



# An Analysis of Logic Rule Dissemination in Sentiment Classifiers

Shashank Gupta<sup>(✉)</sup>, Mohamed Reda Bouadjeneq,  
and Antonio Robles-Kelly

School of Information Technology, Deakin University, Waurn Ponds Campus,  
Geelong, VIC 3216, Australia  
guptashas@deakin.edu.au

**Abstract.** Disseminating and incorporating logic rules in deep neural networks has been extensively explored for sentiment classification. Methods that are proposed for that goal rely on a component that aims to capture and model logic rules, followed by a sequence model to process the input sequence. While these methods claim to effectively capture syntactic structures that affect sentiment, they only show improvement in terms of accuracy to support their claims with no further analysis. Focusing on the *A-but-B* rule, we use the PERCY metric (a recently developed Post-hoc Explanation-based score for logic Rule dissemination ConsistencY assessment) to analyze and study the ability of these methods to identify the *A-but-B* structure, and to make their classification decision based on the *B* conjunct. PERCY proceeds by estimating feature attribution scores using LIME, a model-agnostic framework that aims to explain the predictions of any classifier in an interpretable and faithful manner. Our experiments show that (a) accuracy is misleading in assessing these methods, (b) not all these methods are effectively capturing the *A-but-B* structure, (c) often, the underlying sequence model is what captures the syntactic structure, and (d) the best method classifies less than 25% of test examples based on the *B* conjunct.

**Keywords:** Sentiment Classification · Logic Rules · Explainable AI

## 1 Introduction

Methods of disseminating and incorporating logic rules in Deep Neural Networks have been extensively explored for sentiment classification. The two main methods developed for that purpose are: (i) Iterative Knowledge Distillation method [1] and (ii) the Contextualized Word Embeddings approach [2]. Briefly, these methods rely on a component aimed at capturing and modeling logic rules (e.g., the teacher network in the Iterative Distillation method and the ELMO model [3] in the Contextualized Word Embeddings approach), followed by a sequence model to process the input sequence, (e.g., a RNN).

The authors of these two methods claim that they effectively capture syntactic structures in the input sentence that affect its sentiment, but they have

only used the improvement in terms of accuracy to support their claim with no further analysis. However, achieving a high classification accuracy does not necessarily mean that a method has effectively captured and encoded rules and other textual syntactic structures. For example, let’s consider the sentence “*the casting was not bad but the movie was awful*” that has an *A-but-B* structure – a component *A* followed by *but* which is then followed by a component *B*. In this example, the conjunction is interpreted as an argument for the second conjunct, with the first functioning concessively [4–6]. While a sentiment classifier can correctly identify that this sentence has a negative sentiment, it may fail to infer it’s decision based *exclusively* on the *B* part of the sentence (i.e., “*the movie was awful*”), but instead, it may base it’s decision on individual negative words also present in Part *A* (i.e., “*bad*”).

While focusing on the *A-but-B* syntactic structure and sentiment classification, we propose in this paper to study the ability of the aforementioned methods to: (i) effectively identifying the *A-but-B* structure in an input sentence, and to (ii) make their classification decision based on the *B* conjunct of a sentence. Specifically, we rely on the PERCY metric [7], a recently developed Post-hoc Explanation-based score for logic Rule dissemination Consistency assessment. PERCY estimates feature attribution scores using LIME [8], a model-agnostic framework that aims to explain predictions of any classifier in an interpretable and faithful manner. We validate our findings with an exhaustive experimental evaluation using the SST2 dataset [6] by testing various sentiment classifiers designed for logic rules dissemination. Among numerous findings, we show that: (a) accuracy is misleading in assessing methods for capturing logic rules, (b) not all methods are effectively capturing the *A-but-B* structure, (c) their sequence model is often what captures the syntactic structure, and (d) the best method bases its decision on the *B* conjunct in less than 25% of test examples.

## 2 Logic Rule Dissemination Methods

In this section, we first describe the neural network architecture we use for sequence modeling, before discussing the main methods we analyse for logic rules dissemination in that architecture.

### 2.1 Network Architecture

The backbone neural network [9,10] we use throughout this paper is depicted in Fig. 1. Three 1D CNN sequence layers (kernel size of 3, 4, and 5) process the word embeddings of an input sequence in parallel in order to extract diverse features and pass the concatenated features into a feed-forward binary classification layer with a sigmoid activation to extract the sentiment of the input sentence – 0 for a negative sentiment and 1 for a positive

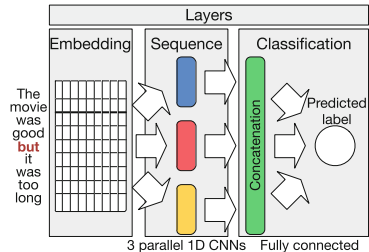


Fig. 1. Neural network.

sentiment. In the next subsections, we will discuss the methods we analyze in this article that aim to incorporate and disseminate logic rules in the neural network architecture depicted in Fig. 1.

