

Big data and migration research

Using social media data for migration research

Jisu Kim, Ph.D



MAX PLANCK INSTITUTE
FOR DEMOGRAPHIC
RESEARCH

Outline

- ▶ Limitations of traditional data and social media data as an alternative data
- ▶ Twitter data format
- ▶ Twitter and migration research
- ▶ Ethics and privacy

Study of migration

Immigrants touch upon multidimensional aspects of both the host country and the home country.

- ▶ Economics:
 - ▶ Jobs
 - ▶ Unemployment
 - ▶ Society:
 - ▶ Integration
 - ▶ Friends
 - ▶ Well-being
 - ▶ Population density
 - ▶ Fertility
 - ▶ Culture:
 - ▶ Food
 - ▶ Language
 - ▶ Inter-marriage
 - ▶ more...
 - ▶ Politics
- and many more...

Immigration

Trump put up walls to immigrants, with stinging rhetoric and barriers made of steel and regulation

The screenshot shows a news article from a German political website. At the top, there's a dark banner with the word "FINISH THE WALL" in large white letters. Below it, the main headline reads "U.S. and Haiti work to address migration challenges". A sub-headline below the main one says "Home | News & Events | U.S. and Haiti work to address migration challenges". At the bottom of the page, there's a navigation bar with links like "GERMAN GENERAL ELECTION", "ELECTIONS", "POLITICS", "REFUGEES IN GERMANY", "ASYLUM", and "MIGRATION POLICY". On the right side, there's a sidebar with a small image of a person and some text. At the very bottom, there's a footer with icons for social media and a page number "3 / 53".

GERMAN GENERAL ELECTION | ELECTIONS | POLITICS | REFUGEES IN GERMANY | ASYLUM | MIGRATION POLICY

3 / 53

Migration research

Who is an **Immigrant**?:

- ▶ “A person who moves to a country other than that of his or her usual residence for a period of at least a year.¹”
- ▶ “Whose movement across borders-whether legal or illegal- is essentially permanent²”
- ▶ “is defined on the ground of the place of birth (foreign-born) or of the citizenship (foreigners)³”

¹United Nation

²World Bank

³OECD

Definitions

Is this your *home* country?

Country of origin:

In the context of migration, the country of origin is the country of nationality where the usual residence was, before migration took place.

Usual residence:

“The geographical place where the enumerated person usually resides” - the concept used in censuses.

Census

Where did this person live one year ago?

Is this person citizen of XX?

Where were you born?

Limitation of traditional data sources

Census, survey, register data

- ▶ Costly
- ▶ Outdated
- ▶ Time consuming
- ▶ Inconsistent
- ▶ Unavailable
- ▶ Lack of data on emigration
- ▶ Incomplete answers/misunderstanding questions etc.
- ▶ Immigrants are often underrepresented traditional data sources.
- ▶ limited in hard-to-reach contexts and societies.

Social media



Big data

is “information assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and process automation”⁴ Laney, 2001. It can be described in 3 Vs which are:

- ▶ Volume: As the name suggests, the size of the data is Big, hence the volume of the data.
- ▶ Velocity: Big data such as Twitter allow us to stream data at real-time. The rate at which we obtain data is faster than the traditional data sources.
- ▶ Variety: Traditional data are mostly structured data. Big data, on the other hand, come in various forms. It can be videos, photos, texts, and audios. It requires a thorough data processing before extracting information/knowledge from it.

⁴<https://www.gartner.com/en/information-technology/glossary/big-data>

Alternative data sources: Social media and Big data

Twitter, Facebook, Yahoo, ...

- ▶ Free
- ▶ Granular data
- ▶ Large scale data
- ▶ Continuously generated
- ▶ Information/opinion shared by users from an uncontrolled environment
- ▶ various forms of data: video, image, text, audio etc.

Required information

- Geo-location
- Information on time

So how do we get this Big data?

CHRIS STOKEWALKER

BUSINESS 28.10.2022 02:28 PM

Elon Musk's Twitter Will Be Chaos

The entrepreneur's laundry list of ideas includes scrapping content moderation, charging subscription fees, and even branching out beyond social media.



PHOTOGRAPH: CARINA JOHANSEN/GETTY IMAGES

5

⁵source:

<https://www.wired.co.uk/article/elon-musk-twitter-deal-chaos>



13 / 53

What is he going to do???



Elon Musk ✅ @elonmusk · 17h

Please note that Twitter will do lots of dumb things in coming months.

...

We will keep what works & change what doesn't.

34.9K

44.3K

435.4K



Elon Musk ✅ @elonmusk · Nov 7

Any name change at all will cause temporary loss of verified checkmark

...

12K

16.9K

142.1K



Elon Musk ✅ @elonmusk · Nov 5

Twitter will soon add ability to attach long-form text to tweets, ending absurdity of notepad screenshots

...

