

Is BERT Robust to Label Noise? A Study on Learning with Noisy Labels in Text Classification

Dawei Zhu, Michael A. Hedderich, Fangzhou Zhai, David Ifeoluwa Adelani & Dietrich Klakow

Saarland University, Saarland Informatics Campus, Germany

{dzhu, mhedderich, didelani, dietrich.klakow}@lsv.uni-saarland.de
fzhai@coli.uni-saarland.de

Abstract

Incorrect labels in training data occur when human annotators make mistakes or when the data is generated via weak or distant supervision. It has been shown that complex noise-handling techniques - by modeling, cleaning or filtering the noisy instances - are required to prevent models from fitting this label noise. However, we show in this work that, for text classification tasks with modern NLP models like BERT, over a variety of noise types, existing noise-handling methods do not always improve its performance, and may even deteriorate it, suggesting the need for further investigation. We also back our observations with a comprehensive analysis.

1 Introduction

For many languages, domains and tasks, large datasets with high-quality labels are not available. To tackle this issue, cheaper data acquisition methods have been suggested, such as crowdsourcing or automatic annotation methods like weak and distant supervision. Unfortunately, compared to gold-standard data, these approaches come with more labeling mistakes, which are known as noisy labels. Noise-handling has become an established approach to mitigate the negative impact of learning with noisy labels. A variety of methods have been proposed that model the noise, or clean and filter the noisy instances (Hedderich et al., 2021; Algan and Ulusoy, 2021). Jindal et al. (2019) show e.g. a 30% boost in performance after applying noise-handling techniques on a CNN-based text classifier.

In a recent work, Tänzer et al. (2021) showed that BERT (Devlin et al., 2019) has an inherent robustness against noisy labels. The generalization performance on the clean distribution drops only slowly with the increase of the mislabeled samples. Also, they show that early-stopping is crucial for learning with noisy labels as BERT will eventually memorize all wrong labels when trained long

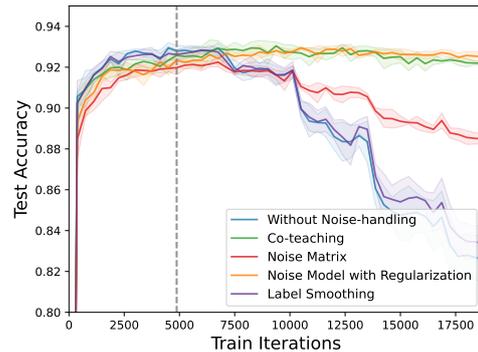


Figure 1: A typical training curve when learning with noise. Learning without noise-handling (blue) will reach a peak accuracy before memorizing the noise. Early-stopping on a noisy validation set (vertical grey line) is often sufficient to find such a peak. Injected uniform noise of 40% on AG-News dataset.

enough. However, their experiments only focus on a single type of noise and a limited range of noise levels. It remains unclear if BERT still performs robustly under a wider range of noise types and with higher fractions of mislabeled samples. Moreover, they perform early-stopping on a clean validation set, which may not be available under low resource settings. Last but not least, they do not compare to any noise-handling methods.

In this work, we investigate the behaviors of BERT on tasks with different noise types and noise levels. We also study the effect of noise-handling methods under these settings. Our main results include (1) BERT is robust against injected noise, but could be vulnerable to noise from weak supervision. In fact, the latter, even at a low level, can be more challenging than high injected noise. (2) Existing noise-handling methods do not improve the peak performance of BERT under any noise settings we investigated; as is shown with further analysis, noise-handling methods rarely render the correct labels more distinguishable from the incorrect ones.¹

¹Our implementation is available on: <https://github.com/uds-lsv/BERT-LNL>.

2 Learning with Noisy Labels

Problem Settings We consider a k -class classification problem. Let D denote the true data generation distribution over $\mathcal{X} \times \mathcal{Y}$ where \mathcal{X} is the feature space and $\mathcal{Y} = \{1, \dots, k\}$ is the label space. In a typical classification task, we are provided with a training dataset $S = \{(x_i, y_i)_{i=1}^n\}$ sampled from D . In learning with noisy labels, however, we have no access to D . Instead, a noisy training set $\hat{S} = \{(x_i, \hat{y}_i)_{i=1}^n\}$ sampled from a label-corrupted data distribution \hat{D} . The goal is to learn a classifier that generalizes well on the clean distribution by only exploiting \hat{S} .

