

# Sentiment Analysis through Conversational Data

Alexander Espronceda Gómez  
Ing. en Tecnologías de Software

Universidad Autónoma de Nuevo León  
Facultad de Ingeniería Mecánica y Eléctrica

November 26, 2021

# Índice

## 1 Introduction

- Justification
- Hypothesis
- Objective

## 2 Background

- Basic Concepts

- Supervised Machine Learning
- Sentiment Analysis
- Tokenization

## 3 Related Work

## 4 Proposed Solution

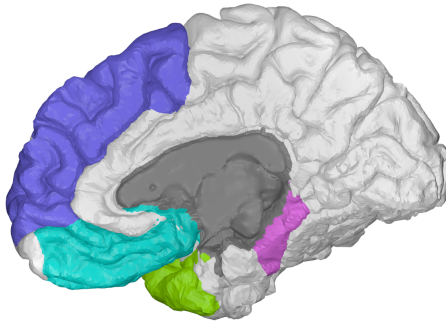
- Evaluation

## 5 Conclusions

## Abstract

In this thesis, open-sourced software is proposed, which interprets the text entered by a person and determines how they are feeling at the moment, with the purpose of being used in tandem with another software or algorithms focused on conversational data.

# Introduction



**Figure:** Lateral brain map of the parts in charge of the empathy processes. Drawing generated using BrainPainter [1].

# Justification

## Justification

This project could prove especially useful towards being used in projects designed for people who have trouble discerning when to console someone or having an idea of how other people or even themselves feel. To this end, the decision was made to work on this project.

# Hypothesis

## Hypothesis

Using supervised machine learning with a neural network could accurately classify the sentiment behind an input text as “Good”, “Neutral” or “Bad”, with the purpose of being implemented in tandem with another software or algorithms focused on conversational data.

# Objective

## Objective

Make software capable of determining how the person that writes the input text is feeling according to the words in it, while keeping the code open-source so it can be used in other projects. This could be achieved thanks to the technology present in machine learning algorithms and an extensive amount of datasets.

# Basic Concepts

**Machine Learning** Also known as ML. The type of algorithm needed for automatic processing, making the machine “learn” (hence the name) over time given enough data.

**Neural Network** A Machine Learning algorithm that uses weights and filters to output data.

**Natural Language Processing** This is the method used for the algorithm to understand the content of the sentences, this is usually achieved by using tokenization but a preset corpus can also be used.



# Basic Concepts

**Sentiment Analysis** This involves a ML algorithm, usually a Neural Network, that is able to analyze sentences and classify them according to the words used.

**Corpus** Preset internal dictionary that the algorithm uses.

**Tokenizing** Process that converts every word in the lexicon to an assigned number for easier processing

# Concept

Supervised ML can be described, broadly and figuratively speaking, as a black box where some data is inserted as an input and numbers come out of it as an output [2].

This output, as opposed to other types of Machine Learning, is later analyzed and compared to real life data.

# Machine Learning Stages

## Training

- Processes the inputs and makes educated guesses.
- Changes weights accordingly.

## Validation

Input is fed to the algorithm and information needs to be compared to the real results to test the accuracy percentage.

# Concept

## Process

- The sentence to analyze is broken down to its component parts, this process is called *tokenization*, and the resulting products are called *tokens*.
- Every token is then tagged, making it part of an internal dictionary or *lexicon*
- A score is assigned to every token depending on the used dataset.

# Tokenization Example

This is an example text

We can tell there are 5 words in the example phrase. So:

1, 2, 3, 4, 5

# Tokenization Example

This is another example

If we used the same tokenizer:

1, 2, 6, 4

## Related Work

**Table:** Comparison between existing literature and the present work:  
✓ indicates the fulfillment of a criterion, otherwise × is used.

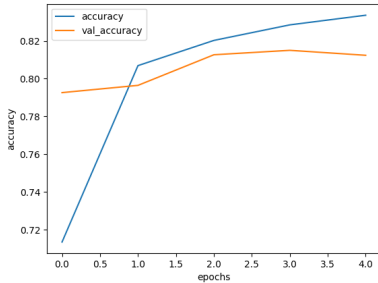
Project	Neural Network	Text Processing	Sentiment Analysis	Open Source	Modular
Blenn et al. [3] Maximum Entropy	✓	✓	✓	×	×
Blenn et al. [3] Support Vector Machines	✓	✓	✓	×	×
Blenn et al. [3] Lingpipe	✓	✓	✓	×	×
Morris et al. [4]	✓	✓	✓	×	×
Bird et al. [5]	✓	✓	×	✓	×
Pang et al. [6]	✓	✓	✓	✓	×
Ahmad et al. [7]	✓	✓	×	✓	×
Wang et al. [8]	✓	✓	✓	✓	×
Capuano et al. [9]	✓	✓	✓	×	×
Chiril et al. [10]	✓	✓	×	✓	×
Rôchert et al. [11]	✓	✓	✓	✓	×
The present work	✓	✓	✓	✓	✓

