
DATE: Detecting Anomalies in Text via Self-Supervision of Transformers

Andrei Manolache^{1 2} Florin Brad¹ Elena Burceanu^{1 2 3}

Abstract

Leveraging deep learning models for Anomaly Detection (AD) has seen widespread use in recent years due to superior performances over traditional methods. Recent deep methods for anomalies in images learn better features of normality in an end-to-end self-supervised setting. These methods train a model to discriminate between different transformations applied to visual data and then use the output to compute an anomaly score. We use this approach for AD in text, by introducing a novel pretext task on text sequences. We learn our DATE model end-to-end, enforcing two independent and complementary self-supervision signals, one at the token-level and one at the sequence-level. Under this new task formulation, we show strong quantitative and qualitative results on the 20NewsGroups and AG News datasets. In the *semi-supervised* setting, we outperform state-of-the-art results by +13.5% and +6.9%, respectively (AUROC). In the *unsupervised* configuration, DATE surpasses all other methods even when 10% of its training data is contaminated with outliers (compared with 0% for the others).

1. Introduction

Anomaly Detection (AD) can be intuitively defined as the task of identifying examples that deviate from the other ones to a degree that arouses suspicion (Hawkins, 1980). Research into AD spans several decades (Chandola et al., 2009; Aggarwal, 2015) and has proved fruitful in several real-world problems, such as intrusion detection systems (Banoth et al., 2017), credit card fraud detection (Dorransoro et al., 1997), and manufacturing (Kammerer et al., 2019).

Our DATE method is applicable in the *semi-supervised* AD setting, in which we only train on clean, labeled normal ex-

amples, as well as the *unsupervised* AD setting, where both unlabeled normal and abnormal data are used for training. Typical deep learning approaches in AD involve learning features of normality using autoencoders (Hawkins et al., 2002; Sakurada & Yairi, 2014; Chen et al., 2017) or generative adversarial networks (Schlegl et al., 2017). Under this setup, anomalous examples lead to a higher reconstruction error or differ compared with generated samples.

Recent deep AD methods for images learn more effective features of visual normality through *self-supervision*, by training a deep neural network to discriminate between different transformations applied to the input images (Golan & El-Yaniv, 2018; Wang et al., 2019). An anomaly score is then computed by aggregating model predictions over several transformed input samples.

We adapt those self-supervised classification methods for AD from vision to learn anomaly scores indicative of text normality. ELECTRA (Clark et al., 2020) proposes an efficient language representation learner, which solves the *Replaced Token Detection* (RTD) task. Here the input tokens are plausibly corrupted with a BERT-based (Devlin et al., 2018) generator, and then a discriminator predicts for each token if it is real or replaced by the generator. In a similar manner, we introduce a complementary sequence-level pretext task called *Replaced Mask Detection* (RMD), where we enforce the discriminator to predict the predefined *mask pattern* used when choosing what tokens to replace. For instance, given the input text ‘They were ready to go’ and the mask pattern $[0, 0, 1, 0, 1]$, the corrupted text could be ‘They were prepared to advance’. The RMD multi-class classification task asks which *mask pattern* (out of K such patterns) was used to corrupt the original text, based on the corrupted text. Our generator-discriminator model solves both the RMD and the RTD task and then computes the anomaly scores based on the output probabilities, as visually explained in detail Fig. 1-2.

We notably simplify the computation of the Pseudo Label (PL) anomaly score (Wang et al., 2019) by removing the dependency on running over multiple transformations and enabling it to work with token-level predictions. This significantly speeds up the PL score evaluation.

To our knowledge, DATE is the first end-to-end deep AD method on text that uses self-supervised classification mod-

¹Bitdefender ²University of Bucharest, Romania ³Institute of Mathematics of the Romanian Academy. Correspondence to: Andrei Manolache <amanolache@bitdefender.com>.

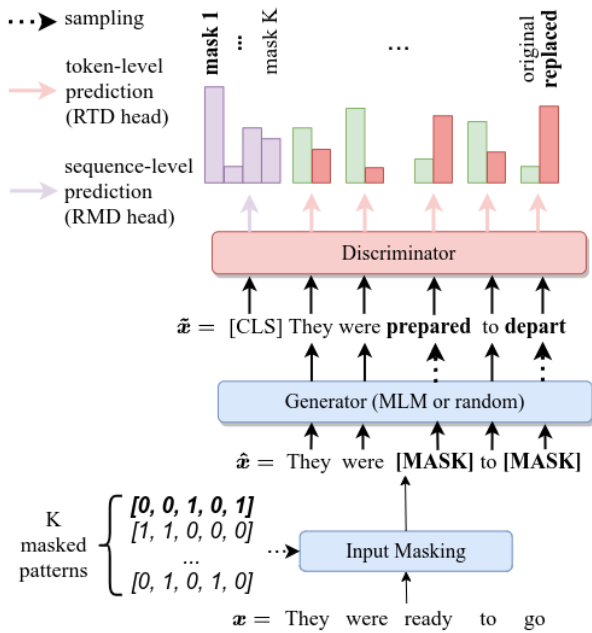


Figure 1. DATE Training. Firstly, the input sequence is masked using a sampled masked pattern and a generator fills in new tokens in place of the masked ones. Secondly, the discriminator receives supervision signals from two tasks: RMD (which mask pattern was applied to the input sequence) and RTD (the per-token status: original or replaced).

els to produce normality scores. Our **contributions** are:

First, we introduce a sequence-level self-supervised task called *Replaced Mask Detection* to distinguish between different transformations applied to a text. Jointly optimizing both sequence and token-level tasks stabilizes training, improving the AD performance.

