

Construction and Performance Analysis of a Groomed Polarity Lexicon Derived from Product Review Source Datasets

Derek Colley
School of Digital, Technologies and Arts
Staffordshire University
Stoke-on-Trent, United Kingdom
derek.colley@staffs.ac.uk

Md Asaduzzaman
School of Digital, Technologies and Arts
Staffordshire University
Stoke-on-Trent, United Kingdom
md.asaduzzaman@staffs.ac.uk

EXTENDED ABSTRACT

Abstract—Using a large, publicly-available dataset [1], we extract over 51 million product reviews. We split and associate each word of each review comment with the review score and store the resulting 3.7 billion word- and score pairs in a relational database. We cleanse the data, grooming the dataset against a standard English dictionary, and create an aggregation model based on word count distributions across review scores. This renders a model dataset of words, each associated with an overall positive or negative polarity sentiment score based on star rating which we correct and normalise across the set. To test the efficacy of the dataset for sentiment classification, we ingest a secondary cross-domain public dataset containing freeform text data and perform sentiment analysis against this dataset. We then compare our model performance against human classification performance by enlisting human volunteers to rate the same data samples. We find our model emulates human judgement reasonably well, reaching correct conclusions in 56% of cases, albeit with significant variance when classifying at a coarse grain. At the fine grain, we find our model is able to track human judgement to within a 7% margin in some cases. We consider potential improvements to our method and further applications, and the limitations of the lexicon-based approach in cross-domain, big data environments.

Keywords—Sentiment analysis, natural language processing, relational databases, language analysis, SQL aggregation.

I. INTRODUCTION

Sentiment analysis involves the assessment of natural language text segments to determine the degree of membership within some nominal classification taxonomy. Challenges in this space include improving the classification success against non-standard phrase forms and short text segments (such as Twitter posts). In this paper, we propose the derivation of a sentiment lexicon from an existing dataset of customer product reviews, and we present, test and validate a method for text segment deconstruction and sentiment calculation against the derived dataset.

II. RELATED WORK

Sentiment analysis is the art of programmatically extracting meaning from often abstract and unstructured segments of text. Applications are varied and include the provision of management information on brand perception in the marketplace (Kaiser et al., 2011), helping to displace risk of reputational damage. Nanli, Ping, Weiguo and Meng (2012) classified the contents of human communication as

objective and subjective and noted that accurate textual sentiment analysis is an unsolved problem.

Algorithmic approaches to addressing sentiment extraction differ widely. These include supervised and unsupervised machine learning methods, lexicon-based methods, use of keywords and concept extraction (Qazi, Rag, Hardaker and Standing, 2017). Our research focuses on sentiment polarity extraction using the lexicon-based method, where opinions are held to be positive, negative or somewhere on a bounded scale between these two finite extremes (Cambria and White, 2014). We attempt this not by using a dictionary-based approach, but by creating a weighted lexicon model, where each word in the lexicon is assigned a score, then applying this model to new text inputs.

We argue that, given a sufficiently-large corpus of balanced data, an algorithmic approach of simple summation and range normalisation can provide performance advantages through lower complexity when a model derived from this approach is used as a quantitative classification mechanism. This is not to denigrate SVM and other ML methods; SVMs (Paltoglou and Thelwall, 2010) are ubiquitous in part because they are able to deal with categorical, hierarchical, semi-structured and ordinal data, whereas simple functions have a brittleness which suits them to low-dimensional, atomic, structured datasets. Shayaa et al. (2018) extend this criticism, noting that the lexicon-based method is vulnerable to bias in the data source and cannot adapt well to unstructured data sources, albeit also noting the significant performance gain of this method over others.

Other challenges in lexicon-based sentiment analysis present themselves. These include the use of non-standard language artefacts such as emojis, abbreviations, misspellings, shortenings of common words and slang. In filtering against a standard English dictionary, potentially interesting sentiment information is lost; an active area of research addressed by, amongst others, Fernández-Gavilanes et al. (2018) who developed an emoji lexicon in mitigation. The evidence in the literature indicates many issues in Sentic computing remain research challenges; Chaturvedi et al. (2018) focus on the difficulties of separating fact and opinion within datasets, or subjectivity detection, noting particular difficulty in classifying weakly-subjective sentences, and in a recent survey, Hussain (2018) analysed 47 previous studies in sentiment analysis noting that certain factors such as world (or

domain) knowledge, 'bi-polar' words and large lexicons have a deleterious effect on overall accuracy rates.