35.4K

67.3K

610.9K



Elon Musk ✅ @elonmusk · Nov 5

Followed by creator monetization for all forms of content

...

9,032

17.3K

220K



Elon Musk ✅ @elonmusk · Nov 5

Trash me all day, but it'll cost \$8

...

112.1K

126K

1.3M



Academic Research track⁷

- ▶ Access to full-archive search and tweet counts
- ▶ Monthly-tweet cap usage: “Certain endpoints (like filtered stream and recent search) have a limit on how many Tweets they can pull per month.” == 10 million tweets⁶
- ▶ Query limit up to 1024 characters

⁶<https://developer.twitter.com/en/docs/twitter-api/rate-limits>

⁷https://blog.twitter.com/developer/en_us/topics/tools/2020/introducing_new_twitter_api

Current API policy⁸

Free) Only "Manage Tweets" and "Users lookup"

Tweets=1,500 tweets per month

Costs=Free

Basic)

Tweets=Retrieve up to 10K Tweets per month

Cost=\$100.00 USD/month

Pro)

Tweets=Retrieve up to 1M Tweets per month Apps

Cost=\$5000.00 USD/month

⁸<https://developer.twitter.com/en/docs/twitter-api>

Don't cry...

- ▶ archive.com
- ▶ [https://www.trackmyhashtag.com/blog/
free-twitter-datasets/](https://www.trackmyhashtag.com/blog/free-twitter-datasets/)
- ▶ <https://data.world/datasets/twitter>
- ▶ ...

Twitter data format

- ▶ Tweet object
- ▶ Entity object
- ▶ User object

Tweet object

“id”: "1050118621198921728"

“text”: "Apply now for our week-long open online course on Digital and Computational Demography hosted by @MPIDRnews scientists ..."

“Context annotations”:

‘name’: ‘Interests and Hobbies Category’, **‘description’**: ‘A grouping of interests and hobbies entities, like Novelty Food or Destinations’, **‘entity’**: ‘id’: '852291840472629248’, ‘name’: ‘Online education’, ‘description’: ‘Online education’

“created_at”: "202x-09-xxTxx:xx:xx.000z"

“conversation_id”: "1435336531519197188"

“lang”: “en”

“geo”: e.g., “coordinates”: [-73.999xx, 40.7416xxx]

“place_id”: "01a9a39529b27f36"

Entity object

```
"hashtags": ["tag": "BuildWhatsNext"],  
"mentions": "tag": "@TwitterDev...",  
"url": "https://t.co/z5RhIVxJFK",
```



Twitter API @TwitterAPI · Feb 9

We have added the sort_order parameter to the search endpoints in the Twitter API v2 which gives developers the option of returning Tweets based on recency or relevancy. Check out the details here 🌟



twittercommunity.com

Introducing the sort_order parameter for search en...
Today, we're sharing a small, but important enhancement to the search functionality in the ...

16

20

73



```
"expanded_url": "https://twittercommunity.com/t/updates-to-retweets-lookup-and...",  
"description": "Thanks for your feedback on the v2 Retweets and Likes endpoints. We've heard you. Starting today, you can retrieve the complete list of accounts that have Liked or Retweeted a Tweet,..."
```

User object

```
"id": "2244994945",
"name": "Twitter Dev",
"username": "TwitterDev",
"location": "127.0.0.1",
"verified": true,
"protected": false,
```

The image shows a Twitter profile card for the account "Twitter Dev" (@TwitterDev). The profile picture is a blue circle with a white Twitter bird icon. The bio reads: "The voice of the #TwitterDev team and your official source for updates, news, and events, related to the #TwitterAPI.". The location is listed as "127.0.0.1". The account is verified, indicated by a blue checkmark. The follower count is 522.3K, and the user has 2,024 following. The profile was joined in December 2013. There are standard Twitter interaction buttons like "Following", "Unfollow", and a bell icon.

Twitter Dev

@TwitterDev

The voice of the #TwitterDev team and your official source for updates, news, and events, related to the #TwitterAPI.

127.0.0.1 developer.twitter.com/en/community Born March 21

Joined December 2013

2,024 Following 522.3K Followers

```
"description": "The voice of the #TwitterDev team and your official source for updates, news, and events, related to the #TwitterAPI.",
"url": "https://t.co/3ZX3TNiZCY",
"profile_image_url": "https://pbs.twimg.com/profile_images/1267175364003901441/tBZNFAgA_normal.jpg",
"created_at": "2013-12-14T04:35:55.000Z"
```

Twitter data and migration research

Information from Twitter data that can be used in migration research:

- ▶ Geo-tagged tweets
- ▶ Self-declared location information from profile
- ▶ Social networks: followers and friends
- ▶ Name
- ▶ Language

Anything else?

Potential issues

- ▶ Tourists
- ▶ Travel blogs
- ▶ Bots
- ▶ ...

