
On the Confidence of Neural Network Predictions for some NLP Tasks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Neural networks are known to produce unexpected results on inputs that are far
2 from the training distribution. One approach to tackle this problem is to detect the
3 samples on which the trained network can not answer reliably. ODIN is a recently
4 proposed method for out-of-distribution detection that does not modify the trained
5 network and achieves good performance for various image classification tasks. In
6 this paper we adapt ODIN for sentence classification and word tagging tasks. We
7 show that the scores produced by ODIN can be used as a confidence measure for
8 the predictions on both in-distribution and out-of-distribution datasets.

9 1 Introduction

10 Neural networks have been shown to perform well on various computer vision and natural language
11 processing tasks. The performance of neural networks is usually measured on the test sets of the
12 corresponding datasets, which does not show how the networks would perform on the samples from
13 other distributions.

14 Complex neural models that perform well on ImageNet dataset produce nonsensical labels for images
15 far from its training and test sets (see [Nguyen et al., 2015] for examples on unrecognizable images
16 and Figure 1 from [Shafaei et al., 2018] for more realistic images). Similar examples can be found
17 for neural models for NLP tasks. Table 1 shows the results of a simple neural POS tagger on different
18 sentences.

Table 1: The output of a neural POS tagger trained on English-LinES dataset on samples from three different datasets.

| | | | | | | | |
|--|-------------|----------|------|--------------|-----------|--------|--------|
| UD English-LinES (Accuracy = 14.3%, PbThreshold(s) = 0.99598, ODIN(s) = 0.05904) | | | | | | | |
| Sentence | Identifying | filters | that | are | currently | in | effect |
| Ground truth | VERB | NOUN | PRON | VERB | ADV | ADP | NOUN |
| Predicted labels | ADJ | NOUN | PRON | VERB | ADV | ADP | NOUN |
| Probabilities | 0.97 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 | 1.00 |
| UD English-EWT (Accuracy = 42.9%, PbThreshold(s) = 0.89237, ODIN(s) = 0.05899) | | | | | | | |
| Sentence | Try | googling | it | for | more | info | :) |
| Ground truth | VERB | VERB | PRON | ADP | ADJ | NOUN | SYM |
| Predicted labels | PRON | VERB | PRON | ADP | ADV | ADV | PUNCT |
| Probabilities | 0.58 | 1.00 | 1.00 | 0.99 | 1.00 | 0.70 | 0.98 |
| UD Dutch-Alpino (Accuracy = 57.1%, PbThreshold(s) = 0.99699, ODIN(s) = 0.05906) | | | | | | | |
| Sentence | Daarbij | is | een | Macedonische | militair | gedood | . |
| Ground truth | ADV | AUX | DET | ADJ | NOUN | VERB | PUNCT |
| Predicted labels | NOUN | AUX | DET | ADJ | ADJ | NOUN | PUNCT |
| Probabilities | 1.00 | 1.00 | 0.98 | 1.00 | 1.00 | 1.00 | 1.00 |

19 This limitation of neural networks makes it hard to deploy them in critical applications. One of the
20 possible directions to solve the problem is to measure the confidence of the prediction and do not
21 output the prediction if the confidence is low.

22 Hendrycks and Gimpel [2017] showed that for various neural models, correctly classified examples
23 tend to have higher maximum softmax probabilities than incorrectly classified examples and out-of-
24 distribution examples. They performed several experiments on deep convolutional networks trained
25 on CIFAR-10 and CIFAR-100 and show that by looking at the maximum predicted probability one
26 can classify images from CIFAR-10 and SUN test sets with 95% AUROC.

27 Liang et al. [2018] improved upon this baseline by using two relatively simple but efficient tricks:
28 (a) by adding temperature to the softmax calculation, (b) by applying adversarial-like perturbation
29 on the inputs. The method is called ODIN. Additionally, the authors showed that the performance
30 of ODIN strongly depends on the actual distance between the test distributions. For example, if the
31 in-distribution and out-of-distribution datasets are non-intersecting subsets of CIFAR-100, then ODIN
32 doesn't work well.

