

# A BERT-Based Ensemble Learning Approach for Sentiment Classification in Twitter<sup>\*</sup>

Shupe Li and Ziyi Xu

LIACS, Leiden University, Leiden, Netherlands  
{s3430863, s3649024}@umail.leidenuniv.nl

**Abstract.** In this project, we present an ensemble learning method for sentiment classification. We use the state-of-the-art BERT as base models and a soft voting classifier as the meta model. The effectiveness of our proposed approach is verified on SemEval-2017 dataset, which quantifies the sentiment in tweets via a three-point ordinal scale. Experimental results show that the overall performance of our proposed model is better than that of baselines.

**Keywords:** Sentiment Analysis · BERT · Ensemble Learning

## 1 Introduction

Sentiment analysis, a growing field today, is the process of analyzing pieces of writing to determine the emotional tone they carry. In simple words, sentiment analysis helps to find the author’s attitude. This method can be used by businesses to analyze product reviews and feedback, especially for social media companies with large information streams because of the wealth of data they generate. Researchers also have a special interest in social media data because of their easy availability and rapid change.

SemEval is an International Workshop on Semantic Evaluation, formerly SenseEval. It is an ongoing series of evaluations of computational semantic analysis systems. The SemEval-2017 Task 4 focuses on classifying and quantifying the sentiment of tweets. The task was included in this workshop in the previous year [6], and Sentiment Analysis in Twitter has been run yearly since 2013 [7]. The subtask A of this task is the problem we try to unravel in this article. It is about deciding the overall sentiment of tweets and marking them on a three-point ordinal scale.

Our proposed model adopts two strategies with the goal of achieving higher accuracy. One is choosing more powerful BERT and its variants as base models instead of CNN, LSTM, or SVM. The other is that we design an ensemble learning approach to integrate the advantages of different BERT models, which encourages learning diversity and might enhance model performance. The classification task is divided into two stages. In the first stage, we train a series of base models to generate corresponding predictions. These predictions are the inputs

---

<sup>\*</sup> Text Mining, Master CS, Fall 2022, Leiden, the Netherlands

of the meta model in the second stage. The final predictions are the outputs of the meta model.

The remainder of the paper is structured as follows. Section 2 discusses previous work related to sentiment classification. Following that, Section 3 introduces datasets provided by SemEval. We explain the data preparation workflow, our proposed approach, and baselines in Section 4. Section 5 describes the procedure of experiments in detail and reports the experimental results. The paper is concluded with a discussion of the results and a brief summary. We also attach the contributions of group members at the end of the paper.

## 2 Related Work

The prevalence of social media has been driving the growth of sentiment analysis research since early 2000 [11]. Earlier studies mainly focused on the lexicon-based approach, which discovers sentimental polarity based either on a predefined dictionary or on some semantic methods [5]. Previous works have also intensively researched the applications of traditional machine learning methods in sentiment analysis, such as SVM, naive Bayes, etc [10]. However, both lexicon-based and machine learning approaches require manual feature selection, which limits their applicability in large and dynamic user-generated content like tweets.

In recent years, deep learning methods have been introduced into sentiment analysis field and have achieved remarkable performance on many benchmark datasets [10]. Popular deep learning frameworks for sentiment classification include CNNs and RNNs. Inspired by the computer vision field, the CNN architecture first constructs a sentence matrix using tokens and word vectors. Then it treats the sentence matrix as an "image" and applies regular convolution operations on the matrix [12]. The major drawback of CNNs is that CNNs can't capture the hierarchical relationship between local features. Besides, the pooling layer may cause the loss of critical information and harm classification accuracy. LSTMs and GRUs are two commonly used RNN variants. Nowak et al. [8] compared LSTM, bidirectional LSTM, and GRU with other algorithms on real-world datasets and found RNNs outperformed baselines in sentiment classification task. However, RNNs are not suitable for processing long sequences and the training process is very slow. There were also attempts to combine CNNs and RNNs, aiming at enhancing model performance further [3][2]. But these hybrid models don't overcome the intrinsic weaknesses of CNNs and RNNs.

In this project, we consider BERT, the state-of-the-art model in several natural language processing tasks. Compared to CNNs and RNNs, BERT leverages the attention mechanism to model the complex relationship between features in short-distance as well as long-distance contexts. And it is efficient on GPUs and TPUs [1]. We describe our method in detail in Section 4.

### 3 Data

The dataset consists of 11 files with tweets from 2013 to 2015. Tweets are marked with sentiment labels on a three-point scale {Positive, Neutral, Negative}. Each tweet corresponds to one row in datasets following a fixed format: [id, sentiment label, text].

