# The 'Arquivo de Opinião' archive

# Miguel Won
**TPDL 2018**

# About me

## *Miguel Won*

2015-2018: FCT Postdoc researcher at INESC-ID in the field of Computational Social Science

miguelwon@tecnico.ulisboa.pt

# 1 Introduction

Political Punditry

# Political commentary

- Political commentary is present in everyday news media:
  - "Experts" in TV broadcasting channels
  - Columnist in newspapers

- This type of opinion plays an important role in the process of the *narrative construction* of the *public realm*:
  - Selection of events
  - Authority position
  - *Deciphers* the political complexities

# Opinion articles

- In this work we consider only opinion articles from newspapers
- Definition: journalistic article, usually about the current  the state of public affairs, authored by one or multiple authors, that expresses the author's personal opinion
- Two-sided role in respect to public opinion
  - They can be interpreted as a mirror of the public opinion
  - But can also be accused of its main influencer
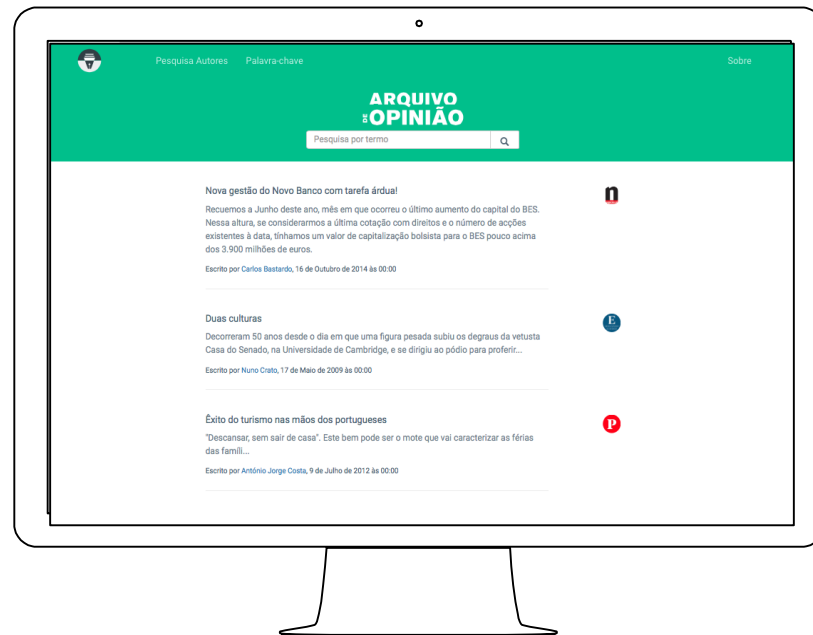
- Essential component of the public debate

# Memory

- Memory of political debates allows the recalling of ideas, main debatable issues, the argumentative logics, as well as the political positions of the various political actors (many political commentators are, were or will be themselves active politicians)
- Memory of political discussion is essential to the proper functioning of democracies
- Archives of this type of memory contributes to a healthy public debate
- Such archives should be digital:
  - Search engine
  - Searches by author, time period or media source
  - Public availability and user friendly

6

# *Arquivo de Opinião*

◉ Digital archive of opinion articles

## 2 Collect & process

Arquivo construction

# Data sources

◉ <mark>Arquivo.pt</mark>: web archive of .pt domain

◉ Opinion section (online)

# Pipeline

| URL identification | Web scraping | Data Cleaning | NLP |
|---|---|---|---|

**URL identification**
- Search for clues such as "opiniao"
- web crawling

**Web scraping**
- Title
- Author
- Publication date
- Body
- ...

**Data Cleaning**
- Name correction
- Remove html code
- Manual inspection

**NLP**
- Part-of-speech tagging
- NER
- Key-phrases extraction

# Pipeline

Tools:

- Python: nltk, re, scikit-learn, etc.



- Scrapy (web scraping)



- Lx-Tagger (pos tagging)



- Stanford NER

Web framework:

- Django 

- MongoDB

# 3 NLP

NLP tasks

# Named Entity Recognition (NER)

◉ Task: given a text as input <mark>identify the entities</mark> within the text
  - ○ Person names
  - ○ Locations
  - ○ Organizations

## NER Examples

**Input**: Vancouver is a coastal seaport city on the mainland of British Columbia. The city's mayor is Gregor Robertson.

Location

**Output**: *Vancouver* is a coastal seaport city on the mainland of *British Columbia*. The city's mayor is *Gregor Robertson*.

