

# Agregadores de Classificadores para Análise de Sentimentos

Alison Pereira Ribeiro  
Prof<sup>a</sup> Dr<sup>a</sup> Nádia F. F. da Silva

Universidade Federal de Goiás  
Instituto de Informática – INF  
Goiânia - Goiás - Brasil

*alisonrib17@gmail.com, nadia@inf.ufg.br*

Dezembro de 2018

- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

# 1 Introdução

## 2 Objetivos

## 3 Metodologia

- Pré-processamento
- Representação vetorial dos tweets
- Abordagens
- Algoritmos

## 4 Datasets

- Sanders
- HCR
- SemEval-2018

## 5 Aplicações dos Métodos

## 6 Resultados

## 7 Conclusões

## 8 Trabalhos Futuros

## 9 Referências

- O que é Análise de Sentimentos?
- Quais os desafios da área?
- Como buscar soluções para os problemas da AS?

## Exemplo 1

- "Vingadores Guerra Infinita foi incrível! Obrigado Marvel! :D", *positivo*.
- "Eu ainda não assisti a nova temporada de *Black Mirror*.", *neutro*.
- "As propostas desses candidatos são péssimas, nunca vi pior! Estamos perdidos!", *negativo*.

- 1 Introdução
- 2 **Objetivos**
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

## Objetivo Geral

Estudar e desenvolver modelos de Inteligência Artificial para Análise de Sentimentos para aplicações reais.

## Objetivo Geral

Estudar e desenvolver modelos de Inteligência Artificial para Análise de Sentimentos para aplicações reais.

## Objetivos Específicos

- Explorar métodos de Análise de Sentimentos;
- Implementar abordagens que aproveitem a variedade dos métodos;
- Divulgar os resultados obtidos.



- 1 Introdução
- 2 Objetivos
- 3 Metodologia**
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

- Pré-processamento;
- Representação vetorial dos *tweets*;
- Abordagens;
- Algoritmos.

- Remoção de *links*;
- Remoção de números;
- Remoção de caracteres especiais;
- *Stop words*;
- *Lowercase*;
- *Stemming*.

## Exemplo 3.1

- "Eu odeio meu computador da @apple. Foi 3500 dólares pelo ralo", *negativo*

## Remoção de números

- "Eu odeio meu computador da @apple. Foi dólares pelo ralo", *negativo*

## Remoção de caracteres especiais

- "Eu odeio meu computador da apple Foi dólares pelo ralo", *negativo*

## *Stop words*

- "Eu odeio computador *apple* dólares ralo", *negativo*

## *Lowercase*

- "eu odeio computador *apple* dólares ralo", *negativo*

## *Stemming*

- Reduz uma palavra em sua forma canônica, em seu morfema;
- Por exemplo casa, casas, casinhas e casarão resultam no mesmo morfema: cas.

- *Bag-of-Words*:
  - Verificação de ocorrência;
  - TF-IDF.
- *Word Embeddings* [Bengio et al. 2003];
- n-gramas:
  - unigrama;
  - bigrama;
  - unigrama + bigrama.

# Representação vetorial dos tweets

## Exemplo 3.2

- *Tweet1*: "Não gosto de *smartphone* muito grande", *negativo*.
- *Tweet2*: "Gostei desse *smartphone* azul", *positivo*.

## Representação

	não	gosto	gostei	smartphone	grande	muito	azul
<i>Tweet1</i>	1	1	0	1	1	1	0
<i>Tweet2</i>	0	0	1	1	0	0	1

**Tabela:** Representação vetorial dos *tweets*.

- Aprendizado de Máquina;
- Dicionários Léxicos;
- *Emoticons*;
- *Part-of-Speech*;
- Combinações de vários métodos.



- *Multinomial Naive Bayes* [Da Silva et al. 2014];
- *Support Vector Machine* [Rosenthal et al. 2014];
- *Random Forest* [Saleiro et al. 2017];
- *Logistic Regression* [Mittal et al. 2012];
- *Ensembles* [Fouad et al. 2018].

- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets**
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

- Sanders;
- HCR;
- SemEval-2018.

## Sanders [Sanders, 2011]

- Possui 5.513 *tweets* classificados como positivos, negativos, neutros e irrelevantes;
- *Tweets* irrelevantes foram desconsiderados;
- Base coletada a partir de quatro tópicos: *@apple*, *#google*, *#microsoft* e *#twitter*.

