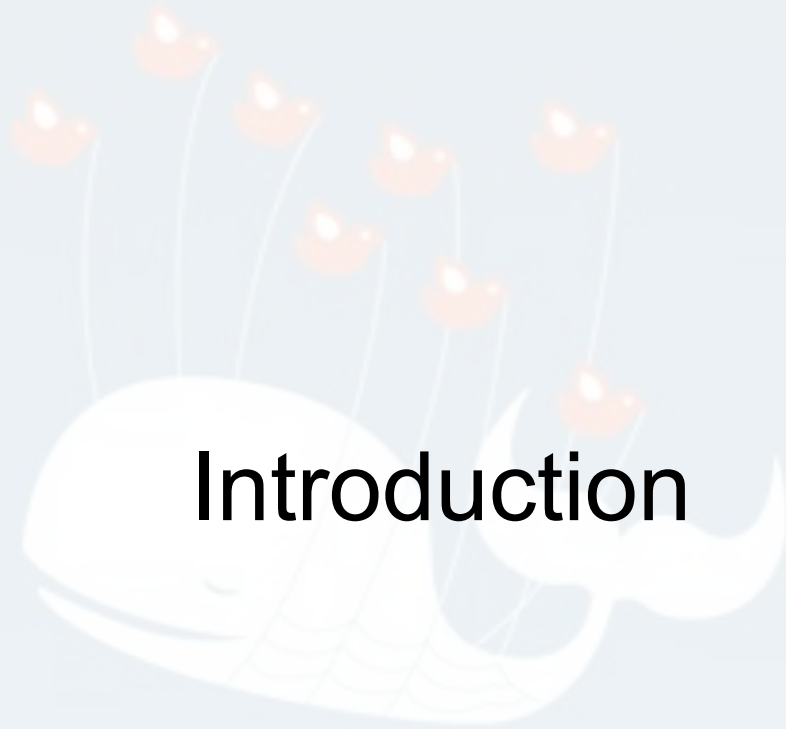# Twitter data as primary sources for historians: a critical approach

Lessons from two projects: the 2015 Greek referendum and the Centenary of the Great War on Twitter

Sofia Papastamkou (European Center for the Humanities and Social Sciences, Lille)
Frédéric Clavert (C2DH, University of Luxembourg)

# Introduction

# Introduction: social networks online, a definition

» We define social network sites as web-based services that allow individuals to (1) construct a public or semi-public profile within a bounded system, (2) articulate a list of other users with whom they share a connection, and (3) view and traverse their list of connections and those made by others within the system. The nature and nomenclature of these connections may vary from site to site. «

(Boyd Danah M et Ellison Nicole B., « Social Network Sites: Definition, History, and Scholarship », *Journal of Computer-Mediated Communication* 13 (1), 01.10.2007, p. 210-230.)

# Introduction: Twitter studies

- Since the beginning of Twitter, its data has been used in Humanities and Social Sciences for different purposes

  See: Williams Shirley A., Terras Melissa M. et Warwick Claire, « What do people study when they study Twitter? Classifying Twitter related academic papers », Journal of Documentation 69 (3), 10.05.2013, pp. 384-410. En ligne: <https://doi.org/10.1108/JD-03-2012-0027>, consulté le 24.10.2018.

- Two kinds of historian's work

  - Collective memory

    Turgeon Alexandre, « Comment Travailler La Mémoire Sur Twitter. Quelques réflexions d'ordreméthodologique à partir de la Grande Noirceur et Révolution Tranquille 2.0 », Études canadiennes / Canadian Studies. Revue interdisciplinaire des études canadiennes en France (76), 01.07.2014, pp. 11-26.

  - Current time events

    Ruest Nick et Milligan Ian, « An Open-Source Strategy for Documenting Events: The Case Study of the 42nd Canadian Federal Election on Twitter », The Code4Lib Journal (32), 25.04.2016. En ligne: <http://journal.code4lib.org/articles/11358>, consulté le 24.10.2018.
    Documenting the Now: https://www.docnow.io/

# Introduction: why Twitter?

Because we can:

- Relatively (though less and less) open APIs
- Several free and one paying APIs
  - Search API (history)
  - Streaming API (what's going on, < 1% of )
  - «If you pay you can get whatever you want» API (but we don't pay, do we?)
- Lots of tools to collect tweets
  - twarc,
  - DMI-TCAT,
  - TAGS,
  - etc.

# Introduction

Plan

I. Two projects, two theoretical backgrounds
II. Two projects, many methods and tools for the creation and the analysis of the corpus
III. Twitter hermeneutics

I. Two projects, two theoretical backgrounds

# #ww1 - The collective memory of the Great War

- Context: Centenary of the Great War
  - First large series of commemorations in the social network online era
  - Multinational(-linguistic) comparisons possible (mainly French and English)

- Collecting tweets related to the Great War
  - **mainly** inductive approach

- Studying collective memory in the digital era
  - Digital memory studies (Andrew Hoskins)
    Hoskins Andrew (éd.), Digital memory studies: media pasts in transition, New York, Routledge, 2017.
  - Will collective memory «change» when confronted to information circulation on social networks online?
    Boullier Dominique, « Big data challenges for the social sciences: from society and opinion to replications », arXiv:1607.05034 [cs], 18.07.2016.

