

**U. PORTO**

**U.  
FLUP**

**CLUP** Centro de  
Linguística da  
Universidade do  
Porto

**UI** UNIVERSIDADE  
BEIRA INTERIOR

**INESCTEC**



# Extending General Sentiment Lexicon to Specific Domains in (Semi-) Automatic Manner

**Pavel Brazdil, Purificação Silvano, Fátima Silva,  
Shamsuddeen Hassan, Fátima Oliveira,  
João Cordeiro, António Leal**

**Workshop SALLD-1@ LDK 2021**

**1 Sept 2021**

# Research Areas, People & Institutions

## AI, ML, Text Mining Sentiment Analysis

- Pavel Brazdil
- Shamsuddeen Hassan (PhD student)
- João Cordeiro

 U. PORTO

 UNIVERSIDADE  
BEIRA INTERIOR

 INESCTEC

 LIAAD  
CENTRO DE LINGUÍSTICA  
APLICADA DA UNIVERSIDADE  
DO PORTO  
INESCTEC

## Linguistics Sentiment Analysis

- Purificação Silvano
- Fátima Silva
- Fátima Oliveira
- António Leal

 U. PORTO

 U.  
FLUP

 CLUP  
Centro de  
Linguística da  
Universidade do  
Porto

Acknowledgements: This work is financed by National Funds through the Portuguese funding agency, *FCT - Fundação para a Ciência e a Tecnologia*, within project UIDB/50014/2020 of INESC TEC and project UIDB/00022/2020 of CLUP.

# Overview

1. Introduction (4-10)
2. Automatic generation of sentiment lexicon (11-18)
3. Experimental results (19)
4. Short phrases and modifier patterns (20-28)
5. Applying modifier patterns and results (29)
6. Conclusions (30-32)

# 1. Introduction (1)

## Our interests:

**Sentiment analysis in specific domains and for specific languages**  
e.g. economics and finance;  
texts in European Portuguese.

Approaches that use **(semi-)automatic methods** for  
the construction of the **domain-sensitive (DS) sentiment lexicon**

## Our aim:

Conceive solutions relatively **easy to explain to the user**  
(i.e. why certain polarity/rating is attributed to the given text)

Lexicon-based approaches are of interest  
as performance is often not the only important criterion

# 1. Introduction (2)

Sentiment analysis (SA) is an important area (product reviews, policy reviews, etc.)

The aim is to predict the sentiment value for a given text (to see what the user thinks about a given product or policy)

## The sentiment values can be

- Categorical (positive, neutral, negative), or
- Numeric (e.g. on scale -3 to 3).

We prefer this option, as it enables to explore more information

# 1. Introduction (3)

Many diverse approaches to SA exist:

- **Manual** annotation of texts (by experts, crowdsourcing), (requires manual effort)
- **Lexicon-based** approaches,
- **Machine learning (ML)** approaches,
  - Classical ML approaches, (e.g. decision trees, random forests, etc.),
  - Deep neural networks (e.g. CNN).

## 1. Introduction (4)

Atteweldt et al. (2021) compared various approaches and showed: DeepNN obtained best performance although they are behind the manual approach.

However, is not easy to see why certain predictions were made.

Lexicon-based approaches offer better explainability, but tend to have rather poor performance in specific domains (e.g., economics and finance, bioinformatics etc.), when a general-purpose lexicon is used.

Hence, we need a domain-specific (DS) lexicon (or extension).

# 1. Introduction (5)

## Strategies for developing DS lexicon:

- Manual - we exclude this, as this requires a lot of manual effort
- Automatic – followed here

## Automatic method for developing DS lexicon for unigrams

Our approach is based on

Almatarneh et al. (2017) and Muhammad et al. (2020).

It requires labelled texts (with numeric ratings).

More details on the method are given in Section 2.



