

# Sentiment Analysis for Mobile App Reviews in Brazilian Portuguese

Larissa F. S. Britto<sup>1</sup>, Luciano D. S. Pacífico<sup>1</sup>

<sup>1</sup>Departamento de Computação (DC) – Universidade Federal Rural de Pernambuco (UFRPE) – Recife – PE – Brazil

{larissa.feliciano, luciano.pacifico@ufrpe.br}

**Abstract.** *Even with the increasing popularity of Sentiment Analysis (SA), the amount of available resources and frameworks are still limited in Brazilian Portuguese language. This work will describe the steps for the development of a Brazilian Portuguese data set in the mobile Apps domain, for SA applications. In addition, the proposed database will be used to compare the main methods used in SA literature, such as Recurrent Neural Networks.*

**Resumo.** *Mesmo com a crescente popularidade da Análise de Sentimentos (AS), a quantidade de recursos e ferramentas disponíveis para o português brasileiro ainda é limitada. Neste trabalho serão descritas as etapas para o desenvolvimento de uma base de dados em português no domínio de aplicativos móveis, para aplicações em AS. Além disso, a base proposta será utilizada para a comparação dos principais métodos utilizados na literatura de AS, como as Redes Neurais Recorrentes.*

## 1. Introdução

A popularização de sites como redes sociais, blogs e fóruns faz com que cada vez mais pessoas exponham suas opiniões na Internet. Essa quantidade crescente de dados textuais faz da Análise de Sentimentos (AS), que consiste no estudo computacional dessas opiniões, uma área popular na Mineração de Texto e Processamento de Linguagem Natural (PLN).

O principal objetivo da AS é extrair informações úteis referentes a sentimentos em textos, informações tais como conhecimentos, críticas e opiniões sobre determinado tema. Há algum tempo a AS tem sido área de interesse de vários pesquisadores [Turney 2002][Ramanathan and Meyyappan 2019][John et al. 2019][Wang et al. 2019], devido a suas aplicações de relevante impacto, como análise do desempenho de produtos [Fang and Zhan 2015], detecção de distúrbios (como depressão) [Wang et al. 2013] e até predição de resultados de disputas eleitorais [Jose and Chooralil 2016].

Existem diversos problemas de interesse da AS, como identificação de sarcasmo ou ironia [Farias and Rosso 2017], adaptação de domínios

[Blitzer et al. 2006][Pan et al. 2010], e o principal deles, a classificação de polaridade [Turney 2002][Zuo 2018]. Na classificação de polaridade, é possível, a partir de um texto, categorizá-lo de acordo com o sentimento que o mesmo expressa, utilizando de algoritmos de aprendizagem supervisionada (*classificadores*). Entre os classificadores mais utilizados na literatura para essa tarefa, podemos citar o Naive Bayes [Zuo 2018], Máquinas de Vetores de Suporte (*Support Vector Machine*) [Lu and Wu 2019][Guan et al. 2018], Regressão Logística [Al Omari et al. 2019][Ramadhan et al. 2017], Árvores de Decisão e Floresta Aleatória (*Random Forest*) [Rathi et al. 2018][Rane and Kumar 2018][Hegde and Padma 2017]. No idioma português, a maioria dos trabalhos existentes abordam a temática de classificação de polaridade, como em [de Aguiar et al. 2018] e [Souza and Vieira 2012], onde algoritmos de aprendizagem de máquina foram utilizados para classificar os sentimentos de postagens feitas por usuários de redes sociais, como o Twitter<sup>1</sup>. Em [Martinazzo and Paraiso 2010], o método estatístico de Análise Semântica Latente (*Latent Semantic Analysis*) foi utilizado para identificar a relação entre palavras em textos, e para a detecção automática de emoções em notícias curtas, como manchetes e subtítulos.