Nevertheless, lexicon-based sentiment analysis has a strong theoretical and practical basis. Our contribution to this field, illustrates our outcomes from building and testing a simple unsupervised sentiment polarity calculator applied to a large and unstructured data set, and the applicability of this set to another domain, with the aim of establishing the limitations of this approach in a modern big data setting.

III. CREATING THE SENTIMENT DICTIONARY

A. Sourcing the ratings data

The ratings data is supplied as a publicly accessible dataset by Ni, Li and McAuley [1] and consists of a record set of 51m Amazon product reviews for books, dating to 2018. We downloaded this data, compromised of a single compressed file in JSON format, and examined the schema. In order to work with the file it was necessary to split it into smaller pieces for data ingestion; we wrote a simple file splitter program to do this, resulting in 1,036 files of approximately 30-50MB per file in size. Each except the last file contained exactly 50,000 documents. Next, we iterated through each file, loading each document into Microsoft SQL Server as a single NVARCHAR(MAX) uniquely identified by a auto-sequence number and timestamp. This resulted in 51m records in SQL server, each numbered. Using the JSON_VALUE() SQL function and the STRING_SPLIT() SQL function nested in a common table expression, we then iterated through the 51m records in batches of 100,000 and extracted the **overall** (review score) and **reviewText** (customer review) elements from each JSON document, splitting **reviewText** into its component words (space-separated) and storing the **overall** numeric, renamed to **score**, alongside each word, renamed to **word**. This resulted in 3,795,765,817 rows in total, each row containing a single **score** and a single **word**, each pair representing an occurrence of the word in the whole review dataset and the associated parent review score. For performance improvements when reading the data for the next steps, we converted the rowstore heap table to a columnar indexed table. Fig. 1 illustrates the preparation process.

FIG. 1. DATA PREPARATION PROCESS (REMOVED FOR ABSTRACT).

B. Data cleansing

We required a reasonably complete English dictionary upon which we could cleanse our collected data of misspellings, orphaned punctuation marks and other errata. We chose Webster's Unabridged Dictionary (2009) hosted at Project Gutenberg [2] which has a plain-text UTF-8 version available for download. This provided 102,668 distinct words.

We filtered our list of 3,795,765,817 words to a shortened list of 2,925,874,997 words by selecting all matching words against the Websters dictionary table into a new table using an inner (predicated) join.

Next, given that case is irrelevant to sentiment, we updated all entries in our table to uppercase format, although this is not strictly necessary as the collation of our database was

case-insensitive. This resulted in a table with rating score ('score') and valid descriptor ('word'). Fig. 2 illustrates the data cleansing process.

FIG. 2. DATA CLEANSING PROCESS (REMOVED FOR ABSTRACT)

C. Aggregation and score normalisation

With a dataset containing one entry for every occurrence of a word and its accompanying score (for the parent review as a whole), we transformed the data into an aggregate grouping over the word column, calculating both the \bar{x} score (per word) and the count of words. This resulted in a new set of distinct words with all duplicates removed, together with average score and frequency count.

To normalise the range of scores, we calculated the final sentiment score x in the range 0-1 using an ordinary normalisation function. Eq. 1 shows this function.