# 1. Introduction (6)

Lexicon that includes single words only (unigrams)  
is insufficient to cover all situations (e.g. phrase “is not good”)

=> **We need to be able to deal with (some) short phrases.**

Alternative solutions:

- **Enumerating** various short phrases and storing their sentiment value in the lexicon.
- **Deriving the sentiment value** of some short phrases using rules (e.g. “*crescimento alto*” (*high growth*))  
This is applicable to short phrases satisfying certain “patterns”
- **Hybrid approach** combining the two approaches above

There can be too many!

Followed here

# 1. Introduction (7)

## Short phrases considered:

We focus on certain type of short phrases that involve intensification, downtoning/attenuation, reversal/inversion.

They are represented with the help of ***modifier patterns***, which are:

- Applicable to various domains;
- Constitute generally useful linguistic knowledge;
- Acquired manually, but **exploit DS lexicon** (acquired automatically);
- The sentiment value of short phrases is obtained in an automatic way.

## 2. Automatic Generation of DS Sentiment Lexicon

The methodology includes:

1. Corpus preparation and annotation
2. Preprocessing
3. Generating the distributions of occurrence of words
4. Generating the sentiment values
5. Combining domain-specific lexicon with general purpose lexicon

## 2.1 Corpus preparation and annotation (1)

Four linguists prepared 23 (1+22) texts from articles on **finance and economy** from different Portuguese online newspapers.

Each text contained between 2 to 30 sentences (or their fractions) (the total was 408).

At least two annotators annotated each sentence (or phrase) with the sentiment value on the **scale of -3 to 3**.


Examples:

Doc	No	Frase	SVal
1	1	Primeiro porque quem está habituado a lidar com a exportação de serviços sabe que a <b>falta de qualificação</b> dos portugueses é uma <b>falsa questão</b> .	1
1	2	porque não só o trabalho dos portugueses se <b>vende</b> como <b>nunca</b> , como também diversas mega empresas europeias estão a mudar para cá os seus <b>serviços mais sofisticados</b>	2
1	3	O <b>saldo positivo</b> das nossas trocas <b>compensa largamente</b> o financiamento das atividades do país.	2

## 2.1 Corpus preparation and annotation (2)

Number of occurrences of different ratings in the corpus:

Rating	-3	-2	-1	1	2	3	Total
Nº of occurrences	21	100	147	81	52	7	408



Somewhat unbalanced data:  
More negative ratings than positive ones.  
Later we show how we deal with this.

A part of this data was used to construct the lexicon;  
another part to generate the predictions and evaluate them.  
(More details later).

## 2.2 Preprocessing

The data was read-in with *Quanteda* package of R.

Further processing was done with *udpipe* package:

POS tagging, lemmatization

Not all lexical categories include sentiment bearing words  
(Martin and White, 2005; Taboada et al., 2011; Liu, 2012, etc.)

So, we focus on four categories of unigrams only:

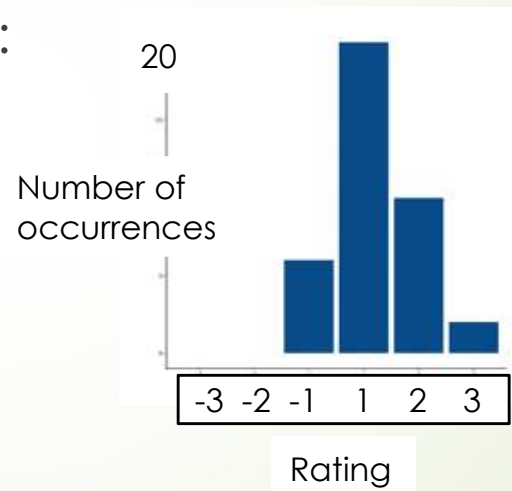
***nouns, verbs, adjectives, adverbs.***

## 2.3 Generating the distributions of occurrences

Consider each word (token) in the labelled sentences in the training data.

Examine the number of occurrences for the given ratings

Ex. word (token) “*bom*” (good):



The occurrences are transformed into probabilities (with Laplace smoothing).

