# An Annotated Bangla Sentiment Analysis Corpus

Fuad Rahman
Apurba Technologies Ltd.
Dhaka, Bangladesh
fuad@apurbatech.com

Habibur Khan
Apurba Technologies Ltd.
Dhaka, Bangladesh
habib@apurbatech.com

Zakir Hossain
Apurba Technologies Ltd.
Dhaka, Bangladesh
zakir@apurbatech.com

Mahfuza Begum
Apurba Technologies Ltd.
Dhaka, Bangladesh
mahfuza@apurbatech.com

Sadia Mahanaz
Apurba Technologies Ltd.
Dhaka, Bangladesh
sadia@apurbatech.com

Ashraful Islam
Apurba Technologies Ltd.
Dhaka, Bangladesh
ashraful@apurbatech.com

Aminul Islam
Apurba Technologies Ltd.
Dhaka, Bangladesh
aminul@apurbatech.com

*Abstract* – **This paper presents a Bangla corpus specifically targeted for sentiment analysis and made available to researchers under an open-source licensing scheme[1]. We have collected and manually annotated over 10,000 sentences with sentiment polarity. We then moved to the Word domain and annotated over 15,000 words derived from these sentences with sentiment polarity. Each entry is the corpus has been cross-annotated by at least two and sometimes three annotators for ensuring quality. Also as a pre-requisite of creating a high quality sentiment analysis corpus, we had to build a secondary corpus for Bangla word stemming, which is also been cross-validated by at least two and sometimes three annotators for ensuring quality.**

*Index Terms* – *Sentiment Analysis, NLP, Bangla Corpus, Annotation, Open Source Corpus*

## I. Introduction

Sentiment analysis is a very important part of natural language processing. While very robust solutions for English already exists both in academic and commercial domains, for Bangla language, work in this area is still in its infancy. As the focus of tools for sentiment analysis has now shifted from rule based to machine learning methods, the need for annotated and ground truth data for training these solutions are of utmost importance. Unfortunately there is almost no serious corpus for Bangla language that is available for sentiment analysis, forcing researchers to stich together their own small corpora which are neither standardized and nor rigorously quality controlled. In this paper, we present a fully annotated corpus for sentiment analysis.

### A. Brief Background

In recent times, there has been some research reported on Bangla sentiment analysis. One common approach seems to be to translate a Bangla word and then use the polarity from the English translated word. [1][6][7]. While it works on straightforward words, it cannot handle the nuances of a language. For example the word "জটিল" means "complex" [2] and it has a negative sense. But in Bangla it is often used in a positive sense, for example, "তামিম আজ জটিল খেলছে". Another example is the word "খাওয়া" which is commonly translated to "eat", the common polarity of which is neutral. But in the sentence "তার খাওয়া নাই", the polarity is distinctly negative.

The reason for the popularity of this type of approach is very simple, a distinct lack of a ground truth corpus suited for training machine learning algorithms. Although there are some existing data set for sentiment analysis, but most of these are not available publicly. Some publicly available data set are small, e.g. [4] has about 4,000 sentences, whereas [1] has about 7,000 sentences.

One of the most significant resources is described in [6]. This corpus size is around 10,000 sentences. These sentences were collected from Facebook, Twitter, YouTube, online news portals and product review pages.

---

[1] See Section VI for details

## II. Sentiment Analysis Corpus

### A. Data Source

The source of this data is the online Sports section of Prothom Alo, as shown in Fig 1 below.



Fig. 1 Daily Prothom Alo Online Edition.