33 Our contributions are the following:

- 34 1. We adapt ODIN for two natural language processing tasks: sentiment analysis and part-of-
35 speech tagging. We show that the two tricks used in ODIN are helpful for sentiment analysis,
36 but we could not get improvement from input perturbations for POS tagging.
- 37 2. We choose the in-distribution (ID) and out-of-distribution (OOD) datasets in a way that the
38 labels used in the datasets are the same. This choice allows us to measure the accuracy of
39 the model on the test set of the OOD dataset. We show that the scores produced by ODIN are
40 relatively higher for the correctly classified examples of the OOD dataset, even when the
41 datasets are close and ODIN fails to properly separate them. For part-of-speech tagging, we
42 show that the scores produced by ODIN have higher rank correlation with the accuracy of
43 the neural network predictions than the simpler baseline on both ID and OOD datasets.
- 44 3. We demonstrate that although character-level embeddings improve POS tagging accuracy,
45 they make it harder to distinguish ID and OOD test sets for the two OOD detection methods.
46 Additionally, the scores produced by these methods have lower rank correlation with the
47 prediction accuracy for all datasets we have tried compared to the models without character-
48 level embeddings.

49 2 Related Work

50 There are various ways to approach out-of-distribution detection problem for high dimensional
51 inputs. Shafaei et al. [2018] made a detailed comparison of different approaches on basic image
52 classification tasks. Hendrycks and Gimpel [2017] set a baseline for these methods by looking at the
53 maximum softmax probabilities ($PbThreshold^1$). Liang et al. [2018] improved upon this baseline
54 using temperature rescaling and perturbations on inputs (ODIN). Gal and Ghahramani [2016] showed
55 that Bayesian interpretation of dropout allows to use it to measure the uncertainty of neural network's
56 prediction (MC-Dropout). Lakshminarayanan et al. [2017] showed that an ensemble of multiple deep
57 networks can be used to capture uncertainty in a non-Bayesian way. It is important to note that these
58 methods require access to the datasets and to the neural networks trained on them. $PbThreshold$
59 and ODIN can be applied to any pretrained neural network. MC-Dropout is applicable to any network
60 with a dropout layer (notably, ResNets do not use dropout), and Lakshminarayanan et al. [2017]
61 requires an ensemble of neural networks.

62 Another class of approaches for OOD detection are based on generative models. Shafaei et al. [2018]
63 evaluated methods based on autoencoders ($AEThreshold$) and $PixelCNN++$ [Salimans et al., 2017],
64 and showed that for basic image classification tasks ODIN outperforms them (in terms of OOD
65 detection accuracy). Recently, Choi and Jang [2018] used an ensemble of generative models to
66 beat ODIN on OOD detection for MNIST and CIFAR-10 datasets (in terms of AUROC). Note that
67 these generative approaches are model-independent, they do not depend on the neural network that
68 performs the actual classification.

¹We use the codenames of the algorithms from [Shafaei et al., 2018]

Table 2: The datasets used in our experiments along with the performance of the state-of-the-art models on these datasets.

| | Labels | Number of sentences | | | Accuracy of SOTA models |
|--------------|--------|---------------------|------|-------|--------------------------------|
| | | Train | Dev. | Test | |
| Yelp Reviews | 5 | 650000 | 0 | 50000 | 70.02 [Howard and Ruder, 2018] |
| SST-5 | 5 | 8544 | 1101 | 2210 | 54.70 [Peters et al., 2018] |
| en-EWT | 17 | 12543 | 2002 | 2077 | 95.94 [Lim et al., 2018] |
| en-LinES | 17 | 2738 | 912 | 914 | 97.06 [Lim et al., 2018] |
| en-GUM | 17 | 2914 | 707 | 769 | 96.44 [Lim et al., 2018] |
| nl-Alpino | 17 | 12269 | 718 | 596 | 96.90 [Straka, 2018] |

69 Most of the research on OOD detection is focused on image classification tasks. Notable exceptions
70 are [Hendrycks and Gimpel, 2017] and [Shalev et al., 2018]. Applications of more advanced methods
71 to NLP tasks remain largely unexplored.

72 Additionally, all experiments described above are performed on ID and OOD datasets with different
73 sets of labels. These experiments are good enough to measure the performance of OOD detection. On
74 the other hand, when the distributions of the two datasets have a significant overlap, OOD detection
75 methods fail to produce high accuracy scores. This phenomenon is demonstrated in Section 4.5 of
76 [Liang et al., 2018]. In the context of measuring the confidence of neural network predictions it is
77 fine to misclassify OOD samples as ID unless the network does not produce incorrect answers for the
78 misclassified examples. In order to measure the accuracy of neural models on misclassified OOD
79 samples we choose the ID and OOD datasets to have the same set of labels.