Location                                                                    Person

Bryan Perozzi  Stony Brook University  Polyglot-NER: Massive Multilingual Named Entity Recognition

Miguel Won
TPDL 2018

# Named Entity Recognition (cont.)

- ◉ <mark>Classification task</mark> (sequential)

- ◉ Many free tools available in the market
  - ○ Stanford NER (CRFs)
  - ○ spaCy (NN)
  - ○ Polyglot (NN)

- ◉ We have trained Stanford NER with an annotated corpus for Portuguese (European): <mark>CINTIL</mark>

# Stanford NER with CINTIL

5-fold Cross-Validation Precision, Recall and F-Measure Results for NER using CINTIL

|       | True Positive | False Positive | False Negative | Precision | Recall | F-Measure |
|-------|---------------|----------------|----------------|-----------|--------|-----------|
| 1     | 952           | 221            | 236            | 0.81      | 0.81   | 0.81      |
| 2     | 846           | 208            | 230            | 0.81      | 0.79   | 0.79      |
| 3     | 921           | 243            | 261            | 0.79      | 0.78   | 0.79      |
| 4     | 939           | 209            | 327            | 0.82      | 0.74   | 0.78      |
| 5     | 892           | 213            | 252            | 0.81      | 0.78   | 0.79      |
| Total | 4550          | 1094           | 1306           | **0.81**  | **0.78** | **0.79** |

# Key-phrase extraction

"Automatic extraction of relevant key-phrases for the study of issue competition",
work in progress with Bruno Martins (INESC-ID) and Filipa Raimundo (ICS)

- Key-phrase: a word or phrase represents a concept, idea, entity,
  etc.
  - Refugee Crisis
  - National Health Service
  - António Costa

- Politicians often guide their speeches using key-phrases
- Key-phrase identification can hint us about the topics addressed in a
  set of speeches

16

# First step: Candidate Selection

- ◉ Part-of-Speech tagging followed by a chunk rule:
  - ○ Crise dos Refugiados: NOUN+PREP+NOUN
  - ○ Sistema Nacional de Saúde: NOUN + ADJ + PREP + NOUN
  - ○ António Costa: NOUN + NOUN

Chunking rule (Portuguese): (<NOUN>+ <ADJ>* <PREP>*)? <NOUN>+

# Second step: rank

- ◉ Several methods: TextRank, Phraseness & Informativeness, EmbedRank, etc.
- ◉ We can achieve state-of-the-art results with simple heuristic rules:
    - ○ Tf-idf
    - ○ Likelihood metric based in the position
    - ○ Length

# 4 Arquivo de Opinião

Opinion in the Portuguese media
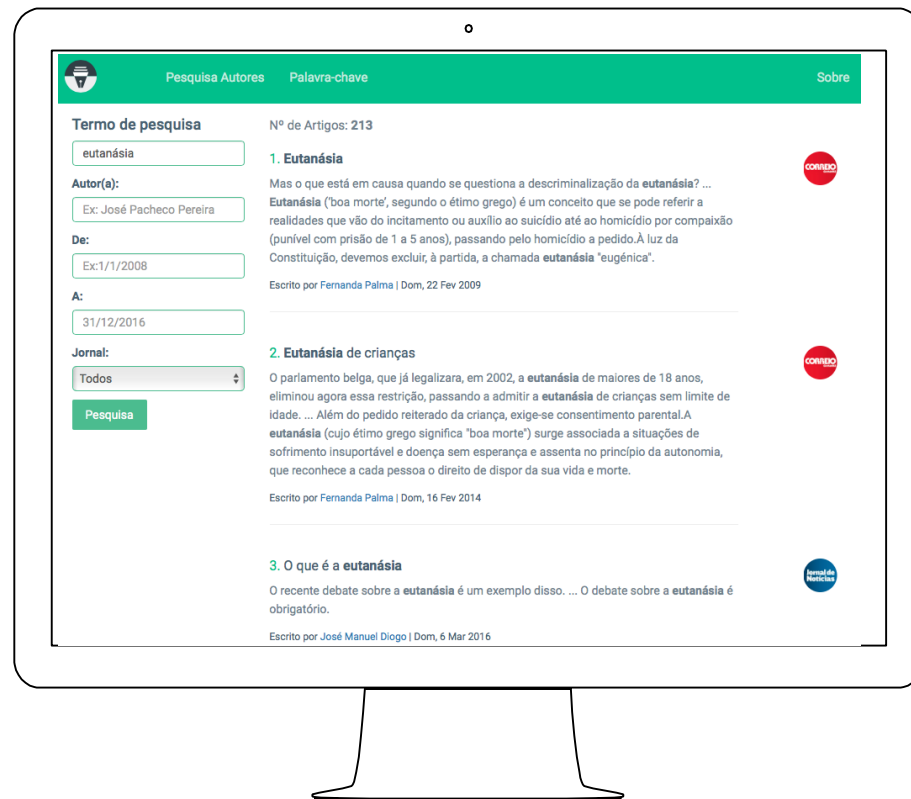
# Arquivo de Opinião