Em comparação a outros idiomas, a língua portuguesa possui, mesmo nos dias atuais, um déficit de recursos que dificultam pesquisas em AS, como por exemplo, bases de dados públicas em domínios variados (ex.: livros, jogos e aplicativos móveis) e ferramentas adequadas para o PLN, sendo tais recursos facilmente encontrados em idiomas como o inglês, por exemplo. O desenvolvimento de *Web Corpus* (coleção estática de vários documentos baixados da Web) [Schäfer and Bildhauer 2015], tem sido empregado pelos pesquisadores para suprir as necessidades por bases de dados. Em [de Souza et al. 2018], um *corpus* de comentários sobre hotéis, obtidos através de *Web Scraping*, foi desenvolvido. Nesse trabalho são analisadas ainda as ferramentas de pré-processamento de texto existentes quando aplicadas ao idioma português. Em [Moraes et al. 2015], uma API para o Twitter foi utilizada para a criação de um *corpus* sobre a Copa do Mundo de Futebol de 2014, onde a anotação de polaridade foi feita manualmente para maior confiabilidade do *corpus*. O Twitter também foi a fonte dos dados em [Brum and das Graças Volpe Nunes 2017], porém, por limitações da API oficial oferecida por essa rede social, a técnica de *Web Scraping* foi utilizada para a obtenção dos documentos, sendo a base de dados anotada manualmente e comparada a outras bases de dados da literatura.

Tendo em vista a escassez de bases de dados textuais em português em alguns domínios, neste trabalho iremos mostrar as etapas do desenvolvimento de um *Web Corpus* de comentários de usuários de aplicativos móveis (*Apps*). Além disso, será feita uma comparação experimental dos classificadores mais comumente utilizados na literatura de AS: Naive Bayes, Árvore de Decisão, Floresta Aleatória, Máquinas de Vetores de Suporte, Regressão Logística e Redes Neurais Recorrentes (*Simples e Long Short-Term Memory*).

As principais contribuições deste trabalho são:

1. Descrição do Desenvolvimento de um *Web Corpus*<sup>2</sup> de aplicativos em Português Brasileiro;

---

<sup>1</sup>[www.twitter.com](http://www.twitter.com)

<sup>2</sup><https://github.com/larifeliciana/Sentiment-Analysis-Portuguese-Datasets/>

2. Avaliação do *Web Corpus* desenvolvido, na tarefa de classificação de polaridade;
3. Comparação do desempenho dos modelos utilizados no idioma português.

O trabalho está dividido como segue. Na próxima seção (Seção 2) será apresentado o processo de criação do *corpus* proposto, sendo suas principais características discutidas. Na Seção 3, os algoritmos utilizados para a avaliação da base de dados proposta serão brevemente descritos. Na Seção 4, os resultados experimentais serão discutidos. Por fim, na Seção 5, as conclusões do trabalho e linhas para possíveis trabalhos futuros serão apresentadas.

## 2. Base de Dados

Com o crescimento da popularidade da Internet, sites como redes sociais, blogs e lojas virtuais têm se tornado uma grande fonte de dados úteis para mineração de texto, principalmente na análise de sentimentos. Os dados obtidos dessas fontes possuem, no entanto, uma série de ruídos e erros que precisam ser tratados. Nesta seção, serão debatidas as etapas para o desenvolvimento da base de dados utilizada neste trabalho (vistas na Figura 1), além de uma discussão das características do *corpus* gerado.

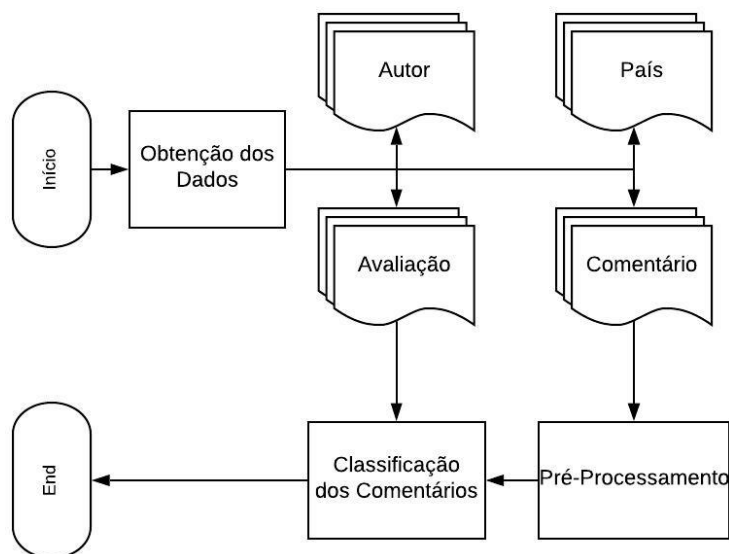


Figura 1. Etapas para o Desenvolvimento do *Corpus*.