Definitions of migrants on Twitter in the literature

- ▶ “A Twitter user has the nationality that others believe you have.” (Huang et al., 2014)
- ▶ “Any individual leaving Venezuela during the time window of observation.” (Mazzoli et al., 2020)
- ▶ “Anyone who tweeted exclusively from Venezuela in the time period between Feb. 1 and April 30 2017.” (Hausmann et al., 2018)
- ▶ “Migrants are users that are identified as people who moved to a different country for at least one of the 4-month periods.” (Zagheni et al., 2014)
- ▶ “A migrant is a person that has the residence different from the nationality.” (Kim et al., 2020)

Identifying migrants (Mazzoli et al., 2020)

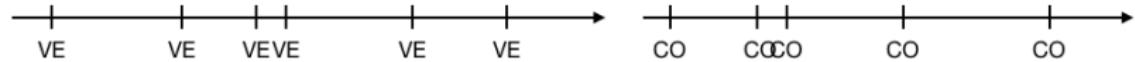


Figure: Mock example of geo-tagged tweet timeline

Inferring nationalities of Twitter users (Huang et al., 2014)



- ▶ Twitter page: xx
- ▶ Homepage: XX
- ▶ Location: Germany
- ▶ Time zone: CET
- ▶ Interface language: English
- ▶ Tweets languages: (1) English: 97% (2) Italian: 2% (3) Korean: 1%
- ▶ Follower locations: (1) US: 151 (2) KR: 72 (3) DE: 3
- ▶ Following locations: (1) US: 80 (2) KR: 40 (3) IT: 30
- ▶ Tweets from (1) IT: 30 (2) KR: 29

Identifying migrants on Twitter (Kim et al., 2020)

"A migrant is a person that has the residence different from the nationality."

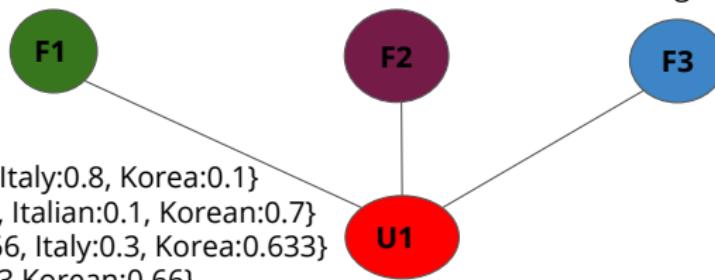
- ▶ Country of residence: "the country with the longest length of stay"
- ▶ Country of nationality: "the ensemble of features that make a person feel like they belong to a certain country"

Identifying migrants on Twitter (Kim et al., 2020)

$\text{loc}^{F1} = \{\text{France}: 0.2, \text{Italy}: 0.8\}$
 $\text{lang}^{F1} = \{\text{Italian}: 1\}$

$\text{loc}^{F2} = \{\text{Italy}: 0.1, \text{Korea}: 0.9\}$
 $\text{lang}^{F2} = \{\text{Korean}: 1\}$

$\text{loc}^{F3} = \{\text{Korea}: 1\}$
 $\text{lang}^{F3} = \{\text{Korean}: 1\}$



$\text{loc}^{U1} = \{\text{France}: 0.1, \text{Italy}: 0.8, \text{Korea}: 0.1\}$
 $\text{lang}^{U1} = \{\text{French}: 0.2, \text{Italian}: 0.1, \text{Korean}: 0.7\}$
 $\text{floc}^{U1} = \{\text{France}: 0.066, \text{Italy}: 0.3, \text{Korea}: 0.633\}$
 $\text{flang}^{U1} = \{\text{Italian}: 0.33, \text{Korean}: 0.66\}$

Validation

Are we correctly identifying migrants?
Are there good ground truth data to compare?

Validation-Gold standard data (Huang et al., 2014)

	Pre.	Rec.	F1
QA	86.67%	95.37%	90.81%
ARA	82.96%	71.16%	76.56%
WES	70.86%	70.62%	70.64%
SA	93.35%	90.48%	91.89%
IN	82.19%	71.13%	76.00%
OTH	78.76%	40.72%	53.54%
UN	30.78%	15.13%	20.16%

Table: The average Precision, Recall, and F1 scores for each nationality group Huang et al., 2014

Validation-Official statistics (Kim et al., 2020)