Injected Label Noise To rigorously evaluate noise-handling methods at different noise levels, researchers in this area often construct noisy datasets from clean ones by injecting noise. This can, e.g., reflect annotation scenarios such as crowdsourcing, where some annotators answer randomly or prefer an early entry in a list of options. Modeling such noise is achieved by flipping the labels of the clean instances according to a pre-defined noise level $\varepsilon \in [0, 1)$ and a noise type. There are two commonly used noise types: the single-flip noise (Reed et al., 2015):

$$p_{\text{flip}}(\hat{y} = j | y = i) = \begin{cases} 1 - \varepsilon, & \text{for } i = j \\ \varepsilon, & \text{for one } i \neq j \\ 0, & \text{else} \end{cases}$$

and uniform-flip (van Rooyen et al., 2015) noise

$$p_{\text{uni}}(\hat{y} = j | y = i) = \begin{cases} 1 - \varepsilon, & \text{for } i = j \\ \frac{\varepsilon}{k-1}, & \text{for } i \neq j \end{cases}.$$

These noise generation processes are feature-independent, i.e. $p(\cdot | y = i, x) = p(\cdot | y = i)$. Therefore, they can be described by a noise transition matrix T with $T_{ij} := p(\hat{y} = j | y = i)$. It is usually assumed that the noise is diagonally-dominant when generating the noisy labels, i.e. $\forall i, T_{ii} > \max_{j \neq i} T_{ij}$.

Label Noise from Weak Supervision Distant and weak supervision (Mintz et al., 2009; Ratner et al., 2016) have become essential methods to acquire labeled data in low-resource scenarios. The resulting noise, unlike injected noise, is often feature-dependent (Lange et al., 2019). We evaluate our methods on two real-world datasets in Hausa and Yorùbá to cover this type of noise.

Dataset	Classes	Average Lengths	Train Samples	Validation Samples	Test Samples	Train Noise Level
IMDB	2	292	21246	3754	25000	various
AG-News	4	44	108000	12000	7600	various
Yorùbá	7	13	1340	189	379	33.28%
Hausa	5	10	2045	290	582	50.37%

Table 1: Statistics of the text classification datasets. The train noise level is the false discovery rate (i.e. 1-precision) of the noisy labels in the training set. The original AG-News has 120k training instances and no validation instances. We therefore held-out 10% of the training samples for validation.

3 Early-Stopping on Noisy Validation Set

When trained on noisy data without noise-handling, BERT reaches a high generalization performance before it starts fitting the noise. Then it memorizes the noise and the performance on clean distribution drops dramatically (the blue curve in Figure 1). Hence, for models without noise-handling, it is crucial to stop training when the generalization performance reaches its maximum.

Tänzer et al. (2021) use a clean validation set to find this point. However, a clean validation set is often unavailable in realistic low-resource scenarios as it requires manual annotation. Therefore, we use a noisy validation set for early-stopping in all of our experiments and we attain models that generalize well on the clean distribution.

In our example in Figure 1, we see that while most noise-handling methods prevent BERT from fitting the noise in the long run, their peak performance is not significantly higher than a vanilla model without noise-handling.

4 Experiments

Dataset Construction We experiment with four text classification datasets: two benchmarks, AG-News (Zhang et al., 2015) and IMDB (Maas et al., 2011), injected with different levels of single-flip or uniform noise; for the weakly supervised noise, we make use of two news topics datasets in two low-resource languages: Hausa and Yorùbá (Hedderich et al., 2020). Hausa and Yorùbá are the second and the third most spoken indigenous language in Africa, with 40 and 35 million native speakers, respectively (Eberhard et al., 2019). The noisy labels were gazetteered. For example, to identify texts for the class ‘‘Africa’’, a labeling rule based on a list of African countries and their capitals is used. Note that while we can vary the noise levels of injected noise, the amount of weak supervision

noise in Hausa and Yorùbá is fixed². We summarize some basic statistics of the datasets in Table 1.

