

# Comparing regression and LSTM sentiment analysis models

Language Processing II exam  
Jens Egholm Pedersen  
xtp778@sc.ku.dk

June 16, 2017

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Background</b>	<b>2</b>
2.1	Language processing . . . . .	2
2.1.1	Sentiment analysis . . . . .	3
2.2	Linear regression . . . . .	3
2.2.1	Error reporting . . . . .	3
2.3	Recurrent neural networks . . . . .	3
<b>3</b>	<b>Data generation</b>	<b>3</b>
3.1	Learning models . . . . .	3
3.1.1	Data set . . . . .	4
3.1.2	K-fold cross validation . . . . .	4
<b>4</b>	<b>Quantitative evaluation</b>	<b>4</b>
<b>5</b>	<b>Qualitative evaluation</b>	<b>4</b>
5.1	Other approaches . . . . .	4
<b>6</b>	<b>Conclusion</b>	<b>4</b>
<b>7</b>	<b>References</b>	<b>5</b>
<b>8</b>	<b>Appendix A: Data</b>	<b>6</b>
<b>9</b>	<b>Appendix B: Software</b>	<b>7</b>
9.1	NLTK . . . . .	7
9.2	Scikit-learn . . . . .	7
9.3	Tensorflow . . . . .	7
9.4	Keras . . . . .	7

# 1 Introduction

Due to the enormous rise in accessible data and processing capabilities, computers have been moving from the restricted formalized domain of mathematics and computability, and into the more complex and seemingly chaotic domain of natural language (Nilsson, 2009; Jurafsky and Martin, 2000). Automating understanding of natural language have a wide range of helpful applications, ranging from translation to recommender systems to generic personal assistants (Cox, 2005; Bird et al., 2009).

As a part of this development, sentiment analysis has been attracting increased attention (Bird et al., 2009; Jurafsky and Martin, 2000). In the present paper, sentiment analysis is applied to a corpus of movie reviews presented by Pang and Lee (2005). By using standard natural language processing tools and machine learning techniques, four models are built and trained to predict the rating of a movie review on a discrete normalized scale.

## 2 Background

This section illuminates the theory and background for language processing and basic methodology for processing text. Sentiment analysis as a fundamental concept for this setting will be introduced second, followed by fundamental statistical metrics. Lastly, recurrent neural networks will be introduced before moving on to the data generation.

### 2.1 Language processing

Following the revolutionizing formalizations of grammar by Chomsky (2002), first written in 1957, much work was put into working with the computability of language (Jurafsky and Martin, 2000). In the beginning of the 21st century, the previous work converged with the statistical machine learning community, giving rise to supervised, semi-supervised and even unsupervised models of language understanding and translation (Jurafsky and Martin, 2000).

There are several approaches to language processing, but it is common to deconstruct text into its constituents and assign meaning to each piece, with the hope that the disassembly gives additional meaning that can aggregates to an understanding of the text as a whole when the pieces are put back together (Jurafsky and Martin, 2000). *Constituents* can range from morphemes to words to whole sentences in this context. This processes is referred to as feature extraction due to the enrichment of the text with meta-data.

A fundamental method for feature extraction is to construct a probabilistic model by simply joining items (ranging from phonemes to word pairs) together in pairs of one, two or more, called  $n$ -grams (Jurafsky and Martin, 2000). The advantage of this approach focuses on speed and the simple, yet efficient, probabilities in for instance word pairs which can be applied to predict likelihood of future occurrences (Jurafsky and Martin, 2000).

Term frequency, inverse document frequency (tf-idf) is another relevant feature that describes the relevance of an item, weighted against the number of times it appears in other documents (Jurafsky and Martin, 2000) (see figure 1). Tf-idf favours words that is used heavily in a single document (assumed to be

$$tf - idf = Nw \cdot \log(\frac{N}{k})$$

Figure 1: Formula for tf-idf.

important) while diminishing the importance of words that often occur in other documents, and is a heavily used metric (Jurafsky and Martin, 2000).

### 2.1.1 Sentiment analysis

Another popular approach for disassembling text called part of speech-tagging (POS tagging), assigns grammatical categories to each word and attempt to build a linguistic model that can be understood by a machine (Jurafsky and Martin, 2000). This works well for information retrieval, where objective truths can be extracted <sup>1</sup>, but fails to capture the complexity and emotional variety of sentences that appear outside the grammar <sup>2</sup>, such as emotions and opinions (Jurafsky and Martin, 2000).