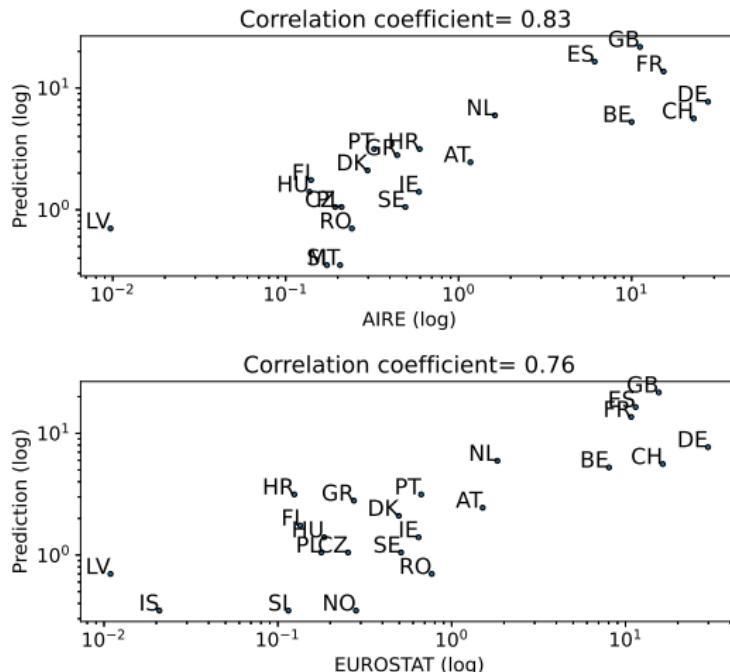


Figure: Correlation between predicted data and official statistics (Kim et al., 2020)

Inferring international and internal migration patterns from Twitter data (Zagheni et al., 2014)

Observe changes in out-migration rate in OECD countries. But the changes could be observed by the changes in behaviours on Twitter, i.e., compositional changes, underline changes in pattern.

- ▶ Correct for overestimation of out-migration rates
 - ▶ Non-representative user set-Twitter
 - ▶ Selection bias that is changing over time
 - ▶ Biases may be different across countries and constant over short period of time.
 - ▶ Assume that the population of Twitter users change in similar ways across all the OECD countries over time.
- Diff-in-diff estimator to evaluate relative changes in trends
(Zagheni et al., 2014)

Difference-in-Differences

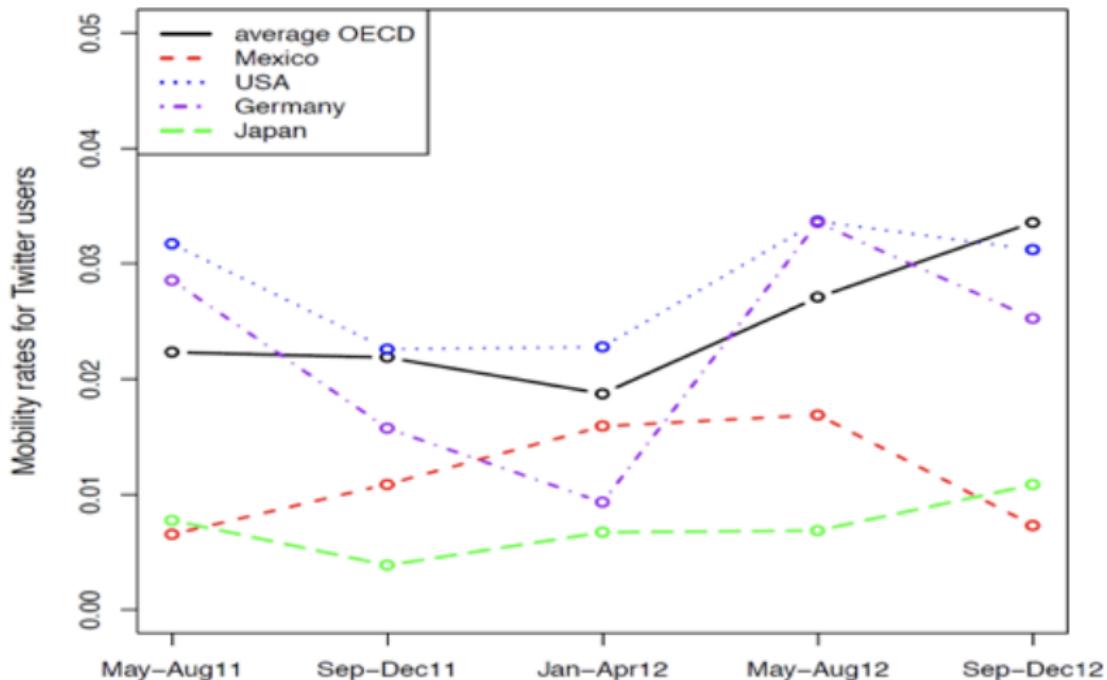


Figure: Mobility rates for Twitter users Zagheni et al., 2014

Cont.