Implementation We use of-the-shelf BERT models for our tasks. Specifically, we apply the BERT-base model for AG-News and IMDB, and the mBERT-base for Yorùbá and Hausa. The fine-tuning approach follows (Devlin et al., 2019). In all settings, we apply early-stopping on a noisy validation set to mimic the realistic low-resource settings, while the test set remains clean. For more implementation details and a discussion on clean and noisy validation sets, see Appendix B and E.

4.1 Baselines

We compare learning without noise-handling with four popular noise-handling methods.³

Without Noise-handling Train BERT on the noisy training set as it was clean. A noisy validation set is used for early-stopping.

No Validation For the sake of comparison, we train the model without noise-handling and until the training loss converges.

Noise Matrix A noise transition matrix is appended after BERT’s prediction to transform the clean label distribution to the noisy one. A variety of methods exists for estimating the noise matrix, i.e. Sukhbaatar et al. (2015); Bekker and Goldberger (2016); Patrini et al. (2017); Hendrycks et al. (2018); Yao et al. (2020). To exclude the effects of estimation errors in the evaluation, we use the ground truth transition matrix as it is the best possible estimation. This matrix is fixed after initialization.

Noise Matrix with Regularization The previous state-of-the-art for text classification with noisy labels (Jindal et al., 2019). Similar to *Noise Matrix*, it appends a noise matrix after BERT’s output. During training, the matrix is learned with an l_2 regularization and is not necessarily normalized to be a probability matrix. In the original implementation they use CNN-based models as backbone, we switch it to BERT for fair comparison.

Co-teaching Han et al. (2018) Train two networks to pick cleaner training subsets for each other. The Co-teaching framework requires an estimation

of the noise level. Similarly to NMat, we use the ground truth noise level to exclude the performance drop caused by estimation error.

Label Smoothing Label smoothing (Szegedy et al., 2016) is a commonly used method to improve model’s generalization and calibration. It mixes the one-hot label with a uniform vector, preventing the model from getting overconfident on the samples. Lukasik et al. (2020) further shows that it improves noise robustness.

4.2 Experimental Results

We evaluate our baselines on both injected noise (on AG-News and IMDB) and weak supervision noise (on Hausa and Yorùbá). The test accuracy is presented Figure 2. On injected noise, our results match and extend the findings by Tänzer et al. (2021) that BERT is noise robust. For example, the test accuracy drops only about 10% after injecting 70% wrong labels (Figure 2(a)). However, we find that BERT is vulnerable under weak supervision noise. The performance can drop up to 35% in a dataset like Hausa with 50% weak supervision noise compared to training with clean labels (Figure 2(c)). This indicates that the experience on injected noise may not be transferable to weak supervision noise.

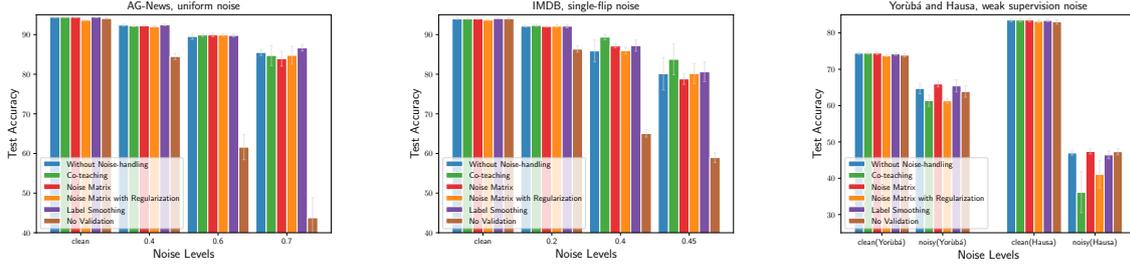
We also observe that noise-handling methods are not always helpful. For injected noise, the benefits from noise-handling become obvious only under high noise levels. But even then, there is no clear winner, meaning that it is hard to decide beforehand which noise method to apply - with the risk that they may even perform worse than BERT without noise-handling. The same applies to weak supervision noise. The maximal performance gap between the best model and BERT without noise-handling is less than 4% and 1.5% under injected noise and weak supervision noise, respectively.

