
Local vs. Global interpretations for NLP

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently, WordsWorth scores have been proposed for calculating feature impor-
2 tance in the context of traditional deep learning models trained for text classification
3 tasks [Anonymous, 2021]. Here, we experiment with the idea behind these scores
4 and present them as a global explanation for a trained model. Interpretability
5 literature shows that delete one method acts as a good explanation for NLP tasks.
6 Since WW scores act as a good proxy to delete one scores for text classification,
7 we extend the argument and utilize them for interpretation. We provide local and
8 global explanations for a CNN trained on the IMDB reviews dataset by comparing
9 these scores with LIME. Similarly to LIME, the global representation is a bag of
10 words representation. Overall, we argue that evaluating a trained neural network on
11 single words, at all possible locations in the input text one by one, gives powerful
12 and valid insights into the workings of these otherwise black box models. This is a
13 work in progress and we are looking for further tests to evaluate the usefulness of
14 our method.

15 1 Introduction

16 Deep learning models are black boxes for the end user. Interpretability is becoming an important
17 concern, for users to trust the model overall as well as individual predictions. Global interpretations
18 try to provide a summary of the model in some form, whereas local interpretations are concerned
19 with explaining model prediction on a specific input. We use the idea of WordsWorth scores which
20 have been recently introduced, to show how to generate global explanations for a CNN, and compare
21 it to LIME. We further argue that evaluating a traditional deep learning model on all the words in the
22 training vocabulary gives a faithful picture of the model. Our method consists of evaluating a neural
23 network on all the words in the training vocabulary one by one, on all possible locations. For any
24 single evaluation, only one location in the input text has a valid word and all other locations are set to
25 zero. This gives us some estimate of how a neural network interprets each word.

26 2 Related Work

27 Interpretability is an exciting area of research in machine learning in general and deep learning in
28 particular. In this section we provide a brief and by no means thorough overview of some existing
29 techniques. This is meant to give the reader a brief idea of the different avenues that are being explored
30 for the problem of interpretability in NLP. An explanation can be global, such that it explains the
31 overall model, or local, such that it explains the decision made by the model on a particular instance.
32 Jacovi et al. [2018] calculate n-gram and word level scores. They calculate word-level score with
33 two methods that they name local and global, but both these methods use leave one out evaluations
34 and require a specific input x . Xiong et al. [2018] compare saliency maps and LRP(Layer-wise
35 Relevance Propagation), and conclude that LRP finds more relevant features. They argue that deleting
36 a particular word does not give a true picture of the contribution of this particular word to output

37 since it could be a part of an important n-gram, for example. Rather than single words, they highlight
38 high level features. Also this requires knowledge of the model parameters and architecture, so this is
39 not a black box explanation.

40 Chen and Ji [2020] place a mask over the word embedding layer during training, and thus propose
41 an interpretation friendly model architecture for text classification tasks. They compare LIME
42 and SampleShapley [Strumbelj and Kononenko, 2010] over simple neural networks and modified
43 networks and show that modified structure is more interpretable. LIME provides local explanations
44 by sampling from the distribution near the example and fitting a local model to it.

45 Lai et al. [2019] compare different methods and conclude that for deep learning model, important
46 features are often different from features for traditional models.

47 Chen et al. [2020] draw inspiration from Shapley values. They generate hierararchical explanations,
48 which might consist of multiple features that interact and drive the classifier to a particular prediction,
49 and argue that these perform better than other popular interpretation methods.

50 Arras et al. [2019] experiment on LSTMs. They compare contextual decomposition, Layer-wise
51 relevance propagation, gradient based methods and occlusion based which involve deleting a particular
52 word. They conclude that LRP performs best, but it is a white box method.

53 Nguyen [2018] examine sentiment analysis and topic classification tasks. They report that leave one
54 out approach and first derivative based generates explanations that correlate with human judgement.
55 For LIME to perform on a comparable level, a large number of samples are required typically around
56 5000. Further they argue that automatic evaluation matrices align reasonably well with human
57 judgement. They examine a logistic regression model and a simple feed forward neural network. We
58 argue that our technique of feature importance attribution aligns closely with omission technique.