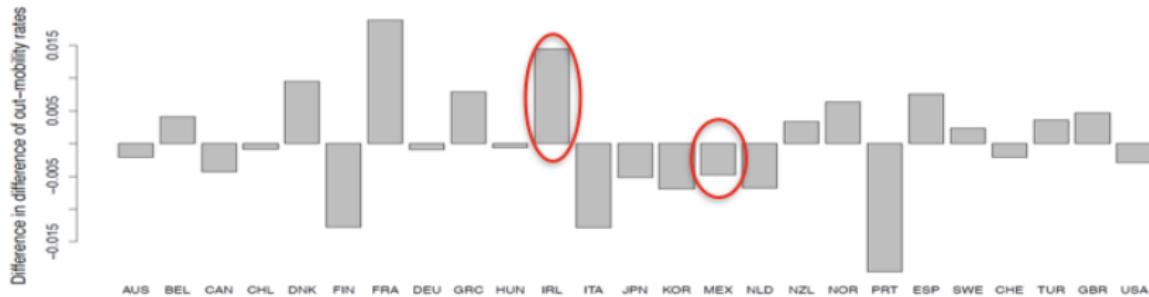


Figure: (Soft) Validation: Ireland out-migration rate grew by 2.2% 2011 – >2012, more than most countries (Irish Central Statistics Office)
Mexico also sees a reduction in out-migration (Pew Research Center)

Other migration related researches

- ▶ Cultural integration (Kim et al., 2021)
- ▶ Spatial integration (Mazzoli et al., 2020; Lamanna et al., 2018)
- ▶ Sentiment analysis (Öztürk et al., 2018; Arcila-Calderón et al., 2021)
- ▶ Social integration (Kim et al., 2023)

Cultural integration (Kim et al., 2021)

Q. *How much do migrants absorb the culture of their destination society? Do they lose connection with their home country?*

Traditional studies focus on elements such as:

- ▶ Language proficiency
- ▶ Marital status
- ▶ Role of media ...

	Low OA	High OA
Low DA	Marginalisation	Separation
High DA	Assimilation	Integration

Table: Theories of integration and their relation to OA and DA Kim et al., 2021

Origin and Destination attachment indexes (Kim et al., 2021)

$$H(h) = \frac{-\sum_c P_h(c) \log P_h(c)}{\log(|P_h(c)|)} \quad (1)$$

$$OA(u) = \frac{\# C_n(u) \text{ hashtags}}{\# \text{ total hashtags}} = \frac{HT(u, C_n(u))}{HT(u)} \quad (2)$$

$$DA(u) = \frac{\# C_r(u) \text{ hashtags}}{\# \text{ total hashtags}} = \frac{HT(u, C_r(u))}{HT(u)} \quad (3)$$

```
'Salvini':  
{'CH': 1,  
'CL': 1,  
'CZ': 2,  
'DE': 2,  
'ES': 3,  
'FR': 2,  
'GB': 1,  
'IT': 544,  
'LU': 1,  
'NL': 2,  
'PL': 1,  
'TH': 2,  
'TR': 1,  
'US': 6}  
/569
```