4.3 Analysis of Loss Distributions

To shed some light on why BERT is robust against injected noise but not weak supervision noise, we track the losses on correctly and wrongly labeled samples during training. Figure 4 depicts typical distributions of losses associated with correctly and incorrectly labeled samples, respectively, when early-stopping is triggered. We see that they have minimal overlap, thus different behaviors throughout the training, potentially allowing the model to distinguish correctly and incorrectly labeled sam-

²refer to Appendix A for detailed noise distribution.

³For a fair comparison, early-stopping on a noisy validation set is applied to all four noise-handling methods.

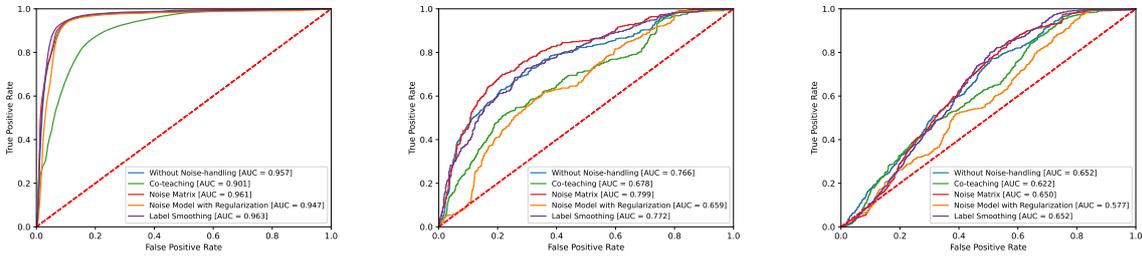


(a) AG-News, uniform noise

(b) IMDB, single-flip noise

(c) Yorùbá & Hausa, weak supervision noise

Figure 2: Test accuracy in different noise settings. a) & b) injected noise with different noise levels c) weak supervision noise, at noise levels of 33.28% and 50.37% in Yorùbá and Hausa, respectively. Noise-handling methods do not always improve peak performances. Further plots in Appendix C.



(a) AG-News - 70% uniform noise

(b) Yorùbá - weak supervision noise

(c) Hausa - weak supervision noise

Figure 3: ROC curves on wrong label detection (binary classification) using the losses. The losses are recorded at the training step when early-stopping is triggered. Noise-handling methods do not make the losses of correct and incorrect labels more distinguishable. Further plots in Appendix D.

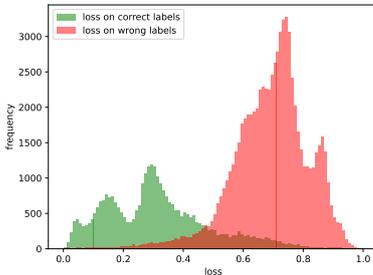


Figure 4: Loss histogram at the training iteration when the early-stopping is triggered. AG-News dataset with 70% uniform noise.

ples from each other. We could further quantify the difference by their separability. Figure 3 presents the receiver operating characteristic (ROC) curves of a threshold-based classifier. We observe that (1) under injected noise, an area under curve (AUC) of more than 90 can be easily achieved without noise-handling (Figure 3(a)), supporting our observation that injected noise has rather a low impact on BERT. (2) Under weak-supervision noise, the AUC score is significantly lower, which means the correct and incorrect labels are less distinguishable. Therefore, BERT fits both labels at similar rates. One reason

could be that the noise in weak supervision is often feature-dependent, it might become easier for BERT to fit them, which in turn deteriorates the generalization. (3) We do not observe a raise in AUC scores when applying noise-handling methods, indicating that noise-handling methods rarely enhance BERT’s ability to further avoid the negative impact of wrong labels. This is consistent with the observation in Section 4.2 that noise-handling methods have little impact on BERT’s generalization performance.