59 For calculating word score matrices, Xu and Du [2020] propose a method which involves using test
60 examples.

61 **Global explanations in computer vision** Wu et al. [2020] outline a global explanations for CNNs.
62 First they find important features, such that occluding these features will flip the prediction on all
63 training instances for a specific class. Then they run different tests, which they call evaluating these
64 important features for semantic tasks.

65 **3 Single word evaluations as global interpretation**

66 **3.1 Direct evaluation on a single word**

67 In contrast to omission techniques, we make a case for directly calculating word influence. Leave
68 one out methods, which seem to be the inspiration for all the feature ranking methods used in NLP
69 [Li et al., 2016], have well grounded backing in statistics as well as intuition. In classical statistics,
70 evaluating the function on just one datapoint or feature does not make much sense.

71 However, the inspiration for single word evaluations comes from the fact that human generated text is
72 a highly compact form of data. Unlike images, where an individual pixel has no significance by itself
73 and the surrounding pixels are always needed to give it context, words carry immense information
74 within themselves. We might expect our deep learning classifier to be able to give meaningful insights
75 for each individual word most of the time.

76 **3.2 The case for position invariance for isolated words**

77 Similarly, we argue that the effect of a word on the prediction would not be highly dependent on its
78 exact position in the input sequence, since the inherent meaning of a word is quite independent of
79 its position in a sentence in which it appears, particularly when the surrounding words are ignored.
80 A well trained CNN, for example, should treat the word 'beautiful' roughly the same, whether it
81 appears near the beginning of the review or the end. If this is the case, we do not need to specifically
82 append words at a particular position to find out their importance.

Table 1: Words with highest WordsWorth scores

Highly positive words	Highly negative words
refreshing	awful
perfect	worst
excellent	waste
superb	poorly
unexpected	fails
perfectly	disappointing
rare	disappointment
enjoyable	forgettable
delightful	unfunny
blake	unwatchable

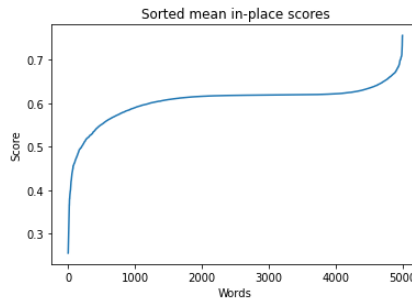


Figure 1: Mean values of in-place scores for IMDB Reviews with CNN

83 3.3 Experiments: IMDB Reviews with a CNN

84 We use the IMDB dataset [Maas et al., 2011], with 25000+25000 training+test examples of variable
 85 length. Each review in the training set has a positive/negative label attached to it. We use a simple
 86 CNN as the starting point of our experiments, with 32 dimensional embedding layer, 32 filters and 64
 87 hidden ReLU units and a single output unit. Test accuracy is 86.7%. Output probabilities above 0.5
 88 correspond to positive sentiment, and lower probabilities correspond to negative sentiment. Training
 89 vocabulary size is 5000 and each review is transformed to a 200 word piece of string by either
 90 appending zeros or removing the trailing part if it exceeds this length.

91 It has been shown in the original work that evaluating a trained model on a single word gives a
 92 score which is somewhat aligned with the sentiment behind that word. We go one step further and
 93 argue that a neural network trained for a sentiment analysis task on an input of size d and a training
 94 vocabulary of size v can be characterized by a $v * d$ matrix. Each row of this matrix corresponds to
 95 the scores for a particular word, and each column represents the score when this word is placed at
 96 that particular position in input. Recall that when a word is being evaluated at a particular position,
 97 all other inputs are set to zero.

98 **Mean and standard deviation of word score matrix** We compute this matrix and plot the mean
 99 and standard deviations for these scores along each row. These scores represent the trained network in
 100 a manner that is both interpretable and faithful to the representations the model has learned. Results
 101 in figures 1,2. The top words associated with each sentiment are presented in Table 2. Notice that the
 102 top 10 high standard deviation words have strong associations with the negative label.

103 For comparison, we present the top 10 words for each sentiment calculated through original
 104 WordsWorth scores in Table 1. Compare to mean scores and notice that the ordering for positive
 105 words is slightly different, and negative word list is exactly the same. This shows that the matrix
 106 we propose is a stable and reliable method of interpreting a neural network.

