORES IN EPORTING ARRONGES WULLES ARE NOT COMMANDE WITH THOSE RELIGIES IN REVOLDS EPORTS. AND ANNOUNCEMENTS SOOM AND AUS EVOLVER TO INFORMET TO COMMANDES TO COMMANDES AND AND RESINCE AND ANNOUNCEMENTS. CHALLES AND ANNOUNCEMENTS AND ANNOUNCEMENTS AND AND ANNOUNCEMENTS AND ANNOUNCEMENTS.

90

we are social <sup>®</sup> Hootsuite





FACEBOOK <sup>1</sup>							2,910
YOUTUBE <sup>2</sup>	ar	e 🥗 Hoo	otsuite <sup>.</sup>			2,562	
WHATSAPP1*	so				2,000		
INSTAGRAM <sup>2</sup>				1,478			
WECHAT			1,263				
тікток			1,000				
FB MESSENGER <sup>2</sup>			988				
DOUYIN <sup>3</sup>	600						
QQ1	574						
SINA WEIBO'	573						
KUAISHOU	573						
SNAPCHAT <sup>2</sup>	557						
TELEGRAM	550						
PINTEREST	444						
TWITTER <sup>2</sup>	436						
REDDIT!*	430						
QUORA1* 300							
99 SOURCES: KEPIOS ANALYS USERS (NOTE THAT MONTI UPDATED USER FIGURES IN		DUNCEMENTS OF MONTHLY A IAY BE HIGHER]. <b>ADVISORY:</b> I FIGURES ARE LESS REPRESENT	ACTIVE USERS; (2) PLATFORMS' SELF-SE USERS MAY NOT REPRESENT UNIQUE I 'ATIVE: BASE CHANGES AND METHOD	RVICE ADVERTISING RESOURCES; [ <b>3</b> ] COMPANY A NOIVIDUAIS, <b>COMPARABILITY:</b> PLATFORMS IDEN OLOGY CHANGES; DATA MAY NOT BE DIRECTLY (		we are social	🥙 Hootsuite

- Twitter, LinkedIn, Facebook, Instagram, TikTok, ...
- Content, but also metadata
- (Used to?) provide some data access
  - $\rightarrow \,\, \text{currently in flux}$

- Twitter, LinkedIn, Facebook, Instagram, TikTok, ...
- Content, but also metadata
- (Used to?) provide some data access
  - $\rightarrow \, {\rm currently} \, {\rm in} \, {\rm flux}$

- Twitter, LinkedIn, Facebook, Instagram, TikTok, ...
- Content, but also metadata
- (Used to?) provide some data access
  - $\rightarrow \,\, {\rm currently}$  in flux

- Twitter, LinkedIn, Facebook, Instagram, TikTok, ...
- Content, but also metadata
- (Used to?) provide some data access
  - $\rightarrow\,$  currently in flux

#### • Facebook Data

 $ightarrow\,$  large community, representative across income distribution

- ightarrow not accessible to users, not representative across age groups
- Twitter data

ightarrow less large community, less representative across income distribution

- Facebook Data
  - $\rightarrow\,$  large community, representative across income distribution
  - $\rightarrow$  not accessible to users, not representative across age groups
- Twitter data
  - ightarrow less large community, less representative across income distribution

- Facebook Data
  - $\rightarrow\,$  large community, representative across income distribution
  - $\rightarrow$  not accessible to users, not representative across age groups
- Twitter data

ightarrow less large community, less representative across income distribution

- Facebook Data
  - $\rightarrow\,$  large community, representative across income distribution
  - ightarrow not accessible to users, not representative across age groups
- Twitter data
  - ightarrow less large community, less representative across income distribution

- Facebook Data
  - ightarrow large community, representative across income distribution
  - ightarrow not accessible to users, not representative across age groups
- Twitter data
  - ightarrow less large community, less representative across income distribution
  - $ightarrow\,$  freely accessible, rich data

- Facebook Data
  - ightarrow large community, representative across income distribution
  - ightarrow not accessible to users, not representative across age groups
- Twitter data
  - ightarrow less large community, less representative across income distribution
  - $\rightarrow$  freely accessible, rich data