## 2.2 Iterative Rule Knowledge Distillation

The Iterative rule knowledge distillation method proposed by Hu et al. [1] aims to transfer the domain knowledge encoded in first order logic rules into a neural network defined by a conditional probability  $p_\theta(y|x)$  where  $\theta$  is a parameter to learn. Specifically, during training, a posterior  $q(y|x)$  is constructed by projecting  $p_\theta(y|x)$  into a subspace constrained by the rules to encode the desirable properties, by using the following loss:

$$\begin{aligned} \min_{q, \xi \geq 0} \quad & KL(q(y|x)||p_\theta(y|x)) + C \sum_{x \in X} \xi_x \\ \text{s.t.} \quad & (1 - \mathbb{E}_{y \leftarrow q(\bullet|x)}[r_\theta(x, y)]) \leq \xi_x \end{aligned}$$

where  $q(y|x)$  denotes the distribution of  $(x, y)$  when  $x$  is drawn uniformly from the train set  $X$  and  $y$  is drawn according to  $q(\bullet|x)$ , and  $r_\theta(x, y) \in [0, 1]$  is a variable that indicates how well labeling  $x$  with  $y$  satisfies the rule. The closed form solution for  $q(y|x)$  is used as soft targets to imitate the outputs of a rule-regularized projection of  $p_\theta(y|x)$ , which explicitly includes rule knowledge as regularization terms.

Next, the rule knowledge is transferred to the posterior  $p_\theta(y|x)$  through knowledge distillation optimization objective:

$$(1 - \pi) \times \mathcal{L}(p_\theta, P_{true}) + \pi \times \mathcal{L}(p_\theta, q)$$

where  $P_{true}$  denotes the distribution implied by the ground truth,  $\mathcal{L}(\bullet, \bullet)$  denotes the cross-entropy function, and  $\pi$  is a hyperparameter that needs to be tuned to calibrate the relative importance of the two objectives. Overall, the Iterative rule knowledge distillation method is agnostic to the network architecture, and thus is applicable to general types of neural models such as the one depicted in Fig. 1.

## 2.3 Contextual Word Embeddings

Traditional word embeddings methods like Word2Vec [11] and Glove [12] do not capture the local context of the word in a sentence. However, language is complex and context can completely change the meaning of a word in a sentence. Hence, contextual word embeddings methods have emerged as a way to capture the different nuances of the meaning of words given the surrounding text. Krishna et al. [2] have advocated that contextualized word embeddings might capture logic rules and thus disseminate that latent information in the 1D CNN sequence models of the neural network in Fig. 1. In the following, we briefly review two of the main contextual word embedding methods we use in our experiments.

**ELMo:** stands for Embeddings from Language Models is a pre-trained model developed by Peters et al. [3]. Instead of using a fixed embedding for each word, ELMo looks at the entire sentence before assigning each word in it an embedding. It uses a bi-directional LSTM trained on a specific task to be able to create those embeddings. Krishna et al. [2] proposed to use ELMo in their method.

**BERT:** stands for Bidirectional Encoder Representations from transformers. This is also a pre-trained model developed by Devlin et al. [13]. Briefly, the BERT is a model based on Encoder Transformer blocks [14], which processes each element of the input sequence by incorporating and estimating the influence of other elements in the sequence to create embeddings.

To further test the hypothesis proposed by Krishna et al. [2], we conduct experiments with two different context-free word embeddings namely Word2vec developed by Mikolov et al. [11] and Glove developed by Pennington et al. [12] in which each token is mapped to a unique vector independent of its context. These word embeddings are used as an ablation study to analyze the effectiveness of the rule knowledge distillation method discussed in the previous section.

### 3 Methodology

As mentioned earlier, our main goal in this paper is to assess each sentiment classifier for its ability to correctly classify a test example with an *A-but-B* structure only on the basis of the *B* conjunct. For this purpose, we use a metric called **PERCY** [7], which stands for *Post-hoc Explanation-based Rule Consistency assessment Score*. Specifically, given a sentence  $S$  which is an ordered sequence of terms  $[t_1 t_2 \dots t_n]$ , PERCY relies on LIME to assign a weight  $w_n$  to each term  $t_n$  in  $S$  where a positive weight indicates that  $t_n$  contributes and supports the positive class, and a negative weight indicates how much  $t_n$  supports the negative class. In order to estimate how much a term  $t_n$  contributes to the final decision of the classifier, PERCY normalizes its weight as follows:

$$\tilde{w}_n = \begin{cases} w_n \times P(y = 1 | S), & \text{if } w_n \geq 0 \\ |w_n| \times P(y = 0 | S), & \text{otherwise} \end{cases} \quad (1)$$