The criterion of tweet selection is covering popular topics at the time of sending tweets. Datasets are downloaded via the Twitter API and have been preprocessed by workshop in advance, where three kinds of data have been removed: repeated tweets, the bag-of-words cosine similarity exceeded 0.6, and topics with less than 100 tweets. CrowdFlower is used to create all annotations on both the training set and the test set. Each tweet is annotated by at least five people to ensure the accuracy of annotations. Another main quality control measure is performing hidden tests to filter out unqualified annotations. There are also manual inspections on pilot runs aiming at adjusting the annotation instructions dynamically. Table 1 summarizes the descriptive statistics of 11 files.

**Table 1.** Descriptive Statistics of Datasets

Dataset	Positive	Neutral	Negative	Total
2013train	3,640	4,586	1,458	9,684
2013dev	575	739	340	1,654
2013test	1,475	1,513	559	3,547
2014sarcasm	20	7	22	49
2014test	982	669	202	1,853
2015train	170	253	66	489
2015test	1,038	987	365	2,390
2016train	3,017	2,001	850	5,868
2016dev	829	746	391	1,966
2016devtest	994	681	325	2,000
2016test	7,059	10,342	3,231	20,632
Total	19,799	22,524	7,809	50,132

## 4 Methodology

### 4.1 Data Preparation

We perform data cleaning and dataset division before modelling. Notice that records in 2016test file have a redundant \t before \n. We remove the unnecessary \t to ensure the dataset would be read into Python appropriately. To create training set and test set, we concatenate all 11 files and randomly sample 80% records as training set while treating the others as test set. As a result, there are 40,105 records in training set and 10,027 records in test set. We also

create the validation set for hyperparameter tuning by randomly sampling 20% records from the training set. The random seed is set to 42 in our experiments to ensure reproducibility.

## 4.2 BERT

BERT [1] is a popular NLP model that has achieved remarkable performance on various tasks, such as text classification, question answering, etc. Compared to traditional RNN models, it encodes sequences in both directions instead of following a left-to-right or right-to-left routine, which is closer to how humans understand the meaning of the text. Its bidirectional encoding ability is accomplished by the transformer architecture, which is a stack of encoders using a multi-head attention mechanism. Specifically, encoders process the entire sequence at once and use layer-wise tensor operations to learn relationships between words in a sentence. This design not only encodes the inputs bidirectionally but also eliminates the possible local bias, for it gives equal importance to the local context and the long-distance context. It is worth mentioning that the training process of BERT is more efficient than RNN due to the feasibility of parallelization.

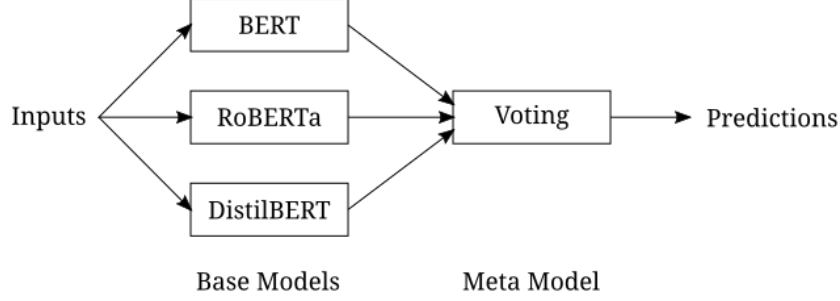
BERT requires a special format of input called WordPiece. The WordPiece tokenizer splits words into tokens and adds special tokens at the beginning as well as the end of the sentence. Preprocessed inputs provide three aspects of information for the BERT model: tokens, sentence segments, and the absolute position. We perform the tokenization operation on both the training set and the test set before modelling.

We also consider two variants of BERT in our project – RoBERTa [4] and DistilBERT [9]. RoBERTa is an optimized version of the BERT model. Authors of RoBERTa find that the original BERT is actually undertrained after reproducing BERT experiments. Their solution is increasing the training epochs of BERT and carefully selecting hyperparameters, which significantly improves the model performance in practice. On the contrary, DistilBERT is a compressed version of BERT. The main idea of DistilBERT is to reduce the model size via the distillation technique. Distillation consists of a teacher model and a student model. The teacher model is trained on a large dataset and is fine-tuned to maximize accuracy. However, many features learned by the teacher model are redundant for a specific task. Therefore, we can train a much smaller student model to focus on key features and imitate the output of the teacher model. Experimental results show that DistilBERT is cheaper to train while maintaining comparable performance to BERT.

## 4.3 Ensemble Learning

Ensemble learning refers to methods that combine multiple models to achieve better performance in machine learning. It encourages base models to learn different aspects of the data to reduce errors and avoid being entrapped in local

optima. In the project, sentiment analysis on Twitter is a multi-class classification problem. And we develop a classification voting ensemble model integrated with BERT and its variants. Figure 1 illustrates the architecture of our model.