### 2.1. Obtenção dos Dados

A primeira etapa para a criação de uma base de dados é a obtenção dos dados que irão compor a mesma. No *Web Corpus*, essa obtenção pode ser feita de duas formas: através de *APIs* disponibilizadas pelos próprios websites (ou desenvolvidas por seus usuários), ou pela aplicação da técnica de *Web Scrapping* a esses websites [Vargiu and Urru 2012]. Os dados utilizados para composição da base utilizada neste trabalho são comentários feitos por usuários de aplicativos da loja de aplicativos da *Apple*<sup>3</sup>, a *App Store*, e foram obtidos utilizando uma *API*<sup>4</sup>. Através dessa *API* foram extraídos, além dos comentários, outras

<sup>3</sup><https://www.apple.com/br/>

<sup>4</sup>[appfollow.io](http://appfollow.io)

informações, como nome do autor, seu país de origem, data do comentário, e a avaliação do usuário para o aplicativo (classificação feita de 1 a 5). Um total de aproximadamente 160 mil comentários foi obtido, dos quais 10 mil foram selecionados para a composição da base de dados deste trabalho.

## 2.2. Pré-Processamento

Atualmente, a maioria dos websites são contruídos em conjunto com a própria comunidade de usuários, ou seja, qualquer pessoa pode se cadastrar e dar suas opiniões sobre os mais diversos assuntos. Essa forma de construção tem a vantagem de que cada vez uma maior quantidade e diversidade de dados é disponibilizada, porém, existem inúmeras desvantagens do ponto de vista da qualidade dos documentos submetidos, que podem conter diversas irregularidades, o que faz do pré-processamento uma das etapas mais importantes no desenvolvimento de um *Web Corpus*. No pré-processamento, técnicas de limpeza e normalização dos comentários são aplicadas, com objetivo de obter um *corpus* padronizado e com menos erros. A ferramenta NLTK<sup>5</sup> e Expressões Regulares foram utilizadas para uniformizar os comentários, de acordo com as seguintes tarefas:

- **Lower Case** - Todas as letras são convertidas em letras minúsculas;
- **Remoção de Letras Repetidas** - Palavras com letras repetidas, utilizadas na Internet para dar ênfase ao que está sendo dito, também são removidas;
- **Correção Ortográfica** - Nesta etapa, erros ortográficos comuns foram corrigidos usando um dicionário, bem como algumas gírias e abreviaturas comumente usadas por usuários na Internet;
- **Remoção de Hashtags e links** - Nesta etapa, comandos comuns usados especialmente na Internet, como *hashtags*, links e menções, também são removidos;
- **Remoção de Stopwords** - Foram removidas todas as *stopwords*, com exceção das que podem representar alguma mudança de polaridade, como 'mas', 'sem' e 'não' [de Souza et al. 2018]. Também foram removidas pontuações e números.

Na Tabela 1 usaremos dois comentários da base de dados, um positivo e outro negativo, para ilustrar os efeitos da etapa de pré-processamento utilizada no desenvolvimento do *corpus* proposto.

## 2.3. Classificação

As avaliações na *App Store* são realizadas através de estrelas, que equivalem a notas, podendo ser dadas de 1 (pior avaliação) até 5 (melhor avaliação). Através dessa avaliação foi feita a anotação da polaridade dos sentimentos contidos nos documentos: avaliações com 4 ou 5 estrelas foram classificadas como positivas, enquanto 1 ou 2 estrelas foram consideradas negativas. Avaliações neutras (3 estrelas) foram removidas.

## 2.4. Estatísticas

Utilizando a biblioteca NLTK, obtemos algumas informações sobre o subconjunto da base utilizado neste trabalho. Tais estatísticas podem ser vistas na Tabela 2.

