

Racism in online interactions

FRANCESCO BAILO (WITH GERARD GOGGIN)

Department of Media and Communications, The University of Sydney

Gale Digital Scholar Lab

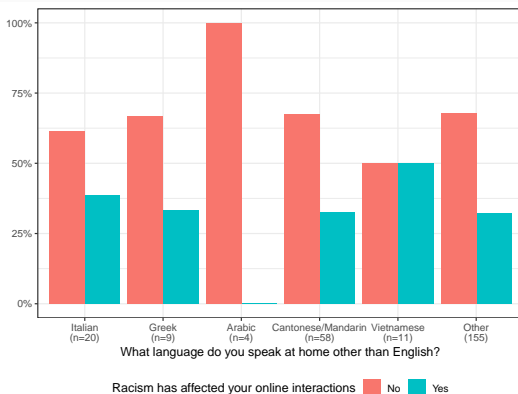
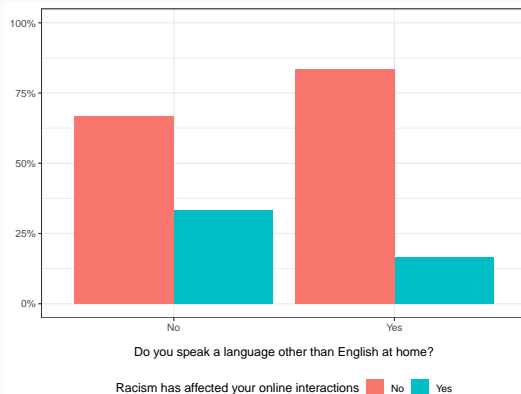
3 October 2018

Justification and research question

Research design

Justification and research question

According to a survey we conducted within the Digital Rights and Governance project among 1,600 Australians, about one-third of respondents who declare to speak a language other than English at home have been subjected to racism in their online interactions. (Goggin et al., 2017)



Research question?

Based on empirical observation, how likely an Australian Twitter user with a non-Anglo-Saxon background is subjected to racism?

Research design

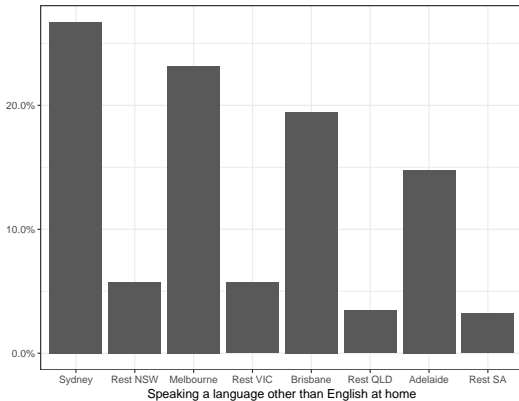
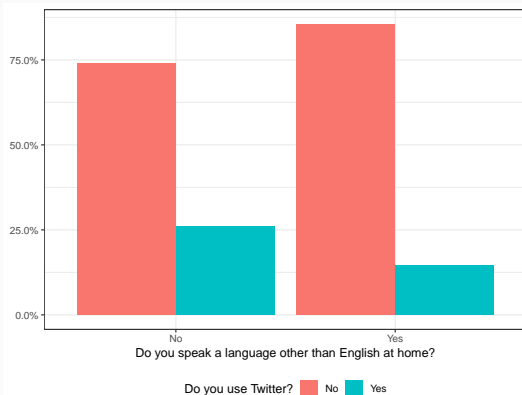
Main challenges

In order to answer to the research question, I need to

1. Create a *representative sample* of Twitter users and tweets;
2. Identify the *ethnic background* of Twitter users;
3. Assess if tweets are *racist* or not.

