

**Management delle informazioni e
gestione della conoscenza
AA 2021-22**

Sentiment Analysis

Roberto Boselli
roberto.boselli@unimib.it

Sentiment Analysis, theory

Sentiment analysis : semantic (qualitative) analysis of messages to identify whether the opinion expressed towards a brand, a product or a service is positive, negative or neutral

Theories:

- ▶ The 5 Factors Model of Personality:
 - ▶ Extraversion (also often spelled extroversion), agreeableness, openness, conscientiousness, and neuroticism
- ▶ Emotions (pleasure, sadness, joy, anger, disgust, love, displeasure, fear, amazement)
 - ▶ e.g., Analisi della Felicità delle province italiane in Twitter (Voices from the blogs)

A limit is the difficulty of recognizing ambiguous attitudes, irony, sarcasm etc.



Why is sentiment analysis so important?

- ▶ The micro-blogging content coming from **Twitter** and **Facebook** poses serious challenges, not only because of the amount of data involved, but also because of the kind of language used in them to express sentiments, i.e., short forms, memes and emoticons
- ▶ Sentiment Analysis enables companies to make sense out of data. Thus they are able to elicit vital insights from a vast unstructured dataset without having to manually indulge with it



Why Sentiment Analysis?

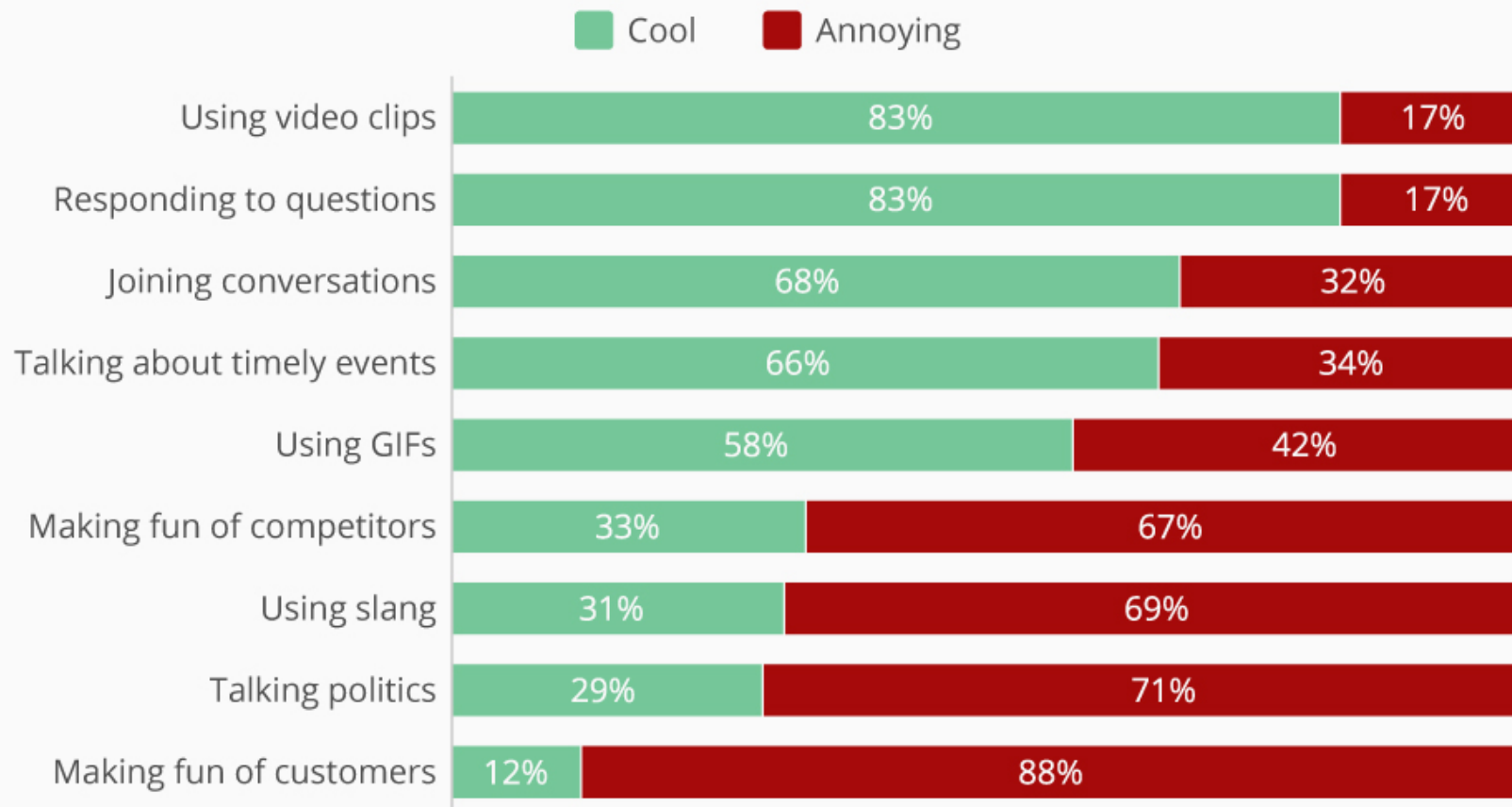
- ▶ Manage critical posts on social media
- ▶ Improve the Customer Experience
- ▶ Analyze product launches
- ▶ Assess the impact of sponsorships and CSR activities
- ▶ Discover new market trends
- ▶ Maintain the quality of the service on a national, international and global scale
- ▶ Monitor the popularity of Management
- ▶ ...