## FACEBOOK DATA



# Social Connectedness: Measurement, Determinants, and Effects

Michael Bailey, Rachel Cao, Theresa Kuchler, Johannes Stroebel, and Arlene Wong

S ocial networks can shape many aspects of social and economic activity: migration and trade, job-seeking, innovation, consumer preferences and sentiment, public health, social mobility, and more. In turn, social networks themselves are associated with geographic proximity, historical ties, political boundaries, and other factors. Traditionally, the unavailability of large-scale and representative data on social connectedness between individuals or geographic regions has posed a challenge for empirical research on social networks. More recently, a body of such research has begun to emerge using data on social connectedness from online social networking services such as Facebook, LinkedIn, and Twitter. To date, most

#### In a nutshell

- Strength of connectedness between two geographic areas as represented by Facebook friendship ties
- Access data thanks to Micheal Bailey (Facebook)
- Validate their Social Connectedness Index (SCI):
  - SCI and geographic distance
  - concentration of social network and socio-economic charcteristics
  - social connectedness and bilateral economic ties (trade, innovation)
  - social connectedness and bilateral social activity (migration)
- SCI is openly available (upon request)

#### Social Conncectedness Index

- 1. Assign people to geographic areas
- 2. Calculate connectedness

$$SCI_{ij} = \frac{n_{ij}}{n_i \times n_j}$$
 (1)

where  $n_{ij}$  are the number of users in country i that are friends with j (friendship is symmetric in FB!),  $n_i$  FB users in i and  $n_j$  users in j

- 3. Drop small counts and add noise: remove all locations with a low number of observations and add random noise to the number of friendships between each set of locations to ensure no one can be re-identified.
- 4. Final sampling: The final SCI is the average scale of friendship ties across 10 random draws from 99% of active Facebook users to further protect privacy.





	Dependent Variable: Log(SCI)					
	(1)	(2)	(3)	(4)	(5)	
log(Distance in Miles)	$-1.483^{***}$ (0.065)	-1.287*** (0.061)	-1.160*** (0.059)	$-1.988^{***}$ (0.043)	$-1.214^{***}$ (0.055)	
Same State		$1.496^{***}$ (0.087)	$1.271^{***}$ (0.083)	$1.216^{***}$ (0.044)	$1.496^{***}$ (0.085)	
$\Delta$ Income (\$1,000)					$-0.006^{***}$ (0.001)	
$\Delta$ Share Population White (%)					$-0.012^{***}$ (0.001)	
$\Delta$ Share Population No High School (%)					$-0.012^{***}$ (0.002)	
$\Delta$ 2008 Obama Vote Share (%)					$-0.006^{***}$ (0.001)	
$\Delta$ Share Population Religious (%)					$-0.002^{***}$ (0.001)	
County Fixed Effects	Y	Y	Y	Y	Y	
Sample			>200 miles	<200 miles		
Number of observations $R^2$	2,961,968 0.907	2,961,968 0.916	2,775,244 0.916	$186,669 \\ 0.941$	2,961,968 0.922	

Note: Table shows results from a regression of the log of the Social Connectedness Index on a number of explanatory variables. The log of the geographic distance between the counties is the explanatory variable in column 1. In column 2, we include an additional control indicating whether both counties are within the same state. In columns 3 and 4, we restrict the sample to county-pairs that are more and less than 200 miles apart, respectively. The unit of observation is a county-pair. Standard errors are given in parentheses. The online Appendix (http://e-jep.org) provides more details on the data sources and exact specifications.

\*, \*\*, and \*\*\* indicate significance levels of p < 0.1, p < 0.05, and p < 0.01, respectively.

#### Network Concentration and County-Level Characteristics



# Table 3 Social Connectedness and Across-Region Economic Interactions

	(1)	(2)	(3)	(4)				
Panel A: Dependent Variable: log(State-Level Trade Flows)								
log(Distance)	$-1.057^{***}$ $(0.071)$		$-0.531^{***}$ (0.084)	$-0.533^{***}$ (0.085)				
log(SCI)		0.999 * * * (0.051)	$0.643^{***}$ (0.071)	$0.637^{***}$ (0.060)				
State Fixed Effects	Y	Y	Y	Y				
Other State Differences	N	N	N	Ŷ				
Observations $R^2$	2,219 0.912	2,220 0.918	2,219 0.926	2,219 0.930				
20	0.512	0.510	0.040	0.550				