Outra informação importante para compreender os comentários que compõe a base são as palavras mais frequentes da mesma. Através delas, podemos observar opiniões frequentes dos usuários sobre os aplicativos. Nas Figuras 2 e 3, podemos ver os unigramas e

---

<sup>5</sup><https://www.nltk.org/book/>

	Positivo	Negativo
<b>Original</b>	<i>“Parabéns #Team99 cada vez melhor Sempre usei e vou usar”</i>	<i>“Sugiro enviar instrução para os taxistas aceitarem todas as bandeiras dos cartões ! #ficadika”</i>
<b>Pré-processado</b>	<i>“parabéns tag cada vez melhor sempre usei vou usar”</i>	<i>“sugiro enviar instrução taxistas aceitarem todas bandeiras cartões tag”</i>

**Tabela 1. Comentários Antes e Depois do Pré-Processamento.**

	Positive	Negative	Total
Nº Comentários	5000	5000	10000
Setenças/ Comentário	≈ 2	≈ 2.5	≈ 2.2
Palavras / Comentário	9.38	19.60	14.49
Palavras / Sentença	17.2	18.2	17.8
Vocabulário	5482	8121	10130

**Tabela 2. Informações sobre o corpus. As informações referente à quantidade de sentenças podem estar incorretas devido aos erros de pontuação por parte dos usuários.**

bigramas mais frequentes de cada classe de comentários. Nos comentários negativos, vemos alguns termos que expressam um sentimento negativo, como “horrrível”, “problema”, “não recomendo” e “não abre”. Já nos comentários positivos, se destacam termos como “excelente”, “melhor”, “funciona bem” e “serviço bom”.



**Figura 2. Unigramas e Bigramas mais frequentes nos comentários negativos.**



Figura 3. Unigramas e Bigramas mais frequentes nos comentários positivos.

## 2.5. Extração e Seleção de Características

A extração de características consiste em transformar o *corpus* obtido em informações úteis e suportadas pelos algoritmos que serão utilizados para a classificação. Neste trabalho, foi utilizado o modelo *Bag-of-Words*, no qual os documentos passam a ser representados pela frequência de cada termo que os compõem. A grande quantidade de termos em documentos pode fazer de uma tarefa de classificação simples um trabalho computacionalmente custoso. Para resolver esse problema, a seleção de características é uma técnica que busca reduzir esse custo, mantendo o bom desempenho dos classificadores pela seleção das características mais representativas da base de dados. Neste trabalho, a frequência dos termos será utilizada para a seleção das características, de modo que apenas os termos mais frequentes serão considerados na composição da base proposta.

## 3. Modelos Selecionados

Nesta seção, serão descritos todos os algoritmos utilizados na análise experimental deste trabalho: Naive Bayes, Árvore de Decisão, Floresta Aleatória, Regressão Logística, Máquinas de Vetores de Suporte, Redes Neurais Recorrentes Simples e *Long Short-Term Memory*.

### 3.1. Naive Bayes

O Naive Bayes (NB) é um modelo probabilístico baseado no teorema de Bayes, que utiliza da probabilidade condicional de cada palavra de um documento para calcular a probabilidade do mesmo pertencer a uma classe. O NB assume que todas as variáveis são estatisticamente independentes, ou seja, ele desconsidera qualquer correlação e contexto que possa haver entre as palavras de um documento.

### 3.2. Árvore de Decisão

A Árvore de Decisão (AD) é um método no qual a função de aprendizagem é representada através de uma estrutura hierarquizada [Bhatnagar 2018]. O processo de decisão é realizado dividindo o conjunto de dados de acordo com questões binárias executadas em uma característica por vez, gerando uma estrutura de árvore formada por nós de decisões. O processo de divisão é guiado por uma medida de satisfação (a entropia), usando uma abordagem gulosa para decidir quais características serão levadas em consideração para dividir o conjunto de dados em uma determinada iteração do método.

### 3.3. Floresta Aleatória

O algoritmo de Floresta Aleatória (*Random Forest* - RF) [Criminisi et al. 2011] é um método de *ensemble* que combina a previsão de várias Árvores de Decisão [Bhatnagar 2018], usando um mecanismo de voto majoritário. O algoritmo de Floresta Aleatória cria um conjunto de árvores de decisão para subconjuntos selecionados aleatoriamente dos dados de treinamento, agregando os votos de diferentes estimadores para decidir a classe final do dado de teste.

### 3.4. Regressão Logística

A Regressão Logística (RL) é um modelo discriminativo utilizado para prever a probabilidade das possíveis saídas de uma variável dependente, dado um conjunto de variáveis independentes. No caso da classificação em textos, é prevista a probabilidade da amostra em análise ser de uma classe ou de outra. O RL assume que a variável dependente pode ser prevista por uma combinação linear dos recursos problemáticos e parâmetros do modelo, representados pelo conjunto de treinamento.