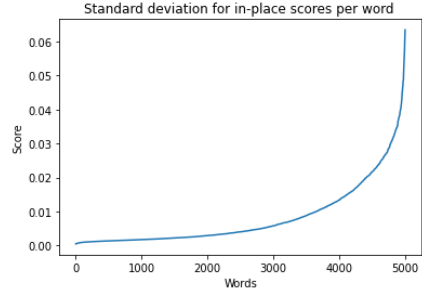


Figure 2: Standard deviation for in-place scores for IMDB Reviews with CNN

Table 2: Most important mean scores and corresponding words

Highly positive word	Scores	Highly negative words	Scores
perfect	0.7564884	awful	0.25465319
excellent	0.75603143	worst	0.2563593
refreshing	0.75512171	waste	0.26369784
superb	0.74893515	poorly	0.27800871
perfectly	0.73354795	fails	0.28250759
unexpected	0.73317402	disappointing	0.2965854
rare	0.72838456	disappointment	0.30099407
enjoyable	0.72051218	forgettable	0.30488086
delightful	0.71923713	unfunny	0.31003236
blake	0.71628927	unwatchable	0.32535568

107 **4 LIME vs WordsWorth: A case of two reviews**

108 Here we present two reviews and their local explanation given by LIME and global explanation given
 109 by WordsWorth scores.

110 **Review 1:** "Naturally in a film who's main themes are of mortality, nostalgia, and loss of innocence
 111 it is perhaps not surprising that it is rated more highly by older viewers than younger ones. However
 112 there is a craftsmanship and completeness to the film which anyone can enjoy. The pace is steady
 113 and constant, the characters full and engaging, the relationships and interactions natural showing that
 114 you do not need floods of tears to show emotion, screams to show fear, shouting to show dispute or
 115 violence to show anger. Naturally Joyce's short story lends the film a ready made structure as perfect
 116 as a polished diamond, but the small changes Huston makes such as the inclusion of the poem fit in
 117 neatly. It is truly a masterpiece of tact, subtlety and overwhelming beauty." Prediction: 0.99668723
 118 Results in figures 3,4,5 and 6. The input review is overwhelmingly positive and this is depicted by

Table 3: Most and least stable words, picked by standard deviation of in-place scores

Most stable words	Least stable words
hugh	awful
face	worst
zombie	waste
thomas	poorly
walls	fails
dickens	disappointing
carter	disappointment
soderbergh	forgettable
clown	unfunny
fairy	unwatchable

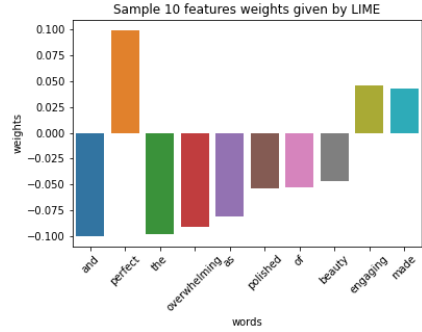


Figure 3: Top ten features highlighted by LIME in review 1

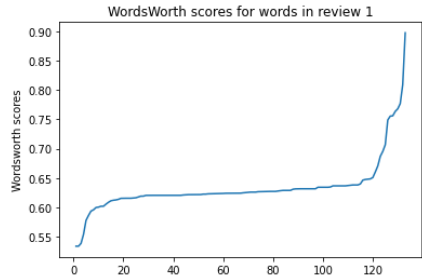


Figure 4: WordsWorth scores for all words in review 1

119 the WordsWorth score distribution. If a more faithful explanation is required, the scores from the
 120 word matrix we have introduced can be utilized for each word.

121 **Review 2:** "This movie is a disaster within a disaster film. It is full of great action scenes, which
 122 are only meaningful if you throw away all sense of reality. Let's see, word to the wise, lava burns
 123 you; steam burns you. You can't stand next to lava. Diverting a minor lava flow is difficult, let alone a
 124 significant one. Scares me to think that some might actually believe what they saw in this movie.

 125 />
Even worse is the significant amount of talent that went into making this film. I mean the
 126 acting is actually very good. The effects are above average. Hard to believe somebody read the scripts
 127 for this and allowed all this talent to be wasted. I guess my suggestion would be that if this movie is
 128 about to start on TV ... look away! It is like a train wreck: it is so awful that once you know what
 129 is coming, you just have to watch. Look away and spend your time on more meaningful content."
 130 Prediction: 0.46570438 Results in figures 7,8,9 and 10. Notice the high number of negative words in
 131 figure 8.