Panel B: Dependent Variable: Indicator for Patent Citation							
log(Distance)	$-0.048^{***}$ (0.002)		-0.011 ** (0.005)	-0.021** (0.009)			
log(SCI)		0.063*** (0.003)	0.049*** (0.006)	0.066*** (0.012)			
Technological Category + County Fixed Effects	Y	Y	Υ	Y			
Cited + Issued Patent Fixed Effects, Other County Differences	Ν	Ν	Ν	Y			
Observations $R^2$	$2,171,754 \\ 0.056$	$2,171,754 \\ 0.059$	$2,171,754 \\ 0.059$	$2,168,285 \\ 0.101$			

Panel C: Dependent Variable: log(County-Level Migration)							
log(Distance)	$-0.973^{***}$ (0.048)		0.023 (0.021)	0.031 (0.021)			
log(SCI)		$1.134^{***}$ (0.019)	$1.148^{***}$ (0.024)	$1.159^{***}$ (0.024)			
County Fixed Effects	Y	Y	Y	Y			
Other County Differences	Ν	Ν	Ν	Y			
Observations	25,305	25,305	25,305	25,287			
$R^2$	0.610	0.893	0.893	0.893			

## Food for thought

- What could one do with SCI data?
- You can access the data at the link https://data.humdata.org/dataset/social-connectedness-index

## Food for thought

- What could one do with SCI data?
- You can access the data at the link https://data.humdata.org/dataset/social-connectedness-index



Contents lists available at ScienceDirect

#### Journal of International Economics

journal homepage: www.elsevier.com/locate/jie



#### International trade and social connectedness

Michael Bailey <sup>a</sup>, Abhinav Gupta <sup>b</sup>, Sebastian Hillenbrand <sup>b</sup>, Theresa Kuchler <sup>b</sup>, Robert Richmond <sup>b,\*</sup>, Johannes Stroebel <sup>b</sup>



<sup>a</sup> Facebook, Inc, United States of America
<sup>b</sup> Stern School of Business, New York University, United States of America

#### ARTICLE INFO

Article history: Received 11 September 2020 Received in revised form 23 December 2020 Accepted 23 December 2020 Available online 29 December 2020

Dataset link: https://data.mendeley.com/ datasets/7wddm84w9r/1

#### JEL codes: F1 F5 F6

#### ABSTRACT

We use de-identified data from Facebook to construct a new and publicly available measure of the pairwise social connectedness between 170 countries and 332 European regions. We find that two countries trade more when they are more socially connected, especially for goods where information frictions may be large. The social connections that predict trade in specific products are those between the regions where the product is produced in the exporting country and the regions where it is used in the importing country. Once we control for social connectedness, the estimated effects of geographic distance and country borders on trade decline substantially.

© 2020 Elsevier B.V. All rights reserved.

#### Table 2

Gravity Regressions - Goods Trade Heterogeneity in 2017.

	Dependent variable: Product-Specific Exports					
	(1)	(2)	(3)	(4)	(5)	
log(SCI)	0.275*** (0.027)	0.299*** (0.028)	0.304*** (0.024)	0.281*** (0.031)	0.287*** (0.025)	
$\log(SCI) \times Share Exchange-Traded$		-0.179** (0.080)	-0.148** (0.070)			
$\log(SCI) \times Rule$ of Law Destination				-0.014 (0.021)	-0.010 (0.019)	
$\log(SCI) \times Rule of Law Origin$				0.000 (0.019)	0.005 (0.015)	
Origin Country × Product FE	Y	Y	Y	Y	Y	
Destination Country × Product FE	Y	Y	Y	Y	Y	
Other Gravity Controls	Y	Y	Y	Y	Y	
$log(Distance) \times Product FE$	Y	Y		Y		
Distance Group × Product FE			Y		Y	
$R^2$	0.932	0.933	0.946	0.932	0.946	
N	2,597,760	2,597,760	2,597,760	2,597,760	2,597,760	
N - Explained by FE	334,186	334,186	334,186	405,093	405,093	