### 3.5. Máquinas de Vetores de Suporte

As Máquinas de Vetores de Suporte (*Support Vector Machine* - SVM) são técnicas de classificação de aprendizado de máquina supervisionada baseada no princípio de Minimização de Risco Estrutural (*Structural Risk Minimization*) de [Vapnik 1991]. O SVM usa uma função de *kernel* para mapear os espaço de características de entrada em um novo espaço onde as classes são linearmente separáveis [Bergsma et al. 2005]. Para isso, o SVM constrói um hiperplano ótimo, que possa separar da melhor forma as instâncias de diferentes classes.

### 3.6. Redes Neurais Recorrentes

Redes Neurais Recorrentes (RNN) são tipos de redes neurais com um ou mais laços de realimentação, que usam suas conexões para armazenar informações sobre entradas recentes. Nessa arquitetura, o conhecimento adquirido de entradas anteriores é utilizado para gerar a saída atual. Desta forma, decisões são sempre influenciadas pelo que já foi aprendido. Isso faz com que as conexões da RNN possuam a capacidade de armazenar informações, o que a torna útil em atividades que levem em conta tempo e sequência, ou seja, onde as entradas são dependentes umas das outras, como por exemplo em textos, séries temporais, falas e vídeos.

### 3.7. Long Short-Term Memory

Redes *Long Short-Term Memory* (LSTM) são expansões das RNNs, propostas por [Hochreiter and Schmidhuber 1997], que ampliam suas memórias permitindo que a rede lembre de suas entradas por um longo período de tempo. Isso acontece porque a estrutura do LSTM funciona como uma memória [Sharfuddin et al. 2018] semelhante à estrutura de um computador. A estrutura do LSTM pode ler, gravar e excluir informações a partir de unidades de portas de entrada e saída, que são utilizadas para proteger os dados armazenados de entradas irrelevantes. A decisão das informações a serem mantidas ou excluídas é feita através dos pesos contidos nas informações.

#### 4. Resultados e Discussão

Visando avaliar a qualidade da base de dados proposta, alguns classificadores selecionados da literatura de AS foram aplicados, e suas performances comparadas. Nesta seção, os resultados experimentais para o *Web Corpus* desenvolvido serão apresentados. Um esquema de validação cruzada com 10 *folds* foi utilizado. Para a avaliação dos modelos, as seguintes métricas de classificação foram utilizadas: Acurácia, Revocação, Precisão e *F-Measure*. Na Tabela 3 os resultados obtidos podem ser visualizados.

Métrica \ Modelo	Acurácia		Precisão		Revocação		F-Measure	
	Média	Std.	Média	Std.	Média	Std.	Média	Std.
AD	0.7969	0.0126	0.7914	0.0223	0.8071	0.0192	0.7988	0.0128
RF	0.8456	0.0117	0.855	0.0208	0.8321	0.0112	0.8433	0.0137
NB	0.8311	0.0139	0.8622	0.0189	0.7896	0.0299	0.8237	0.0134
SVM	0.8498	0.0114	0.8312	0.0204	0.8778	0.0119	0.8537	0.0126
RL	0.8618	0.0062	0.8502	0.0153	0.8785	0.0136	0.864	0.0075
RNN	0.856	0.0106	0.873	0.0245	0.8338	0.0213	0.8525	0.0123
LSTM	<b>0.8699</b>	<b>0.01</b>	<b>0.89</b>	<b>0.0184</b>	<b>0.8445</b>	<b>0.0159</b>	<b>0.8665</b>	<b>0.01</b>

**Tabela 3. Resultados dos Experimentos.**

A Rede Neural Recorrente LSTM obteve o melhor resultado de acordo com as métricas adotadas, alcançando uma acurácia de 87%, seguida pela Regressão Logística e RNN Simples, que alcançaram 86.1% e 85.6% de acurácia média, respectivamente. Os três melhores modelos tiveram pouca diferença (menos de 1%), não se podendo assumir nenhum deles como o melhor. A Árvore de Decisão obteve um resultado inferior aos demais modelos, obtendo apenas 79.6% de acurácia média.

Todos os resultados condiziam com o que era esperado de cada modelo de acordo com o que é descrito na literatura, atestando assim a qualidade da base de dados proposta.