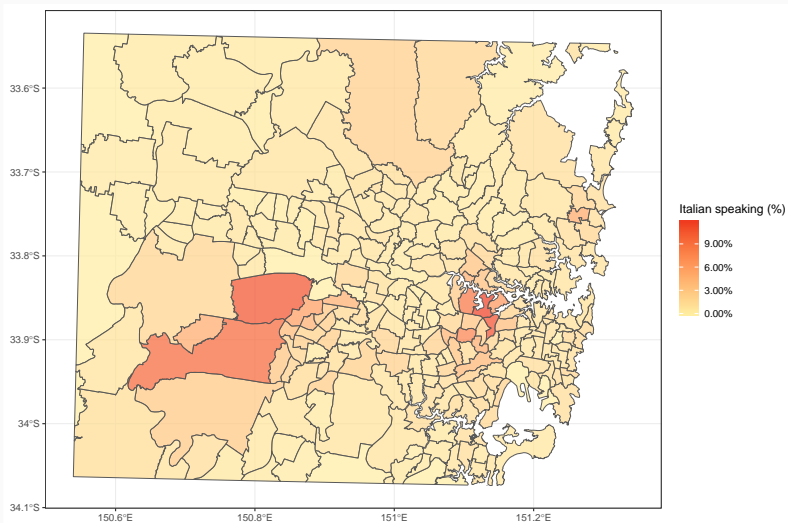
Note: The research project is under revision by the research and ethics committee of the University.

Creating a representative sample of Twitter users and tweets

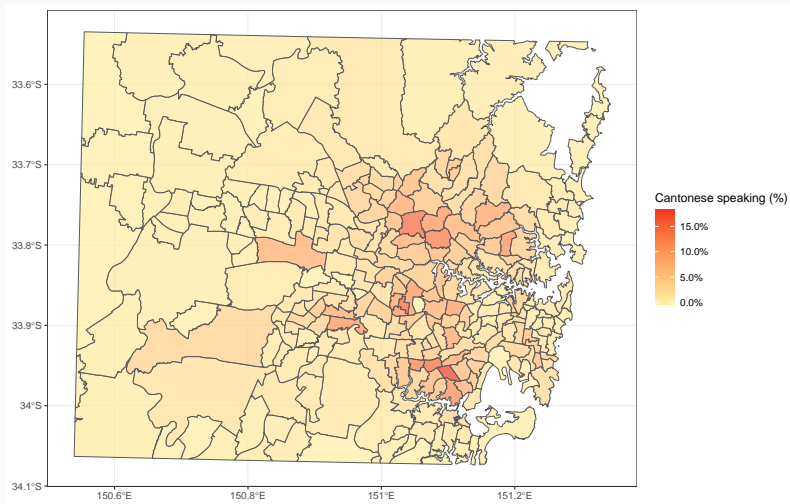


Twitter is used by less than 15% of Internet users who speaks a language other than English but they are highly concentrated in cities.

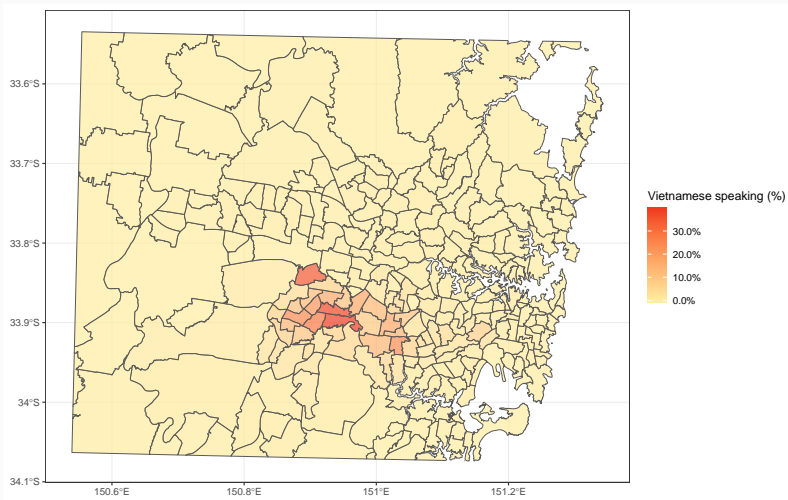
Distribution of Italian speaking population (Sydney, 2016 census)