80 3 Experiments

81 3.1 Datasets

82 We performed experiments for two tasks: sentiment analysis and part-of-speech tagging. For
83 sentiment analysis we used the five-class Yelp reviews dataset from [Zhang et al., 2015] and Stanford
84 Sentiment Treebank (SST) [Socher et al., 2013]. SST has labels for every sentence and for every
85 phrase produced from the dependency trees of the sentences. The labels are real numbers between
86 0 and 1. We created a five-class version of the labels by splitting [0,1] into five equal intervals.
87 For part-of-speech tagging we used two English and one more Dutch treebanks from Universal
88 Dependencies v2.2 [Zeman et al., 2018a] used in the CoNLL Shared Task 2018 [Zeman et al., 2018b].
89 The majority of the sentences in English-LinES treebank are from literature. English-EWT dataset is
90 larger and is more diverse. The datasets are described in Table 2².

91 3.2 Neural models

92 In contrast to Liang et al. [2018], we did not use state-of-the-art neural networks. Instead, we trained
93 simple recurrent models for both tasks. For simplicity, we did not use pretrained word embeddings.
94 Similar to [Joulin et al., 2017], we used *hashing trick* to map the words to hashes and embed the
95 hashes using D -dimensional vectors. We randomly initialized the embedding matrix and made it
96 trainable.

97 Let s be a sentence from one of the datasets, and w_1, \dots, w_M be the words. The embedding of
98 the m -th word of the sentence will be $x_m = \text{Hash}(w_m)$. We apply bidirectional LSTM on the
99 embeddings: $\vec{h}_m, \overleftarrow{h}_m = \text{BiLSTM}(x_m)$. For sentiment analysis we apply a dense layer on the
100 concatenation of the last states of the two LSTMs: $f_{sc}(s) = W[\vec{h}_M, \overleftarrow{h}_1] + b$. The loss function
101 is a cross-entropy: $\text{loss}(s) = \text{ce}(S(f_{sc}(s)), 1)$, where $S_i(\mathbf{z}, T) = \frac{\exp(z_i/T)}{\sum_{j=1}^C \exp(z_j/T)}$ is the modified
102 softmax function, T is the temperature scaling parameter, and C is the number of classes.

²State-of-the-art results for sentiment analysis datasets are retrieved from <http://nlpprogress.com>
on 27.10.2018. UD POS tagging results are from [http://universaldependencies.org/conll118/
results-upos.html](http://universaldependencies.org/conll118/results-upos.html).

Table 3: Out-of-distribution detection performance on part-of-speech tagging tasks. Performance is reported in AUROC scores (higher is better). Char. column indicates whether character-level embeddings were used in the model. ϵ , T column lists the best hyperparameters for ODIN found using grid search on validation sets. All scores are in percents.

| ID | OOD | Char. | Accuracy | | PbThreshold | | | ODIN | | | |
|----------|-----------|-------|----------|------|--------------------|------|------|--------------------|-------------|--------------|------|
| | | | ID | OOD | AUROC | | | ϵ , T | AUROC | | |
| Yelp | SST-5 | No | 59.6 | 24.5 | 79.9 | | | 0.011, 2K | | 90.68 | |
| | | | | | Mean / Med. / Tok. | | | Mean / Med. / Tok. | | | |
| en-LinES | en-EWT | Yes | 89.6 | 75.9 | 66.8 | 70.7 | 57.1 | 0, 5K | 75.1 | 71.3 | 58.2 |
| en-LinES | en-EWT | No | 87.8 | 73.4 | 72.4 | 73.9 | 57.9 | 0, 10K | 77.0 | 74.0 | 59.1 |
| en-LinES | nl-Alpino | Yes | 89.6 | 31.7 | 95.2 | 97.4 | 76.5 | 0, 2K | 99.1 | 98.3 | 79.8 |
| en-LinES | nl-Alpino | No | 87.8 | 28.6 | 97.5 | 97.4 | 77.4 | 0, 5 | 98.2 | 97.9 | 81.7 |

103 For POS tagging we apply a dense layer on every hidden state: $f_{st}(w_m) = W([\vec{h}_m, \overleftarrow{h}_m]) + b$. The
 104 loss function is the average of word-level cross entropies $loss(s) = \frac{1}{M} \sum ce(S(f_{st}(w_m), 1))$.