Most of source sentences for the corpus were collected from the comments Section as shown in Fig 2 below.
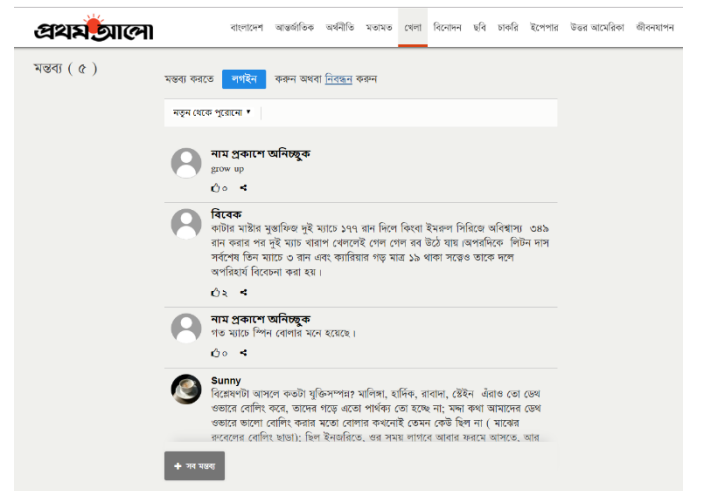


Fig. 2 Comments in the Sports Section.

### B. Methodology

The corpus was prepared in a combination of manual and automated steps. Initially the sentences appearing in the Sports Section were copied manually. These data were labeled by hand into three categories "positive", "negative" and "neutral", by a "Content Team" and then crosschecked by a "QA Team".

The truth labeling is done at two levels.

Sentence Level: The first level is the sentence level, as seen in Fig 3 below. In this case, the polarity applies to the overall sentence.



Fig. 3 Sentence level polarity.

- Manual collection and preprocessing: We collected more than 10,000 sentences from comments from online Bangla newspaper (https://www.prothomalo.com/), primarily from cricket sports news. We only included valid and complete single sentences. The task was distributed among the 5 members team members, who are all native Bangla speakers. Each member tagged sentiment polarity for the sentences allocated to him/her. Another member then crosschecked this. A third member then re-assigned the polarity sentences if the first and second members disagreed. The same methodology was applied to crosschecking sentence validation. In case there are disagreements, the assigned final polarity is at least assigned by two team members. If all three team members disagreed, the sentence was considered to be too ambiguous and dropped from the corpus.

- Automatic processing: We removed unwanted characters, words and symbols from the sentences, such as:
  - ['১','২','৩','৪','৫','৬','৭','৮','৯','1','2','3','4','5','6','7','8','9']
  - [' ','?',',',':',';','(',')','-','_','!','/','_','*','%','।','\','+','<','>','—','o','=']
  - ['"','"','|',':','…',')','`','@','#',',','','&','-','_',😊,'💋','💅',😀,😂]
  - [A-Z]
  - [a-z]

  We also removed duplicate sentences.

Word Level: The second is polarity on the word level, as shown in Fig 4 below.

- Automatic processing: We tokenized all collected sentences and removed numbers, digits, and symbols from the tokenized word list. We then identified the unique words from the word list of nearly 15,000

words. We stemmed the word list using two different stemmers, StemmerR[10] and StemmerP[11]. We then identified the words that produced the same result from these two stemmers.



Fig. 4 Word level polarity.

- Manual collection and preprocessing: We checked whether the un-stemmed word is already a root or not and manually corrected the roots for those words that were stemmed wrongly. Once a clean word list was created, we then tagged the polarity of each word manually, using the same three-tiered approach as described before. This step also resulted in identifying some words that were ambiguous. These are then dropped for the final corpus.

A snapshot of this corpus is shows in Fig 5 below.



Fig. 5 Corpus after stemming using two different algorithms.

III. CORPUS STATISTICS
This Section qualifies the corpus.

TABLE I
RAW DATA COLLECTED FOR THE SENTIMENT CORPUS

| Total number of sentences | 10,008 |
|---|---|
| Total number of words before filtering | 19,731 |

| Total number of words after filtering | 14,874 |
|---|---|
| How many ambiguous | 2,140 |
| How many words accepted | 12,734 |

Table 1 shows the statistics of raw data collected from the sentiment corpus.

TABLE 2
COMPARING THE TWO STEMMERS

| | StemmerP | StemmerR |
|---|---|---|
| How many words were stemmed | 14,874 | 14,874 |
| How many times the stemmer was able to stem a word | 14,662 | 14,874 |
| How many times the stemmer was not able to stem a word | 212 | 14,874 |

Table 2 compares the performance of the two stemmers. This step was completely manual and three-level peer-viewed.

