# Sentiment Analysis on Online Automotive Forums

*Candidate*: Giuseppe Ravagnani
*Supervisor*: Prof. Nicola Ferro

September 30, 2019

# Outline

*Sentiment Analysis* is an active area of research in Natural Language Processing, motivated to improve the automated recognition of sentiment expressed in text
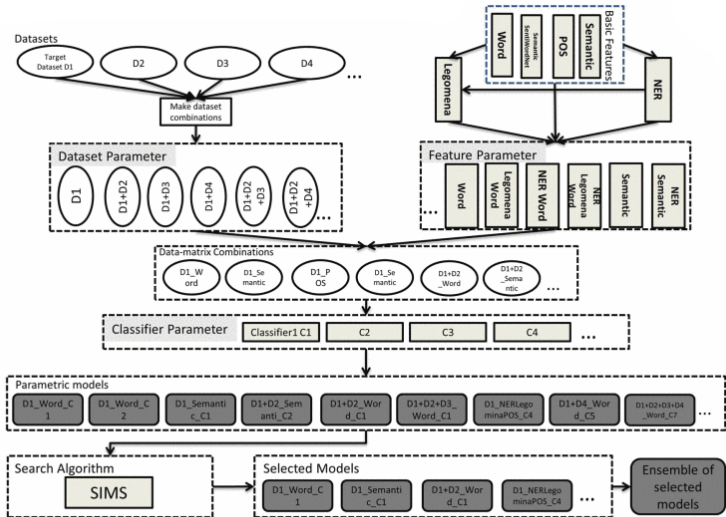
- From the whole scenic of online social media, Twitter is the most used to achieve contents

- Most existing techniques for sentiment classification involve supervised learning

- Used for collecting people's sentiment information about a given topic

$\rightarrow$ **Dataset**: the need of a set of already annotated data for applying machine learning models

$\rightarrow$ **Development**:
- Twitter Sentiment Analysis Algorithm
- Relevance Detection for class "Engine"
- Sentiment Classification for class "Engine"
- Cascade Classifier

$\rightarrow$ **Goal**: Develop an aspect-based sentiment analysis tool in automotive field for a specific brand

$\rightarrow$ **Collaboration**: Reply Technology for commission of Porsche

From: Zimbra, David & Abbasi, Ahmed & Zeng, Daniel. (2018).
The State-of-the-Art in Twitter Sentiment Analysis:
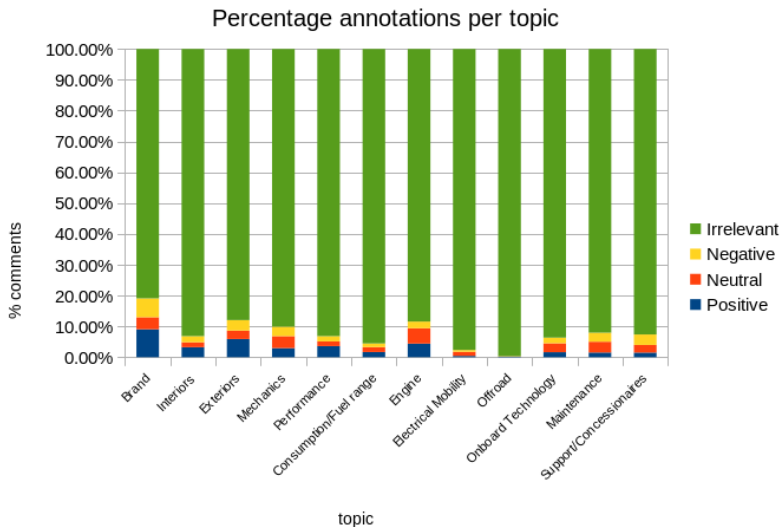A Review and Benchmark Evaluation.

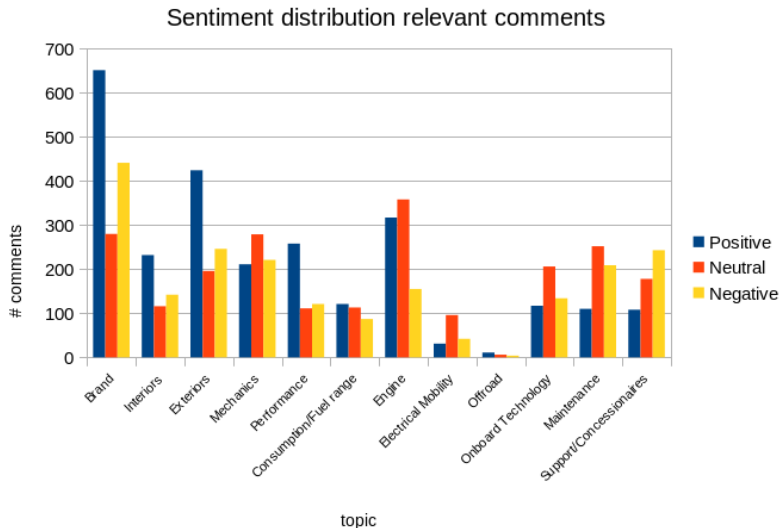| | Average | Pharma | Retail | Security | Tech | Telco | Ensemble |
|---|---|---|---|---|---|---|---|
| **BPEF** | 71.38 | 67.81 | 65.24 | 75.32 | 76.30 | 72.21 | yes |
| **NRC** | 71.33 | 75.26 | 64.93 | 76.39 | 64.96 | 75.08 | no |
| **Webis** | 71.41 | 76.16 | 64.40 | 77.37 | 63.68 | 75.46 | yes |

# Dataset Gathering

The creation of a suitable dataset is mandatory to train machine learning algorithm.