5 Conclusion

On several text classification datasets and for different noise types, we showed that BERT is noise resistant under injected noise, but not necessarily under weak supervision noise. In both cases, the improvement obtained by applying noise-handling methods are limited. Our analysis on the separability of losses corresponding to correct and incorrect labeled samples provides evidence to this argument. Our analysis offers both motivation and insights to further improve label noise-handling methods and make them useful on more realistic types of noise.

6 Broader Impact Statement and Ethics

Noisy labels are a cheaper source of supervision. This could make it easier to use machine learning for improper use cases. However, it also opens up NLP methods for low-resource settings such as under-resourced languages or applications developed by individuals or small organizations. It can, therefore, be a step towards the democratization of AI.

Acknowledgments

This work has been partially funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – Project-ID 232722074 – SFB 1102 and the EU Horizon 2020 projects ROX-ANNE under grant number 833635 and COM-PRISE under grant agreement No. 3081705.

References

- Görkem Algan and Ilkay Ulusoy. 2021. [Image classification with deep learning in the presence of noisy labels: A survey](#). *Knowl. Based Syst.*, 215:106771.
- Alan Joseph Bekker and Jacob Goldberger. 2016. [Training deep neural-networks based on unreliable labels](#). In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2016, Shanghai, China, March 20-25, 2016*, pages 2682–2686. IEEE.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- David M. Eberhard, Gary F. Simons, and Charles D. Fennig (eds.). 2019. [Ethnologue: Languages of the world. twenty-second edition](#).
- Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and Masashi Sugiyama. 2018. [Co-teaching: Robust training of deep neural networks with extremely noisy labels](#). In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 8536–8546.
- Michael A. Hedderich, David Ifeoluwa Adelani, Dawei Zhu, Jesujoba O. Alabi, Udia Markus, and Dietrich Klakow. 2020. [Transfer learning and distant supervision for multilingual transformer models: A study on african languages](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 2580–2591. Association for Computational Linguistics.
- Michael A. Hedderich, Lukas Lange, Heike Adel, Jan-nik Strötgen, and Dietrich Klakow. 2021. [A survey on recent approaches for natural language processing in low-resource scenarios](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2545–2568, Online. Association for Computational Linguistics.
- Dan Hendrycks, Mantas Mazeika, Duncan Wilson, and Kevin Gimpel. 2018. [Using trusted data to train deep networks on labels corrupted by severe noise](#). In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 10477–10486.
- Ishan Jindal, Daniel Pressel, Brian Lester, and Matthew S. Noleby. 2019. [An effective label noise model for DNN text classification](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 3246–3256. Association for Computational Linguistics.
- Lukas Lange, Michael A. Hedderich, and Dietrich Klakow. 2019. [Feature-dependent confusion matrices for low-resource NER labeling with noisy labels](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3554–3559, Hong Kong, China. Association for Computational Linguistics.
- Michal Lukasik, Srinadh Bhojanapalli, Aditya Menon, and Sanjiv Kumar. 2020. [Does label smoothing mitigate label noise?](#) In *International Conference on Machine Learning*, pages 6448–6458. PMLR.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. [Learning word vectors for sentiment analysis](#). In *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA*, pages 142–150. The Association for Computer Linguistics.
- Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. 2009. [Distant supervision for relation extraction without labeled data](#). In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1003–1011, Suntec, Singapore. Association for Computational Linguistics.

- Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. 2017. [Making deep neural networks robust to label noise: A loss correction approach](#). In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 2233–2241. IEEE Computer Society.
- Alexander J Ratner, Christopher M De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. 2016. [Data programming: Creating large training sets, quickly](#). In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Scott E. Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. 2015. [Training deep neural networks on noisy labels with bootstrapping](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Workshop Track Proceedings*.
- Sainbayar Sukhbaatar, Joan Bruna, Manohar Paluri, Lubomir Bourdev, and Rob Fergus. 2015. [Training convolutional networks with noisy labels](#). In *3rd International Conference on Learning Representations, ICLR 2015*.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2016. [Rethinking the inception architecture for computer vision](#). *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826.
- Michael Tănzer, Sebastian Ruder, and Marek Rei. 2021. [BERT memorisation and pitfalls in low-resource scenarios](#). *CoRR*, abs/2105.00828.
- Brendan van Rooyen, Aditya Krishna Menon, and Robert C. Williamson. 2015. [Learning with symmetric label noise: The importance of being unhinged](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 10–18.
- Yu Yao, Tongliang Liu, Bo Han, Mingming Gong, Jiankang Deng, Gang Niu, and Masashi Sugiyama. 2020. [Dual T: reducing estimation error for transition matrix in label-noise learning](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 649–657.