105 3.3 Out-of-distribution detection method

106 We use ODIN to detect out-of-distribution samples and compare it with the PbThreshold baseline. For
 107 every sentence s we compute the scores for each of the methods: $PbThreshold(s) = \max S(f_{sc}(s))$
 108 and $ODIN(s) = \max S(f_{sc}(\tilde{\mathbf{x}}), T)$, where $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} S_{\hat{y}}(\mathbf{x}))$, where $\hat{y} = \text{argmax} S(\mathbf{x}, 1)$.
 109 Here ϵ (perturbation magnitude) and T (temperature) are hyperparameters, which are chosen based
 110 on the OOD detection performance on the development sets. For POS tagging, the gradient in the
 111 ODIN score formula is applied to the mean of word-level probability maximums. The final ODIN
 112 score of the sentence is some aggregate of the word-level ODIN scores. We tried mean and median
 113 as the aggregate function. Additionally, we tried to do out-of-distribution detection at the level of
 114 tokens, which is expected to be more challenging, as the vocabularies of the datasets have a significant
 115 overlap.

116 Hendrycks and Gimpel [2017] and Liang et al. [2018] report multiple scores for OOD detection
 117 performance, while Shafaei et al. [2018] reports only accuracy. We follow [Choi and Jang, 2018]
 118 and choose AUROC, as it does not require to tune a threshold. We report the scores on the test sets.
 119 Appendix H of [Liang et al., 2018] demonstrates that the choice of the OOD distribution is not critical
 120 for the hyperparameter tuning.

121 3.4 Evaluating confidence scores

122 For every sentence we produce the scores $PbThreshold(s)$ and $ODIN(s)$ and attempt to interpret
 123 them as confidence scores for the prediction of the neural network. To measure how well these scores
 124 can perform as a confidence measure, we calculate Spearman’s rank correlation coefficient between
 125 the scores and the accuracy numbers.

126 4 Results and Discussions

127 Table 3 shows the results for OOD detection and Table 4 shows the rank correlation coefficients for
 128 PbThreshold and ODIN methods.

129 **The role of the temperature scaling and input perturbations** All our experiments confirm the
 130 observation from [Liang et al., 2018] that temperature scaling improves out-of-distribution detection.
 131 The effect of higher temperatures saturates when T reaches thousands (Figure 1). The positive effect
 132 of the perturbations on the inputs is visible for sentiment analysis, but not for POS tagging. We
 133 have noticed that when we try to perform OOD detection at the level of tokens, small perturbations
 134 ($\epsilon = 0.005$) bring tiny improvements to the performance (by less than 1 AUROC percent point). But
 135 for sentence level OOD detection we get worse AUROC scores for $\epsilon > 0$ (or even when $\epsilon < 0$) for
 136 both mean and median averaging.

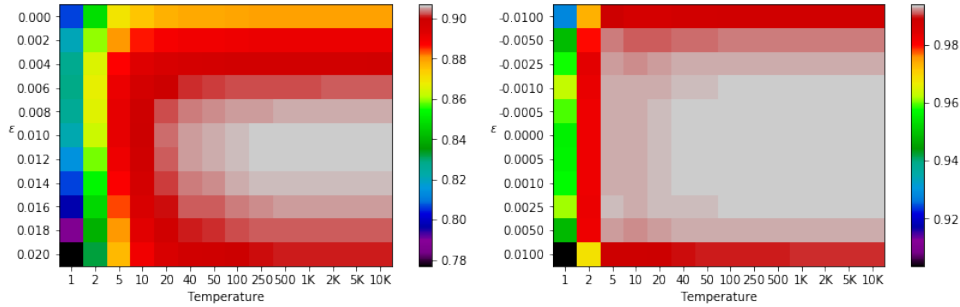


Figure 1: AUROC scores for ODIN out-of-distribution detection method for different values of hyperparameters. The left image is for sentiment analysis (Yelp dataset as ID, SST-5 as OOD), the right one is for POS tagging (English-LinES as ID, Dutch-Alpino as OOD, the neural model uses character-level embeddings).

137 **Sentence-level scores** In POS tagging experiments, median of the token-level scores is usually
 138 a better score for sentence-level OOD detection than the mean when $T = 1$, while for higher
 139 temperatures mean is consistently better. As expected, token-level OOD detection doesn't work well
 140 when both ID and OOD datasets are in English (the vocabularies overlap significantly).

141 **Ranking of the sentences** ODIN is clearly better than PbThreshold according to Spearman's rank
 142 correlation coefficient for POS tagging tasks (Table 4). For a neural network trained on en-LinES,
 143 ODIN scores are a good indicator how the network will perform on OOD samples. It is a much
 144 better indicator for the closer English dataset than the further Dutch dataset. On the other hand,
 145 when we consider the union of ID and OOD datasets, it works better for the union of en-LinES and
 146 nl-Alpino datasets.