Note: Table shows results from regression 3. The dependent variable is exports of product category k from country i to country j in 2017. Product-level trade data are aggregated up to the first 2 digits of the HS96 product classification. Other gravity controls include a common border dummy, a common official language dummy, a dummy indicating whether the two countries had a common colonizer post-1945, and a dummy indicating whether the pair of countries was in a colonial relationship post-1945. We also separately control for the logarithm of distance interacted with product categories in columns 1, 2, 4 and for distance groups (dummise) and existibution j interacted with product categories in columns 5 and 5. Share Exchange-Traded refers to the proportion of exchange-traded products—based on the conservative classification scheme in Rauch (1999)—within a product category. Rule of law is obtained from the World Governance Indicators published by the World Bank. All specifications include fixed effects for the importer country interacted with product categories. Standard errors are clustered by exporter and importer country. The data include 165 countries and 96 product categories. (which amounts to 2,597,760 observations. Observations that are fully explained by the fixed effects are dropped before the PPML estimation. Significance levels: (p<0.01), \*\*(p<0.05).

## TWITTER DATA

- Twitter Streaming API: 1 % random sample of all tweets
  - ightarrow filters: keyword, geolocation
  - ightarrow between 40 and 60 per second
- 42 variables: text, username, user\_lang, lang, followers, timezone, latitude, longitude, place, source,...

- Twitter Streaming API: 1 % random sample of all tweets
  - $\rightarrow\,$  filters: keyword, geolocation
  - ightarrow between 40 and 60 per second
- 42 variables: text, username, user\_lang, lang, followers, timezone, latitude, longitude, place, source,...

- Twitter Streaming API: 1 % random sample of all tweets
  - $\rightarrow\,$  filters: keyword, geolocation
  - $\rightarrow\,$  between 40 and 60 per second
- 42 variables: text, username, user\_lang, lang, followers, timezone, latitude, longitude, place, source,...

- Twitter Streaming API: 1 % random sample of all tweets
  - $\rightarrow\,$  filters: keyword, geolocation
  - $\rightarrow~$  between 40 and 60 per second
- 42 variables: text, username, user\_lang, lang, followers, timezone, latitude, longitude, place, source,...




```
"created_at": "Tue Apr 18 15:22:19 +0000 2017",
"id": 854354410041991168
"id str": "854354410041991168".
"text": "@ichmandasnicht offenbar nicht Mathematik .....
"source": "<a href=\"http://taphots.com/tweethot\" rel=\"nofollow\">Tweethot for iOS</a>".
"in_reply_to_status_id": 854247992186073088,
"in reply to status id str": "854247992186073088".
"in reply to screen name": "ichmagdasnicht".
 "id": 19030252.
 "name": "Timo Zander".
 "screen_name": "tinkengil",
 "url": "http://about.me/timozander".
  "description": "PhD-Student | Podcastet bei plavtogether-podcast.de | bloggt gelegentlich bei insulinaspekte.de und http://tinkengil.com | http://instagram.com/tinkengil".
 "friends count": 344.
 "favourites count": 1830.
  "created at": "Thu Jan 15 17:40:27 +0000 2009".
  "utc_offset": 7200.
  "time_zone": "Bern".
  "geo_enabled": true.
  "profile_background_color": "EBEBEB",
  "profile background image url": "http://pbs.twimg.com/profile background images/590786545/5vvvvdxrk528xhz91w86.ipeg".
  "profile background image url https": "https://pbs.twimg.com/profile background images/590786545/5vvvvdxrk528xhz91w86.ipeg".
  "profile link color": "990000".
 "profile use background image": false.
  "profile_image_url": "http://pbs.twimg.com/profile_images/549318880876048384/zag6999H_normal.jpeg",
  "profile image url https": "https://pbs.twimq.com/profile images/549318880876048384/zag6999H normal.jpeg",
```

```
"profile background image url": "http://pbs.twimg.com/profile background images/590786545/5vvvvdxrk528xhz91w86.ipeg".
  "profile background image url https": "https://pbs.twimg.com/profile background images/590786545/5vyvydxrk528xhz91w86.jpeg",
  "profile link color": "990000".
 "profile_image_url": "http://pbs.twimg.com/profile_images/549318880876048384/zag6999H_normal.jpeg",
 "profile_image_url_https": "https://pbs.twimq.com/profile_images/549318880876048384/zag6999H_normal.jpeg",
 "default profile image": false.
"aeo": {
 "type": "Point".
   54.32436928,
    10.12301066
"coordinates": {
   10.12301066.
    54.32436928
 "id": "1b9b5e83e647a7ed".
 "upl": "https://api.twitter.com/1.1/geo/id/1b9b5e83e647a7ed.ison".
 "country": "Germany".
    "type": "Polygon".
         10.032937.
         54.250693
         10.032937.
         54,432916
```

```
54.250693
         10.032937.
         54,432916
         10.218568
          54,432916
          54,250693
     "screen name": "ichmagdasnicht".
     "name": "Marvin || Runaways".
"timestamp_ms": "1492528939148"
```