We found that 1.5% times the stemmers did not agree with each other. In order to correct that, we manually fixed the stemming.

TABLE 3
MANUAL STEMMING

| Category | Number |
|---|---|
| Accuracy of StemmerP | 58.08% |
| Accuracy of StemmerR | 59.65% |
| How many words were manually stemmed | 5,138 |
| How many spelling were corrected manually | 682 |

Table 3 shows the results of manually fixing stemming issues.

After all this cleanup and manual fixes, the final corpus has the following entries, as seen in Table 4.

TABLE I
RAW DATA COLLECTED FOR THE SENTIMENT CORPUS

| | Positive | | Negative | | Neutral | |
|---|---|---|---|---|---|---|
| | Actual | % | Actual | % | Actual | % |
| Total Number of sentences | 3,183 | 33.0 | 4,110 | 42.67 | 2,337 | 24.27 |
| Total Number of words | 824 | 6.47 | 1,068 | 8.38 | 10,804 | 84.80 |



Fig. 6 WordCloud of corpus sentences with both positive and negative polarity.

## IV. CORPUS PROPERTIES

### A. WordCloud of Collected Sentences

Fig. 6 shows a WordCloud of the collected sentences with positive and negative polarity. Please note that we have used a stop word list to filter the words before this was created. We have identified 398 stop words. The same applies to building WordClouds for negative and neutral sentences.

Fig. 7 shows a WordCloud of the collected sentences with positive polarity.



Fig. 7 WordCloud of corpus sentences with a positive polarity.

Fig. 8 shows a WordCloud of the collected sentences with negative polarity.



Fig. 8 WordCloud of corpus sentences with a negative polarity.



Fig. 9 Word frequency of top 20 words

### B. Frequency Distribution of Top 20 Words

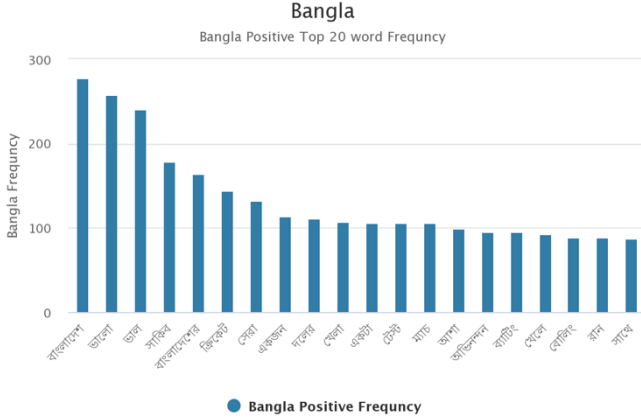Fig. 9 shows the frequency of the top 20 words in our corpus.



Fig. 10 Word frequency of top 20 positive words

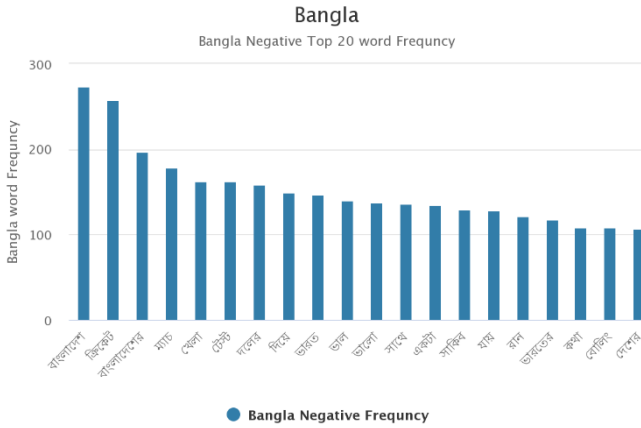Fig. 10 shows the frequency of the top 20 positive words in our corpus.



Fig. 11 Word frequency of top 20 negative words

Fig. 11 shows the frequency of the top 20 negative words in our corpus.