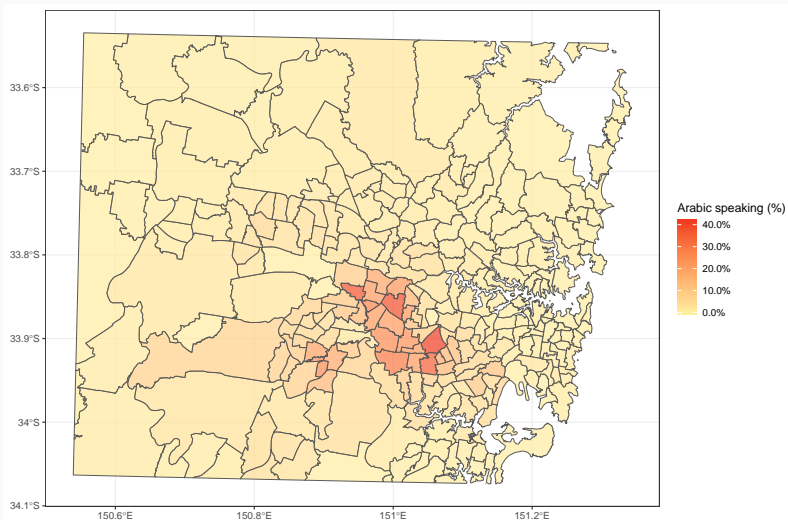
Distribution of Cantonese speaking population (Sydney, 2016 census)



Distribution of Vietnamese speaking population (Sydney, 2016 census)



Distribution of Arabic speaking population (Sydney, 2016 census)



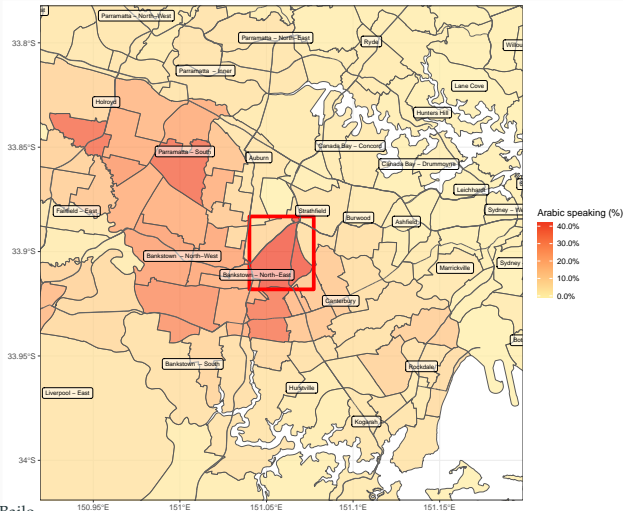
Creating a representative sample of Twitter users and tweets

- 23,401,892 (100%) people living in Australia on census night 2016;
- 6,388,717 (27%) people speaking a language other than English (PLOE);
- 4,400,000 (70%, estim.) PLOE on the Internet;
- 650,000 (14%) PLOE on Twitter.

Residents speaking a language other than English tend to cluster in specific areas.

I can use the filtering functions of the Twitter API to collect tweets published in those areas.

Distribution of Arabic speaking population (Sydney, 2016 census)



Limits to filter realtime Tweets via API for *Greenacre*

statuses/filter &

```
locations=151.04014,
```

-33.91810, 151.07729,

-33.88318

Testing Greenacre API query

Over 4 days (long weekend):

- 35,584 tweets;
- 41% of tweets are a reply;
- 6,338 Twitter users.

Creating a representative sample of Twitter users and tweets

Defining 30 bounding boxes in Sydney, Melbourne, Brisbane and Adelaide and filtering results from the Streaming API based on the diffusion of languages.

- Italian;
- Greek;
- Arabic;
- Cantonese;
- Mandarin;
- Vietnamese.