**Fig. 1.** The Architecture of the Proposed Model

Each base model outputs the class prediction for each record. Here we consider two strategies for the meta model — hard voting and soft voting. Hard voting means the voting classifier adopts the plurality voting strategy to generate the final prediction. In other words, the final prediction is the class label received the majority votes from base models. If all classes have the same votes, the voting classifier will choose a class label randomly as the output. However, hard voting ignores the probability information related to each model’s prediction and is likely to be affected by relatively weaker base models. Thus, we set the hard voting ensemble learning as one of our baselines and adopt the soft voting strategy.

The soft voting strategy averages the output probabilities of base models and selects the class with the highest probability. However, the outputs of base models may contain negative numbers and are not normalized. We convert the outputs into a probability distribution by applying the softmax function:

$$P(c_i) = \frac{\exp(\hat{y}_i)}{\sum_i \exp(\hat{y}_i)}$$

where  $P(c_i)$  is the estimated probability of class  $i$  and  $\hat{y}_i$  is the output value of the base model for class  $i$ . In summary, our proposed model is an ensemble of three BERT base learners and a soft voting classifier.

#### 4.4 Baselines

We compare the proposed method against four baselines. The first three benchmarks are BERT, RoBERTa, and DistilBERT alone. We set the ensemble BERT model with hard voting strategy as the fourth baseline, as mentioned in Section 4.3. Accuracy, macro recall, and macro F1 are selected as metrics in experiments.

## 5 Experiments

### 5.1 Experimental Setup

We replace three patterns in tweets and tokenize the textual data before feeding them into models. The first pattern is `@user_name`, which indicates other users on Twitter. We substitute it with a single space. The second is to replace `&amp;` with `&`. And the last is to delete redundant spaces. We also map the text labels into categorical variables to meet the requirements of APIs, i.e., setting Negative as 0, Neutral as 1, and Positive as 2. We apply Hugging Face’s implementation of BERT models. Table 2 summarizes the software environment of our experiments.

**Table 2.** Summary of Software Environment

Item	Version	Function
Python	3.9	Programming language.
Transformers	4.25.1	Transformer models in Hugging Face.
TensorFlow	2.11.0	Deep learning APIs.
scikit-learn	1.2.0	Preprocessing and evaluation.

Our experimental process consists of hyperparameter tuning, training, and evaluation. Grid search technique is used to optimize hyperparameters. The best hyperparameter setting is `batch size= 16`, `initial learning rate=  $10^{-5}$` . We combine early stopping and model checkpoint saving strategies to enhance the model performance further. Specifically, we set the number of epochs as 50, stop the training if the accuracy doesn’t improve in 3 epochs, and load the best model when evaluating. All experiments are deployed on a cloud server with an Intel Xeon(R) Gold 6330 CPU and an RTX A5000 GPU.

### 5.2 Results

Table 3 reports the experimental results. The detailed evaluation report for each class is attached as Appendix.

**Table 3.** Summary of Results

Model	Precision	Macro Recall	Macro F1
BERT	0.7338	0.7169	0.7181
RoBERTa	0.7473	0.7236	0.7310
DistilBERT	0.7226	0.6920	0.7002
Ensemble BERT with Hard Voting	0.7458	0.7229	0.7289
Ensemble BERT with Soft Voting	<b>0.7510</b>	<b>0.7305</b>	<b>0.7347</b>

Our proposed model outperforms benchmarks on all metrics, which supports the effectiveness of our method. For the base model, standard BERT performs worse than RoBERTa but better than DistilBERT, which is consistent with the theory. It is worth noting that the ensemble BERT model with hard voting achieves slightly worse scores on all metrics than RoBERTa alone. We think the reason is that the performance of hard voting in this task is affected by DistilBERT which is a weaker model compared to the other two.

## 6 Discussion

Our proposed model is inspired by BERT and ensemble learning. In experiments, we fine-tune BERT models and investigate the performance of the individual BERT as well as the ensemble one. The overall computational cost is reasonable, for we can train all models with only one available GPU by leveraging pre-trained BERT from Hugging Face. And the runtime of fine-tuning doesn't have much difference among BERT models. Considering that most tweets are related to topics in daily life, pre-trained BERT models generalize well in this task. However, if the dataset involves lots of terminologies in a specific domain, we need to search for the domain-specific BERT instead of directly downloading the BERT model trained on general domain text. In rare cases, there may be no suitable pre-trained BERT model. Then we are supposed to train BERT from scratch, which is time-consuming and computationally expensive.

Ensemble learning is the key to improving the performance of our proposed method. Generally, different models have different strengths and weaknesses. The idea of ensemble learning is that combining predictions from different models may increase learning diversity and encourage models to complement each other. In experiments, applying the proper ensemble learning technique enhances the model performance significantly. Besides, a comparison between the experimental result of hard voting and that of soft voting illustrates the importance of model combination scheme selection. An ideal combination strategy should minimize the information loss during model aggregation.