Second, we compute an efficient Pseudo Label score for anomalies, by removing the need for evaluating multiple transformations, allowing it to work directly on individual tokens probabilities. This makes our model faster and its results more interpretable.

Third, we outperform existing state-of-the-art semi-supervised AD methods on text by a large margin (AUROC) on two datasets: 20Newsgroups (+13.5%) and AG News (+6.9%). Moreover, in unsupervised AD settings, even with 10% outliers in training data, DATE surpasses all other methods trained with 0% outliers.

2. Our Approach

Our method is called **DATE** for 'Detecting Anomalies in Text using ELECTRA'. We propose an end-to-end AD approach for the discrete text domain that combines our novel self-supervised task (Replaced Mask Detection), a power-

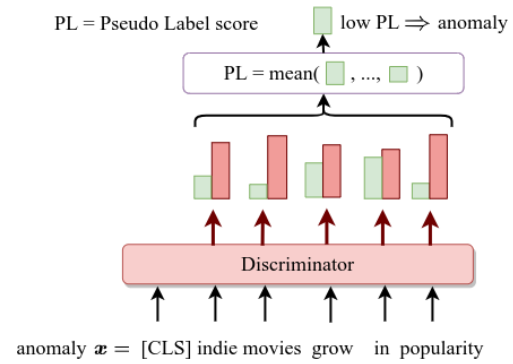


Figure 2. DATE Testing. The input text sequence is fed to the discriminator, resulting in token-level probabilities for the normal class, which are further aggregated into an anomaly score, as detailed in Appendix A. For deciding whether a sample is either normal or abnormal, we aggregate over all of its tokens.

ful representation learner for text (ELECTRA), and an AD score tailored for sequential data. We present in Fig. 1-2 a visual representation for the training and testing pipeline, explained in detail in Appendix A.

Replaced Mask Detection task. We introduce a novel self-supervised task for text, called Replaced Mask Detection (RMD). This discriminative task creates training data by transforming an existing text using one out of K given operations. It further asks to predict the correct operation, given the transformed text. The transformation over the text consists of two steps: 1) *masking* some of the input words using a predefined mask pattern and 2) replacing the masked words with alternative ones (e.g. 'car' with 'taxi').

Connecting RMD and RTD tasks. RTD (Replaced Token Detection) is a binary sequence tagging task, where some tokens in the input are corrupted with plausible alternatives, similarly to RMD. The discriminator must then predict for each token if it's the *original* token or a *replaced* one. Distinctly from RTD, which is a *token-level* discriminative task, RMD is a *sequence-level* one, where the model distinguishes between a fixed number of predefined transformations applied to the input. As such, RMD can be seen as the text counterpart task for the self-supervised classification of geometric alterations applied to images (Golan & El-Yaniv, 2018; Wang et al., 2019). While RTD predictions could be used to sequentially predict an entire mask pattern, they can lead to masks that are not part of the predefined K patterns. But the RMD constraint overcomes this behaviour. We thus train DATE to solve both tasks simultaneously, which increases the AD performance compared to solving one task only, as shown in Sec. C. Furthermore, this approach also improves training stability.

We solve RMD and RTD by jointly training a generator, G , and a discriminator, D . G is an MLM used to *replace* the masked tokens with plausible alternatives. We also consider a setup with a *random generator*, in which we sample tokens uniformly from the vocabulary. D is a deep neural network with two prediction heads used to distinguish between *corrupted* and original tokens (RTD) and to predict which mask pattern was applied to the corrupted input (RMD). At test time, G is discarded and D 's probabilities are used to compute an anomaly score. Both G and D models are based on a BERT encoder, which consists of several stacked Transformer blocks (Vaswani et al., 2017).

Anomaly Detection score. We adapt the Pseudo Label (PL) based score from the E^3 *Outlier* framework (Wang et al., 2019) in a novel and efficient way. In its general form, the PL score aggregates responses corresponding to multiple transformations of x . This approach requires k input transformations over an input x and k forward passes through a discriminator. It then takes the probability of the ground truth transformation and averages it over all k transformations.

To compute PL for our RMD task, we take x to be our input text and the K mask patterns as the possible transformations. We corrupt x with mask $m^{(i)}$ and feed the resulted text to the discriminator. We take the probability of the i -th mask from the RMD head. We repeat this process k times and average over the probabilities of the correct mask pattern. This formulation requires k feed-forward steps through the DATE network, which slows down inference. We propose a more computationally efficient approach next.

PL over RTD classification scores. Instead of aggregating *sequence-level* responses from multiple transformations over the input, we can aggregate *token-level* responses from a single model over the input to compute an anomaly score. More specifically, we can discard the generator and feed the original input text to the discriminator directly. We then use the probability of each token being *original* (not *corrupted*) and then average over all the tokens in the sequence.

3. Experimental analysis

We detail next the empirical validation of our method by presenting the semi-supervised and unsupervised experimental setup and comparison with state-of-the-art methods. DATE does not use any form of pre-training or knowledge transfer (from other datasets or tasks), learning all the embeddings from scratch. We observe that when fine-tuning pre-trained models with our task formulation, the OD performance does not improve during training. We posit that using pre-trained models introduces unwanted prior knowledge about the outliers, making our model considering them known (normal).

Anomaly Detection setup. We use a semi-supervised setting in Sec. C-3.1 and an unsupervised one in Sec. 3.2. In the semi-supervised case, we successively treat one class as normal (*inliers*) and all the other classes as abnormal (*outliers*). In the unsupervised AD setting, we add a fraction of outliers to the inliers training set, thus contaminating it. We compute the Area Under the Receiver Operating Curve (AUROC) for comparing our method with the previous state-of-the-art.