# #greferendum - Studying the 2015 greek referendum

- Context: Greek debt crisis, Eurozone crisis,
  - A rich, born-digital (SNS), transnational documentation
  - A personal experimentation: archiving and analysing an event
- Collecting the #greferendum tweets
  - An *ad hoc* collect
  - Holistic approach (by hashtag)
- Studying the event
  - An important concept for historians (Seignobos 1898; Nora 1972; Le Goff 1999; Sirinelli 2002)
  - Twitter: the medium of the event
- Studying Twitter as a source for historians
  - non-institutional; decentralised; wild; born-digital

II. Two projects
Many methods and tools

# #ww1

- Collecting tweets
  - 140dev [abandonned] and [the incredible] [DMI-TCAT](DMI-TCAT)
    Rieder Bernhard et Borra Erik, « Programmed method: developing a toolset for capturing and analyzing tweets », Aslib Journal of Information Management 66 (3), 19.05.2014, p. 262-278.
  - a regularly updated tool, that can manage the many and regular changes in Twitter API
  - 5 Millions+ tweets, 1 million users (and GPDR headache) stored in a mariaDB database
- Preparing Twitter data for analysis
  - spreadsheets, OpenRefine, Dataïku DSS, SQL query, etc.
  - «in-between» tools that we don't always talk about (but we should)
- Analysing tweets
  - IRaMuTeQ (iramuteq.org) = data mining
  - Gephi = social networks analysis
  - Dataïku DSS / spreadsheets for simple stats

# #greferendum: collecting tweets

- Dates: 6-16 July 2015 (the "international" phase)

- Holistic collect - main hashtag: #greferendum

- NodeXL: ≈ 20,000 tweets per day

204 714 tweets:

- 139 945 retweets (68,36 %)

- 8 686 replies (4,24 %)

- 56 086 unique tweets (27,39 %)

# #greferendum: preparing data for analysis

OpenRefine => data (:hashtag) cleaning (:clustering)

TEI P5 XML => (a very basic) text encoding of data (:tweet text) - subcorpus par date

# #greferendum: analysing tweet data

- Hashtags
  Qualitative work : a typology of the most frequent hashtags (frq>99, 158 words)
  R (wordcloud package) => textual data visualisation (:hashtags)
- Tweets

  text statistical analysis  (:cooccurrences)  => TXM textométrie

- Users
  network metrics and visualisation => Gephi
  Qualitative work on most central accounts
- Domains: simple statistics with Voyant tools

# Typology 1



Hashtag type

# Typology 2: hashtag function

Commentary: 14/158 (8,861 %)

Tag: 144/158 (91,139 %)

# III. Twitter Hermeneutics

# Hermeneutics of APIs

Twitter APIs constraints: choosing an API as the first step to interpretation

- Search API: 7 days in the past, around 3000 tweets per hour
  (some workarounds: https://github.com/taspinar/twitterscraper)
  Either sampling / or small corpus
- Streaming API: anticipation of what will be the past
  Limitation: 1% of the tweets that are being published
  Progressive construction of massive corpus

# Hermeneutics of keywords: hashtags

Most research on twitter are based on keywords/hashtags which means that:

- The studied object must be quite well-known by the researcher to find the best keywords
- Hashtags / keywords are not conversation

    D'heer Evelien, Vandersmissen Baptist, Neve Wesley De et al., « What are we missing? An empirical exploration in the structural biases of hashtag-based sampling on Twitter », First Monday 22 (2), 16.01.2017.

- Therefore
    - collecting massive data != collecting exhaustive data
    - sampling data can be better than massively collecting data
- Numerous ways to understand what a hashtag is

# … some thoughts from the #greferendum corpus

Hashtags

- tell the big story (quantitative + relational analysis)
- reveal different temporalities related to connected histories of the Eurozone crisis
- a common conversation? a European space? (cf. works of Camille Roth)

# Hermeneutics of networks

See:
http://theconversation.com/four-more-years-that-obama-tweet-and-the-politics-of-intimacy-10606

# … some thoughts from the #greferendum corpus

- main actors in a graphe - main actors irl?
- what is a Twitter network?

# Hermeneutics of tools

Hypothesis: a tool = a method = a theory = a specific way to interpret data

- **Gephi**
  Visualizing social networks => sociology of social networks != sociology of field and *habitus*

- **IRaMuTeQ**
  *Théorie des mondes lexicaux*

  Reinert Max, « Une méthode de classification descendante hiérarchique: application à l'analyse lexicale par contexte », Les cahiers de l'analyse des données 8 (2), 1983, pp. 187-198.
  - French School of Data Analysis (yes, there is one)
  - Mondes lexicaux: one point of view = one coherent set of words = social representations (= Émile Durkheim)

# … some thoughts from the #greferendum corpus

- Dataviz is useful… metrics are important
- How to be comfortable with the algorithm? (transparency, stability issues)
- Need for tools that behave well with multilingual corpora (TXM is fine)
- Preservation and sharing issues

# Hermeneutics of Twitter: Twitter as a primary source

- A primary source in the historian's point of view
  Traditionally something that is fixed within a set framework (= the Archive)

- Twitter is always moving, is a "source" in the original meaning (source of water): something that is endlessly flowing, that cannot by definition be fixed
  Ex: The "four more years" Obama tweet

- What we do
  - transforming something that is not supposed to stay still into an archive, something that is fixed
  - What do we lose in this process?

# Hermeneutics of Metadata

Information embedded in the metadata are crucial for the interpretation of a / numerous tweet(s)

Ex: timestamps and the interpretation of temporalities

- Timestamps in tweet metadata correspond to the unending (well…) and continuous feed of tweets that is the essence of Twitter
  - Western vision of time
- Many more artifacts of temporalities are embedded in the text of a tweet
  - How to deal with other kind of temporalities whereas collective memory, for instance, is the result of an interlacing of temporalities

Conclusion: what is the allure of born digital archive?

# What is a tweet?

[{
  "created_at": "Thu...
2017",
  "id": 87799460456...
  "id_str": "8779946...
  "text": "Creating a...
Angular, Part 1: Add...
https://t.co/xFox78ju...
  "truncated": false,
  "entities": {
    "hashtags": [{
      "text": "Angular"...
      "indices": [103,...
    }],