• Obvious: Text-mining

- Not so obvious: Metadata
  - $\rightarrow$  Language distribution
  - ightarrow Migration

• Obvious: Text-mining

- Not so obvious: Metadata
  - ightarrow Language distribution
  - ightarrow Migration

• Obvious: Text-mining

- Not so obvious: Metadata
  - $\rightarrow$  Language distribution
  - ightarrow Migration

• Obvious: Text-mining

- Not so obvious: Metadata
  - $\rightarrow$  Language distribution
  - $\rightarrow$  Migration

### • Spatial distribution of languages in Europe

- Geolocation from "coordinates", and "user\_lang" or "lang"
  - ightarrow large heterogeneity across and within countries
- Coordinates provided either by the user's device's GPS coordinates, or a self-assigned location
  - ightarrow Barratt, J. Cheshire, and E. Manley (2013) use similar data for NY boroughs

- Spatial distribution of languages in Europe
- Geolocation from "coordinates", and "user\_lang" or "lang"
  - ightarrow large heterogeneity across and within countries
- Coordinates provided either by the user's device's GPS coordinates, or a self-assigned location
  - ightarrow Barratt, J. Cheshire, and E. Manley (2013) use similar data for NY boroughs

- Spatial distribution of languages in Europe
- Geolocation from "coordinates", and "user\_lang" or "lang"
  - $\rightarrow\,$  large heterogeneity across and within countries
- Coordinates provided either by the user's device's GPS coordinates, or a self-assigned location
  - ightarrow Barratt, J. Cheshire, and E. Manley (2013) use similar data for NY boroughs

- Spatial distribution of languages in Europe
- Geolocation from "coordinates", and "user\_lang" or "lang"
  - ightarrow large heterogeneity across and within countries
- Coordinates provided either by the user's device's GPS coordinates, or a self-assigned location
  - ightarrow Barratt, J. Cheshire, and E. Manley (2013) use similar data for NY boroughs

- Spatial distribution of languages in Europe
- Geolocation from "coordinates", and "user\_lang" or "lang"
  - $\rightarrow\,$  large heterogeneity across and within countries
- Coordinates provided either by the user's device's GPS coordinates, or a self-assigned location
  - ightarrow Barratt, J. Cheshire, and E. Manley (2013) use similar data for NY boroughs























- Bots: an issue, Chu et al. (2012) suggest only taking those sent from smart phones and official app
- 6.6 million unique human Twitter users
- 481,720 unique human Twitter users in Europe
- 73 different languages
- 25 % tweet in more than 1 language, in Germany 31 %
- 958,071 unique language-user observations

- Bots: an issue, Chu et al. (2012) suggest only taking those sent from smart phones and official app
- 6.6 million unique human Twitter users
- 481,720 unique human Twitter users in Europe
- 73 different languages
- 25 % tweet in more than 1 language, in Germany 31 %
- 958,071 unique language-user observations

- Bots: an issue, Chu et al. (2012) suggest only taking those sent from smart phones and official app
- 6.6 million unique human Twitter users
- 481,720 unique human Twitter users in Europe
- 73 different languages
- 25 % tweet in more than 1 language, in Germany 31 %
- 958,071 unique language-user observations

- Bots: an issue, Chu et al. (2012) suggest only taking those sent from smart phones and official app
- 6.6 million unique human Twitter users
- 481,720 unique human Twitter users in Europe
- 73 different languages
- 25 % tweet in more than 1 language, in Germany 31 %
- 958,071 unique language-user observations

- Bots: an issue, Chu et al. (2012) suggest only taking those sent from smart phones and official app
- 6.6 million unique human Twitter users
- 481,720 unique human Twitter users in Europe
- 73 different languages
- 25 % tweet in more than 1 language, in Germany 31 %
- 958,071 unique language-user observations

- Bots: an issue, Chu et al. (2012) suggest only taking those sent from smart phones and official app
- 6.6 million unique human Twitter users
- 481,720 unique human Twitter users in Europe
- 73 different languages
- 25 % tweet in more than 1 language, in Germany 31 %
- 958,071 unique language-user observations

# Twitter and UK Census Population



#### Twitter and UK Census Main Language



Language use on Twitter and UK census, correlation = 0.49.