◉ Frontpage with a search engine

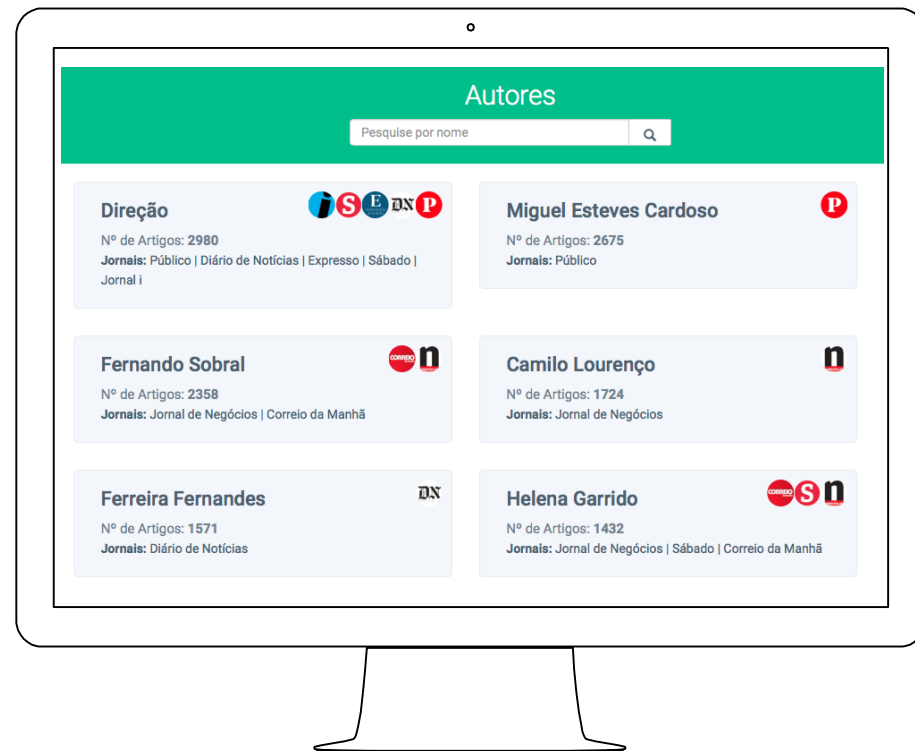# Search engine (mongo)

◉ Search for text of phrase

◉ Filters:
  ○ author
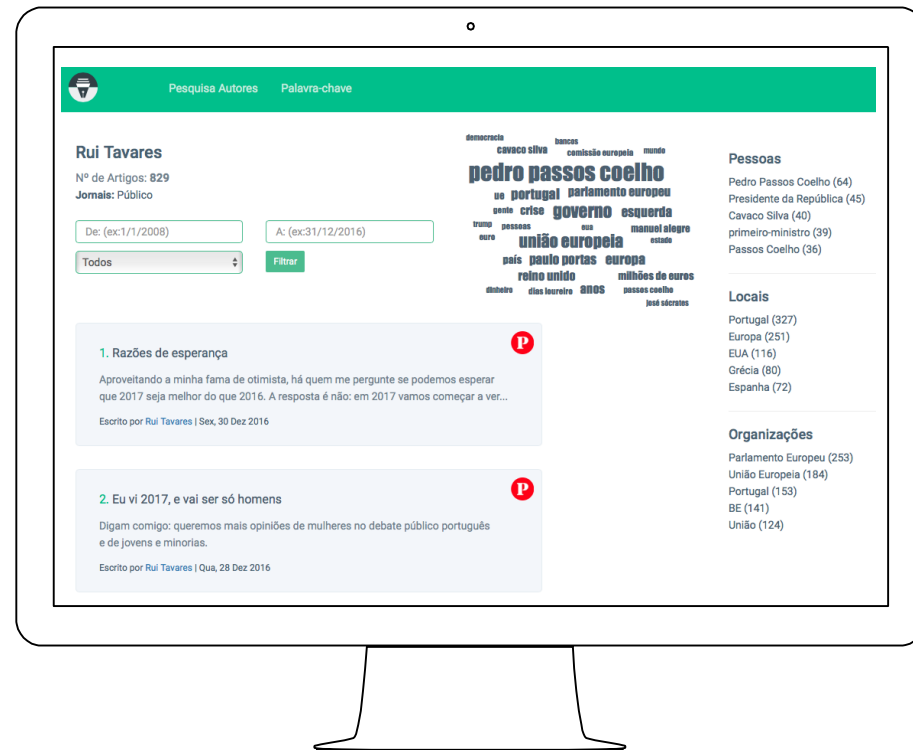  ○ time interval
  ○ source