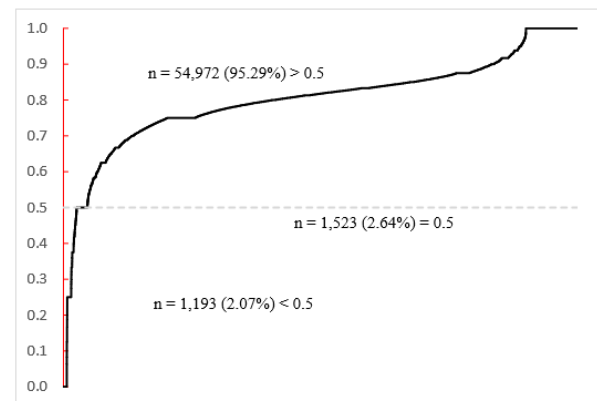
$$f(x) = \left(\frac{x - \min(X)}{\max(X) - \min(X)} \right), \forall x \text{ in } X \quad (1)$$

This resulted in a set of distinct words and associated normalized sentiment scores which we round to 3 d.p., representing a single scalar sentiment score $0 \leq x \leq 1$ for each word. At this point, the original ratings are discarded.

This process leaves a total of 57,688 distinct words, each with a single normalised sentiment score in the range 0-1, where 0 represents negative and 1 represents positive sentiments.

We note a potential issue; the resulting range of values is not uniformly distributed. Instead, we find 95.29% of all scores claim a positive sentiment (above the neutral value of 0.5); 2.64% of all scores are exactly neutral at 0.5; the remaining 2.07% of all scores claim a negative sentiment (below the neutral value of 0.5). This is indicative of the parent reviews, where similar percentages are rated 3, 4 or 5 'stars'. The consequence of this skew is that negative-sentiment words may be unrepresented. This skew, identical to the z-score distribution except bounded between 0.0-1.0, is shown in Fig. 3.

FIG. 3. SCORE DISTRIBUTION (NORMALISED)



We overcome this issue by resetting the median point at which we define a 'negative' or a 'positive' sentiment to the

median of the normalised scores, which is 0.817 for our lexicon. In doing so, we rebalance the population of words on either side of this dividing line. The function for midpoint correction for all scores s is thus:

CODE LISTING 1. USER-DEFINED SENTIMENT SCORE CALCULATION FUNCTION IN TRANSACT-SQL (REMOVED FOR ABSTRACT)

CODE LISTING 2. USAGE EXAMPLES FOR THE SENTIMENT SCORE CALCULATION UDF (REMOVED FOR ABSTRACT)

$$\begin{aligned} f(s) &= 0.817 - s + 0.5 \therefore \\ f(s) &= 1.317 - s \end{aligned} \quad (2)$$

Next, we demonstrate the creation of a user-defined function to take advantage of this lexicon for sentiment analysis.

D. Programmatic Access

To implement a mechanism to access the data programmatically, we created a SQL user-defined function (UDF), returning the sentiment score (a numeric x bounded to $0 \leq x \leq 1$ to 3 d.p.) given a string expression e of input words with no upper bound on length. We use \bar{x} smoothing to calculate the similarity score for any given expression e . Code Listing 1 provides this UDF in Transact-SQL form compatible with Microsoft SQL Server 2016 or above, although the code is easily portable to any ANSI-SQL compatible platform. Code Listing 2 provides example of uses. Both code samples rely on the existence of the ‘aggregations’ table as described, containing the lexicon and associated sentiment scores.

IV. TESTING AND VALIDATION

To test the efficacy of our implementation, we sourced Twitter data from the Internet Archive [3], extracting all tweets for a single day. We considered all tweets with substantial English-language content (defined as at least 8 separate words). We anonymised the metadata and removed identifiable information, such as user handles, then prepared this set as a survey instrument which we administered to 22 participants. We then calculated the average score given by each participant across each tweet to set a sentiment score which we treat as a truthful reflection of the tweet sentiment and divided by 10 to normalise. We prepared the same ordered set of 100 tweets and applied our UDF function against them, setting the neutral point at the calculated balance point of 0.817 (see Eq. 2), resulting in 100 sentiment scores which we annotated with the appropriate classification. Table 1 illustrates this data.

TABLE I. AVERAGE SENTIMENT SCORES (REMOVED FOR ABSTRACT)

We observe that our UDF function is notably cautious in score allocation, with a score range of 0.086, a maximum score of 0.561 and a minimum score of 0.475. Compare this against the human participant results with a range of 0.286, maximum of 0.691 and minimum of 0.4; a range increase of 332%.