## 2.4 Generating the Sentiment Value

The sentiment value is derived from the distribution of probabilities and ratings -3, -2, -1, 1, 2, 3

$$SV_{ti} = \sum P''_{ti} * W$$

Examples of some entries induced:

<u>Ecolex</u>	
<u>Word</u>	<u>Sent.Val.</u>
<u>investimento</u> (investment)	<b>1.000</b>
<u>respeito</u> (respect)	<b>1.000</b>
<u>importante</u> (important)	<b>0.955</b>
<u>crescer</u> (grow)	<b>0.931</b>
<u>bom</u> (good)	<b>0.909</b>

Comparison with one existing sentiment lexicon:

<u>Sentilex-PT</u>	
<u>Word/Idiom</u>	<u>Sent.Val.</u>
<u>investimento</u>	-
<u>faltar ao respeito</u>	-1
<u>importante</u>	-
<u>crescer</u>	-
<u>bom</u>	1



## 2.5 Combining DS+GP Lexicons and Evaluation

### Combination of lexicons:

- Use preferentially the DS lexicon (Ecolex)
- Add entries from the GP lexicon (Sentilex-PT) (those that do not appear in Ecolex)

## 2.5 Combining DS+GP Lexicons and Evaluation

### Evaluation using 5-fold cross-validation

Repeat 5 times:

- Use 4 partitions (folds) as the “train data” to construct the lexicon.
- Balance the training data by adding some duplicates of positively rated sentences (or their fractions).  
This way we obtained approx. 318 cases + 100 duplicates
- Use 1 partition (fold) to as test data ( $\approx 81$  cases).

### 3 Experimental Results

The accuracy of default system (random prediction) is 50% for balance data.

	Accuracy (%)					
	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean
<u>Ecolex</u>	60.37	54.21	65.71	47.06	55.66	56.60
<u>Sentilex-PT</u>	62.26	39.25	65.71	<b>70.58</b>	<b>58.49</b>	59.60
<u>Ecolex+Sentilex</u>	<b>76.41</b>	<b>77.57</b>	<b>73.33</b>	66.67	55.66	<b>69.93</b>

The combined lexicon (Ecolex+Sentilex) is better than the constituents (about 10% improvement of accuracy)!

## 4 Short phrases and modifier patterns

We focus on certain short phrases only that include:

- **intensification** ex. *melhorar muito* (improve **greatly**)
- **downtoning/attenuation** ex. *melhorar pouco* (improve **a little**)
- **reversal/inversion** ex. *não é bom* (is **not** good)

modifier M

focal element F

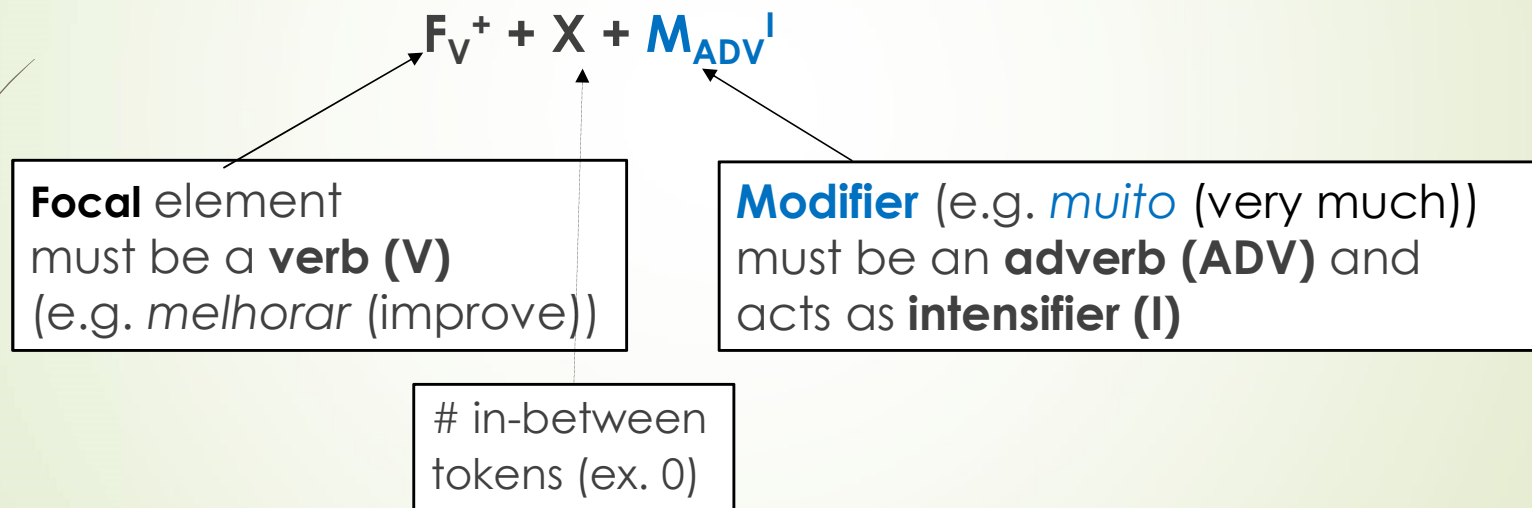
Modifier and focal element form part of **modifier patterns**.