## HCR [Speriosu et al. 2011]

- Construído a partir da *hashtag* #hcr;
- Dados de treinamento, desenvolvimento e teste;
- *Tweets* classificados como positivos, negativos, neutro e irrelevantes;
- *Tweets* irrelevantes foram desconsiderados.

## SemEval-2018 [Barbieri et al. 2018]

- *Tweets* coletados com a API do Twitter;
- Geolocalização do *tweets* é no EUA;
- 500 mil tweets para treinamento, 50 mil para desenvolvimento e 50 mil para teste;
- Problema multiclasse: 20 classes (*emojis*).

- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

- Aprendizado de Máquina;
- Dicionários léxicos:
  - Opinion Lexicon;
  - SenticNet;
  - SemEval-2015 Lexicon.
- *Emoticons*;
- *Part-of-Speech*;
- Dicionário léxico + *Emoticons*;
- Combinação dos três dicionários léxicos;
- Dicionário léxico + *Emoticons* + *Part-of-Speech*.



## Aprendizado de Máquina (*Machine Learning*)

- No método baseado em ML, representa-se os *tweets* através do *Bag-of-Words* (ou *Word Embeddings*), utilizando *n*-gramas;
- Posteriormente aplica-se algum algoritmo de classificação.

	termo <sub>1</sub>	termo <sub>2</sub>	...	termo <sub>n</sub>
<i>tweet</i> <sub>1</sub>	1	0	...	0
<i>tweet</i> <sub>2</sub>	0	0	...	1
<i>tweet</i> <sub>3</sub>	0	1	...	0
...	...	...	...	...
<i>tweet</i> <sub>n</sub>	1	0	...	0

**Tabela:** Representação dos *tweets* com método de ML.

## Opinion Lexicon

- Possui 4.783 léxicos positivos 2.006 negativos;
- Estratégia proposta por [Mohammad et al. 2013];
- Conta as palavras do dicionário;
- Se o número de termos positivos for maior que negativos, então o *tweet* é positivo;
- Caso contrário, o *tweet* é negativo;
- No caso de empate, o *tweet* é neutro.

## SenticNet

- Conta com 50.000 palavras classificadas como positivas e negativas;
- Estratégia proposta por [Mohammad et al. 2013];
- Conta as palavras do dicionário;
- Se o número de termos positivos for maior que negativos, então o *tweet* é positivo;
- Caso contrário, o *tweet* é negativo;
- No caso de empate, o *tweet* é neutro.

## SemEval-2015 Lexicon

- Possui 1515 termos de sentimentos;
- Cada palavra possui uma pontuação (número real);
- Soma a pontuação das palavras encontradas nos *tweets*:

$$n = \sum_{i=1}^N K_i \quad (1)$$

- Rotulação dos *tweets* por meio da proposta apresentada:

$$x = \begin{cases} \textit{positivo}, & \textit{se } n > 0 \\ \textit{negativo}, & \textit{se } n < 0 \\ \textit{neutro}, & \textit{se } n = 0 \end{cases} \quad (2)$$

## Part-of-Speech (POS)

- Categoriza cada palavra na respectiva classe sintática, como: verbo, pronome, advérbio, entre outros;
- Pacote utilizado de Stanford [Manning et al. 2014];
- Rotula as palavras dos *tweets* e conta o número de *tags*;

	CC	JJ	VB	...	NN	classe
<i>tweet<sub>1</sub></i>	0	2	0	...	0	positivo
<i>tweet<sub>2</sub></i>	1	0	3	...	0	negativo
<i>tweet<sub>3</sub></i>	0	2	0	...	0	neutro
...	...	...	...	...	...	...
<i>tweet<sub>n</sub></i>	1	0	2	...	0	positivo

**Tabela:** Representação dos *tweets* com método de POS.

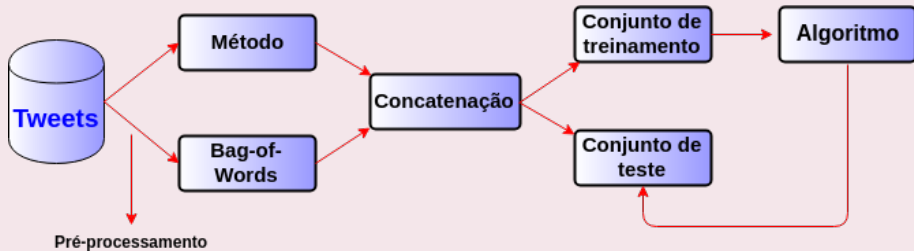
## *Emoticons*

- Emoticon Sentiment Lexicon;
- Contém 476 *emoticons*, distribuídos em 179 com sentimento positivo, 278 com sentimento negativo e 20 com sentimento neutro;
- Se é encontrado um *emoticon* positivo no *tweet*, então o *tweet* é positivo, essa lógica se repete no caso de encontrar um *emoticon* negativo ou neutro;
- Nesse método pode haver que um *tweet* não tenha *emoticon*.