To establish the extent of correlation between the UDF-generated scores and the human-generated scores, we first chart the data points on a scatter diagram (as shown in Fig. 4).

We then calculate the bivariate correlation co-efficient (PCC) in the normal way as 0.39, illustrating a weak but present correlation between the two variable sets.

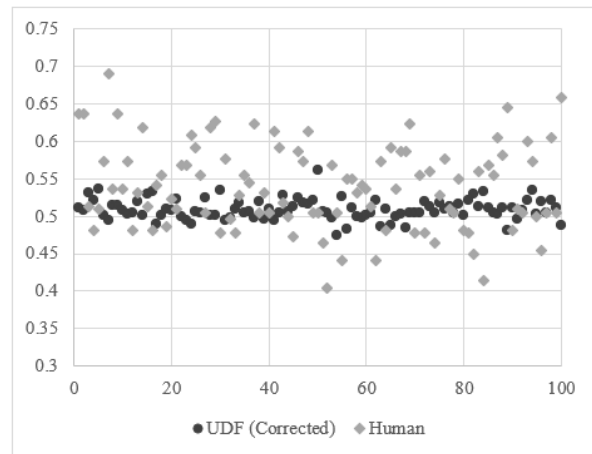


FIG. 4. COMPARISON OF UDF- VS. HUMAN SENTIMENT SCORES, PER TWEET

Another way we may establish relationships is by considering the midpoint 0.5 as the neutral point and classifying all scores above 0.5 as positive and all scores below 0.5 as negative regardless of fine-grained classification. Thus, we obtain the results shown in Table 2.

TABLE 2. UDF VS. HUMAN SCORE PERFORMANCE METRICS

UDF: Total positive tweets ($s > 0.5$)	78
UDF: Total neutral tweets ($s = 0.5$)	0
UDF: Total negative tweets ($s < 0.5$)	22
Human: Total positive tweets ($s > 0.5$)	76
Human: Total neutral tweets ($s = 0.5$)	0
Human: Total negative tweets ($s < 0.5$)	24
UDF vs. human s-score agreement	56%
UDF vs. human s-score disagreement	44%

We can observe that the overall volume of classification was similar for both UDF vs. human performance – the UDF classified 78 tweets as positive vs. human classification of 76 tweets as positive, and likewise negative; however, there was significant variance on a per-item level. In 56% of cases there was agreement on positive vs. negative between UDF and human scores. The UDF classifies very near the 0.5 midpoint boundary and the variance inequality means there is significantly more room for individual classification error.

As per our PCC correlation calculation of 0.39, there is a weak but present relationship between human- and UDF-driven score generation using our lexicon and method; that in general, the UDF method is conservative in range and consequently the outcome of text classification using this data

source and this method yields a 6% advantage over random chance when using coarse, or discrete, classification.

Finally, we can examine the average difference between human- and UDF-driven scoring by examining the deltas between the mean of the human scores and the UDF-generated score, per text item. These deltas were illustrated in Table 1. We note the average difference between them is just -0.031, with a range of 0.315 (31.5% potential swing) and a low standard deviation of +/-6.3% (0.063, bounded from 0.0-1.0) meaning that UDF-driven scoring generally tracks human-driven scoring to within a 7% margin when considering a continuous score range.

V. CONCLUSIONS AND FUTURE WORK

It is evident from the literature review that Sentic computing continues to present challenges for researchers and industry practitioners in ensuring accuracy in the face of the ever-changing variety of data available to mine, the constant evolution of language to incorporate new vocabulary and phrasings, the difficulty of classifying words, phrases and larger text blocks into quantised groupings and, not least, the difficulties in telling fact from fiction. In this research, we investigated lexicon-based sentiment analysis using a polarity weighting technique and applied this to a big data set derived from Amazon product reviews to create a training dataset; we then attempted to use this training set against another domain of data, a selection of random posts on a social media platform. We found a middling degree of success when compared against human performance at the same task, with 56% of posts classified correctly by the algorithm; we found a better degree of accuracy when looking at particular cases, tracking to within 7% of human performance. We conclude that although lexicon-based methods are applicable to big data sets, as evidenced in our case study, there remain challenges to be solved in applying such rigid polarity lexicons across domains, and these challenges are exacerbated by different vocabularies, intentions and inconsistencies naturally present within informal human language.