- Crawled some of the most visited Italian automotive forums (Quattroruote, Autopareri, Bmwpassion, HDmotori, Porschemania, Forumelettrico)
- Comments have been annotated with respect to the classes "Brand", "Interiors", "Exteriors", "Mechanics", "Performance", "Consumption", "Engine", "Electrical mobility", "Off-road" and "Technology", picking the labels "positive", "negative", "neutral" or "irrelevant".
- From a total of 1,200,000 crawled comments, 7,183 have been manually annotated

# Dataset Statistics



Percentage annotations per topic

# Dataset Statistics

Sentiment distribution relevant comments

# Comparison with Twitter

### Example block

@united I do not see where it talks about military baggage fees. Can you please guide me. Thanks #usairline

### Example block

Sono reali calcolati nel arco del tutto anno nel estate qualcosa in più causa gomme di 17" e climatizzatore nel inverno un po di meno. Per quanto riguarda le autostrade quelle che percorro io principalmente la A4 e molto congestionata cosi spesso la media e 110-115 km/h che ovviamente influisce positivamente a i consumi. Ma quello che mi piace di più è assenza dei guasti. Sulla vecchia Accord il primo guasto lo ho avuto a 200000 km si è rotto il termostato della clima. Ogni tanto faccio giro di altri forum e leggo delle turbine rotte catene di distribuzione progettate male iniettori fatti male mah nel 2015 per me sono le cose incomprensibili. Con tutti gli difetti che può avere preferisco la Honda.

$\rightarrow$ Twitter preprocessing:

1. lowercase, punctuation and stopwords removal
2. "http://someurl" $\rightarrow$ "URL"
3. "#hashtag" $\rightarrow$ "hashtag"
4. happy emoticons $\rightarrow$ "EMO_POS"
   sad emoticons $\rightarrow$ "EMO_NEG"
5. stemming

$\rightarrow$ Support Vector Machine (SVM) Classifier with Term Frequency - Inverse Document Frequency (TF-IDF) features vectorization

$$tf_{i,j} = \frac{n_{i,j}}{|d_j|}, \quad idf_i = \log \frac{|D|}{|\{d : i \in d\}|}, \quad tf\text{-}idf_{i,j} = tf_{i,j} \times idf_i$$

# Baseline Approach (2)

$\rightarrow$ Our dataset preprocessing:

1. encoding correction
2. lowercase, punctuation and stopwords removal
3. "http://someurl" $\rightarrow$ "URL"
4. replacing domain-specific tokens with common string (distances, speed, consumption, weight, power, ...)
5. stemming

$\rightarrow$ Support Vector Machine (SVM) Classifier with Term Frequency - Inverse Document Frequency (TF-IDF) features vectorization

$$tf_{i,j} = \frac{n_{i,j}}{|d_j|}, \quad idf_i = \log \frac{|D|}{|\{d : i \in d\}|}, \quad tf\text{-}idf_{i,j} = tf_{i,j} \times idf_i$$
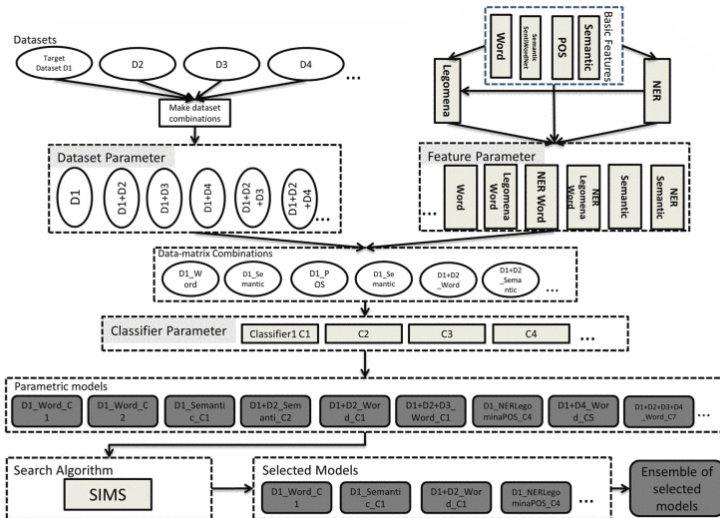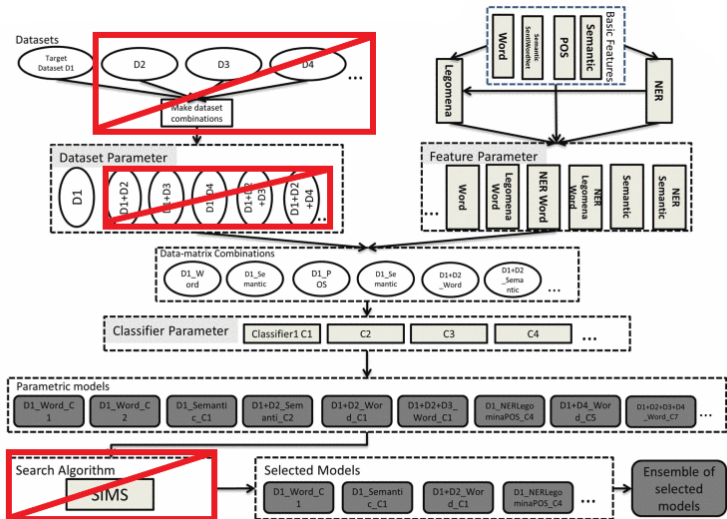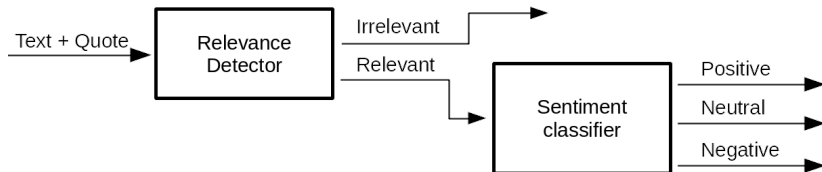
Implementation of a cascade classifier for the four-label classification.