#### 5. Conclusão

Neste trabalho, foram comparados, no idioma português, o desempenho de alguns dos principais métodos utilizados na análise de sentimentos. Nos testes realizados com a base de dados proposta, os algoritmos de RNN Simples, LSTM e o Regressão Logística obtiveram os melhores resultados médios. O trabalho apresenta ainda a descrição detalhada das etapas realizadas para o desenvolvimento do *Web Corpus* no domínio de aplicativos proposto. Devido aos ruídos e erros encontrados em textos produzidos pelos usuários da Internet, uma intensa etapa de pré-processamento foi necessária. Através de uma análise na



base desenvolvida, pudemos obter informações úteis sobre os comentários feitos na *App Store*, como os maiores problemas encontrados pelos usuários em aplicativos, através dos termos mais frequentes encontrados nos documentos, como por exemplo, problemas com atualizações e com formas de pagamento. Também foi possível observar os pontos fortes desses aplicativos (pontos considerados importantes para a maioria dos usuários), tais como rapidez, segurança e promoções. Os resultados da análise experimental realizada atestaram a qualidade da base desenvolvida, tendo os modelos testados obtido resultados condizentes com a literatura da área de AS. Acreditamos que o *corpus* desenvolvido pode ser útil para a elaboração de outros trabalhos e pesquisas na área de análise de sentimentos e classificação de textos no idioma português. Como trabalhos futuros, pretendemos fazer a aplicação do *corpus* em outras tarefas da análise de sentimentos, buscando aumentar a quantidade de pesquisas desenvolvidas no nosso idioma.

### Agradecimentos

Os autores gostariam de agradecer ao CNPq e a CAPES pelo suporte financeiro.

### Referências

- Al Omari, M., Al-Hajj, M., Hammami, N., and Sabra, A. (2019). Sentiment classifier: Logistic regression for arabic services' reviews in lebanon. In *2019 International Conference on Computer and Information Sciences (ICCIS)*, pages 1–5.
- Bergsma, S., Jung, D., Lau, R., Wang, Y., and Wang, S. (2005). Machine learning approaches to sentiment classification cmut 551 : Course project winter , 2005.
- Bhatnagar, R. (2018). *Machine Learning and Big Data Processing: A Technological Perspective and Review*, pages 468–478.
- Blitzer, J., McDonald, R., and Pereira, F. (2006). Domain adaptation with structural correspondence learning. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing, EMNLP '06*, pages 120–128, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Brum, H. B. and das Graças Volpe Nunes, M. (2017). Building a sentiment corpus of tweets in brazilian portuguese. *CoRR*, abs/1712.08917.
- Criminisi, A., Konukoglu, E., and Shotton, J. (2011). Decision forests for classification, regression, density estimation, manifold learning and semi-supervised learning.
- de Aguiar, E. J., Faiçal, B. S., Ueyama, J., Silva, G. C., and Menolli, A. (2018). Análise de sentimento em redes sociais para a língua portuguesa utilizando algoritmos de classificação. In *Anais do XXXVI Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos*, Porto Alegre, RS, Brasil. SBC.
- de Souza, J. G. R., de Paiva Oliveira, A., and Moreira, A. (2018). Development of a brazilian portuguese hotel's reviews corpus. In *PROPOR*.
- Fang, X. and Zhan, J. (2015). Sentiment analysis using product review data. *J Big Data*, 2.
- Farias, D. H. and Rosso, P. (2017). Chapter 7 - irony, sarcasm, and sentiment analysis. In Pozzi, F. A., Fersini, E., Messina, E., and Liu, B., editors, *Sentiment Analysis in Social Networks*, pages 113 – 128. Morgan Kaufmann, Boston.