where  $P(y = c | S)$  is the probability to predict class  $c$  given sentence  $S$ . Hence, every sentence in our test set is mapped to a vector  $[\tilde{w}_1 \tilde{w}_2 \dots \tilde{w}_n]$  with  $\tilde{w}_n$  indicating how much the word  $t_n$  contributed to the final decision of the classifier. Next, given a sentence that contains an *A-but-B* structure, PERCY defines the normalized weights  $\tilde{W}(A) = [\tilde{w}_0 \dots \tilde{w}_{i-1}]$  and  $\tilde{W}(B) = [\tilde{w}_{i+1} \dots \tilde{w}_n]$  as respectively the left and right sub-sequences w.r.t the word “*but*” indexed by  $i$ . Finally, PERCY computes an expectation over weights as follows:  $\mathbb{E}_A(W) = \sum_{\tilde{w}_k \in \tilde{W}(A)} \tilde{w}_k$  and  $\mathbb{E}_B(W) = \sum_{\tilde{w}_k \in \tilde{W}(B)} \tilde{w}_k$ , and concludes that a classifier has based its classification prediction by relying on the *B* conjunct if:  $\mathbb{E}_B(W) > \mathbb{E}_A(W)$  **and**  $p$ -value  $\leq 0.05$  – this condition aims to make sure that the observed difference is statistically significant.

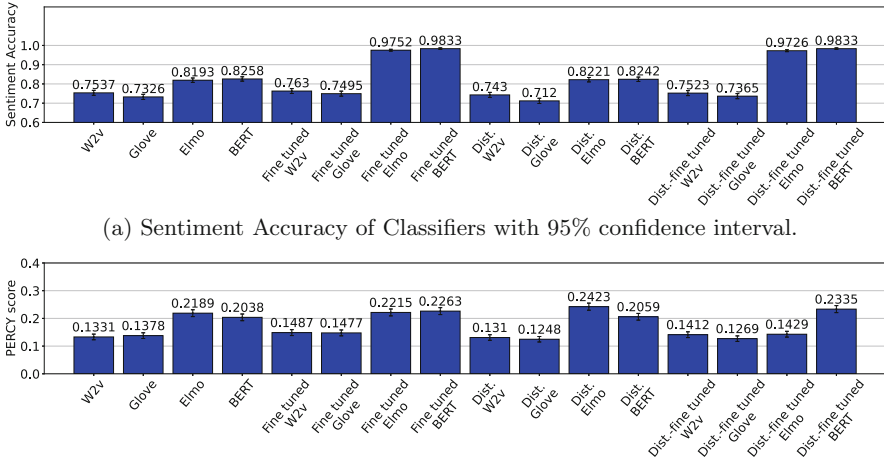


Fig. 2. Results obtained on SST2 dataset.

## 4 Experimental Evaluation

In this section, we first describe the dataset we have used in our evaluation before discussing the obtained results.

### 4.1 Dataset

Our experiments (as well as those presented in Hu et al. [1] and Krishna et al. [2]) are based on the Stanford Sentiment Treebank (SST2) dataset [6], which is a binary sentiment classification dataset. The dataset consists of 9,613 single sentences extracted from movie reviews, where sentences are labelled as either positive or negative each accounting for about 51.6% and 48.3%. A total of 1,078 sentences contain the  $A$ -but- $B$  syntactic structure which accounts for about 11.2% of the dataset. We report our results only on test examples that contain an  $A$ -but- $B$  syntactic structure to demonstrate the ability of a classifier to capture  $A$ -but- $B$  pattern. Hence, all classifier are trained, tuned, and tested using stratified nested  $k$ -fold cross-validation and evaluated primarily according to accuracy. These sentences are identified simply by searching for the word “but” as proposed in [1, 2, 6].

### 4.2 Performance Evaluation

In this section, we discuss the results of our analysis of logic rules dissemination methods in sentiment classifiers. The configuration options that were considered

are the following:  $\{\text{Word2vec, Glove, ELMo, BERT}\} \times \{\text{Static, Fine-tuning}\} \times \{\text{no distillation, distillation}\}$ , which gives a total of 16 classifier analysed on sentences with an *A-but-B* structure. To summarize all the results obtained over all the above configurations, Figs. 2a and 2b show the accuracy and the ability of the methods to base their classification decisions on the *B* conjunct. From these results, we make the following observations:

**Accuracy Analysis:** In Fig. 2a, we observe that the distillation model described in Hu et al. [1] is ineffective as it gives almost no improvement in terms of accuracy as also noted in [2]. Second, we note that fine-tuning all embeddings provides a statistically significant improvement of accuracy for almost all methods. Finally, it is clear that the best method is BERT, followed by ELMo, followed by either Glove or Word2vec.

