

Identifying Entity Aspects in Microblog Posts



damiano@lsi.uned.es, {edgar.meij, derijke}@uva.nl, {oghina,mbui,mbreuss}@science.uva.nl



UNED NLP & IR Group
Madrid, Spain

35th International ACM SIGIR Conference on Research and Development in Information Retrieval
Portland, Oregon, USA. August 12–16, 2012.

ISLA, University of Amsterdam
Amsterdam, The Netherlands

Introduction

- ❖ Scenario: **Online reputation management**

What do people say about an *entity*?

entity = { company, organization, individual, product }

Identification of aspects discussed on microblog posts

- ❖ **Aspects** \approx products, services, competitors, key people related to the entity

Dataset

- ❖ 94 entities, 17,775 tweets, \approx 177 tweets/entity

- ❖ Built upon WePS-3 ORM Task Dataset.
Disambiguated company names (e.g. apple fruit vs. Apple Inc.)



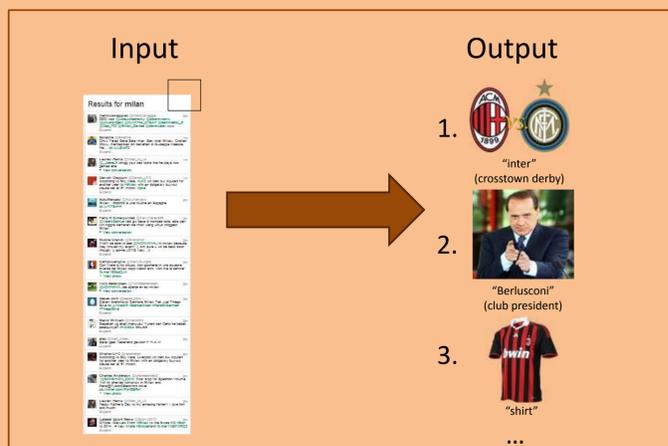
- ❖ Pooling methodology.
Three annotators with substantial agreement
(Cohen/Fleiss' kappa > 0.6)

<http://bit.ly/profilingTwitter>

- ❖ Annotated 2455 terms, 1304 aspects (54.11%)
Most of the true aspects are nouns (89.72%).

Entity	Aspects
Apple Inc.	ipad, iphone, prototype, store, gizmodo, employee
Sony	advertising, headphones, digital, pro, music, xperia, dsc, x10, bravia, camera, vegas, battery, ericsson, playstation
Starbucks	coffee, latte, tea, frappuccino, barista, drink, mocha

Identifying Entity Aspects



Main idea

- ❖ Comparing a **pseudo-document** D built from entity-specific tweets with a **background corpus** C .
- ❖ Score of a term $t = s(t, D, C)$

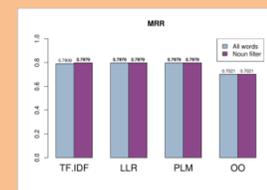
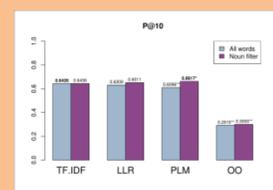
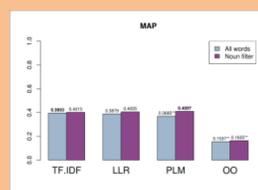
Four models tested:

- ❖ **TF.IDF**
- ❖ **LLR**: Log-Likelihood Ratio
- ❖ **PLM**: Parsimonious Language Models
- ❖ **OO**: Opinion-oriented. Opinion target extraction using topic-specific subjective lexicons

Method	Scoring function
TF.IDF	$s(t, D, C) = tf(t, D) \cdot \log \frac{N}{df(t)}$ $N = \text{number of pseudo-documents } D_i \text{ in } C$ $df(t) = \sum_i^N tf(t, D_i) > 0, D_i \in C$
LLR	$s(t, D, C) = 2 \cdot ((a \cdot \log(\frac{a}{E_1})) + (b \cdot \log(\frac{b}{E_2})))$ $E_1 = \frac{c(a+b)}{c+d}, E_2 = \frac{d(a+b)}{c+d}$ $a = tf(t, D), b = cf(t)$ $c = \sum_i^N tf(t_i, D), d = \sum_i^N cf(t_i)$
PLM	$s(t, D, C) = P(t D)$ (when model converges) E-step: $e_t = tf(t, D) \cdot \frac{\lambda \cdot P(t D)}{(1-\lambda) \cdot P(t C) + \lambda \cdot P(t D)}, \lambda = 0.1$ M-step: $P(t D) = \frac{e_t}{\sum_i e_i}$ initial $P(t D) = \frac{tf(t, D)}{\sum_i tf(t_i, D)}, P(t C) = \frac{cf(t)}{\sum_i cf(t_i)}$
OO	$s(t, D, C) = \chi^2(\text{target}(t, D), \text{target}(t, C))$ $\chi^2(o, e) = \frac{o - e}{(o + e)}$ $\text{target}(t, D) = \text{freq. of potential target } t \text{ in tweets } D$ $\text{target}(t, C) = \text{freq. of potential target } t \text{ in background } C$

Experiments and Results

- ❖ **All Words**:
 - **TF.IDF** significantly **outperforms PLM and OO in precision**.
 - The results for **OO** are **much lower** than for the other methods.
- ❖ **Noun filter**:
Applying a **part-of-speech filter** and only consider terms tagged as **nouns**
 - For all methods, MAP and precision values are slightly higher than in the all words condition: considering only nouns helps to identify aspects.
 - PLM best method using Noun filter.
- ❖ **Observations** (after manually inspecting the results):
 - Results for **TF.IDF, LLR and PLM** are **very similar**.
 - **OO** tends to return more **subjective terms as aspects** (syntactic parsing errors)
Examples: haha, pls, xd, safety, win
 - **OO** has more **difficulty to filter out generic terms**.
Examples: new, use, today, come



Student's t-test statistic significance w.r.t to TF.IDF All words baseline with $\alpha=0.05$ (*) and $\alpha=0.01$ (**).

Conclusions

- ❖ Simple statistical methods such as **TF.IDF** are a **strong baseline** for the task of **identifying entity aspects**, significantly **outperforming opinion-oriented methods**.
- ❖ **Only considering terms tagged as nouns improves the results for all the methods analyzed.**