REFERENCES

- [1] J. Ni, J. Li and J. McAuley, "Justifying recommendations using distantly-labeled reviews and fine-grained aspects", in *Empirical Methods in Natural Language Processing (EMNLP)*, 2019. Available at: <http://cseweb.ucsd.edu/~jmcauley/pdfs/emnlp19a.pdf> [Accessed 22 Apr. 2021].
- [2] Project Gutenberg, "Webster's Unabridged Dictionary by Various", 2009. Available at: <http://www.gutenberg.org/ebooks/29765> [Accessed 22 Apr. 2021].
- [3] Internet Archive, "Archive Team Twitter Stream 2020-10", 2020. Available at: <https://archive.org/details/archiveteam-twitter-stream-2020-10> [Accessed 23 Apr. 2021].
- [4] Nanli, Z., Ping, Z., Weiguang, L.I. and Meng, C., 2012, November. Sentiment analysis: A literature review. In *2012 International Symposium on Management of Technology (ISMOT)* (pp. 572-576). IEEE.
- [5] Al-Moslemi, T., Omar, N., Abdullah, S. and Albared, M., 2017. Approaches to cross-domain sentiment analysis: A systematic literature review. *Ieee access*, 5, pp.16173-16192.
- [6] Qazi, A., Raj, R.G., Hardaker, G. and Standing, C., 2017. A systematic literature review on opinion types and sentiment analysis techniques. *Internet Research*.
- [7] Kumar, A. and Jaiswal, A., 2020. Systematic literature review of sentiment analysis on Twitter using soft computing techniques. *Concurrency and Computation: Practice and Experience*, 32(1), p.e5107.
- [8] Kaiser, Carolin; Schlick, Sabine; Bodendorf, Freimut. Warning system for online market research - Identifying critical situations in online opinion formation[J]. *Knowledge-Based Systems*, v 24, n 6, August 2011, pp. 824-836
- [9] Cambria, E. and White, B. (2014), "Jumping NLP curves: a review of natural language processing research", *IEEE Computational Intelligence Magazine*, Vol. 9 No. 2, pp. 48-57.
- [10] Das, A., 2017, July. Sentiment analysis. In *2017 8th International Conference on Computing, Communication and Networking Technologies (ICCCNT)* (pp. 1-6). IEEE.
- [11] Paltoglou, G. and Thelwall, M., 2010, July. A study of information retrieval weighting schemes for sentiment analysis. In *Proceedings of the 48th annual meeting of the association for computational linguistics* (pp. 1386-1395).
- [12] Shayaa, S., Jaafar, N.I., Bahri, S., Sulaiman, A., Wai, P.S., Chung, Y.W., Piprani, A.Z. and Al-Garadi, M.A., 2018. Sentiment analysis of big data: Methods, applications, and open challenges. *IEEE Access*, 6, pp.37807-37827.
- [13] Chaturvedi, I., Cambria, E., Welsch, R.E. and Herrera, F., 2018. Distinguishing between facts and opinions for sentiment analysis: Survey and challenges. *Information Fusion*, 44, pp.65-77.
- [14] Hussein, D.M.E.D.M., 2018. A survey on sentiment analysis challenges. *Journal of King Saud University-Engineering Sciences*, 30(4), pp.330-338.
- [15] Mohammad, S.M., 2017. Challenges in sentiment analysis. In *A practical guide to sentiment analysis* (pp. 61-83). Springer, Cham.
- [16] Fernández-Gavilanes, M., Juncal-Martínez, J., García-Méndez, S., Costa-Montenegro, E. and González-Castaño, F.J., 2018. Creating emoji lexica from unsupervised sentiment analysis of their descriptions. *Expert Systems with Applications*, 103, pp.74-91.
- [17] Taboada, M., Brooke, J., Tofiloski, M., Voll, K. and Stede, M., 2011. Lexicon-based methods for sentiment analysis. *Computational linguistics*, 37(2), pp.267-307.