Entropy score=0.11
∴ Italian specific

OA & DA indexes

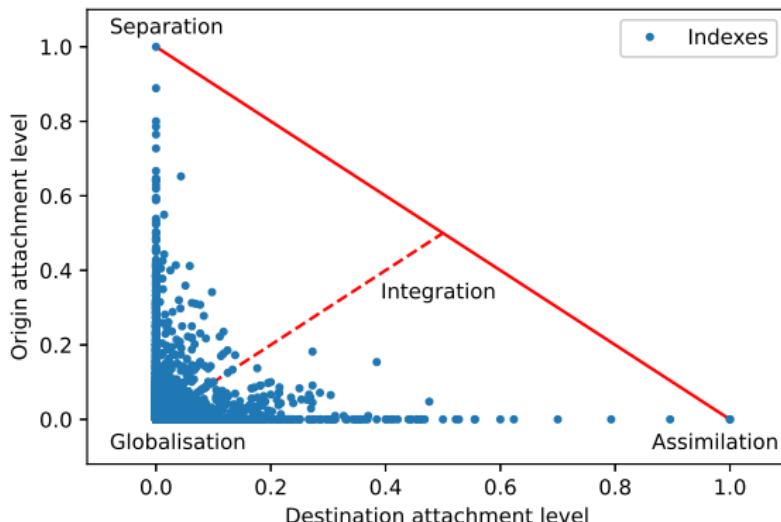


Figure: Relationship between OA & DA

Italian emigrants in overseas

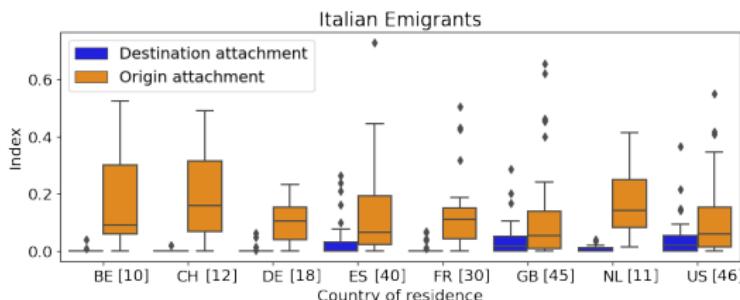
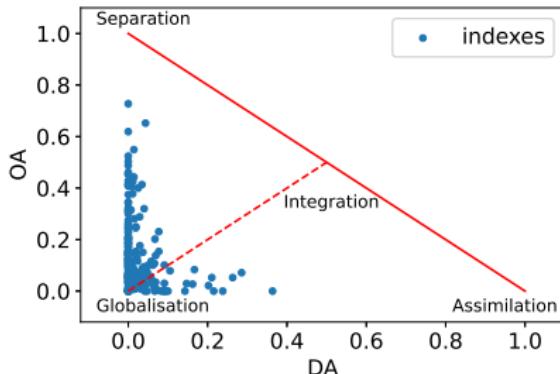
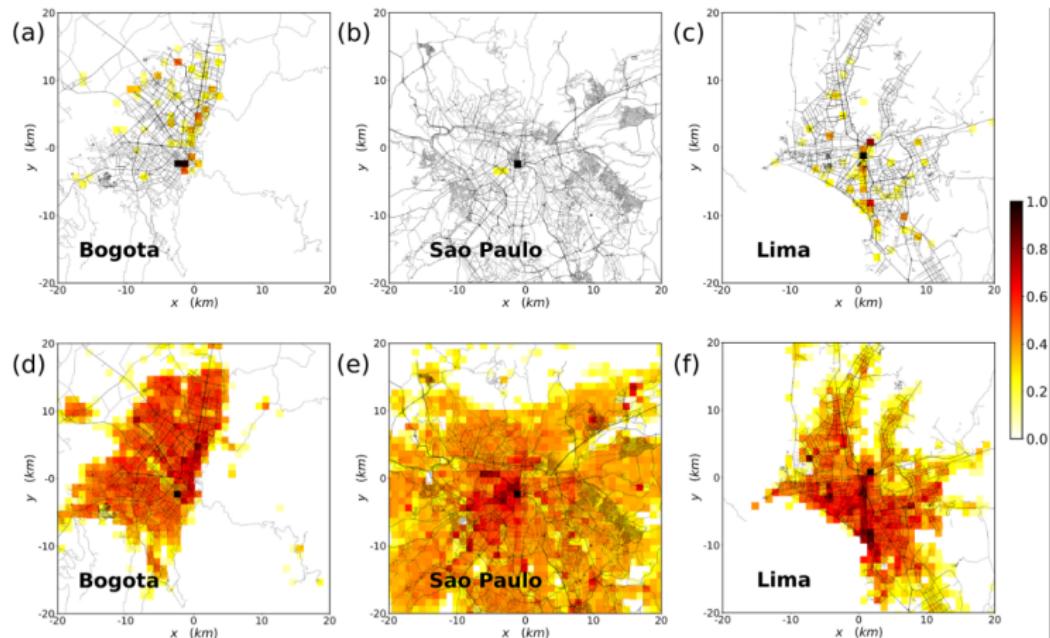


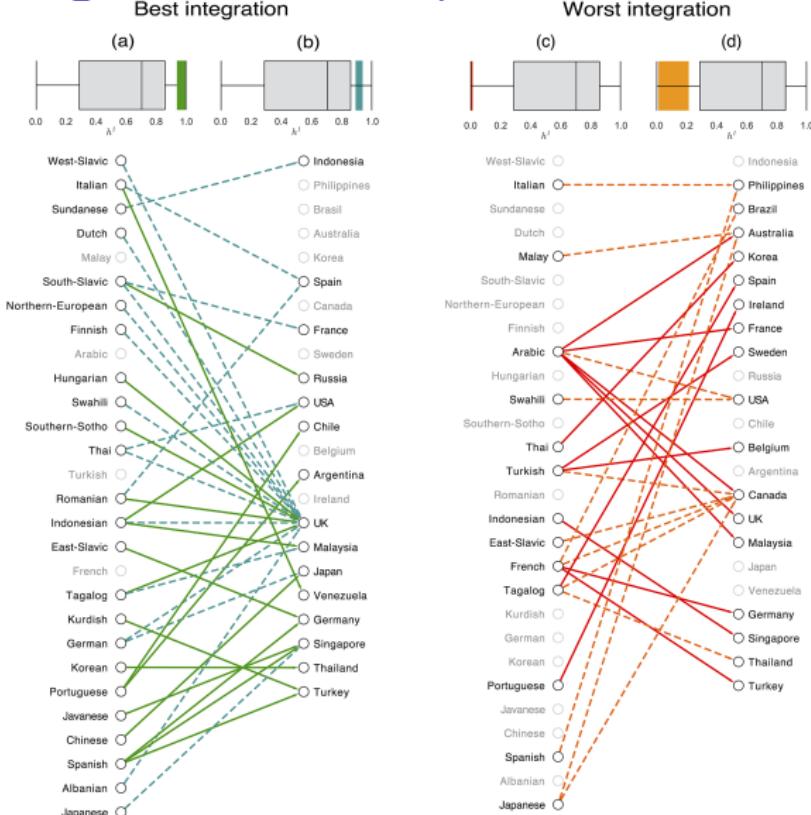
Figure: Case study of Italian emigrants in overseas