A Noise Matrix on Yorùbá and Hausa

The training and validation sets of Yorùbá and Hausa have two sets of labels: the human-annotated (clean) labels and labels obtained from weak supervision. This makes it possible to compute the ground truth noise matrix in the training set. The noise matrices in Yorùbá and Hausa are shown in Figure 5. The Noise Matrix method evaluated in section 4 uses these two matrices for initialization. The labeling rules in the weak supervision are described in (Hedderich et al., 2020). The Yorùbá dataset has a rather low noise level, and the diagonally-dominant noise assumption holds in the training set. Oppositely, the Hausa training set is quite noisy. For the label “nigeria” the wrong labels is overwhelming, violating the diagonally-dominant noise assumption. Label “politics” is often misrecognized as “nigeria”. Moreover, many labels are misrecognized as the label “world”, making an unbalanced classification dataset. These factors make it very challenging to conquer the noise in this dataset.

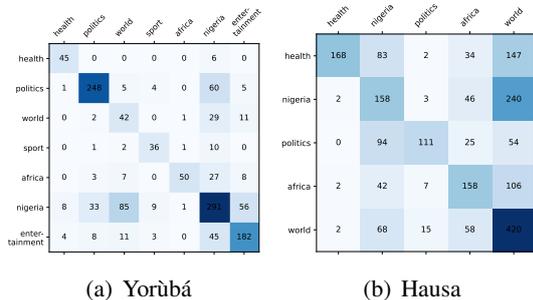


Figure 5: Noise matrix constructed from the Yorùbá (Hausa) training set.

B Comparing Early-stopping on Clean and Noisy Validation Sets

We compare the difference in model performance when using a noisy validation set rather than the clean one. Table 2 presents the results on datasets with injected noise. For a noise level below 60% uniform noise or 40% single flip-noise, we see the difference is often less than 0.5%, indicating that a noisy validation set can already serve as a good estimator for the generalization error. In an even higher noise level, the difference can be up to 2.14%. As for the datasets obtained from weak supervision, the difference is higher in general. Table 3 summarizes the difference on the Yorùbá and Hausa.

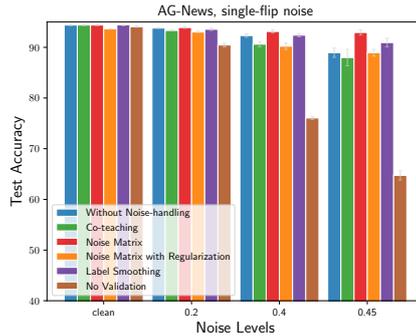


Figure 6: Test accuracy (%) on AG-News dataset with single-flip noise.

C BERT Performance on Different Datasets and Noise Settings

We evaluate the baselines under different noise settings and different datasets. The full result is shown in Table 4 and Table 5. A visualization of the result on AG-News with single-flip noise can be found in Figure 6 (other plots can be found in the main paper). BERT clearly shows its robustness against injected noise. Although noise-handling methods do help under a high noise level, the effect is limited (less than 4%). Compared to injected noise, the noise from weak supervision is much more challenging for BERT, especially on the Hausa dataset. For both noise types, there is no single noise-handling method that outperforms the simple baseline method without noise-handling in all settings.

D More ROC Curves

We present additional ROC curves under different settings with injected noise in Figure 7. It is obvious that the AUC decreases when the noise levels increase. However, the absolute AUC score remains at a high level even under extremely high noise levels of injected noise.