## Combinações dos métodos

- Mantidas as mesmas abordagens descritas anteriormente;
- Dicionário léxico + *Emoticons*;
- Combinação dos três dicionários léxicos;
- Dicionário léxico + *Emoticons* + *Part-of-Speech*.

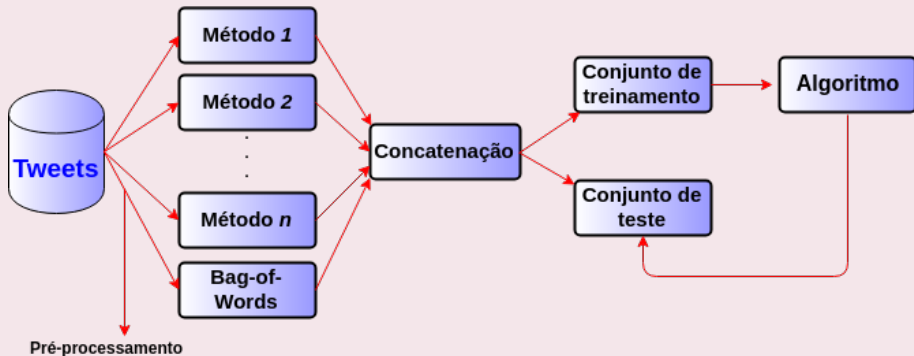
## Diagrama de Aplicação



**Figura:** Diagrama de aplicação dos métodos, individualmente.

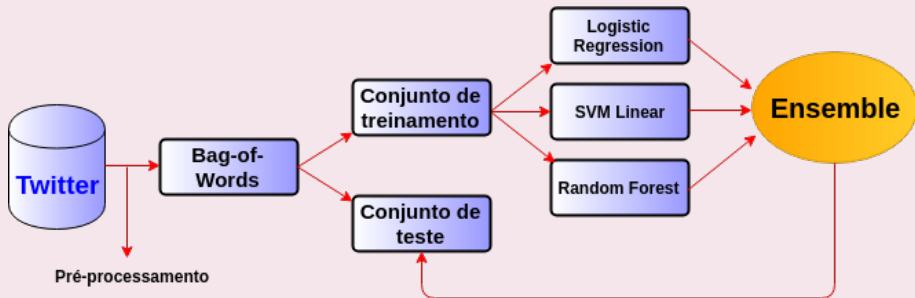


## Diagrama de Aplicação



**Figura:** Diagrama de aplicação dos métodos combinados.

## Diagrama de Aplicação



**Figura:** Diagrama de aplicação dos métodos em agregadores de classificadores.

- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados**
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

## Resultados - Sanders

Método	Algoritmo	Acc. (%)
Aprendizado de Máquina	LR	78,10
<i>Part-of-Speech (POS)</i>	LR	77,40
Opinion Lexicon	<b>LR</b>	<b>78,85</b>
SemEval-2015 Lexicon	LR	77,67
SenticNet Lexicon	SVM	77,35
<i>Emoticons</i>	LR	77,91
Opinion Lex. + Emoticons	LR	78,56
Léxicos combinados	SVM	78,40
Opinion Lex. + Emoticons + POS	LR	78,48

## Resultados - HCR

Método	Algoritmo	Acc. (%)
Aprendizado de Máquina	SVM	64,47
<i>Part-of-Speech (POS)</i>	LR	63,20
Opinion Lexicon	SVM	65,70
SemEval-2015 Lexicon	LR	64,03
SenticNet Lexicon	SVM	64,86
<i>Emoticons</i>	LR	62,99
Opinion Lex. + Emoticons	LR	65,07
Léxicos combinados	SVM	65,90
Opinion Lex. + Emoticons + POS	<b>SVM</b>	<b>66,32</b>