Spatial integration (Mazzoli et al., 2020)



Spatial integration of Venezuelan migrants in Colombia, Brazil and Peru
(Top figures: spatial distribution of migrants, bottom figures: spatial distribution of local population)

Language integration network (Lamanna et al., 2018)



Population distribution in function of language spoken by the users. UK sets its dominant role in uniformly integrating several communities. Arabic arises as the most common spatially segregated community followed by French speaking communities that appear to be spatially concentrated in countries like Germany and Turkey.



Sentiment analysis (Arcila-Calderón et al., 2021)

Hate speech towards migrants/refugees? What are the main topics behind the messages in Twitter in Spain with hate speech against refugees?

Topic	Most Common Words and Their Frequency within the Topic (the Underlined Words are the Most Determinant for the Labelling of the Topic)
Pull effect and consequences	"immigrants" (0.008) + "spain" (0.007) + "boat" (0.005) + " <u>effect</u> " (0.005) + " <u>pull</u> " (0.004) + "goes" (0.004) + "valencia" (0.004) + " <u>consequences</u> " (0.004) + "europe" (0.004) + " <u>concentration</u> " (0.004)
Pull effect and not welcoming "illegals"	"immigrants" (0.013) + "people" (0.008) + "spain" (0.008) + " <u>pull</u> " (0.007) + " <u>effect</u> " (0.007) + "illegal" (0.007) + " <u>protectyourborders</u> " (0.007) + " <u>spaniards</u> " (0.006) + "harbor" (0.006) + " <u>aquariusnotwelcome</u> " (0.006)
Not welcoming and terrorism	"go" (0.008) + "spain" (0.008) + "people" (0.006) + " <u>aquariusnotwelcome</u> " (0.005) + "country" (0.005) + " <u>boko</u> " (0.005) + " <u>haram</u> " (0.005) + "immigrants" (0.005) + "boat" (0.004) + "have" (0.004)
Smugglers and NGOs	"immigrants" (0.023) + "spain" (0.010) + "spaniards" + " <u>come</u> " (0.008) + " <u>mafias</u> " (0.008) + "people" (0.007) + "go" (0.006) + " <u>illegal</u> " (0.006) + " <u>ngos</u> " (0.006) + "boat" (0.006)
Money and jobs	"refugees" (0.011) + "immigrants" (0.009) + " <u>pay</u> " (0.007) + "spain" (0.007) + "countries" (0.005) + "boat" (0.005) + "spaniards" (0.005) + "people" (0.005) + " <u>solution</u> " (0.005) + "work" (0.005)
Entrance to Europe	"spain" (0.018) + "immigrants" (0.009) + "europe" (0.006) + "boat" (0.006) + "valencia" (0.005) + "spaniards" (0.005) + "government" (0.005) + "immigration" (0.005) + " <u>north</u> " (0.004) + " <u>millions</u> " (0.005)

Social integration (Kim et al., 2023)

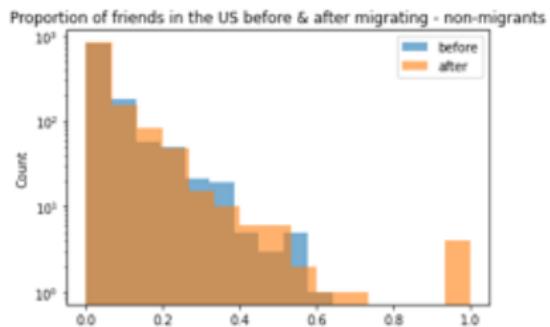
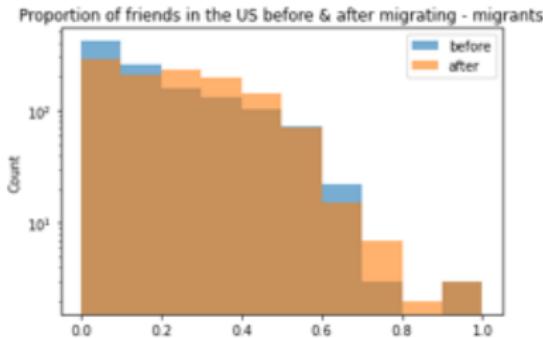


Figure: Distribution of fraction of friends residing in the United States before and after migration for migrants (left) and for propensity-score matched non-migrants (right).