Datasets. We test our solution using two text classification datasets, after stripping headers and other metadata. For **20Newsgroups**¹ dataset we keep the exact setup, splits, and preprocessing (lowercase, removal of: punctuation, number, stop word and short words) as in (Ruff et al., 2019). The second dataset, **AG News**² (Zhang et al., 2015) is significantly larger, better suited for deep learning methods. Our code, including preprocessing steps, is publicly available³.

3.1. Comparison with other AD methods

We compare our method against classical AD baselines like Isolation Forest (Liu et al., 2008) and existing state-of-the-art OneClassSVMs (Schölkopf et al., 2001b) and CVDD (Ruff et al., 2019). We outperform all previously reported performances on all *20Newsgroups* splits by a large margin: 13.5% over the best reported CVDD and 11.7% over the best OCSVM, as shown in Tab. 1. In contrast, DATE uses the same set of hyper-parameters per dataset.

3.2. Unsupervised AD

We further analyse how our algorithm works in a fully unsupervised scenario, namely when the training set contains some anomalous samples (which we treat as normal ones). By definition, the quantity of anomalous events in the training set is significantly lower than the normal ones. In this experiment, we show how our algorithm performance is influenced by the percentage of anomalies in training data. Our method proves to be extremely robust, surpassing state-of-the-art, which is a semi-supervised solution, trained over a clean dataset (with 0% anomalies), even at 10% contamination, with +0.9% in AUROC (see Fig. 3). The reported scores are the mean over all AG News splits.

3.3. Qualitative results

We show in Fig. 4 how DATE performs in identifying anomalies. Each token is colored based on its PL score.

Separating anomalies. We see in Fig. 5 how our anomaly score (PL) is distributed among normal vs abnormal samples.

¹<http://qwone.com/~jason/20Newsgroups/>

²http://groups.di.unipi.it/~gulli/AG_corpus_of_news_articles.html

³<https://github.com/bit-ml/date>

	Inlier class	IsoForest best	OCSVM best	CVDD best	DATE (Ours)
20 News	comp	66.1	78.0	74.0	92.1
	rec	59.4	70.0	60.6	83.4
	sci	57.8	64.2	58.2	69.7
	misc	62.4	62.1	75.7	86.0
	pol	65.3	76.1	71.5	81.9
	rel	71.4	78.9	78.1	86.1
AG News	business	79.6	79.9	84.0 [‡]	90.0
	sci	76.9	80.7	79.0 [‡]	84.0
	sports	84.7	92.4	89.9 [‡]	95.9
	world	73.2	83.2	79.6 [‡]	90.1

Table 1. Semi-supervised performance (AUROC%). We test on the 20Newsgroups and AG News datasets, by comparing DATE against several strong baselines and state-of-the-art solutions (with multiple variations, choosing the best score per split): IsoForest, OCSVM, and CVDD. We largely outperform competitors with an average improvement of 13.5% on 20Newsgroups and 6.9% on AG News compared with the next best solution. Note that DATE uses the same set of hyper-parameters per dataset.

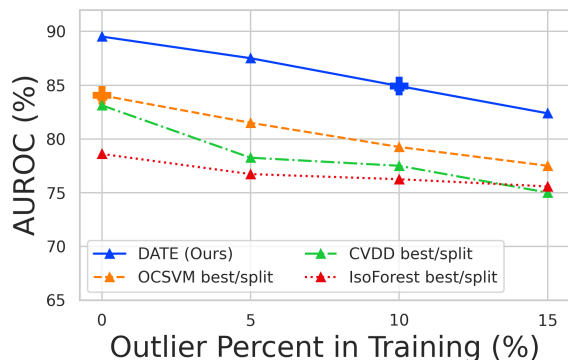


Figure 3. Unsupervised AD. We test the performance of our method when training on impure data, which contains anomalies in various percentages: 0%-15%. The performance slowly decreases when we increase the anomaly percentage, but even at 10% contamination, it is still better than state-of-the-art results on self-supervised anomaly detection in text (Ruff et al., 2019), which trains on 0% anomalous data, proving the robustness of our method. Experiments were done on all AG News splits.

For visualization, we chose two splits from AG News and report the scores from the beginning of the training to the end. At the beginning, the outliers’ distribution of scores fully overlaps with the one for inliers, but in the end the two are well separated, proving the effectiveness of our method.

[‡]Experiments done using the CVDD published code <https://github.com/lukasruff/CVDD-PyTorch>.

Inlier	Label	Pred	Sample (BERT tokens)
Sports	Outlier (World)	Outlier	jail #ing democrat china politically motivated af #p af #p hong kong democrats accused china jail #ing one members trump #ed prostitution charges bid disgrace political movement beijing feud #ing seven years
Sci	Outlier (World)	Outlier	panama flooding kills nine people least nine people seven children died flooding capital panama authorities say least people still missing heavy rainfall caused rivers break banks
Business	Inlier	Inlier	motorola cut jobs new york reuters telecommunications equipment maker motorola inc hr #ef http www investor reuters com full #qu #ote as #p #x tick #er mo #t target stocks quick #in #fo full #qu #t mo #t said tuesday would cut jobs take related charges million focus wireless business

Figure 4. Qualitative examples. Lower scores are shown in a more intense red, and point to anomalies. In the 1st example, words from politics are flagged as anomalous for sports. In the 2nd one, words describing natural events are outliers for technology. In the 3rd row, while few words have higher anomaly potential for the business domain, most of them are appropriate.

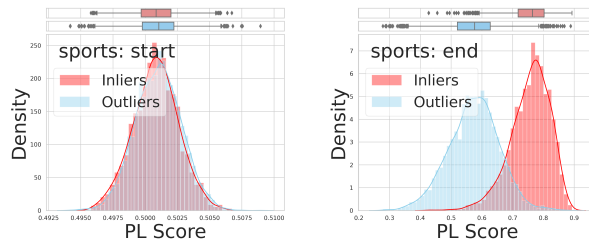


Figure 5. Normalized histogram for anomaly score. We see how the anomaly score (PL) distribution varies among inliers and outliers, from the beginning of the training (1st column) to the end (2nd column), where the two become well separated, with relatively low interference between classes.