## 7 Conclusion

In this paper, we proposed a BERT-based ensemble learning method to classify tweets according to sentiment. We follow the standard sentiment analysis workflow — hyperparameter tuning, training, and evaluation. To improve model performance, we also adopt training strategies like early stopping. Experimental results support the effectiveness of our model. Due to the limited computational resources, we search a relatively small hyperparameter space and don't try other ensemble learning methods such as bagging and boosting, which are possible directions for future works.

## 8 Contributions

Authors contribute equally to the programming and the report.

## References

1. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. ArXiv **abs/1810.04805** (2019)
2. Hassan, A., Mahmood, A.: Deep learning approach for sentiment analysis of short texts. 2017 3rd International Conference on Control, Automation and Robotics (ICCAR) pp. 705–710 (2017)
3. Huang, Q., Chen, R., Zheng, X., Dong, Z.: Deep sentiment representation based on cnn and lstm. 2017 International Conference on Green Informatics (ICGI) pp. 30–33 (2017)
4. Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., Stoyanov, V.: Roberta: A robustly optimized bert pretraining approach. ArXiv **abs/1907.11692** (2019)
5. Medhat, W., Hassan, A., Korashy, H.: Sentiment analysis algorithms and applications: A survey. Ain Shams Engineering Journal **5**(4), 1093–1113 (2014)
6. Nakov, P., Ritter, A., Rosenthal, S., Sebastiani, F., Stoyanov, V.: SemEval-2016 task 4: Sentiment analysis in Twitter. In: Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016). pp. 1–18. Association for Computational Linguistics, San Diego, California (Jun 2016)
7. Nakov, P., Rosenthal, S., Kozareva, Z., Stoyanov, V., Ritter, A., Wilson, T.: SemEval-2013 task 2: Sentiment analysis in Twitter. In: Second Joint Conference on Lexical and Computational Semantics (\*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013). pp. 312–320. Association for Computational Linguistics, Atlanta, Georgia, USA (Jun 2013)
8. Nowak, J., Taspinar, A., Scherer, R.: Lstm recurrent neural networks for short text and sentiment classification. In: Artificial Intelligence and Soft Computing. pp. 553–562. Springer International Publishing, Cham (2017)
9. Sanh, V., Debut, L., Chaumond, J., Wolf, T.: Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. ArXiv **abs/1910.01108** (2019)
10. Yadav, A., Vishwakarma, D.K.: Sentiment analysis using deep learning architectures: A review. Artif. Intell. Rev. **53**(6), 4335–4385 (Aug 2020)
11. Zhang, L., Wang, S., Liu, B.: Deep learning for sentiment analysis: A survey. WIREs Data Mining and Knowledge Discovery **8**(4), e1253 (2018)
12. Zhang, Y., Wallace, B.: A sensitivity analysis of (and practitioners’ guide to) convolutional neural networks for sentence classification (2015)

## Appendix: Detailed Evaluation Report



**Table 4.** Classification Report of BERT

Class	Precision	Recall	F1
Negative	0.6628	0.6464	0.6545
Neutral	0.7233	0.7175	0.7204
Positive	0.7723	0.7868	0.7795
Macro Avg.	0.7195	0.7169	0.7181
Weighted Avg.	0.7332	0.7338	0.7335

**Table 5.** Classification Report of RoBERTa

Class	Precision	Recall	F1
Negative	0.7129	0.6280	0.6678
Neutral	0.7287	0.7394	0.7340
Positive	0.7795	0.8034	0.7913
Macro Avg.	0.7404	0.7236	0.7310
Weighted Avg.	0.7464	0.7473	0.7463

**Table 6.** Classification Report of DistilBERT

Class	Precision	Recall	F1
Negative	0.6584	0.5754	0.6141
Neutral	0.6974	0.7479	0.7218
Positive	0.7771	0.7526	0.7647
Macro Avg.	0.7110	0.6920	0.7002
Weighted Avg.	0.7229	0.7226	0.7218

**Table 7.** Classification Report of Ensemble BERT with Hard Voting

Class	Precision	Recall	F1
Negative	0.6945	0.6324	0.6620
Neutral	0.7257	0.7475	0.7364
Positive	0.7875	0.7888	0.7881
Macro Avg.	0.7359	0.7229	0.7289
Weighted Avg.	0.7453	0.7458	0.7452

**Table 8.** Classification Report of Ensemble BERT with Soft Voting

Class	Precision	Recall	F1
Negative	0.6967	0.6464	0.6706
Neutral	0.7396	0.7374	0.7385
Positive	0.7827	0.8077	0.7950
Macro Avg.	0.7397	0.7305	0.7347
Weighted Avg.	0.7499	0.7510	0.7502