Identifying the ethnic background of Twitter users

The identification of the ethnic background of Twitter users is conducted by comparing user names to a dictionary of common first names derived from en.wiktionary.org.

The dictionary is compiled by parsing all given names/surnames and their variations in the categories:

- Italian male given names;
- Italian female given names;
- Greek male given names;
- Greek female given names;
- ...
- English surnames from Chinese;
- Vietnamese male given names.
- Vietnamese female given names.

Not logged in

TalkContributionsPreferencesCreate accountLog in

ReadEditHistory

Search Wiktionary

シ

𐄎

𐄏

𐄑

𐄒

𐄓

𐄔

𐄕

𐄖

𐄗

𐄘

𐄙

𐄚

𐄛

𐄜

𐄝

𐄞

𐄟

𐄠

𐄡

𐄢

𐄣

𐄤

𐄥

𐄦

𐄧

𐄨

𐄩

𐄪

𐄫

𐄬

𐄭

𐄮

𐄯

𐄰

𐄱

𐄲

𐄳

𐄴

𐄵

𐄶

𐄷

𐄸

𐄹

𐄺

𐄻

𐄼

𐄽

𐄾

𐄿

𐅀

𐅁

𐅂

𐅃

𐅄

𐅅

𐅆

𐅇

𐅈

𐅉

𐅊

𐅋

𐅌

𐅍

𐅎

𐅏

𐅐

𐅑

𐅒

𐅓

𐅔

𐅕

𐅖

𐅗

𐅘

𐅙

𐅚

𐅛

𐅜

𐅝

𐅞

𐅟

𐅠

𐅡

𐅢

𐅣

𐅤

𐅥

𐅦

𐅧

𐅨

𐅩

𐅪

𐅫

𐅬

𐅭

𐅮

𐅯

𐅰

𐅱

𐅲

𐅳

𐅴

𐅵

𐅶

𐅷

𐅸

𐅹

𐅺

𐅻

𐅼

𐅽

𐅾

𐅿

𐆀

𐆁

𐆂

𐆃

𐆄

𐆅

𐆆

𐆇

𐆈

𐆉

𐆊

𐆋

𐆌

𐆍

𐆎

𐆏

𐆐

𐆑

𐆒

𐆓

𐆔

𐆕

𐆖

𐆗

𐆘

𐆙

𐆚

𐆛

𐆜

𐆝

𐆞

𐆟

𐆠

𐆡

𐆢

𐆣

𐆤

𐆥

𐆦

𐆧

𐆨

𐆩

𐆪

𐆫

𐆬

𐆭

𐆮

𐆯

𐆰

𐆱

𐆲

𐆳

𐆴

𐆵

𐆶

𐆷

𐆸

𐆹

𐆺

𐆻

𐆼

𐆽

𐆾

𐆿

𐇀

𐇁

𐇂

𐇃

𐇄

𐇅

𐇆

𐇇

𐇈

𐇉

𐇊

𐇋

𐇌

𐇍

𐇎

𐇏

𐇐

𐇑

𐇒

𐇓

𐇔

𐇕

𐇖

𐇗

𐇘

𐇙

𐇚

𐇛

𐇜

𐇝

𐇞

𐇟

𐇠

𐇡

𐇢

𐇣

𐇤

𐇥

𐇦

𐇧

𐇨

𐇩

𐇪

𐇫

𐇬

𐇭

𐇮

𐇯

𐇰

𐇱

𐇲

𐇳

𐇴

𐇵

𐇶

𐇷

𐇸

𐇹

𐇺

𐇻

𐇼

𐇽

𐇾

𐇿

𐈀

𐈁

𐈂

𐈃

𐈄

𐈅

𐈆

𐈇

𐈈

𐈉

𐈊

𐈋

𐈌

𐈍

𐈎

𐈏

𐈐

𐈑

𐈒

𐈓

𐈔

𐈕

𐈖

𐈗

𐈘

𐈙

𐈚

𐈛

𐈜

𐈝

𐈞

𐈟

𐈠

𐈡

𐈢

𐈣

𐈤

𐈥

𐈦

𐈧

𐈨

𐈩

𐈪

𐈫

𐈬

𐈭

𐈮

𐈯

𐈰

𐈱

𐈲

𐈳

𐈴

𐈵

𐈶

𐈷

𐈸

𐈹

𐈺

𐈻

𐈼

𐈽

𐈾

𐈿

𐉀

𐉁

𐉂

𐉃

𐉄

𐉅

𐉆

𐉇

𐉈

𐉉

𐉊

𐉋

𐉌

𐉍

𐉎

𐉏

𐉐

𐉑

𐉒

𐉓

𐉔

𐉕

𐉖

𐉗

𐉘

𐉙

𐉚

𐉛

𐉜

𐉝

𐉞

𐉟

𐉠

𐉡

𐉢

𐉣

𐉤

𐉥

𐉦

𐉧

𐉨

𐉩

𐉪

𐉫

𐉬

𐉭

𐉮

𐉯

𐉰

𐉱

𐉲

𐉳

𐉴

𐉵

𐉶

𐉷

𐉸

𐉹

𐉺

𐉻

𐉼

𐉽

𐉾

𐉿

𐊀

𐊁

𐊂

𐊃

𐊄

𐊅

𐊆

𐊇

𐊈

𐊉

𐊊

𐊋

𐊌

𐊍

𐊎

𐊏

𐊐

𐊑

𐊒

𐊓

𐊔

𐊕

𐊖

𐊗

𐊘

𐊙

𐊚

𐊛

𐊜

𐊝

𐊞

𐊟

𐊠

𐊡

𐊢

𐊣

𐊤

𐊥

𐊦

𐊧

𐊨

𐊩

𐊪

𐊫

𐊬

𐊭

𐊮

𐊯

𐊰

𐊱

𐊲

𐊳

𐊴

𐊵

𐊶

𐊷

𐊸

𐊹

𐊺

𐊻

𐊼

𐊽

𐊾

𐊿

𐋀

𐋁

𐋂

𐋃

𐋄

𐋅

𐋆

𐋇

𐋈

𐋉

𐋊