4. Conclusion

We propose DATE, a model for tackling Anomaly Detection in Text, and formulate an innovative self-supervised task, based on masking parts of the initial input and predicting which mask pattern was used. After masking, a generator reconstructs the initially masked tokens and the discriminator predicts which mask was used. We optimize a loss composed of both token and sequence-level parts, which stabilizes learning and improves the AD performance. For computing the anomaly score, we alleviate the burden of aggregating predictions from multiple transformations by introducing an efficient variant of the Pseudo Label score, which is applied per token, only on the original input. We show that this score separates very well the abnormal entries from normal ones, leading DATE to outperform state-of-the-art results on all AD splits from 20Newsgroups and AG News datasets, by a large margin, both in the semi-supervised and unsupervised AD settings.

References

- Aggarwal, C. C. Outlier analysis. In *Data mining*, pp. 237–263. Springer, 2015.
- Banoth, L., Teja, M., Saicharan, M., and Chandra, N. A survey of data mining and machine learning methods for cyber security intrusion detection. *International Journal of Research*, 4:406–412, 2017.
- Beltagy, I., Lo, K., and Cohan, A. Scibert: A pretrained language model for scientific text. In *EMNLP/IJCNLP*, 2019.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3(null):993–1022, March 2003. ISSN 1532-4435.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, Dec 2017. ISSN 2307-387X. doi: 10.1162/tacl_a.00051. URL http://dx.doi.org/10.1162/tacl_a.00051.
- Chandola, V., Banerjee, A., and Kumar, V. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3): 1–58, 2009.
- Chen, J., Sathe, S., Aggarwal, C. C., and Turaga, D. S. Outlier detection with autoencoder ensembles. In *SDM*, 2017.
- Clark, K., Luong, M.-T., Le, Q. V., and Manning, C. D. Electra: Pre-training text encoders as discriminators rather than generators. *ArXiv*, abs/2003.10555, 2020.
- de Vries, W., van Cranenburgh, A., Bisazza, A., Caselli, T., van Noord, G., and Nissim, M. Bertje: A dutch bert model. *ArXiv*, abs/1912.09582, 2019.
- Deng, L. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE Signal Processing Magazine*, 29:141–142, 2012.
- Devlin, J., Chang, M., Lee, K., and Toutanova, K. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018. URL <http://arxiv.org/abs/1810.04805>.
- Dorransoro, J., Ginel, F., Sanchez, C., and Cruz, C. Neural fraud detection in credit card operations. *IEEE transactions on neural networks*, 84:827–34, 1997.
- Golan, I. and El-Yaniv, R. Deep anomaly detection using geometric transformations. *CoRR*, abs/1805.10917, 2018. URL <http://arxiv.org/abs/1805.10917>.
- Hawkins, D. M. *Identification of outliers*, volume 11. Springer, 1980.
- Hawkins, S., He, H., Williams, G. J., and Baxter, R. A. Outlier detection using replicator neural networks. In *Proceedings of the 4th International Conference on Data Warehousing and Knowledge Discovery, DaWaK 2000*, pp. 170–180, Berlin, Heidelberg, 2002. Springer-Verlag. ISBN 3540441239.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <https://doi.org/10.1162/neco.1997.9.8.1735>.
- Joulin, A., Grave, E., Bojanowski, P., and Mikolov, T. Bag of tricks for efficient text classification. *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, 2017. doi: 10.18653/v1/e17-2068. URL <http://dx.doi.org/10.18653/v1/E17-2068>.
- Kammerer, K., Hoppenstedt, B., Pryss, R., Stöckler, S., Allgaier, J., and Reichert, M. Anomaly detections for manufacturing systems based on sensor data—insights into two challenging real-world production settings. *Sensors (Basel, Switzerland)*, 19, 2019.
- Kannan, R., Woo, H., Aggarwal, C. C., and Park, H. Outlier detection for text data : An extended version. *CoRR*, abs/1701.01325, 2017. URL <http://arxiv.org/abs/1701.01325>.
- Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., and Soricut, R. Albert: A lite bert for self-supervised learning of language representations, 2019.
- Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C. H., and Kang, J. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 2020.
- Liu, F. T., Ting, K. M., and Zhou, Z.-H. Isolation forest. In *2008 Eighth IEEE International Conference on Data Mining*, pp. 413–422. IEEE, 2008.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *ICLR 2019*, 2019.
- Mahapatra, A., Srivastava, N., and Srivastava, J. Contextual anomaly detection in text data. *Algorithms*, 5(4):469–489, 2012.
- Manevitz, L. and Yousef, M. One-class svms for document classification. *J. Mach. Learn. Res.*, 2:139–154, 2001.
- Manevitz, L. and Yousef, M. One-class document classification via neural networks. *Neurocomputing*, 70:1466–1481, 2007.