147 To visualize what these correlation coefficients imply in practice, we split the samples into 20 equal
 148 buckets according to the scores, and compute accuracy on each of the buckets. We expect to see that
 149 the accuracy numbers are monotonically increasing, and that the accuracy numbers on different test
 150 sets are close to each other for each bucket. Figure 2 shows that ODIN performs slightly better than
 151 PbThreshold.

152 For sentiment analysis, we could not get improvement in rank correlation with ODIN, although
 153 the performance of OOD detection is improved. The reasons of this phenomenon are yet to be
 154 investigated.

155 **The role of character-level embeddings** Character-level embeddings improve the accuracy of the
 156 neural POS tagger for both ID and OOD datasets (consistent with the results reported by Reimers
 157 and Gurevych [2017]). On the other hand, they make it harder for PbThreshold and ODIN methods
 158 to separate ID and OOD datasets. Additionally, the scores from the models with character-level
 159 embeddings have lower rank correlation with the prediction accuracy. This implies, that the usage of
 160 character-level embeddings can be a tradeoff between the accuracy of the model and the reliability of
 161 the confidence scores.

Table 4: Spearman's ranking correlation coefficients between sample-level accuracies and the scores produced by PbThreshold and ODIN (higher is better). For each method we report the coefficient for ID, OOD test sets and the union of both.

| ID | OOD | Char. | PbThreshold | | | ODIN | | |
|----------|-----------|-------|-------------|-------|-------|-------|-------|-------|
| | | | ID | OOD | Both | ID | OOD | Both |
| en-LinES | nl-EWT | Yes | 0.365 | 0.571 | 0.565 | 0.441 | 0.641 | 0.636 |
| en-LinES | nl-EWT | No | 0.408 | 0.632 | 0.616 | 0.476 | 0.684 | 0.676 |
| en-LinES | nl-Alpino | Yes | 0.365 | 0.265 | 0.793 | 0.441 | 0.308 | 0.828 |
| en-LinES | nl-Alpino | No | 0.408 | 0.312 | 0.803 | 0.476 | 0.445 | 0.833 |

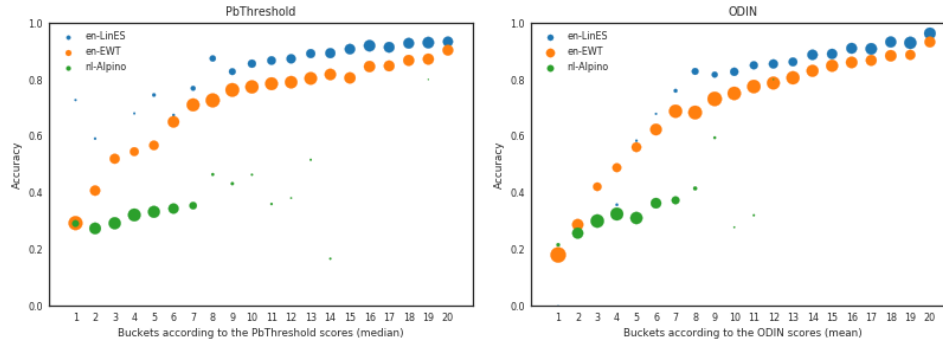


Figure 2: Accuracy of the POS tagger trained on en-LinES (without character-level embeddings) on 20 equal buckets of the union of three test sets. The buckets are computed according to the scores produced by each method (PbThreshold and ODIN). ϵ and T for ODIN are determined based on the development sets of en-LinES (ID) and en-EWT (OOD). The size of a circle is proportional to the number of samples that fall into that bucket. Ideally, accuracy scores for the i -th bucket should be higher than for the $(i - 1)$ -th bucket, and y coordinates of the three circles for each bucket should be the same.

162 5 Conclusions and Future Work

163 In this work we have adapted ODIN out-of-distribution detection method on sentence classification
 164 and sequence tagging tasks. We showed that as an OOD detector it performs consistently better than
 165 for the PbThreshold baseline. Additionally, we attempted to quantify how well the scores produced
 166 by these methods can be used as confidence scores for the predictions of neural models.