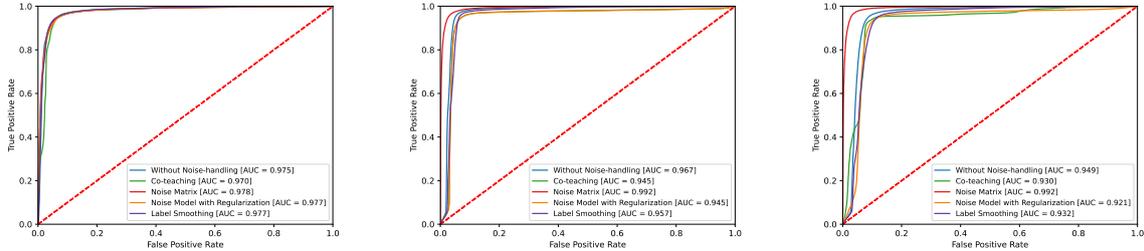
E Implementation Details

Datasets We experiment with the following four datasets: AG-News, IMDB, Yorùbá and Hausa.

1. AG-News: originates from AG, which is a large collection of news articles. Zhang et al. (2015) constructed the AG-News dataset from the AG collection and it is used as a benchmark dataset for text classification.
2. IMDB: consists of movie reviews with binary labels. It is a commonly used benchmark

Performance Difference (%)	AG-News						IMDB		
	uniform			single-flip			single-flip		
	40%	60%	70%	20%	40%	45%	20%	40%	45%
	0.10±0.09	0.56±0.50	1.96±0.97	0.06±0.07	0.29±0.19	2.00±0.60	0.14±0.19	1.71±2.05	1.76±2.79

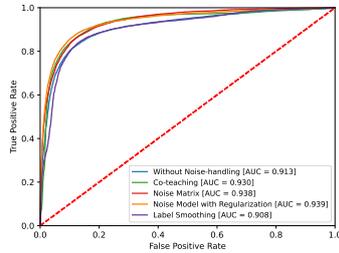
Table 2: Average performance difference (in %) and standard deviation (5 trials) between the test accuracy based on early-stopping with the clean validation and with the noisy validation set.



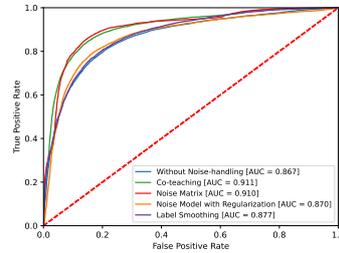
(a) AG-News - 60% uniform noise

(b) AG-News - 40% single-flip noise

(c) AG-News - 45% single-flip noise



(d) IMDB - 40% single-flip noise



(e) IMDB - 45% single-flip noise

Figure 7: The losses are record at the training step when early-stopping is triggered. noise-handling methods do not make the losses of correct and incorrect labels more distinguishable.

	Yorùbá		Hausa	
	FT	TP+FT	FT	TP+FT
Performance Difference (%)	1.93±1.71	1.00±0.70	1.08±0.79	1.92±1.64

Table 3: Average difference (in %) and standard deviation (10 trials) between the test accuracy based on early-stopping on the clean validation and on the noisy validation set.

dataset used for text classification.

3. Yorùbá: The dataset was created from BBC Yorùbá news titles along with the noisy dataset (Hedderich et al., 2020).
4. Hausa: Similar to Yorùbá, the Hausa dataset and the corresponding noisy dataset were created by Hedderich et al. (2020) from VOA Hausa news titles and distant supervision using keywords.

Models We use the official BERT-base model (Devlin et al., 2019) for text classification on AG-News and IMDB. It consists of an embedding layer, a 12-layer encoder, and a pooling layer. It contains

110M parameters in total. We use the multilingual version of BERT-base for text classification on Yorùbá and Hausa. It has the same architecture as the original BERT-base model. It also has 110M training parameters.

Fine-tuning on Text Classification Task For the vanilla models (Without Noise-handling and No Validation models in the paper), we pass the final layer of the [CLS] token representation (\mathbb{R}^{768}) to a feed forward layer for prediction. Noise Matrix and Noise Matrix with Regularization append a noise matrix $N \in \mathbb{R}^{k \times k}$ after the model’s prediction. For Noise Matrix we initialize the matrix with the ground truth information. Following (Jindal et al., 2019), when applying Noise Matrix with Regularization, we initialize the noise matrix using an identity matrix. The hyper-parameters for Noise Matrix with Regularization, Co-teaching and Label Smoothing are chosen so that the model performs the best on the noisy validation set.