## Resultados - Sanders

Método	Acc. (%)	P (%)	R (%)	F1 (%)
Aprendizado de Máquina	76,89	69,91	65,64	67,50
<i>Part-of-Speech</i>	76,14	68,54	65,35	66,78
Opinion Lex.	79,09	<b>76,20</b>	64,22	68,41
SemEval-2018 Lex.	76,25	68,40	64,11	65,95
SenticNet Lexicon	76,11	68,47	64,26	66,06
<i>Emoticons</i>	75,71	67,83	63,65	65,44
Opinion Lex. + Emoji	76,01	68,47	64,26	66,06
Léx. comb.	76,68	68,97	<b>66,30</b>	67,52
Op. Lex. + Emoji + POS	78,96	74,81	65,21	68,77
<i>Ensemble</i>	78,37	74,40	63,83	67,50
GloVe (WE)	<b>79,36</b>	75,82	64,89	<b>68,85</b>

## Resultados - HCR

Método	Acc. (%)	P (%)	R (%)	F1 (%)
Aprendizado de Máquina	65,44	62,83	59,36	60,47
<i>Part-of-Speech</i>	63,61	62,03	57,97	59,29
Opinion Lex.	<b>69,11</b>	66,78	62,67	63,90
SemEval-2018 Lex.	65,44	63,54	58,17	59,53
SenticNet Lexicon	64,83	62,99	58,00	59,34
<i>Emoticons</i>	62,69	63,32	57,03	58,86
Opinion Lex. + Emoji	64,83	62,31	58,52	59,57
Léx. comb.	64,53	61,46	61,59	61,52
Op. Lex. + Emoji + POS	62,69	61,25	59,60	60,23
<i>Ensemble</i>	65,44	61,89	59,29	59,97
GloVe (WE)	68,22	<b>70,31</b>	<b>63,34</b>	<b>65,34</b>

## Observações

- Base de dados muito grande;
- Apenas 10% utilizada para o treinamento;
- 30ª colocação;
- 1º Colocado utilizou SVM e RNN [Çöltekin et al. 2018].

## Resultados

Modelo	F1	P	R	Acc
<b>WE+BoW-LR</b>	<b>21.497</b>	<b>26.208</b>	<b>20.843</b>	<b>31.588</b>
WE+BoW-SVM	21.023	27.034	21.403	32.570
BoW-LR	20.351	24.923	19.824	30.830
BoW-SVM	20.194	26.659	20.518	31.966
BoW-RF	15.793	19.890	15.310	25.842
Tübingen-Oslo	35.991	36.551	36.222	47.094



- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências

## Considerações finais

- Explorar métodos de Análise de Sentimentos;
- Implementar abordagens que aproveitem a variedade dos métodos;
- Divulgar os resultados obtidos.

## Considerações finais

- Explorar métodos de Análise de Sentimentos;
- Implementar abordagens que aproveitem a variedade dos métodos;
- Divulgar os resultados obtidos.

## Publicações

- Artigo publicado na 5ª Escola Regional de Informática;
- Artigo publicado na Revista de Sistemas de Informação da Faculdade Salesiana Maria Auxiliadora;
- Artigo publicado no SemEval-2018 Task 2 - Emoji Prediction in Tweets.

- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros**
- 9 Referências

## Objetivo de pesquisa

Estudar algoritmos de *Deep Learning*, especificamente sobre Redes Neurais Recorrentes.

## Objetivo de pesquisa

Estudar algoritmos de *Deep Learning*, especificamente sobre Redes Neurais Recorrentes.

## Tema do SemEval-2019 Task 5

- Detecção multilíngue de discurso de ódio contra imigrantes e mulheres no Twitter.

FIM.

Perguntas?

- 1 Introdução
- 2 Objetivos
- 3 Metodologia
  - Pré-processamento
  - Representação vetorial dos tweets
  - Abordagens
  - Algoritmos
- 4 Datasets
  - Sanders
  - HCR
  - SemEval-2018
- 5 Aplicações dos Métodos
- 6 Resultados
- 7 Conclusões
- 8 Trabalhos Futuros
- 9 Referências





B. Pang e L. Lee, “Opinion mining and sentiment analysis”, Found.Trends Inf.Reptr., vol.2,no1-2, pp.1–135,jan. de 2008, issn:1554-0669.  
doi:10.1561/15000000011.  
endereço:<http://dx.doi.org/10.1561/15000000011>.