## V. SOME OBSERVATIONS

We have created this corpus from a very focused source, the sports news domain and have incorporated sentences, unique words and stemmed words. We extensively cleaned the data using a combination of filtering and stop word list — employing both manual and automated process. Every entry is cross-validated using at least two, and sometimes three annotators. We have manually corrected misspelling and stemming errors. So this is not just an annotated and ground truth corpus on sentiment analysis, it is also a corpus for training stemming engines.

It was also very important for us to build in auditability in the corpus. That is why every word and root word is cross-referenced against the source sentences. This way this corpus can be adopted for other NLP solutions with ease.

The other aspect of our corpus design is the transparency of the data collection process. It is a natural extension of the auditability of the data mentioned above.

## VI. OPEN SOURCE LICENSING

Not-For-Profit and academic organizations and government agencies may use this corpus for noncommercial linguistic research and education only. For-profit organizations may use this corpus after signing a commercial technology development contract. Not-For-Profit and academic organizations and government agencies cannot use this corpus to develop or test products for commercialization, nor can they use this in any commercial product or for any commercial purpose.

## VII. CONCLUSIONS AND FUTURE WORK

We have presented a Bangla corpus specifically targeted for sentiment analysis. We described the methodology, source and clean up process. In the future, we plan to extend the corpus to support aspect based sentiment analysis for Bangla, where different clauses of a single sentence may have different sentiments. We also plan to extend this by adding sentiments for phrases, idioms and clauses. In addition, we plan to offer a set of machine learning models that can use this corpus.

## REFERENCES

[1] Adrija Roy and Abhishek Anand Singh. Sentimental Analysis (Bengali) https://github.com/abhie19/Sentiment-Analysis-Bangla-Language.

[2] English & Bengali Online Dictionary & Grammar. http://www.english-bangla.com/bntoen/index/%E0%A6%9C%E0%A6%9F%E0%A6%BF%E0%A6%B2

[3] Tazim Hoque. Word and Doc2Vec file for Bengali Sentiment Analysis. https://www.kaggle.com/tazimhoque/bengali-sentiment-text/.

[4] Atik Rahman. Bangla Aspect Based Sentiment Analysis Dataset. https://github.com/AtikRahman/Bangla_ABSA_Datasets.

[5] Md. Atikur Rahman and Emon Kumar Dey. Datasets for Aspect-Based Sentiment Analysis in Bangla and its Baseline Evaluation. 4 May 2018. Institute of Information Technology, University of Dhaka, Dhaka 1000, Bangladesh. https://res.mdpi.com/data/data-03-00015/article_deploy/data-03-00015.pdf?filename=&attachment=1

[6] Asif Hassan, Mohammad Rashedul Amin, Abul Kalam Al Azada, Nabeel Mohammed. Sentiment Analysis on Bangla and Romanized Bangla Text (BRBT) using Deep Recurrent models. 24 Nov 2016. Dept. of Computer Science and Engineering University of Liberal Arts Bangladesh (ULAB), Bangladesh. https://docs.google.com/viewerng/viewer?url=http://resources.apurbatech.com/publication/upload/1610.00369.pdf

[7] D. Das and S. Bandyopadhyay, Developing Bengali WordNet Affect for Analyzing Emotion. Int. Conf. on the Computer. Processing of Oriental Languages, pp. 35-40, 2010.

[8] Cliff Goddard. Natural Language Processing, Edition: 2nd edition, Chapter: 5, Publisher: CRC Press, Taylor & Francis, Editors: Nitin Indurkhya, Fred J. Damerau, pp.92-120

[9] Mohammad Daniul Huq. Semantic values in Translating from English to Bangla, The Dhaka University Journal of Linguistics: Vol. 1 No.2 August, 2008 Pages: 45-66.

[10] Rafi Kamal. Bangla Stemmer. https://github.com/rafi-kamal/Bangla-Stemmer.

[11] Sazedul Islam. Rule based Bengali Stemmer written in python. https://pypi.org/project/py-bangla-stemmer/.