167 There are many other OOD detection methods that have yet to be tested on NLP tasks. On the other
 168 hand, our analysis notably doesn't cover sequence-to-sequence tasks. We have shown that the usage
 169 of character-level embeddings makes OOD detection harder for both PbThreshold and ODIN. The
 170 role of the pretrained word vectors, the size of the embeddings, the choice of the neural architecture
 171 (recurrent, convolutional or Transformer-like) on OOD detection performance is left for future work.

172 References

- 173 Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence
 174 predictions for unrecognizable images. In *Proceedings of the IEEE Conference on Computer
 175 Vision and Pattern Recognition*, pages 427–436, 2015.
- 176 Alireza Shafaei, Mark Schmidt, and James J. Little. Does Your Model Know the Digit 6 Is Not a
 177 Cat? A Less Biased Evaluation of Outlier Detectors. *ArXiv e-prints*, 2018.
- 178 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
 179 examples in neural networks. 2017. URL <https://openreview.net/forum?id=Hkg4TI9x1>.
- 180 Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image
 181 detection in neural networks. In *International Conference on Learning Representations*, 2018.
 182 URL <https://openreview.net/forum?id=H1VGkIxRZ>.
- 183 Yarín Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model
 184 uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059,
 185 2016.
- 186 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive
 187 uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing
 188 Systems*, pages 6402–6413, 2017.
- 189 Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. Pixelcnn++: Improving the
 190 pixelcnn with discretized logistic mixture likelihood and other modifications. In *International*

- 191 *Conference on Learning Representations*, 2017. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=BJrFC6ceg)
192 [BJrFC6ceg](https://openreview.net/forum?id=BJrFC6ceg).
- 193 Hyunsun Choi and Eric Jang. Generative ensembles for robust anomaly detection. *arXiv preprint*
194 *arXiv:1810.01392*, 2018.
- 195 Gabi Shalev, Yossi Adi, and Joseph Keshet. Out-of-distribution detection using multiple semantic
196 label representations. *arXiv preprint arXiv:1808.06664*, 2018.
- 197 Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text
198 classification. In *Advances in neural information processing systems*, pages 649–657, 2015.
- 199 Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Ng, and
200 Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank.
201 In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages
202 1631–1642, 2013.
- 203 Dan Zeman et al. Universal Dependencies 2.2 – CoNLL 2018 shared task development and test data,
204 2018a. URL <http://hdl.handle.net/11234/1-2184>. LINDAT/CLARIN digital library at
205 the Institute of Formal and Applied Linguistics, Charles University, Prague, <http://hdl.handle.net/11234/1-2184>.
206
- 207 Daniel Zeman, Jan Hajič, Martin Popel, Martin Potthast, Milan Straka, Filip Ginter, Joakim Nivre,
208 and Slav Petrov. CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal
209 Dependencies. In *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw*
210 *Text to Universal Dependencies*, pages 1–20, Brussels, Belgium, October 2018b. Association for
211 Computational Linguistics. ISBN 978-1-948087-82-7.
- 212 Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification.
213 In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*
214 *(Volume 1: Long Papers)*, pages 328–339. Association for Computational Linguistics, 2018. URL
215 <http://aclweb.org/anthology/P18-1031>.
- 216 Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and
217 Luke Zettlemoyer. Deep contextualized word representations. In *Proc. of NAACL*, 2018.
- 218 KyungTae Lim, Cheoneum Park, Changki Lee, and Thierry Poibeau. SEx BiST: A multi-source
219 trainable parser with deep contextualized lexical representations. In *Proceedings of the CoNLL*
220 *2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 143–
221 152, Brussels, Belgium, October 2018. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/K18-2014>.
222
- 223 Milan Straka. UDPipe 2.0 prototype at CoNLL 2018 UD shared task. In *Proceedings of the*
224 *CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages
225 197–207, Brussels, Belgium, October 2018. Association for Computational Linguistics. URL
226 <http://www.aclweb.org/anthology/K18-2020>.
- 227 Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text
228 classification. In *Proceedings of the 15th Conference of the European Chapter of the Association*
229 *for Computational Linguistics: Volume 2, Short Papers*, volume 2, pages 427–431, 2017.
- 230 Nils Reimers and Iryna Gurevych. Reporting score distributions makes a difference: Performance
231 study of lstm-networks for sequence tagging. In *Proceedings of the 2017 Conference on Empirical*
232 *Methods in Natural Language Processing*, pages 338–348. Association for Computational Linguistics,
233 2017. doi: 10.18653/v1/D17-1035. URL <http://aclweb.org/anthology/D17-1035>.