Identifying the ethnic background of Twitter users

Labelling approaches:

- Simple match (given name/surname is present in Twitter user name);
- Supervised learning: training an algorithm by manually coding a selected sample.

Issues:

- False positives (e.g. the Arabic name Ali will match Ali short for *Alice*);
- Supervised learning requires significant time investment for manual labelling and assessment of predictions.

Creating a representative sample of Twitter users and tweets

When user names are identified as belonging to PLOE, replies directed to them will be pulled from TrISMA (trisma.org).

TrISMA have been tracking and storing tweets of the 'most active Australian Twitter accounts'.

TrISMA data will allow to expand the dataset longitudinally.

Assessing if tweets are racist or not

Labelling approaches:

Dictionary-based approach : A dictionary was created from the Wiktionary category English ethnic slurs (containing 332 terms).

Supervised learning approach based on three categories *Toxicity*, *Aggression* and *Personal Attacks* based on the manually coded dataset released by the 'Wikipedia Detox project'.

Assessing if tweets are racist or not

The **Dictionary-based approach** will return a binary answer to the question if a tweet contains racist terms.

But the **Supervised learning approach** will return a more nuanced answer based (racism is a subjective experience!) on three dimensions:

Toxicity How strongly the message is 'perceived as likely to make [the recipient] want to leave the discussion'.

Aggression 'How strongly the message is perceived as aggressive' towards the recipient.

Personal Attack '[W]hether it contains a personal attack' towards the recipient.

References



Goggin, G., Vromen, A., Weatherall, K., Martin, F., Adele, W., Sunman, L. & Bailo, F. (2017). *Digital Rights in Australia*. Sydney: The University of Sydney. Retrieved September 6, 2018, from <https://ses.library.usyd.edu.au/handle/2123/17587>