- Logistic Regression relevance detector
- BPEF Sentiment classifier

Baseline

BPEF

**Predicted value**

| | Positive | Neutral | Negative |
|---|---|---|---|
| Positive | 235 | 44 | 21 |
| Neutral | 36 | 196 | 68 |
| Negative | 27 | 42 | 231 |

Actual value

| **F-macro** | 0.735 |
|---|---|
| **Accuracy** | 0.736 |

**Predicted value**

| | Positive | Neutral | Negative |
|---|---|---|---|
| Positive | 230 | 55 | 15 |
| Neutral | 39 | 221 | 40 |
| Negative | 20 | 72 | 208 |

Actual value

| **F-macro** | 0.734 |
|---|---|
| **Accuracy** | 0.732 |

**Relevance detection**

SVM

Logistic Regression

**Predicted value**

Irrelevant    Relevant

| | Irrelevant | Relevant |
|---|---|---|
| **Irrelevant** | 962 | 55 |
| **Relevant** | 60 | 73 |

**Actual value**

| | |
|---|---|
| **F1-macro** | 0.559 |
| **Recall** | 0.570 |
| **Precision** | 0.549 |

**Predicted value**

Irrelevant    Relevant

| | Irrelevant | Relevant |
|---|---|---|
| **Irrelevant** | 933 | 84 |
| **Relevant** | 50 | 83 |

**Actual value**

| | |
|---|---|
| **F1-macro** | 0.553 |
| **Recall** | 0.624 |
| **Precision** | 0.497 |

**Sentiment classification**

SVM

BPEF



**Predicted value**

|  | Positive | Neutral | Negative |
|---|---|---|---|
| Positive | 35 | 15 | 1 |
| Neutral | 24 | 33 | 0 |
| Negative | 13 | 9 | 3 |

Actual value

| **F1-macro** | 0.451 |
|---|---|



**Predicted value**

|  | Positive | Neutral | Negative |
|---|---|---|---|
| Positive | 29 | 22 | 0 |
| Neutral | 15 | 42 | 0 |
| Negative | 7 | 15 | 3 |

Actual value

| **F1-macro** | 0.467 |
|---|---|

## 4-labels classification

SVM

Cascade classifier



**Predicted value**

| | Irrelevant | Positive | Neutral | Negative |
|---|---|---|---|---|
| Irrelevant | 1000 | 11 | 6 | 0 |
| Positive | 32 | 15 | 4 | 0 |
| Neutral | 38 | 10 | 9 | 0 |
| Negative | 20 | 3 | 2 | 0 |

| **F1-macro** | 0.378 |
|---|---|



**Predicted value**

| | Irrelevant | Positive | Neutral | Negative |
|---|---|---|---|---|
| Irrelevant | 933 | 18 | 61 | 5 |
| Positive | 19 | 23 | 9 | 0 |
| Neutral | 20 | 4 | 33 | 0 |
| Negative | 11 | 3 | 3 | 8 |

| **F1-macro** | 0.556 |
|---|---|

New data involving some other information have been crawled for test a use case of the classifier.

# Conclusions

- BPEF model overcomes baseline approach
- Implemented model can be considered reliable for sentiment classification for Italian automotive forums
- Expanding the dataset, same algorithms should improve their scores
- Design a more sophisticated relevance detector for identifying the topic
- Integration in a production system: scheduled crawled, database and business intelligence software