- Martin, L., Muller, B., Suárez, P. J. O., Dupont, Y., Romary, L., de la Clergerie, E., Seddah, D., and Sagot, B. Camembert: a tasty french language model. *ArXiv*, abs/1911.03894, 2020.
- Mikolov, T., Chen, K., Corrado, G. S., and Dean, J. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
- Miller, G. A. Wordnet: A lexical database for english. *Commun. ACM*, 38(11):39–41, November 1995. ISSN 0001-0782. doi: 10.1145/219717.219748. URL <https://doi.org/10.1145/219717.219748>.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch, 2017.
- Pennington, J., Socher, R., and Manning, C. D. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1532–1543, 2014. URL <http://www.aclweb.org/anthology/D14-1162>.
- Peters, M., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. Deep contextualized word representations. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2018. doi: 10.18653/v1/n18-1202. URL <http://dx.doi.org/10.18653/v1/N18-1202>.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners, 2019.
- Ruff, L., Zemlyanskiy, Y., Vandermeulen, R., Schnake, T., and Kloft, M. Self-attentive, multi-context one-class classification for unsupervised anomaly detection on text. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4061–4071, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1398. URL <https://www.aclweb.org/anthology/P19-1398>.
- Sakurada, M. and Yairi, T. Anomaly detection using autoencoders with nonlinear dimensionality reduction. In *Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis*, MLSDA’14, pp. 4–11, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450331593. doi: 10.1145/2689746.2689747. URL <https://doi.org/10.1145/2689746.2689747>.
- Schlegl, T., Seeböck, P., Waldstein, S. M., Schmidt-Erfurth, U., and Langs, G. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *IPMI*, 2017.
- Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A., and Williamson, R. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13:1443–1471, 2001a.
- Schölkopf, B., Platt, J. C., Shawe-Taylor, J., Smola, A. J., and Williamson, R. C. Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7): 1443–1471, 2001b.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *CoRR*, abs/1706.03762, 2017. URL <http://arxiv.org/abs/1706.03762>.
- Wang, S., Zeng, Y., Liu, X., Zhu, E., Yin, J., Xu, C., and Kloft, M. Effective end-to-end unsupervised outlier detection via inlier priority of discriminative network. In *NeurIPS*, 2019.
- Zhang, X., Zhao, J. J., and LeCun, Y. Character-level convolutional networks for text classification. In *NeurIPS*, 2015.

A. DATE Architecture

Input masking. Let $\mathbf{m} \in \{0, 1\}^T$ be a mask pattern corresponding to the text input $\mathbf{x} = [x_1, x_2, \dots, x_T]$. For training, we generate and fix K mask patterns $\mathbf{m}^{(1)}, \mathbf{m}^{(2)}, \dots, \mathbf{m}^{(K)}$ by randomly sampling a constant number of ones. Instead of masking random tokens on-the-fly as in ELECTRA, we first sample a mask pattern from the K predefined ones. Next we apply it to the input, as in Fig. 1. Let $\hat{\mathbf{x}}(\mathbf{m}) = [\hat{x}_1, \hat{x}_2, \dots, \hat{x}_T]$ be the input sequence \mathbf{x} , masked with \mathbf{m} , where:

$$\hat{x}_i = \begin{cases} x_i, & \mathbf{m}_i = 0 \\ [\text{MASK}], & \mathbf{m}_i = 1 \end{cases}$$

For instance, given an input $\mathbf{x} = [\text{bank}, \text{hikes}, \text{prices}, \text{before}, \text{election}]$ and a mask pattern $\mathbf{m} = [0, 0, 1, 0, 1]$, the masked input is $\hat{\mathbf{x}}(\mathbf{m}) = [\text{bank}, \text{hikes}, [\text{MASK}], \text{before}, [\text{MASK}]]$.

Replacing [MASK]s. Each masked token can be replaced with a word token (*e.g.* by sampling uniformly from the vocabulary). For more plausible alternatives, masked tokens can be sampled from a Masked Language Model (MLM) generator such as BERT, which outputs a probability distribution P_G over the vocabulary, for each token. Let $\tilde{\mathbf{x}}(\mathbf{m}) = [\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_T]$ be the plausibly corrupted text, where:

$$\tilde{x}_i = \begin{cases} x_i, & \mathbf{m}_i = 0 \\ w_i \sim P_G(x_i | \hat{\mathbf{x}}(\mathbf{m}); \theta_G), & \mathbf{m}_i = 1 \end{cases}$$

For instance, given the masked input $\hat{\mathbf{x}}(\mathbf{m}) = [\text{bank, hikes, [MASK], before, [MASK]}]$, a plausibly corrupted input is $\tilde{\mathbf{x}}(\mathbf{m}) = [\text{bank, hikes, fees, before, referendum}]$.

The BERT encoder transforms an input token sequence $\mathbf{x} = [x_1, x_2, \dots, x_T]$ into a sequence of contextualized word embeddings $h(\mathbf{x}) = [h_1, h_2, \dots, h_T]$.

Generator. G is a BERT encoder with a linear layer on top that outputs the probability distribution P_G for each token. The generator is trained using the MLM loss:

$$\mathcal{L}_{MLM} = E \left[\sum_{\substack{i=1; \\ s.t. m_i=1}}^T -\log P_G(x_i | \hat{\mathbf{x}}(\mathbf{m}); \theta_G) \right] \quad (1)$$

Discriminator. D is a BERT encoder with two prediction heads applied over the contextualized word representations:

i. RMD head. This head outputs a vector of logits for all mask patterns $\mathbf{o} = [o_1, \dots, o_K]$. We use the contextualized hidden vector $h_{[\text{CLS}]}$ (corresponding to the $[\text{CLS}]$ special token at the beginning of the input) for computing the mask logits \mathbf{o} and P_M , the probability of each mask pattern:

$$P_M(\mathbf{m} = \mathbf{m}^{(k)} | \tilde{\mathbf{x}}(\mathbf{m}^{(k)}); \theta_D) = \frac{\exp(o_k)}{\sum_{i=1}^K \exp(o_i)} \quad (2)$$

ii. RTD head. This head outputs scores for the two classes (*original* and *replaced*) for each token x_1, x_2, \dots, x_T , by using the contextualized hidden vectors h_1, h_2, \dots, h_T .