These can be transformed into rules and used to derive the sentiment value of the phrase.

## 4.1 Intensification (1)

Intensification includes **modifier M** that **increases** the sentiment value of the focal element F

Example of a *modifier pattern*:



## 4.1 Intensification (2)

More details

$F_V^+ + X + M_{ADV}^I$

No need to specify the word  
(just its class, e.g. "v")

Admissible adverbs are given.  
The list was elaborated manually.  
Currently it has 26 elements:

- *muito* (very much),
- *bastante* (rather a lot),
- *mais* (more),
- *poderosamente* (with a great power),
- *muitíssimo* (very, very much),
- *incrivelmente* (incredibly),
- *extraordinariamente* (extraordinarily),
- ...

## 4.1 Intensification (3)

Determining the sentiment value of modifier patterns

$$SV(F_V^+ + X + M_{ADV}^I) = C_{MI} * SV(F_V^+)$$

Sentiment value of  
the modifier pattern

Sentiment value of the  
focal element  
(retrieved from the lexicon)

Constant (e.g. 2)

$$\begin{aligned} \text{Ex. } SV(\text{"melhorar"}, \text{"muito"}) &= 2 * SV(\text{"melhorar"}) \\ &= 2 * 0.6 = 1.2 \end{aligned}$$

## 4.1 Intensification (4)

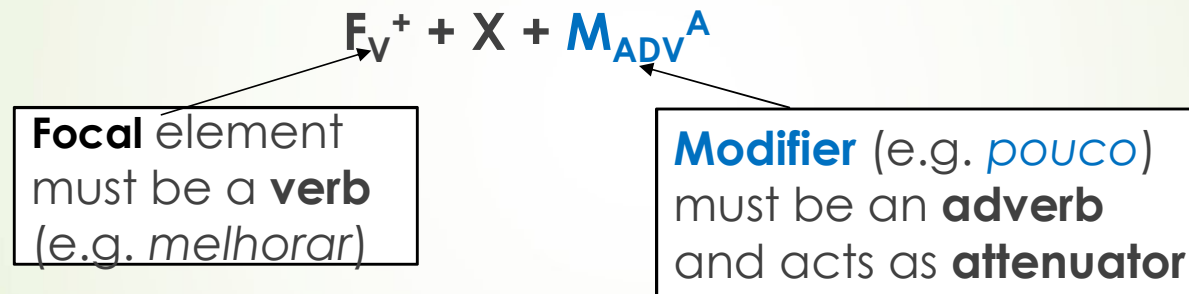
<b>Modifier pattern</b>	<b>F (domain-specific)</b>	<b>M</b>
$F_V^+ + X + M_{ADV}^I$	compensar, melhorar, crescer	muito, tanto, largamente, como nunca
$F_N + X + M_{ADJ}^I$	saldo, números de economia	positivo, lisonjeiros
$F_N + X + M_{ADJ}^I$	portugueses	produtivos
$F_V^+ + X + M_{QUANT}^I$	crescer	mais que X
$F_V^- + X + M_{QUANT}^I$	perder	mais que X
$F_N^+ + X + M_V^I$	crescimento, exportações	suplantou, evoluíram
<b>Modifier pattern</b>	<b>M</b>	<b>F</b>
$M_{ADV}^I + X + F_{ADJ}^+$	muito, bastante, mais	interessante, lisonjeiro, sofisticado
$M_{ADV}^I + X + F_{ADJ}^-$	muito, bastante, mais	negativo
$M_{ADJ}^I + X + F_N^+$	maior	credibilidade
$M_V^I + X + F_N^+$	aumentaram	qualificação
$M_V^I + X + F_N^+$	aumentará	economia
$M_N^I + X + F_{NP}^+$	o acelerar	crescimento económico