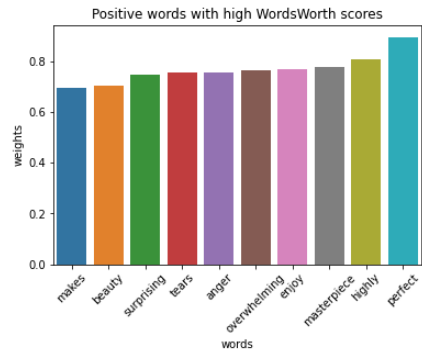


Figure 5: Most positive words for review 1

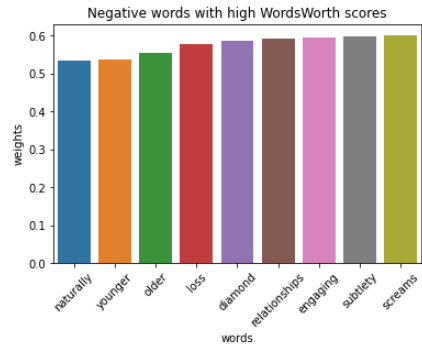


Figure 6: Most negative words for review 1

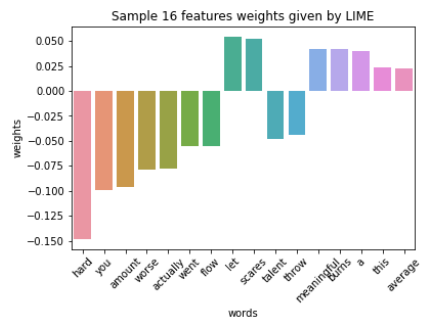


Figure 7: Top sixteen features highlighted by LIME in review 2

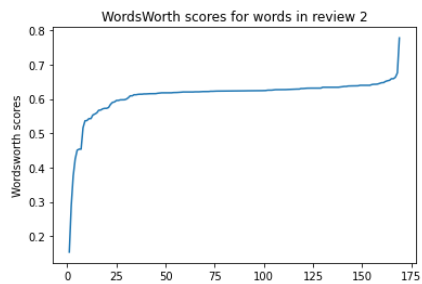


Figure 8: WordsWorth scores for all words in review 2

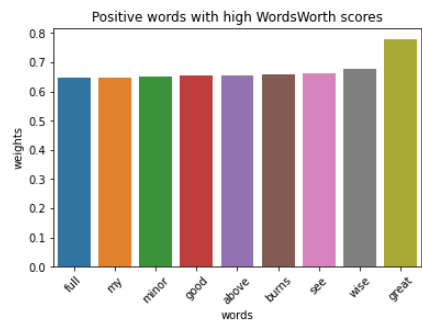


Figure 9: Most positive words for review 2

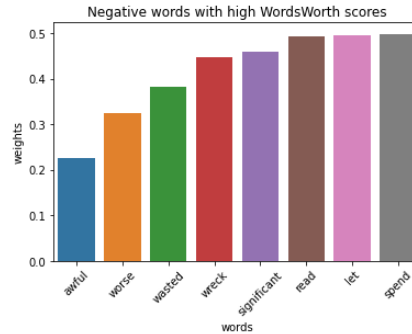


Figure 10: Most negative words for review 2

132 5 Conclusion

133 We provide a word matrix for evaluating a trained neural network on a sentiment analysis task.
 134 Further we show how to provide explanations for a specific input using WordsWorth scores. These
 135 explanations are not entirely local but they still provide a useful summary of the input. We compare
 136 our results to LIME. Leave-one-out scores have been shown to be effective at explaining classifier
 137 decisions and WordsWorth scores provide a good proxy to these local explanations. Further, since
 138 these scores have been shown to be effective at attacking LSTMS as well as for topic classification
 139 taskd, they might serve as a faithful explanation in scenarios which have not been explored in this
 140 paper.