Loss. We train the DATE network in a maximum-likelihood fashion using the \mathcal{L}_{DATE} loss:

$$\min_{\theta_D, \theta_G} \sum_{\mathbf{x} \in \mathcal{X}} \mathcal{L}_{DATE}(\theta_D, \theta_G; \mathbf{x}) \quad (3)$$

The loss contains both the token-level losses in ELECTRA, as well as the sequence-level mask detection loss \mathcal{L}_{RMD} :

$$\mathcal{L}_{DATE}(\theta_D, \theta_G; \mathbf{x}) = \mu \mathcal{L}_{RMD}(\theta_D; \mathbf{x}) + \mathcal{L}_{MLM}(\theta_G; \mathbf{x}) + \lambda \mathcal{L}_{RTD}(\theta_D; \mathbf{x}), \quad (4)$$

where the discriminator losses are:

$$\mathcal{L}_{RMD} = \mathbb{E} \left[-\log P_M(\mathbf{m} | \tilde{\mathbf{x}}(\mathbf{m}); \theta_D) \right], \quad (5)$$

$$\mathcal{L}_{RTD} = \mathbb{E} \left[\sum_{\substack{i=1; \\ x_i \neq [\text{CLS}]}}^T -\log P_D(m_i | \tilde{\mathbf{x}}(\mathbf{m}); \theta_D) \right], \quad (6)$$

where P_D is the probability distribution that a token was replaced or not.

The ELECTRA loss enables D to learn good feature representations for language understanding. Our RMD loss puts the representation in a larger sequence-level context. After pre-training, G is discarded and D can be used as a general-purpose text encoder for downstream tasks. Output probabilities from D are further used to compute an anomaly score for new examples.

PL over RTD classification scores (detailed) We show next the equation for the Pseudo Label score:

$$PL_{RTD}(x) = \frac{1}{T} \sum_{i=1}^T P_D(m_i = 0 | \tilde{\mathbf{x}}(\mathbf{m}^{(0)}); \theta_D), \quad (7)$$

where $\mathbf{m}^{(0)} = [0, 0, \dots, 0]$ effectively leaves the input unchanged. As can be seen in Fig. 2, the RTD head will be less certain in predicting the *original* class for *outliers* (having a probability distribution unseen at training time), which will lead to lower PL scores for *outliers* and higher PL scores for *inliers*. We use PL at testing time, when the entire input is either normal or abnormal. Our method also speeds up inference, since we only do one feed-forward pass through the discriminator instead of k passes. Moreover, having a per token anomaly score helps us better understand and visualize the behavior of our model, as shown in Fig. 4.

B. Experimental setup details

In this section, we detail the empirical validation of our method by presenting: the semi-supervised and unsupervised experimental setup, a comprehensive ablation study on DATE, and the comparison with state-of-the-art on the semi-supervised and unsupervised AD tasks. DATE does not use any form of pre-training or knowledge transfer (from other datasets or tasks), learning all the embeddings from scratch. Using pre-training would introduce unwanted prior knowledge about the outliers, making our model considering them known (normal).

Model and Training. For training the DATE network we follow the pipeline in Fig. 1. We detail next the modules in our model: for the ablation experiments in Tab. 2, **Generator (small)**: 1 Transformer layer, with 4 self-attention heads, token and positional embeddings of size 128, hidden layer of size 16, feedforward layer of sizes 1024 and 16; **Generator (large)**: 1 Transformer layer, with 4 self-attention heads, token and positional embeddings of size 128, hidden layer of size 64, feedforward layer of sizes 1024 and 64; As empirical experiments showed us, we choose a *random Generator* (samples were drawn from a uniform distribution over the vocabulary) in our final model. **Discriminator**: 4 Transformer layers, each with 4 self-attention heads, hidden layers of size 256, feedforward layers of sizes of 1024 and 256, 128-dimensional token and positional

embeddings, which are *tied* with the generator. For other unspecified hyper-parameters we use the ones in ELECTRA-Small model. **Prediction Heads:** both heads have 2 linear layers separated by a non-linearity, ending in a classification. **Loss weights:** We set the RTD λ weight to 50 as in (Clark et al., 2020), and the RMD μ weight to 100.

We train the networks with AdamW with amsgrad (Loshchilov & Hutter, 2019), $1e^{-5}$ learning rate, using sequences of maximum length 128 for AG News, and 498 for 20Newsgroups. We use $K = 50$ predefined masks, covering 50% of the input for AG News and $K = 25$, covering 25% for 20Newsgroups. The training converges on average after 5000 update steps and the inference time is 0.005 sec/sample in PyTorch (Paszke et al., 2017), on a single GTX Titan X.

B.1. Comparison with other AD methods

We compare DATE against the methods bellow:

OCSVM. We use the One-Class SVM model implemented in the CVDD work [‡]. For each split, we choose the best configuration (fastText vs Glove, rbf vs linear kernel, $\nu \in [0.05, 0.1, 0.2, 0.5]$).

Isolation Forest. We apply it over fastText or Glove embeddings, varying the number of estimators (64, 100, 128, 256), and choosing the best model per split. In the unsupervised AD setup, we manually set the percent of outliers in the train set.

CVDD. This model (Ruff et al., 2019) is the current state-of-the-art solution for AD on text. For each split, we chose the best column out of all reported context sizes (r). The scores reported using the c^* context vector depends on the ground truth and it only reveals "the potential of contextual anomaly detection", as the authors mention.

C. Ablation studies

To better understand the impact of different components in our model and making the best decisions towards a higher performance, we perform an extensive set of experiments (see Tab. 2). Note that we successively treat each AG News split as inlier and report the mean and standard deviations over the four splits. The results show that our model is robust to domain shifts.