BUT

Bias in the data:

“Twitter users are younger, more educated and more likely to be Democrats than general public”⁹

What about in other countries?

Sample size: Started BIG, ended small...

⁹<https://www.pewresearch.org/internet/2019/04/24/sizing-up-twitter-users/>

Pros and cons

Pros:

- ▶ Access
- ▶ Open box
- ▶ Longitudinal: access to historical data

Cons:

- ▶ Bias
- ▶ Sample size

Issue:

- ▶ Ethics and Privacy

What can we do to mitigate the issues related to Ethics and Privacy?

Ethics and Privacy

Twitter API is an open-access

Eavesdropping someone's conversation in a bus → can I expose this information?

Does this mean that we can do anything with the data?

→ Quick answer= NO!

GDPR

General Data Protection Regulation was enacted in 2016: An EU law that aims to protect and regulate data and privacy.

7 principles¹⁰:

- ▶ **Lawfulness, fairness and transparency**
- ▶ Purpose limitation
- ▶ Data minimisation
- ▶ Accuracy
- ▶ **Storage limitation**
- ▶ **Integrity and confidentiality (security)**
- ▶ Accountability

For more info: <https://gdpr.eu>

¹⁰<https://www.onetrust.com/blog/gdpr-principles/>

Consent

Article 4(11) defines consent: “Consent of the data subject means any freely given, specific, informed and unambiguous indication of the data subject’s wishes by which he or she, by a statement or by a clear affirmative action, signifies agreement to the processing of personal data relating to him or her.”¹¹

Consent also means that “the data subject knows your identify, what data processing activities you intend to conduct, the purpose of it and that they can withdraw their consent at any time.”¹²

¹¹<https://gdpr.eu/gdpr-consent-requirements/>

¹²Idem.

Confidentiality

- ▶ Pseudonymization: jisu → DKjxid73
- ▶ Anonymization: jisu → *****
- ▶ Data transformation: By modifying the format, value, or structure of data

Data storage limitation

How, where and for how long can one store sensitive and personal data?

- Before storing the data, think of pseudonymizing, anonymising or transforming the data
- Secure data storage: need to store data where it is safe and well protected where it can prevent possible attacks or loss of data
- The length of the time that you are planning to store the data is required by the GDPR.

Take home message

“Social big data can be proposed to fill some of the gaps and complement traditional data types but cannot replace them”.

-  Arcila-Calderón, Carlos et al. (2021). "Refugees Welcome? Online Hate Speech and Sentiments in Twitter in Spain during the Reception of the Boat Aquarius". In: *Sustainability* 13.5, p. 2728.
-  Hausmann, Ricardo et al. (2018). *Measuring venezuelan emigration with twitter*. Tech. rep. Kiel Working Paper.
-  Huang, Wenyi et al. (2014). "Inferring nationalities of twitter users and studying inter-national linking". In: *Proceedings of the 25th ACM conference on Hypertext and social media*, pp. 237–242.
-  Kim, Jisu et al. (2020). "Digital footprints of international migration on twitter". In: *International Symposium on Intelligent Data Analysis*. Springer, pp. 274–286.
-  Kim, Jisu et al. (2021). "Home and destination attachment: study of cultural integration on Twitter". In: *arXiv preprint arXiv:2102.11398*.
-  Kim, Jisu et al. (2023). *Online social integration of migrants: evidence from Twitter*. Tech. rep. Max Planck Institute for Demographic Research, Rostock, Germany.

-  Lamanna, Fabio et al. (2018). "Immigrant community integration in world cities". In: *PLoS one* 13.3, e0191612.
-  Laney, Doug (2001). *3D data management: Controlling data volume, velocity and variety*. URL:
<http://blogs.gartner.com/doug-laney/files/2012/01/ad949-3D-Data-%20Management-Controlling-Data-Volume-Velocity-and-Variety.pdf>.
-  Mazzoli, Mattia et al. (2020). "Migrant mobility flows characterized with digital data". In: *PLoS one* 15.3, e0230264.
-  Öztürk, Nazan et al. (2018). "Sentiment analysis on Twitter: A text mining approach to the Syrian refugee crisis". In: *Telematics and Informatics* 35.1, pp. 136–147.
-  Zagheni, Emilio et al. (2014). "Inferring international and internal migration patterns from twitter data". In: *Proceedings of the 23rd International Conference on World Wide Web*, pp. 439–444.

Supplement materials

F1-score measures the accuracy;

$$F1score = 2 * \frac{precision * recall}{precision + recall}$$

- ▶ Precision: How many selected items are relevant?

$$\frac{Truepositive}{Truepositive + Falsepositive}$$

- ▶ Recall: How many relevant items are selected?

$$\frac{Truepositive}{Truepositive + Falsenegative}$$

The score ranges from its worst score of 0 to its best score 1.