## 4.2 Downtoning/Attenuation (1)

Downtoning/Attenuation Includes **modifier M** that **decreases** the sentiment value of the focal element F

Example of a *modifier pattern*:



Determining the sentiment value:

$$SV(F_{V^+} + X + M_{ADV^A}) = C_{MA} * SV(F_{V^+})$$

Constant (e.g. 0.5)

## 4.2 Downtoning/Attenuation (2)

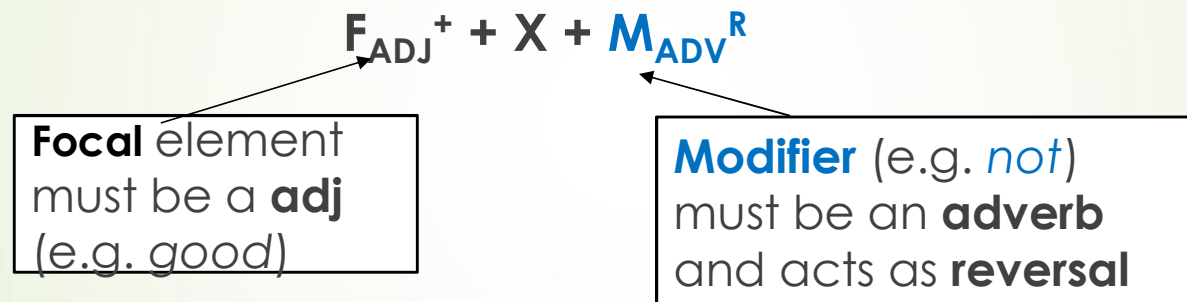
Our study includes various modifier patterns for downtoning/attenuation, e.g.:

Modifier pattern	F (domain-specific)	M
$F_N^+ + X + M_{ADJ}^A$	saldo	negativo
$F_N^- + X + M_{ADJ}^A$	déficit	controlado
$F_N^+ + X + M_{PREP}^A$	crescimento	apenas X
$F_N^- + X + M_{PREP}^A$	desemprego	abaixo de X
	<b>M</b>	<b>F</b>
$M_{ADV}^A + X + F_V$	pouco	alterou
$M_{PREP}^A + X + F_{ADJ}^-$	apesar	negativo
$M_N^A + X + F_N^+$	falta	qualificação

## 4.3 Reversal/Inversion (1)

Reversal/Inversion includes **modifier M** that **inverts** the sentiment value of the focal element F

Example of a *modifier pattern*:



Determining the sentiment value:

$$SV(F_{ADJ}^+ + X + M_{ADV}^R) = \bar{c}_{MR} * SV(F_{ADJ}^+)$$

Constant (e.g. 1)

## 4.3 Reversal/Inversion (2)

Our study includes various modifier patterns for reversal/inversion, e.g.:

Modifier pattern	M	F
$M_{ADV}^R + X + F_N^+$	não	solução
$M_{ADV}^R + X + F_{ADJ}^-$	não, nem, tudo menos	desajustado, mau, problemático
$M_V^R + X + F_{NP}^-$	inverta	ciclo negativo
	F	M
$F_N^- + X + M_{ADJ}^R$	falta de $N^+$	disparatado
$F_N + X + M_{NP}^R$	$F_N$	falsa questão

## 5 Applying Modifier Patterns and Results (1)

So far, we have conducted experiments with some modifier patterns only.