A. Anomaly score. We explore three anomaly scores introduced in the E^3 Outlier framework (Wang et al., 2019) on semi-supervised and unsupervised AD tasks in Computer Vision: Maximum Probability (MP), Negative Entropy (NE) and our modified Pseudo Label (PL_{RTD}). These scores are computed using the softmax probabilities from the fi-

Abl.	Method	Variation	AUROC(%)
	CVDD	best	83.1 \pm 4.4
	OCSVM	best	84.0 \pm 5.0
	ELECTRA	adapted for AD	84.6 \pm 4.5
	DATE	(Ours)	90.0 \pm 4.2
A.	Anomaly score	MP	72.4 \pm 3.7
		NE	73.1 \pm 3.9
B.	Generator	small	89.3 \pm 4.2
		large	89.8 \pm 4.4
C.	Loss func	RTD only	89.4 \pm 4.4
		RMD only	85.9 \pm 4.1
D.	Masking patterns	5 masks	87.5 \pm 4.5
		10 masks	89.2 \pm 4.3
		25 masks	89.8 \pm 4.3
		100 masks	89.8 \pm 4.3
E.	Mask percent	15%	89.5 \pm 4.1
		25%	89.5 \pm 4.1

Table 2. Ablation study. We show results for the competition and report ablation experiments which are only one change away from our best **DATE** configuration: **A.** PL_{RTD} ; **B.** Rand **C.** RTD + RMD; **D.** 50 masks; **E.** 50%. For the ELECTRA line, we use: A. PL_{RTD} ; B. Rand; C. RTD only; D. Unlimited; E. 15%. A. The Anomaly Score used over classification probabilities shows that PL_{RTD} (used in DATE) is the best in predicting anomalies, meaning that our self-supervised classification task is well defined, with few ambiguous samples; B. A learned Generator does not justify its training cost; C. RMD Loss proved to be complementary with RTD Loss, their combination (in DATE) increasing the score and stabilizes the training; D+E.

nal classification layer of the discriminator. PL is an ideal score if the self-supervised task manages to build and learn well separated classes. The way we formulate our mask prediction task enables a very good class separation, as theoretically proved in detail in the Appendix E. Therefore, PL_{RTD} proves to be significantly better in detecting the anomalies compared with MP and NE metrics, which try to compensate for ambiguous samples.

B. Generator performance. We tested the importance of having a learned generator, by using a one-layer Transformer with hidden size 16 (small) or 64 (large). The *random generator* proved to be better than both parameterized generators.

C. Loss function. For the final loss, we combined RTD (which sanctions the prediction per token) with our RMD (which enforces the detection of the mask applied on the entire sequence). We also train our model with RTD or RMD only, obtaining weaker results. This proves that combining losses with supervisions at different scales (locally:

token-level and globally: sequence-level) improves AD performance. Moreover, when using only the RTD loss, the training can be very unstable (AUROC score peaks in the early stages, followed by a steep decrease). With the combined loss, the AUROC is only stationary or increases with time.

D. Masking patterns. The mask patterns are the root of our task formulation, hiding a part of the input tokens and asking the discriminator to classify them. As experimentally shown, having more mask patterns is better, encouraging increased expressiveness in the embeddings. Too many masks on the other hand can make the task too difficult for the discriminator and our ablation shows that having more masks does not add any benefit after a point. We validate the percentage of masked tokens in **E. Mask percent** ablation.

D. Relation to prior work

Our work relates to self-supervision for language representation as well as self-supervision for learning features of normality in AD.

D.1. Self-supervision for NLP

Self-supervision has been the bedrock of learning good feature representations in NLP. The earliest neural methods leveraged shallow models to produce static word embeddings such as word2vec (Mikolov et al., 2013), GloVe (Pennington et al., 2014) or fastText (Bojanowski et al., 2017; Joulin et al., 2017). More recently, contextual word embeddings have produced state-of-the-art results in many NLP tasks, enabled by Transformer-based (Vaswani et al., 2017) or LSTM-based (Hochreiter & Schmidhuber, 1997) architectures, trained with language modeling (Peters et al., 2018; Radford et al., 2019) or masked language modeling (Devlin et al., 2018) tasks.

Many improvements and adaptations have been proposed over the original BERT, which address other languages (Martin et al., 2020; de Vries et al., 2019), domain specific solutions (Beltagy et al., 2019; Lee et al., 2020) or more efficient pre-training models such as ALBERT (Lan et al., 2019) or ELECTRA (Clark et al., 2020). ELECTRA pre-trains a BERT-like generator and discriminator with a Replacement Token Detection (RTD) Task. The generator substitutes masked tokens with likely alternatives and the discriminator is trained to distinguish between the original and masked tokens.

D.2. Self-supervised classification for AD

Typical representation learning approaches to deep AD involve learning features of normality using autoencoders (Hawkins et al., 2002; Sakurada & Yairi, 2014; Chen et al., 2017) or generative adversarial networks (Schlegl

et al., 2017). More recent methods train the discriminator in a self-supervised fashion, leading to better normality features and anomaly scores. These solutions mostly focus on image data (Golan & El-Yaniv, 2018; Wang et al., 2019) and train a model to distinguish between different transformations applied to the images (*e.g.* rotation, flipping, shifting). An interesting property that justifies self-supervision under unsupervised AD is called *inlier priority* (Wang et al., 2019), which states that during training, inliers (normal instances) induce higher gradient magnitudes than outliers, biasing the network’s update directions towards reducing their loss. Due to this property, the outputs for *inliers* are more consistent than for *outliers*, enabling them to be used as anomaly scores.

