

# Detecting Real-World Influence Through Twitter

Jean-Valère Cossu  
NLP team @ University of Avignon, France

ENIC 2015

22 September 2015

## Participants

- **J.-V. Cossu\***, N. Dugue\*\* and V. Labatut\*
  - \* Université d'Avignon
  - \*\* Université d'Orléans

## Get started with RepLab Online Reputation Monitoring on Twitter

### What ?

- Assist brand managers in their daily work

### How ?

- Learn experts behaviour about entities concerns
- Automatically propagate these assessments
- Justify the hypothesis

### 2012

- Monitoring in unknown entity case
- New challenge, is it possible to compute the issue ?
- Features discovery

### 2013

- Machine Learning
- Messages selection/ranking
- Raises new issues

### 2014

- **User profiling**
- Age, gender, occupation, influence, traits and so on

## Data

- Large profiles collection (7,000 profiles manually labelled)
- 2 economic domains Automotive, Banking
- Each profile comprises the last 600 posted tweets

## Main objectives

- **Identify Opinion-Makers**  
70% are not influential - Search/IR
- Categorize profile according to their activity
- See PAN for gender and age identification

## Piotr's profile

**Piotr Bródka**  
@BrodkaPiotr

I am assistant professor of Computer Science at the Department of Computational Intelligence, Wroclaw University of Technology, Poland.

Wroclaw, Poland  
[il.pwr.edu.pl/~brodka/index\\_...](http://il.pwr.edu.pl/~brodka/index_...)

TWEETS 14 ABONNEMENTS 45 ABONNÉS 22 FAVORIS 3

Tweets Tweets & réponses Photos & vidéos

**Piotr Bródka** @BrodkaPiotr · 25 août  
#SNAA2015 is over We had 8 great papers, awesome audience and discussion.  
It was great pleasure See you next week during #SNAA2016

**Piotr Bródka** @BrodkaPiotr · 25 août  
This year Best Paper Award goes to Sarka Zehnalova, Milos Kudelka and Zdenek Horak. Congrats. #SNAA2015 #ASONAM2015

- Is Piotr an influent Twitter user (in his community-domain) ?
- What about Gender ? Age ? Traits ?

## RepLab Author Profiling

- (Social) Network approaches
- Content-based approaches

## Features

- Public profile (description and personal data)
- (S)N features (followers, followees, etc.)
- Writing behaviour (retweets, hashtag, links)
- External data (Klout, Kred, Google)

## Features issues

- Features names
- Features relevance to the problem
- Data crawling, impossible to collect the complete network

## Profile representation

- A profile is a **Bag of Tweets** where each tweet is processed  
Tweets can be bag-of-words or set of features
- A profile is a vocabulary or a set of features (**User as Document**)

## Related content-based features

- Bag of words binary representation (word presence)
- Tweet length, Special characters, Hashtags, Links

## Extra-features

- Part-of-speech tagging, Named entities
- Tweet enrichment



## Proposed content-based features

- Bag of  $N$ -grams with TF, IDF
- Purity index (word distribution through classes) :

$$G(i) = \sum_{c \in \mathbb{C}} \mathbb{P}^2(i|c) = \sum_{c \in \mathbb{C}} \left( \frac{DF_i(c)}{DF_T(i)} \right)^2 \quad (1)$$

## Strategies

- Monolingual and domain specific model VS global model
- Majority vote over the profile

# What about performances ?

## Is tweet content sufficient ?

- Cosine distance between document and influents class' vocabulary
- No specific training, parameters tuning, or features selection

## Author ranking performances

System	Automotive	Banking	Avg MAP
<b>Cosine</b>	<b>.803</b>	.626	<b>.714</b>
LIA_Participation	.764	<b>.652</b>	.708
<i>REPLAB1</i>	.721	.410	.565
<i>Baseline</i>	.370	.385	.378
<i>Klout</i>	.304	.275	.289

- LIA\_Participation is a manually tuned KNN\* (not replicable)
- REPLAB1 used hot topics information
- Baseline rank according to the Followers number
- Raw features (Tweets, followees etc.) are under Klout

## Influence detection classification

System	Automotive	Banking	Avg F
<b>Cosine</b>	<b>.833</b>	<i>.751</i>	<b>.792</b>
LIA_new	.702	.726	.714
<i>Best_System</i>	.696	.693	.694
<i>Baseline</i>	.500	.500	.500

- Same configuration ... same results

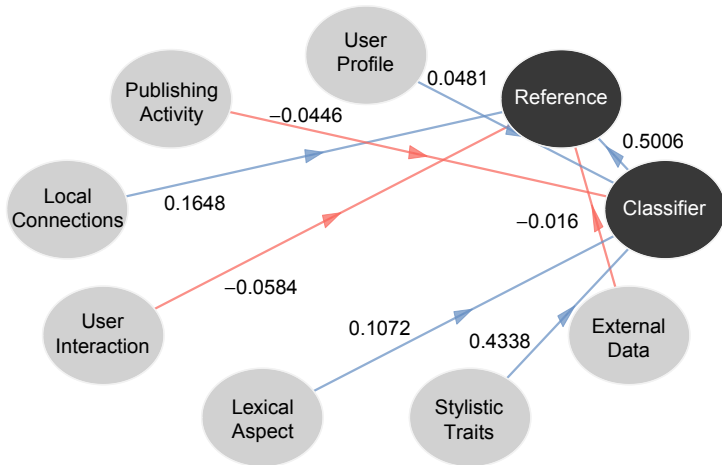
## What we have done here ?

- **Distinguish influencers using their account characteristics**
- **Distinguish influencers using contents they produce**
- Ranking them
- Real influence is not node importance
  - \* *Results limited to this dataset and the annotation quality*

## What's next ?

- Profile summarization in progress
  - How much we need ?*
  - Does it work with only 50 tweets ?*
- Visual modelling

## Features behaviour regarding influence for banking



Thank you !

Contact:

- [jvcossu@gmail.com](mailto:jvcossu@gmail.com)
- [www.jeanvalerecossu.fr](http://www.jeanvalerecossu.fr)

## References

- **NLP-based classifiers to generalize experts assessments in E-Reputation monitoring**  
*Cossu J-V., Ferreira E., Gaillard J., Janod K. and El-Bèze M. @CLEF 2015*
- **Automatic Classification and PLS-PM Modeling for Profiling Reputation of Corporate Entities on Twitter**  
*Cossu J-V., SanJuan E., Torres J-M. and El-Bèze M. @NLDB 2015*
- **Overview of the 3rd Author Profiling task at PAN 2015.**  
*Rangel F., Rosso P., Potthast M., Stein B., and Daelemans W., D. @CLEF 2015*
- **Overview of RepLab 2014: Evaluating Online Reputation Monitoring Systems.**  
*Amigó, E. and Carrillo de Albornoz, J. and Chugur, I. and Corujo, A. and Gonzalo, J. and Martín, T. and Meij, E. and de Rijke, M. and Spina, D. @CLEF 2014*