In all settings, a batch size of 32 is used, and the

	AG-News							IMDB			
	clean	uniform			single-flip			clean	single-flip		
		40%	60%	70%	20%	40%	45%		20%	40%	45%
NV	94.07±0.13	84.48±0.78	61.61±3.18	43.78±5.07	90.46±0.37	76.06±0.33	64.74±0.94	94.03±0.13	86.34±0.77	65.05±0.90	58.97±1.26
CT	-	92.18±0.21	89.90±0.38	84.74±2.56	93.33±0.12	90.62±0.53	87.99±1.64	-	92.32±0.27	89.36±0.67	83.77±3.88
NMat	-	92.25±0.14	89.91±0.48	83.9±1.87	93.91±0.15	93.13±0.31	92.93±0.51	-	92.07±0.21	87.13±0.44	78.82±1.37
NMwR	93.64±0.06	92.02±0.20	89.91±0.33	84.77±2.24	93.03±0.17	90.23±0.65	88.93±0.68	93.68±0.14	92.12±0.35	85.94±0.86	80.17±2.57
LS	94.43±0.19	92.45±0.21	89.79±0.38	86.64±0.78	93.56±0.23	92.40±0.33	90.94±0.86	94.06±0.09	92.13±0.43	87.22±1.39	80.61±2.48
WN	94.40±0.13	92.40±0.25	89.53±0.75	85.49±0.76	93.80±0.08	92.33±0.35	88.94±0.92	93.98±0.15	92.13±0.21	85.88±2.78	80.12±4.09

Table 4: Average test accuracy (%) and standard deviation (5 trials) on AG-News and IMDB with uniform and single-flip noise. NV: without noise-handling and no validation set, i.e. train the model without noise-handling and until the training loss converges. CT: Co-teaching. NMat: Noise Matrix. NMwR: Noise Matrix with Regularization. LS: Label Smoothing. CT and NMat are equivalent to WN in the clean setting. Note that as IMDB is a binary-classification task, single-flip noise is equivalent to the uniform noise in this case.

	Yorùbá		Hausa	
	clean	noisy	clean	noisy
NV	74.11±0.26	63.88±1.59	83.02±0.45	46.98±1.01
CT	-	61.37±1.58	-	31.65±2.71
NMat	-	65.96±0.81	-	46.58±0.88
NWwR	73.78±0.32	61.32±0.71	83.21±0.40	35.36±3.60
LS	74.22±0.37	65.44±1.67	83.44±0.35	46.44±0.78
WN	74.45±0.32	64.72±1.45	83.55±0.47	46.97±0.81

Table 5: Average test accuracy (%) and standard deviation (10 trials) on Yorùbá and Hausa with noise from weak supervision. NV: without noise-handling and no validation set, i.e. train the model without noise-handling and until the training loss converges. CT: Co-teaching. NMat: Noise Matrix. NMwR: Noise Matrix with Regularization. LS: Label Smoothing.

	Average Runtime (Hours)			
	AG-News	IMDB	Yorùbá	Hausa
CT	5	4.5*	0.1*	0.1*
NMat	2.5	8	0.1*	0.1*
NMwR	3	8	0.1*	0.1*
LS	2.5	8	0.1*	0.1*
WN	2.5	8	0.1*	0.1*

Table 6: Average runtime (in hours) of each method. The Numbers with “*” indicates that the experiment was run on a Nvidia Tesla V100. Other experiments were run on a Nvidia GeForce GTX TITAN X.

learning rate is set to $2e-5$. We train all models until the training loss converges. However, we report the score where the model performs the best on the validation set during training except for the No Validation baseline where we report the last-epoch performance.

Hardware and Average Runtime We use Nvidia Tesla V100 and Nvidia GeForce GTX TITAN X to accelerate training. The average runtime for each method and dataset is summarized in Table 6.