Sentiment analysis focuses on extracting information about the emotions or opinions of a text (Pang and Lee, 2008). Defining what a *sentiment* is not a simple task, and have been scrutinized extensively (Jurafsky and Martin, 2000; Pang and Lee, 2008). A much used approach was introduced by Ekman (1992) where he divided six basic emotions, which could be conveniently operationalized by computers. Since the dataset for this paper revolves around opinions in a single dimension, *sentiment* will in this paper simply refer to a subjective experience, discretized to a number (see section 3.1.1).

## 2.2 Linear regression

### 2.2.1 Error reporting

Mean absolute error (MAE) Root mean square error (RMSE)

## 2.3 Recurrent neural networks

# 3 Data generation

This section

## 3.1 Learning models

The four models model is built using the Scikit Learn (Sklearn) package and second is based on the deep learning library Keras (see appendix in section 9). Two The corpus is evaluated against four different implementations of sentiment analysis, based on two different models. The first model uses the

---

<sup>1</sup>Such as the objective statements about an object in the sentence "*The cat likes tuna fish*".

<sup>2</sup>Such as a heavily sarcastic sentence like "*I'd really truly love going out in this weather!*". For POS-tagging sarcasm is difficult to capture because it is not implicitly encoded in the grammatical categories.

Model	MAE	RMSE
Linear	0.1649	0.2075
Linear + tf-idf	0.1494	0.1855
LSTM	0.1667	0.2045
LSTM + tf-idf	0.1641	0.2014

Figure 2: Comparison of metrics for the four different models using mean absolute error (MAE) and root mean squared error (RMSE).

### 3.1.1 Data set

### 3.1.2 K-fold cross validation

The dataset

N-grams Jurafsky and Martin (2000).

max features for bow (10'000)

this section must report on actual test-runs utilizing the movie review corpus and applying at least two machine learning algorithms (Naïve-Bayes being one example)

## 4 Quantitative evaluation

demonstrating how SA results can be evaluated automatically

## 5 Qualitative evaluation

evaluating and discussing the quantitative results wrt. validity, reliability and/or relevance in a “real -world” perspective

The very definition of *sentiment* in section ?? as a subjective opinion makes comparison difficult. As stated in (Pang and Lee, 2005) the entries in the present dataset is difficult to compare upon, since ratings are relative between reviewers; a rating of 0.8 might be high for one reviewer, but low for another. Unfortunately a simple normalization does not solve the problem because the scale of one reviewer might not even be linearly comparable to another.

### 5.1 Other approaches

Not taken AFINN, leksika word2vec Nielsen (2011)

## 6 Conclusion

presenting in a condensed form your results and observations, e.g. pointing to strengths, weaknesses, and future directions

## 7 References

- Bird, S., Klein, E., and Loper, E. (2009). *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*. O'Reilly, Beijing.
- Chomsky, N. (2002). *Syntactic Structures*. Mouton classic. Bod Third Party Titles.
- Cox, M. T. (2005). Metacognition in computation: A selected research review. *Artificial Intelligence*, 169(2):104 – 141. Special Review Issue.
- Ekman, P. (1992). An argument for basic emotions. *Cognition and Emotion*, pages 169–200.
- Jurafsky, D. and Martin, J. H. (2000). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1st edition.
- Nielsen, F. Å. (2011). AFINN.
- Nilsson, N. J. (2009). *The quest for artificial intelligence - A history of ideas and achievements*. Cambridge University Press.
- Pang, B. and Lee, L. (2005). Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of ACL*, pages 115–124.
- Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2:1–135.

## 8 Appendix A: Data

The data used in this report is collected from Pang and Lee (2005). The data is available in the software repository (see 9). The readme with further descriptions on the data source is available at:

<http://www.cs.cornell.edu/people/pabo/movie-review-data/scaledata.README.1.0.txt>

## 9 Appendix B: Software

All software used to produce the results presented in this report is open-source and available through GitHub at:

<https://github.com/Jegp/langprocexam>

Below follows a list of the software used. An in-depth description on how to reproduce the results are available from the GitHub repository above.

### 9.1 NLTK

A toolkit for natural language processing in python. Visited 15th of June 2017.

<http://www.nltk.org/api/nltk.sentiment.html>

### 9.2 Scikit-learn

Machine learning library for Python. Visited 15th of June 2017.

<http://scikit-learn.org/stable/index.html>

### 9.3 Tensorflow

A machine learning library initially developed by Google Inc. Visited 15th of June 2017.

<https://www.tensorflow.org/>

### 9.4 Keras

A deep learning library which builds on top of Tensorflow. Visited 4th of June 2017.

<https://keras.io>