# Tools

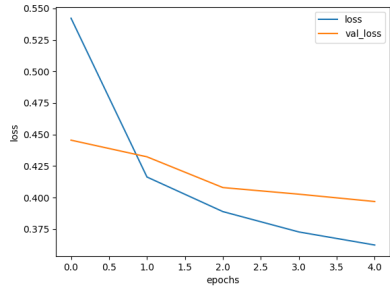
## Tools

This project is built on Python v3.8.10, The libraries used for this project to come to fruition are TensorFlow v2.6.0 and Keras v2.6.0 for the Neural Network section and Natural Language Toolkit v3.5 (also known as NLTK) for the tokenization and stemming process.





**Figure:** Accuracy values of the finished project



**Figure:** Loss values of the finished project

# Evaluation

The purpose of these experiments is to determine if the parameters chosen for this project are optimal and, if not, correct them and know the reason behind the improvement.  
Lower loss and higher accuracy are preferred.

# Experiment Results

Table: Experiment results

	Training		Cross-Validation	
	Loss	Accuracy	Loss	Accuracy
Experiment 1	0.6916	0.7130	<b>0.8709</b>	0.6234
Experiment 2	0.5956	0.7576	0.7649	0.6821
Experiment 3	0.5829	0.7564	0.7373	0.6780
Experiment 4	0.5455	0.7741	0.6704	0.7110
Experiment 5	0.6222	0.6550	0.7186	0.5357
Experiment 6	0.6041	0.7451	0.6555	0.7097
Experiment 7	0.6030	0.7421	0.6579	0.7156
Experiment 8	0.3624	0.8337	0.3871	<b>0.8124</b>

# Conclusions

## Conclusion

In the end, the conclusion reached is that the hypothesis was correct; it is possible for a Machine Learning algorithm to predict how a person is feeling based on an input, albeit some faults can be caused by the datasets used in training it. This can be changed using some quality control on them, or using personalized data specifically catered to this project.

## Future Work

This project would greatly benefit from a dataset that takes into consideration sentences that can be said in any context and still be correctly classified. And, of course, the less ortographical errors there are, the better.

# References

- [1] Razvan Marinescu, Arman Eshaghi, Daniel Alexander, and Polina Golland. Brainpainter: A software for the visualisation of brain structures, biomarkers and associated pathological processes. *arXiv preprint arXiv:1905.08627*, 2019.
- [2] Xian-Da Zhang. Machine learning. In *A Matrix Algebra Approach to Artificial Intelligence*, pages 223–440. Springer, 2020.
- [3] Norbert Blenn, Cassandra Charalampidou, and Christian Doerr. Context-sensitive sentiment classification of short colloquial text. In Robert Bestak, Lukas Kencl, Li Erran Li, Joerg Widmer, and Hao Yin, editors, *NETWORKING 2012*, pages 97–108, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. ISBN 978-3-642-30045-5.
- [4] Robert R Morris, Kareem Kouddous, Rohan Kshirsagar, and Stephen M Schueller. Towards an artificially empathic text sentiment analysis. *Ecological Informatics*, 64, 2021. ISSN 1574-9541. doi:<https://doi.org/10.1016/j.ecoinf.2021.10148>.
- [5] Jordan J Bird, Anikó Ekárt, and Diego R Faria. Chatbot interaction with artificial intelligence: human data augmentation with t5 and language transformer ensemble for text classification. *Journal of Ambient Intelligence and Humanized Computing*, pages 1–16, 2021.
- [6] Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. Thumbs up? sentiment classification using machine learning techniques. *arXiv preprint cs/0205070*, 2002.
- [7] Shakeel Ahmad, Muhammad Zubair Asghar, Fahad M Alotaibi, and Irfanullah Awan. Detection and classification of social media-based extremist affiliations using sentiment analysis techniques. *Human-centric Computing and Information Sciences*, 10(1):1–15, 2020. ISSN 1438-8871.
- [8] Nicola Capuano, Luca Greco, Pierluigi Ritrovato, and Mario Vento. Sentiment analysis for customer relationship management: An incremental learning approach. *Applied Intelligence*, 51:3339–3352, 2021.
- [9] Patricia Chiril, Endang Wahyu Pamungkas, Farah Benamara, Véronique Moriceau, and Viviana Patti. Emotionally informed hate speech detection: a multi-target perspective. *Cognitive Computation*, pages 1–31, 2021.
- [10] Daniel Röcher, German Neubaum, and Stefan Stieglitz. Identifying political sentiments on youtube: A systematic comparison regarding the accuracy of recurrent neural networks and machine learning models. In *Multidisciplinary*