#### Twitter and Eurobarometer



#### Language use on Twitter and Eurobarometer, correlation = 0.74.



• Economic crisis in Venezuela: Large (?) number of refugees  $\rightarrow$  lack of official numbers

- Dataset of geolocalized Tweets of people that tweeted from Venezuela between February 2017 and May 2018
   → 5.4 million tweets
   → 490.000 tweets from 30.000 human Twitter users
- Idea: What location(s) do they tweet from over time?

• Economic crisis in Venezuela: Large (?) number of refugees  $\rightarrow$  lack of official numbers

- Dataset of geolocalized Tweets of people that tweeted from Venezuela between February 2017 and May 2018
   → 5.4 million tweets
   → 490.000 tweets from 30.000 human Twitter users
- Idea: What location(s) do they tweet from over time?

- Economic crisis in Venezuela: Large (?) number of refugees  $\rightarrow$  lack of official numbers
- Dataset of geolocalized Tweets of people that tweeted from Venezuela between February 2017 and May 2018
  - ightarrow 5.4 million tweets
  - ightarrow 490.000 tweets from 30.000 human Twitter users
- Idea: What location(s) do they tweet from over time?

- Economic crisis in Venezuela: Large (?) number of refugees  $\rightarrow$  lack of official numbers
- Dataset of geolocalized Tweets of people that tweeted from Venezuela between February 2017 and May 2018  $\rightarrow$  5.4 million tweets
  - ightarrow 490.000 tweets from 30.000 human Twitter users
- Idea: What location(s) do they tweet from over time?

- Economic crisis in Venezuela: Large (?) number of refugees  $\rightarrow$  lack of official numbers
- Dataset of geolocalized Tweets of people that tweeted from Venezuela between February 2017 and May 2018
  - $\rightarrow$  5.4 million tweets
  - ightarrow 490.000 tweets from 30.000 human Twitter users
- Idea: What location(s) do they tweet from over time?

- Economic crisis in Venezuela: Large (?) number of refugees  $\rightarrow$  lack of official numbers
- Dataset of geolocalized Tweets of people that tweeted from Venezuela between February 2017 and May 2018
  - $\rightarrow$  5.4 million tweets
  - ightarrow 490.000 tweets from 30.000 human Twitter users
- Idea: What location(s) do they tweet from over time?




## Distribution of countries



#### Distribution of countries of last recorded locations of users outside Venezuela

#### • Hawelka (2014): global mobility patterns, tourism flows

#### • Jurdak (2015) city-to-city travel in Australia

- Morstatter (2013): random sample creates an accurate picture of the entire population of geolocated Tweets
- Question: How representative are geolocalized tweets?

- Hawelka (2014): global mobility patterns, tourism flows
- Jurdak (2015) city-to-city travel in Australia
- Morstatter (2013): random sample creates an accurate picture of the entire population of geolocated Tweets
- Question: How representative are geolocalized tweets?

- Hawelka (2014): global mobility patterns, tourism flows
- Jurdak (2015) city-to-city travel in Australia
- Morstatter (2013): random sample creates an accurate picture of the entire population of geolocated Tweets
- Question: How representative are geolocalized tweets?

- Hawelka (2014): global mobility patterns, tourism flows
- Jurdak (2015) city-to-city travel in Australia
- Morstatter (2013): random sample creates an accurate picture of the entire population of geolocated Tweets
- Question: How representative are geolocalized tweets?