141 References

- 142 Anonymous. Wordsworth scores for attacking {cnn}s and {lstm}s for text classification. In *Submitted to*
 143 *International Conference on Learning Representations*, 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=e6hMkY6MFcU)
 144 [id=e6hMkY6MFcU](https://openreview.net/forum?id=e6hMkY6MFcU). under review.
- 145 Leila Arras, Ahmed Osman, Klaus-Robert Müller, and Wojciech Samek. Evaluating recurrent neural network
 146 explanations. *CoRR*, abs/1904.11829, 2019. URL <http://arxiv.org/abs/1904.11829>.
- 147 Hanjie Chen and Yangfeng Ji. Learning variational word masks to improve the interpretability of neural text
 148 classifiers, 2020.
- 149 Hanjie Chen, Guangtao Zheng, and Yangfeng Ji. Generating hierarchical explanations on text classification
 150 via feature interaction detection. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault,
 151 editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL*
 152 *2020, Online, July 5-10, 2020*, pages 5578–5593. Association for Computational Linguistics, 2020. URL
 153 <https://www.aclweb.org/anthology/2020.acl-main.494/>.
- 154 Alon Jacovi, Oren Sar Shalom, and Yoav Goldberg. Understanding convolutional neural networks for text
 155 classification. *CoRR*, abs/1809.08037, 2018. URL <http://arxiv.org/abs/1809.08037>.
- 156 Vivian Lai, Zheng Cai, and Chenhao Tan. Many faces of feature importance: Comparing built-in and post-
 157 hoc feature importance in text classification. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan,
 158 editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and*
 159 *the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong*
 160 *Kong, China, November 3-7, 2019*, pages 486–495. Association for Computational Linguistics, 2019. doi:
 161 [10.18653/v1/D19-1046](https://doi.org/10.18653/v1/D19-1046). URL <https://doi.org/10.18653/v1/D19-1046>.
- 162 Jiwei Li, Will Monroe, and Dan Jurafsky. Understanding neural networks through representation erasure. *CoRR*,
 163 abs/1612.08220, 2016. URL <http://arxiv.org/abs/1612.08220>.
- 164 Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts.
 165 Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association*
 166 *for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA,
 167 June 2011. Association for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/P11-1015)
 168 [P11-1015](https://www.aclweb.org/anthology/P11-1015).
- 169 Dong Nguyen. Comparing automatic and human evaluation of local explanations for text classification. In
 170 Marilyn A. Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North*

- 171 *American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-*
172 *HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, pages 1069–1078.
173 Association for Computational Linguistics, 2018. doi: 10.18653/v1/n18-1097. URL [https://doi.org/10.](https://doi.org/10.18653/v1/n18-1097)
174 [18653/v1/n18-1097](https://doi.org/10.18653/v1/n18-1097).
- 175 Erik Strumbelj and Igor Kononenko. An efficient explanation of individual classifications using game theory. *J.*
176 *Mach. Learn. Res.*, 11:1–18, 2010. URL <https://dl.acm.org/citation.cfm?id=1756007>.
- 177 Weibin Wu, Yuxin Su, Xixian Chen, Shenglin Zhao, Irwin King, Michael R. Lyu, and Yu-Wing Tai. Towards
178 global explanations of convolutional neural networks with concept attribution. In *2020 IEEE/CVF Conference*
179 *on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 8649–
180 8658. IEEE, 2020. doi: 10.1109/CVPR42600.2020.00868. URL [https://doi.org/10.1109/CVPR42600.](https://doi.org/10.1109/CVPR42600.2020.00868)
181 [2020.00868](https://doi.org/10.1109/CVPR42600.2020.00868).
- 182 Wenting Xiong, Iftitahu Ni'mah, Juan M. G. Huesca, Werner van Ipenburg, Jan Veldsink, and Mykola Pech-
183 enizkiy. Looking deeper into deep learning model: Attribution-based explanations of textcnn. *CoRR*,
184 [abs/1811.03970](http://arxiv.org/abs/1811.03970), 2018. URL <http://arxiv.org/abs/1811.03970>.
- 185 Jincheng Xu and Qingfeng Du. On the interpretation of convolutional neural networks for text classification. In
186 Giuseppe De Giacomo, Alejandro Catalá, Bistra Dilkina, Michela Milano, Senén Barro, Alberto Bugarín,
187 and Jérôme Lang, editors, *ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8*
188 *September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference*
189 *on Prestigious Applications of Artificial Intelligence (PAIS 2020)*, volume 325 of *Frontiers in Artificial*
190 *Intelligence and Applications*, pages 2252–2259. IOS Press, 2020. doi: 10.3233/FAIA200352. URL
191 <https://doi.org/10.3233/FAIA200352>.