**Rule Dissemination Analysis:** In Fig. 2b we show the proportion of test examples that have been *correctly* classified based on the *B* conjunct using PERCY score described in Sect. 3. Briefly, we first observe that for all methods, less than 25% of the test examples are effectively classified based on the *B* conjunct, which shows that the intent of these methods as described by their authors in [1, 2] is far from being achieved. This suggests that there is still a lot of research to be done on this NLP topic. Second, we again note that there is almost no improvement between for instance Word2vec with and without distillation (Figs. 2a and 2b), which simply suggests that in [1] it is the 1D CNN sequence model that is capturing to some extent the *A-but-B* structure. Finally, we note that some models although have higher sentiment accuracy perform poorly on rule dissemination performance and vice-versa. For example, Dist. Elmo and Dist. BERT have similar sentiment accuracy in Fig. 2a but Dist. Elmo outperforms Dist. BERT by a statistically significant margin on rule dissemination performance in Fig. 2b. Similar phenomenon can be observed for Dist. fine-tuned Elmo and BERT models where later outperforms former even though having similar sentiment accuracy. This indicates that accuracy is misleading and there is no correlation between sentiment accuracy and actual rule dissemination performance.

## 5 Conclusion

This paper gives an analysis and a study of logic rules dissemination methods on their ability to identify *A-but-B* structures while making their classification decision based on the *B* conjunct. We use a rule consistency assessment metric called PERCY for that goal. Our experimental evaluation shows that (a) accuracy is misleading to assess whether the classifier based its decision as per *B* conjunct (b) not all methods are effectively capturing *A-but-B* structure, (c) that their underlying sequence model is often the one that captures to some

extent the syntactic structure, and (d) that for the best method, less than 25% of test examples are effectively classified based on the *B* conjunct, indicating that a lot of research needs to be done in this topic.

## References

1. Hu, Z., Ma, X., Liu, Z., Hovy, E., Xing, E.: Harnessing deep neural networks with logic rules. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, Berlin, Germany, August 2016, vol. 1: Long Papers, pp. 2410–2420. Association for Computational Linguistics (2016)
2. Krishna, K., Jyothi, P., Iyyer, M.: Revisiting the importance of encoding logic rules in sentiment classification. In: Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October–November 2018, pp. 4743–4751. Association for Computational Linguistics (2018)
3. Peters, M.E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., Zettlemoyer, L.: Deep contextualized word representations. In Proceedings of NAACL (2018)
4. Lakoff, R.: If’s, and’s and but’s about conjunction. In: Fillmore, C.J., Langendoen, D.T. (eds.) Studies in Linguistic Semantics, Irvington, pp. 3–114 (1971)
5. Blakemore, D.: Denial and contrast: a relevance theoretic analysis of “but”. *Linguist. Phil.* **12**(1), 15–37 (1989)
6. Socher, R., et al.: Recursive deep models for semantic compositionality over a sentiment treebank. In: Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, Seattle, Washington, USA, October 2013, pp. 1631–1642. Association for Computational Linguistics (2013)
7. Gupta, S., Bouadjenek, M.R., Robles-Kelly, A.: PERCY: a post-hoc explanation-based score for logic rule dissemination consistency assessment in sentiment classification. Technical report, Deakin University, School of Information Technology (2022)
8. Ribeiro, M.T., Singh, S., Guestrin, C.: “why should i trust you?”: explaining the predictions of any classifier. In: Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD 2016, New York, NY, USA, pp. 1135–1144. Association for Computing Machinery (2016)
9. Kim, Y.: Convolutional neural networks for sentence classification. In: Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), Doha, Qatar, October 2014, pp. 1746–1751. Association for Computational Linguistics (2014)
10. Zhang, Y., Wallace, B.: A sensitivity analysis of (and practitioners’ guide to) convolutional neural networks for sentence classification. In: Proceedings of the Eighth International Joint Conference on Natural Language Processing, Taipei, Taiwan, November 2017, vol. 1: Long Papers, pp. 253–263. Asian Federation of Natural Language Processing (2017)
11. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: Burges, C.J.C., Bottou, L., Welling, M., Ghahramani, Z., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems, vol. 26. Curran Associates Inc. (2013)
12. Pennington, J., Socher, R., Manning, C.D.: Glove: global vectors for word representation. In: Empirical Methods in Natural Language Processing (EMNLP), pp. 1532–1543 (2014)

13. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, vol. 1 (Long and Short Papers), Minneapolis, Minnesota, June 2019, pp. 4171–4186. Association for Computational Linguistics (2019)
14. Vaswani, A., et al.: Attention is all you need. In: Guyon, I., et al. (eds.) Advances in Neural Information Processing Systems, vol. 30. Curran Associates Inc. (2017)