## Population and Tweets



"Gridded Population of the World" and number of Tweets by location

### Population and Users



#### "Gridded Population of the World" and number of Twitter users by location

## Representativeness of Twitter users in Venezuela

- "Digital in 2017 Global Overview report": 44% of Venezuelans social media, 35% from mobile device
- "Tendencias Digitales": 56% of internet users in Venezuela use Twitter or comparable social media services
- Twitter: penetration in Venezuela 26 %

## Representativeness of Twitter users in Venezuela

- "Digital in 2017 Global Overview report": 44% of Venezuelans social media, 35% from mobile device
- "Tendencias Digitales": 56% of internet users in Venezuela use Twitter or comparable social media services
- Twitter: penetration in Venezuela 26 %

## Representativeness of Twitter users in Venezuela

- "Digital in 2017 Global Overview report": 44% of Venezuelans social media, 35% from mobile device
- "Tendencias Digitales": 56% of internet users in Venezuela use Twitter or comparable social media services
- Twitter: penetration in Venezuela 26 %

#### Tweets per users



Number of tweets per user in the dataset

### Days per users



Number of days a user is observed in the dataset

#### • narrow sample to users who

 $\rightarrow$  tweeted from Venezuela exclusively between Feb and May '17 (Period 1)  $\rightarrow$  tweeted from *a country* exclusively between Feb and May '18 (Period 2)

Everyone who is not in Venezuela in period 2: migrant

reduces sample to 818 (!)

 $\rightarrow$  Problem: Large heterogeneity in tweet frequency

- narrow sample to users who
  - ightarrow tweeted from Venezuela exclusively between Feb and May '17 (Period 1)
  - $\rightarrow$  tweeted from *a country* exclusively between Feb and May '18 (Period 2)
- Everyone who is not in Venezuela in period 2: migrant
- reduces sample to 818 (!)
  - $\rightarrow$  Problem: Large heterogeneity in tweet frequency

- narrow sample to users who
  - ightarrow tweeted from Venezuela exclusively between Feb and May '17 (Period 1)
  - $\rightarrow$  tweeted from *a country* exclusively between Feb and May '18 (Period 2)
- Everyone who is *not* in Venezuela in period 2: migrant
- reduces sample to 818 (!)
  - $\rightarrow$  Problem: Large heterogeneity in tweet frequency

- narrow sample to users who
  - ightarrow tweeted from Venezuela exclusively between Feb and May '17 (Period 1)
  - $\rightarrow$  tweeted from *a country* exclusively between Feb and May '18 (Period 2)
- Everyone who is not in Venezuela in period 2: migrant
- reduces sample to 818 (!)
  - $\rightarrow$  Problem: Large heterogeneity in tweet frequency

- narrow sample to users who
  - ightarrow tweeted from Venezuela exclusively between Feb and May '17 (Period 1)
  - $\rightarrow$  tweeted from *a country* exclusively between Feb and May '18 (Period 2)
- Everyone who is not in Venezuela in period 2: migrant
- reduces sample to 818 (!)

 $\rightarrow$  Problem: Large heterogeneity in tweet frequency

- narrow sample to users who
  - ightarrow tweeted from Venezuela exclusively between Feb and May '17 (Period 1)
  - $\rightarrow$  tweeted from *a country* exclusively between Feb and May '18 (Period 2)
- Everyone who is not in Venezuela in period 2: migrant
- reduces sample to 818 (!)
  - $\rightarrow$  Problem: Large heterogeneity in tweet frequency



#### Tweets by migrants and non-migrants in two periods

### • Need weight to correct for sampling bias

 Suppose probability of individual i tweeting exactly x tweets in three-month period given by

$$p_{i,x} = Pr\{tw_i = x\}$$

•  $tw_i$  random variable denoting tweets  $i \rightarrow$  assume this probability distribution constant across periods

- Need weight to correct for sampling bias
- Suppose probability of individual *i* tweeting exactly *x* tweets in three-month period given by

$$p_{i,x} = \Pr\{tw_i = x\}$$

•  $tw_i$  random variable denoting tweets  $i \rightarrow$  assume this probability distribution constant across periods

- Need weight to correct for sampling bias
- Suppose probability of individual *i* tweeting exactly *x* tweets in three-month period given by

$$p_{i,x} = Pr\{tw_i = x\}$$

•  $tw_i$  random variable denoting tweets  $i \rightarrow$  assume this probability distribution constant across periods

- Twitter provides s = 0.01 of all tweets, independent of user  $\rightarrow q = (1 - s) = 99\%$  of Tweets not reported
- Denote  $U^1$  ( $U^2$ ) set all users observed at least once in period 1 (2)

- Twitter provides s = 0.01 of all tweets, independent of user  $\rightarrow q = (1 s) = 99\%$  of Tweets not reported
- Denote  $U^1$  ( $U^2$ ) set all users observed at least once in period 1 (2)