Consumer sentiment

What Brands Should and Shouldn't Do on Social Media

Consumer sentiment on brand behavior on social media channels



How does Sentiment work?

First approach:

1. **Detection and evaluation of keywords:**

A set of words identified as keywords in messages, a positive or negative value is assigned to the words. Each time the tool detects these keywords, it assigns a positive or negative value to the entire message

Many limits: **it is not possible to contextualize the message**, the level of accuracy of the analysis generally varies from **50% to 80%**, **sarcasm not detected**

=> Adding rules to create wider association patterns

- ▶ It provides valid results only if aimed at identifying the macro trend of highly discussed topics or if applied on a large scale



Example 1° approach

☰ Results

↶ Legend

I ¹ really enjoyed using the ¹ Canon Ixus in Madrid on March 4. The ² Panasonic Lumix ² is a bit disappointing, but the ³ Canon camera is ³ not bad at all. All I want when taking photos is point it and then just press the button. For only 200 dollars, a ⁴ really fair ⁴ price, this ⁵ camera is ⁵ perfect for me. Besides, I have had a ⁶ good ⁶ customer service experience.



How does Sentiment work?

Second approach

2. **Customized categories:**

- ▶ Users can establish categories to manually classify a few results, which constitute the **training set**, and the rules that algorithms will then have to follow
- ▶ It offers a higher level of precision of results
- ▶ Limits: requires a significant investment in time, and very strict parameters for the classification of results, this approach returns a limited number of results



Third Approach: Sentiment & AI

- ▶ Algorithms classify the contents on the base of entire sentence analysis, e.g. they are able to contextualize a tweet, a post or an article, and to accurately interpret the opinion of consumers
- ▶ They use **deep learning models** capable of simulating the cognitive functions of the human brain, technology is able to distinguish and understand complex linguistic structures, as well as entire sentences and simple forms of sarcasm and irony
- ▶ The accuracy of the results increases as well as the training set increases

Size of Training Set	Sentiment Accuracy
10K	55%
100K	62%
1M	71%
10M	84%
Multiple 10M	90%



Why is Sentiment Analysis a Hard to perform Task?

- ▶ Understanding emotions through text are not always easy. Sometimes even humans can get misled
- ▶ A text may contain multiple sentiments all at once
- ▶ Computers aren't too comfortable in comprehending Figurative Speech. Figurative language uses words in a way that deviates from their conventionally accepted definitions in order to convey a more complicated meaning or heightened effect. Use of similes, metaphors, hyperboles etc
- ▶ Heavy use of emoticons and slangs with sentiment values in social media texts like that of Twitter and Facebook also makes text analysis difficult



Other Issues

- ▶ A sentence containing positive or negative words could be neutral, that is not to express any opinion. In the questions or in the conditional sentences:
 - ▶ “Can you tell me which Sony camera is good?”
 - ▶ “If I can find a good camera in the shop, I will buy it”
- ▶ The use of sarcasm is difficult to grasp:
 - ▶ “What a great car! It stopped working in two days.”
- ▶ Some phrases do not have sentiment words but indicate anyway an implicit opinion:
 - ▶ “After two days of normal usage, the screen became black on the bottom”



Sentiment Analysis Tools

- ▶ The most common application of sentiment analysis is in consumer products and services reviews
- ▶ There are lexicon-based analysis methods, rule-based analysis methods, and machine learning techniques
- ▶ **Liu and Hu opinion lexicon** is a list of positive and negative words and is one of most used lexicons in sentiment tools. It contains around 6800 opinion words or sentiment words for English language. This list was composed over many years
- ▶ **VADER** is a **rule-based method** (next slides)
- ▶ These two tools are included in Orange DM



VADER Sentiment Analysis

- ▶ **VADER (Valence Aware Dictionary and sEntiment Reasoner)** is a lexicon and rule-based sentiment analysis tool that is specifically **attuned to sentiments expressed in social media**
- ▶ VADER uses a combination of a sentiment lexicon, is a list of lexical features (e.g., words) which are generally labelled according to their semantic orientation as either positive or negative
- ▶ VADER not only tells about the Positivity and Negativity score, but also tells us about how positive or negative a sentiment is

(Hutto, Gilbert, VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text, 2014)