D.3. AD for text

There are a few shallow methods for AD on text, usually operating on traditional document-term matrices. One of them uses one-class SVMs (Schölkopf et al., 2001a) over different sparse document representations (Manevitz & Yousef, 2001). Another method uses non-negative matrix factorization to decompose the term-document matrix into a low-rank and an outlier matrix (Kannan et al., 2017). LDA-based (Blei et al., 2003) clustering algorithms are augmented with semantic context derived from WordNet (Miller, 1995) or from the web to detect anomalies (Mahapatra et al., 2012).

D.4. Deep AD for text

While many deep AD methods have been developed for other domains, few approaches use neural networks or pre-trained word embeddings for text anomalies. Earlier methods use autoencoders (Manevitz & Yousef, 2007) to build document representations. More recently, pre-trained word embeddings and self-attention were used to build contextual word embeddings (Ruff et al., 2019). These are jointly optimized with a set of *context vectors*, which act as topic centroids. The network thus discovers relevant topics and transforms normal examples such that their contextual word embeddings stay close to the topic centroids. Under this setup, anomalous instances have contextual word embeddings which on average deviate more from the centroids.

E. Disjoint patterns analysis

We start from two observations regarding the performance of DATE, our Anomaly Detection algorithm. First, a discriminative task performs better if the classes are well separated (Deng, 2012) and there is a low probability for confusions. Second, the PL score for anomalies achieves best performance when the probability distribution for its input is clearly separated. Intuitively, for three classes, $PL([0.9, 0.05, 0.05])$ is better than $PL([0.5, 0.3, 0.2])$ because it allows PL to give either near 1 score if the class is correct,

either near 0 score if it is not, avoiding the zone in the middle where we depend on a well chosen threshold.

Since the separation between the mask patterns greatly influences our final performance, we next analyze our AD task from the mask pattern generation point of view. Ideally, we want to have a sense of how disjoint our randomly sampled patterns are and make an informed choice for the pattern generation hyper-parameters.

First, we start by computing an upper bound for the probability of having two patterns with at least p common masked points. We have $\binom{S}{M}$ patterns, where S is the sequence length and M is the number of masked tokens. We fix the first p positions that we want to mask in any pattern. Considering those fixed masks, the probability of having a sequence with M masked tokens, with p tokens in the first positions is r :

$$r = \frac{\binom{S-p}{M-p}}{\binom{S}{M}}. \quad (8)$$

Next, the probability that two sequences mask the first p tokens is r^2 . But we can choose those two positions in a $\binom{S}{p}$ ways. So the probability that any two sequences have at least p common masked tokens is lower than UB_2 :

$$UB_2 = \binom{S}{p} r^2 \quad (9)$$

Next, out of our generated patterns, we sample N masks, so the probability becomes less than the upper bound UB_N :

$$\begin{aligned} UB_N &= \binom{N}{2} UB_2 = \binom{N}{2} \binom{S}{p} r^2 \\ &= \binom{N}{2} \binom{S}{p} \left(\frac{\binom{S-p}{M-p}}{\binom{S}{M}} \right)^2. \end{aligned} \quad (10)$$

In our experiments, the sequence length is $S = 128$ and we chose the number of masked tokens to be between 15% and 50% (M between 19 and 64). We consider that two patterns are disjoint when they have less than p masked tokens in common, for N sampled patterns.

The probability that any two patterns collide (have more than p masked tokens in common) is very low. We compute several values for its upper bound: $UB_{N=100,p=12} = 5e - 4$, $UB_{N=100,p=15} = 1e - 9$, $UB_{N=10,p=15} = 1e - 11$, $UB_{N=10,p=13} = 1e - 7$.

In conclusion, for our specific setup, the probability for two masks to largely overlap (large p compared with S) is extremely small, ensuring us a good performance in the discriminator. We take advantage of this property of our pre-

Inlier	Label	Pred	Sample (BERT tokens)
Sports	Outlier (Business)	Outlier	airways watch bankruptcy court judge stephen mitchell hear arguments today asking rec ##ons ##ider four month percent pay cut imposed many union ##ized workers last month
Sci	Outlier (World)	Outlier	bush defend ##s decision invade iraq president bush went skeptical hall world leaders tuesday mount vigorous defense war iraq telling united nations iraqi people
Business	Outlier (Sci)	Outlier	nokia nec test new multimedia sub ##sy ##ste ##m two high tech communications players completed first phase series tests show next generation data communications infrastructure works
World	Outlier (Sports)	Outlier	man plan athens four years ago sydney gymnast ##s gone medal free olympics first time years federation president bob cola ##ross ##i sitting table explaining turn ##around already begun women moved sixth fourth world one year men sixth fifth
Sports	Inlier	Inlier	dolphins steelers play sunday night reuters reuters miami dolphins pittsburgh steelers play scheduled game sunday night [SEP]
Sci	Inlier	Inlier	microsoft launch new search engine software giant microsoft corp decided release beta versions updated ms ##n search engine earlier today company hopes compete leading search engines google yahoo [SEP]

Figure 6. More qualitative examples.

Subset	business	sci	sports	world
AUPR-in	74.8	62.4	88.8	81.9
AUPR-out	96.1	93.5	98.5	95.5

Table 3. We report AUPR metric for AG News splits, on inliers and outliers since this is a more relevant metric for unbalanced classes (which is the case for all splits in text AD, as explained in Anomalies setup).

text task by combining the discriminator output probabilities with the PL score.

F. More qualitative and quantitative Results

In Fig. 6 we show more qualitative results, trained on different inliers. To encourage further more detailed comparisons, we report the AUPR metric on AG News for inliers and outliers (see Tab. 3). When all the other metrics are almost saturated, we notice that AUPR-in better captures the performance on a certain split.

Acknowledgments. This work has been supported in part by UEFISCDI, under Project PN-III-P2-2.1-PTE-2019-0532.