- Twitter provides s = 0.01 of all tweets, independent of user  $\rightarrow q = (1 s) = 99\%$  of Tweets not reported
- Denote  $U^1$  ( $U^2$ ) set all users observed at least once in period 1 (2)

• Probability of observing an individual who tweeted  $x_i$  times in period 1

$$Pr\{i \in U^1 | tw_i^1 = x\} = 1 - q^x.$$

• Probability of observing same individual who tweeted  $y_i$  times in period 2

$$Pr\{i \in U^2 | tw_i^2 = y\} = 1 - q^y.$$

• Probability of observing an individual who tweeted  $x_i$  times in period 1

$$Pr\{i \in U^1 | tw_i^1 = x\} = 1 - q^x.$$

• Probability of observing same individual who tweeted  $y_i$  times in period 2

$$Pr\{i \in U^2 | tw_i^2 = y\} = 1 - q^y.$$

 Assuming independence between the two sample, probability to be observed in both periods

$$\begin{split} Pr\{i \in U^1 \text{ and } i \in U^2\} &= \sum_{x=0}^{\infty} \sum_{y=0}^{\infty} Pr\{i \in U^1 | tw_i^1 = x\} Pr\{tw_i^1 = x\} \times \\ Pr\{i \in U^2 | tw_i^2 = y\} Pr\{tw_i^2 = y\} \\ &= \sum_{x=0}^{\infty} p_{i,x}(1-q^x) \sum_{y=0}^{\infty} p_{i,y}(1-q^y) \\ &= (1-E_i[q^x])^2 = (1-G_i(q))^2 \end{split}$$

•  $G_i(q)$  probability generating function

- Model the individuals' tweeting behavior as a Poisson process
- Assume each individual has Poisson tweet rate in a three month period  $\lambda_i$
- With Poisson distribution, rewrite the probability generating function as

$$G_i(q) = e^{-\lambda_i(1-q)} = e^{-\lambda_i s}.$$

- Model the individuals' tweeting behavior as a Poisson process
- Assume each individual has Poisson tweet rate in a three month period  $\lambda_i$
- With Poisson distribution, rewrite the probability generating function as

$$G_i(q) = e^{-\lambda_i(1-q)} = e^{-\lambda_i s}.$$

- Model the individuals' tweeting behavior as a Poisson process
- Assume each individual has Poisson tweet rate in a three month period  $\lambda_i$
- With Poisson distribution, rewrite the probability generating function as

$$G_i(q) = e^{-\lambda_i(1-q)} = e^{-\lambda_i s}.$$

• Hence probability of being observed in both periods

$$Pr\{i\in U^0 ext{ and } i\in U^1\}=(1-e^{-\lambda_i s})^2$$

with s = 0.01 in our case.

(2)

## Net outflow over time

		(1)	(2)	(3)	(4)	(5)	(6)	(7)
		Venezuela	Colombia	Argentina	Brazil	Germany	Venezuela	Colombia
Emigration	unweighted	6,76%	7,78%	7,62%	3,88%	11,59%	6,99%	6,06%
	weighted	9,59%	7,84%	7,92%	3,97%	13,18%	7,98%	6,10%
Immigration	unweighted	2,01%	5,21%	10,48%	3,59%	11,27%	1,77%	5,21%
	weighted	2,22%	5,48%	10,70%	3,67%	12,41%	1,70%	5,37%
Difference	unweighted	-4,75%	-2,57%	2,86%	-0,29%	-0,32%	-5,22%	-0,85%
	weighted	-7,37%	-2,36%	2,78%	-0,30%	-0,77%	-6,28%	-0,73%
Annualized weighted perc.		-9,7%	-3,1%	3,7%	-0,4%	-1%	-12,1%	-1,4%
Period 1		02-04/17	02-04/17	02-04/17	02-04/17	02-04/17	12/16-04/17	12/16-04/17
Period 2		02-04/18	02-04/18	02-04/18	02-04/18	02-04/18	12/17-04/18	12/17-04/18

Source: Authors' calculations.

Computed emigration and immigration numbers

## Distribution of countries



Distribution of countries of users between February and April '18

#### Conclusion

- Social media data allows researchers to observe people, revealed preferences
- Design of exercise important: Endogeneity, sampling, ...
## Social Media Data

DSIER [/dɪˈzaɪər/]

Julian Hinz

**Bielefeld University**