N. F. F. SILVA, “Análise de sentimentos em textos curtos provenientes de redes sociais”, Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, 2016.



Da Silva, N. F., Hruschka, E. R., and Hruschka Jr, E. R. (2014). Tweet sentiment analysis with classifier ensembles. Decision Support Systems, 66:170–179.



Félix, Nádia. and P. Ribeiro, Alison. (2018). #TeamINF at SemEval-2018 Task 2: Emoji Prediction in Tweets.



Alison Ribeiro and Nádia Silva. 2017. Métodos para análise de sentimentos em tweets: um estudo comparativo. In V ERI-GO 2017.



Wang, Y., Huang, M., & Zhao, L. (2016). Attention-based lstm for aspect-level sentiment classification. In Proceedings of the 2016 conference on empirical methods in natural language processing (pp. 606-615).



Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8), 1735-1780.



Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014). Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555.



Joulin, A., Grave, E., Bojanowski, P., & Mikolov, T. (2016). Bag of tricks for efficient text classification. arXiv preprint arXiv:1607.01759.



M. Speriosu, N. Sudan, S. Upadhyay e J. Baldridge, “Twitter polarity classification with label propagation over lexical links and the follower graph”, em Proceedings of the First Workshop on Unsupervised Learning in NLP, sér. EMNLP '11, Edinburgh, Scotland: Association for Computational Linguistics, 2011, pp. 53–63, isbn:978-1-937284-13-8.  
endereço:<http://dl.acm.org/citation.cfm?id=2140458.2140465>.



N. J. Sanders, “Sanders-twitter sentiment corpus”, Sanders Analytics LLC, 2011.



Francesco Barbieri, Jose Camacho-Collados, Francesco Ronzano, Luis Espinosa-Anke, Miguel Ballesteros, Valerio Basile, Viviana Patti, and Horacio Saggion. 2018. SemEval-2018 Task 2: Multilingual Emoji Prediction. In Proceedings of the 12th International Workshop on Semantic Evaluation (SemEval-2018), New Orleans, LA, United States. Association for Computational Linguistics.



S. M. Mohammad, S. Kiritchenko e X. Zhu, “Nrc-canada: Building the state-of-the-art in sentiment analysis of tweets”, CoRR, vol. abs/1308.6242, 2013. endereço: <http://arxiv.org/abs/1308.6242>.



C. D. Manning, M. Surdeanu, J. Bauer, J. Finckel, S. J. Bethard e D. McClosky, “The StanfordCoreNLP natural language processing toolkit”, em Association for Computational Linguistics (ACL) System Demonstrations, 2014, pp.55–60. endereço: <http://www.aclweb.org/anthology/P/P14/P14-5010>.

# Referências



Çöltekin, Çağrı, and Taraka Rama. "Tübingen-Oslo at SemEval-2018 Task 2: SVMs perform better than RNNs in Emoji Prediction." Proceedings of The 12th International Workshop on Semantic Evaluation. 2018.



S. Rosenthal, A. Ritter, P. Nakov e V. Stoyanov, "Semeval-2014 task 9: Sentiment analysis in twitter", em Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014), 2014, pp. 73–80.



P. Saleiro, E. M. Rodrigues, C. Soares e E. Oliveira, "Feup at semeval-2017 task 5: Predicting sentiment polarity and intensity with financial word embeddings", ArXiv preprint arXiv:1704.05091, 2017.



M. Fouad, T. Gharib e A. Mashat, Efficient twitter sentiment analysis system with feature selection and classifier ensemble, jan. de 2018.



Mittal, Anshul, and Arpit Goel. "Stock prediction using twitter sentiment analysis." Stanford University, CS229 (2011  
<http://cs229.stanford.edu/proj2011/GoelMittal-StockMarketPredictionUsingTwitterSentimentAnalysis.pdf>) 15 (2012).



Y. Bengio, R. Ducharme, P. Vincent e C. Jauvin, "A neural probabilistic language model", *Journal of machine learning research*, vol. 3, n o Feb, pp. 1137–1155, 2003.

# Agregadores de Classificadores para Análise de Sentimentos

Alison Pereira Ribeiro  
Profª Drª Nádia F. F. da Silva

Universidade Federal de Goiás  
Instituto de Informática – INF  
Goiânia - Goiás - Brasil

*alisonrib17@gmail.com, nadia@inf.ufg.br*

Dezembro de 2018