



# Rozpoznávání ironie na Twitteru

Francesco Barbieri  
& Horacio Saggion

Univerzita Pompeua Fabry,  
Španělsko

# Co je to ironie?

- Na definici nepanuje shoda.
- Je řečen je opak toho, co má být sděleno. **Jde o negaci bez explicitních negačních znaků.**
- Dochází k záměrnému porušení očekávání. **Důležité je tedy překvapení.**
- Opakování cizího výroku, které má za úkol ukázat jeho nepravdivost.

# Trénovací data

- 40 000 příspěvků z Twitteru (v angličtině)
- Strojové rozřazení podle #hashtagů
- Čtyři kategorie: **ironické příspěvky, vzdělání, humor a politika**

# Použité charakteristiky

- Obvyklost použitých slov
- Kontrast psaného a mluveného jazyka
- Intenzita hodnotících slov
- Struktura příspěvku (délka, interpunkce, ...)
- Emoční povaha slov
- Volba synonym
- Používání víceznačných výrazů

# Obvyklost (frekvence)

- Překvapivý dojem může být vyvolán **kombinací obvyklých a velmi vzácných slov.**
- Frekvence jsou převzaty z Amerického národního korpusu.
- Měří se průměrná frekvence slov v příspěvku, frekvence nejneobvyklejšího slova a rozdíl těchto hodnot.

# Psaný versus mluvený jazyk

- Myšlenka: překvapivě působí **použití hovorového výrazu ve formálním textu**, nebo také naopak.
- Měří se průměrná frekvence slov v psaném jazyce a v mluveném jazyce a rozdíl těchto dvou údajů.

# Struktura příspěvku

- Ironické příspěvky jsou v průměru delší, slova v nich použitá jsou ale naopak kratší.
- **Zkoumá se interpunkce:** čárky, výpustky, vykřičníky atd.
- Ironické příspěvky používají méně smajlíků, s jednou výraznou výjimkou – **mrknutí ;)**
- Měří se také četnost citoslovečných výrazů typu „hahaha“, „lol“ apod.

# Intenzita hodnotících slov

- Měří se intenzita adjektiv a adverbii.  
(Každému slovu je bez hlubšího zkoumání přiřazen nejčastější slovní druh.)
- Ohodnocení slov je založeno na korelaci komentářů s uděleným počtem hvězdiček v uživatelských hodnoceních produktů na internetu.
- **Počítá se absolutní součet hodnotících slov a průměrné hodnocení na slovo.**



# Synonyma

- Překvapení může být vyvoláno užitím **neobvyklého synonyma ve společnosti obvyklých**, nebo také naopak.
- Množiny synonym jsou převzaty z WordNetu.

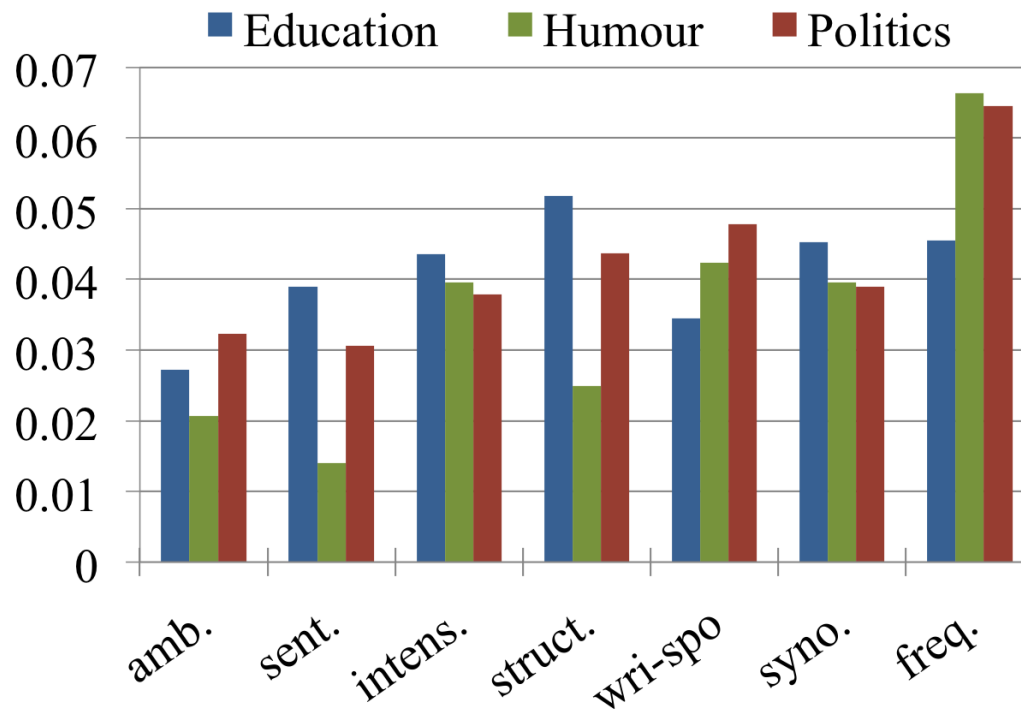
# Víceznačnost

- Ironie může využít nejednoznačných slov.
- Postup je podobný jako u synonym – zjišťuje se **užití velmi mnohoznačného slova v prostředí slov jednoznačných.**
- Opět se využívá WordNet.

# Nálada

- Tato charakteristika je podobná intenzitě hodnotících slov.
- Zjišťuje se užití **výrazně jinak pozitivně/negativně zabarvených slov, než je průměr příspěvku**, tedy jakási náladová nevyváženost.
- Využívá se slovník SentiWordNet.

# Přínos jednotlivých atributů



# Výsledky

- Přesnost rozpoznání ironie se pohybuje v rozmezí 60–75 %, podle použité trénovací množiny.
- Nejužitečnější charakteristiky jsou: **neobvyklost slov, volba neobvyklých synonym a zvláštní interpunkce.**

# Analýza na základě použitých slov (bag of words)

- Přesnost dosahuje 50–85 %.
- Bag of words přináší lepší výsledky při klasifikaci do jednotlivých tříd.
- Problém ale nastává při testování na jiné třídě, než byla použita pro trénink (pak úspěšnost klesá k 50 %).
- Při práci s obecným textem se proto jeví jako výhodnější použít popsané lingvistické charakteristiky.

# Zdroj

- Student Research Workshop,  
14. konference EACL, duben 2014
- <http://aclweb.org/anthology//E/E14/E14-3007.pdf>

Zpracoval: Josef Plch, Masarykova univerzita, 2014