VADER score

- ▶ VADER uses a vocabulary in text format inside which the recognized words are present, an average sentiment value that varies in the range $[-4; 4]$ and a standard deviation value based on how much the estimated sentiment of that word varies in the texts collected
- ▶ VADER provides four different elements of sentiment once inserted a word or phrase: the **negative** polarity, the **positive** polarity, the **neutrality** and a value called **compound** which calculates the sum of all the lexicon ratings which have been normalized between -1 (most extreme negative) and $+1$ (most extreme positive)



Compound

- ▶ The value of this element is calculated as:

$$\text{compound} = \frac{\text{score}}{\sqrt{\text{score}^2 + \alpha}}$$

- Where the score is the normalized sentiment value of a sentence and α a smoothing value which is selected by default at 15
- Essentially the compound value is a normalization of the score of a sentence that is calculated by the sum of the sentiments present in the vocabulary keeping in mind their standard deviation



Key points

- ▶ VADER analyses sentiments primarily based on certain key points:
 - ▶ **Punctuation:** the use of an exclamation mark (!) increases the magnitude of the intensity without modifying the semantic orientation
 - ▶ **Capitalization:** using upper case letters to emphasize a sentiment-relevant word in the presence of other non-capitalized words, increases the magnitude of the sentiment intensity
 - ▶ **Degree modifiers:** also called intensifiers, they impact the sentiment intensity by either increasing or decreasing the intensity (e.g. “extremely”, “marginally”)



Key points (2)

- ▶ **Conjunctions:** use of conjunctions like “but” signals a shift in sentiment polarity, with the sentiment of the text following the conjunction being dominant. “The food here is great, but the service is horrible” has mixed sentiment, with the latter half dictating the overall rating
- ▶ **Preceding Tri-gram:** by examining the tri-gram preceding a sentiment-laden lexical feature, we catch nearly 90% of cases where negation flips the polarity of the text. A negated sentence would be “The food here isn’t really all that great”
- ▶ VADER performs very well with emojis, slangs, and acronyms in sentences

I am 😊 today {'neg': 0.0, 'neu': 0.476, 'pos': 0.524, 'compound': 0.6705}



Comparison of VADER results: Orange vs. Python

Text	Orange				Python			
	pos	neg	neu	compound	pos	neg	neu	compound
The food here is good	0,42	0	0,58	0,4404	0,42	0	0,58	0,44
The food here is good!	0,444	0	0,556	0,4926	0,44	0	0,55	0,49
The food here is good!!	0,466	0	0,534	0,5399	0,46	0	0,53	0,53
The food here is good!!!	0,486	0	0,514	0,5826	0,48	0	0,51	0,58
The food here is GREAT!	0,562	0	0,438	0,729	0,56	0	0,43	0,72
The food here is great!	0,523	0	0,477	0,6588	0,52	0	0,47	0,65
The service here is extremely good	0,39	0	0,61	0,4927	0,39	0	0,61	0,49
The service here is marginally good	0,343	0	0,657	0,3832	0,34	0	0,65	0,38
The food here is great, but the service is horrible	0,167	0,31	0,523	-0,4939	0,16	0,31	0,52	-0,49
I am 😊 today	0	0	1	0	0,52	0	0,47	0,67
😊	0	0	0	0	0,66	0	0,33	0,71
😭	0	0	0	0	0,45	0,27	0,26	0,32
😞	0	0	1	0	0	0,7	0,29	-0,34
💕	0	0	0	0	0	0	1	0
Today SUX!	0	0,779	0,221	-0,5461	0	0,77	0,21	-0,54
Today only kinda sux! But I'll get by, lol	0,251	0,179	0,569	0,2228	0,31	0,12	0,55	0,52
Make sure you :) or :D today!	0,706	0	0,294	0,8633	0,7	0	0,29	0,86
The food here isn't really all that great	0,384	0	0,616	0,6557 /	/	/	/	/



Other tools

- ▶ **NLTK SentimentAnalyzer:** machine learning approach with several modules and functions using both Liu and Vader lexicons
- ▶ **TextBlob:** text processing Python library. The sentiment property returns a named tuple of the form Sentiment (polarity, subjectivity)
- ▶ **Stanford CoreNLP:** deep learning model computes the sentiment based on how individual words change the meaning of longer phrases
- ▶ **R packages:** Syuzhet (4 dictionaries), Rsentiment (sarcasm), Sentiment Analysis (customized dictionaries)
- ▶ **Google Cloud Prediction API:** machine learning models



Resources

- ▶ <https://github.com/cjhutto/vaderSentiment/blob/master/vaderSentiment/vaderSentiment.py>
- ▶ <http://www.nltk.org/api/nltk.sentiment.html>
- ▶ <http://textblob.readthedocs.io/en/dev/quickstart.html#sentiment-analysis>
- ▶ <https://nlp.stanford.edu/sentiment/index.html>
- ▶ <https://github.com/mjockers/syuzhet>
- ▶ <https://cran.r-project.org/web/packages/RSentiment/index.html>
- ▶ <https://github.com/sfeuerriegel/SentimentAnalysis>
- ▶ https://cloud.google.com/prediction/docs/sentiment_analysis