So far, we have noted a modest improvement.

We plan to complete this study shortly.

We expect that we can gain further 5-10% on accuracy.

## 6 Conclusions - Note on *Linguistic Linked Data*

The topic is of interests to us, because we are:

- interested in collaborating with others on different SA approaches,
- willing to share/exchange data.

Currently, our data is stored in CSV format, so it can be interchanged.

This applies both to:

- labelled texts with sentiment score,
- modifier patterns (these are transformed into rules by our program).

Transforming this data into other formats could be done.

## 6 Conclusions – Contributions of this Work

We have shown a method for **automatic construction of a domain-specific (DS) lexicon** (requires modest number of labelled examples)

It is advantageous to use **labelled data with numeric values**, rather than with just categorical labels.

It is useful to use a **combined DS + GP lexicon**, i.e., combine the induced DS lexicon with an existing general-purpose (GP) lexicon.

Experiments show that this is a promising line to follow. We have obtained about **10% increase in accuracy** when compared with the existing GP lexicon.

## 6 Conclusions – Contributions of this Work

We have designed a method (based on previous work) for **representing and applying modifier patterns** (for intensification, downtoning/attenuation, reversal/inversion).

We have shown how these can be combined with the lexicon-based approach.

Experiments are in progress and we expect that we can obtain substantial gains in accuracy (5-10%).

### **Additional advantage – explainability:**

The predictions of sentiment values can be justified to the user (i.e., we can show how they were derived).



## References (selection)

- Almatarneh, S., Gamallo, P. (2017). Automatic construction of domain-specific sentiment lexicons for polarity classification. In: *Advances in Intelligent Systems and Computing*, pp. 175–182.
- van Atteweld, W., van der Velden, M. A. C. G., Bokes, M. (2021). *The Validity of Sentiment Analysis: Comparing Manual Annotation, Crowd-Coding, Dictionary Approaches, and Machine Learning Algorithms*, *Communication Methods and Measure*
- Forte, A. C., Brazdil, P. (2016). Determining the Level of Clients' Dissatisfaction from their Commentaries, *Proc. PROPOR 2016*, Tomar, Springer, Portugal, 74-85
- Liu, (2012). *Sentiment Analysis and Opinion Mining: Synthesis Lectures on Human Language Technologies*. California: Morgan & Claypool Publishers.
- Martin, J.R., White, P.R.R. (2005). *The Language of Evaluation*. New York: Palgrave.
- Muhammad, S. H., Brazdil, P., Jorge, A. (2020). Incremental Approach for Automatic Generation of Domain-Specific Sentiment Lexicon. In: *Advances in Information Retrieval*, LNCS, vol. 12036. *Proc. ECIR 2020*, 619-623.

## References (selection)

- Silva, F., Silvano, P., Leal, A., Oliveira, F., Brazdil, P., Cordeiro, C., Oliveira, D. (2018). Análise de sentimento em artigos de opinião, *Linguística: Revista de estudos linguísticos da Universidade do Porto*, 13, 79-114.
- Trnavac, R., Das, D., Taboada, M. (2016). Discourse relations and evaluation, *Corpora* 11 (2), 169-190.
- Taboada, M., Brooke, J., Tofiloski, M., Voll, K., Stede, M. (2011). Lexicon-Based Methods for Sentiment Analysis. *Association for Computational Linguistics*. 37(2): 267-307

**U. PORTO**

**U.  
FLUP**

**CLUP** Centro de  
Linguística da  
Universidade do  
Porto

**UI** UNIVERSIDADE  
BEIRA INTERIOR

**INESCTEC**



# Extending General Sentiment Lexicon to Specific Domains in (Semi-) Automatic Manner

**Pavel Brazdil, Purificação Silvano, Fátima Silva,  
Shamsuddeen Hassan, Fátima Oliveira,  
João Cordeiro, António Leal**

**Workshop SALLD-1@ LDK 2021**

**1 Sept 2021**