# Author

◉ Search for author

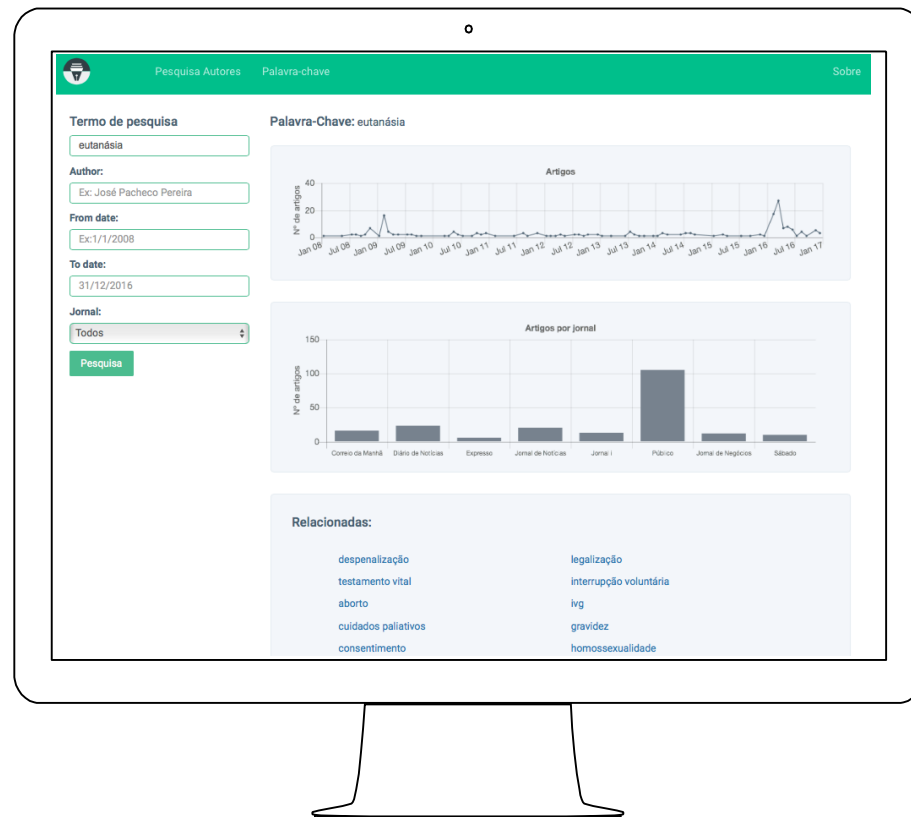◉ ~3500 available authors

# Author (cont.)

◉ Each author has its page

◉ Key-phrases cloud

◉ Mentioned entities:
  ○ Persons names
  ○ Locations
  ○ Organizations

# Key-phrases

- Search indexed key-phrases (with autocomplete)

- Outputs
  - No. articles by time and source
  - Related (word embeddings)

**85 530**
Articles

**3571**
Authors

**30 000**
key-phrases

**9**
Years of publications (2008-2016)

# 5 Next steps and final remarks

# Version 2.0

- Add additional sources: Observador, O Jornal Económico
- 2016-Present
- Social Media:
  - Authors pages
  - Shares, likes, etc.
  - Networks
- Real time monitoring (daily, weekly?)
- Add more NLP metrics: topic modeling, sentiment, etc.

# Final remarks

- Political commentary is an important section of newspaper media
- A digital archive of this type of memory contributes to a better public debate
- *Arquivo de Opinião* main objective is to offer a digital online archive of the political opinion published in the main Portuguese newspapers
- All data was processed in order to extract additional information (NLP)
- Future work will be carried out towards the inclusion of external data, in particular from social media

# Arquivo.pt awards (3rd place)

# Acknowledgements