



UNIVERSITÀ DEGLI STUDI DI MILANO - BICOCCA

Scuola di Scienze

Dipartimento di Informatica, Sistemistica e Comunicazione

Corso di laurea in Informatica

Enhancing irony detection with affective information

Relatore: Prof. Elisabetta Fersini

Relazione della prova finale di:

Gianluca Giudice

Matricola 830694

Anno Accademico 2019-2020

Contents

1	Stato dell'arte	4
1.1	Approccio supervisionato	4
1.2	Approccio semi supervisionato	4
1.3	Approccio non supervisionato	4
2	Sistema realizzato	5
2.1	Descrizione del sistema proposto	5
2.2	Dataset	5
2.2.1	Metodi di labeling	5
2.2.2	Dataset utilizzato	5
2.3	Rappresentazione del testo	6
2.3.1	Rappresentazione booleana	6
2.3.2	Rappresentazione mediante transformer	6
2.4	Caratteristiche linguistiche	6
2.4.1	PP	6
2.4.2	POS	6
2.4.3	Onomatopeic	6
2.4.4	...Tutto il resto	6
2.4.5	EMOT	6
2.5	Modelli supervisionati	6
2.5.1	Alberi di decisione	6
2.5.2	Support Vector Machine	6
2.5.3	Naive Bayes	6
2.5.4	Bayesian Network	6
2.6	Strumenti utilizzati	6
2.6.1	Scikit-learn	6
2.6.2	Weka	6
3	Campagna sperimentale	7
3.1	Dataset	7
3.2	Misure di performance	7
3.3	BOW + Caratteristiche linguistiche	7
3.4	BERT + Caratteristiche linguistiche	7
3.5	SBERT + Caratteristiche linguistiche	7
3.6	Analisi lessico con PCA	7
4	Conclusioni e sviluppi futuri	8

Introduzione

onsiste nell'affermare il contrario di ciò che si pensa con lo scopo di ridicolizzare o sottolineare concetti provocando, a volte, una risata e finendo, in quei casi, nel sarcasmo, ma ha assunto anche significati più profondi. Di essi si possono definire tre accezioni:

Descrizione del problema

Il riconoscimento automatico dell'ironia nei contenuti generati da utenti, è uno dei compiti più complessi per quanto riguarda l'elaborazione del linguaggio naturale. Questo è molto importante per tutti i sistemi di sentiment analysis, in quanto facendo uso dell'ironia è possibile invertire completamente la polarità di una propria opinione, facendola passare da positiva a negativa e viceversa, penalizzando in questo modo le performance dei sistemi. Diventa pertanto cruciale sviluppare irony-aware sentiment analysis systems, ovvero sistemi che sono in grado di riconoscere questo fenomeno.

L'ironia è un tema studiato in diverse discipline, come la linguistica, filosofia e psicologia, ma è difficile da definire formalmente, soprattutto per questo motivo ne è difficile il riconoscimento. Tuttavia ci sono basi teoriche che suggeriscono il ruolo importante della sfera emotiva nell'uso dell'ironia, quindi un fattore chiave per riconoscerlo. Con questo si intende anche un uso indiretto e non esplicito del carico emotivo in ciò che si vuole comunicare.

I social network in generale, e twitter nello specifico, sono ampiamente utilizzati come fonte di informazione per sperimentare con modelli computazionali per il riconoscimento dell'ironia, essendo di fatti una grande risorsa per quanto riguarda i dati testuali generati da utenti.

CITAZIONE: 2 Paper

Approccio al problema

Si può considerare il riconoscimento dell'ironia come un problema di classificazione. Una frase potrà quindi essere classificata come appartenente alla classe ironica o non ironica.

Come già detto, twitter è una risorsa che fornisce moltissimi contenuti generati da utenti. Viene sfruttato questo aspetto creando dei modelli supervisionati di machine learning in grado di apprendere dai dati. A questo scopo, si estraggono varie caratteristiche dal testo, le quali permettono di distinguere le due classi. Le features e i modelli utilizzati saranno meglio spiegati nei capitoli successivi, ma tra tutte verranno usate caratteristiche relative alla sfera emotiva, così da cercare di migliorare le performance dei modelli.

Sintesi dei risultati

Devo mettere le tabelle i grafici? Quanto deve essere sintetico?

Chapter 1

Stato dell'arte

1.1 Approccio supervisionato

1.2 Approccio semi supervisionato

1.3 Approccio non supervisionato

Chapter 2

Sistema realizzato

In questo capitolo vengono spiegate ed analizzate tutte le operazioni di preprocessing e le features estratte dal testo, così da preparare il dataset per essere fornito in input ai vari classificatori nella fase di training.

2.1 Descrizione del sistema proposto

Grafico con workflow

2.2 Dataset

2.2.1 Metodi di labeling

Il dataset a disposizione con le relative etichette (ironico/non ironico) sono una parte cruciale per identificare l'ironia e quindi la costruzione dei modelli. Per etichettare i messaggi degli utenti si possono seguire due strade:

- Self-Tagging

Twitter mette a disposizione l'utilizzo degli hashtag nei messaggi. Assumendo che un utente che utilizza l'hashtag *#irony*, voglia esprimere ironia, è facile collezionare una serie di tweet etichettati come ironici.

- Crowdsourcing

I vari tweet vengono etichettati manualmente da alcune persone

2.2.2 Dataset utilizzato

Nel caso specifico, viene utilizzato il dataset *TwReyes2013*, composto da 40,000 tweet accumulati usando la tecnica self-tagging. Vengono quindi considerati 4 hashtag diversi:

Numero	Hashtag	Label assocaita
10,000	<i>#education</i>	non ironico
10,000	<i>#humor</i>	non ironico
10,000	<i>#politics</i>	non ironico
10,000	<i>#irony</i>	ironico

Table 2.1: Hashtag e label associate

2.3 Rappresentazione del testo

2.3.1 Rappresentazione booleana

2.3.2 Rappresentazione mediante transformer

2.4 Caratteristiche linguistiche

2.4.1 PP

2.4.2 POS

2.4.3 Onomatopeic

2.4.4 ...Tutto il resto

2.4.5 EMOT

2.5 Modelli supervisionati

2.5.1 Alberi di decisione

2.5.2 Support Vector Machine

2.5.3 Naive Bayes

2.5.4 Bayesian Network

2.6 Strumenti utilizzati

2.6.1 Scikit-learn

2.6.2 Weka

Chapter 3

Campagna sperimentale

3.1 Dataset

3.2 Misure di performance

3.3 BOW + Caratteristiche linguistiche

3.4 BERT + Caratteristiche linguistiche

3.5 SBERT + Caratteristiche linguistiche

3.6 Analisi lessico con PCA

Chapter 4

Conclusioni e sviluppi futuri