- Guan, X., Li, Y., Gong, H., Sun, H., and Zhou, C. (2018). An improved svm for book review sentiment polarity analysis. In *2018 International Conference on Transportation Logistics, Information Communication, Smart City (TLICSC 2018)*. Atlantis Press.
- Hegde, Y. and Padma, S. K. (2017). Sentiment analysis using random forest ensemble for mobile product reviews in kannada. In *2017 IEEE 7th International Advance Computing Conference (IACC)*, pages 777–782.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9:1735–80.
- John, A., John, A., and Sheik, R. (2019). Context deployed sentiment analysis using hybrid lexicon. In *2019 1st International Conference on Innovations in Information and Communication Technology (ICHICT)*, pages 1–5.
- Jose, R. and Chooralil, V. S. (2016). Prediction of election result by enhanced sentiment analysis on twitter data using classifier ensemble approach. In *2016 International Conference on Data Mining and Advanced Computing (SAPIENCE)*, pages 64–67.
- Lu, K. and Wu, J. (2019). Sentiment analysis of film review texts based on sentiment dictionary and svm. In *Proceedings of the 2019 3rd International Conference on Innovation in Artificial Intelligence, ICAI 2019*, pages 73–77, New York, NY, USA. ACM.
- Martinazzo, B. and Paraiso, E. C. (2010). Identificação de emoções em notícias curtas.
- Moraes, S. M. W., Manssour, I. H., and Silveira, M. S. (2015). 7x1-PT: um corpus extraído do twitter para análise de sentimentos em língua portuguesa (7x1-PT: a corpus extracted from twitter for sentiment analysis in Portuguese language). In *Proceedings of the 10th Brazilian Symposium in Information and Human Language Technology*, pages 21–25, Natal, Brazil. Sociedade Brasileira de Computação.
- Pan, S. J., Ni, X., Sun, J.-T., Yang, Q., and Chen, Z. (2010). Cross-domain sentiment classification via spectral feature alignment. In *Proceedings of the 19th International Conference on World Wide Web, WWW '10*, pages 751–760, New York, NY, USA. ACM.
- Ramadhan, W. P., Novianty, S. T. M. T. A., and Setianingsih, S. T. M. T. C. (2017). Sentiment analysis using multinomial logistic regression. In *2017 International Conference on Control, Electronics, Renewable Energy and Communications (ICCREC)*, pages 46–49.
- Ramanathan, V. and Meyyappan, T. (2019). Twitter text mining for sentiment analysis on people’s feedback about oman tourism. In *2019 4th MEC International Conference on Big Data and Smart City (ICBDSC)*, pages 1–5.
- Rane, A. and Kumar, A. (2018). Sentiment classification system of twitter data for us air-line service analysis. In *2018 IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC)*, volume 01, pages 769–773.
- Rathi, M., Malik, A., Varshney, D., Sharma, R., and Mendiratta, S. (2018). Sentiment analysis of tweets using machine learning approach. In *2018 Eleventh International Conference on Contemporary Computing (IC3)*, pages 1–3.

- Schäfer, R. and Bildhauer, F. (2015). Web corpus construction roland schäfer and felix bildhauer (freie universität berlin) morgan claypool (synthesis lectures on human language technologies, edited by graeme hirst, volume 22), 2013, 145 pages, paper-bound, isbn 9781608459834, doi:10.2200/s00508ed1v01y201305hlt022. *Computational Linguistics*, 41:161–163.
- Sharfuddin, A. A., Tihami, M. N., and Islam, M. S. (2018). A deep recurrent neural network with bilstm model for sentiment classification. *2018 International Conference on Bangla Speech and Language Processing (ICBSLP)*, pages 1–4.
- Souza, M. and Vieira, R. (2012). Sentiment analysis on twitter data for portuguese language. In Caseli, H., Villavicencio, A., Teixeira, A., and Perdigão, F., editors, *Computational Processing of the Portuguese Language*, pages 241–247, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Turney, P. D. (2002). Thumbs up or thumbs down?: Semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, pages 417–424, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Vapnik, V. (1991). Principles of risk minimization for learning theory. In *Proceedings of the 4th International Conference on Neural Information Processing Systems, NIPS'91*, pages 831–838, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Vargiu, E. and Urru, M. (2012). Exploiting web scraping in a collaborative filtering- based approach to web advertising. *Artif. Intell. Research*, 2:44–54.
- Wang, L., Niu, J., and Yu, S. (2019). Sentidiff: Combining textual information and sentiment diffusion patterns for twitter sentiment analysis. *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1.
- Wang, X., Zhang, C., Ji, Y., Sun, L., Wu, L., and Bao, Z. (2013). A depression detection model based on sentiment analysis in micro-blog social network. In *PAKDD Workshops*.
- Zuo, Z. (2018). Sentiment analysis of steam review datasets using naive bayes and decision tree classifier.